

TarDis: Achieving robust and structured disentanglement of multiple covariates

Highlights

- *TarDis* disentangles multiple covariates in single-cell data into isolated latents
- Handles both categorical and continuous covariates with interpretable latent spaces
- Generates ordered, interpretable latent axes for continuous covariates
- Superior disentanglement enables robust data integration and improved OOD predictions

Authors

Kemal Inecik, Aleyna Kara,
Antony Rose, Muzlifah Haniffa,
Fabian J. Theis

Correspondence

kemal.inecik@helmholtz-munich.de (K.I.),
fabian.theis@
helmholtz-munich.de (F.J.T.)

In brief

Inecik et al. introduce *TarDis*, a deep generative model that disentangles multiple covariates in single-cell genomics into structured latent spaces. By jointly modeling categorical and continuous factors with ordered, interpretable representations, *TarDis* enhances data integration, supports counterfactual and hypothesis-driven analyses, and enables the exploration of complex cellular dynamics across diverse conditions.

Methods

TarDis: Achieving robust and structured disentanglement of multiple covariates

Kemal Inecik,^{1,2,8,*} Aleyna Kara,^{1,3} Antony Rose,^{4,5} Muzlifah Haniffa,^{4,5,6} and Fabian J. Theis^{1,7,*}

¹Institute of Computational Biology, Helmholtz Center Munich, Neuherberg 85764, Germany

²School of Life Sciences, Technical University of Munich, Freising 85354, Germany

³Department of Computer Science, Technical University of Munich, Garching 85748, Germany

⁴Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SA, UK

⁵Biosciences Institute, Newcastle University, Newcastle upon Tyne NE2 4HH, UK

⁶Department of Dermatology, Newcastle Hospitals NHS Foundation Trust, Newcastle upon Tyne NE1 4LP, UK

⁷Department of Mathematics, Technical University of Munich, Garching 85748, Germany

⁸Lead contact

*Correspondence: kemal.inecik@helmholtz-munich.de (K.I.), fabian.theis@helmholtz-munich.de (F.J.T.)

<https://doi.org/10.1016/j.cels.2026.101573>

SUMMARY

Addressing challenges in domain invariance within single-cell genomics necessitates innovative strategies for managing the heterogeneity of multi-source datasets while maintaining the integrity of biological signals. We introduce targeted disentanglement (*TarDis*), an end-to-end deep generative model designed to disentangle intricate covariate structures across diverse biological datasets, distinguishing technical artifacts from true biological variations. By employing tailored covariate-specific loss components and a self-supervised approach, *TarDis* effectively generates multiple latent-space representations that capture each continuous and categorical target covariate separately, along with unexplained variation. Our extensive evaluations demonstrate that *TarDis* outperforms existing methods in data integration, covariate disentanglement, and robust out-of-distribution predictions. The model's capacity to produce interpretable and structured latent spaces, including its introduction of ordered latent representations for continuous covariates, markedly enhances its utility in hypothesis-driven research. Consequently, *TarDis* offers a promising analytical platform for advancing scientific discovery, providing insights into cellular dynamics, and enabling targeted therapeutic interventions.

INTRODUCTION

Domain invariance tackles the challenge of learning from datasets that, while representing the same physical phenomena, originate from disparate sources such as different users, acquisition devices, or locations.¹ As the data source often lacks direct relevance to the task, the objective is to develop a model that maintains performance robustness by being invariant to these domain variations. This invariance not only enhances model reliability across shifts, whether subpopulational² or distributional,³ but also is an end in itself where the source is obscured to comply with data protection requirements.⁴ Such shifts, frequently observed in practical machine learning scenarios, necessitate that models be resilient to variations in multi-domain datasets by learning to minimize the disparity in data distributions within the representation space, ideally achieving a low metric distance between them. This concept is closely aligned with distributionally robust optimization strategies, promoting the development of universally applicable machine learning models that withstand out-of-distribution (OOD) variations.^{5–8}

The identification of spurious correlations within these multi-domain datasets can provide critical insights for certain down-

stream applications, enriching the interpretive scope beyond mere domain invariance. Moreover, models leveraging data representations or predictors derived from true correlations, including domain-specific attributes or nuisance factors, more effectively discern causal relationships, thereby enhancing their generalization capabilities.^{9,10} This recognition has spurred interest in disentangled representation learning, aiming to segregate and independently model spurious and invariant characteristics within the data.^{10–12} Developing invariant representation learning models is a complex multi-objective optimization problem, frequently necessitating linear constraints on the data representations and classifiers^{9,12,13} or the incorporation of conditional priors within the variational autoencoder (VAE) framework.^{10,14}

Existing invariant representation learning methods often fail for continuous domain problems, an area that is largely underexplored yet critically important.^{15–17} Examples include patient monitoring systems, where physiological spurious data vary daily and across activities¹⁸; finance, where models predicting stock prices or market trends must generalize across varying economic conditions and times¹⁹; and climate modeling, where models use invariant learning to forecast weather or long-term climate changes across diverse locations and time periods.²⁰ Existing methods are generally designed for discrete categorical

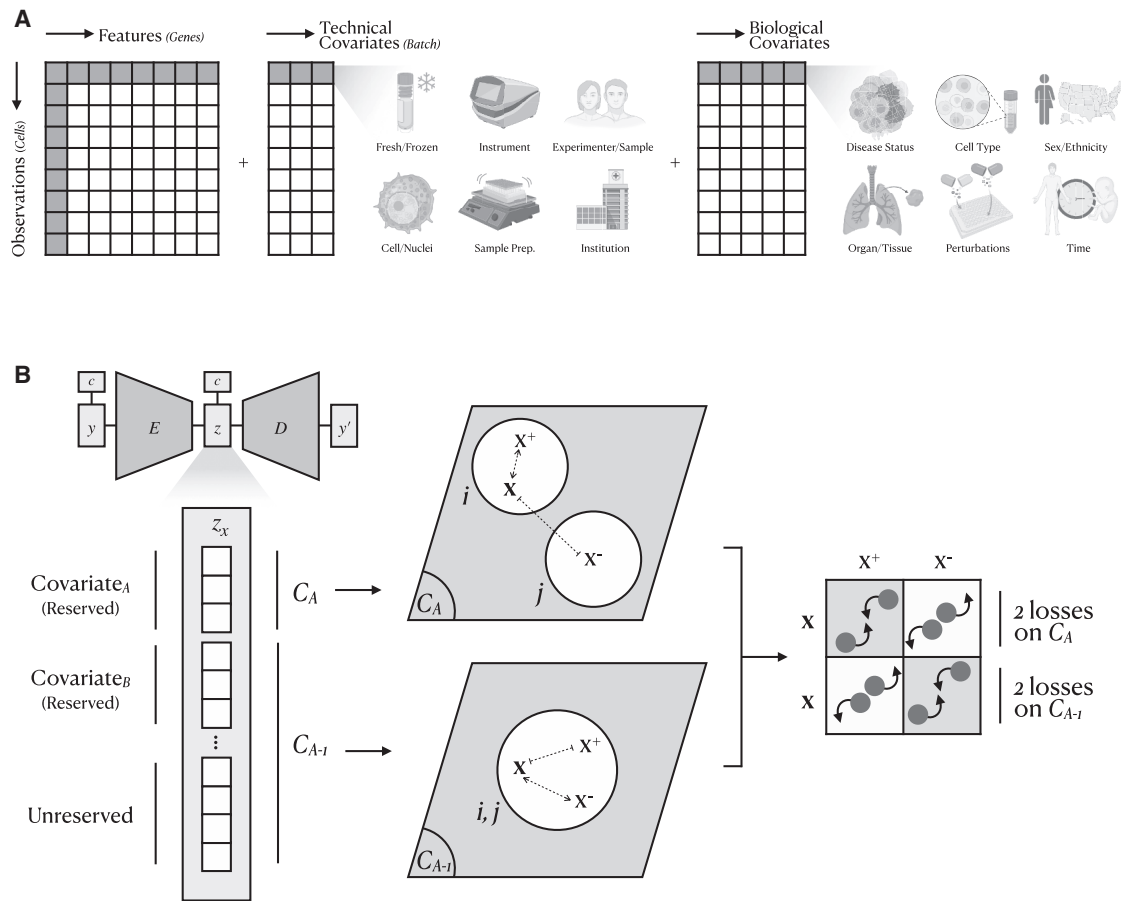


Figure 1. Overview of single-cell covariates and the TarDis framework

(A) A typical single-cell dataset can contain thousands or even millions of cells (observations), each with measured gene expression profiles (features). Alongside these high-dimensional measurements, each cell is characterized by a host of technical covariates (e.g., instrumentation or sample preparation) and biological covariates (e.g., disease status, cell type, sex, ethnicity, organ, perturbations, and temporal factors). Collectively, these covariates influence gene expression, making their accurate representation and disentanglement crucial for extracting meaningful biological insights from complex single-cell datasets.

(B) A conceptual illustration of our proposed method. The figure depicts data containing multiple covariates, denoted as A , B , and so on. Covariate A has two categories, i and j , and the current training iteration uses a data point x associated with category i . For example, covariate A may represent the disease status of cells, with i corresponding to diseased cells and j to healthy cells. Here, x^+ is a randomly chosen data point in category i , while x^- is a randomly chosen data point in categories other than i . The C_A (short for covariate_A, a reserved latent subset of A) is designated to capture covariate A -specific information, whereas its complement C_{A-1} (short for covariate_{A-1}, the remaining part of the latent space) encodes all remaining variation, including other covariates and unexplained factors. The objective of the covariate-specific loss component of A , L_{covA} , is to remove A -specific information from C_{A-1} and to confine it only in C_A . L_{covA} is composed of four loss terms: two operate on C_A to pull data points of the same category closer and push different categories apart, and two operate on C_{A-1} , which acts opposite of the losses on C_A . By constructing and applying such sets of losses for each covariate, TarDis produces distinct latent subsets dedicated to individual covariates, alongside a separate space for residual variation.

domains and struggle with the continuous nature of many real-world tasks. This leads to challenges such as sparse data in each domain, making it hard to accurately estimate invariant correlations, and segmentation of continuous data into discrete blocks, which can misrepresent true data distributions. Addressing these issues is crucial for advancing model robustness and ensuring applicability in dynamic environments.¹⁵

In the context of domain invariance, multi-domain and multi-condition single-cell genomics datasets present a critical testbed where the integration of data representations confronts complex challenges in biological and pharmaceutical research.²¹ Single-cell genomics offers a granular view of individual cells' genetic diversity, highlighting the variability among

cells and being essential for understanding cellular and molecular processes.^{22–24} However, the data often come from a range of labs and varied experimental setups, incorporating batch effects and technical artifacts that can mask true biological signals (Figure 1A).^{25,26} These challenges are compounded when data include cells affected by chemical or genetic perturbations, sourced from diseased states, or differing in their origin, such as specific organs, organisms, developmental stages, ethnicity, age, sex, and other factors that further contribute to variability.^{27–31} Effective data integration is vital for separating technical artifacts from relevant biological signals, facilitating a robust comparison of biological landscapes across various domains and enhancing our understanding of

the underlying cellular dynamics, with broad implications for advancing disease research and therapeutic development.^{32,33}

Hence, it becomes essential to disentangle invariant and spurious correlations for single-cell data integration, where spurious correlations often obscure biological signals. The disentanglement of these elements not only enhances data integration by clarifying underlying biological processes but also bolsters OOD prediction capabilities.^{10,34} Furthermore, there is a compelling need for researchers to explore the potential effects of one covariate on another, whether categorical or continuous, by manipulating such disentangled latent representations. For instance, adjusting the continuous “drug dose” representation while holding the representations of “disease state,” “patient,” and continuous “age” constant could reveal the dose-dependent effects on gene expression independent of the disease’s progression or patient characteristics. Such analyses would deepen our understanding of the interactions between various factors at the cellular level, thereby unlocking new avenues for complex, hypothesis-driven research with single-cell genomics data.

