

MCluster-VAEs: an end-to-end variational deep learning-based clustering method for subtype discovery using multi-omics data

Zhiwei Rong

Peking University

Jiali Song

Peking University

Lei Cao

Harbin Medical University

Zhilin Liu

Peking University

Yipei Yu

Peking University

Mantang Qiu

Peking University People's Hospital

Yan Hou (✉ houyan@bjmu.edu.cn)

Peking University

Research Article

Keywords: cancer subtype discovery, multi-omics data integration, cluster, deep learning, variational bayes

Posted Date: May 25th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1668552/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **MCluster-VAEs: an end-to-end variational deep**
2 **learning-based clustering method for subtype**
3 **discovery using multi-omics data**

4
5 Zhiwei Rong^a, Jiali Song^a, Lei Cao^b, Zhilin Liu^a, Yipei Yu^a, Mantang Qiu^{c,*} and Yan
6 Hou^{a,d,*}

7
8 ^a Department of Biostatistics, School of Public Health, Peking University, Beijing
9 100191, China

10 ^b Department of Epidemiology and Biostatistics, School of Public Health, Harbin
11 Medical University, Harbin 150086, China

12 ^c Department of Thoracic Surgery, Peking University People's Hospital, No.11
13 Xizhimen South Street, Xicheng District, Beijing 100044, China

14 ^d Peking University Clinical Research Center, Peking University, 38 Xueyuan Rd.,
15 Haidian District, Beijing, China

16 * To whom correspondence should be addressed.

17 Yan Hou, Tel: +86 (10)82805952; E-mail: houyan@bjmu.edu.cn.

18 Mantang Qiu, Tel: +86 (10)88325983; E-mail: qiumantang@163.com.

1 **Abstract**

2 **Background:** The discovery of cancer subtype based on unsupervised clustering helps
3 provide precise diagnoses, guide treatment and improve patients' prognoses. Instead of
4 single-omics data, multi-omics data can improve performance of the clustering because
5 it obtains a comprehensive landscape for understanding biological systems and
6 mechanisms. However, heterogeneous data from multiple sources raises high
7 complexity and different kinds of noise, which will be detrimental to the extraction of
8 clustering information.

9 **Methods:** We propose an end-to-end deep learning-based method, Multi-omics
10 Clustering Variational Autoencoders (MCluster-VAEs), that can extract cluster-friendly
11 representations on multi-omics data. First, unified network architecture with an
12 attention mechanism is developed for modeling multi-omics data precisely. Then, using
13 a novel objective function built from the Variational Bayes technique, the model is
14 trained to effectively obtain the posterior estimation of clustering assignments.

15 **Results:** Compared with twelve other state-of-the-art multi-omics clustering methods,
16 MCluster-VAEs achieved outstanding performance on benchmark datasets from the
17 TCGA database. On the Pan Cancer dataset, MCluster-VAEs achieved adjusted Rand
18 index of around 0.78 for cancer category recognition, an increase of more than 18%
19 compared with other methods. Furthermore, the survival analysis and clinical
20 parameters enrichment tests on ten cancer datasets demonstrate that MCluster-VAEs
21 delivered comparable or even better results than many typical integrative methods.

22 **Conclusions:** These results demonstrate that MCluster-VAEs is a new powerful tool

- 1 for dissecting complex multi-omics relationships and providing new insights for cancer
- 2 subtype discovery.
- 3 Key words: cancer subtype discovery; multi-omics data integration; cluster; deep
- 4 learning; variational bayes
- 5

1 **1. Background**

2 Cancer being a highly complex and heterogeneous genomics disease, showcases
3 variability in tumor responsiveness to the therapy (1). This problem forms the basis of
4 one of the critical areas of cancer research, i.e. the development of excellent subtype
5 discovery methods. These methods are designed to decouple the heterogeneity of cancer,
6 and accordingly, divide cancer patients into different groups to better understand the
7 pathogenesis of cancer and to boost clinical treatment (2). Since the occurrence and
8 development of cancer is related to several different biological layers and molecular
9 systems, it is more pertinent to identify cancer subtypes based on multi-omics data
10 rather than single-omics data (3).

11 The most commonly used methods to recognize cancer subtypes are clustering, such as
12 hierarchical clustering (4), spectral clustering (5) and K-means (6). These methods,
13 designed for single source data, cannot address the dimensional redundancy and
14 heterogeneous information integration brought by multi-omics data. Lately, the advent
15 of high-throughput technologies has proliferated multi-omics data, leading to the
16 development of multi-omics clustering methods (7). These methods are based on either
17 additional regularization controlling the problem of dimensionality (3,8,9), the fusion
18 of similarity networks (10–14), the factorization of multiple matrices (15,16), or the
19 simple linear probabilistic model (17,18). However, all these methods can only deal
20 with simple hypotheses to describe molecular systems, and cannot accurately depict the
21 complex regulation of multi-omics data. They still face significant challenges of data

1 complexity.

2 More recently, deep learning is emerging rapidly in many fields, exhibiting innovative
3 performance in processing images (19), texts (20) and graphs data (21). Several deep
4 learning-based methods have been developed to try to solve multi-omics clustering
5 tasks, such as autoencoder (AE) (22,23), variational autoencoder (VAE) (24,25) and
6 Subtype-GAN (26). They use non-linear neural networks to learn an integrated
7 representation of multi-omics data by the unsupervised framework and then apply a
8 traditional clustering algorithm to this representation. We call them “two-steps”
9 methods (Figure 1A (a)), because their entire process includes two separate steps:
10 representation learning and clustering. Due to the substantial capacity of neural
11 networks, the integrated representations contain rich high-level information. Thus, the
12 clustering step is expected to obtain improved performance. However, representation
13 learning is independent of the following clustering step, and the representation from the
14 first step is not guaranteed to be suitable for the second. For example, most
15 representation learning approaches try to squeeze all reconstruct-friendly information
16 into lower-dimensional representations, and these representations contain a lot of
17 information unrelated to clustering (27). The cluster-friendly information will be
18 overwhelmed and interfered by cluster-unfriendly information in the following
19 clustering step. Thus, if we can integrate these two steps as one step (Figure 1A (b)),
20 the representation learning will be guided by the clustering target, and cluster-friendly
21 representation will be extracted to improve clustering performance. This “one-step”
22 idea has been studied in the field of single-modal deep learning (28–30), but not in the

1 field of multi-omics.

2 In this research, we propose Multi-omics Clustering Variational Autoencoders
3 (MCluster-VAEs), an end-to-end deep learning-based method for clustering multi-
4 omics data. It describes multi-omics data using a new probabilistic model with a global
5 discrete latent variable considered as the clusters. Using Variational Bayes approach
6 (31), we derive a unified end-to-end architecture and a novel objective function, to fit
7 the multi-omics data and infer the posterior probability of the clustering assignments.
8 They form a “one-step” framework, so the clustering target can guide representation
9 learning and improve clustering performance. In addition, an attention module is used
10 to effectively incorporate multiple omics and reveal the contribution of each omics to
11 clustering results. To evaluate the performance of MCluster-VAEs, we collected the Pan
12 Cancer dataset with 8,314 samples and the ten cancer datasets with a total of 4,154
13 samples from TCGA database. The results generated by MCluster-VAEs were fully
14 validated, which showed great potential for novel cancer subtype discovery from deep
15 learning model.

16 **Figure 1.** Overview of MCluster-VAEs. (A) Two deep learning-based multi-omics clustering
17 pipelines. (B) The data generative model. (C) The posterior variational inference path. (D) The
18 architecture of MCluster-VAEs model for the experiments using mRNA expression, miRNA
19 expression, copy number alterations (CNA) and DNA methylation.

20 **2. Materials and methods**

21 **2.1 Probabilistic model with discrete latent variables**

22 Let us consider some multi-omics dataset $\mathbf{X} = \{\mathbf{X}^1, \dots, \mathbf{X}^M\}$ consisting of N i.i.d.

1 samples, and each sample contains M kinds of omics data which can be expressed as
2 $\mathbf{x}_i = \{\mathbf{x}_i^1, \dots, \mathbf{x}_i^M\}$, $i = 1, \dots, N$. The vector \mathbf{x}_i^m represents the features of omics m of
3 sample i . We assume that the data are generated by some random process, involving
4 an unobserved discrete random variable \mathbf{y} with C categories and M unobserved
5 continuous random variables $\mathbf{z} = \{\mathbf{z}^1, \dots, \mathbf{z}^M\}$. The process consists of three steps: (1)
6 a categorical value y_i is generated from a prior distribution $p(y)$; (2) a value \mathbf{z}_i^m is
7 assumed to follow a mixture of C Gaussian distributions and generate from
8 distribution $p_\theta(\mathbf{z}^m | y_i)$ corresponding to omics m ; (3) a value \mathbf{x}_i^m is generated
9 from a conditional distribution $p_\theta(\mathbf{x}_i^m | \mathbf{z}_i^m, y_i)$. This process can be expressed by the

10 Figure 1B and the following formula:

$$11 \quad p_\theta(\mathbf{x}_i) = \int_{\mathbf{z}_i^1, \dots, \mathbf{z}_i^M, y_i} \prod_{m=1}^M p_\theta(\mathbf{x}_i^m | \mathbf{z}_i^m, y_i) p_\theta(\mathbf{z}_i^m | y_i) p(y_i) d\mathbf{z}_i^m dy_i \quad (1)$$

12 We assume that $p_\theta(\mathbf{x}_i^m | \mathbf{z}_i^m, y_i)$ and $p_\theta(\mathbf{z}_i^m | y_i)$ come from parametric families of
13 distributions $p_{\theta'}(\mathbf{x} | \mathbf{z}, y)$ and $p_{\theta'}(\mathbf{x} | \mathbf{z})$, and that their probability density functions
14 (PDFs) are differentiable almost everywhere w.r.t both θ' , \mathbf{z} and y . The above
15 process describes a realistic data generative model and has enough flexibility to fit the
16 complex multi-omics data. The unobserved y_i can be considered as the cluster
17 assignments of sample i with C categories, and the clustering task can be equivalent
18 to inference of the unknown parameters θ and posterior distribution
19 $p_\theta(y_i | \mathbf{x}_i^1, \dots, \mathbf{x}_i^M)$.

