

DeepBioGen: Generalizing predictions to unseen sequencing profiles via visual data augmentation

Min Oh¹ and Liqing Zhang^{1,*}

¹Department of Computer Science, Virginia Tech, Blacksburg, VA, USA

*Correspondence: lqzhang@cs.vt.edu

Abstract

Predictive models trained on sequencing profiles often fail to achieve expected performance when externally validated on unseen profiles. While many factors such as batch effects, small data sets, and technical errors contribute to the gap between source and unseen data distributions, it is a challenging problem to generalize the predictive models across studies without any prior knowledge of the unseen data distribution. Here, this study proposes DeepBioGen, a sequencing profile augmentation procedure that characterizes visual patterns of sequencing profiles, generates realistic profiles based on a deep generative model capturing the patterns, and generalizes the subsequent classifiers. DeepBioGen outperforms other methods in terms of enhancing the generalizability of the prediction models on unseen data. The generalized classifiers surpass state-of-the-arts methods, evaluated on RNA sequencing tumor expression profiles for anti-PD1 therapy response prediction and WGS human gut microbiome profiles for type 2 diabetes diagnosis.

Introduction

Predictive models relying on genomic signatures and biomarkers often suffer significantly inferior performance in the independent validation on external data sets in biomedical research such as disease

diagnostics, prognostics, drug discovery, and precision medicine, resulting in a contribution to reproducibility crisis¹⁻⁴. Irreproducible models can lead to not only invalid conclusions misleading subsequent studies but also a substantial waste of time and effort for researchers trying to commercialize the models to benefit patients⁵. A major factor behind these failures is the lack of generalizability across studies, in each of which the number of the heterogeneous data points is insufficient to obtain statistical power to overcome the generalization barrier. In addition to the sample size, usually, there is a significant gap between source data that are used to train classifiers and target data that are used to evaluate the classifiers. One possible cause of the gap is the batch effect such as different sample cohorts, different lab environments, and differences in experimental protocols across studies^{4,6}, which violates the assumption that source and target data are drawn from the same distribution.

In many real-world applications, trained systems fail to produce accurate predictions for unseen data with the shifted distribution. For example, illumination or viewpoint changes in data acquisition for an object detection system and noisier environments for a speech-to-text translation system could easily disrupt the desired outcome. To address this issue, domain adaptation algorithms have been proposed to better align source and target data in a domain-invariant feature space when knowledge of target domains is available during the training phase⁷⁻⁹. However, in practice, it is common that no clue on the target domain is provided. As a more ambitious goal, domain generalization studies focus on training a model generalizing to the unseen domain without any foreknowledge of the unseen domain. Recent studies proposed different ways of domain generalization such as extracting domain-invariant features¹⁰⁻¹², leveraging self-supervised tasks to guide and learn robust representation¹³, simulating domain shift in meta-learning¹⁴, and adding perturbed samples^{15,16}. Although these methods achieved promising performance on benchmark data sets, their requirements, such as having datasets from multiple source domains or sufficient enough for splitting and simulating domain shift, are often not satisfied in biomedical research where only a limited number of heterogeneous data points in a single source domain is available.

Data augmentation techniques in the computer vision field show promising potential in improving classifiers by reducing overfitting to source data¹⁷⁻¹⁹. Especially, recent advances in deep generative models such as generative adversarial networks (GAN)²⁰ allow generating visual contents that are indistinguishable from real ones and also augmenting image data to guide in finding better decision boundaries^{17,19}. More recently, generative models have been utilized to augment medical images, including Magnetic Resonance Images (MRI)²¹, computed tomography (CT)²², and X-ray images²³. However, there has been little effort in transferring the success in computer vision to biomedical sequencing data²⁴. Furthermore, it is unclear whether augmentation of sequencing data could overcome the generalization barrier across different studies.

In this study, DeepBioGen, a data augmentation procedure that establishes visual patterns from sequencing profiles and generates new sequencing profiles capturing the visual patterns based on conditional Wasserstein GAN, is proposed to enhance the generalizability of the prediction models to unseen data. DeepBioGen outperforms other augmentation methods in generalizing classifiers to unseen data. Also, the classifiers generalized by DeepBioGen surpass state-of-the-art classifiers that are designed to work on unseen profiles when tested on two scenarios: devising a prediction model for immune checkpoint blockade (anti-PD1) responsiveness in melanoma patients based on RNA sequencing (RNA-seq) data and building a diagnostic model for type 2 diabetes based on whole-genome metagenomic sequencing data. DeepBioGen source code is free and available at <https://github.com/minoh0201/DeepBioGen>.

Results

Formation and augmentation of visual patterns of sequencing profiles

Sequencing profiles, such as RNA-seq measurements of gene expression levels, consist of numerical values that indicate the activity of thousands of genes in different samples or patients. While many

statistical methods such as multivariate linear regression assume that variables are independent of one another, in reality, genes' activities are highly correlated²⁵. In DeepBioGen, to take into account and visually formalize the interactivity of related genes, similar features in the profiles were clustered together, presenting visible patterns after converting numerical values to colors (Figure 1a). Subsequently, a conditional Wasserstein GAN equipped with convolutional layers to capture the local visual patterns was implemented to augment the sequencing profiles conditioned on class labels. During the augmentation phase, multiple GANs were initialized and trained with different random seeds to promote diversity in the augmented data points (Figure 1b).

To inspect the visual quality of augmented data, two different sequencing profiles were used to train the generative models: one is RNA-seq expression profiles of melanoma patients, and the other is gut microbiome profiles of type 2 diabetes patients. Visual assessment showed that the augmented profiles preserved the boundaries of the clustered features and within-cluster color patterns in the same manner as source data. It is also difficult to distinguish an augmented profile from source data without the original tag (Supplementary Figure S1 and S2).