To address the complexities inherent in multi-domain and multi-condition datasets, we introduce *TarDis*, an end-to-end deep generative model specifically designed for the targeted disentanglement of multiple covariates, such as those encountered in extensive single-cell genomics data.¹ *TarDis* employs covariate-specific loss functions through a self-supervision strategy, enabling the learning of disentangled representations that achieve accurate reconstructions and effectively preserve essential biological variations across diverse datasets. It eschews additional architectural complexities, enabling straightforward application to large datasets. *TarDis* ensures the independence of invariant signals from noise, enhancing interpretability that is crucial for extracting biological insights obscured by spurious data correlations. *TarDis* handles both categorical and, notably, continuous variables, demonstrating its adaptability to diverse data characteristics and allowing for a granular understanding and representation of underlying data dynamics within a coherent and interpretable latent space. This capability is instrumental for exploring complex biological phenomena and conducting hypothesis-driven research. Empirical benchmarking across multiple datasets highlights *TarDis*’s superior performance in covariate disentanglement, data integration, and OOD predictions, outperforming existing models (Box 1).^{2,3}

RESULTS

TarDis achieves robust disentanglement of covariates into isolated latent spaces

We assessed the *TarDis* model’s ability to disentangle covariates using the *Afriat* single-cell genomics dataset, which includes three distinct covariates: age, zone status, and time (see the “dataset insights” section). Experiments were conducted with two methodologies: disentangling all covariates simultaneously, *TarDis*_{multiple},⁴ and disentangling each covariate individually, followed by concatenating the reserved latent spaces, *TarDis*_{single}. The disentanglement performance was benchmarked using the maximum mutual information gap (maxMIG), as detailed in the “evaluation metrics” section, demonstrating that both configurations of *TarDis* surpassed existing models^{35–40} and achieved

nearly 0.9 maxMIG scores on validation sets (Figure 2A). These results underscore the efficacy of *TarDis* in handling multiple covariates simultaneously without compromising disentanglement quality. Further analysis using the mutual information (MI) metric reveals minimal differences in the preservation of information within the unreserved and reserved latent spaces between the two training strategies, indicating the model’s effective scalability for disentanglement tasks (Figure 2B). For all subsequent experiments detailed in this paper, we have exclusively employed the multiple-covariate disentanglement approach.

An ablation study was performed to evaluate the model’s robustness against feature reduction, where varying percentages of input features were systematically removed. Results in Figure 2C show that *TarDis* maintained high maxMIG and R² reconstruction scores, above 0.65 and 0.94, respectively, affirming its resilience to input variability. Additionally, modifying the auxiliary loss weight, λ_C , systematically influenced the clustering quality and disentanglement accuracy, as indicated by the increased maxMIG score and mean centroid distance with higher λ_C values (Figures 2D and S2). Moreover, the silhouette scores, calculated on the unreserved latent space \mathbf{z}_{n0} using cell-type annotations as the labels, provided empirical evidence that effective disentanglement correlates with enhanced biological signal representation, as further investigated in the “results” section. Overall, these results not only validate the robustness of *TarDis* in disentangling complex covariate structures but also highlight its utility in preserving essential biological variations, which are pivotal for advancing single-cell genomic data analysis.

TarDis achieves superior performance in single-cell genomics data integration

To probe the efficacy of invariant representation learning, we turned our attention to the *Suo* dataset, a massive single-cell genomics dataset capturing human embryonic development. This dataset includes about 850k cells from various organs and time points, using multiple methods, instruments, samples, and platforms, as well as a wide range of cell types (see the “dataset insights” section). Its complexity makes it an ideal testbed for evaluating model performance in integrating intricate datasets. We assessed the data integration quality using the scIB package metrics,⁴¹ which are recognized benchmarks in the single-cell genomics community for evaluating the balance between biological signal preservation and batch effect mitigation (see the “evaluation metrics” section). This balance is crucial, as inadequate correction can lead to data clustering by batch, obscuring true biological variance, while overcorrection may suppress biological signals, reducing the biological relevance of the outcomes.

In our experimental setup, we tested two configurations of the *TarDis* model. *TarDis*-1 focuses on covariates typically considered as batch keys in single-cell data integration tasks, such as library platform, donor, sample status, and instruments. *TarDis*-2 extends this disentanglement to additional covariates, including sex, age, and, notably, organ. The comparative results, detailed in Table 1, show that *TarDis*, particularly *TarDis*-2, outperforms state-of-the-art models⁵ and maintains an optimal balance between biological conservation and batch correction. By effectively disentangling various spurious correlations from invariant biological signals,

Box 1. Progress and potential

Modern single-cell genomics provides an unprecedented view into cellular heterogeneity, yet the very richness that propels new discoveries also complicates downstream analysis. Gene expression patterns emerge from overlapping biological processes (e.g., differentiation programs and disease progression) and extrinsic factors (e.g., laboratory protocols and technical artifacts). *Disentanglement*, in this context, aims to parse these intertwined influences into interpretable latent representations, a crucial step for elucidating how complex covariates shape cellular states. While methods that correct for batch effects have become standard, these strategies often fall short in achieving the deeper objective of capturing subtle, high-dimensional biological dynamics. In single-cell experiments, cells navigate intricate developmental trajectories, respond nonlinearly to environmental or pharmaceutical perturbations, and exhibit myriad context-specific behaviors. Without disentanglement, these diverse signals frequently remain intermingled, limiting biological interpretability and hindering hypothesis-driven research.

Disentangling biological covariates is particularly vital for addressing nuanced questions in single-cell research. For example, in a disease model involving multiple genetic variants and variable drug dosing, researchers may wish to examine the effect of each variant independently or investigate how dosage influences a specific mutant background. Similarly, in developmental biology, uncovering how cells evolve across a continuum of pseudotime (e.g., from pluripotent to fully differentiated states) is critical for identifying the genes that orchestrate fate decisions while isolating the influence of developmental time from tissue-specific contexts, along with other confounding factors such as culture conditions, sample preparation, or donor genetic characteristics. Alternatively, disentangling lineage commitment signals from spatial patterning cues enables the identification of master regulators driving fate decisions. Moreover, by explicitly isolating and representing each covariate as an independent latent dimension, one can systematically navigate and interrogate a rich *multidimensional covariate space*. This approach extends beyond merely observing biological states—it enables exploration of novel or unmeasured cellular conditions through latent-space manipulations. For instance, disentangled latent spaces could allow researchers to computationally predict cellular responses at drug dosages or developmental stages that were never experimentally observed, substantially broadening the scope and predictive power of experimental datasets. Such analyses yield testable hypotheses for unexplored biological phenomena and enable informed planning of subsequent experimental validations.

The challenge of covariate disentanglement stems fundamentally from the complexity of modeling joint distributions of gene expression conditioned simultaneously on multiple covariates, both categorical (e.g., tissue type and disease condition) and continuous (e.g., pseudotime and dosage). This is inherently an underdetermined problem because single-cell measurements represent only sparse snapshots within a vast combinatorial space of covariate conditions. Conventional modeling approaches often conflate correlated covariates, collapsing biological variability into ambiguous latent factors, and typically fail to explicitly create separate latent representations for disentangled covariates. Moreover, continuous covariates introduce an additional layer of complexity, yet discretizing them artificially imposes arbitrary boundaries, obscuring subtle transitions and hindering accurate capture of biological gradients. Therefore, preserving the continuous nature of such covariates in disentangled representations is critical, as it maintains their intrinsic ordering and enables researchers to discern nuanced biological shifts—such as identifying thresholds in dose-response relationships or characterizing gradual developmental transitions—in a naturally interpretable manner. The key idea in this paper is to devise a tailored deep generative model for systematically separating both categorical and continuous covariates into independent latent dimensions, while still ensuring coherent integration of the underlying gene expression data. By explicitly *targeting* these covariates and preserving continuous variables as smooth, ordered latent axes, our approach clarifies complex interactions and uncovers nuanced patterns that remain concealed under standard analyses. The resulting disentangled representations can then support robust OOD generalizations, refined differential analyses, and more principled hypotheses about how diverse factors interact to drive cellular variation.

TarDis has demonstrated its robust capability to manage the complexities inherent in vast and heterogeneous datasets.

TarDis generates ordered latent representation for continuous covariates

In addressing the challenge of learning the representation of disentangled *continuous* covariates, *TarDis* provides a solution that captures data variations without reducing them to mere categorical approximations. Continuous covariates such as age or treatment dosage are critical for understanding gradients in biological processes, cellular behavior, and disease progression. To manage the subtleties associated with these variables, *TarDis* employs a distance-based loss function for each auxiliary loss component. The model employs negative pair losses weighted by the distance between the values of the continuous covariates, omitting positive pair losses due to the continuous

nature of the covariate, which results in generating an ordered and interpretable latent space (Figure 3).

In our studies, we focused on two primary continuous covariates, age and drug dosage, which present distinct challenges due to their variability and substantial impact on cellular phenotypes. We employed two datasets to evaluate the effectiveness of *TarDis* in producing ordered latent representations of these covariates. The first dataset, named *Sciplex* (see the “[dataset insights](#)” section), involves drug perturbation experiments and helps in analyzing the structured response of cells to varying drug dosages. The second dataset, referred to as *Braun*, comprises 1.6 million cells from human embryonic brain development, providing a complex scenario for assessing the impact of time as a continuous variable. Through *TarDis*, we managed to produce ordered latent representations of these covariates within isolated latent subsets while concurrently disentangling

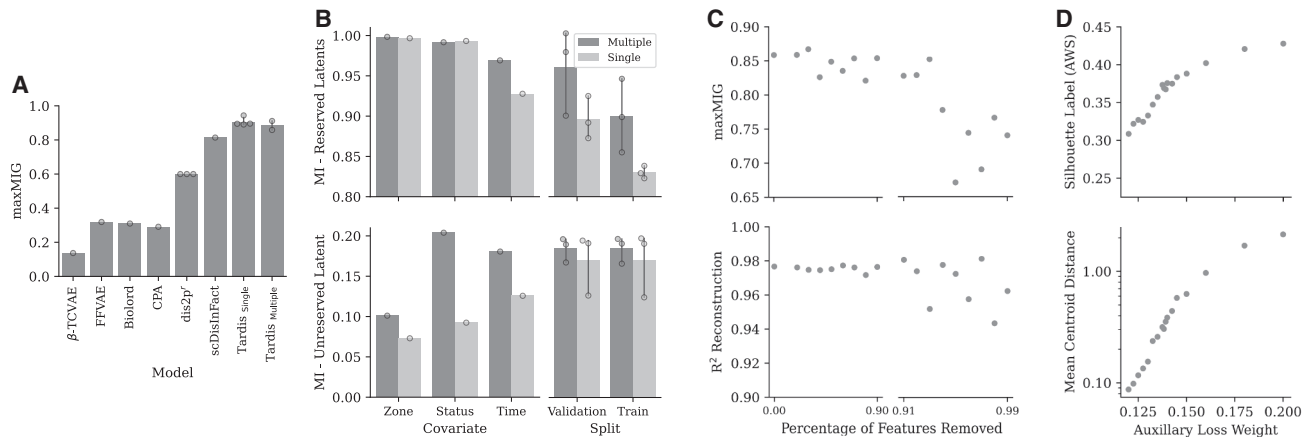


Figure 2. TarDis variants outperform baselines and effectively isolate covariate information within latent subspaces

(A) Comparison of disentanglement performance using maximum mutual information gap (maxMIG) showing that both *TarDis* variants outperform existing models, where r stands for the reported score by the authors. Experiments in (A) and (B) assessed *TarDis*_{multiple}, disentangling all covariates simultaneously, and *TarDis*_{single}, disentangling individually and concatenating latent spaces, with *TarDis*_{single} performing marginally superior.

(B) MI in the reserved, \mathbf{z}_{nr} , and unreserved, \mathbf{z}_{no} , latent spaces for *TarDis* under multiple and single covariate training conditions (as light and dark gray bars) across various covariates and data splits. Overlaid dots in (A) and (B) show individual data points (n visible). Bars show the mean, and whiskers show the standard deviation.

(C) Relationship between the percentage of input features removed and the corresponding maxMIG and R^2 reconstruction scores, indicating robustness to feature removal.

(D) Impact of auxiliary loss weight (λ_C) on mean centroid distance in reserved latents, \mathbf{z}_{nr} , and average silhouette width (ASW) scores at the unreserved latent, \mathbf{z}_{no} . Each dot is one ablation (n visible) in (C) or one λ_C setting in (D). The *Afriat* dataset ($n = 19,053$ cells) is used for all panels, and all axes are linear except (D) bottom, which is log.

other variables such as the type of library platform, donor characteristics, sample status, instrumentation used, and tissue types (Figures 3 and 4).

This representation has enabled previously unfeasible hypothesis-driven biological analyses. For example, *TarDis* allows for the exploration of organ-specific developmental

gene expression patterns for specific cell types, an analysis that previously was not optimal with non-batch-corrected input spaces. Unlike existing models such as single-cell variational inference (scVI) and single-cell annotation using variational inference (scANVI), which address batch effects but often fail to retain essential biological information like age or organ

Table 1. Benchmarking data integration performance by scIB package metrics, organized into biological signal conservation and batch correction categories (refer to the “evaluation metrics” section)

Metric group	Metric	PCA	Harmony	scVI	scANVI	inVAE	<i>TarDis</i> -1	<i>TarDis</i> -2
Bio conservation	isolated labels	0.610	0.563	0.638	0.774	0.798	0.662	0.767
	K-means normalized mutual information (NMI)	0.691	0.620	0.649	0.792	0.651	0.634	0.713
	K-means adjusted Rand index (ARI)	0.226	0.182	0.209	0.360	0.191	0.185	0.228
	silhouette label (AWS)	0.504	0.482	0.496	0.576	0.508	0.497	0.508
	cell-type local inverse Simpson’s index (LISI) (cLISI)	0.999	0.997	0.999	1.000	0.999	0.998	0.999
Batch correction	silhouette batch	0.851	0.862	0.867	0.861	0.840	0.903	0.896
	integration LISI (iLISI)	0.057	0.100	0.098	0.093	0.040	0.094	0.080
	kBET per label	0.309	0.475	0.487	0.526	0.194	0.448	0.430
	graph connectivity	0.793	0.671	0.866	0.912	0.836	0.866	0.879
	PCR comparison	0.000	0.350	0.699	0.222	0.000	0.931	0.850
Aggregate score	bio conservation	0.606	0.569	0.598	0.701	0.629	0.595	0.643
	batch correction	0.402	0.492	0.603	0.523	0.382	0.648	0.627
	total	0.524	0.538	0.600	0.629	0.530	0.616	0.637

Data are from the *Suo* dataset ($n = 841,922$ cells). Quantification employed a comprehensive set of metrics, with aggregate scores derived according to scIB standards. Cell-type annotations are incorporated in the metrics where labels are necessary. Available technical covariates—such as library platform, donor, sample status, and instrument—are concatenated and used as batch keys for model training, except for principal-component analysis (PCA). For this experiment, two configurations of the *TarDis* model were tested: *TarDis*-1 targets the batch keys defined in this experiment, while *TarDis*-2 expands disentanglement to also include sex, age, and organ. For both *TarDis* variants, the unreserved latent subsets are used for benchmarking with scIB metrics.