20 **2.2 Variational Bayes**

21 In order to solve the above inference task, we introduce a trainable parametric

1 approximation $q_\phi(y_i, \mathbf{z}_i | \mathbf{x}_i)$ of posterior distribution $p_\theta(y_i, \mathbf{z}_i | \mathbf{x}_i)$. The marginal
 2 likelihood is composed of a sum over the marginal likelihoods of individual samples
 3 $\log p_\theta(\mathbf{x}_1, \dots, \mathbf{x}_N) = \sum_{i=1}^N \log p_\theta(\mathbf{x}_i)$, which can each be rewritten as:

$$\begin{aligned}
 \log p_\theta(\mathbf{x}_i) &= D_{KL}(q_\phi(y_i, \mathbf{z}_i | \mathbf{x}_i) \| p_\theta(y_i, \mathbf{z}_i | \mathbf{x}_i)) + L(\theta, \phi; \mathbf{x}_i) \\
 &\geq L(\theta, \phi; \mathbf{x}_i) \\
 &= \mathbb{E}_{q_\phi(y_i, \mathbf{z}_i | \mathbf{x}_i)}(\log p_\theta(\mathbf{x}_i | y_i, \mathbf{z}_i)) - D_{KL}(q_\phi(y_i, \mathbf{z}_i | \mathbf{x}_i) \| p_\theta(y_i, \mathbf{z}_i))
 \end{aligned}
 \tag{2}$$

5 The term $L(\theta, \phi; \mathbf{x}_i)$ is called (variational) evidence lower bound (ELBO), which
 6 strictly less than the marginal likelihood which can be maximized indirectly by
 7 updating generative parameters θ and variational parameters ϕ . This approach was
 8 proposed first by Kingma and Welling (31), and has been applied in several generative
 9 models, including the famous VAE (31).

10 We can assume that $q_\phi(y_i, \mathbf{z}_i | \mathbf{x}_i)$ has the mean field form and can be decomposed into
 11 product form:

$$q_\phi(y_i, \mathbf{z}_i | \mathbf{x}_i) = q_\phi(y_i | \mathbf{x}_i) \prod_{m=1}^M q_\phi(\mathbf{z}_i^m | \mathbf{x}_i^m, y_i)
 \tag{3}$$

13 Figure 1C illustrates this inference path. Therefore, the ELBO can be derived in the
 14 following form:

$$\begin{aligned}
 L(\theta, \phi; \mathbf{x}_i) &= \underbrace{\sum_{m=1}^M \left(\mathbb{E}_{q_\phi(y_i | \mathbf{x}_i)} \left(\mathbb{E}_{q_\phi(\mathbf{z}_i^m | \mathbf{x}_i^m, y_i)} \left(\log(p_\theta(\mathbf{x}_i^m | y_i, \mathbf{z}_i^m)) \right) \right) \right)}_{\text{reconstruction term}} \\
 &\quad - \underbrace{\sum_{m=1}^M \left(\mathbb{E}_{q_\phi(y_i | \mathbf{x}_i)} \left(D_{KL}(q_\phi(\mathbf{z}_i^m | \mathbf{x}_i^m, y_i) \| p_\theta(\mathbf{z}_i^m | y_i)) \right) \right)}_{\text{conditional prior term}} \\
 &\quad - \underbrace{D_{KL}(q_\phi(y_i | \mathbf{x}_i) \| p(y_i))}_{\text{conditional entropy term}}
 \end{aligned}
 \tag{4}$$

16 The detailed derivation is given in Supplementary Note S1. The first Right-Hand-Side

1 (RHS) term is actually a measurement of reconstruction by samples of posterior
2 distributions. The last two terms are the Kullback-Leibler (KL) divergences of
3 approximated posterior and prior distributions. We refer to all RHS terms as the
4 reconstruction term, conditional prior term and conditional entropy term, respectively.
5 Now, with some reasonable and loose assumptions for \mathbf{z}_i^m and y_i , equation (4) can
6 be used as loss function and to optimize parameters by stochastic gradient descent (SGD)
7 with Monte Carlo gradient estimator. After the training completed, $q_\phi(y_i|\mathbf{x}_i)$, as the
8 approximation of posterior $p_\theta(y_i|\mathbf{x}_i)$, will be used to obtain the clustering assignment
9 of each sample.

10 2.3 Multi-omics Clustering Variational Autoencoders (MCluster-VAEs)

11 Just like VAE, we use some neural networks for $p_\theta(\mathbf{x}_i^m | y_i, \mathbf{z}_i^m)$, $p_\theta(\mathbf{z}_i^m | y_i)$, $q_\phi(y_i | \mathbf{x}_i)$
12 and $q_\phi(\mathbf{z}_i^m | \mathbf{x}_i^m, y_i)$ in equation (4). The prior over the latent variables y_i is assumed
13 as uniform categorical distribution $p_\theta(y_i) = 1/M$. We also assume that $p_\theta(\mathbf{z}_i^m | y_i)$
14 and $q_\phi(\mathbf{z}_i^m | \mathbf{x}_i^m, y_i)$ are multivariate Gaussian distributions with diagonal covariance,
15 whose distribution parameters (means and covariances) are computed from \mathbf{x}_i^m and
16 y_i with two multilayer perceptrons (MLP, fully-connected neural networks with some
17 hidden layers), respectively. $q_\phi(y_i | \mathbf{x}_i)$ is also a MLP with M output nodes with
18 softmax activation function to get categorical probabilities of each clustering label.
19 $p_\theta(\mathbf{x}_i^m | y_i, \mathbf{z}_i^m)$ is considered as multivariate Gaussian distribution for real-valued data,
20 whose means are calculated from \mathbf{z}_i^m and y_i through an additional MLP. The whole
21 model architecture is shown in figure 1D.

22 There are two expectations in the above ELBO term. The first is w.r.t discrete

1 distribution $q_\phi(y_i|\mathbf{x}_i)$, which can be calculated by iterating over all possible values of
 2 y_i . The second is w.r.t continuous distribution $q_\phi(\mathbf{z}_i^m|\mathbf{x}_i^m, y_i)$, which should be
 3 estimated by reparameterization sampling trick (31). To summarize all of the above, the
 4 loss function can be rewritten as the following form (the derivation is given in
 5 Supplementary Note S2):

$$\begin{aligned}
 \text{Loss} &= -\sum_{i=1}^N L'(\theta, \phi; \mathbf{x}_i) = L_{rec} + L_{cprior} + L_{centropy} \\
 L_{rec} &= \frac{1}{N} \sum_{i=1}^N \sum_{m=1}^M \sum_{c=1}^C \sum_{j=1}^{d_m^x} (x_{ijc}^m - x'_{ijc}{}^m)^2 \pi_{ic} \\
 L_{cprior} &= -\frac{1}{2N} \sum_{i=1}^N \sum_{m=1}^M \sum_{c=1}^C \sum_{j=1}^{d_m^z} \left(\log \left(\frac{(\sigma_{ijc}^m)^2}{(\sigma'_{ijc}{}^m)^2} \right) - \frac{(\sigma_{ijc}^m)^2}{(\sigma'_{ijc}{}^m)^2} - \frac{(\mu_{ijc}^m - \mu'_{ijc}{}^m)^2}{(\sigma'_{ijc}{}^m)^2} \right) \pi_{ic} \quad (5) \\
 L_{centropy} &= \frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C \pi_{ic} \log \pi_{ic}
 \end{aligned}$$

7 where $L'(\theta, \phi; \mathbf{x}_i)$ represents an estimator of $L(\theta, \phi; \mathbf{x}_i)$ using Gaussian
 8 reparameterization trick. d_m^x and d_m^z are the dimensions of original features and
 9 latent continuous embeddings of the omics m , respectively. $x'_{ijc}{}^m$ is the reconstruction
 10 of x_{ijc}^m , which is the value of feature j of omics m of sample i if its cluster
 11 assignment is c . μ_{ijc}^m and $\mu'_{ijc}{}^m$ are means of $q_\phi(\mathbf{z}_i^m|\mathbf{x}_i^m, y_i = c)$ and $p_\theta(\mathbf{z}_i^m|y_i = c)$,
 12 respectively. $(\sigma_{ijc}^m)^2$ and $(\sigma'_{ijc}{}^m)^2$ are variances of $q_\phi(\mathbf{z}_i^m|\mathbf{x}_i^m, y_i = c)$ and
 13 $p_\theta(\mathbf{z}_i^m|y_i = c)$, respectively. π_{ic} is the value of $q_\phi(y_i = c|\mathbf{x}_i)$.

14 2.4 Gated attention mechanism

15 $q_\phi(y_i|\mathbf{x}_i)$ is the main part of MCluster-VAEs, which merges all omics data together to
 16 infer the posterior probability of clustering assignments. The most common deep
 17 learning-based approach for multi-omics information integration is the *intermediate*
 18 *integration* (32), which first processes each omics data using dedicated layers,

1 concatenates the outputs of dedicated layers and then uses further layers to integrate the
 2 features extracted from each omics data. This approach enables the most suitable
 3 dedicated layers to be used for each omics data and can hence extract more predictive
 4 features.