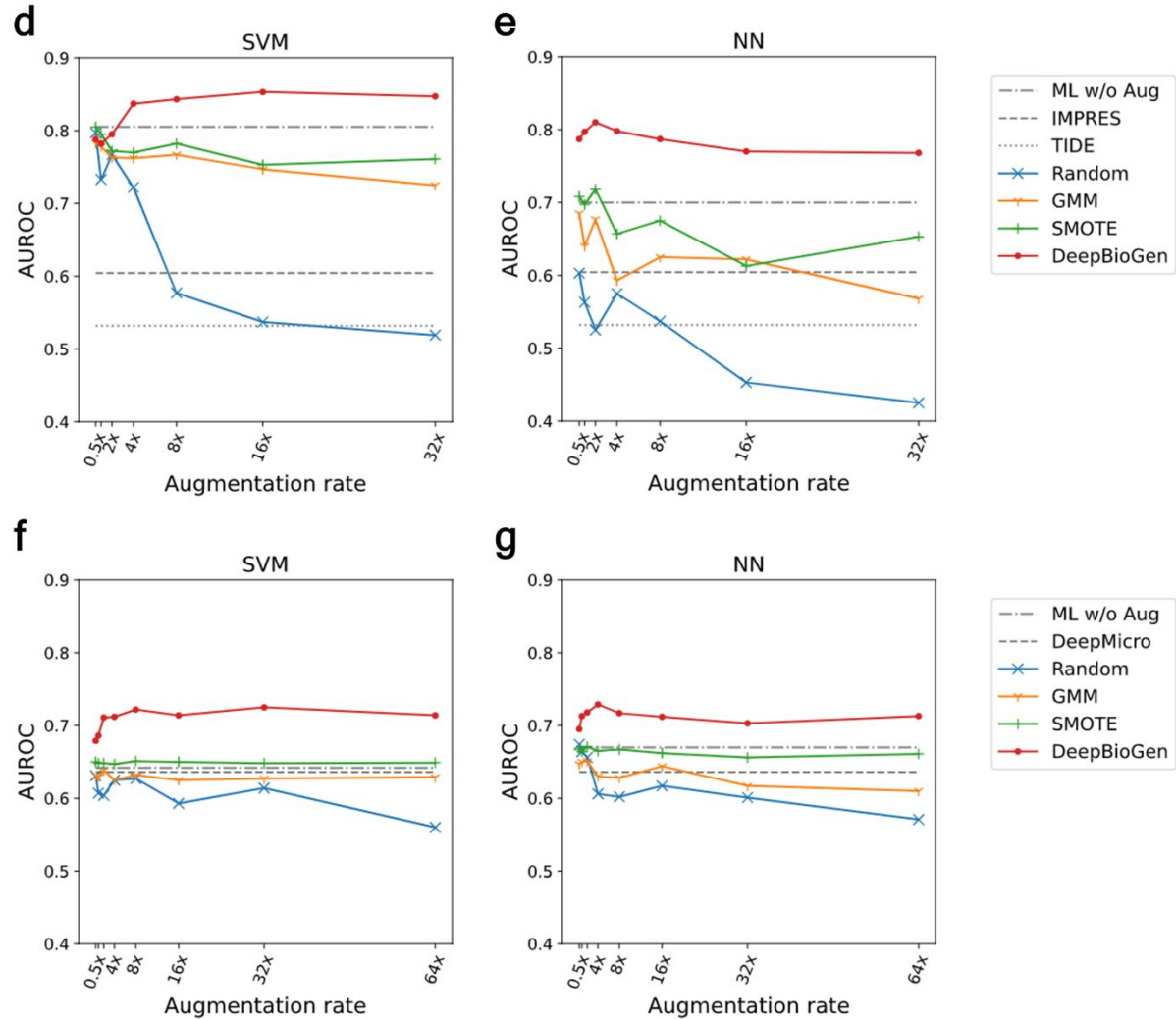
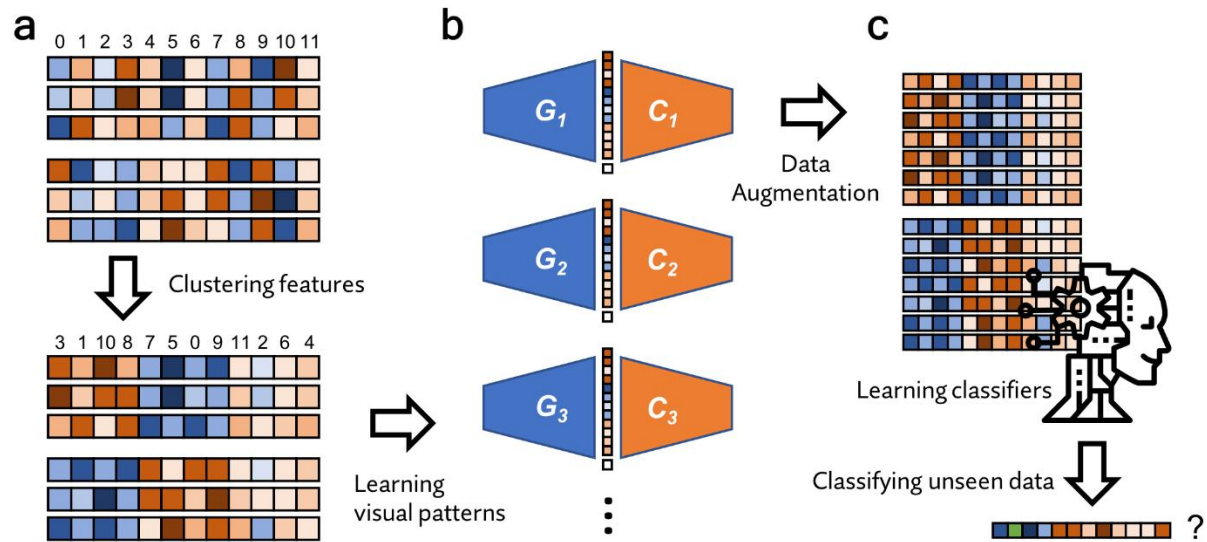