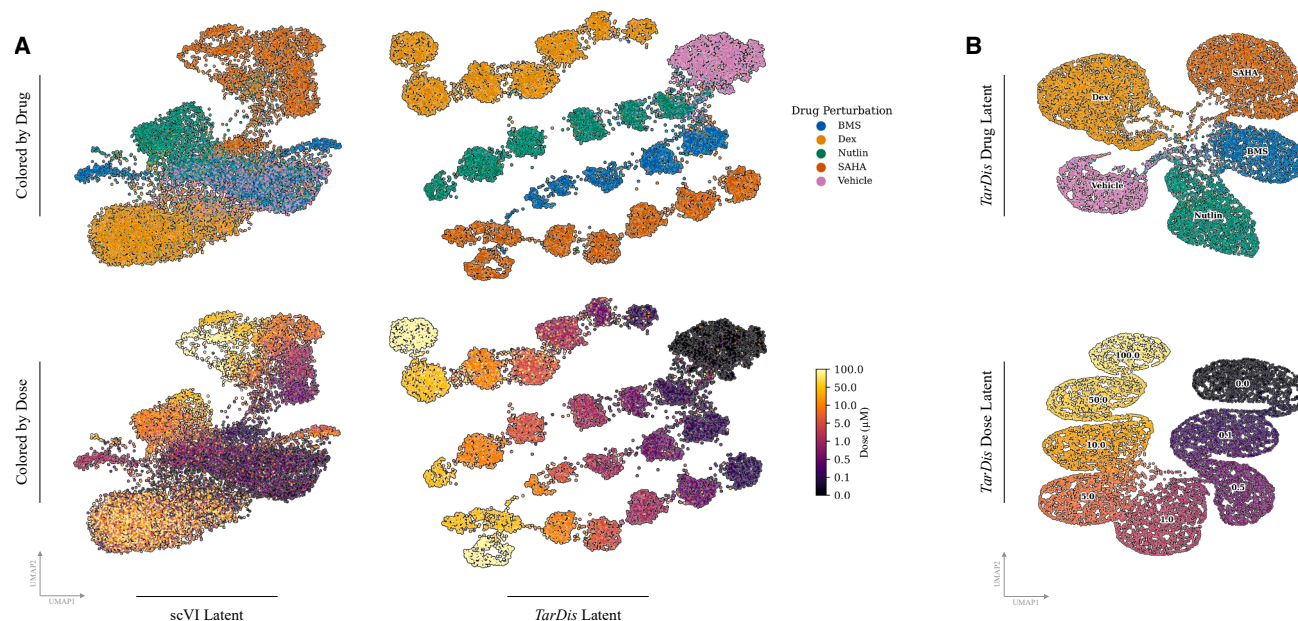


Figure 3. UMAP visualization of *TarDis* latent-space representations from the *Sciplex* dataset (GEO: GSE139944; $n = 14,811$ cells)

(A) Comparing the performance of scVI (trained using default parameters from the `scvi_tools` package,⁴² with batch keys included) and *TarDis* models in capturing drug responses and dosage effects. The upper row displays clusters differentiated by drug types (categorical palette; no numeric scale), while the bottom row illustrates the ordered representation of dosage (μM , arbitrary scale), showcasing the ability of *TarDis* to structurally organize cellular responses across different drug concentrations.

(B) *TarDis* model training generates three distinct latent spaces: unreserved, dose, and drug. Displayed uniform manifold approximation and projections (UMAPs⁴³) are the dose and drug latent subspaces, demonstrating structured separation and ordered representation. Points are cells, descriptive visualization, no hypothesis testing.

specifics—either being overly corrected by batch keys or inadequately accounted for^{26,45}—*TarDis* allows researchers to isolate cells from two different organs using the organ-specific latent subset and, for a given cell type, compare expression patterns across developmental stages in a massive multi-organ developmental single-cell dataset. This analysis benefits from a batch-corrected latent space, thanks to a set of other latent subsets that disentangle batch effects.⁶ In Figure 4, bottom right, *TarDis* enabled to identify genes including *EGR2-3-4*, *KLF2-4*, *RTL1*, *SPRY4-AS1*, and *FOSB*, which decrease in expression through the embryonic development of human forebrain neurons within the *Braun* dataset, which were shown to be associated with brain development, aging, and diseases including Down syndrome and bipolar disorder.^{46–51} In a parallel experiment using the *Sciplex* perturbation dataset, *TarDis* effectively disentangled the influences of drug type and dosage (Figures 3 and 4). Using the data points corresponding to the *Nutlin* cluster in drug latency, we analyzed how gene expression responds to increasing doses. As shown in Figure 4, upper right, this approach allowed us to pinpoint the expression patterns of genes such as *TP53I3*, *CDKN1A*, *GDF15*, *MDM2*, *FDXR*, and *NUPR1*, which are known to be responsive to escalating doses of *Nutlin*.^{7,8,52,53}

***TarDis* predicts counterfactual gene expressions accurately under OOD conditions**

The capacity of predictive models to generate accurate gene expressions under OOD conditions is pivotal for extrapolating

research findings to new or novel environments. In evaluating this capacity, *TarDis* was systematically tested using two distinct datasets to gauge its effectiveness in predicting counterfactual gene expressions. Using the *Afriat* dataset, previously introduced, multiple models were trained, each excluding a different combination of three covariates to create respective OOD sets. Additionally, the *Miller* dataset, which comprises samples from human developmental embryonic lungs, was utilized to disentangle the effects of age and donor covariates (see the “dataset insights” section). Similar to the *Afriat* dataset, combinations of two covariates were systematically omitted during training to simulate various OOD conditions.

TarDis demonstrated superior performance in predicting gene expressions under OOD conditions, outperforming compositional perturbation autoencoder (CPA),⁵ another model that concurrently disentangles multiple covariates,³⁸ in both the *Afriat* and *Miller* datasets. In the *Afriat* dataset, *TarDis* achieved higher R^2 reconstruction scores, showcasing its strong capability for accurate reconstruction under varied and unseen conditions. In the *Miller* dataset, the challenge intensified with the evaluation focusing on differentially expressed genes (DEGs) (see the “evaluation metrics” section). *TarDis* excelled, achieving significantly better OOD predictions for DEGs compared with CPA. Here, our inclusion of both R^2 and R^2 -DEG provides a dual-resolution view, where broad data distribution reconstructions and global patterns across the entire gene set are highly accurate. Simultaneously, the subset of

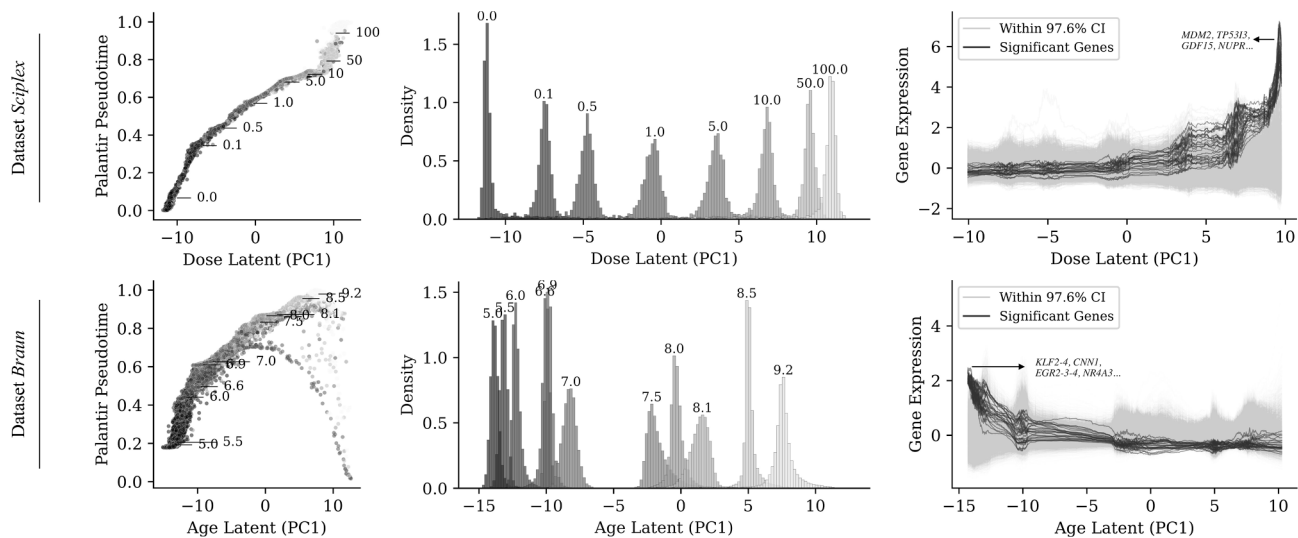


Figure 4. Ordered latent spaces for dose (μM) and age (pcw) in the *Sciplex* ($n = 14,811$ cells) and *Braun* datasets ($n = 1,661,498$ cells), respectively

(Left) Principal component 1 (PC1) of the continuous covariate latent space plotted against Palantir pseudotime,⁴⁴ which uses a k-nearest neighbor graph to infer cell pseudotime trajectories. (Middle) Density distribution of the continuous covariate in the respective latent subset, illustrating ordered peaks corresponding to varying levels of the covariate. (Right) Differential gene expression profiles plotted against the continuous covariate latent space, identifying genes that show variation in expression levels associated with changes in the covariate, indicative of underlying cellular processes. Gene expression patterns are highlighted with (upper right) increasing doses of *Nutlin* and (bottom right) through human embryonic developmental stages of *forebrain neuron*.^{6,7} Latent PC1 values are arbitrary units (dimensionless), Palantir pseudotime is unitless on a linear scale, and gene expression curves represent Z scores on a linear scale. For both the Braun dataset (ages, post-conception week [pcw]) and the *Nutlin* dose (μM), values are color-coded from dark gray to light gray with increasing magnitude, and no color bar is shown.

genes undergoing biologically meaningful expression shifts remains well captured, ensuring that critical biological signals crucial for downstream analyses are preserved. These results, shown in Figure 5, affirm the utility of *TarDis* not only in disentangling complex covariate interactions within datasets but also in its capability to generalize across novel, unseen domains, which is key to advancing the precision and reliability of predictive models in single-cell genomics.

TarDis produces interpretable latent representations of disentangled covariates

In exploring the capabilities of *TarDis* to yield interpretable latent representations, we utilized the *Norman* dataset, a comprehensive collection comprising 108,000 cells subjected to single or combinatorial gene perturbations (see the “dataset insights” section). This dataset is particularly challenging due to the diversity and complexity of its perturbations, with a total of 284 distinct perturbation conditions included in this analysis. In this experiment, the inference model in *TarDis* relied solely on input features without the introduction of covariate information, \mathbf{s}_n . This approach ensured that the learning process was purely driven by the data’s inherent structure rather than external annotations. Our results indicate that *TarDis* effectively disentangles these perturbations, with each perturbation distinctly isolated in the latent space. Notably, perturbations that share a common cellular program, as identified in the original publication of the dataset,⁵⁴ were found to cluster closely. The results support *TarDis*’ ability to capture interpretable and biologically meaningful patterns, as

the clustering is not random qualitatively but reflects the underlying biological relationships (Figure 6A).

A particularly rigorous test of the model’s interpretability involved the relabeling of certain perturbations in the dataset. Specifically, the labels were altered to appear as two distinct entities: “ $X + 0$ ” and “ $0 + X$,” despite originating from the same perturbation. This was designed to test whether *TarDis* could recognize and reconcile these as identical despite their nominal differences. The results were in line with our expectations: *TarDis* successfully overlapped these perturbations in the latent space, affirming its capability to generate biologically coherent and interpretable latent representations, even under challenging conditions (Figure 6B). This analysis not only confirms the robustness of *TarDis*’s disentanglement capabilities but also highlights its potential in generating actionable insights from complex genomic data, where interpretability is crucial for meaningful biological inference.