5 However, the intermediate integration approach cannot consider the contribution of
 6 multiple omics data. In order to make more effective use of multi-omics data, we add
 7 an attention mechanism into the encoder $q_\phi(y_i|\mathbf{x}_i)$. The attention module is shown in
 8 figure 1D, which is implemented by a simple gated attention mechanism. Each omics
 9 data first maps into the hidden representation with the same dimension, then the
 10 representations are weighted sum according to attention scores calculated from the
 11 representations:

$$\begin{aligned}
 \mathbf{h}_i &= \sum_m^M s_i^m \mathbf{h}_i^m \\
 s_i^m &= \frac{1}{1 + \exp(-W^m \mathbf{h}_i^m)}
 \end{aligned}
 \tag{6}$$

13 where \mathbf{h}_i^m and s_i^m are the hidden representation and attention scores of omics m
 14 of sample i , respectively. \mathbf{h}_i represents the integrated information of all types of
 15 omics and then is used to predict clustering assignments.

16 **2.5 Conditional entropy annealing trick**

17 The most unusual term in the ELBO is the conditional entropy term. It tries to minimize
 18 the KL divergence of y-posterior and uniform prior. It seems to be in the opposite
 19 direction of clustering. However, it is the existence of this term which makes reasonable
 20 clustering possible. Actually, this term is a Bayesian regularization term against bias of

1 maximum likelihood, which enables every category of y_i to have enough chance to
2 join the training process and the model does not quickly fall into a very bad local
3 optimum.

4 In practice, model training often falls into collapse mode when the equation (5) is
5 directly considered as loss function. In this case, all clustering assignments tend to
6 become a same value. This problem may come from the enormous contribution of large
7 dimensional multi-omics data to the training of reconstruction term L_{rec} . To solve it,
8 we need to increase the weight of the conditional entropy term. Therefore, we propose
9 a conditional entropy annealing trick for experiments on real data, which multiplies the
10 conditional entropy term $L_{entropy}$ by a large value of γ at the beginning of training,
11 and then gradually reduces the γ during training to take into account all categories of
12 y_i :

$$13 \quad \text{Loss} = L_{rec} + L_{cprior} + \gamma L_{entropy} \quad (7)$$

14 The value of γ is described in Supplementary Note 3. In the following experiments,
15 we will show that this trick is very helpful in training of MCluster-VAEs.

16 **2.6 Gumbel softmax reparameterization**

17 The training speed of MCluster-VAEs is relatively slow, because all categories of y_i
18 need to be passed in model when calculating the loss function. One way to solve this
19 problem is to use reparameterization trick on the discrete variable y_i as well.

20 According to the techniques introduced by Jang and Maddison (33,34), we use Gumbel
21 softmax distribution to approximate the categorical distribution to achieve
22 reparameterization. The Gumbel softmax distribution can be sampled by the following

1 process:

$$2 \quad y_{ic} = \frac{\exp\left(\frac{(\log \pi_{ic} + g_c)}{\tau}\right)}{\sum_{k=1}^C \exp\left(\frac{(\log \pi_{ik} + g_k)}{\tau}\right)} \quad (8)$$
$$3 \quad g_k = -\log(-\log(u_k)), u_k \sim U(0,1), k = 1, \dots, C$$

4 where π_{ic} is the posterior probability of category c of sample i (the value of
5 $q_{\phi}(y_i = c | \mathbf{x}_i)$). τ is an annealing factor that decreases gradually during whole
6 training. The pseudo-code of MCluster-VAEs with and without the gumbel softmax
7 trick is shown in Supplementary Note S3.

7 **2.7 Datasets and preprocessing**

8 The mRNA expression, miRNA expression, DNA methylation (450K) and copy
9 number alterations (short for mRNA, miRNA, methylation, CNA) from TCGA
10 database were used in this study. MCluster-VAEs and comparison methods were tested
11 in the following two settings.

12 The Pan Cancer dataset consisted of a total of 8,314 samples across 32 types of cancers.
13 These cancer types originated from different types of tissues and dissection positions.
14 Therefore, the Pan Cancer dataset had natural categorical structure (32 cancer types)
15 and could be used to test the performance of clustering algorithms. We then
16 preprocessed this dataset using a 6-steps approach: 1. to filter samples to ensure that
17 they exist in all four omics; 2. to perform log transformation of mRNA and miRNA; 3.
18 to remove duplicated regions data of CNA and convert it into gene-level form; 4. to
19 select the features using Yang's approaches (26), retaining 3105 CNA features, 3217

1 mRNA features, 383 miRNA features and 3139 methylation features; 5. to perform
2 missing data imputation by the samples' mean value; 6. To standardize features in each
3 omics by removing the mean and scaling the feature to unit variance. The sample sizes
4 and abbreviations of 32 cancer types were included in the Supplementary Table S1.
5 Secondly, ten cancer types in the TCGA dataset were applied in this study (including
6 4154 tumors): 1031 BRCA tumors, 399 BLCA tumors, 488 KIRC tumors, 127 GBM
7 tumors, 490 LUAD tumors, 176 PAAD tumors, 446 SKCM tumors, 407 STAD tumors,
8 510 UCEC tumors and 80 UVM tumors. These datasets were used to test survival and
9 clinical differences of clustering assignments, which have been used in Yang's study
10 (26). These datasets with the above-mentioned four omics were preprocessed by the
11 above 6-steps approach.

12 **2.8 Evaluation criteria and comparison algorithms**

13 Since the cancer types were known, the Pan Cancer dataset was tested to determine
14 whether the clustering algorithms could correctly identify the labels. It was measured
15 by unsupervised clustering accuracy (ACC), adjusted Rand index (ARI), normalized
16 mutual information (NMI), and clustering F measure (F1). They were defined in
17 Supplementary Note S4.

18 On the ten single cancer datasets, the significance of survival analysis and the number
19 of enriched clinical parameters were used to measure performance of MCluster-VAEs
20 and the comparison algorithms. It assumed that if clusters of samples exhibit significant
21 difference in clinical parameters and survival, they are different in a biologically
22 meaningful way. The clinical parameters included gender, age at diagnosis, pathologic

1 T, pathologic M, pathologic N and pathologic stage. The four latter parameters were
2 discrete pathological parameters, measuring the progression of the tumor (T),
3 metastases (M) and cancer in lymph nodes (N), and the total progression (stage). We
4 measured significance of survival analysis among the obtained clustering assignments
5 using the log-rank test (35), enrichment for discrete parameters using the chi-square
6 test, and enrichment for numeric parameters using Kruskal-Wallis test. According to
7 Rappoport and Shamir (36), the accuracy of pure log-rank test, chi-square test and
8 Kruskal-Wallis test might be influenced by small sample size and unbalanced cluster
9 assignments. Therefore, we followed their advice and used the corresponding
10 permutation tests to estimate empirical P values. More details on the permutation tests
11 can be found in their study (36). After that, the P values were corrected for multiple
12 hypotheses using Bonferroni correction for each cancer and corresponding method.
13 Note that cancer subtypes that are biologically different may have similar survival, and
14 this is also true for enrichment of clinical parameters. However, these measures are
15 widely used for clustering assessment, including in the most multi-omics clustering
16 studies (12,26,36–38).

17 To evaluate the clustering performance of MCluster-VAEs, we compared its
18 performance with the performances of twelve state-of-the-art clustering methods (i.e.
19 *k-means* (6), *spectral clustering* (5), *MCCA* (15), *SNF* (10,11), *COCA* (39), *ANF* (40),
20 *iClusterBayes* (3), *CIMLR* (41), *NEMO* (12), *MAUI(VAE)* (24,25), *DCAP(AE)* (22,23),
21 *SubtypeGAN* (26)). *MUAI* and *DCAP* used VAE and AE with *intermediate integration*,
22 respectively. Hence, *MUAI(VAE)* and *DCAP(AE)* were used to represent them. These

1 methods were chosen to represent diverse approaches to multi-omics clustering. The
 2 detailed information and the parameter setting of all methods were shown in
 3 Supplementary Note S5, S6, Table S2 and S3.

4 **3. Results**

5 **3.1 Comparison of MCluster-VAEs with state-of-the-art clustering methods on the** 6 **Pan Cancer dataset**

7 **Table 1.** Clustering performance on the Pan Cancer dataset using MCluster-VAEs and twelve other
 8 methods.

Method	ACC	ARI	F1	NMI
iCluster	0.0569±0.0022	0.0014±0.0004	0.0386±0.0004	0.0231±0.0012
MCCA	0.1275±0.0026	0.0557±0.0032	0.1067±0.0053	0.1320±0.0031
SNF	0.4809	0.3223	0.3659	0.6050
Spectral	0.5635	0.4423	0.4690	0.6857
COCA	0.5459±0.0099	0.5051±0.0030	0.5264±0.0029	0.6732±0.0042
DCAP(AE)	0.6656±0.0256	0.5267±0.0347	0.5482±0.0324	0.7725±0.0112
SubtypeGAN	0.6334±0.0374	0.5335±0.0505	0.5526±0.0484	0.7311±0.0265
CIMLR	0.6342±0.0100	0.5355±0.0209	0.5561±0.0197	0.7337±0.0081
MAUI(VAE)	0.6750±0.0241	0.5380±0.0528	0.5587±0.0494	0.7669±0.0079
K-means	0.6507±0.0308	0.5738±0.0275	0.5922±0.0264	0.7494±0.0124
ANF	0.6924±0.0672	0.6134±0.0991	0.6314±0.0965	0.7944±0.0843
NEMO	0.7475	0.6642	0.6782	0.7966
MCluster-VAEs	0.8200±0.0152	0.7826±0.0291	0.7924±0.0279	0.8789±0.0088

9 The proposed MCluster-VAEs and the state of art methods were applied to the Pan
 10 Cancer dataset labeled with known cancer types, with clusters set equal to the number
 11 of cancer types (i.e., 32). The clustering performance (ACC, ARI, F1, and NMI) was
 12 shown in Table 1. Except for SNF, Spectral, and NEMO, which produced deterministic
 13 results, all methods were repeated five times to avoid the influence of randomness. The
 14 results obtained by MCluster-VAEs were 0.8200 (ACC), 0.7826 (ARI), 0.7924 (F1),

1 and 0.8789 (NMI), which showed the best performance compared with the other
2 methods. Even when compared to the second-placed NEMO, there was about 10-18%
3 improvement. Moreover, MCluster-VAEs showed more minor variation than other deep
4 learning-based methods (DCAP(AE), MAUI(VAE), SubtypeGAN) and were at the
5 same level as k-means. The correspondence between clustering assignments and cancer
6 types labels was shown in Supplementary Figures S1 and S2, demonstrating that the
7 performance improvement of MCluster-VAEs is reflected in almost all cancer types,
8 not just in individual cancer types. These findings indicated that the performance of
9 MCluster-VAEs was better than those of the twelve state-of-the-art clustering methods
10 on the Pan Cancer dataset.

