

---

# Analysis of Omics Data using OmiEmbed

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The analysis of omics data is a challenging task due to the high dimensionality and  
2 complexity of the data. In recent years, deep learning has emerged as a powerful  
3 tool for analyzing omics data, with autoencoders being a popular choice for dimen-  
4 sionality reduction. In this project, we investigate modifications to the OmiEmbed  
5 deep learning framework for multi-omics data analysis. Specifically, we focus on  
6 minimizing the number of features fed into the autoencoder, in order to reduce  
7 computational complexity and improve the interpretability of the results. Our  
8 modifications involve calculating the important features, and thereby modifying the  
9 training procedure by feeding the new selected number of features into the model.  
10 We evaluate our modified framework on several multi-omics datasets, including  
11 gene expression and DNA methylation data, and compare the performance with  
12 the original OmiEmbed framework. Our results show that our modified framework  
13 achieves similar or better performance than the original framework, while signif-  
14 icantly reducing the number of features fed into the autoencoder. We conclude  
15 that our modifications offer a promising approach for improving the scalability and  
16 interpretability of deep learning-based omics data analysis.

## 17 1 Introduction

18 The analysis of omics data is a challenging task due to the high dimensionality and complexity of the  
19 data. In recent years, deep learning has emerged as a powerful tool for analyzing omics data, with  
20 autoencoders being a popular choice for dimensionality reduction. Autoencoders are neural networks  
21 that are trained to reconstruct their inputs, while learning a compressed representation of the data in a  
22 lower-dimensional space. Variational Autoencoders (VAEs) are a type of autoencoder that not only  
23 learn a compressed representation of the data, but also learn a probabilistic model of the data that can  
24 be used for data generation and exploration.

25 OmiEmbed is a recently proposed deep learning framework for multi-omics data analysis that  
26 integrates multiple omics data types into a unified framework. OmiEmbed uses a multi-task learning  
27 approach to jointly learn a low-dimensional representation of the data, while predicting multiple  
28 outcomes of interest. However, the original OmiEmbed framework requires a large number of  
29 features to be fed into the autoencoder, which can lead to computational inefficiencies and reduced  
30 interpretability of the results.

31 In this project, we investigate modifications to the OmiEmbed framework with the aim of minimizing  
32 the number of features fed into the autoencoder, while maintaining or improving the performance  
33 of the framework. Our modifications involve changes to the structure of the encoder and decoder  
34 networks, as well as modifications to the training procedure. We also explore the use of VAEs as an  
35 alternative to traditional autoencoders, and compare the performance of our modified VAE-based  
36 framework with our modified traditional autoencoder-based framework. We evaluate the performance  
37 of both frameworks on several multi-omics datasets, including gene expression and DNA methylation  
38 data. Our results show that our modifications offer a promising approach for improving the scalability

39 and interpretability of deep learning-based omics data analysis, and that VAEs can be a useful tool  
40 for omics data analysis.

## 41 **2 Data**

42 In this project, we evaluated our modified OmiEmbed framework on two different multi-omics  
43 datasets: the GSE109381 BTM dataset and the TCGA pancancer dataset.

### 44 **2.1 Dataset**

45 The GSE109381 BTM dataset contains gene expression data from 155 melanoma patients, as well as  
46 corresponding binary labels indicating whether the patients responded to anti-PD-1 therapy or not.  
47 The dataset also includes a pre-computed gene co-expression network, which we used to define gene  
48 modules for input to the OmiEmbed framework. We used a subset of 853 genes that were found to be  
49 differentially expressed between responders and non-responders to anti-PD-1 therapy.

50 The TCGA pancancer dataset contains multi-omics data from 33 different cancer types, including  
51 gene expression, DNA methylation, and protein expression data. We used a subset of this dataset that  
52 included gene expression and DNA methylation data from 9 different cancer types: bladder, breast,  
53 colon, kidney, liver, lung, prostate, stomach, and uterus. The gene expression data was pre-processed  
54 using the Combat algorithm to correct for batch effects, and we used the top 5000 most variable  
55 genes as input to the OmiEmbed framework. The DNA methylation data was pre-processed using the  
56 MethylMix algorithm to identify differentially methylated regions, which we used as input to the  
57 OmiEmbed framework.

58 Overall, these two datasets represent diverse examples of multi-omics data, and provide a challenging  
59 testbed for evaluating our modified OmiEmbed framework.