DISCUSSION

In this study, we presented *TarDis*, a deep generative model designed for the *TarDis* of covariates in complex multi-domain and multi-condition datasets, particularly focusing on the challenges presented by single-cell genomics data. Our approach leverages a series of covariate-specific loss functions to facilitate robust disentanglement and invariant representation of both continuous and categorical variables, thus enhancing data integration capabilities and enabling more insightful biological analyses. Through rigorous benchmarking against existing models and diverse

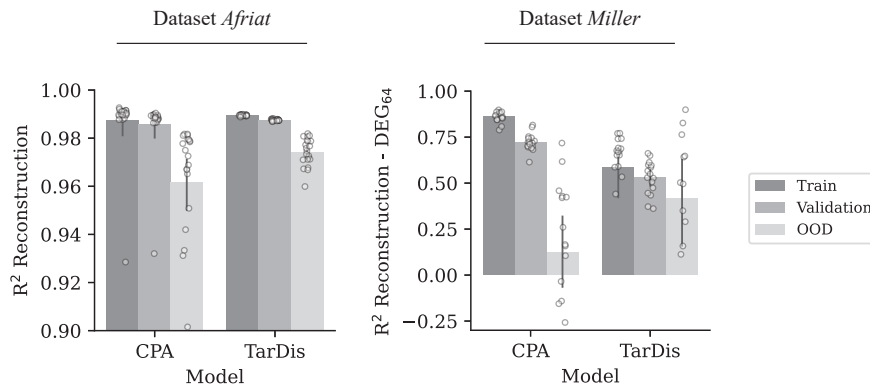


Figure 5. Performance comparison of TarDis and CPA in predicting counterfactual gene expressions under OOD conditions using the Afriat ($n = 19,053$ cells) and Miller datasets ($n = 7,405$ cells)

R^2 scores for reconstructed gene expressions and differentially expressed genes (DEGs) across varying unseen covariate combinations highlight TarDis's superior predictive capabilities. Overlaid dots show different OOD scenarios ($n = 20$ for condition). Bars show the mean, and the whiskers show standard deviation.

datasets, TarDis has demonstrated superior performance not only in its capacity to disentangle complex covariate structures but also in maintaining essential biological signals crucial for accurate data interpretation and analysis and generating robust predictions under OOD conditions. Moreover, TarDis's ability to generate ordered latent representations of continuous covariates enhances differential analyses across varying conditions. The model performs robustly in generating interpretable and biologically meaningful latent representations, which could empower researchers to conduct advanced hypothesis-driven research, potentially unveiling novel insights and therapeutic targets.

Including all covariates may appear simpler *a priori*, but it often induces conflation of correlated factors, leakage into the unreserved latent space, and diminished performance and interpretability (Figure 2; “limitations” under “STAR Methods” section). Consequently, *hypothesis-driven selection* of a subset of covariates—guided by the overarching research question—remains the most direct and scientifically productive strategy for employing TarDis. By aligning the targeted covariates with biological or clinical priorities, researchers can extract clearer insights, avoid overcorrection, and preserve critical aspects of the data that remain unreserved for other analyses.

TarDis establishes a robust approach for exploring complex biological questions, offering researchers comprehensive clarity in dissecting the nuanced interactions between diverse covariates. This capability is instrumental in advancing personalized medicine, supporting the development of customized therapeutic strategies grounded in a profound understanding of individual responses to different treatments. Considering the expansion of TarDis applications beyond genomics, for instance, into neuromarketing using electroencephalogram (EEG) event-related potential (ERP) data, it becomes crucial to acknowledge that modifications to the model may be necessary to accommodate different types of data. We are actively investigating these potential applications, aiming to extend the reach and impact of TarDis across various scientific and applied fields.⁹

RESOURCE AVAILABILITY

Lead contact

Further information and requests for resources should be directed to and will be fulfilled by the lead contact, Kemal Inecik (kemal.inecik@helmholtz-munich.de).

Materials availability

This study did not generate any new materials.

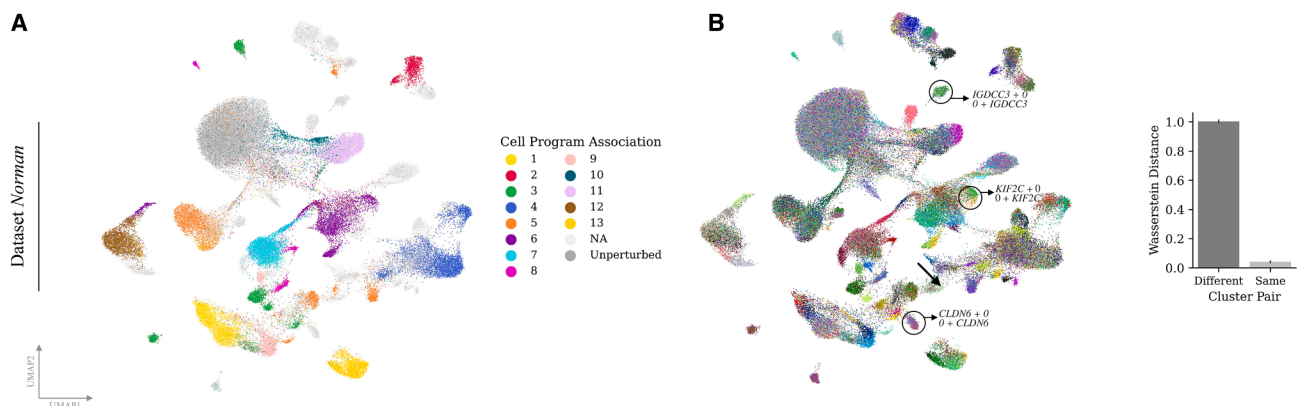


Figure 6. UMAP visualization of the TarDis perturbation latent space derived from the Norman dataset (GEO: GSE133344; $n = 108,497$ cells)
(A) Clusters corresponding to sets of perturbations associated with similar cell programs as identified in the original publication,⁵⁴ demonstrating the model's ability to capture underlying biological patterns. The categorical palette contains 15 categories and represents distinct groups without a numeric scale.
(B) UMAP visualization of TarDis latent space, colored by 284 perturbation identities using a categorical palette without a numeric scale. Representative clusters are highlighted, illustrating the model's capability to align identical perturbations accurately despite nominal labeling differences, thus confirming label reconciliation. Wasserstein distances are computed to quantitatively confirm the close, often overlapping, clustering of identical perturbations.⁵⁵

Data and code availability

- This manuscript utilizes pre-existing datasets that are publicly accessible, and the specific accession numbers pertinent to these datasets are provided in the [key resources table](#). Detailed descriptions of the dataset contents, along with alternative sources for data acquisition, are presented in the “[dataset insights](#)” section. Comprehensive documentation of the preprocessing procedures applied to the majority of these datasets is available in the GitHub repository of *TarDis*.
- All original code has been deposited at the GitHub repository (theislabs/tardis) and is publicly available as of the date of publication. DOIs are listed in the [key resources table](#).
- Any additional information required to reanalyze the data reported in this paper is available from the [lead contact](#) upon request.

ACKNOWLEDGMENTS

We extend our sincere gratitude to the reviewers of the RECOMB 2025 conference version as well as the additional reviewers from *Cell Systems* for their insightful and constructive feedback, which substantially improved the manuscript's quality. We also gratefully acknowledge the members of the Theis lab for their diligent proofreading and valuable suggestions, which significantly enhanced the clarity and coherence of the text. Open access funding provided by “Helmholtz Zentrum München—Deutsches Forschungszentrum für Gesundheit und Umwelt (GmbH).”

AUTHOR CONTRIBUTIONS

K.I. conceptualized and conceived the ideas of the study, formulated the underlying hypotheses, developed the methodologies, implemented the software and the analyses, applied data preprocessing, conducted all experimental procedures, and authored the manuscript. These were carried out under the supervision of F.J.T., who also provided guidance throughout the research process and secured funding for the project. A.R. was instrumental in the collection of metadata for human developmental datasets. A.K. contributed to the manuscript by assisting with the writing process. All authors have thoroughly read and given their approval for the final version of the manuscript.

DECLARATION OF INTERESTS

F.J.T. consults for Immunai Inc., CytoReason Ltd, Cellarity, BioTuring Inc., and Genbio.AI Inc. and has an ownership interest in Dermagnostix GmbH and Cellarity.

DECLARATION OF GENERATIVE AI AND AI-ASSISTED TECHNOLOGIES IN THE WRITING PROCESS

During the preparation of this work, the authors used ChatGPT in order to improve language and readability. After using this tool/service, the authors reviewed and edited the content as needed and take full responsibility for the content of the publication.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- [KEY RESOURCES TABLE](#)
- [METHOD DETAILS](#)
 - Overview
 - VAE skeleton
 - Auxiliary loss
 - Additional details
- [QUANTIFICATION AND STATISTICAL ANALYSIS](#)
 - Loss functions
 - Evaluation metrics

SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.cels.2026.101573>.

Received: April 1, 2025

Revised: September 16, 2025

Accepted: March 11, 2026

REFERENCES

1. Andéol, L., Kawakami, Y., Wada, Y., Kanamori, T., Müller, K.-R., and Montavon, G. (2023). Learning domain invariant representations by joint Wasserstein distance minimization. *Neural Netw.* 167, 233–243. <https://doi.org/10.1016/j.neunet.2023.07.028>.
2. Koh, P.W., Sagawa, S., Marklund, H., Xie, S.M., Zhang, M., Balsubramani, A., Hu, W., Yasunaga, M., Phillips, R.L., Gao, I., et al. (2021). WILDS: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning (PMLR)*, pp. 5637–5664. <https://doi.org/10.48550/arXiv.2012.07421>.
3. Goel, K., Gu, A., Li, Y., and Ré, C. (2020). Model patching: Closing the subgroup performance gap with data augmentation. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2008.06775>.
4. Hajihassnai, O., Ardakanian, O., and Khzaei, H. (2021). ObscureNet: Learning attribute-invariant latent representation for anonymizing sensor data. In *Proceedings of the International Conference on Internet-of-Things Design and Implementation (ACM)*, pp. 40–52. <https://doi.org/10.1145/3450268.3453534>.
5. Lu, C., Wu, Y., Hernández-Lobato, J.M., and Schölkopf, B. (2022). Invariant Causal Representation Learning for Out-of-Distribution Generalization. In *International Conference on Learning Representations* <https://openreview.net/forum?id=e4EXDWXnSn>.
6. Yin, M., Wang, Y., and Blei, D.M. (2021). Optimization-based causal estimation from heterogenous environments. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2109.11990>.
7. Guo, S., Tóth, V., Schölkopf, B., and Huszar, F. (2024). Causal de Finetti: On the identification of invariant causal structure in exchangeable data. *Advances in Neural Information Processing Systems* 36. <https://doi.org/10.48550/arXiv.2203.15756>.
8. Sturma, N., Squires, C., Drton, M., and Uhler, C. (2024). Unpaired multi-domain causal representation learning. *Advances in Neural Information Processing Systems* 36. <https://doi.org/10.48550/arXiv.2302.00993>.
9. Ahuja, K., Wang, J., Dhurandhar, A., Shanmugam, K., and Varshney, K.R. (2020). Empirical or invariant risk minimization: a sample complexity perspective. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2010.16412>.
10. Aliee, H., Kapl, F., Hediye-Zadeh, S., and Theis, F.J. (2023). Conditionally Invariant Representation Learning for Disentangling Cellular Heterogeneity. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2307.00558>.
11. Arjovsky, M., Bottou, L., Gulrajani, I., and Lopez-Paz, D. (2019). Invariant risk minimization. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1907.02893>.
12. Kong, L., Xie, S., Yao, W., Zheng, Y., Chen, G., Stojanov, P., Akinwande, V., and Zhang, K. (2022). Partial disentanglement for domain adaptation. In *Proceedings of the 39th International Conference on Machine Learning*, 162, K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, eds. (*Proceedings of Machine Learning Research*. PMLR), pp. 11455–11472. <https://proceedings.mlr.press/v162/kong22a.html>.
13. Ahuja, K., Caballero, E., Zhang, D., Gagnon-Audet, J.-C., Bengio, Y., Mitiagkas, I., and Rish, I. (2021). Invariance principle meets information bottleneck for out-of-distribution generalization. *Advances in Neural Information Processing Systems* 34, 3438–3450. <https://doi.org/10.48550/arXiv.2106.06607>.
14. Lu, C., Wu, Y., Hernández-Lobato, J.M., and Schölkopf, B. (2021). Nonlinear invariant risk minimization: A causal approach. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2102.12353>.

