



OPEN

A novel interpretable deep transfer learning combining diverse learnable parameters for improved T2D prediction based on single-cell gene regulatory networks

Sumaya Alghamdi^{1,2} & Turki Turki¹✉

Accurate deep learning (DL) models to predict type 2 diabetes (T2D) are concerned not only with targeting the discrimination task but also with learning useful feature representation. However, existing DL tools are far from perfect and do not provide appropriate interpretation as a guideline to explain and promote superior performance in the target task. Therefore, we provide an interpretable approach for our presented deep transfer learning (DTL) models to overcome such drawbacks, working as follows. We utilize several pre-trained models including SEResNet152, and SEResNeXT101. Then, we transfer knowledge from pre-trained models via keeping the weights in the convolutional base (i.e., feature extraction part) while modifying the classification part with the use of Adam optimizer to deal with classifying healthy controls and T2D based on single-cell gene regulatory network (SCGRN) images. Another DTL models work in a similar manner but just with keeping weights of the bottom layers in the feature extraction unaltered while updating weights of consecutive layers through training from scratch. Experimental results on the whole 224 SCGRN images using five-fold cross-validation show that our model (TFeSEResNeXT101) achieving the highest average balanced accuracy (BAC) of 0.97 and thereby significantly outperforming the baseline that resulted in an average BAC of 0.86. Moreover, the simulation study demonstrated that the superiority is attributed to the distributional conformance of model weight parameters obtained with Adam optimizer when coupled with weights from a pre-trained model.

Keywords Deep transfer learning, Optimizers, Explainable AI, T2D prediction, Single-cell gene regulatory network

Type 2 diabetes (T2D) is a common condition that over time when left untreated can cause damage reaching various organs, including kidney, eye, and heart, to just name a few^{1,2}. Patients with diabetes incur an overall average medical expenditure more than two times that of those without diabetes. Therefore, diabetes is considered as a burden associated with higher medical costs, and increased mortality rates³. Obtaining a highly accurate tool to discriminate between healthy and T2D subjects can aid in disease diagnosis, management, prevention and understanding⁴. Therefore, scientific efforts have been made contributing to detect T2D using computational methods^{5–8}.

Pyrros et al.⁹ employed deep learning (DL)-based approach to identify T2D using chest X-ray (CXR) images pertaining to healthy and T2D patients working as follows. They employed the ResNet34 DL model incorporating typical data augmentation with the use of Adam optimizer¹⁰. The dataset to develop (i.e., train from scratch to induce) the ResNet34 model consisted of 271,065 CXR training images in which 45,961 were designated as T2D CXR images and the remaining 225,104 were as CXR images for healthy control subjects. The trained model was then applied for a testing set consisting of 9943 CXR images. Results demonstrated that the DL model achieved an AUC of 0.84 when compared to an AUC of 0.79 for the baseline linear regression (LR) incorporating only clinical information data. An ensemble of both LR and ResNet34 generated an AUC of 0.85, considered as a marginal improvement in the prediction performance. These results demonstrate the feasibility of DL in

¹Department of Computer Science, King Abdulaziz University, 21589 Jeddah, Saudi Arabia. ²Department of Computer Science, Albaha University, 65799 Albaha, Saudi Arabia. ✉email: tturki@kau.edu.sa

screening T2D patients using CXR images. Wachinger et al.¹¹ presented a DL approach to predict T2D based on neck-to-knee MRI images and clinical information. The MRI images dataset consisted of 3406 MRI images in which the class distribution is uniformed (i.e., 1703 as T2D and 1703 as MRI images for healthy control subjects). The DL approach consists of convolutional layers, maxpooling layers, batch normalization layers, dropout layer and incorporation of dynamic affine feature map transform (DAFT) within convolutional layers to concatenate features obtained from clinical and MRI image data. Five-fold cross-validation was utilized to assess the performance of the whole data. Results demonstrated the superiority of CNN-DAFT achieving an AUC of 0.871, significantly outperforming CNN using only MRI images and linear regression using only clinical information.

Das¹² et al. presented a learning-based approach combining deep and machine learning for the diagnosis of T2D based on DNA sequences working as follows. First, they transformed the DNA sequence pertaining to healthy and T2D to images, provided as input to ResNet¹³ and VGG19 DL models to extract features. Then, providing the extracted features along with corresponding class labels to machine learning algorithms, namely support vector machines(SVM) and KNN. Experimental results using cross-validation on the whole image dataset demonstrate the good performance of SVM when coupled with extracted features using ResNet DL model. Naveed et al.¹⁴ employed DL to predict T2D. The dataset consisted of 19,181 patient records data in which 7715 records were for diabetic patients while the remaining 11,466 records were for non-diabetic patients. The dataset was divided into training and testing in which training composed of 80% of the dataset while the remaining 20% was assigned for testing. DL models included CNN, LSTM, and CNN-LSTM¹⁵. Experimental results demonstrated the superior performance of CNN-LSTM achieving the highest performance results when compared to other models including decision tree and SVM. Specifically, CNN-LSTM generated an accuracy of 91.6, and F1-Score of 89.2. These results demonstrate the feasibility of DL in early predicting T2D. Other AI-driven computational methods have been proposed to aid in predicting T2D^{8,16}.

As inferred single-cell gene regulatory networks (SCGRNs) encode the molecular interactions pertaining to components of specific cell types and thereby can aid in characterizing cellular differentiation in healthy and disease subjects^{17,18}, Turki et al.¹⁹ presented a novel DL approach to discriminate between healthy controls and T2D based on SCGRN images working as follows. Because rapid progress in single-cell technologies has contributed to the availability of biological experiments pertaining to gene regulatory networks (GRNs), single-cell gene expression data from the ArrayExpress repository was processed with the use of bigScale and NetBioV packages^{20–22}, generating 224 SCGRN images. The class distribution was distributed evenly in terms of healthy controls and T2D images. Then, utilizing RMSprop optimizer¹⁵ with the following DL models: VGG16²³, VGG19²³, Xception²⁴, ResNet50¹³, ResNet101¹³, DenseNet121²⁵, and DenseNet169²⁵ to discriminate between healthy controls and T2D SCGRN images. Experimental results demonstrated the VGG19 performed better than studied DL models. However, no interpretation was provided to back the prediction performance.

Although these recently developed methods aimed to address the task of prediction T2D, these methods are still far from perfect and do not provide interpretation for practical deep transfer learning (DTL) models aiding in the explanation of the performance superiority. Therefore, this study is unique in the following aspects: (1) we present highly accurate DTL models working by combining weight parameters from pre-trained models and weights obtained with the use of Adam optimizer; (2) we provide, to the best of our knowledge, the first interpretation behind DTL models inspecting and quantifying that conformance of pre-trained model weight parameters with weight parameters obtained with the incorporation of Adam and RMSprop Optimizers. This interpretation framework can guide in the process of designing highly efficient DTL models applicable to wide range of problems; and (3) we conduct experimental study to report the prediction performance and computational running time for the task of predicting single-cell gene regulatory network images pertaining to healthy controls and T2D. Experimental results demonstrate the superiority of our DTL model, TFeSEResNeXT101, performing better than the baseline with 11% improvements. In terms of the running time, our DL models exhibited a significant reduction in training time attributed to transfer learning, which reduced the number of trainable weight parameters. In addition, simulation study unveiled the conformance of parameter weights of both transfer weights from pre-trained models with weights obtained from Adam optimizer as compared to RMSprop that was used by the baseline and resulted in inferior prediction performance, attributed to the divergence of its weight parameters from the weight parameters of pre-trained models.

Materials and methods

Biological networks

We provide an illustration in Fig. 1 for the biological network images used in this study, which were downloaded from¹⁹ and consisted of 224 SCGRN images pertaining to healthy and T2D. The class distribution for these biological network images is balanced (i.e., 224 divided evenly into the two classes). These biological network images were produced with the help of bigScale package to process the single-cell gene expression data and build regulatory networks, then visualizing networks via the NetBioV package. In terms of the single-cell gene expression data pertaining to healthy controls and T2D patients, it was obtained from ArrayExpress repository under accession number E-MTAB-5061²⁶.

Deep transfer learning

Figure 1 demonstrates how our deep transfer learning (DTL) approach is performed. First, we adapt the following pre-trained models: VGG19, DenseNet201²⁵, InceptionV3²⁷, ResNet50V2²⁸, ResNet101V2²⁸, SEResNet152²⁹, and SEResNeXT101²⁹. Each pre-trained model has a feature extraction part (i.e., series of convolutional and pooling layers) for feature extraction and a densely connected classifier for classification. Then, we keep the weights unchanged for the feature extraction part of a pre-trained model and change the densely connected classifier to deal with binary classification instead of 1000 classes. Therefore, when feeding the SCGRN image dataset, we