Figure 1. DeepBioGen, a sequencing profile augmentation procedure that generalizes classifiers to enhance prediction performance on unseen data. **a**, Feature-wise clustering of sequencing profiles to form perceptible visual patterns. **b**, Training multiple conditional Wasserstein GANs equipped with up-convolutional and convolutional layers. **c**, Generating augmented data from the multiple generators of GAN models and learning classifiers based on the augmented data along with source data to predict unseen data. **d-e**, Results of anti-PD1 therapy response prediction on unseen data by the state-of-the-art and baseline classifiers (gray) and by classifiers generalized with DeepBioGen (red), SMOTE (green), GMM (yellow), and Random augmentation (blue); Classification algorithms: Support Vector Machine (SVM) and Neural network (NN) which is a multi-layer perceptron; Evaluation metric: Area under the receiver operating characteristics (AUROC). **f-g**, Results of type 2 diabetes prediction on unseen data.

Generalized classification on unseen sequencing profiles

The augmented data derived from the multiple generators of GANs were injected into training data along with the source data. The training data was used to train three machine learning classifiers, support vector machine (SVM), an artificial neural network (NN), and random forest (RF) (Figure 1c). The classifiers were trained to predict non-responders of cancer immunotherapy (anti-PD1) based on RNA-seq gene expression profiles or type 2 diabetes based on human gut microbiome profile.

To validate the generalizability of the classifiers, test (unseen) data were secured from studies that are independent of the source studies. Classification performances on test data were evaluated using an area under the receiver operating characteristics (AUROC) and an area under the precision-recall curve (AUPRC). State-of-the-art predictors, TIDE²⁶ and IMPRES²⁷ for predicting patient response to anti-PD1 therapy, and DeepMicro²⁸ for using deep representations of microbiome data to predict disease states, were compared to DeepBioGen. Besides, widely-used data augmentation techniques, such as Gaussian Mixture Model (GMM)²⁹ and Synthetic Minority Over-sampling Technique (SMOTE)³⁰, were used to

generate augmented data for comparison. The classifiers trained only on source data were used as the baseline comparison.

Remarkably, DeepBioGen-based classifiers surpass not only state-of-the-art classifiers but also classifiers that are trained on augmented data generated by different augmentation methods in both immunotherapy response (Figure 1d-e and Supplementary Figure S3) and diabetes predictions (Figure 1f-g and Supplementary Figure S3). Notably, even though DeepBioGen-based classifiers have no clue of test data, it outperforms Gide et al.'s immune marker classifier (AUROC=0.77) that directly leverages the test data through differential expression analysis³¹. Especially, DeepBioGen provides a stable performance boost to SVM and NN classifiers for both problems as the augmentation rate increases. RF classifiers partially benefit from DeepBioGen, showing generally worse performance than SVM and NN classifiers (Supplementary Figure S4). Consistently, DeepBioGen reduces \mathcal{H} -divergence between the source data and the test data more than other augmentation methods (Table 1).

Table 1. \mathcal{H} -divergence between source and test data

| Data type | DeepBioGen | SMOTE | GMM | Random |
|----------------------------------|--------------|-------|-------|--------|
| RNA-seq tumor expression profile | 0.368 | 0.688 | 0.512 | 0.888 |
| WGS human gut microbiome profile | 0.268 | 0.288 | 0.352 | 0.858 |

Impact of visual clusters and multiple generators

DeepBioGen uses the elbow method³² to estimate the optimal number of visual clusters and GANs. To assess the ability of the approach in inferring the ideal parameters based on source data only, DeepBioGen models with a varying number of visual clusters or GANs were used to generate the augmented data for training classifiers. The classification results of unseen data show that the elbow method elicits an optimal or nearly optimal number of clusters and GANs in both immunotherapy response and diabetes prediction problems (Supplementary Figure S5-S8).

Notably, the number of clusters has more impact on classification performance than the number of GANs, suggesting that how sequencing data are clustered and thus presented visually plays a major role in improving the generalizability of DeepBioGen (Supplementary Figure S5-S8). Results also show that diverse generators of multiple Wasserstein GANs are more effective in diversifying the augmented sequencing data than a single generator, thus leading to better generalizability (Supplementary Table S3).