15. Yong, L., Zhou, F., Tan, L., Ma, L., Liu, J., He, Y., Yuan, Y., Liu, Y., Zhang, J.Y., Yang, Y., et al. (2024). Continuous Invariance Learning. The Twelfth International Conference on Learning Representations. <https://doi.org/10.48550/arXiv.2310.05348>.
16. Azzam, M., Gnanha, A.T., Wong, H.-S., and Wu, S. (2021). Adversarially Constrained Interpolation for Unsupervised Domain Adaptation. In 25th International Conference on Pattern Recognition (ICPR) (IEEE), pp. 2375–2381. <https://doi.org/10.1109/ICPR48806.2021.9412471>.
17. Zhang, Y., and Davison, B.D. (2021). Adversarial continuous learning in unsupervised domain adaptation. In Pattern Recognition. ICPR International Workshops and Challenges: Virtual Event, January 10–15, 2021, Proceedings, Part II, A. Del Bimbo, R. Cucchiara, S. Sclaroff, G.M. Farinella, T. Mei, M. Bertini, H.J. Escalante, and R. Vezzani, eds. (Springer), pp. 672–687. https://doi.org/10.1007/978-3-030-68790-8_52.
18. Cao, Z., Yu, H., Yang, H., and Sano, A. (2023). Pirl: participant-invariant representation learning for healthcare using maximum mean discrepancy and triplet loss. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2302.09126>.
19. Huang, H., Chen, M., and Qiao, X. (2024). Generative Learning for Financial Time Series with Irregular and Scale-Invariant Patterns. In The Twelfth International Conference on Learning Representations <https://openreview.net/forum?id=CdjnzWsQax>.
20. Beucler, T., Gentine, P., Yuval, J., Gupta, A., Peng, L., Lin, J., Yu, S., Rasp, S., Ahmed, F., O’Gorman, P.A., et al. (2024). Climate-invariant machine learning. *Sci. Adv.* 10, ead7250. <https://doi.org/10.1126/sciadv.adj7250>.
21. Heumos, L., Schaar, A.C., Lance, C., Litinetskaya, A., Drost, F., Zappia, L., Lücken, M.D., Strobl, D.C., Henao, J., Curion, F., et al. (2023). Best practices for single-cell analysis across modalities. *Nat. Rev. Genet.* 24, 550–572. <https://doi.org/10.1038/s41576-023-00586-w>.
22. Inecik, K., and Theis, F.J. (2023). scARE: Attribution regularization for single cell representation learning. Preprint at bioRxiv. <https://doi.org/10.1101/2023.07.05.547784>.
23. Perez, R.K., Gordon, M.G., Subramaniam, M., Kim, M.C., Hartoularos, G.C., Targ, S., Sun, Y., Ogorodnikov, A., Bueno, R., Lu, A., et al. (2022). Single-cell RNA-seq reveals cell type-specific molecular and genetic associations to lupus. *Science* 376, eabf1970. <https://doi.org/10.1126/science.abf1970>.
24. Bergen, V., Lange, M., Peidli, S., Wolf, F.A., and Theis, F.J. (2020). Generalizing RNA velocity to transient cell states through dynamical modeling. *Nat. Biotechnol.* 38, 1408–1414. <https://doi.org/10.1038/s41587-020-0591-3>.
25. Eraslan, G., Simon, L.M., Mircea, M., Mueller, N.S., and Theis, F.J. (2019). Single-cell RNA-seq denoising using a deep count autoencoder. *Nat. Commun.* 10, 390. <https://doi.org/10.1038/s41467-018-07931-2>.
26. Lopez, R., Regier, J., Cole, M.B., Jordan, M.I., and Yosef, N. (2018). Deep generative modeling for single-cell transcriptomics. *Nat. Methods* 15, 1053–1058. <https://doi.org/10.1038/s41592-018-0229-2>.
27. Srivatsan, S.R., McFaline-Figueroa, J.L., Ramani, V., Saunders, L., Cao, J., Packer, J., Pliner, H.A., Jackson, D.L., Daza, R.M., Christiansen, L., et al. (2020). Massively multiplex chemical transcriptomics at single-cell resolution. *Science* 367, 45–51. <https://doi.org/10.1126/science.aax6234>.
28. Sikkema, L., Ramírez-Suástegui, C., Strobl, D.C., Gillett, T.E., Zappia, L., Madisson, E., Markov, N.S., Zaragosi, L.-E., Ji, Y., Ansari, M., et al. (2023). An integrated cell atlas of the lung in health and disease. *Nat. Med.* 29, 1563–1577. <https://doi.org/10.1038/s41591-023-02327-2>.
29. Hrovatin, K., Moifar, A.A., Zappia, L., Lapuerta, A.T., Lengerich, B., Kellis, M., and Theis, F.J. (2024). Integrating single-cell RNA-seq datasets with substantial batch effects. Preprint at bioRxiv. <https://doi.org/10.1101/2023.11.03.565463>.
30. Haniffa, M., Taylor, D., Linnarsson, S., Aronow, B.J., Bader, G.D., Barker, R.A., Camara, P.G., Camp, J.G., Chédotal, A., Copp, A., et al. (2021). A roadmap for the human developmental cell atlas. *Nature* 597, 196–205. <https://doi.org/10.1038/s41586-021-03620-1>.
31. Muus, C., Luecken, M.D., Eraslan, G., Sikkema, L., Waghray, A., Heimberg, G., Kobayashi, Y., Vaishnav, E.D., Subramanian, A., Smillie, C., et al. (2021). Single-cell meta-analysis of SARS-CoV-2 entry genes across tissues and demographics. *Nat. Med.* 27, 546–559. <https://doi.org/10.1038/s41591-020-01227-z>.
32. Regev, A., Teichmann, S.A., Lander, E.S., Amit, I., Benoist, C., Birney, E., Bodenmiller, B., Campbell, P., Carninci, P., Clatworthy, M., et al. (2017). The human cell atlas. *eLife* 6, e27041. <https://doi.org/10.7554/eLife.27041>.
33. Rood, J.E., Maartens, A., Hupalowska, A., Teichmann, S.A., and Regev, A. (2022). Impact of the Human Cell Atlas on medicine. *Nat. Med.* 28, 2486–2496. <https://doi.org/10.1038/s41591-022-02104-7>.
34. Liu, R., Qian, K., He, X., and Li, H. (2024). Integration of scRNA-seq data by disentangled representation learning with condition domain adaptation. *BMC Bioinformatics* 25, 116. <https://doi.org/10.1186/s12859-024-05706-9>.
35. Piran, Z., Cohen, N., Hoshen, Y., and Nitzan, M. (2024). Disentanglement of single-cell data with biolord. *Nat. Biotechnol.* 42, 1678–1683. <https://doi.org/10.1038/s41587-023-02079-x>.
36. Zhang, Z., Zhao, X., Bindra, M., Qiu, P., and Zhang, X. (2024). scDisInFact: disentangled learning for integration and prediction of multi-batch multi-condition single-cell RNA-sequencing data. *Nat. Commun.* 15, 912. <https://doi.org/10.1038/s41467-024-45227-w>.
37. Shamsaie, K., Megas, S., Asadollahzadeh, H., Teichmann, S.A., and Lotfollahi, M. (2024). Disentangling Covariates to Predict Counterfactuals for single-cell data. OpenReview. <https://openreview.net/forum?id=YeOUqnPVwM>.
38. Lotfollahi, M., Susmelj, A.K., De Donno, C., Ji, Y., Ibarra, I.L., Wolf, F.A., Yakubova, N., Theis, F.J., and Lopez-Paz, D. (2021). Learning interpretable cellular responses to complex perturbations in high-throughput screens. Preprint at bioRxiv. <https://doi.org/10.1101/2021.04.14.439903>.
39. Chen, R.T., Li, X., Grosse, R.B., and Duvenaud, D.K. (2018). Isolating sources of disentanglement in variational autoencoders. *Advances in Neural Information Processing Systems* 31. <https://doi.org/10.48550/arXiv.1802.04942>.
40. Creager, E., Madras, D., Jacobsen, J.-H., Weis, M., Swersky, K., Pitassi, T., and Zemel, R. (2019). Flexibly fair representation learning by disentanglement. In *International Conference on Machine Learning (PMLR)*, pp. 1436–1445. <https://doi.org/10.48550/arXiv.1906.02589>.
41. Luecken, M.D., Büttner, M., Chaichoompu, K., Danese, A., Interlandi, M., Mueller, M.F., Strobl, D.C., Zappia, L., Dugas, M., Colomé-Tatché, M., et al. (2022). Benchmarking atlas-level data integration in single-cell genomics. *Nat. Methods* 19, 41–50. <https://doi.org/10.1038/s41592-021-01336-8>.
42. Gayoso, A., Lopez, R., Xing, G., Boyeau, P., Valiollah Pour Amiri, V., Hong, J., Wu, K., Jayasuriya, M., Mehlman, E., Langevin, M., et al. (2022). A Python library for probabilistic analysis of single-cell omics data. *Nat. Biotechnol.* 40.2, 163–166. <https://doi.org/10.1038/s41587-021-01206-w>.
43. McInnes, L., Healy, J., and Melville, J. (2018). UMAP: Uniform manifold approximation and projection for dimension reduction. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1802.03426>.
44. Setty, M., Kisieliovas, V., Levine, J., Gayoso, A., Mazutis, L., and Pe’er, D. (2019). Characterization of cell fate probabilities in single-cell data with Palantir. *Nat. Biotechnol.* 37, 451–460. <https://doi.org/10.1038/s41587-019-0068-4>.
45. Xu, C., Lopez, R., Mehlman, E., Regier, J., Jordan, M.I., and Yosef, N. (2021). Probabilistic harmonization and annotation of single-cell transcriptomics data with deep generative models. *Mol. Syst. Biol.* 17, e9620. <https://doi.org/10.15252/msb.20209620>.
46. Manning, C.E., Eagle, A.L., Kwiatkowski, C.C., Achargui, R., Woodworth, H., Potter, E., Ohnishi, Y., Leininger, G.M., and Robison, A.J. (2019).