11 **Figure 2.** The t-SNE visualization on the latent variables generated by four deep learning-based
12 methods (MCluster-VAEs, AE, VAE, SubtypeGAN) used the Pan Cancer dataset. The color of each
13 point indicated the cancers. MCluster-VAEs had a greater degree of separation.

14 As mentioned earlier, the excellent performance of MCluster-VAEs may come from its
15 “one-step” framework, which guides the model to learn cluster-friendly representation.
16 In order to prove it, we visualized the representations learned of MCluster-VAEs,
17 DCAP(AE), MAUI(VAE), and SubtypeGAN through t-SNE (42). For MCluster-VAEs,
18 the representations were the output of the second last layer of encoder $q_\phi(y_i|\mathbf{x}_i)$. For
19 DCAP(AE), MAUI(VAE), and SubtypeGAN, the representations were the outputs of
20 the encoder. As shown in Figure 2, the representations of MCluster-VAEs had more
21 remarkable dissociation among different cancer types. For example, although data
22 points representing HNSC, LUSC, CESC, and BLCA (red circle) tended to mix in
23 subfigures of DCAP(AE), MAUI(VAE), and SubtypeGAN, MCluster-VAEs had a

1 greater degree of separation. In addition, the points representing STAD tended to mix
2 with the points representing COAD for these three “two-steps” methods (light blue
3 circle). The confusion could be caused by the fact that these two kinds of cancers are
4 all digestive tract cancer. However, the representations learned by MCluster-VAEs
5 separated STAD from COAD. MCluster-VAEs also made some improvements in LIHC
6 and CHOL (green circle). The above results suggested that MCluster-VAEs could learn
7 cluster-friendly representations and better clusters than “two-step” methods.

8 **Figure 3.** Distribution of metrics and attention scores of MCluster-VAEs using single omics data of
9 the Pan Cancer datasets. **A.** ACC, ARI, NMI and F1 of MCluster-VAEs based on all four omics
10 (MOmics) or single-omics data. **B.** The distribution of attention scores for each omics data. Here,
11 mRNA represents mRNA expression, methy denotes DNA methylation (450K), miRNA represents
12 miRNA expression and CNA represents copy number alterations.

13 To determine the necessity of usage of multi-omics data, we compared the clustering
14 results of MCluster-VAEs based on four single-omics data and multi-omics
15 data. MCluster-VAEs removed the attention module when processing single-omics data
16 because the fusion of multiple omics features was no longer required, while other
17 parameter settings remained unchanged. As shown in Figure 3A, CNA alone was
18 insufficient for accurate cancer classification. In comparison to methylation, mRNA,
19 and miRNA data alone, the combined use of multi-omics data increased ARI of
20 approximately 18%, 19%, and 33%, respectively. In addition, the ARI of the multi-
21 omics data had a smaller distribution range than the ARI values of single-omics data.
22 ACC, F1, and NMI's conclusions are consistent with ARI's. These results indicate that
23 MCluster-VAEs can give more accurate and stable clustering results using multi-omics
24 information compared with single-omics datasets.

1 3.2 The role and influence of attentional mechanism

2 How to integrate multiple source data is a major challenge in multi-omics studies. Deep
3 learning-based methods have an inherent advantage, because the neural networks can
4 contain multiple independent layers to extract appropriate information for each omics,
5 and the shared layer then fuse these information. This architecture, called *intermediate*
6 *integration*, was easy to implement and has been used in many studies (22,25,26).
7 However, different omics data contained varying amounts of clustering information,
8 making it difficult for the *intermediate integration* architecture to identify these
9 differences and lack interpretation of omics contribution. The solution provided by this
10 study was the attention mechanism, which adaptively learned the weight of each omics
11 data for each sample.

12 The distribution of the attention scores of each omics was shown in Figure 3B. It
13 indicated that when clustering Pan Cancer dataset, MCluster-VAEs place a greater
14 emphasis on mRNA, miRNA, and methylation while giving CNA less weight. It is
15 worth noting that the distribution of attention scores matched that of clustering metrics
16 on single omics data (Figure 3A) and that the metrics can be used to measure the amount
17 of clustering information in each omics data to some extent. This similarity showed that
18 the attention module allowed the model to focus more on the omics with more
19 clustering information, allowing it to extract information more efficiently. We also
20 implemented the non-attention version of MCluster-VAEs on the Pan Cancer dataset,
21 compared with the standard version (with attention mechanism) of MCluster-VAEs.
22 The results were shown in Supplementary Figure S3 and indicated that the performance

1 of the attention version was slightly better than the non-attention version (ARI:
2 0.7647 ± 0.0257 vs. 0.7523 ± 0.0188 , ACC: 0.7978 ± 0.0154 vs. 0.7704 ± 0.0167). Thus,
3 using an attention module does not affect clustering performance. Simultaneously, the
4 attention module improved the interpretability of the clustering results, which helped
5 investigate the biological significance of cancer subtypes. For example, methylation
6 has a more negligible effect on PRAD and BRCA than on other cancers, whereas CNA
7 has a more significant effect on OV (Supplementary Figure S4). In conclusion, the
8 above results showed that the attention mechanism could help with more explainable
9 and accurate multi-omics data integration.

10 **3.3 Gumbel Softmax and conditional entropy weighted annealing**

11 The exact form of MCluster-VAEs was time-consuming, especially when the number
12 of clustering was huge (such as 32 of the Pan Cancer dataset). The reason was that
13 MCluster-VAEs must traverse all possible clustering assignments and run
14 $q_{\phi}(\mathbf{z}_i^m | \mathbf{x}_i^m, y_i)$ and $p_{\theta}(\mathbf{z}_i^m | y_i)$ repeatedly. The application of the Gumbel Softmax
15 reparameterization trick would make the model just run one or few times to calculate
16 the loss function. The running time of the Gumbel version and exact version was
17 reported in Supplementary Figure S5, and we could see that the Gumbel version was
18 more than 5 times faster than the exact version in the Pan Cancer dataset. The Gumbel
19 version also had a better performance than the exact version (Figure 4). This
20 improvement of ACC might be because this reparameterization trick improved gradient
21 variance and gave more space for exploration.

22 Another challenge of multi-omics data analysis was the curse of dimensionality, which

1 also appeared in single-omics analysis but would be more severe in multi-omics
2 analysis because of the integration of multiple data sources. One of its influences for
3 MCluster-VAEs was that the effect of the reconstruction term overwhelmed the KL
4 divergence regularization. As a result, it quickly led $q_\phi(y_i | \mathbf{x}_i)$ to converge into a
5 unique clustering assignment during the early stage of training. This is a “Lazy behavior”
6 of model training, resulting in convergence to a bad local minimum.

7 **Figure 4.** The monitoring records of the training process of MCluster-VAEs on the Pan Cancer
8 dataset. “WA” means using conditional entropy weight decay annealing trick. “Gumbel” means
9 using Gumbel Softmax reparameterization trick. “none” means neither used.

10 To solve this problem, we proposed the conditional entropy weight annealing trick. The
11 γ factor was used to provide obvious gradient signals for each clustering category
12 during backpropagation learning. Figure 4 showed the entropy, conditional entropy and
13 ACC during the whole training period in the Pan Cancer dataset. When the weight
14 annealing was not applied, conditional entropy would decline sharply to zero in first
15 10-20 epochs. Correspondingly, the entropy was soon stop at a lower level until training
16 end. When using the weight annealing trick, conditional entropy would still decline to
17 zero but more slowly, the entropy would rise to a higher level, and the ACC would also
18 rise up to a higher value. Unexpectedly, the use of Gumbel softmax reparameterization
19 also eased the convergence problem, which may be due to the reparameterization that
20 improved efficiency of gradient backpropagation. Overall, by using the conditional
21 entropy weight annealing trick and Gumbel Softmax trick, the convergence problem of
22 MCluster-VAEs could be solved effectively.

1 3.4 Integrative clustering across ten cancer types with thirteen methods

2 We applied MCluster-VAEs and the other twelve state-of-the-art clustering methods on
3 the ten TCGA multi-omics datasets. The clustering assignments and clinical parameters
4 were used to calculate corrected empirical P values of survival analysis and enrichment
5 test (36). To ensure the fairness of comparison, we followed the setup of Yang's study
6 (26) and used the reasonable number of clusters obtained from the previous studies
7 (2,43–48) as the setting for all methods (BRCA: 5, BLCA: 5, KIRC: 4, SKCM: 4,
8 UCEC: 4, UVM: 4, GBM: 3, LUAD: 3, STAD: 3, PAAD: 2).

9 **Figure 5.** Performance of the algorithms on the ten cancer datasets. The x-axis was the number of
10 clinical parameters enriched in the clusters, and the y-axis measured the differential survival
11 between clusters ($-\log_{10}$ of permuted logrank's test P values). Colors indicated clustering methods.