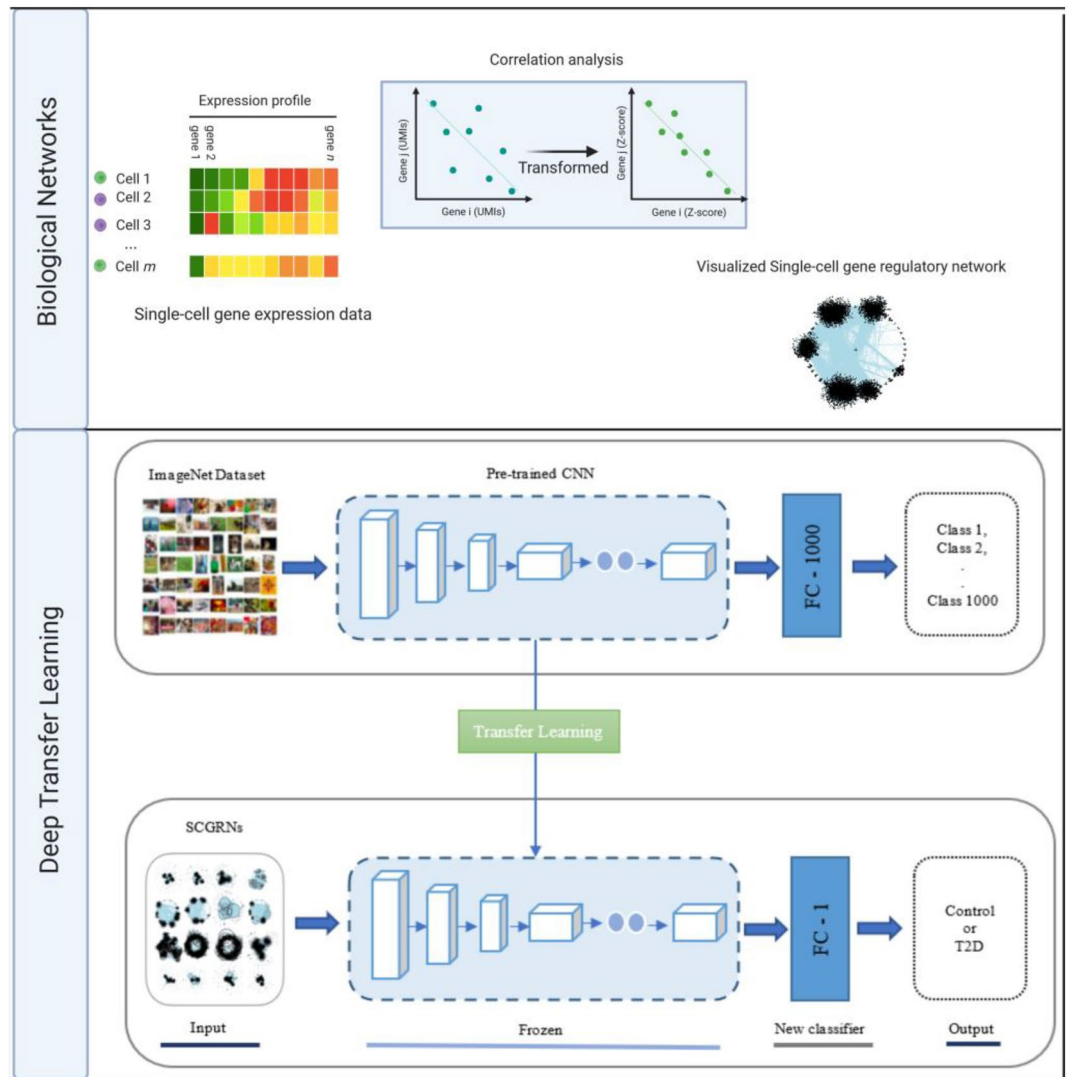


Figure 1. Flowchart of the deep transfer learning-based approach for the predicting T2D using SCGRNs. **Biological Networks:** To infer single-cell gene regulatory network (SCGRN), gene expression data are provided to bigScale (performing clustering and differential expression analysis) changing measured correlation between genes from expression values to Z-score, followed by retaining significant correlations to guide in building a regulatory network. A visualization is performed using NetBioV. **Deep Transfer Learning:** Transfer learning applying feature extraction with new classifier (TFe) to distinguish between T2D and healthy control SCGRNs.

extract features using weights of pre-trained models while training the densely connected classifier from scratch and performing prediction. We refer to models using this type of DTL computations as TFeVGG19, TFeDenseNet201, TFeInceptionV3, TFeResNet50V2, TFeResNet101V2, TFeSEResNet152, and TFeSEResNeXT101 (see Fig. 1). For the other DTL computations, we keep weights of the bottom layers unchanged in the feature extraction part while performing training from scratch to change weights of top layers in feature extraction part and densely connected layers.

As in TFe-based models, we modify the densely connected classifier dealing with binary classification problem before performing the training phase. As seen in Fig. 2, we refer to models employing this type of deep transfer learning as TFtVGG19, TFtDenseNet201, TFtInceptionV3, TFtResNet50V2, TFtResNet101V2, TFtSEResNet152, and TFtSEResNeXT101.

When changing weights during training, we employed three optimizers: Adam, RMSprop, and SGD³⁰. When weights are kept unchanged referring to the transfer of knowledge from pre-trained models using SGD optimizer. In terms of predictions of unseen SCGRN images, predictions are mapped to healthy control subjects if the predicted values are greater than 0.5. Otherwise, predictions are mapped to T2D.

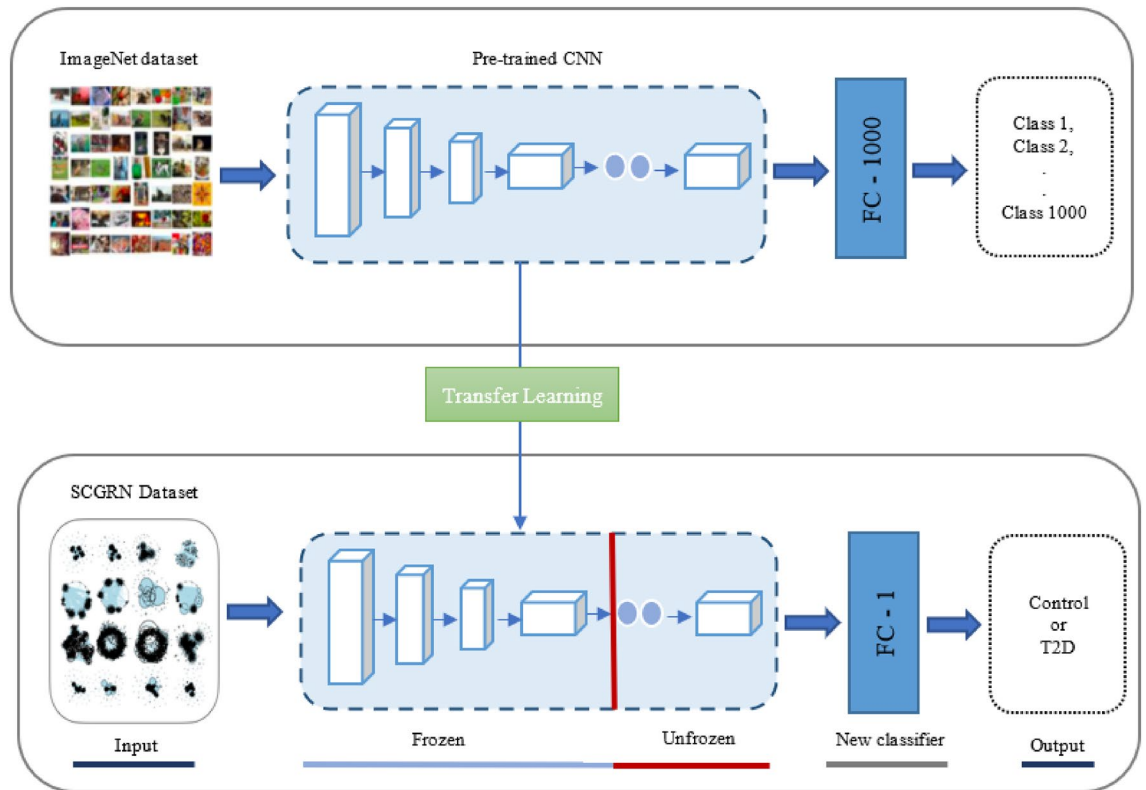


Figure 2. Transfer learning applying fine tuning with new classifier (TfT) to distinguish between T2D and healthy control SCGRNs.

Results

Classification methodology

In this study, we considered seven pre-trained models, namely VGG19, DenseNet201, InceptionV3, ResNet50V2, ResNet101V2, SEResNet152, and SEResNeXT101. Each of the pre-trained models was trained on 1.28 million images from ImageNet database to classify images into 1000 different categories. In terms of TFe-based models, we used the feature extraction part of pre-trained models in which weights were kept unchanged and were used to extract feature from SCGRN images. Moreover, the densely connected classifier was trained from scratch to handle the binary class classification problem. Regarding the TfT-based models, we trained the top layers and densely connected classifier from scratch while retaining the weights of bottom layers unchanged in the feature extraction part. For both TfT-based and TFe-based models, we employed Adam optimizer when updating weights of layers. Moreover, we compared the performance of our deep transfer learning approaches using different optimizers including the baseline (i.e., RMSprop optimizer) as well as against training models from scratch. We set optimization parameters as follows: 0.00001 for the learning rate, 10 for the number of epochs, and 32 for the batch size. In terms of the loss function, we utilized categorical cross-entropy³¹.

To assess the performance of studied models, we employed Balanced Accuracy (BAC), Accuracy (ACC), Precision (PRE), Recall (REC), and F1 computed as follows:

$$BAC = \frac{1}{2} \left(\frac{TP}{TP + FN} + \frac{TN}{TN + FP} \right) \tag{1}$$

$$ACC = \frac{TP + TN}{TP + TN + FP + FN} \tag{2}$$

$$PRE = \frac{TP}{TP + FP} \tag{3}$$

$$REC = \frac{TP}{TP + FN} \tag{4}$$

$$F1 = \frac{2 * PRE * REC}{PRE + REC} \tag{5}$$

where *TN* designates true negative, corresponding to the number of T2D images that were correctly predicted as T2D. *FP* designates false positive, corresponding to the number of T2D images that were incorrectly predicted as healthy controls. *TP* designates true positive, corresponding to the number of healthy control images that were correctly predicted as healthy controls. *FN* designates false negative, corresponding to the number of healthy control images that were incorrectly predicted as T2D.

To evaluate the results on the whole SCGRN image dataset, we employed five-fold cross-validation as follows. We partitioned the SCGRN image datasets and randomly assigned images into 5 folds. During the first run of five-fold cross-validation, we used 4 of the folds to train our deep learning models and perform predictions to the remaining fold for testing and record the performance results. Such a process was repeated for an additional 4 runs in which performance results were recorded. Finally, we report the average performance results corresponding to the results obtained from five-fold cross-validation.

Implementation details

All experiments were run on a machine equipped with central processing unit (CPU) of Google Colab. The specifications of CPU runtime offered by Google Colab were Intel Xeon Processor with two cores with 2.30 GHz and 13 GB RAM where the installed version of Python is 3.10.11. For the analysis of models, we used R statistical software³² to run the experiments and utilized the *optim* package in R to run Adam optimizer³³. All plots were performed using Matplotlib package in python³⁴.