- Hippocampal subgranular zone FosB expression is critical for neurogenesis and learning. *Neuroscience* 406, 225–233. <https://doi.org/10.1016/j.neuroscience.2019.03.022>.
47. Chou, M.-Y., Hu, M.-C., Chen, P.-Y., Hsu, C.-L., Lin, T.-Y., Tan, M.-J., Lee, C.-Y., Kuo, M.-F., Huang, P.-H., Wu, V.-C., et al. (2022). RTL1/PEG11 imprinted in human and mouse brain mediates anxiety-like and social behaviors and regulates neuronal excitability in the locus coeruleus. *Hum. Mol. Genet.* 31, 3161–3180. <https://doi.org/10.1093/hmg/ddac110>.
48. Kitazawa, M., Sutani, A., Kaneko-Ishino, T., and Ishino, F. (2021). The role of eutherian-specific RTL1 in the nervous system and its implications for the Kagami-Ogata and Temple syndromes. *Genes Cells* 26, 165–179. <https://doi.org/10.1111/gtc.12830>.
49. Yin, K.-J., Hamblin, M., Fan, Y., Zhang, J., and Chen, Y.E. (2015). Krüppel-like factors in the central nervous system: novel mediators in Stroke. *Metab. Brain Dis.* 30, 401–410. <https://doi.org/10.1007/s11011-013-9468-1>.
50. Palmer, C.R., Liu, C.S., Romanow, W.J., Lee, M.-H., and Chun, J. (2021). Altered cell and RNA isoform diversity in aging Down syndrome brains. *Proc. Natl. Acad. Sci. USA* 118, e2114326118. <https://doi.org/10.1073/pnas.2114326118>.
51. Poirier, R., Cheval, H., Mailhes, C., Charnay, P., Davis, S., and Laroche, S. (2007). Paradoxical role of an Egr transcription factor family member, Egr2/Krox20, in learning and memory. *Front. Behav. Neurosci.* 1, 6. <https://doi.org/10.3389/neuro.08.006.2007>.
52. Voltan, R., Secchiero, P., Corallini, F., and Zauli, G. (2014). Selective induction of TP53/p53-inducible gene 3 (PIG3) in myeloid leukemic cells, but not in normal cells, by Nutlin-3. *Mol. Carcinog.* 53, 498–504. <https://doi.org/10.1002/mc.21985>.
53. Huang, B., and Vassilev, L.T. (2009). Reduced transcriptional activity in the p53 pathway of senescent cells revealed by the MDM2 antagonist Nutlin-3. *Aging (Albany, NY)* 1, 845–854. <https://doi.org/10.18632/aging.100091>.
54. Norman, T.M., Horlbeck, M.A., Replogle, J.M., Ge, A.Y., Xu, A., Jost, M., Gilbert, L.A., and Weissman, J.S. (2019). Exploring genetic interaction manifolds constructed from rich single-cell phenotypes. *Science* 365, 786–793. <https://doi.org/10.1126/science.aax4438>.
55. Vallender, S.S. (1974). Calculation of the Wasserstein distance between probability distributions on the line. *Theory Probab. Appl.* 18, 784–786. <https://doi.org/10.1137/1118101>.
56. Afriat, A., Zuzarte-Luís, V., Bahar Halpern, K.B., Buchauer, L., Marques, S., Chora, Á.F., Lahree, A., Amit, I., Mota, M.M., and Itzkovitz, S. (2022). A spatiotemporally resolved single-cell atlas of the Plasmodium liver stage. *Nature* 611, 563–569. <https://doi.org/10.1038/s41586-022-05406-5>.
57. Suo, C., Dann, E., Goh, I., Jardine, L., Kleshchevnikov, V., Park, J.-E., Botting, R.A., Stephenson, E., Engelbert, J., Tuong, Z.K., et al. (2022). Mapping the developing human immune system across organs. *Science* 376, eabo0510. <https://doi.org/10.1126/science.abo0510>.
58. Braun, E., Danan-Gotthold, M., Borm, L.E., Lee, K.W., Vinsland, E., Lönnerberg, P., Hu, L., Li, X., He, X., Andrusivová, Ž., et al. (2023). Comprehensive cell atlas of the first-trimester developing human brain. *Science* 382, eadf1226. <https://doi.org/10.1126/science.adf1226>.
59. Miller, A.J., Yu, Q., Czerwinski, M., Tsai, Y.-H., Conway, R.F., Wu, A., Holloway, E.M., Walker, T., Glass, I.A., Treutlein, B., et al. (2020). In vitro and in vivo development of the human airway at single-cell resolution. *Dev. Cell* 53, 117–128.e6. <https://doi.org/10.1016/j.devcel.2020.01.033>.
60. Weinberger, E., Lin, C., and Lee, S.-I. (2023). Isolating salient variations of interest in single-cell data with contrastiveVI. *Nat. Methods* 20, 1336–1345. <https://doi.org/10.1038/s41592-023-01955-3>.
61. Oh, C., Won, H., So, J., Kim, T., Kim, Y., Choi, H., and Song, K. (2022). Learning fair representation via distributional contrastive disentanglement. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (ACM)*, pp. 1295–1305. <https://doi.org/10.1145/3534678.3539232>.
62. Kingma, D.P., and Welling, M. (2013). Auto-encoding variational bayes. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1312.6114>.
63. Buettner, F., and Theis, F.J. (2012). A novel approach for resolving differences in single-cell gene expression patterns from zygote to blastocyst. *Bioinformatics* 28, i626–i632. <https://doi.org/10.1093/bioinformatics/bts385>.
64. De Donno, C., Hediye-Zadeh, S., Moinfar, A.A., Wagenstetter, M., Zappia, L., Lotfollahi, M., and Theis, F.J. (2023). Population-level integration of single-cell datasets enables multi-scale analysis across samples. *Nat. Methods* 20, 1683–1692. <https://doi.org/10.1038/s41592-023-02035-2>.
65. Tran, H.T.N., Ang, K.S., Chevrier, M., Zhang, X., Lee, N.Y.S., Goh, M., and Chen, J. (2020). A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome Biol.* 21, 12. <https://doi.org/10.1186/s13059-019-1850-9>.
66. Lample, G., Zeghidour, N., Usunier, N., Bordes, A., Denoyer, L., and Ranzato, M. (2018). Fader Networks: Manipulating Images by Sliding Attributes. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1706.00409>.
67. Bunne, C., Stark, S.G., Gut, G., Del Castillo, J.S., Levesque, M., Lehmann, K.-V., Pelkmans, L., Krause, A., and Rättsch, G. (2023). Learning single-cell perturbation responses using neural optimal transport. *Nat. Methods* 20, 1759–1768. <https://doi.org/10.1038/s41592-023-01969-x>.
68. Lotfollahi, M., Klimovskaia Susmelj, A., De Donno, C., Hetzel, L., Ji, Y., Ibarra, I.L., Srivatsan, S.R., Naghipourfar, M., Daza, R.M., Martin, B., et al. (2023). Predicting cellular responses to complex perturbations in high-throughput screens. *Mol. Syst. Biol.* 19, e11517. <https://doi.org/10.15252/msb.202211517>.
69. Inecik, K., Uhlmann, A., Lotfollahi, M., and Theis, F. (2022). MultiCPA: Multimodal compositional perturbation autoencoder. Preprint at bioRxiv. <https://doi.org/10.1101/2022.07.08.499049>.
70. Butler, A., Hoffman, P., Smibert, P., Papalexi, E., and Satija, R. (2018). Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.* 36, 411–420. <https://doi.org/10.1038/nbt.4096>.
71. Virshup, I., Rybakov, S., Theis, F.J., Angerer, P., and Wolf, F.A. (2021). anndata: Annotated data. Preprint at bioRxiv. <https://doi.org/10.1101/2021.12.16.473007>.
72. Biology, C.S.-C., Abdulla, S., Aevermann, B., Assis, P., Badajoz, S., Bell, S.M., Bezzi, E., Cakir, B., Chaffer, J., Chambers, S., et al. (2023). CZ CELLxGENE Discover: A single-cell data platform for scalable exploration, analysis and modeling of aggregated data. *Nucleic Acids Res.* 53, D886–D900. <https://doi.org/10.1093/nar/gkac1142>.
73. Tong, Q., and Kobayashi, K. (2021). Entropy-Regularized Optimal Transport on Multivariate Normal and q-normal Distributions. *Entropy (Basel)* 23, 302. <https://doi.org/10.3390/e23030302>.
74. Kim, M., Wang, Y., Sahu, P., and Pavlovic, V. (2019). Relevance factor VAE: Learning and identifying disentangled factors. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1902.01568>.
75. Silva, W.B., Freitas, C.C., Sant’Anna, S.J.S., and Frery, A.C. (2013). Classification of segments in PolSAR imagery by minimum stochastic distances between Wishart distributions. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* 6, 1263–1273. <https://doi.org/10.1109/JSTARS.2013.2248132>.
76. Cao, D., Zhou, S., Liu, H., Liu, J., and Zang, H. (2017). Signal censoring and fusing with system-level communication constraints in multistatic radar: a J-divergence and Bhattacharyya distance-based approach. *IET Radar Sonar Navig.* 11, 1802–1814. <https://doi.org/10.1049/iet-rsn.2017.0159>.
77. Sintini, L., and Kunze, L. (2020). Unsupervised and Semi-supervised Novelty Detection using Variational Autoencoders in Opportunistic Science Missions. *British Machine Vision Conference*. <https://doi.org/10.5244/C.34.149>.

78. Zhang, S., Xie, L., Cui, Y., Carone, B.R., and Chen, Y. (2022). Detecting Fear-Memory-Related Genes from Neuronal scRNA-seq Data by Diverse Distributions and Bhattacharyya Distance. *Biomolecules* *12*, 1130. <https://doi.org/10.3390/biom12081130>.
79. Baker, D.N., Dyjack, N., Braverman, V., Hicks, S.C., and Langmead, B. (2021). Fast and memory-efficient scRNA-seq k-means clustering with various distances. In Proceedings of the 12th ACM Conference on Bioinformatics, Computational Biology, and Health Informatics, BCB '21 (ACM). <https://doi.org/10.1145/3459930.3469523>.
80. Moon, K.R., Stanley, J.S., Burkhardt, D., van Dijk, D., Wolf, G., and Krishnaswamy, S. (2018). Manifold learning-based methods for analyzing single-cell RNA-sequencing data. *Curr. Opin. Syst. Biol.* *7*, 36–46. <https://doi.org/10.1016/j.coisb.2017.12.008>.
81. Wang, J., Agarwal, D., Huang, M., Hu, G., Zhou, Z., Ye, C., and Zhang, N.R. (2019). Data denoising with transfer learning in single-cell transcriptomics. *Nat. Methods* *16*, 875–878. <https://doi.org/10.1038/s41592-019-0537-1>.
82. Lin, E., Mukherjee, S., and Kannan, S. (2020). A deep adversarial variational autoencoder model for dimensionality reduction in single-cell RNA sequencing analysis. *BMC Bioinformatics* *21*, 64. <https://doi.org/10.1186/s12859-020-3401-5>.
83. Shree, A., Pavan, M.K., and Zafar, H. (2023). scDREAMER for atlas-level integration of single-cell datasets using deep generative model paired with adversarial classifier. *Nat. Commun.* *14*, 7781. <https://doi.org/10.1038/s41467-023-43590-8>.
84. Tschannen, M., Bachem, O., and Lucic, M. (2018). Recent advances in autoencoder-based representation learning. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1812.05069>.
85. Khemakhem, I., Kingma, D., Monti, R., and Hyvarinen, A. (2020). Variational autoencoders and nonlinear ICA: A unifying framework. In International Conference on Artificial Intelligence and Statistics (PMLR), pp. 2207–2217. <https://doi.org/10.48550/arXiv.1907.04809>.
86. Rousseeuw, P.J. (1987). Silhouettes: A Graphical Aid to the Interpretation and Validation of Cluster Analysis. *J. Comput. Appl. Math.* *20*, 53–65. [https://doi.org/10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7).
87. Haghverdi, L., Lun, A.T.L., Morgan, M.D., and Marioni, J.C. (2018). Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nat. Biotechnol.* *36*, 421–427. <https://doi.org/10.1038/nbt.4091>.
88. Sokolova, M., and Lapalme, G. (2009). A systematic analysis of performance measures for classification tasks. *Inf. Process. Manag.* *45*, 427–437. <https://doi.org/10.1016/j.ipm.2009.03.002>.
89. Gayoso, A., Kim, M., Kronfeld, O., Hong, J., and Nir, Y. (2026). scib-metrics: Accelerated, Python-only, single-cell integration benchmarking metrics. GitHub. <https://github.com/yoseflab/scib-metrics>.
90. Duncan, T.E. (1970). On the calculation of mutual information. *SIAM J. Appl. Math.* *19*, 215–220. <https://doi.org/10.1137/0119020>.
91. Kraskov, A., Stögbauer, H., and Grassberger, P. (2004). Estimating mutual information. *Phys. Rev. E Stat. Nonlin. Soft Matter Phys.* *69*, 066138. <https://doi.org/10.1103/PhysRevE.69.066138>.
92. Vinh, N., Epps, J., and Bailey, J. (2010). Information Theoretic Measures for Clusterings Comparison: Variants, Properties, Normalization and Correction for Chance. *J. Mach. Learn. Res.* *11*, 2837–2854. <https://doi.org/10.5555/1756006.1953024>.
93. Higgins, I., Matthey, L., Pal, A., Burgess, C.P., Glorot, X., Botvinick, M.M., Mohamed, S., and Lerchner, A. (2016). beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework. International Conference on Learning Representations. <https://api.semanticscholar.org/CorpusID:46798026>.
94. Wu, Y., Price, L.C., Wang, Z., Ioannidis, V.N., Barton, R.A., and Karypis, G. (2022). Variational causal inference. Preprint at arXiv. <https://doi.org/10.48550/arXiv.2209.05935>.
95. Kumar, A., Sattigeri, P., and Balakrishnan, A. (2017). Variational inference of disentangled latent concepts from unlabeled observations. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1711.00848>.
96. Kim, H., and Mnih, A. (2018). Disentangling by factorising. In International Conference on Machine Learning (PMLR), pp. 2649–2658. <https://doi.org/10.48550/arXiv.1802.05983>.
97. Seplarskaia, A., Kiseleva, J., and de Rijke, M. (2019). How to not measure disentanglement. Preprint at arXiv. <https://doi.org/10.48550/arXiv.1910.05587>.
98. Hubert, L.J., and Arabie, P. (1985). Comparing partitions. *J. Classif.* *2*, 193–218. <https://doi.org/10.1007/BF01908075>.
99. Büttner, M., Miao, Z., Wolf, F.A., Teichmann, S.A., and Theis, F.J. (2018). A test metric for assessing single-cell RNA-seq batch correction. *Nat. Methods* *16*, 43–49. <https://doi.org/10.1038/s41592-018-0254-1>.
100. Hetzel, L., Boehm, S., Kilbertus, N., Günemann, S., Lotfollahi, M., Theis, F., et al. (2022). Predicting cellular responses to novel drug perturbations at a single-cell resolution. *Advances in Neural Information Processing Systems* *35*, 26711–26722. <https://doi.org/10.48550/arXiv.2204.13545>.
101. Korsunsky, I., Millard, N., Fan, J., Slowikowski, K., Zhang, F., Wei, K., Baglaenko, Y., Brenner, M., Loh, P.-R., and Raychaudhuri, S. (2019). Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat. Methods* *16*, 1289–1296. <https://doi.org/10.1038/s41592-019-0619-0>.