12 The $-\log_{10}$ P -values of differential survival and the number of enriched clinical
13 parameters of all methods on the ten datasets were shown in Figure 5 (the values were
14 further reported detailly in Supplementary Tables S4 and S5). On all nine datasets
15 except LUSC, MCluster-VAEs had the smallest log-rank P value and the most clinical
16 parameters that were enriched. Even on the LUSC dataset, MCluster-VAEs still had the
17 best survival analysis results and the second-highest number of clinical variables after
18 SubtypeGAN. The mean of $-\log_{10}$ log-rank P values was 5.850 ($-\log_{10} 0.05 = 1.301$)
19 obtained by MCluster-VAEs, while the second-best method (SubtypeGAN) achieved
20 3.511 and the third-best method (SNF) achieved 3.485. Especially on the KIRC and
21 SKCM datasets, the $-\log_{10}$ P values obtained by MCluster-VAEs were bigger than 10
22 ($P < 1 \times 10^{-10}$). Figure 6 showed all the Kaplan-Meier survival plots of MCluster-VAEs,
23 indicating that MCluster-VAEs can reveal cancer datasets' survival differences.

1 **Figure 6.** Kaplan-Meier survival plots of MCluster-VAEs on the ten cancer datasets.

2 For the number of enriched clinical parameters, the MCluster-VAEs achieved the mean
3 of 4.0, while the second-best method (SubtypeGAN) achieved 3.4 and the third-best
4 method (SNF) achieved 2.9. For example, on the STAD datasets, four clinical
5 parameters' enrichment tests of MCluster-VAEs were significant, followed by three
6 significant results of COCA. On the PAAD dataset, MCluster-VAEs had four enriched
7 clinical parameters, whereas the other methods had no more than one. MCluster-VAEs
8 also obtained better results than the comparison methods on the other datasets.
9 Therefore, we can conclude that our method outperformed most of the typical methods
10 on many TCGA datasets.

11 The MCluster-VAEs was also applied to single-omics data of the ten single cancer
12 datasets. The results are shown in Figure 7. On five datasets (PAAD, UVM, GBM,
13 STAD, UCEC), MCluster-VAEs had the smallest log-rank P value, while the remaining
14 five datasets ranked second. The means of $-\log_{10} P$ values of survival analysis obtained
15 by MCluster-VAEs were 3.845 (methylation-alone), 2.601 (mRNA-alone), 2.367
16 (miRNA-alone) and 1.407 (CNA-alone). They were worse than that of multi-omics data
17 (5.850). In addition, MCluster-VAEs had the most clinical parameters which were
18 enriched on nine datasets (except UVM). The means were 4.00 (multi-omics), 2.80
19 (methylation-alone), 2.78 (miRNA-alone), 1.6 (CNA-alone) and 1.3 (mRNA-alone).
20 These results again proved that MCluster-VAEs could effectively use multiple
21 information from multi-omics data. Combining the results of the Pan Cancer dataset,
22 we had reason to believe that MCluster-VAEs can achieve better integration of multi-

1 omics data for clustering tasks.

2 **Figure 7.** Performance of MCluster-VAEs based on four omics or single-omics data on the ten
3 cancer datasets. The x-axis was the number of clinical parameters enriched in the clusters, and the
4 y-axis measured the differential survival between clusters ($-\log_{10}$ of permuted logrank's test P
5 values). Colors indicated the omics data applied. Here, MOmics represents four omics data, mRNA
6 represents mRNA expression, methy denotes DNA methylation (450K), miRNA represents miRNA
7 expression and CNA represents copy number alterations.

8 **3.5 Identification of marker genes of BRCA dataset**

9 It is crucial to demonstrate that the subtypes identified by MCluster-VAEs are
10 biologically interpretable as an application in medicine. To identify essential genes
11 among the five subtypes of breast carcinoma (BRCA) identified by MCluster-VAEs
12 (i.e., C1, C2, C3, C4, and C5), we calculated the signal-to-noise ratio (SNR) of each
13 gene, proposing a one-vs-rest differential procedure for one subtype as one group with
14 the other four subtypes as another group. The gene with a high absolute value of SNR
15 indicates significant differential expression in its subtype. We selected ten genes with
16 the highest SNR and ten genes with the lowest SNR for each BRCA subtype, such that
17 a total of one hundred marker genes were filtered and listed in Supplementary Table S6.
18 To interpret the biological role and potential functions of the marker genes identified
19 by MCluster-VAEs, functional enrichment analysis of the marker genes was conducted
20 with a gene list of GO (Gene Ontology) Molecular Functions, GO Biological Processes
21 and GO Cellular Components in Metascape (<https://metascape.org>). All genes in the
22 genome have been used as the enrichment background. Terms with a P-value <0.01 , a
23 minimum overlap of 3, and an enrichment factor >1.5 are collected and grouped into
24 clusters based on their membership similarities. The results of GO analyses are

1 displayed in Figure 8A and 8B. In addition, the top 18 biological process items are listed
2 in Supplementary Table S7.

3 Biological processes include epithelial cell differentiation, RAGE receptor binding, cell
4 maturation, female sex differentiation, lipid binding, multicellular organismal
5 homeostasis, cellular process involved in reproduction in multicellular organism,
6 negative regulation of cell population proliferation and growth factor activity were
7 significantly regulated by these genes. A protein–protein interaction (PPI) network was
8 created with STRING, BioGrid, OmniPath, InWeb_IM. Only physical interactions in
9 STRING (physical score > 0.132) and BioGrid are used. Moreover, significant modules
10 of PPI were identified in Figures 8C. In addition, we utilized the molecular complex
11 detection (MCODE) algorithm to analyze clusters of the PPI networks. The MCODE
12 component extracted was mainly associated with RAGE receptor binding, sequestering
13 of metal ion and epithelial cell differentiation (Figure 8 D–E). We also found that the
14 molecular subtyping results of MCluster-VAEs are partially consistent with the
15 previous analysis results obtained by gene expression data (Supplementary Table S8),
16 confirming the ability of MCluster-VAEs to acquire reasonable cancer subtypes from
17 the original multi-omics data without cancer-related prior information.

18 **Figure 8.** Enrichment analysis of the marker genes identified by MCluster-VAEs (Metascape). **A.**
19 Pathway and process enrichment analysis has been carried out with Gene Ontology (GO) Biological
20 Processes. Bar graph of enriched terms across the marker genes, colored by P-values; **B.** Network
21 of enriched terms across the marker genes: colored by P-value, where terms containing more genes
22 tend to have a more significant P-value; **C.** Protein–protein interaction network and the MCODE
23 component identified in the marker genes; **D–E.** Pathway and process enrichment analysis has been
24 applied to the MCODE component (**D**), and the three best-scoring terms by P-value have been
25 retained as the functional description of the MCODE components, shown in the table (**E**).

1 **4. Discussion**

2 We presented a brand-new deep learning-based multi-omics clustering algorithm,
3 called MCluster-VAEs. As far as we know, MCluster-VAEs is the first “one-step” deep
4 learning-based clustering algorithm for multi-omics data. The experiments on the Pan
5 Cancer dataset proved that MCluster-VAEs could better identify the intrinsic
6 categorical information than other comparison methods, while the results on the ten
7 cancer datasets showed that it could find biologically significant clusters.

8 The main insight of MCluster-VAEs is that it is better to use only cluster-friendly
9 information than all information in multi-omics data. We believe that MCluster-VAEs’
10 better performance than “two-steps” deep learning-based methods stems considerably
11 from this insight. Therefore, it is important that the probabilistic model contains both
12 global categorical latent variable y and omic-special latent variables \mathbf{z}^m . \mathbf{z}^m will
13 be used to bear cluster-unfriendly information, which makes y focus on the cluster-
14 friendly information. The attention module further enhances this ability, enabling the
15 model to make trade-offs and selections among multiple omics data. The main
16 differences between MCluster-VAEs and the "one-step" single-modal clustering
17 methods are the omics-special \mathbf{z}^m and attention module, which makes MCluster-
18 VAEs more suitable for multi-omics data.

19 In the result section, the number of clusters k of all methods is set to a fair value
20 obtained from known labels or previous large-scale studies for each tumor type.
21 However, determining the optimal k is also a challenging task. For ten cancer types,

1 we explored a simple strategy to determine the most appropriate k based on
2 MCluster-VAEs: (1) choose a possible k as the MCluster-VAEs setting and fit model;
3 (2) extract the fusion hidden representation from the attention module's output hidden
4 layer; (3) calculate silhouette score using this representation, and cluster assignments;
5 (4) perform (1)-(3) on the all possible k and choose the k with the largest silhouette
6 score. This strategy was applied to 10 cancer datasets, and the results were shown in
7 Supplementary Figure S7. We found that MCluster-VAEs determined a reasonable k
8 for all datasets. On six of the ten datasets (STAD, PAAD, LUAD, SKCM, GBM, UVM),
9 the k values obtained by MCluster-VAEs are precisely equal to the values used in the
10 benchmark. For other datasets, MCluster-VAEs obtained relatively conservative (i.e.,
11 fewer) k . Thus, we believe that even when the number of clusters is unknown, the
12 subtype results obtained by this strategy based on MCluster-VAEs and silhouette score
13 are reliable.

14 There are still some limitations of the MCluster-VAEs. The first limitation is that the
15 attention module is a little crude and cannot do gene-level evaluation. The second
16 limitation is the instability in small sample size data, which could be seen in PAAD,
17 GBM, and UVM datasets. Even if using the weight annealing trick and Gumbel softmax
18 reparameterization, there is still a slight chance to train into the collapse. For this
19 problem, our experience is that the entropy and conditional entropy should be
20 monitored. The training must guarantee the gradual rise of the entropy and gradual fall
21 of the conditional entropy by controlling the hyperparameters (such as learning rate, the
22 weight of conditional entropy items, etc.). Simultaneously, the entropy finally needs to

1 be as large as enough. When the above conditions are satisfied, MCluster-VAEs will
2 perform well in most cases.

3 **5. Conclusions**

4 Clustering cancer patients into subgroups has the potential to be used for personalized
5 diagnosis and therapy. The increasing diversity of omics data, as well as their reduced
6 cost, creates an opportunity to use multi-omics data to discover such subgroups.
7 MCluster-VAEs' ability to handle complex relationships on multi-omics data makes it
8 a valuable method. Its importance will become apparent with the further exploration of
9 omics data.