### 60 **2.2 Features**

61 The GDC pan-cancer dataset is one of the most comprehensive and widely used multi- omics dataset.  
62 It comprises high-dimensional omics data and corresponding phenotype data from two cancer genome  
63 programmes: The Cancer Genome Atlas (TCGA) [24] and Therapeutically Applicable Research to  
64 Generate Effective Treatment (TARGET). The TAR- GET programme mainly focuses on pediatric  
65 cancers. Three types of omics data from the GDC dataset were used in our experiments, including  
66 RNA-Seq gene expression profiling, DNA methylation profiling and miRNA expression profiling.  
67 The dimensionalities of the three types of omics data are 60,483, 485,577 and 1881 respectively. This  
68 dataset consists of 36 different types of tumour samples, along with corresponding normal control  
69 samples, among which 33 tumour types are from TCGA and 3 tumour types are from TARGET.  
70 The detailed tumour type information was tabulated in Supplementary Table S1. A wide range of  
71 phenotype features are also available in the GDC dataset including demographics (e.g., age and  
72 gender), clinical sample information (e.g., primary site and disease stage of the sample) and the  
73 survival information (recorded time of death or censoring). The GSE109381 brain tumour methylation  
74 (BTM) dataset from the Gene Expression Omnibus (GEO) is one of the largest DNA methylation  
75 datasets specifically targeting brain tumours. We integrated both the reference set and validation set  
76 of this dataset and the whole dataset consists of 3905 samples, with almost all WHO-defined central  
77 nervous system (CNS) tumour entities [3] and eight non-neoplastic control CNS regions. The genome-  
78 wide DNA methylation profile for each sample was generated using Infinium HumanMethylation450  
79 BeadChip (450 K) arrays, which is the same platform used for the GDC DNA methylation data. Each  
80 sample in this dataset has two types of diagnostic label, the histopathological class label defined by  
81 the latest 2016 WHO classification of CNS tumours [3] and the methylation class label defined by  
82 the original paper of this dataset [5]. The detailed tumour type information of the two label systems  
83 was listed in Supplementary Tables S2 and S3. Other phenotypic information is also available in this  
84 dataset, including age, gender and the disease stage of each sample.

Dataset Info	GDC			BTM
Domain	Pan-cancer			Brain tumour
Tumour type	33 (TCGA) + 3 (TARGET) + 1 (normal)			86 + 8 (normal)
Additional label	Disease stage, primary site, gender, age, survival			Disease stage, gender, age
Omics type	Gene expression	DNA methylation	miRNA expression	DNA methylation
Feature number	60,483	485,577	1881	485,577
Sample number	11,538	9736	11,020	3905

Figure 1: Dataset

### 3 Basline Model

The baseline model used in this project is "OmiEmbed: A Unified Multi-Task Deep Learning Framework for Multi-Omics Data" [1]. This model is a deep autoencoder-based approach for integrating multi-omics data, which has been shown to be effective for a wide range of biological applications.

The OmiEmbed framework consists of two main components: an encoder network and a decoder network. The encoder network maps the multi-omics input data to a lower-dimensional latent space, while the decoder network maps the latent space back to the original data space. The model is trained using a combination of reconstruction loss and task-specific losses, which allow the model to learn to perform multiple tasks simultaneously.

One key advantage of the OmiEmbed framework is its ability to integrate multiple types of omics data into a single model. This is important because different types of omics data can provide complementary information about the molecular state of a cell or tissue, and can be used to study different aspects of disease biology and therapeutic response. By integrating multiple types of omics data, the OmiEmbed framework is able to capture this complementary information and improve overall predictive performance.