STAR★METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Deposited data		
Afriat dataset, single-cell RNA sequencing data of host and parasite gene expression during the liver stage of the rodent malaria parasite Plasmodium berghei ANKA.	Afriat et al. ⁵⁶	GSE181725 (GEO)
Suo dataset, a multi-organ, single-cell transcriptomic perspective, capturing dynamic immune system developments across nine prenatal human tissues during embryonic stages	Suo et al. ⁵⁷	E-MTAB-11343 (ArrayExpress)
Braun dataset, a comprehensive single-cell transcriptomic analysis of the human brain during the crucial first trimester	Braun et al. ⁵⁸	EGAS00001004107 (EGA)
Miller dataset, provides a detailed single-cell mRNA sequencing atlas of human lung development from 11.5 to 21 weeks, integrated with studies on homogeneous human bud tip organoid cultures	Miller et al. ⁵⁹	E-MTAB-8221 (ArrayExpress)
Sciplex dataset, derived from the sci-Plex technology using nuclear hashing, quantifies transcriptional responses to chemical perturbations at single-cell resolution	Srivatsan et al. ²⁷	GSE139944 (GEO)
Norman dataset, leverages high-content Perturb-seq technology to explore cellular and organismal complexity through combinatorial gene expression	Norman et al. ⁵⁴	GSE133344 (GEO)
Software and algorithms		
TarDis	This study	GitHub: theislab/tardis; DOI: 10.5281/zenodo.17135254
β -TCVAE	Chen et al. ³⁹	GitHub: rtqichen/beta-tcvae
FFVAE	Creager et al. ⁴⁰	GitHub: nomnomnonono/FFVAE
Biolord	Piran et al. ³⁵	GitHub: nitzanlab/biolord
CPA	Lotfollahi et al. ⁶⁰	GitHub: theislab/cpa
scIB	Luecken et al. ⁴¹	GitHub: YosefLab/scib-metrics
scVI	Lopez et al. ²⁶	GitHub: scverse/scvi-tools
scANVI	Xu et al. ⁴⁵	GitHub: scverse/scvi-tools
Harmony	Korsunsky et al. ⁵¹	GitHub: lilab-bcb/harmony-pytorch
scDisInFact	Zhang et al. ³⁶	GitHub: ZhangLabGT/scDisInFact
inVAE	Aliee et al. ¹⁰	GitHub: theislab/inVAE; arXiv:2307.00558
Other		
Mamba environment for reproducibility, including all used software and their versions	This study	GitHub: theislab/tardis/environment

METHOD DETAILS

Overview

Let \mathcal{D} represent a single-cell genomics dataset containing N_C cells, where each cell, denoted as n , is characterized by its gene expression (\mathbf{x}_n) and associated covariates (\mathbf{s}_n). The gene expression is represented by a count vector $\mathbf{x}_n = [x_{ng}]_{g=1}^{N_G}$, where $x_{ng} \in \mathbb{Z}_{\geq 0}$ is the expression count of gene g , and N_G is the total number of genes in the dataset. Additionally, each cell n is associated with a vector of covariates $\mathbf{s}_n = [s_{nk}]_{k=1}^{N_K}$, which may be either continuous or discrete, and N_K indicates the number of covariates. *TarDis* constructs a latent representation \mathbf{z}_n for gene expression \mathbf{x}_n , organized as $\mathbf{z}_n = (\mathbf{z}_{n0}, [\mathbf{z}_{nk}]_{k \in J_k})$, where $J_k \subseteq \{1, \dots, N_K\}$ denotes the subset of covariates targeted for disentanglement. Note that since indexing of targeted covariates uses $k \in J_k$ starting from 1, the subscript '0' in \mathbf{z}_{n0} does not correspond to any index in J_k but instead represents the portion of the latent space reserved for information not explained by the targeted covariates. Specifically, \mathbf{z}_{nk} is a latent vector constructed for each targeted covariate, while \mathbf{z}_{n0} captures residual information independent of targeted covariates. During model training, *TarDis* employs a specific strategy

to foster disentanglement by generating pairs of additional latent vectors $(\mathbf{z}_n^{(k)})^-$ and $(\mathbf{z}_n^{(k)})^+$ corresponding to two data points $(\mathbf{x}_n^{(k)})^-$ and $(\mathbf{x}_n^{(k)})^+$. These data points are selected *randomly* and differ in the k_{th} covariate value, such that $(\mathbf{s}_{nk}^{(k)})^+ = \mathbf{s}_{nk}$ and $(\mathbf{s}_{nk}^{(k)})^- \neq \mathbf{s}_{nk}$. Observe that, for each categorical covariate k , a *positive* pair for covariate k consists of the original data point and another randomly selected data point from the dataset whose category matches exactly for covariate k ; formally, $(\mathbf{s}_{nk}^{(k)})^+ = \mathbf{s}_{nk}$. On the other hand, a *negative* pair for covariate k comprises the original data point and another randomly selected data point whose category differs (i.e., belongs to any category other than the original) for covariate k . Formally, $(\mathbf{s}_{nk}^{(k)})^- \neq \mathbf{s}_{nk}$. These definitions are consistently maintained in all subsequent formulations involving pair-wise terms in the loss function.

The primary objective of *TarDis* training is to optimize the latent vectors based on a distance measure F . While F is defined conceptually as a real-valued function, $F: \mathbb{R}^{|\mathbf{z}_{nk}|} \rightarrow \mathbb{R}_{\geq 0}$, here just to illustrate the underlying concept, practical implementation typically employ multiple loss terms instead of a single function for optimizing latent vectors, as will be discussed in further detail. For each covariate $k \in \mathcal{J}_k$, F should satisfy $F(\mathbf{z}_{nk}, (\mathbf{z}_{nk}^{(k)})^-) \geq F(\mathbf{z}_{nk}, (\mathbf{z}_{nk}^{(k)})^+)$, implying that latent vector \mathbf{z}_{nk} should be more similar to another vector that shares the same covariate value, $(\mathbf{z}_{nk}^{(k)})^+$, than to a vector with a different covariate value, $(\mathbf{z}_{nk}^{(k)})^-$. Furthermore, the latent vector \mathbf{z}_{n0} should show equal similarity to any other vectors regardless of their covariate values, whether $(\mathbf{z}_{n0}^{(k)})^-$ and $(\mathbf{z}_{n0}^{(k)})^+$, thus fulfilling the condition: $F(\mathbf{z}_{n0}, (\mathbf{z}_{n0}^{(k)})^-) = F(\mathbf{z}_{n0}, (\mathbf{z}_{n0}^{(k)})^+)$. This equality ensures that \mathbf{z}_{n0} remains unaffected by covariate-specific information, thereby providing a covariate-neutral representation of the cell's gene expression. Ultimately, the aim of *TarDis* is to produce a latent representation in which \mathbf{z}_{nk} reflects the influence of its corresponding covariate \mathbf{s}_{nk} , while \mathbf{z}_{n0} offers a covariate-neutral representation of the cell's gene expression profile, unaffected by any covariate-specific variations.¹⁰

VAE skeleton

TarDis builds upon a variational autoencoder (VAE) to construct a high-fidelity generative model that underpins our disentanglement objectives. The VAE component optimization guided by the Evidence Lower Bound (ELBO), a surrogate for the intractable marginal log-likelihood as shown in Equation 1.⁶² Here, the covariates, \mathbf{s}_n , are pivotal for capturing factors that might influence the observed data, such as batch effects. *TarDis* incorporates the target covariates as \mathbf{s}_n , and also allows inclusion of non-target covariates, providing flexibility in managing different types of data impacts. The first term of \mathcal{L}_{VAE} represents the reconstruction loss, \mathcal{L}_R , which quantifies the expected negative log-likelihood of the observed data \mathbf{x}_n , conditioned on the latent variables, \mathbf{z}_n . It aims to minimize the discrepancy between the observed data and its reconstruction from the latent space. The reconstruction loss is formally expressed using the negative binomial (NB) distribution, ideal for capturing the count variability inherent in data types like single-cell genomics (Equation 2). In this equation, Γ denotes the gamma function, μ and θ refer to the mean and inverse dispersion parameters of the negative binomial distribution, respectively.²² The second term measures the Kullback-Leibler divergence (KL), \mathcal{L}_{KL} , penalizing deviations of the learned posterior distribution $q_\psi(\mathbf{z}_n|\mathbf{x}_n, \mathbf{s}_n)$ from the prior distribution $p(\mathbf{z}_n)$. In Equation 3, the approximate posterior distribution is assumed to be Gaussian distribution with mean μ_n and diagonal covariance matrix Σ_n , and the prior distribution $p(\mathbf{z}_n)$ is typically a standard normal distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$ where \mathbf{I} is the identity matrix in $\mathbb{R}^{|\mathbf{z}_n| \times |\mathbf{z}_n|}$. This acts as a regularizer, aligning the latent embeddings with a predefined distribution, to ensure the statistical robustness and generalization capability of the model.⁶²

$$\mathcal{L}_{\text{VAE}}(\theta, \psi; \mathbf{x}_n, \mathbf{s}_n) = -\mathbb{E}_{q_\psi(\mathbf{z}_n|\mathbf{x}_n, \mathbf{s}_n)}[\log p_\theta(\mathbf{x}_n|\mathbf{z}_n)] + D_{\text{KL}}(q_\psi(\mathbf{z}_n|\mathbf{x}_n, \mathbf{s}_n) \| p(\mathbf{z}_n)) \quad (\text{Equation 1})$$

$$\mathcal{L}_R = \text{NB}(\mathbf{x}_n; \mu_n, \theta_n) = \frac{\Gamma(\mathbf{x}_n + \theta_n)}{\Gamma(\mathbf{x}_n + 1)\Gamma(\theta_n)} \left(\frac{\theta_n}{\theta_n + \mu_n}\right)^{\theta_n} \left(\frac{\mu_n}{\theta_n + \mu_n}\right)^{\mathbf{x}_n} \quad (\text{Equation 2})$$

$$\mathcal{L}_{\text{KL}} = D_{\text{KL}}(\mathcal{N}(\mu_n, \Sigma_n) \| \mathcal{N}(\mathbf{0}, \mathbf{I})) \quad (\text{Equation 3})$$

Auxiliary loss

In *TarDis* model training, the VAE optimization is intertwined with a specifically designed auxiliary loss component introduced, \mathcal{L}_C , to construct $\mathbf{z}_n = (\mathbf{z}_{n0}, [\mathbf{z}_{nk}]_{k \in \mathcal{J}_k})$ with $\mathbf{z}_{nk} \sim \mathcal{N}(\mu_{nk}, \Sigma_{nk})$. The overall loss function of *TarDis* integrates these components through a weighted sum, controlled by hyperparameters λ_C , λ_{KL} , and λ_R (Equation 8). Specifically, \mathcal{L}_C is a composite loss function that incorporates four distinct loss components for each covariate.⁶³ For each target covariate \mathbf{s}_{nk} , the loss function, $\mathcal{L}_C^{(k)}$, includes $(N_{\mathcal{L}}^{(k)})^+$ positive and $(N_{\mathcal{L}}^{(k)})^-$ negative loss terms. Similarly, for the covariate-free representation \mathbf{z}_{n0} , it includes $(N_{\mathcal{L}}^{(k_0)})^+$ positive and $(N_{\mathcal{L}}^{(k_0)})^-$ negative terms. The losses for positive pairs and negative pairs given in Equations 4 and 5, respectively. Here, the λ values are hyperparameters that determine the weight of each loss component, while the \mathcal{L} loss functions encompass metrics such as KL divergence and mean squared error (MSE).¹¹ Thus, the overall covariate loss, $\mathcal{L}_C^{(k)}$, is computed as the sum of these two pair losses, as specified in Equation 6. By aggregating these individual covariate losses, the total auxiliary loss, \mathcal{L}_C , is expressed in Equation 7.