10 **6. List of abbreviations**

- 11 ● MCluster-VAEs: Multi-omics Clustering Variational Autoencoders
- 12 ● AE: autoencoder
- 13 ● VAE: variational autoencoder
- 14 ● PDFs: probability density functions
- 15 ● ELBO: evidence lower bound
- 16 ● KL: Kullback-Leibler
- 17 ● RHS: Right-Hand-Side
- 18 ● SGD: stochastic gradient descent
- 19 ● MLP: multilayer perceptrons
- 20 ● ACC: unsupervised clustering accuracy

- 1 ● ARI: adjusted Rand index
- 2 ● NMI: normalized mutual information
- 3 ● F1: clustering F measure
- 4 ● GO: Gene Ontology
- 5 ● PPI: protein-protein interaction
- 6 ● The abbreviations of each cancer types are shown in Supplementary Table S1

7 **7. Declarations**

8 **7.1 Ethics approval and consent to participate**

9 Not applicable.

10 **7.2 Consent for publication**

11 Not applicable.

12 **7.3 Availability of data and materials**

13 The datasets used during the current study are available in the UCSC Xena data portal,

14 <https://xenabrowser.net/datapages/>.

15 Code has been made publicly available on GitHub

16 <https://github.com/luyiyun/MCluster-VAEs>.

17 **7.4 Competing interests**

18 The authors declare that they have no competing interests.

19 **7.5 Funding**

20 This work was supported by the National Natural Science Foundation of China

21 [81773550, 82173615].

1 **7.6 Authors' contributions**

2 ZR proposed the main idea and implemented the codes of MCluster-VAEs, and was a
3 major contributor in writing the manuscript. JS and LC collected and analyzed the
4 benchmark datasets used in current research. ZL tested the comparison methods in the
5 benchmark datasets. YY corrected the language of the manuscript. YH were the funder
6 of the current research. All authors read and approved the final manuscript.

7 **7.7 Acknowledgements**

8 Not applicable.

9 **8. References**

- 10 1. Pavía-Jiménez A, Tcheuyap VT, Brugarolas J. Establishing a human renal cell carcinoma
11 tumorgraft platform for preclinical drug testing. *Nature Protocols*. 2014;9(8):1848–59.
- 12 2. Berger AC, Korkut A, Kanchi RS, Hegde AM, Lenoir W, Liu W, et al. A Comprehensive
13 Pan-Cancer Molecular Study of Gynecologic and Breast Cancers. *Cancer Cell*. 2018 Apr
14 9;33(4):690–705.e9.
- 15 3. Mo Q, Shen R, Guo C, Vannucci M, Chan KS, Hilsenbeck SG. A fully Bayesian latent
16 variable model for integrative clustering analysis of multi-type omics data. *Biostatistics*.
17 2018 Jan 1;19(1):71–86.
- 18 4. Bridges CC. Hierarchical Cluster Analysis. *Psychol Rep*. 1966;18(3):851–4.
- 19 5. Shi J, Malik J. Normalized cuts and image segmentation. 2000.

- 1 6. MacQueen J, others. Some methods for classification and analysis of multivariate
2 observations. In: Proceedings of the fifth Berkeley symposium on mathematical statistics
3 and probability. Oakland, CA, USA; 1967. p. 281–97.
- 4 7. Cantini L, Zakeri P, Hernandez C, Naldi A, Thieffry D, Remy E, et al. Benchmarking joint
5 multi-omics dimensionality reduction approaches for the study of cancer. *Nature*
6 *Communications*. 2021;12(1):124.
- 7 8. Wu D, Wang D, Zhang MQ, Gu J. Fast dimension reduction and integrative clustering of
8 multi-omics data using low-rank approximation: application to cancer molecular
9 classification. *BMC Genomics*. 2015;16(1):1022.
- 10 9. Shen R, Olshen AB, Ladanyi M. Integrative clustering of multiple genomic data types
11 using a joint latent variable model with application to breast and lung cancer subtype
12 analysis. *Bioinformatics*. 2009;25(22):2906–12.
- 13 10. Wang B, Jiang J, Wang W, Zhou ZH, Tu Z. Unsupervised metric fusion by cross diffusion.
14 In: 2012 IEEE Conference on Computer Vision and Pattern Recognition. 2012. p. 2997–
15 3004.
- 16 11. Wang B, Mezlini AM, Demir F, Fiume M, Tu Z, Brudno M, et al. Similarity network fusion
17 for aggregating data types on a genomic scale. *Nature Methods*. 2014;11(3):333–7.
- 18 12. Rappoport N, Shamir R. NEMO: cancer subtyping by integration of partial multi-omic
19 data. *Bioinformatics*. 2019 Sep 15;35(18):3348–56.

- 1 13. Zhu J, Oh JH, Deasy J, Tannenbaum A. vWCluster: A Network Based Clustering of Multi-
2 omics Breast Cancer Data Based on Vector-Valued Optimal Transport. 2021.
- 3 14. Lemsara A, Ouadfel S, Fröhlich H. PathME: pathway based multi-modal sparse
4 autoencoders for clustering of patient-level multi-omics data. BMC Bioinformatics. 2020
5 Apr 16;21(1):146.
- 6 15. Witten DM, Tibshirani RJ. Extensions of sparse canonical correlation analysis with
7 applications to genomic data. Stat Appl Genet Mol Biol. 2009;8(1):Article28.
- 8 16. Lock EF, Hoadley KA, Marron JS, Nobel AB. JOINT AND INDIVIDUAL VARIATION
9 EXPLAINED (JIVE) FOR INTEGRATED ANALYSIS OF MULTIPLE DATA TYPES. Ann Appl
10 Stat. 2013 Mar 1;7(1):523–42.
- 11 17. Lock EF, Dunson DB. Bayesian consensus clustering. Bioinformatics. 2013;29(20):2610–
12 6.
- 13 18. Kirk P, Griffin JE, Savage RS, Ghahramani Z, Wild DL. Bayesian correlated clustering to
14 integrate multiple datasets. Bioinformatics. 2012;28(24):3290–7.
- 15 19. He K, Zhang X, Ren S, Sun J. Deep Residual Learning for Image Recognition. In 2016
16 [cited 2022 May 7]. p. 770–8. Available from:
17 [https://openaccess.thecvf.com/content_cvpr_2016/html/He_Deep_Residual_Learning_](https://openaccess.thecvf.com/content_cvpr_2016/html/He_Deep_Residual_Learning_CVPR_2016_paper.html)
18 [CVPR_2016_paper.html](https://openaccess.thecvf.com/content_cvpr_2016/html/He_Deep_Residual_Learning_CVPR_2016_paper.html)
- 19 20. Devlin J, Chang MW, Lee K, Toutanova K. BERT: Pre-training of Deep Bidirectional

1 Transformers for Language Understanding. arXiv e-prints. 2018;arXiv:1810.04805.

2 21. Kipf TN, Welling M. Semi-supervised classification with graph convolutional networks.
3 arXiv preprint arXiv:160902907. 2016;

4 22. Chaudhary K, Poirion OB, Lu L, Garmire LX. Deep Learning–Based Multi-Omics
5 Integration Robustly Predicts Survival in Liver Cancer. Clin Cancer Res. 2018;24(6):1248.

6 23. Chai H, Zhou X, Zhang Z, Rao J, Zhao H, Yang Y. Integrating multi-omics data through
7 deep learning for accurate cancer prognosis prediction. Comput Biol Med. 2021
8 Jul;134:104481.

9 24. Ronen J, Hayat S, Akalin A. Evaluation of colorectal cancer subtypes and cell lines using
10 deep learning. Life Sci Alliance. 2019 Dec 1;2(6):e201900517.

11 25. Hira MT, Razzaque MA, Angione C, Scrivens J, Sawan S, Sarker M. Integrated multi-
12 omics analysis of ovarian cancer using variational autoencoders. Scientific Reports.
13 2021;11(1):6265.

14 26. Yang H, Chen R, Li D, Wang Z. Subtype-GAN: a deep learning approach for integrative
15 cancer subtyping of multi-omics data. Bioinformatics. 2021 Feb 18;

16 27. Yang B, Fu X, Sidiropoulos ND, Hong M. Towards K-Means-Friendly Spaces:
17 Simultaneous Deep Learning and Clustering. In: Proceedings of the 34th International
18 Conference on Machine Learning - Volume 70. JMLR.org; 2017. p. 3861–70. (ICML'17).

- 1 28. Min E, Guo X, Liu Q, Zhang G, Cui J, Long J. A Survey of Clustering With Deep Learning:
2 From the Perspective of Network Architecture. *IEEE Access*. 2018;6:39501–14.
- 3 29. Diallo B, Hu J, Li T, Khan GA, Liang X, Zhao Y. Deep embedding clustering based on
4 contractive autoencoder. *Neurocomputing*. 2021;433:96–107.
- 5 30. Guo X, Gao L, Liu X, Yin J. Improved Deep Embedded Clustering with Local Structure
6 Preservation. 2017;1753–9.
- 7 31. Kingma DP, Welling M. Auto-encoding variational bayes [J]. 2013;
- 8 32. Eraslan G, Avsec Ž, Gagneur J, Theis FJ. Deep learning: new computational modelling
9 techniques for genomics. *Nature Reviews Genetics*. 2019;1.
- 10 33. Jang E, Gu S, Poole B. Categorical Reparameterization with Gumbel-Softmax.
11 2016;arXiv:1611.01144.
- 12 34. Maddison C, Mnih A, Teh Y. The Concrete Distribution: A Continuous Relaxation of
13 Discrete Random Variables. 2016 02;
- 14 35. Mantel N. Evaluation of survival data and two new rank order statistics arising in its
15 consideration. *Cancer Chemother Rep*. 1966 Mar;50(3):163–70.
- 16 36. Rappoport N, Shamir R. Multi-omic and multi-view clustering algorithms: review and
17 cancer benchmark. *Nucleic Acids Research*. 2018;46(20):10546–62.
- 18 37. Zhang C, Chen Y, Zeng T, Zhang C, Chen L. Deep latent space fusion for adaptive