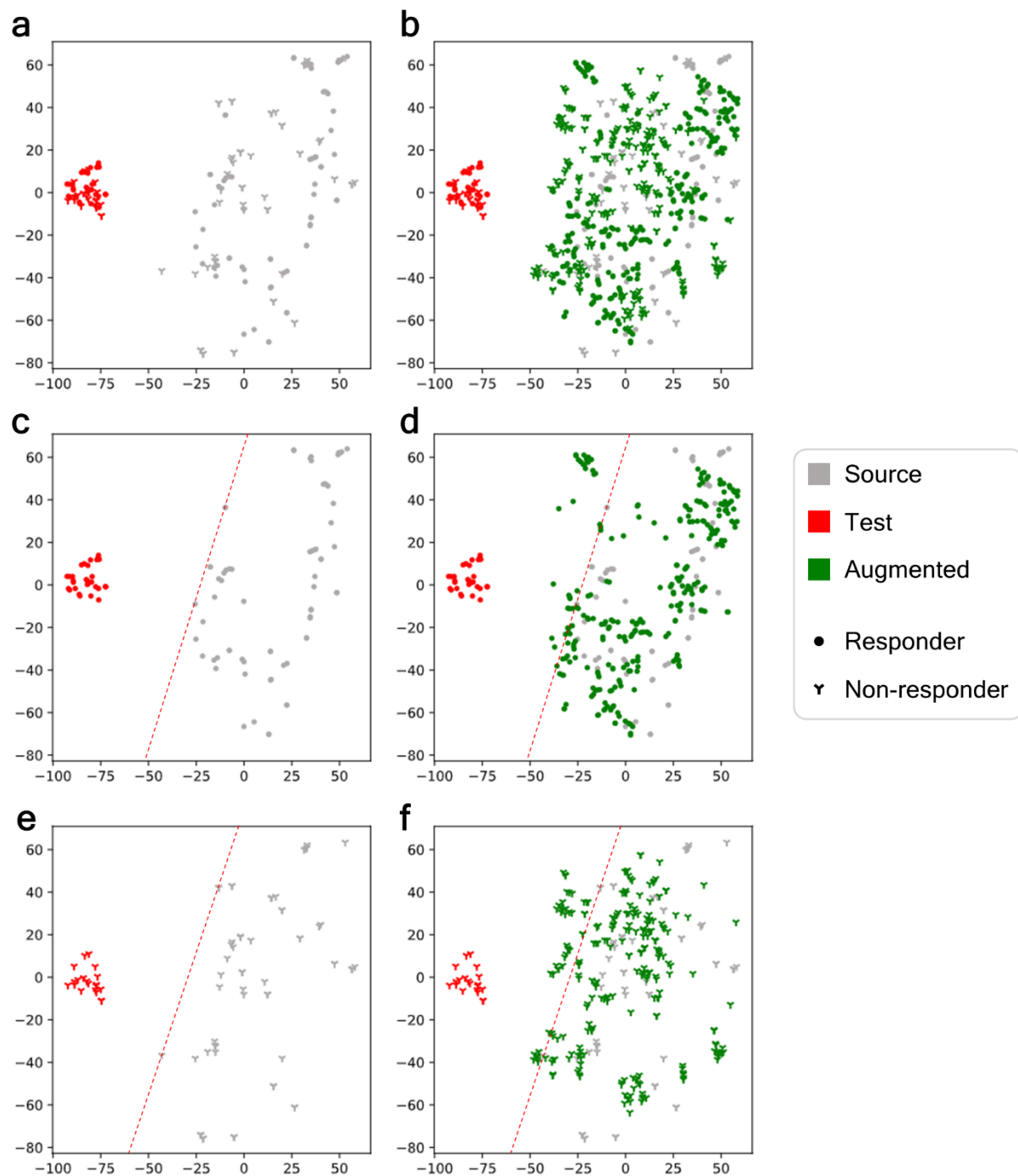
Augmentations beyond the boundary of source data

To visualize how DeepBioGen augmented data to generalize classifiers, the source, augmented and test data were embedded to 2-dimensional space with t-SNE algorithm³³. In melanoma patient profiles, the source and test data are placed distantly, while within-cluster data points with different anti-PD1 responses are located closely in both data clusters (Figure 2a). The data embeddings were plotted separately for two classes, and an empirical outer boundary of the source data based on the outermost data points heading toward the test data was drawn with a red dotted line (Figure 2c and 2e). Interestingly, DeepBioGen generated data points beyond the outer boundaries of the source data cluster (Figure 2d and 2f), whereas other augmentation methods rarely produced data points that cross the boundaries (Supplementary Figure S7-S9).

In microbiome profiles of healthy controls and diabetic patients, the test data cluster resides in the side region of the source data cluster, thus depicting a moderately shifted distribution (Supplementary Figure S12). DeepBioGen produced augmented microbiome profiles across boundaries of the source data cluster. Particularly, the outermost augmented data points beyond the source boundaries are closely placed with test data points that cross the border (Supplementary Figure S12), while other methods rarely generate data points overpassing the boundaries (Supplementary Figure S13-S15).

155

156



157

158 **Figure 2.** t-SNE visualization of augmented tumor expression profiles derived from DeepBioGen along
159 with the source (grey), augmented (green), and test (unseen, red) data of melanoma patients treated with
160 anti-PD1 therapy. **a**, The source and test data. **b**, The source, test, and augmented data. **c**, Responders of

the source and test data; An empirical boundary of responders of source data (red dotted line). **d**, Responders of the source, test, and augmented data. **e**, Non-responders of the source and test data; An empirical boundary of non-responders of source data (red dotted line). **f**, Non-responders of the source, test, and augmented data.

Progression-free survival analysis of predicted anti-PD1 treatment responders

For the predicted responder (PR) and non-responder (PNR) patients to anti-PD1 treatment determined by DeepBioGen-supported SVM classifier, progression-free survival analysis was conducted to estimate the clinical outcome. For comparison, state-of-the-art classifiers based on genomic signatures, IMPRES and TIDE, were evaluated with the same analysis. With the DeepBioGen classifier or IMPRES, the PR group has a significantly longer progression-free survival rate compared to the PNR group (Figure 3a and 3b), whereas the two TIDE predicted groups do not show a significant difference.

Importantly, the median survival time of PRs classified by the DeepBioGen classifier was 755 days (95% CI [335, N/A]), compared to 440 days (95% CI [125, N/A]) for the IMPRES classified PRs . Also, the DeepBioGen classifier tends to be more sensitive in predicting responders than IMPRES, likely posing a lower risk of unnecessary treatment suggestions often accompanied by unnecessary side effects (Figure 3 and Table 2).