Classification results

Training results

In Fig. 3, we illustrate the training accuracy performance results when running five-fold cross-validation. It can be seen that our models outperformed all other models trained from scratch. Specifically, TFeVGG19 and TFtVGG19 achieved average accuracies of 0.976 and 0.962, respectively, while VGG19 achieved an average accuracy of 0.530. TFeDenseNet201 outperformed DenseNet201 via achieving an average accuracy of 0.988 while DenseNet201 performed better than TFtDenseNet201 via achieving an average accuracy of 0.982 compared to 0.946. For TFe- and TFt-based models when coupled with ResNet101V2, SEResNet152 and SEResNetXT101, they outperformed their counterparts when not applying deep transfer learning (DTL) models. These superior performance results are attributed to the learned representation using transfer learning.

Testing results

Figures 4 and 5 report the generalization (i.e., test) accuracy performance results and combined confusion matrices, respectively, when five-fold cross-validation is utilized. TFeSEResNeXT101 achieved the highest average accuracy of 0.968.

The second-best model is TFeDenseNet201, achieving an average accuracy of 0.958, followed by TFeVGG19, TFeResNet50V2, TFeSEResNet152, TFeInceptionV3, and TFeResNet101V2 (generating average accuracies of 0.946, 0.940, 0.936, 0.930, and 0.918, respectively). TFt-based models also outperformed all models trained from scratch (see Fig. 4b,c). Particularly, TFt-based models generated average accuracies lower and upper bounded by 0.864 and 0.916, respectively, while models trained from scratch were lower and upper bounded by average accuracies of 0.468 and 0.590. These results demonstrate the superior performance of models employing our DTL computations.

In terms of reporting testing performance results using different metrics, our model TFeSEResNeXT101 outperforms all other models (see Table 1) via achieving an average BAC of 0.97, average PRE of 0.97 (tie with our model TFeSEResNet152), and average F1 of 0.97. Moreover, TFeVGG19 and TFtVGG19 perform better than VGG19. Similarly, TFeDenseNet201, TFeInceptionV3, TFeResNet50V2, TFeResNet101V2, and TFeSEResNet152 performed better than DenseNet201, Inception V3, ResNet50V2, ResNet101V2, and SEResNet152, respectively. The same holds true for TFt-based models outperforming their counterparts (i.e., VGG19, DenseNet201, InceptionV3, ResNet50V2, ResNet101V2, SEResNet152, and SEResNetXT101).

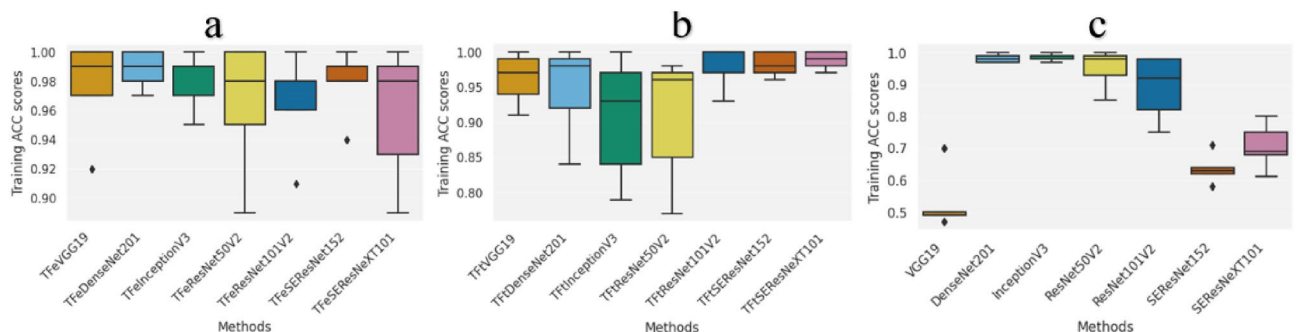


Figure 3. The boxplots presenting the average five-fold cross-validation results using the ACC measure for the training folds. (a) Deep transfer learning models using feature extraction (referred with the prefix TFe). (b) Deep transfer learning models using fine tuning (referred with the prefix TFt). (c) Deep learning models trained from scratch. ACC is accuracy.

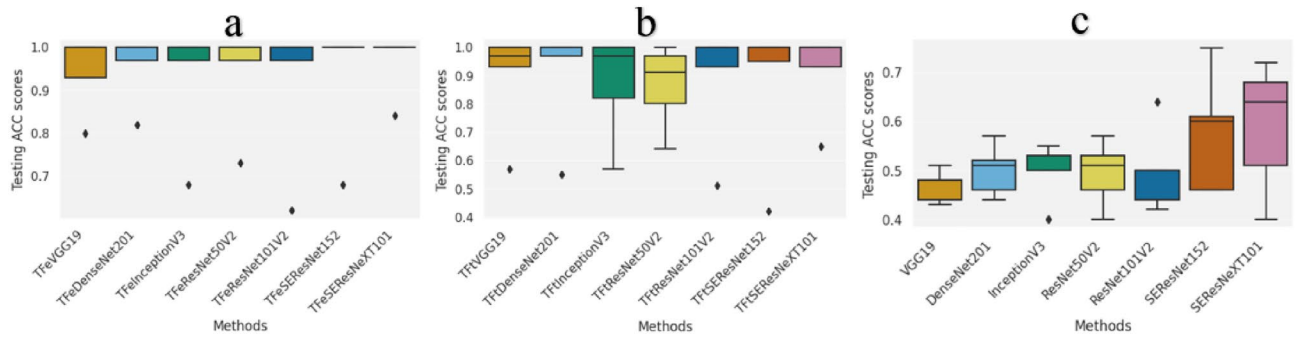


Figure 4. The boxplots presenting the average five-fold cross-validation results using the ACC measure for the testing folds. (a) Deep transfer learning models using feature extraction (referred with the prefix TFe). (b) Deep transfer learning models using fine tuning (referred with the prefix TFt). (c) Deep learning models trained from scratch. ACC is accuracy.

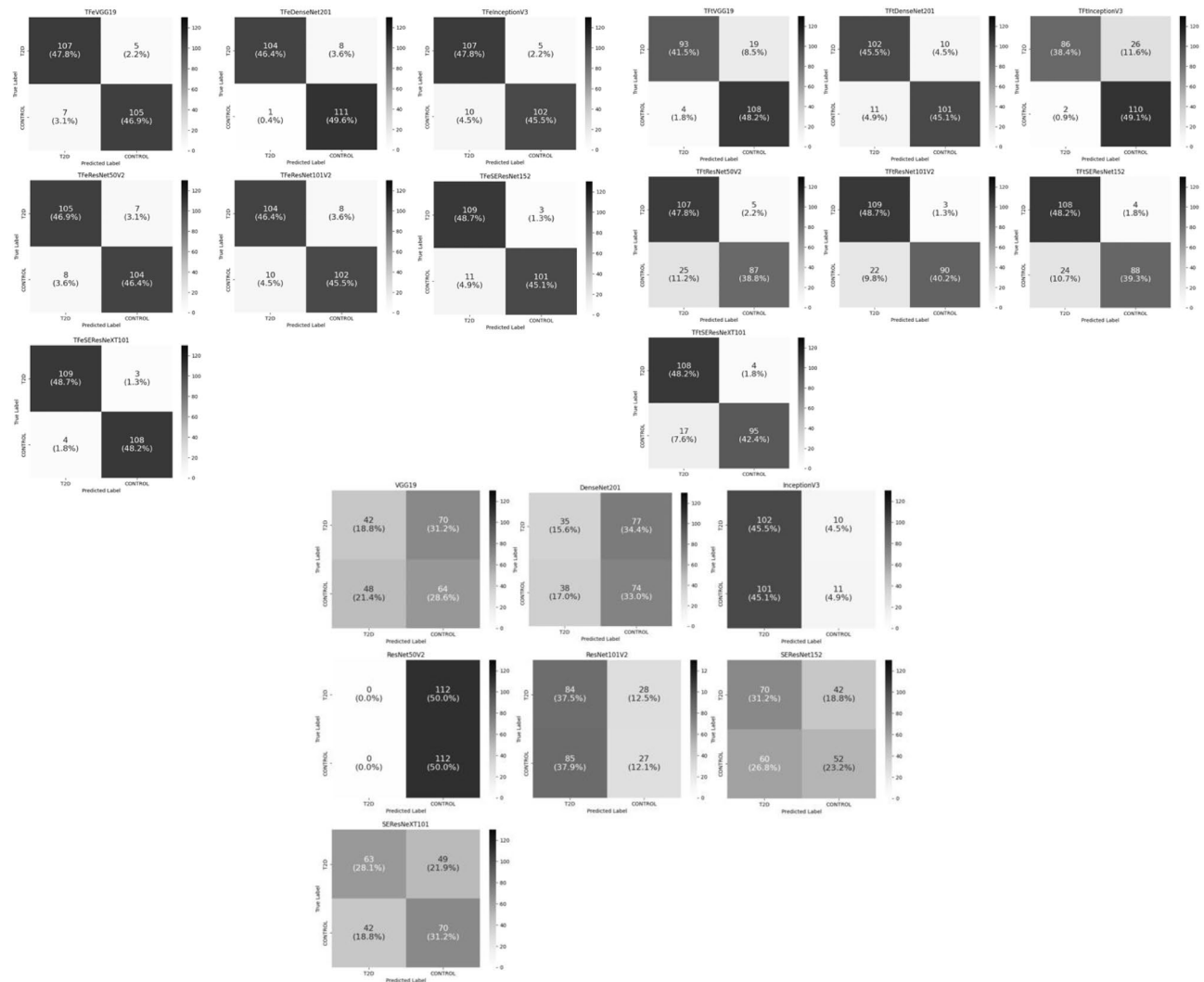


Figure 5. Combined confusion matrices for all methods during the running of five-fold cross-validation.