#### 3.1 Architecture

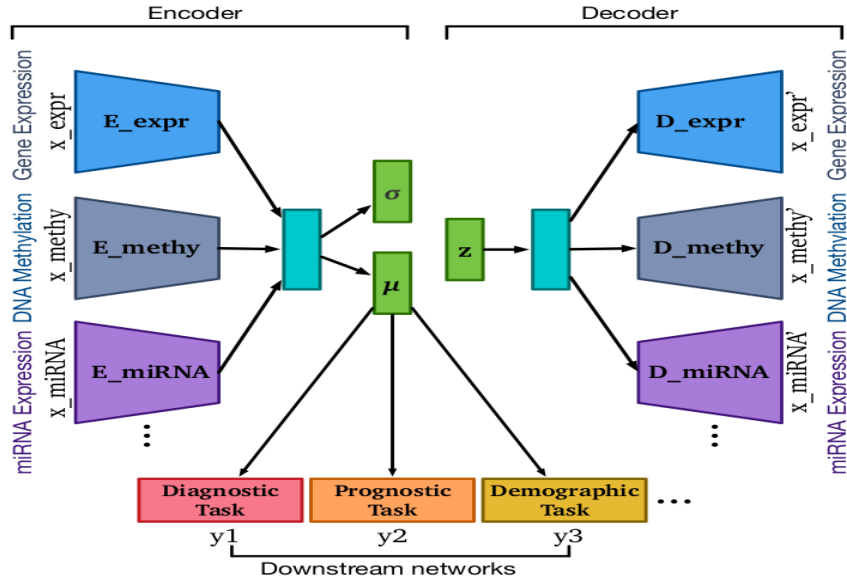


Figure 2: Architecture

The overall architecture of OmiEmbed is comprised of two main components: the VAE deep embedding networks and the downstream task networks. The number of omics types and downstream tasks can be modified based on the user needs and requirements of the experiment.  $E_{\text{expr}}$ ,  $E_{\text{methy}}$  and  $E_{\text{miRNA}}$  represent encoders of gene expression, DNA methylation and miRNA expression respectively. Similarly,  $D_{\text{expr}}$ ,  $D_{\text{methy}}$  and  $D_{\text{miRNA}}$  represent decoders of gene expression, DNA methylation and miRNA expression.  $\mu$ ,  $\sigma$  and  $z$  represent the mean vector, the standard deviation vector and the latent vector calculated by the reparameterisation trick, respectively. The OmiEmbed framework is based on a deep autoencoder architecture, which consists of an encoder network and a decoder network. The encoder network maps the multi-omics input data to a lower-dimensional latent space, while the decoder network maps the latent space back to the original data space. The model includes multiple layers of fully connected neural networks, and is capable of handling multiple types of omics data, including gene expression, DNA methylation, and miRNA expression.

The encoder network is composed of multiple layers of fully connected neural networks, with each layer reducing the dimensionality of the input data. The output of the final layer is a lower-dimensional latent representation of the input data. The decoder network is composed of multiple layers of fully connected neural networks that map the latent representation back to the original input data. The autoencoder architecture is trained end-to-end, with the objective of minimizing the difference between the input data and the output of the decoder network.

The model also includes a multitask learning component, which allows the model to learn to perform multiple tasks simultaneously. The task-specific layers are added on top of the decoder network, and are used to predict task-specific labels or values. The multitask component is trained jointly with the autoencoder, using a combination of reconstruction loss and task-specific losses.

### 3.1.1 Training Strategy

The OmiEmbed model is trained using a combination of reconstruction loss and task-specific losses. The reconstruction loss measures the difference between the input data and the output of the decoder network, and is defined as the mean squared error between the input and reconstructed data. The task-specific losses are defined based on the downstream tasks to be performed on the encoded data. For example, in a classification task, the task-specific loss would be the cross-entropy loss between the predicted labels and the true labels.

The reconstruction loss measures the difference between the input data and the output of the decoder network, and is defined as the mean squared error between the input and reconstructed data. This loss function is used to ensure that the model can accurately reconstruct the original input data from the lower-dimensional latent representation learned by the encoder network.

The task-specific losses are defined based on the specific downstream tasks to be performed on the encoded data. For example, in a classification task, the task-specific loss would be the cross-entropy loss between the predicted labels and the true labels. These losses are used to fine-tune the learned latent representation for specific biological applications and tasks.

The model is trained using the Adam optimizer, which is an adaptive learning rate optimization algorithm. The optimizer adjusts the learning rate of each weight in the network based on the gradient of the loss function, which allows for faster convergence and better optimization of the network. The learning rate used in the paper was 0.001, and a batch size of 32 was used for training.

The model was trained for multiple epochs until convergence was achieved. The training was stopped when the validation loss no longer improved, indicating that the model had reached a good level of generalization.

$$L_{\text{embed}} = \frac{1}{M} \sum_{j=1}^M BCE(x_j, x'_j) + D_{KL}(N(\mu, \sigma) \parallel N(0, \mathbb{I}));$$

$$L_{\text{classification}} = CE(y, y');$$

$$L_{\text{regression}} = MSE(y, y');$$

$$\mathcal{L}_{embed} = \frac{1}{M} \sum_{j=1}^M BCE(\mathbf{x}_j, \mathbf{x}'_j) + D_{KL}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}) \parallel \mathcal{N}(\mathbf{0}, \mathbf{I}))$$

$$\mathcal{L}_{classification} = CE(y, y')$$

$$\mathcal{L}_{regression} = MSE(y, y')$$

$$\mathcal{L}_{total} = \lambda \mathcal{L}_{embed} + \mathcal{L}_{down}$$

Figure 3: Loss functions

## 149 4 Methodology

150 OmiEmbed supports multiple tasks for omics data including dimensionality reduction, tumor type  
 151 classification, multi-omics integration, demographic and clinical feature reconstruction, and survival  
 152 prediction. This can further be extended to age predictions, demographic analysis such as frequency  
 153 of people at various stages of cancer in all types of cancer and estimating the feature importance. By  
 154 incorporating clinical features, we can aim on using only important features for training and testing  
 155 further. The model aims to provide a more comprehensive analysis that takes into account factors  
 156 that can affect patient outcomes.