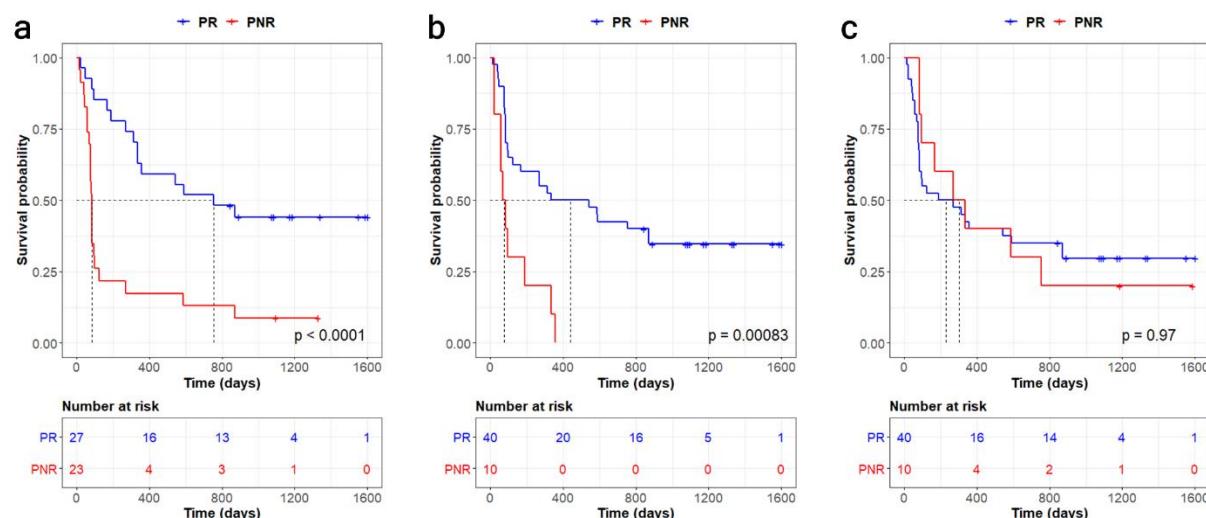


Figure 3. Kaplan-Meier plots of progression-free survival for predicted responder (PR) and non-responder (PNR) patients determined by three classifiers. **a**, generalized SVM classifier with DeepBioGen augmentations. **b**, IMPRES. **c**, TIDE.

Table 2. Summary statistics for progression-free survival analysis

| Classifier | Prediction | N | Median survival time (days) | 95% CI | MR* | HR** | 95% CI | P-value |
|----------------|------------|----|-----------------------------|-----------|------|------|--------------|---------|
| DeepBioGen-SVM | PR | 27 | 755 | [335, NA] | 9.21 | 3.72 | [1.88, 7.36] | < 0.001 |
| | PNR | 23 | 82 | [76, 125] | | | | |
| IMPRES | PR | 40 | 440 | [125, NA] | 5.71 | 3.47 | [1.66, 7.49] | 0.002 |
| | PNR | 10 | 77 | [58, NA] | | | | |
| TIDE | PR | 40 | 231 | [82, 870] | 0.76 | 0.99 | [0.45, 2.17] | > 0.9 |
| | PNR | 10 | 303 | [96, NA] | | | | |

*Median ratio; **Hazard ratio

Discussion

DeepBioGen provides a framework for effective data augmentation in sequencing profiles that can be used to boost the training data and improve the performance of prediction models on unseen data. It adversarially learns multiple generative models that capture visual signals from source data. With multiple generators, DeepBioGen generates realistic augmented data beyond the boundary of the source domain. The augmented data can be used to amplify training data and train classifiers resilient to unknown domain shifts. Consequently, DeepBioGen can improve the transferability and reproducibility of the prediction models without any knowledge of unseen data.

DeepBioGen is unique as it takes input sequencing profiles in machine-understandable visual form, while visualization of sequencing data (e.g. heatmap of differentially expressed genes) has been typically used to present findings in a human-understandable manner. One potential advantage of feeding DeepBioGen with visually recognizable data is that visual patterns difficult to be identified with human eyes may be captured and characterized in embedding space.

Even with a limited amount of source data, DeepBioGen can alleviate batch effects of independent studies without details for batch correction such as sample cohorts, lab environments, and experimental protocol, by reducing the gap between the source and unseen data. Also, DeepBioGen is highly extensible to other biological data whose feature dependency is not negligible.

In the future study, it is envisioned that the process of forming visual patterns from sequencing profiles can be learned with cutting-edge machine learning models toward the better formation of machine-understandable patterns.

Methods

Sequencing profiles and pre-processing

Clinical genomic data containing RNA-seq tumor expression profiles of melanoma patients and their responsiveness to anti-PD1 therapy were secured from three independent studies^{31,34,35} (Supplementary Table S1). Fifty samples in the most recent study³¹ were used as test data and the others were used as source data. RNA-seq read counts were normalized to transcripts per million (TPM) and then log2-transformed. To focus on genes related to primary mechanisms of tumor immune evasion, recently identified T cell signature genes²⁶, such as regulators of T cell dysfunction and suppressors of T cell infiltration into the tumor, were selected out of 18,570 common genes across the studies. In total, 702 genes were considered as features of initial inputs.

Human gut metagenomic sequencing reads of type 2 diabetic patients and healthy controls were acquired from two independent studies: one on the Chinese cohort³⁶ and the other on the European women cohort³⁷ (Supplementary Table S1). Using MetaPhlAn2³⁸, strain-level marker profiles were extracted from the metagenomic samples. In total, the number of common strain-level markers that are considered as initial features was 74,240. The European samples in the more recent study were used as test data and Chinese samples as source data.