Table 2 reports our best DTL models with different optimizers. It can be shown that TFeSEResNeXT101 and TFtSEResNeXT101 generate the highest performance results when coupled with Adam optimizer method. Specifically, TFeSEResNeXT101 with Adam optimizer generates the highest average BAC (and F1) of 0.97 (and 0.97). TFtSEResNeXT101 with Adam optimizer archives the highest average BAC of 0.91, highest average F1 of 0.90 (tie with SGD optimizer). When TFeSEResNeXT101 and TFtSEResNeXT101 are coupled with SGD optimizer, they generate inferior performance results.

Model	Optimizer	BAC	PRE	REC	F1
TFeVGG19	Adam	<u>0.94</u>	<u>0.95</u>	0.94	<u>0.94</u>
TFtVGG19	Adam	0.90	0.82	0.97	0.89
VGG19	Adam	0.50	0.37	0.46	0.41
TFeDenseNet201	Adam	<u>0.96</u>	<u>0.95</u>	<u>0.95</u>	<u>0.96</u>
TFtDenseNet201	Adam	0.90	0.91	0.90	0.90
DenseNet201	Adam	0.50	0.20	0.51	0.29
TFeInceptionV3	Adam	<u>0.93</u>	<u>0.95</u>	0.91	<u>0.92</u>
TFtInceptionV3	Adam	0.87	0.76	<u>0.92</u>	0.83
InceptionV3	Adam	0.50	0.91	0.50	0.64
TFeResNet50V2	Adam	<u>0.93</u>	0.93	<u>0.92</u>	<u>0.93</u>
TFtResNet50V2	Adam	0.87	<u>0.95</u>	0.81	0.87
ResNet50V2	Adam	0.50	–	–	–
TFeResNet101V2	Adam	<u>0.92</u>	<u>0.92</u>	<u>0.91</u>	<u>0.92</u>
TFtResNet101V2	Adam	0.89	0.90	0.86	0.88
ResNet101V2	Adam	0.52	0.75	0.49	0.59
TFeSEResNet152	Adam	<u>0.94</u>	0.97	<u>0.90</u>	<u>0.93</u>
TFtSEResNet152	Adam	0.88	0.96	0.81	0.88
SEResNet152	Adam	0.55	0.67	0.50	0.57
TFeSEResNeXT101	Adam	0.97	0.97	<u>0.96</u>	0.97
TFtSEResNeXT101	Adam	0.90	0.96	0.90	0.91
SEResNeXT101	Adam	0.60	0.56	0.60	0.58

Table 1. Reported average performance results during the running of five-fold cross-validation on testing using studied models. BAC is balanced accuracy. PRE is precision. REC is recall. The best overall result is underlined and is shown in bold. The method outperforming its counterparts is just underlined.

Model	Optimizer	BAC	F1
TFeSEResNeXT101	Adam	0.97	0.97
	RMSprop	0.92	0.92
	SGD	0.77	0.77
TFtSEResNeXT101	Adam	0.91	0.90
	RMSprop	0.90	0.89
	SGD	0.90	0.90

Table 2. Performance comparison of our best deep transfer learning model under different optimizers during the five-fold cross-validation. BAC is balanced accuracy. Best performance result is shown in bold.

In Table 3, we compare our model TFeVGG19 with Adam optimizer against the best performing baseline TFeVGG19 with RMSprop optimizer, named VGG19 in¹⁹. It is evident that our model TFeVGG19 with Adam optimizer achieves the highest average BAC of 0.94 while the baseline obtained an average BAC of 0.86. Moreover, when F1 performance measure is considered, TFeVGG19 with Adam optimizer attains the highest average F1 of 0.94 while the baseline achieved an average F1 of 0.88. The same holds true for TFtVGG19, which achieved the highest average BAC of 0.91, highest average F1 of 0.90.

In Fig. 6, we report the running time in seconds for the process of running five-fold cross-validation when utilizing our best model (TFeSEResNeXT101) and TFtSEResNeXT101 compared to their peer SEResNeXT101.

Model	Optimizer	BAC	F1
TFeVGG19	Adam (ours)	0.94	0.94
	RMSprop (baseline)	0.86	0.88
TFtVGG19	Adam (ours)	0.91	0.90
	RMSprop (baseline)	0.89	0.89

Table 3. Performance comparison of our deep transfer learning model against recent baseline methods when five-fold cross-validation is employed. BAC is balanced accuracy. Best performance result is shown in bold.

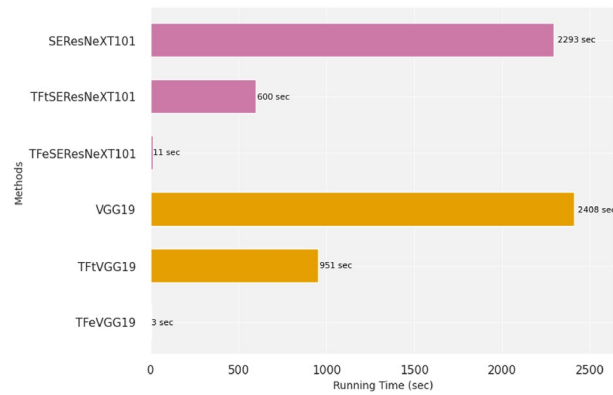


Figure 6. Running time comparisons in seconds for selected models when running five-fold cross-validation.

Our model TFeSEResNeXT101 is $208.45 \times$ faster than SEResNeXT101. Also, our model TFtSEResNeXT101 is $3.82 \times$ faster than SEResNeXT101. Moreover, TFeVGG19 and TFtVGG19 are $802.67 \times$ and $2.53 \times$, respectively, faster than VGG19. These results demonstrate the computational efficiency of the DTL models, in addition to the highly achieved performance results.

Models introspection

Stochastic gradient descent (SGD)

To minimize the objective function $Q(\theta_0, \theta_1)$ for parameters θ_0 and θ_1 of model $H(x_i)$, we employ gradient descent optimization algorithms to find θ_0 and θ_1 minimizing the objective function. The optimization problem can be formulated as follows:

$$\begin{aligned} \min_{\theta_0, \theta_1} Q(\theta_0, \theta_1) &= \frac{1}{m} \sum_{i=1}^m Q_i(\theta_0, \theta_1) \\ &= \frac{1}{m} \sum_{i=1}^m (H(x_i) - y_i)^2 \\ &= \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i)^2 \end{aligned} \quad (6)$$

We utilize SGD, RMSprop, and Adam optimization algorithm to minimize the objective function and estimate the model parameters. For SGD, we initialize the parameters θ_0 and θ_1 according to the uniform distribution $U(0, 1)$ and setting the learning rate $\eta = 0.001$, maximum number of iterations to 3000. Then, in each time, shuffling the data of m examples followed by looping m times over the following to update model parameters. After the end of looping, the algorithm stops if the maximum number of iterations is reached or $\|\nabla Q(\theta_0, \theta_1)\| \leq 0.001$:

$$\begin{aligned} \theta_0 &:= \theta_0 - \eta \cdot \frac{\partial Q_i}{\partial \theta_0} \\ &:= \theta_0 - \eta \cdot 2(\theta_0 + \theta_1 x_i - y_i) \\ \theta_1 &:= \theta_1 - \eta \cdot \frac{\partial Q_i}{\partial \theta_1} \\ &:= \theta_1 - \eta \cdot 2x_i(\theta_0 + \theta_1 x_i - y_i) \end{aligned} \quad (7)$$

Root mean square propagation (RMSprop)

For RMSprop, we initialize model parameters according to uniform distribution $U(0, 1)$, set the learning rate $\eta = 0.001$, maximum number of iterations to 3000, $\beta = 0.9$, $\epsilon = 10^{-6}$, BatchSize = 16, referred in the following as $[S]$. Then, looping to update model parameters according to each selected batch. After ending of looping over all selected batches, the algorithm terminates when $\|\nabla Q(\theta_0, \theta_1)\| \leq 0.001$ or the maximum number of iterations is reached:

$$\begin{aligned}
v_{\theta_0} &:= \beta \cdot v_{\theta_0} + (1 - \beta) \left(\frac{\partial Q}{\partial \theta_0} \right)^2 \\
&:= \beta \cdot v_{\theta_0} + (1 - \beta) \left(\frac{2}{|S|} \sum_i (\theta_0 + \theta_1 x_i - y_i) \right)^2 \\
v_{\theta_1} &:= \beta \cdot v_{\theta_1} + (1 - \beta) \left(\frac{\partial Q}{\partial \theta_1} \right)^2 \\
&:= \beta \cdot v_{\theta_1} + (1 - \beta) \left(\frac{2}{|S|} \sum_i (\theta_0 + \theta_1 x_i - y_i) \cdot x_i \right)^2 \\
\theta_0 &:= \theta_0 - \frac{\eta}{\sqrt{v_{\theta_0} + \epsilon}} \frac{\partial Q}{\partial \theta_0} \\
&:= \theta_0 - \frac{\eta}{\sqrt{v_{\theta_0} + \epsilon}} \left(\frac{2}{|S|} \sum_i (\theta_0 + \theta_1 x_i - y_i) \right) \\
\theta_1 &:= \theta_1 - \frac{\eta}{\sqrt{v_{\theta_1} + \epsilon}} \frac{\partial Q}{\partial \theta_1} \\
&:= \theta_1 - \frac{\eta}{\sqrt{v_{\theta_1} + \epsilon}} \left(\frac{2}{|S|} \sum_i (\theta_0 + \theta_1 x_i - y_i) x_i \right)
\end{aligned} \tag{8}$$

Adaptive moment estimation (Adam)