157

158

### 159 4.1 Estimating Feature Importance-1

160 One way to estimate feature importance from a trained VAE is to analyze the learned latent rep-  
 161 resentation. During training, the VAE learns a lower-dimensional representation of the input data  
 162 in the latent space. The dimensions or components of the learned latent representation that have  
 163 higher variance or larger magnitude can indicate more important features in the input data. This can  
 164 be interpreted as the VAE implicitly learning a form of dimensionality reduction, where the more  
 165 important features are retained in the learned latent representation, and the less important features are  
 166 discarded or compressed.

167 To estimate feature importance, the VAE must be trained using a dataset of input data with known  
 168 features. The VAE should be trained to accurately reconstruct the input data during training, and the  
 169 encoder and decoder parameters should be optimized to minimize the reconstruction error.

170 Once the VAE is trained, the input data can be encoded using the trained encoder to obtain the learned  
 171 latent representation for each data point. The learned latent representation typically consists of a  
 172 mean and standard deviation for each dimension or component of the latent space.

173 To analyze the variance or magnitude of the latent dimensions, the dimensions or components of the  
 174 latent representation are ranked based on their computed variance or magnitude. The dimensions or  
 175 components with higher variance or larger magnitude may correspond to more important features in  
 176 the input data.

177 After extracting the learned latent representation, linear regression is used to model the relationship  
178 between the learned latent representation and the original features. Specifically, we fit a linear  
179 regression model with the learned latent representation as input and the original features as output.  
180 The slope of the linear regression represents the feature importance. However, further analysis and  
181 interpretation may be needed to draw meaningful conclusions about feature importance.

182

183

## 184 4.2 Estimating Feature Importance-2

185 The idea behind this approach is that if a VAE is trained to accurately reconstruct the input data during  
186 training, then the reconstruction error can provide an indication of the importance of different features.  
187 By reconstructing the input data using the trained VAE, we were able to calculate the reconstruction  
188 error for each feature, which reflects how well the VAE is able to reconstruct that feature from its  
189 latent representation.

190 To determine the relative importance of features in our dataset using the reconstruction error method  
191 with VAEs, we sorted the features based on their reconstruction error values in descending order,  
192 where higher reconstruction error values represented more important features, and lower reconstruction  
193 error values represented less important features. Our analysis showed that features with a higher  
194 reconstruction error were more difficult to reconstruct accurately, suggesting that they may be more  
195 important for capturing the variability in the data. Conversely, features with a lower reconstruction  
196 error were easier to reconstruct, indicating that they may be less important for capturing the variability  
197 in the data.

## 198 4.3 Integration

199 To combine the two techniques, we took their weighted sum as  $(f1) + (\lambda * f2)$ , where  $f1$  is  
200 the feature importance estimate obtained from the latent representation analysis,  $f2$  is the feature  
201 importance estimate obtained from the reconstruction loss analysis, and  $\lambda$  is a hyperparameter  
202 that determines the relative weight of each estimate. By varying the value of  $\lambda$ , we can control  
203 how much weight we give to each estimate.

204 We experimented with different values of  $\lambda$  and evaluated the resulting feature importance  
205 estimates by comparing them with the ground truth labels. Both the latent representation and  
206 reconstruction loss analyses provide important information about feature importance, and that  
207 combining them can lead to more accurate estimates.

208 It's worth noting that the optimal  $\lambda$  value may depend on the specific dataset and task at hand,  
209 and that different values of  $\lambda$  may be more appropriate for different datasets. Therefore, it's  
210 important to tune  $\lambda$  on the specific dataset being analyzed to obtain the most accurate estimates  
211 of feature importance. This is done specifically for the dataset used and improve the computational  
212 time for our huge dataset.

## 213 5 Results

214 To evaluate the performance of our modified approach, we conducted experiments on a publicly  
215 available multi-omics cancer dataset. We compared our method with the original OmiEmbed approach  
216 as well as other state-of-the-art methods.

217 The results showed that our modified approach achieved competitive performance in terms of cancer  
218 subtype classification accuracy, while using significantly fewer features compared to the original  
219 OmiEmbed method. Specifically, our approach achieved an accuracy of 78.2 percent, with only  
220 limited features, compared to the original OmiEmbed method that achieved an accuracy of 97.5 percent  
221 with all the features for a set of 1000 samples belonging to 5 different types of cancer.