Formation of visual patterns from sequencing profiles

Each measurement in source data was standardized by subtracting the mean and dividing by the standard deviation. The same standardization was applied to test data using the mean and standard deviation of source data. To meet the dimensional requirement of the pre-defined input layer, the extremely randomized trees³⁹ feature selection algorithm was applied to the source data to select 256 features. The k-means clustering algorithm was used to cluster features. Based on the elbow point where the within-cluster sum of squared errors (WSS) starts to decrease significantly, the optimal number of clusters was

determined to be 4 for RNA-seq tumor expression profiles and 6 for human gut microbiome profiles (Supplementary Figure S16). The selected features were then sorted and rearranged by cluster labels so that similar features are placed nearby. The features of test data were also rearranged in the same order.

Augmentation of sequencing profiles based on their visual patterns

DeepBioGen captures local visual patterns of sequencing profiles by training conditional Wasserstein GAN, whose generator and critic networks are composed of up-convolutional and convolutional layers, respectively. The generator tries to generate realistic images enough to fool the critic, whereas the critic tries to assign higher values for real images than for generated images. During training, the generator and the critic progressively become better at their jobs by competing against each other. This adversarial training can be conducted by optimizing a minimax objective. Wasserstein distance (or Earth Mover) formulated by Kantorovich-Rubinstein duality is used in the objective term for better reaching Nash equilibrium⁴⁰. Also, the gradient penalty is applied to the objective function to enforce the Lipschitz constraint, alleviating potential instability in the critic⁴¹. Generator function G and critic function C are conditioned on the class label y and the final objective function of conditional Wasserstein GAN is as follows:

$$\min_G \max_C \mathbb{E}_{z \sim p(z)} [C(G(z|y))] - \mathbb{E}_{x \sim P_r} [C(x|y)] - \mathbb{E}_{\hat{x} \sim P_{\hat{x}}} [(\|\nabla_{\hat{x}} C(\hat{x}|y)\|_2 - 1)^2]$$

where z denotes a random noise vector derived from random noise distribution $p(z)$, x a real profile derived from the real data distribution P_r , and $\hat{x} \sim P_{\hat{x}}$ sampling uniformly along straight lines connecting the real data distribution P_r and the output distribution of generator $P_g = G(z|y)$. The gradient penalty term directly constrains the norm of the critic's output concerning its input, enforcing the Lipschitz constraint along the straight lines.

The architecture of neural networks that approximate generator function G and critic function C is illustrated in Supplementary Figure S17. The generator begins with two input layers, one for receiving a

random noise vector and the other for a class label, followed by dense and embedding layers. Embedded random noise vector and label vector are reshaped and concatenated. Subsequently, two up-convolutional blocks, composed of an up-convolutional layer, batch normalization layer, and Leaky ReLU activation layer, perform inverse convolution operations. Lastly, the final up-convolutional layer produces generated sequencing profile. Note that each sequencing profile is considered as a 1x256 pixel image in a single channel. Similarly, the critic has two input layers, one for sequencing profile and the other for a class label, which is embedded, reshaped, and concatenated onto the sequencing profile vector. The two consecutive convolutional blocks, each of which consists of a convolutional layer, Leaky ReLU activation, and dropout layer, are followed by the output layer with a single unit. Across the generator and critic, the alpha value of Leaky ReLU is set as 0.3, and the dropout rate is set at 0.3.

To achieve better generalization, multiple clones of the GAN are trained in the same way except for initial weights in the neural networks. The number of desired GANs is estimated by approximating modes of samples with the elbow method under the assumption that most modes are generated if the number of generators is at least as many as the number of modes in source data (Supplementary Figure S18). Individual generators produce the same number of augmented data points.

Generalized predictions on unseen sequencing profiles

To generalize classifiers predicting clinical outcomes or disease states to unseen data, three classifiers, SVM, NN, and RF, were built on training data composed of source and augmented sequencing profiles. Hyper-parameters of the classifiers were optimized based only on source data with a 5-fold cross-validation scheme. Grid search was applied to explore hyper-parameter space (see details in Supplementary Table S2). With the best hyper-parameters, prediction models were trained on the pooled source and augmented data. The generalizability and performance of the prediction models were

evaluated on the unseen test data using AUROC and AUPRC. The performance evaluation was repeated by gradually changing the augmentation rate indicating how many times the size of augmented data is of the source data.

For comparison, state-of-the-art classifiers designed to work on unseen data, including TIDE²⁶, IMPRES²⁷, and DeepMicro²⁸, were evaluated on test data. TIDE predicts anti-PD1 responsiveness of melanoma patients based on genome-wide expression signatures of T cell dysfunction and exclusion. To satisfy its requirement, the test data without filtering out any genes from the original data was submitted to the TIDE response prediction web service. IMPRES is a predictor of anti-PD1 response in melanoma patients, which is a rule-based classifier manually built based on gene expression relationships between immune checkpoint gene pairs. Its source code was utilized to evaluate the performance of IMPRES on the test data. DeepMicro is a deep representation learning framework for improving predictors based on microbiome profiles. The source data was utilized to learn a low-dimensional representation of the microbiome data, and classifiers were then trained on the representation and evaluated on the test data. Furthermore, as an alternative to DeepBioGen, widely-used data augmentation approaches, including GMM²⁹ and SMOTE³⁰, as well as statistics-based random augmentation were evaluated. An independent GMM model was fitted for each class label, and the optimal number of components in the GMM model was estimated with the Bayesian information criterion (BIC). SMOTE derives the generated samples from linear combinations of nearest neighboring samples. Random augmentation draws data points from the normal distribution whose mean and standard deviation are the same as those of the source data, assigning an arbitrary class label. Also, as a baseline comparison, machine learning classifiers that are trained only on source data (i.e., no augmented data) were evaluated on test data.