In terms of Adam, m_{θ_0} , m_{θ_1} , v_{θ_0} , v_{θ_1} , β_1 , β_2 , ϵ , and η are initialized as in³⁵. Then, looping to update model parameters according to all selected batches of examples. After ending of looping over all selected batches, the algorithm stops if the maximum number of iterations (i.e., 3000) is reached or $\|\nabla Q(\theta_0, \theta_1)\| \leq 0.001$:

$$\begin{aligned}
m_{\theta_0} &:= \beta_1 \cdot m_{\theta_0} + (1 - \beta_1) \left(\frac{\partial Q}{\partial \theta_0} \right) \\
&:= \beta_1 \cdot m_{\theta_0} + (1 - \beta_1) \left(\frac{2}{|S|} \sum_i (\theta_0 + \theta_1 x_i - y_i) \right) \\
m_{\theta_1} &:= \beta_1 \cdot m_{\theta_1} + (1 - \beta_1) \left(\frac{\partial Q}{\partial \theta_1} \right) \\
&:= \beta_1 \cdot m_{\theta_1} + (1 - \beta_1) \left(\frac{2}{|S|} \sum_i (\theta_0 + \theta_1 x_i - y_i) x_i \right) \\
v_{\theta_0} &:= \beta_2 \cdot v_{\theta_0} + (1 - \beta_2) \left(\frac{\partial Q}{\partial \theta_0} \right)^2 \\
&:= \beta_2 \cdot v_{\theta_0} + (1 - \beta_2) \left(\frac{2}{|S|} \sum_i (\theta_0 + \theta_1 x_i - y_i) \right)^2 \\
v_{\theta_1} &:= \beta_2 \cdot v_{\theta_1} + (1 - \beta_2) \left(\frac{\partial Q}{\partial \theta_1} \right)^2 \\
&:= \beta_2 \cdot v_{\theta_1} + (1 - \beta_2) \left(\frac{2}{|S|} \sum_i (\theta_0 + \theta_1 x_i - y_i) \cdot x_i \right)^2 \\
\hat{m}_{\theta_0} &:= \frac{m_{\theta_0}}{(1 - \beta_1)} \\
\hat{m}_{\theta_1} &:= \frac{m_{\theta_1}}{(1 - \beta_1)} \\
\hat{v}_{\theta_0} &:= \frac{v_{\theta_0}}{(1 - \beta_2)} \\
\hat{v}_{\theta_1} &:= \frac{v_{\theta_1}}{(1 - \beta_2)} \\
\theta_0 &:= \theta_0 - \frac{\eta}{\sqrt{\hat{v}_{\theta_0} + \epsilon}} \hat{m}_{\theta_0} \\
\theta_1 &:= \theta_1 - \frac{\eta}{\sqrt{\hat{v}_{\theta_1} + \epsilon}} \hat{m}_{\theta_1}
\end{aligned} \tag{9}$$

Simulated data

To demonstrate the efficiency of the proposed deep transfer learning (DTL) models incorporation mixed parameters derived from both SGD and Adam, we conducted simulation studies to explain the superiority behind the proposed models as well as imitate the numerical behavior. Particularly, we consider the following four predictive models:

$$F_1(x) = 0.5 + 0.79x + \epsilon \tag{10}$$

$$F_2(x) = 1 - \exp\left(\frac{-1}{2x}\right) \tag{11}$$

$$F_3(x) = 0.7 + 3x^2 + \epsilon \tag{12}$$

$$F_4(x) = (6x - 2)^2 + \sin(12x - 4) \tag{13}$$

where $X \sim U(0, 1)$ and $\epsilon \sim N(0, 0.2)$ in which $U()$ and $N()$ are uniform and normal distributions, respectively. For F_1 , when we have X and $Y = F_1(X)$, we perform the following steps. Let $H_{SGD}(x_i) = \theta_0 + \theta_1 x_i$ (for $i = 1..m$) be the model in which we want to estimate parameters using (X, Y) data from F_1 coupled with Eq. (7). Similarly, let $H_{RMSprop}(x_i)$ and $H_{Adam}(x_i)$ be models in which we want to estimate their parameters using Eqs. (8) and (9), respectively, coupled with (X, Y) data from F_1 . Then, we provide each $x_i \in X$ to perform predictions corresponding to y'_i . In Fig. 7, we report 2D plots for X and predicted $Y' = \{y'_1, \dots, y'_m\}$ via each model using data generated according to Eq. (10), where SGD refers to plotting $(x_i, H_{SGD}(x_i))$ while Adam and RMSprop refer to plotting results obtained via $(x_i, H_{Adam}(x_i))$ and $(x_i, H_{RMSprop}(x_i))$, respectively, and $i = 1..m$. We then repeat this process for an additional 8 runs. Therefore, we have 9 runs in total.

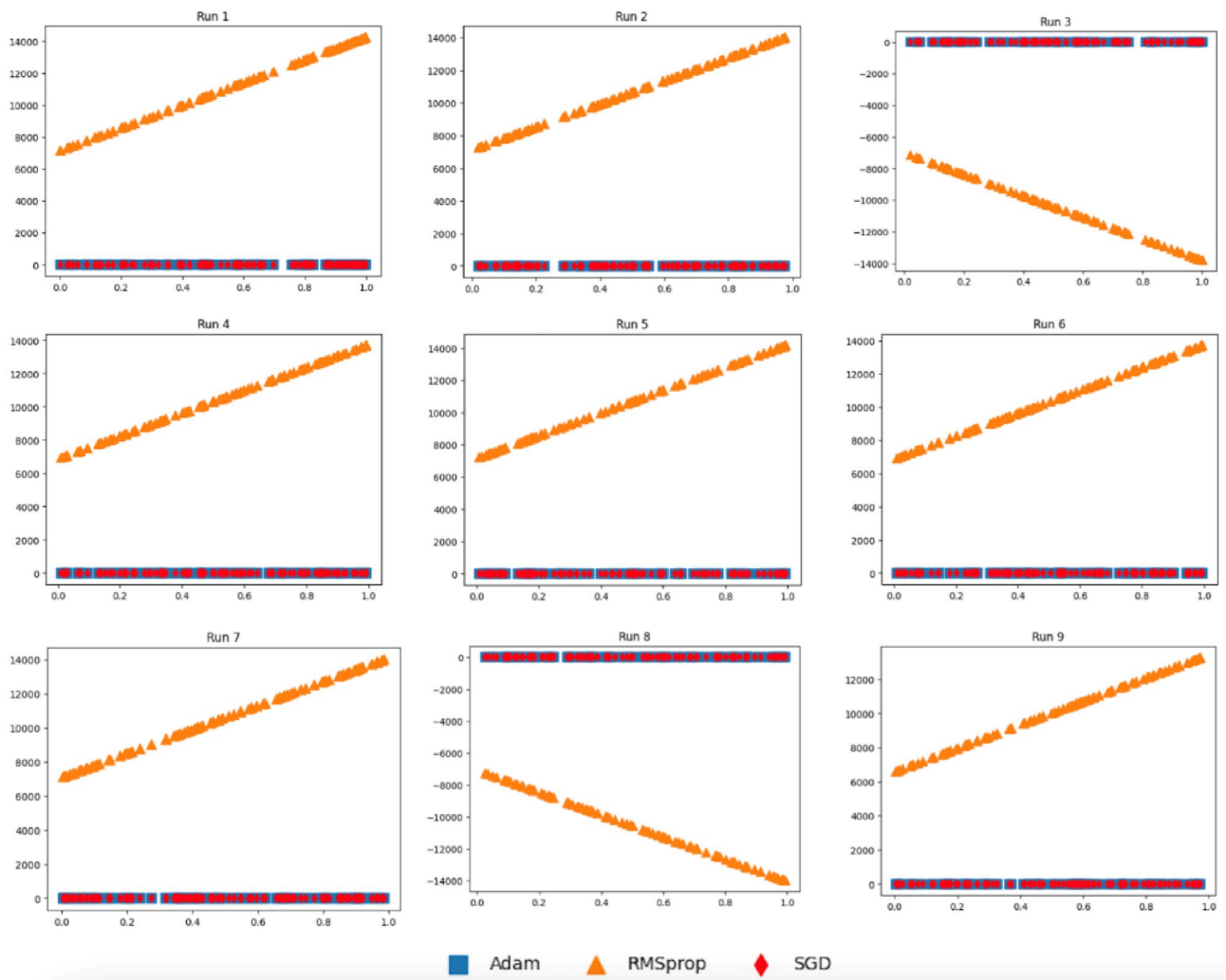


Figure 7. Plots for the three models as $(x_i, H_{SGD}(x_i))$, $(x_i, H_{Adam}(x_i))$, and $(x_i, H_{RMSprop}(x_i))$ for $i = 1..m$ according to x_i generated using F_1 .

It can be seen from Fig. 7 that model induced via RMSprop has more distributional differences compared to those obtained via SGD and Adam. To quantify distributional differences between SGD and Adam against SGD and RMSprop, we perform the following computations:

$$d_{SA} = \sqrt{\sum_i \left((x_i^{SGD} - x_i^{Adam})^2 + (z_i^{SGD} - z_i^{Adam})^2 \right)} \tag{14}$$

$$d_{SR} = \sqrt{\sum_i \left((x_i^{SGD} - x_i^{RMSprop})^2 + (z_i^{SGD} - z_i^{RMSprop})^2 \right)} \tag{15}$$

where d_{SA} measuring the distance between data associated with SGD and Adam. Similarly, d_{SR} measures the distance between data associated with SGD and RMSprop. The lower the distance value, the less the distributional difference is. Figure 11a plots d_{SA} and d_{SR} for the 9 runs. It can be seen that $H_{Adam}(x_i)$ is closer to $H_{SGD}(x_i)$ than $H_{RMSprop}(x_i)$ to $H_{SGD}(x_i)$ in most runs. Moreover, the distributional differences are statistically significant (P -value = 7.28×10^{-14} from t -test).

These results demonstrate conformance of the weight parameters of models utilizing Adam and SGD optimizers. Figure 8 reports the 2D plots of three induced models as $(x_i, H_{SGD}(x_i))$, $(x_i, H_{Adam}(x_i))$, and $(x_i, H_{RMSprop}(x_i))$ for $i = 1..m$ using data generated according to Eq. (11) (i.e., F_2)³⁶. It can be clearly seen that the data distributional difference of results via SGD is closer to that of Adam when compared to results obtained with the help of RMSprop. In Fig. 11b, we quantify distributional differences using Eqs. (14) and (15). It can be shown that Adam is closer to SGD as shown from AdamSGD when compared to that of RMSprop to SGD (i.e., RMSpropSGD) over the 9 runs. The quantification of AdamSGD is attributed to d_{SA} while RMSpropSGD is attributed to d_{SR} .