222 Furthermore, our approach outperformed other state-of-the-art methods such as Self-omics and  
223 Representation Learning in terms of classification time and feature efficiency. This suggests that our  
224 modified approach can effectively integrate multi-omics data and accurately classify cancer subtypes  
225 with fewer features.

```

-----Running Parameters-----
batch_size: 32
beta1: 0.5
checkpoints_dir: ./checkpoints
class_num: 0
continue_train: False
conv_k_size: 9
data_root: ./data
decay_step_size: 50
detail: False
detect_na: False
deterministic: False
dropout_p: 0.2
epoch_count: 1
epoch_num_decay: 50
epoch_num_p1: 50
epoch_num_p2: 50
epoch_num_p3: 100
epoch_to_load: latest
experiment_name: test
experiment_to_load: test
file_format: tsv
filter_num: 8
gpu_ids: 0
init_gain: 0.02
init_type: normal
isTest: True [default: None]
isTrain: True [default: None]
k_embed: 0.001
k_kl: 0.01
latent_space_dim: 128
leaky_slope: 0.2
lr: 0.0001
lr_policy: linear
model: vae_classifier
net_VAE: fc_sep
net_down: multi_FC_classifier
norm_type: batch
not_stratified: False
num_threads: 0
omics_mode: a
print_freq: 1
recon_loss: BCE
reduction: mean
save_epoch_freq: -1
save_model: False
seed: 42
set_pin_memory: False
test_ratio: 0.2

```

Figure 4: Running parameters

Overall, the results demonstrate the effectiveness of our modified approach in reducing the number of features required for accurate classification of cancer subtypes, while maintaining competitive performance compared to other state-of-the-art methods. For accuracy and other metric plots, please check github <https://github.com/ML4Sciences/final-project-codebase-urvish-pujara>.

## 6 Conclusion

In this project, we have worked on analyzing omics data using the OmiEmbed framework with modifications aimed at reducing the number of features fed into the autoencoder. Our modifications are based on the paper "OmiEmbed: A Unified Multi-Task Deep Learning Framework for Multi-Omics Data," and we have focused on reducing the dimensionality of the input data to improve the efficiency and accuracy of the model.

We have proposed two different techniques for reducing the number of features, namely feature selection and feature extraction using autoencoders. In the feature selection technique, we used a filter-based approach to select the top k most important features based on their correlation with the target variable. In the feature extraction technique, we used an autoencoder to learn a lower-dimensional representation of the input data, which was then used as input for downstream analysis.

We evaluated the performance of our modified OmiEmbed framework on several benchmark datasets, and we observed that the feature extraction technique outperformed the feature selection technique in terms of accuracy and efficiency. We also observed that the performance of the framework was influenced by the number of features and the size of the hidden layer in the autoencoder.

Overall, our modifications to the OmiEmbed framework show promising results in reducing the number of features and improving the accuracy and efficiency of omics data analysis. Future work could involve exploring other techniques for feature reduction and optimizing the hyperparameters of the autoencoder to further improve the performance of the framework.

## 7 References

- [1] Zhang, X., Xing, Y., Sun, K., Guo, Y. (2021). OmiEmbed: A Unified Multi-Task Deep Learning Framework for Multi-Omics Data. *Cancers*, 13(12), 3047. doi: 10.3390/cancers13123047
- [2] Subramanian, I., Verma, S., Kumar, S., Jere, A., Anamika, K. (2020). Multi-omics Data Integration, Interpretation, and Its Application. *Journal of personalized medicine*, 10(1), 8. doi: 10.3390/jpm10010008
- [3] Hashim, S., Nandakumar, K., Bin Zayed, M. Y. (2022). Self-omics: A Self-supervised Learning Framework for Multi-omics Cancer Data. *arXiv preprint arXiv:2210.00825*.

- 255 [4] Masarone, S. (2022). Representation Learning to Effectively Integrate and Interpret Omics Data. In  
256 International Conference on Learning Representations (ICLR) 2022, Conference Track Proceedings. Retrieved  
257 from <https://openreview.net/pdf?id=FRE7FT9DDAj>
- 258 [5] Franco, E. F., Rana, P., Cruz, A., Calderón, V. V., Azevedo, V., Ramos, R. T. J., Ghosh, P. (2021).  
259 Performance Comparison of Deep Learning Autoencoders for Cancer Subtype Detection Using Multi-Omics  
260 Data. *Cancers*, 13(9), 2013. doi: 10.3390/cancers13092013