To understand the impact of generalization on reducing the discrepancy between the source and test data, a classifier-induced divergence measure, \mathcal{H} -divergence, was determined with various classifiers. For a given set of binary hypotheses $\mathcal{H} \subseteq \{h: X \rightarrow \{0,1\}\}$, \mathcal{H} -divergence is the largest possible difference

between probabilities of being classified as 1 in source and test distributions^{42,43}. More formally, the empirical \mathcal{H} -divergence can be written as:

$$d_{\mathcal{H}}(D_S, D_T) = 2 \sup_{h \in \mathcal{H}} |P_{x \sim D_S}[h(x) = 1] - P_{x \sim D_T}[h(x) = 1]|$$

where D_S and D_T are the source and test data, respectively, and

$$P_{x \sim D}[h(x) = 1] = \frac{|\{x: x \in D, h(x) = 1\}|}{|D|}$$

As a proxy of \mathcal{H} for each augmentation method, all classifiers trained on the augmented training data by varying an augmentation rate and classification algorithms were included in a set of binary hypotheses.

Impact of multiple generators on the diversity of generated sequencing profiles

Wasserstein GAN may suffer less from mode collapse than infant GAN relying on Jensen-Shannon divergence in its loss term⁴⁰. However, a single Wasserstein GAN may not be able to produce all modes of data, and it can be hypothesized that multiple Wasserstein GANs may increase the diversity of augmented sequencing profiles. To evaluate the diversity of the augmented profiles generated with multiple Wasserstein GANs, the adapted inception score is used. Originally, the inception score was introduced to evaluate the quality and diversity of generated images based on the predicted class probability distributions derived from a pre-trained Inception v3 model⁴⁴. More recently, Gurumurthy et al. suggested a modified inception score considering within-class diversity of the generated data⁴⁵, and this scoring method is used in the current evaluation. Also, according to the note that non-ImageNet data generator should not be evaluated by the Inception v3 classifier⁴⁶, it is replaced with the best performing baseline-classifier trained only on source data. Consequently, the adapted inception score ranges from 1 to 2, and the higher the score, the better the diversity and quality of the augmented profiles.

t-SNE visualization of the augmented data

To visualize how augmented data is arranged in a high-dimensional space, the augmented data along with source and test data was embedded into a 2-dimensional space using t-SNE. Also, a class-specific boundary of the source data cluster facing the test data cluster in the embedded space was drawn with one or two straight lines through the outermost data points of the source data cluster.

Progression-free survival analysis

The Kaplan-Meier plots were drawn to conduct progression-free survival analysis for predicted responder and non-responder patients. For each classifier, a receiver operating characteristic (ROC) curve was used to determine the cut-off value of predictions. The closest point from (0, 1) on the ROC curve was chosen, at which the threshold well balancing true positive rate and false-positive rate is identified. The log-rank test was used to validate statistical significance.

References

- 1 Baker, M. 1,500 scientists lift the lid on reproducibility. *Nature* **533** (2016).
- 2 Bernau, C. *et al.* Cross-study validation for the assessment of prediction algorithms. *Bioinformatics* **30**, i105-i112 (2014).
- 3 Castaldi, P. J., Dahabreh, I. J. & Ioannidis, J. P. An empirical assessment of validation practices for molecular classifiers. *Briefings in bioinformatics* **12**, 189-202 (2011).
- 4 Collins, F. S. & Tabak, L. A. Policy: NIH plans to enhance reproducibility. *Nature* **505**, 612-613 (2014).

352 5 Mattsson-Carlgrén, N., Palmqvist, S., Blennow, K. & Hansson, O. Increasing the reproducibility
353 of fluid biomarker studies in neurodegenerative studies. *Nature communications* **11**, 1-11 (2020).

354 6 Leek, J. T. *et al.* Tackling the widespread and critical impact of batch effects in high-throughput
355 data. *Nature Reviews Genetics* **11**, 733-739 (2010).

356 7 Ganin, Y. *et al.* Domain-adversarial training of neural networks. *The Journal of Machine*
357 *Learning Research* **17**, 2096-2030 (2016).

358 8 Hoffman, J. *et al.* in *Proceedings of the International Conference on Machine Learning* 1989-
359 1998 (ICML, 2018).

360 9 Saenko, K., Kulis, B., Fritz, M. & Darrell, T. in *Proceedings of the European Conference on*
361 *Computer Vision*. 213-226 (ECCV, 2010).

362 10 Li, H., Jialin Pan, S., Wang, S. & Kot, A. C. in *Proceedings of the IEEE Conference on Computer*
363 *Vision and Pattern Recognition*. 5400-5409 (CVPR, 2018).

364 11 Li, Y. *et al.* in *Proceedings of the European Conference on Computer Vision*. 624-639 (ECCV,
365 2018).

366 12 Matsuura, T. & Harada, T. in *Proceedings of the Thirty-Fourth AAAI Conference on Artificial*
367 *Intelligence*. 11749-11756 (AAAI, 2020).

368 13 Carlucci, F. M., D'Innocente, A., Bucci, S., Caputo, B. & Tommasi, T. in *Proceedings of the*
369 *IEEE Conference on Computer Vision and Pattern Recognition*. 2229-2238 (CVPR, 2019).

370 14 Li, D., Yang, Y., Song, Y.-Z. & Hospedales, T. in *Proceedings of the Thirty-Second AAAI*
371 *Conference on Artificial Intelligence*. 3490-3497 (AAAI, 2018).