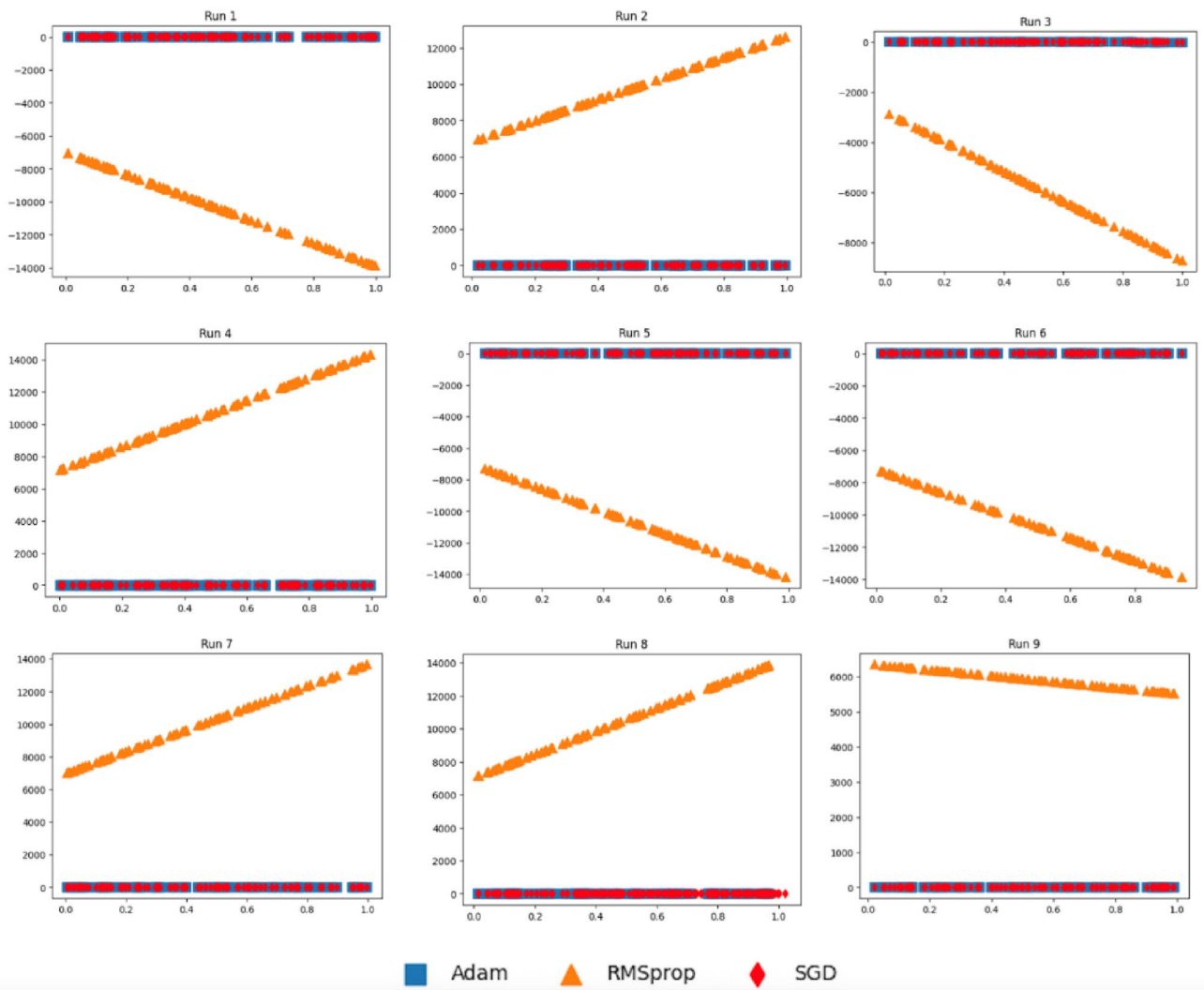


Figure 8. Plots for the three models as $(x_i, H_{SGD}(x_i))$, $(x_i, H_{Adam}(x_i))$, and $(x_i, H_{RMSprop}(x_i))$ for $i = 1..m$ according to x_i generated using F_2 .

In addition, the distributional differences between AdamSGD and RMSpropSGD are statistically significant (P -value = 7.01×10^{-7} from t -test).

Figures 9 and 10 report 2D plots of $(x_i, H_{SGD}(x_i))$, $(x_i, H_{Adam}(x_i))$, and $(x_i, H_{RMSprop}(x_i))$ for $i = 1..m$ using generated data of Eqs. (12) (F_3) and (13) (F_4)³⁷, where models were induced with SGD, Adam, and RMSprop optimizers. It can be seen from the alignment of Adam with SGD that Adam has a closer data representation to SGD compared to RMSprop to SGD. When quantifying the data distributional differences in Fig. 11c and d, it can be clearly shown that the distributional differences of SGD and Adam (referred to AdamSGD) are closer than SGD to RMSprop over the 9 runs.

These quantified results for AdamSGD and RMSprop are attributed to d_{SA} and d_{SR} , respectively. Additionally, the distributional differences of between AdamSGD and RMSprop were statistically significant (P -value = 3.48×10^{-12} from t -test when F_3 is used while P -value = 1.49×10^{-3} from t -test when F_4 is used). These results demonstrate the stable performance when SGD is coupled with Adam.

Discussion

Our deep transfer learning (DTL) models work as follows. In the TFe-based models, the convolutional base (also called the feature extraction part) in the pre-trained model is left unchanged while the densely connected classifier is modified to deal with the binary class classification at hand. Therefore, we applied the features extraction part of pre-trained models to the SCGRN images to extract features followed by a flattening step to train densely connected classifier from scratch. It can be noted that only weights of densely connected classifier are changed according to Adam optimizer while we transferred knowledge (i.e., weights) of the feature extraction part from pre-trained models. In terms of the T Ft-based models, we keep weights of the bottom layers in the feature extraction part of pre-trained models unchanged while modifying weights in the proceeding layers including the densely connected classifier according to the Adam optimizer. Moreover, the densely connected

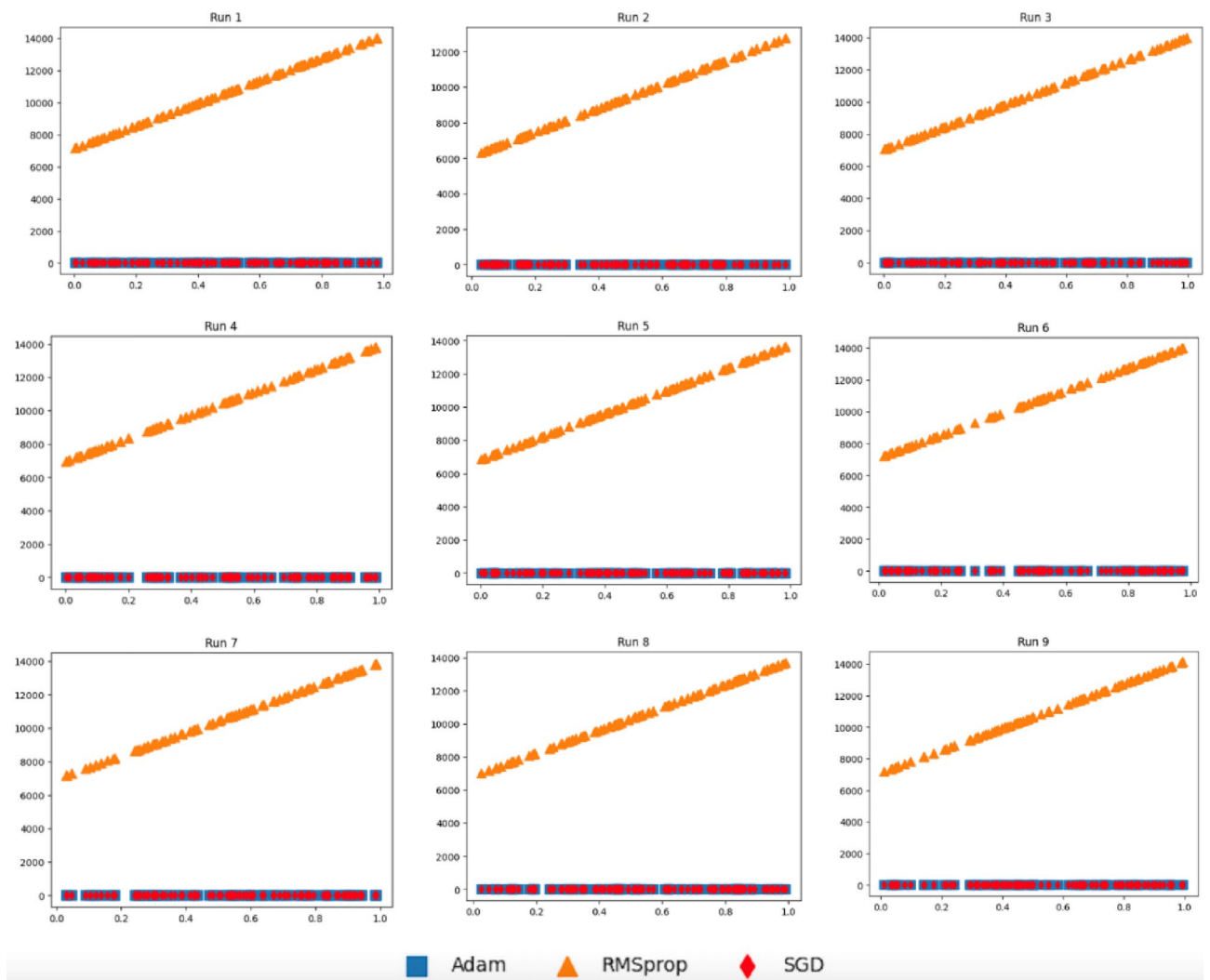


Figure 9. Plots for the three models as $(x_i, H_{SGD}(x_i))$, $(x_i, H_{Adam}(x_i))$, and $(x_i, H_{RMSprop}(x_i))$ for $i = 1..m$ according to x_i generated using F_3 .