- 1 representation of heterogeneous multi-omics data. *Briefings in Bioinformatics*. 2022
2 Mar 1;23(2):bbab600.
- 3 38. Yang Y, Tian S, Qiu Y, Zhao P, Zou Q. MDICC: novel method for multi-omics data
4 integration and cancer subtype identification. *Briefings in Bioinformatics*. 2022 Apr
5 18;bbac132.
- 6 39. Hoadley KA, Yau C, Wolf DM, Cherniack AD, Tamborero D, Ng S, et al. Multiplatform
7 analysis of 12 cancer types reveals molecular classification within and across tissues of
8 origin. *Cell*. 2014 Aug 14;158(4):929–44.
- 9 40. Ma T, Zhang A. Integrate multi-omic data using affinity network fusion (ANF) for cancer
10 patient clustering. In: 2017 IEEE International Conference on Bioinformatics and
11 Biomedicine (BIBM). 2017. p. 398–403.
- 12 41. Ramazzotti D, Lal A, Wang B, Batzoglou S, Sidow A. Multi-omic tumor data reveal
13 diversity of molecular mechanisms that correlate with survival. *Nature Communications*.
14 2018 Oct 26;9(1):4453.
- 15 42. Maaten L van der, Hinton G. Visualizing Data using t-SNE. *Journal of Machine Learning*
16 *Research*. 2008;9(86):2579–605.
- 17 43. Genomic Classification of Cutaneous Melanoma. *Cell*. 2015 Jun 18;161(7):1681–96.
- 18 44. Collisson EA, Campbell JD, Brooks AN, Berger AH, Lee W, Chmielecki J, et al.
19 Comprehensive molecular profiling of lung adenocarcinoma. *Nature*.

1 2014;511(7511):543–50.

2 45. Creighton CJ, Morgan M, Gunaratne PH, Wheeler DA, Gibbs RA, Gordon Robertson A,
3 et al. Comprehensive molecular characterization of clear cell renal cell carcinoma.
4 Nature. 2013;499(7456):43–9.

5 46. Levine DA, Getz G, Gabriel SB, Cibulskis K, Lander E, Sivachenko A, et al. Integrated
6 genomic characterization of endometrial carcinoma. Nature. 2013;497(7447):67–73.

7 47. Robertson AG, Kim J, Al-Ahmadie H, Bellmunt J, Guo G, Cherniack AD, et al.
8 Comprehensive Molecular Characterization of Muscle-Invasive Bladder Cancer. Cell.
9 2017 Oct 19;171(3):540-556.e25.

10 48. Verhaak RGW, Hoadley KA, Purdom E, Wang V, Qi Y, Wilkerson MD, et al. Integrated
11 genomic analysis identifies clinically relevant subtypes of glioblastoma characterized by
12 abnormalities in PDGFRA, IDH1, EGFR, and NF1. Cancer Cell. 2010 Jan 19;17(1):98–110.

13

14

15

16

Figures

Figure 1

Overview of MCluster-VAEs. (A) Two deep learning-based multi-omics clustering pipelines. (B) The data generative model. (C) The posterior variational inference path. (D) The architecture of MCluster-VAEs model for the experiments using mRNA expression, miRNA expression, copy number alterations (CNA) and DNA methylation.

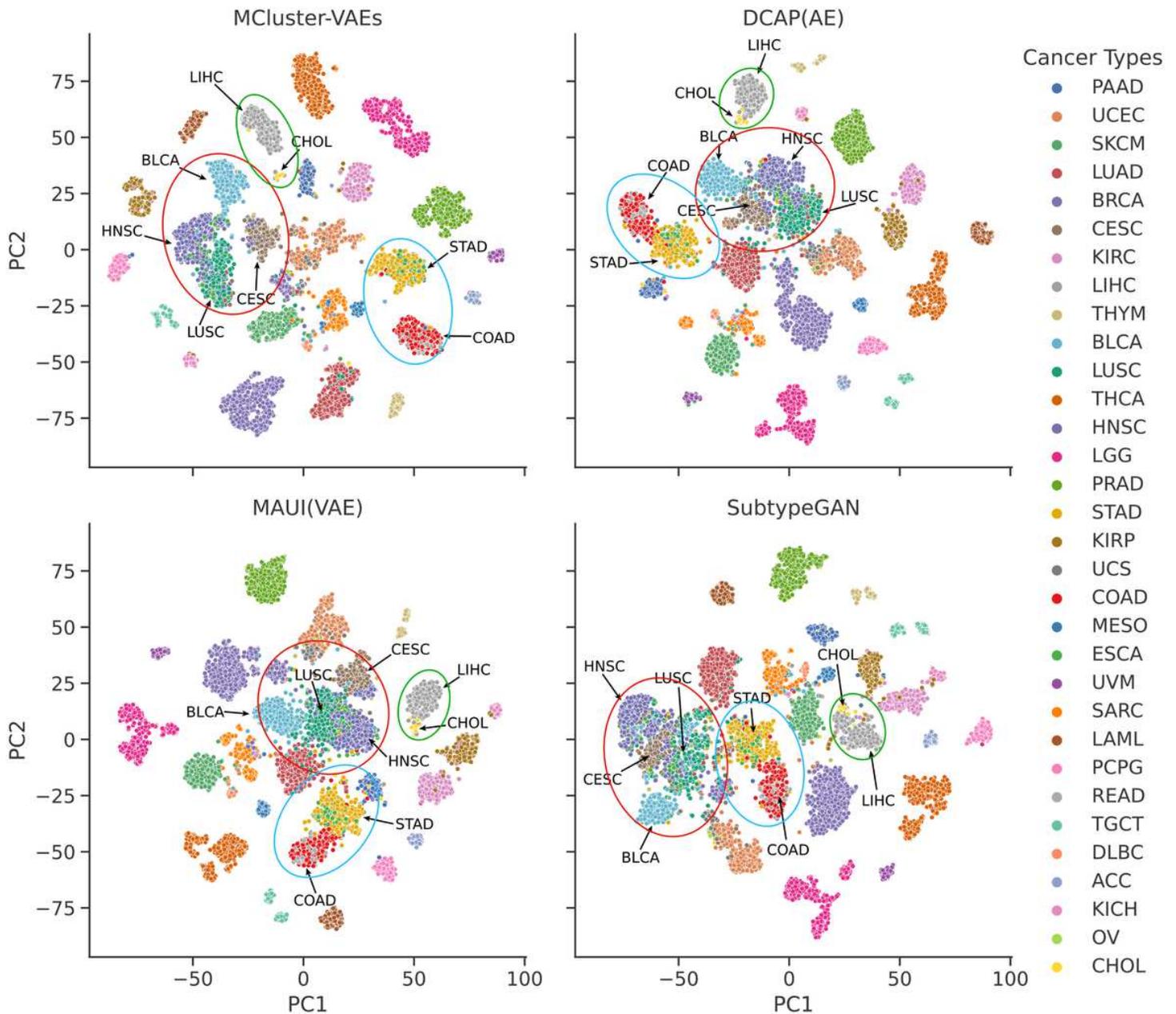


Figure 2

The t-SNE visualization on the latent variables generated by four deep learning-based methods (MCluster-VAEs, AE, VAE, SubtypeGAN) used the Pan Cancer dataset. The color of each point indicated the cancers. MCluster-VAEs had a greater degree of separation.

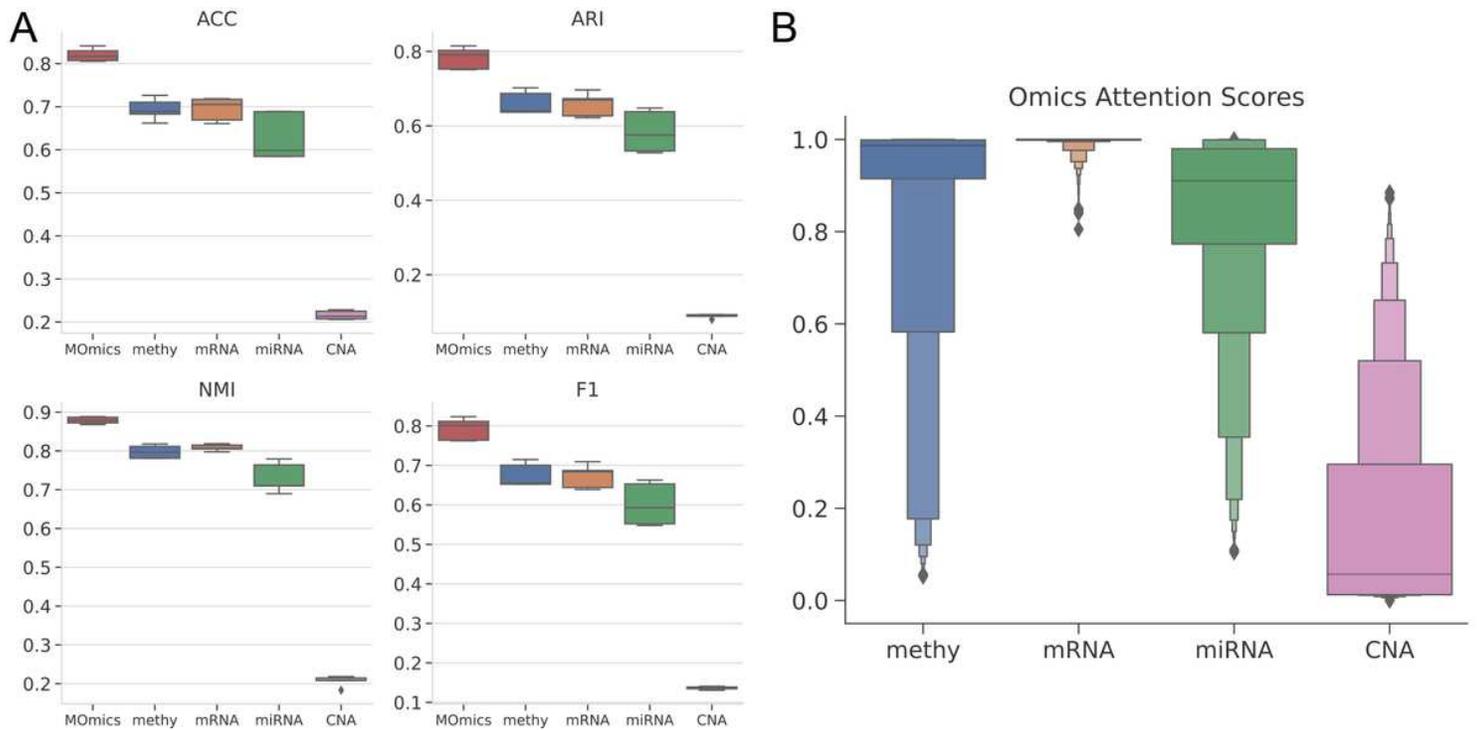


Figure 3

Distribution of metrics and attention scores of MCluster-VAEs using single omics data of the Pan Cancer datasets. **A.** ACC, ARI, NMI and F1 of MCluster-VAEs based on all four omics (MOmics) or single-omics data. **B.** The distribution of attention scores for each omics data. Here, mRNA represents mRNA expression, methy denotes DNA methylation (450K), miRNA represents miRNA expression and CNA represents copy number alterations.