372 15 Shankar, S. *et al.* in *Proceedings of the International Conference on Learning Representations*.
373 (ICLR, 2018).

374 16 Volpi, R. *et al.* in *Proceedings of the 32nd International Conference on Neural Information*
375 *Processing Systems*. 5339-5349.

376 17 Antoniou, A., Storkey, A. & Edwards, H. Data augmentation generative adversarial networks.
377 *arXiv preprint arXiv:1711.04340* (2017).

378 18 Wong, S. C., Gatt, A., Stamatescu, V. & McDonnell, M. D. in *Proceedings of the International*
379 *Conference on Digital Image Computing: techniques and applications.* 1-6 (IEEE DICTA,
380 2016).

381 19 Zhang, X., Wang, Z., Liu, D. & Ling, Q. in *Proceedings of the International Conference on*
382 *Acoustics, Speech and Signal Processing.* 2807-2811 (IEEE ICASSP).

383 20 Goodfellow, I. *et al.* Generative adversarial nets. *Advances in neural information processing*
384 *systems* **27**, 2672-2680 (2014).

385 21 Calimeri, F., Marzullo, A., Stamile, C. & Terracina, G. in *International conference on artificial*
386 *neural networks.* 626-634 (Springer).

387 22 Sandfort, V., Yan, K., Pickhardt, P. J. & Summers, R. M. Data augmentation using generative
388 adversarial networks (CycleGAN) to improve generalizability in CT segmentation tasks.
389 *Scientific reports* **9**, 1-9 (2019).

390 23 Madani, A., Moradi, M., Karargyris, A. & Syeda-Mahmood, T. in *Proceedings of the*
391 *International Society for Optics and Photonics.* 105741M (2018).

392 24 Marouf, M. *et al.* Realistic in silico generation and augmentation of single-cell RNA-seq data
393 using generative adversarial networks. *Nature communications* **11**, 1-12 (2020).

394 25 Emilsson, V. *et al.* Genetics of gene expression and its effect on disease. *Nature* **452**, 423-428
395 (2008).

396 26 Jiang, P. *et al.* Signatures of T cell dysfunction and exclusion predict cancer immunotherapy
397 response. *Nature medicine* **24**, 1550-1558 (2018).

398 27 Auslander, N. *et al.* Robust prediction of response to immune checkpoint blockade therapy in
399 metastatic melanoma. *Nature medicine* **24**, 1545-1549 (2018).

400 28 Oh, M. & Zhang, L. DeepMicro: deep representation learning for disease prediction based on
401 microbiome data. *Scientific reports* **10**, 1-9 (2020).

402 29 Reynolds, D. A., Quatieri, T. F. & Dunn, R. B. Speaker verification using adapted Gaussian
403 mixture models. *Digital signal processing* **10**, 19-41 (2000).

404 30 Chawla, N. V., Bowyer, K. W., Hall, L. O. & Kegelmeyer, W. P. SMOTE: synthetic minority
405 over-sampling technique. *Journal of artificial intelligence research* **16**, 321-357 (2002).

406 31 Gide, T. N. *et al.* Distinct immune cell populations define response to anti-PD-1 monotherapy and
407 anti-PD-1/anti-CTLA-4 combined therapy. *Cancer cell* **35**, 238-255. e236 (2019).

408 32 Thorndike, R. L. Who belongs in the family? *Psychometrika* **18**, 267-276 (1953).

409 33 Maaten, L. v. d. & Hinton, G. Visualizing data using t-SNE. *Journal of machine learning*
410 *research* **9**, 2579-2605 (2008).

411 34 Hugo, W. *et al.* Genomic and transcriptomic features of response to anti-PD-1 therapy in
412 metastatic melanoma. *Cell* **165**, 35-44 (2016).

413 35 Riaz, N. *et al.* Tumor and microenvironment evolution during immunotherapy with nivolumab.
414 *Cell* **171**, 934-949. e916 (2017).

415 36 Qin, J. *et al.* A metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature*
416 **490**, 55-60 (2012).

417 37 Karlsson, F. H. *et al.* Gut metagenome in European women with normal, impaired and diabetic
418 glucose control. *Nature* **498**, 99-103 (2013).

419 38 Truong, D. T. *et al.* MetaPhlAn2 for enhanced metagenomic taxonomic profiling. *Nature*
420 *methods* **12**, 902-903 (2015).

421 39 Geurts, P., Ernst, D. & Wehenkel, L. Extremely randomized trees. *Machine learning* **63**, 3-42
422 (2006).

423 40 Arjovsky, M., Chintala, S. & Bottou, L. in *International conference on machine learning*. 214-
424 223 (PMLR).

425 41 Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V. & Courville, A. Improved training of
426 wasserstein gans. *arXiv preprint arXiv:1704.00028* (2017).

427 42 Ben-David, S., Blitzer, J., Crammer, K. & Pereira, F. Analysis of representations for domain
428 adaptation. *Advances in neural information processing systems* **19**, 137 (2007).

429 43 Kifer, D., Ben-David, S. & Gehrke, J. in *VLDB*. 180-191 (Toronto, Canada).

430 44 Salimans, T. *et al.* in *Proceedings of the 30th International Conference on Neural Information*
431 *Processing Systems*. 2234-2242.
432 45 Gurumurthy, S., Kiran Sarvadevabhatla, R. & Venkatesh Babu, R. in *Proceedings of the IEEE*
433 *conference on computer vision and pattern recognition*. 166-174.
434 46 Barratt, S. & Sharma, R. A note on the inception score. *arXiv preprint arXiv:1801.01973* (2018).
435