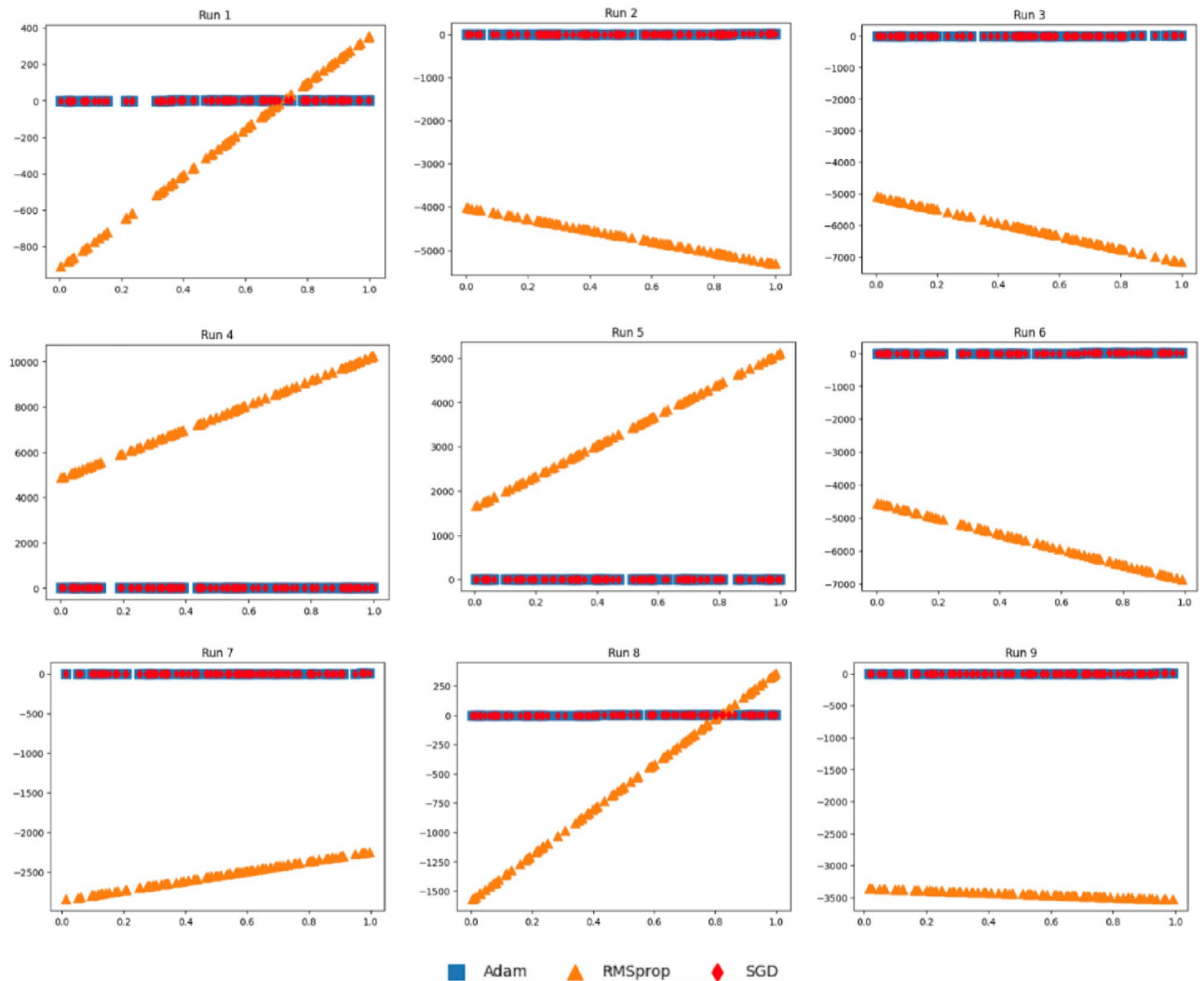


Figure 10. Plots for the three models as $(x_i, H_{\text{SGD}}(x_i))$, $(x_i, H_{\text{Adam}}(x_i))$, and $(x_i, H_{\text{RMSprop}}(x_i))$ for $i = 1..m$ according to x_i generated using F_4 .

classifier was altered to deal with the binary class classification problem pertaining to distinguishing between healthy controls and T2D SCGRN images. It can be seen that updating model weight parameters is done through the training with the use of Adam optimizer.

When conducting experimental study to assess deep transfer learning models, we used different optimizers, including SGD, RMSprop, and Adam. For SGD optimizer used in our DTL models, the model weight parameters were different than DTL models coupled with RMSprop and Adam optimizers. Therefore, we induced three sets of different models attributed to the three optimizers. Experimental results demonstrate the superiority of DTL models utilizing SGD and Adam optimizers when compared to that using SGD and RMSprop optimizers. We reported training loss for epochs pertaining to TFe-based and T Ft-based models when running five-fold cross-validation in Supplementary Fig. S1. For each DTL model, the number of layers including frozen and unfrozen layers is reported in Supplementary Table S1.

In our study, mitigation of overfitting is attributed to (1) transfer learning in which many layers in DTL models are frozen and thereby reducing the number of trainable parameters; and (2) applying the dropout to the last fully-connected layer in which we set the dropout rate to 0.5³⁸. It is worth mentioning that we assessed the performance of other deep learning (DL) models such as ConvNeXtTiny and ConvNeXtLarge³⁹. Although ConvNeXtLarge outperformed ConvNeXtTiny, they didn't exhibit superior performance when compared to TFeSEResNeXT101. Therefore, we include their performance results in Supplementary Tables S2 and S3 (and Supplementary Fig. S2). In terms of the running time, TFeSEResNeXT101 was $1.54 \times$ faster than TFeConvNeXtLarge. We report running time for ConvNeXtTiny-based and ConvNeXtLarge-based models in Supplementary Fig. S3.

In our DTL models, we have transferred weights from pre-trained models coupled with weights obtained with the help of Adam optimizer. If weight parameters obtained using two optimizers are close, then the two models almost behave the same. On the other hand, when the weight parameters resulted from two optimizers are not close, then the two models behave differently. To mimic the real scenario and investigate the effects of coupling

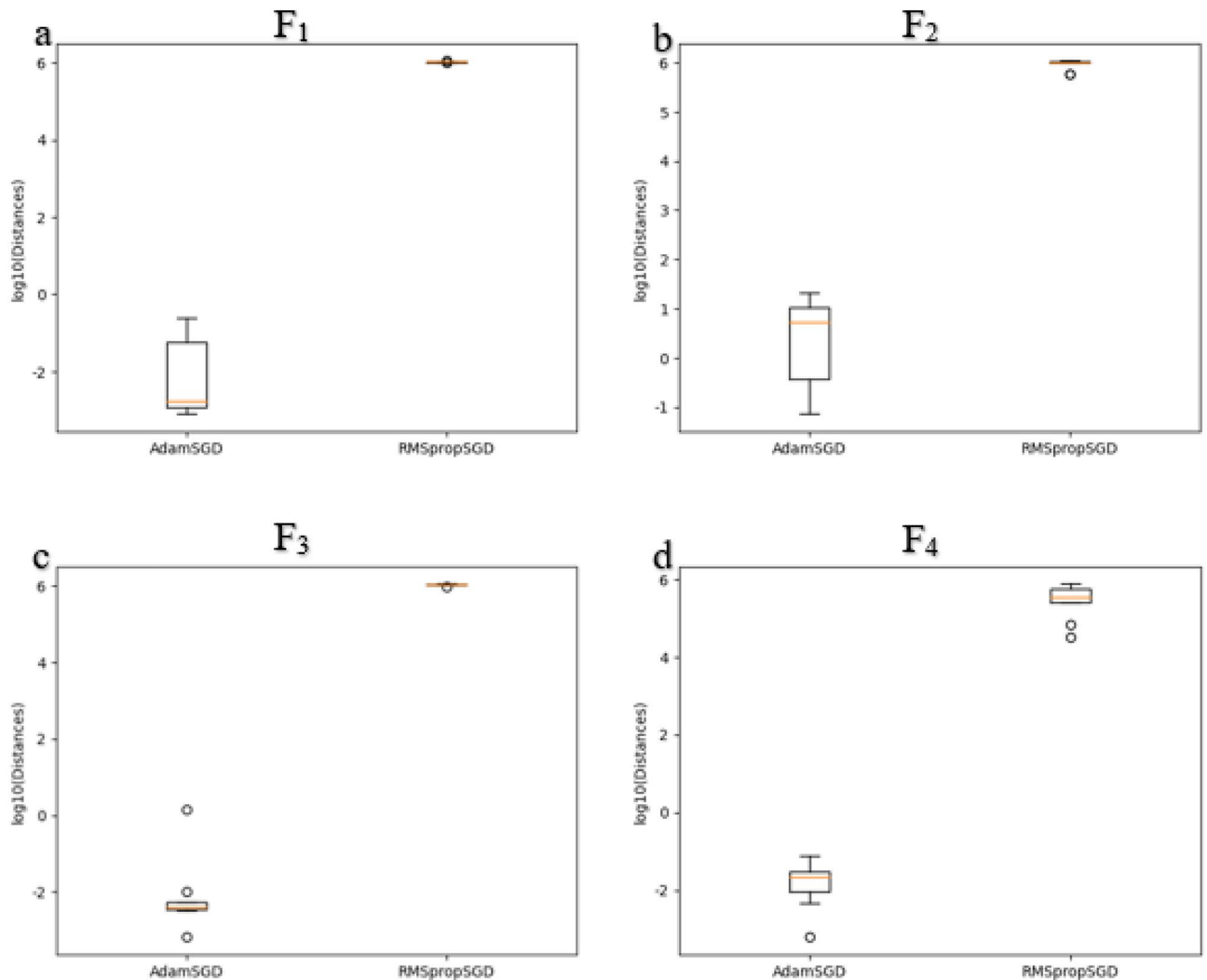


Figure 11. Boxplots of the four studied models, F_1 – F_4 , showing the distance distribution over nine runs for AdamSGD and RMSpropSGD. (a) results for F_1 . (b) results for F_2 . (c) results for F_3 . (d) results for F_4 .

different model weight parameters, we performed a simulated study. In Figs. 7, 8, 9, 10 and 11, we showed that a model induced with the help of SGD optimizer is closer to a model induced with Adam optimizer when compared to a model induced with the help of RMSprop optimizer. It can be evident from visualized results in our study that SGD and Adam had less distributional differences than that of SGD and RMSprop. That resembles the case of having two related datasets for SGD and Adam against unrelated datasets for SGD and RMSprop. As a result, inferior performance results for models utilizing RMSprop are attributed to the high distributional differences in model weight parameters.