Figure 4

The monitoring records of the training process of MCluster-VAEs on the Pan Cancer dataset. “WA” means using conditional entropy weight decay annealing trick. “Gumbel” means using Gumbel Softmax reparameterization trick. “none” means neither used.

Figure 5

Performance of the algorithms on the ten cancer datasets. The x-axis was the number of clinical parameters enriched in the clusters, and the y-axis measured the differential survival between clusters (-

log10 of permuted logrank's test P values). Colors indicated clustering methods.

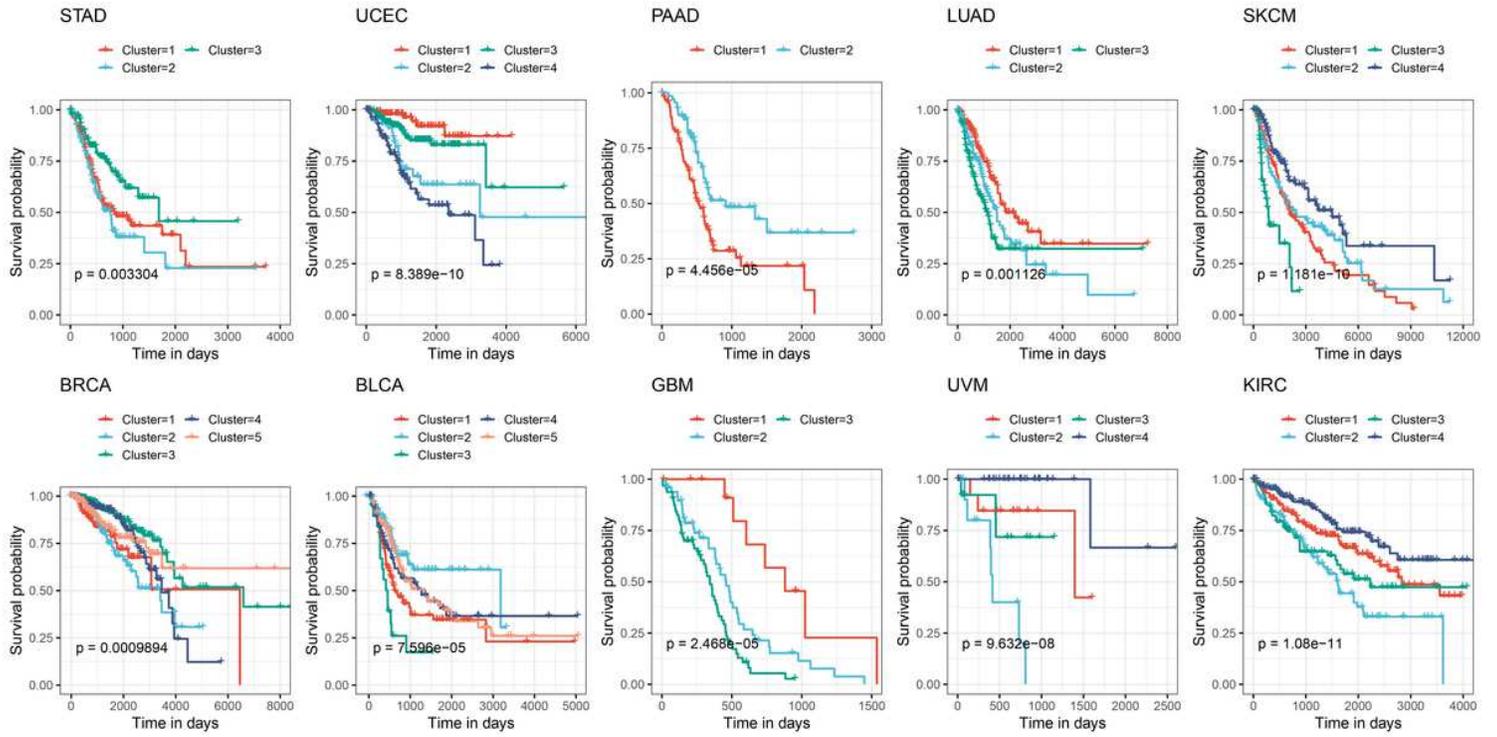


Figure 6

Kaplan-Meier survival plots of MCluster-VAEs on the ten cancer datasets.

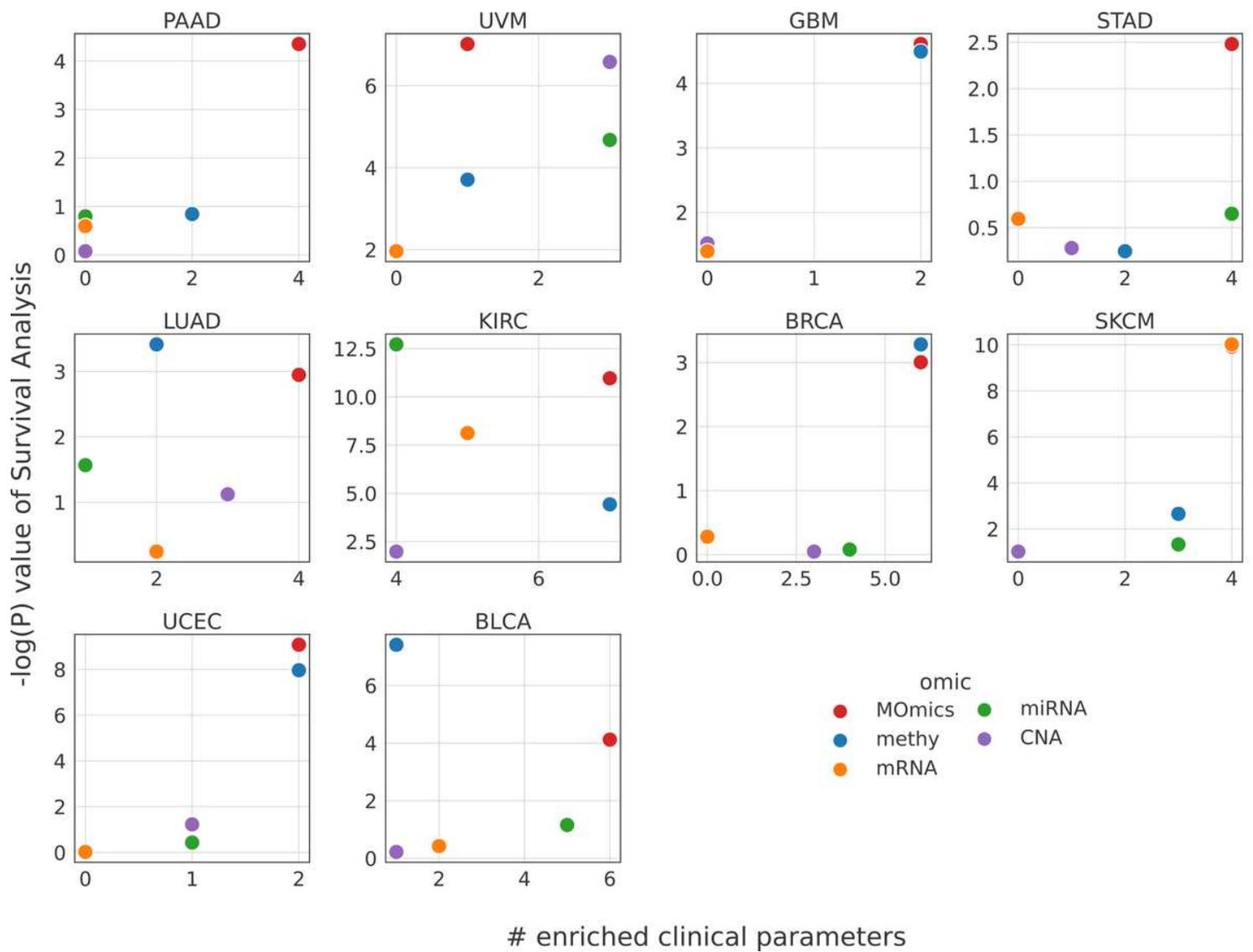


Figure 7

Performance of MCluster-VAEs based on four omics or single-omics data on the ten cancer datasets. The x-axis was the number of clinical parameters enriched in the clusters, and the y-axis measured the differential survival between clusters ($-\log_{10}$ of permuted logrank's test P values). Colors indicated the omics data applied. Here, MOmics represents four omics data, mRNA represents mRNA expression, methy denotes DNA methylation (450K), miRNA represents miRNA expression and CNA represents copy number alterations.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [supplementary.docx](#)