It is worth noting that our DTL models keep weights of many layers unchanged. Therefore, when we trained our models, we had fewer number of updated weights compared to updated weights in models trained from scratch. It can be seen from Fig. 6 that our DTL models are fast and can be adopted into mobile applications. It can be noticed from Tables 2 and 3 that leveraging source task knowledge contributed to improved prediction performance when coupled with updated weight parameters in the target task using Adam optimizer. On the other hand, the transferred knowledge from the source task contributed to degraded performance when coupled with updated weight parameters in the target task using RMSprop and SGD optimizers. Also, when we assessed additional DL models such as ConvNeXtLarge and ConvNeXtTiny, the knowledge transfer contributed to maintain the same performance behavior in which leveraging source domain knowledge when coupled with updated weights in the target task remained to be the best (see Supplementary Tables S2 and S3).

Conclusions and future work

In this paper, we present and analyze deep transfer learning (DTL) models for the task of classifying 224 SCGRN images pertaining to healthy controls and T2D patients. First, we utilized seven pre-trained models (including SEResNet152 and SEResNeXT101) already trained on more than million images from the ImageNet dataset. Then, we left weights in the convolutional base (i.e., feature extraction part) unchanged and thereby transferring knowledge from pre-trained models while modifying the densely connected classifier with the use of Adam optimizer to discriminate healthy and T2D SCGRN images. Another presented DTL models work as follows.

We kept weights of bottom layers in the feature extraction part of pre-trained model unchanged while modifying consequent layers including the densely connected classifier with the use of Adam optimizer. Experimental results on the whole 224 SCGRN image dataset using five-fold cross-validation demonstrate the superiority of TFeSEResNeXT101, achieving the highest average BAC of 0.97 and therefore significantly surpassing the performance of the baseline resulted in an average BAC of 0.86. Furthermore, our simulation study showed that the highly accurate performance in our models is attributed to the distributional conformance of weights obtained with the use of Adam optimizer when coupled with weights of pre-trained models.

Future work includes (1) adopting our computational framework to analyze DTL models with different network topologies and thereby identifying the best practice for DTL; (2) incorporating multi-omics datasets with images to improve the prediction performance using DTL models; (3) developing a boosting mechanism to improve the performance of DTL models in different biological problems^{40,41}; (4) incorporating feature representation obtained via our DTL models with machine learning algorithms for the task of inferring SCGRNs; and (5) utilizing our framework to speed up the learning process, e.g., TFeVGG19 was $802.67 \times$ faster than VGG19, trained from scratch.

Data availability

The dataset analyzed during the current study is available in the dataset folder within supplementary material at <https://www.biorxiv.org/content/10.1101/2020.08.30.273839v1.supplementary-material>. The single-cell gene expression data is available in the ArrayExpress repository under accession number E-MTAB-5061 (<https://www.ebi.ac.uk/biostudies/arrayexpress/studies/E-MTAB-5061>).

Received: 12 September 2023; Accepted: 18 February 2024

Published online: 24 February 2024

References

- Hemerich, D. *et al.* Effect of tissue-grouped regulatory variants associated to type 2 diabetes in related secondary outcomes. *Sci. Rep.* **13**(1), 3579 (2023).
- Xie, D. *et al.* Global burden and influencing factors of chronic kidney disease due to type 2 diabetes in adults aged 20–59 years, 1990–2019. *Sci. Rep.* **13**(1), 20234 (2023).
- Parker, E. D. *et al.* Economic costs of diabetes in the US in 2022. *Diabetes Care* **47**(1), 26–43 (2024).
- Mohsen, F. *et al.* A scoping review of artificial intelligence-based methods for diabetes risk prediction. *NPJ Dig. Med.* **6**(1), 197 (2023).
- Su, X. *et al.* Ten metabolites-based algorithm predicts the future development of type 2 diabetes in Chinese. *J. Adv. Res.* <https://doi.org/10.1016/j.jare.2023.11.026> (2023).
- He, Y. *et al.* Comparisons of polyexposure, polygenic, and clinical risk scores in risk prediction of type 2 diabetes. *Diabetes Care* **44**(4), 935–943 (2021).
- Edlitz, Y. & Segal, E. Prediction of type 2 diabetes mellitus onset using logistic regression-based scorecards. *Elife* **11**, e71862 (2022).
- Kokkorakis, M. *et al.* Effective questionnaire-based prediction models for type 2 diabetes across several ethnicities: A model development and validation study. *EClinicalMedicine* **64**, 102235 (2023).
- Pyrros, A. *et al.* Opportunistic detection of type 2 diabetes using deep learning from frontal chest radiographs. *Nat. Commun.* **14**(1), 4039 (2023).
- Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. in *3rd International Conference on Learning Representations* (2015).
- Wachinger, C., Wolf, T. N. & Pölsterl, S. Deep learning for the prediction of type 2 diabetes mellitus from neck-to-knee Dixon MRI in the UK biobank. *Heliyon* **9**(11), e22239 (2023).
- Das, B. A deep learning model for identification of diabetes type 2 based on nucleotide signals. *Neural Comput. Appl.* **34**(15), 12587–12599 (2022).
- He, K. *et al.* Deep residual learning for image recognition. in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2016).
- Naveed, I. *et al.* Artificial intelligence with temporal features outperforms machine learning in predicting diabetes. *PLOS Dig. Health* **2**(10), e0000354 (2023).
- Bengio, Y., Goodfellow, I. & Courville, A. *Deep Learning* Vol. 1 (MIT Press, 2017).
- Wu, D. *et al.* Multi-feature map integrated attention model for early prediction of type 2 diabetes using irregular health examination records. *IEEE J. Biomed. Health Inform.* **65**, 1–10 (2023).
- Shu, H. *et al.* Modeling gene regulatory networks using neural network architectures. *Nat. Comput. Sci.* **1**(7), 491–501 (2021).
- Badia-i-Mompel, P. *et al.* Gene regulatory network inference in the era of single-cell multi-omics. *Nat. Rev. Genet.* **24**, 739–754 (2023).
- Turki, T. & Taguchi, Y. H. Discriminating the single-cell gene regulatory networks of human pancreatic islets: A novel deep learning application. *Comput. Biol. Med.* **132**, 104257 (2021).
- Iacono, G. *et al.* bigSCale: An analytical framework for big-scale single-cell data. *Genome Res.* **28**(6), 878–890 (2018).
- Iacono, G., Massoni-Badosa, R. & Heyn, H. Single-cell transcriptomics unveils gene regulatory network plasticity. *Genome Biol.* **20**(1), 110 (2019).
- Tripathi, S., Dehmer, M. & Emmert-Streib, F. NetBioV: An R package for visualizing large network data in biology and medicine. *Bioinformatics* **30**(19), 2834–2836 (2014).
- Simonyan, K. & Zisserman, A. *Very Deep Convolutional Networks for Large-Scale Image Recognition*, in *3rd International Conference on Learning Representations (ICLR)*. (2015).
- Chollet, F. Xception: Deep learning with depthwise separable convolutions. in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2017).
- Huang, G. *et al.* Densely connected convolutional networks. in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2017).
- Segerstolpe, Å. *et al.* Single-cell transcriptome profiling of human pancreatic islets in health and type 2 diabetes. *Cell Metab.* **24**(4), 593–607 (2016).
- Szegedy, C. *et al.* Rethinking the inception architecture for computer vision. in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2016).
- He, K. *et al.* Identity mappings in deep residual networks. in *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*. (Springer, 2016).

29. Hu, J., Shen, L. & Sun, G. Squeeze-and-excitation networks. in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2018).
30. Bottou, L. Stochastic gradient descent tricks. In *Neural Networks: Tricks of the Trade* 2nd edn 421–436 (Springer, 2012).
31. Chollet, F. *Deep Learning with Python* (Manning Publications Co., 2017).
32. RC Team. R: A language and environment for statistical computing. *J. Stat. Softw.* **25**(1), 1–10 (2008).
33. Franco, V. R. *Optim: General-Purpose Gradient-Based Optimization*. (2021).
34. Hunter, J. D. Matplotlib: A 2D graphics environment. *Comput. Sci. Eng.* **9**(03), 90–95 (2007).
35. Ruder, S. *An Overview of Gradient Descent Optimization Algorithms*. [arXiv:1609.04747](https://arxiv.org/abs/1609.04747) (2016).
36. Currin, C. *et al.* *A Bayesian Approach to the Design and Analysis of Computer Experiments* (Oak Ridge National Lab, 1988).
37. Forrester, A., Sobester, A. & Keane, A. *Engineering Design Via Surrogate Modelling: A Practical Guide* (Wiley, 2008).
38. Gao, H., Pei, J. & Huang, H. Demystifying dropout. in *International Conference on Machine Learning*. (PMLR, 2019).
39. Liu, Z. *et al.* A convnet for the 2020s. in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. (2022).
40. Turki, T. & Wei, Z. Boosting support vector machines for cancer discrimination tasks. *Comput. Biol. Med.* **101**, 236–249 (2018).
41. Turki, T. & Wei, Z. Improved deep convolutional neural networks via boosting for predicting the quality of in vitro bovine embryos. *Electronics* **11**(9), 1363 (2022).

Author contributions

T.T. conceived and designed the study. S.A. performed the deep learning experiments and the visualization of results. T.T. performed the analysis. T.T. and S.A. wrote the manuscript. T.T. supervised the study. All authors have read and agreed to the revised version of the manuscript.

Funding

This study received no funding.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-024-54923-y>.

Correspondence and requests for materials should be addressed to T.T.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2024