The configuration of $\mathcal{L}_C^{(k)}$ is meticulously designed to meet several critical objectives within the *TarDis* framework. First, by minimizing the distance between $(\mathbf{z}_{nk}^{(k)})^+$ and \mathbf{z}_{nk} , the model ensures that the latent representations of positive examples closely align with their corresponding covariate within respective latent subset, accurately reflecting specific characteristics. In contrast, it maximizes the distance between $(\mathbf{z}_{nk}^{(k)})^-$ and \mathbf{z}_{nk} , thereby promoting clear differentiation in the latent representations of negative examples and enhancing the distinction between different covariates. Additionally, the model strategy involves maximizing the distance between $(\mathbf{z}_{n0}^{(k)})^+$ and \mathbf{z}_{n0} , while minimizing the distance between $(\mathbf{z}_{n0}^{(k)})^-$ and \mathbf{z}_{n0} . This approach ensures that \mathbf{z}_{n0} remains free from covariate-specific influences, maintaining its role as a covariate-neutral representation. These operations collectively ensure that covariate information is precisely captured in the respective targeted latent subsets, \mathbf{z}_{nk} , and effectively isolated from \mathbf{z}_{n0} .

$$\begin{aligned} (\mathcal{L}_C^{(k)})^+(\varphi; \mathbf{x}_n, \mathbf{s}_n) &= \frac{1}{(N_{\mathcal{L}}^{(k)})^+} \sum_{i=1}^{(N_{\mathcal{L}}^{(k)})^+} [(\lambda_C^{(k)})^+ (\mathcal{L}_C^{(k)})^+(\varphi; \mathbf{x}_n, \mathbf{s}_n)] \\ &+ \frac{1}{(N_{\mathcal{L}}^{(k_0)})^+} \sum_{i=1}^{(N_{\mathcal{L}}^{(k_0)})^+} [(\lambda_C^{(k_0)})^+ (\mathcal{L}_C^{(k_0)})^+(\varphi; \mathbf{x}_n, \mathbf{s}_n)] \end{aligned} \quad (\text{Equation 4})$$

$$\begin{aligned} (\mathcal{L}_C^{(k)})^-(\varphi; \mathbf{x}_n, \mathbf{s}_n) &= \frac{1}{(N_{\mathcal{L}}^{(k)})^-} \sum_{i=1}^{(N_{\mathcal{L}}^{(k)})^-} [(\lambda_C^{(k)})^- (\mathcal{L}_C^{(k)})^-(\varphi; \mathbf{x}_n, \mathbf{s}_n)] \\ &+ \frac{1}{(N_{\mathcal{L}}^{(k_0)})^-} \sum_{i=1}^{(N_{\mathcal{L}}^{(k_0)})^-} [(\lambda_C^{(k_0)})^- (\mathcal{L}_C^{(k_0)})^-(\varphi; \mathbf{x}_n, \mathbf{s}_n)] \end{aligned} \quad (\text{Equation 5})$$

$$\mathcal{L}_C^{(k)}(\varphi; \mathbf{x}_n, \mathbf{s}_n) = (\mathcal{L}_C^{(k)})^+(\varphi; \mathbf{x}_n, \mathbf{s}_n) + (\mathcal{L}_C^{(k)})^-(\varphi; \mathbf{x}_n, \mathbf{s}_n) \quad (\text{Equation 6})$$

$$\mathcal{L}_C(\varphi; \mathbf{x}_n, \mathbf{s}_n) = \frac{1}{|J_K|} \sum_{k=1}^{|J_K|} \mathcal{L}_C^{(k)}(\varphi; \mathbf{x}_n, \mathbf{s}_n) \quad (\text{Equation 7})$$

$$\mathcal{L}_{\text{TarDis}}(\theta, \varphi; \mathcal{D}) = \frac{1}{N_C} \sum_{(\mathbf{x}_n, \mathbf{s}_n) \in \mathcal{D}} [\lambda_R \mathcal{L}_R(\theta, \varphi; \mathbf{x}_n, \mathbf{s}_n) + \lambda_{\text{KL}} \mathcal{L}_{\text{KL}}(\varphi; \mathbf{x}_n, \mathbf{s}_n) + \lambda_C \mathcal{L}_C(\varphi; \mathbf{x}_n, \mathbf{s}_n)] \quad (\text{Equation 8})$$

To clarify the construction of the observation-specific loss, note that the complete objective is built hierarchically by aggregating per-observation contributions across data points and covariates. Each data point $(\mathbf{x}_n, \mathbf{s}_n)$ generates partial losses for each targeted covariate k , comprising a negative-pair term $(\mathcal{L}_C^{(k)})^-$ (Equation 5) and a positive-pair term $(\mathcal{L}_C^{(k)})^+$ (Equation 4), where i indexes the respective pairs. For covariate k , the negative-pair loss in Equation 5 sums over $i \in \{1, \dots, (N_{\mathcal{L}}^{(k)})^-\}$ to aggregate contributions from all negative pairs, while the positive-pair loss in Equation 4 analogously sums over $i \in \{1, \dots, (N_{\mathcal{L}}^{(k)})^+\}$. These covariate-specific components are combined in Equation 6 as $\mathcal{L}_C^{(k)} = \sum_i (\mathcal{L}_C^{(k)})^- + \sum_i (\mathcal{L}_C^{(k)})^+$, ensuring each data point's contribution is preserved. The global auxiliary loss (Equation 7) then averages $\mathcal{L}_C^{(k)}$ over all targeted covariates $k \in J_K$, which is finally combined with the VAE's reconstruction and KL terms in Equation 8 to form the total loss $\mathcal{L}_{\text{TarDis}}$. This hierarchical formulation explicitly decomposes how individual observations, through their positive and negative pairs, influence the training objective at both covariate-specific and dataset-wide levels.

Various options for the loss functions used in $\mathcal{L}_C^{(k)}$ are explored. Although the theoretical framework primarily employs KL divergence as the loss metric, the principle is also applicable to various losses between anchor points and negative or positive samples with minor adjustments.¹¹ Also note that, for continuous covariates, we adopt an analogous strategy but weigh the negative-pair penalty based on the distance between covariate values rather than using discrete category mismatches, while positive-pair loss components disappears. The optimization of various hyperparameters, including the individual loss weights, is conducted once and uniformly applied across all experiments, unless explicitly stated otherwise.¹² The experiments and benchmarking processes utilize a diverse array of datasets to ensure comprehensive testing and validation of the model. Each dataset is selected to represent different types and scales of data challenges.¹³ Various evaluation metrics are used to assess the model's performance, with a full discussion provided in 'evaluation metrics' section.¹⁴ The assumptions behind the theoretical framework are discussed in 'theoretical assumptions' section.¹⁵

Additional details

Related work

Models like single-cell Variational Inference (scVI) facilitate data integration by incorporating environmental variables such as experimental batches or sequencing protocols alongside gene data, using one-hot vectors processed through a conditional variational autoencoder (cVAE) to reduce technical noise.²⁶ Its extension, single-cell ANnotation using Variational Inference (scANVI), builds on this by introducing cell annotations in a semi-supervised approach, thus enhancing cell integration across diverse environments and adeptly capturing cell type variations.^{45,64} Despite their effective integration, these models may over-correct, adjusting biological signals while targeting technical noise, which can obscure subtle biological variations such as inter-patient differences or treatment effects. Moreover, these methodologies tend to aggregate all sources of spurious correlations indiscriminately, failing to discern the unique characteristics of each source.^{26,64,65} This approach inadequately addresses the nuanced interactions between these sources and biological signals, particularly problematic with continuous spurious covariates such as age or drug dosage. Models equipped to continuously adapt to these subtle variations are thus essential, ensuring that biological insights derived from single-cell genomics are not confounded by these varying conditions.

Several models in single-cell genomics have explored creating multiple latent spaces to handle different sources of variability distinctly. For instance, contrastiveVI models each covariate separately, developing a shared latent space for the common variability across covariates and an exclusive latent space for the target covariate's unique variability.⁶⁰ Similarly, single cell disentangled Integration preserving condition-specific Factors (scDisInFact) develops a shared latent space specifically designed to account for and eliminate batch effects, while simultaneously maintaining separate latent spaces for other covariates, isolating and preserving the variations from batch influences.³⁶ Yet, none of these approaches offer a control latent space dedicated to retaining batch effects while filtering out the influences of other covariates, essential for accurately distinguishing between variations caused by batch effects and those arising from true biological differences. Such methods draw inspiration from broader approaches focused on fair and disentangled representation, such as Flexibly Fair VAE (FFVAE) and Fader networks, and unsupervised disentanglement techniques such as Total Correlation VAE (β -TCVAE).^{37,39,61,66} The cell optimal transport model (CellOT) uses optimal transport (OT) methods to align cells from control and perturbed conditions, but its non-generative, single-covariate focus limits broader applicability.⁶⁷ Biolord offers a unique approach to supervised disentanglement, yet it faces scalability issues due to per-cell optimization.³⁵ The invariant VAE (inVAE) method introduces conditional priors within the VAE framework to effectively disentangle spurious and invariant correlations. While it offers nuanced disentanglement, inVAE faces optimization challenges, particularly in large datasets, and does not separate latent representations for individual covariates, and does not support continuous covariates naively limiting its ability to analyze complex interactions between various biological conditions in detail.¹⁰ On the other hand, Compositional Perturbation Autoencoder (CPA) handles drug perturbations and produce latent embedding but their assumption of linearity in the latent space limits capturing complex, non-linear biological interactions.^{68,69}

While existing approaches in single-cell genomics have notably advanced the disentanglement of spurious and invariant correlations, they predominantly excel within narrowly defined scenarios. Many models, however, simplify continuous covariates by categorizing them, which undermines the granularity of biological insights and limits their applicability in precision medicine. Beyond this, there's a critical need for models that not only handle the diversity of single-cell data but also scale efficiently and train effectively given the heterogeneity inherent in these datasets. Despite the innovative nature of these methods, they are often tailored to specific experimental conditions rather than offering a universal solution across the diverse landscape of single-cell analysis. There remains an unmet need for a comprehensive model that excels in data integration, out-of-distribution prediction, and serves as a robust platform for addressing intricate biological questions across various conditions and experimental setups.

Limitations

While *TarDis* introduces significant advancements in disentangling complex covariate structures in single-cell genomics, it is important to acknowledge several inherent limitations. *TarDis* operates under a supervised learning paradigm, which necessitates access to pre-labeled covariates. This requirement limits its applicability to datasets where such labels are readily available and accurately annotated, constraining its utility in less structured environments.

A notable limitation of *TarDis* is the potential for overfitting. Although rigorous validation protocols and robust regularization strategies, including elevated dropout rates and weight decay—more aggressive than those utilized in generic VAE models like scVI—are employed, the risk remains. In our study, the hyperparameters were carefully optimized at the onset of all experiments, ensuring consistent conditions across all tests, which mitigated the concerns of overfitting. It is important to note that our successful one-time optimization and the avoidance of overfitting in single-cell genomics data do not guarantee similar outcomes across other data types, hence users must conduct cautious benchmarking on validation splits to ensure the model's generalizability.

Moreover, the disentanglement of interdependent covariates introduces unique challenges. For example, accurately disentangling *age* and *donor* in a single-cell genomics data as covariates requires the presence of multiple donors of varying ages to prevent the model from conflating these factors. Without such diversity, the model risks inaccurately attributing the influence of one covariate to another, thereby undermining the reliability of the disentanglement, particularly evident in our validation splits.

Additionally, the implementation of *TarDis* introduces computational overhead, slightly slowing down the processing speed. Nevertheless, this does not significantly impact performance, even with large datasets like the *Braun* dataset, which comprises

1.6 million cells. The primary bottleneck arises from the selection of counteractive minibatches for each covariate during training, which is quantified to increase the average training time by approximately 1.8 times in comparison to scVI, when three covariates were targeted.

The encoding of covariates in a one-hot format, \mathbf{s}_n , while optional as mentioned in ‘STAR Methods’ section, generally fosters better disentanglement in the validation splits. However, the dependency of disentanglement on the input space may necessitate further optimization. This adjustment is crucial for enhancing the model’s utility in specific downstream tasks, as demonstrated in our analysis using the *Norman* dataset in ‘results’ section.

Lastly, *TarDis* necessitates numerous hyperparameters, especially concerning the loss weights for each of the four terms associated with every covariate. This complexity was manageable in our experiments through our aforementioned one-time optimization, and it did not present issues for single-cell data. However, adapting the model to new datasets could necessitate further tuning, potentially complicating its application across varied contexts. It is also important to underscore the model assumptions in ‘theoretical assumptions’ section, as these foundational assumptions highlight potential limitations and areas where *TarDis* might encounter challenges.

Model training details

Compute Resources and System Configuration: For the computational tasks in our research, we employed *NVIDIA Tesla A100* GPUs, which feature 40 GB of high-bandwidth HBM2 memory each. This GPU architecture is specifically designed for accelerating machine learning and high-performance computing applications, providing substantial throughput for both single and mixed-precision computations. We allocated 64 GB of GPU memory for processing large training datasets, which facilitated efficient handling of extensive computational operations without the need for frequent data swapping, thereby minimizing I/O overhead. For smaller datasets, a reduced memory allocation of 16 GB was used, which optimized resource utilization without compromising performance. On the CPU side, our computational nodes were equipped with dual *Intel Xeon Gold 6230* processors. Each processor offers 20 cores operating at a base frequency of 2.1 GHz, which can boost up to 3.9 GHz. This setup provided a robust and responsive environment for handling non-GPU-intensive tasks and managing the preprocessing and postprocessing stages of our experiments. The system’s main memory configuration included 256 GB of DDR4 RAM per node, which was crucial for supporting the high-throughput demands of data-intensive operations, particularly when dealing with large-scale datasets and complex computational models. Computational experiments were orchestrated using an internal *SLURM* (Simple Linux Utility for Resource Management) compute cluster. We configured *SLURM* to efficiently allocate resources based on the demands of queued jobs, with dynamic adjustments based on priority and current load. It should be noted that the computational resources described here sufficed for all phases of the research project; the full project did not require more compute resources than those reported for the experiments.

Hyperparameter optimization: We initiated hyperparameter optimization using a randomized grid search approach, applied specifically to the *Afriat* single-cell genomics dataset due to its moderate size, facilitating efficient and thorough exploration of the hyperparameter space. The optimization targeted two primary evaluation metrics: the maximum Mutual Information Gap (maxMIG), serving as a quantitative measure for assessing covariate disentanglement efficiency, and DEG-R^2 , focusing explicitly on reconstruction accuracy for differentially expressed genes. The optimization was conducted through the Weights & Biases (wandb) platform, leveraging *TarDis*’s integrated logging capabilities. Initial random sampling broadly explored hyperparameters across multiple architectural and training-related dimensions, including latent dimensionalities, number of neurons per layer in multilayer perceptrons (MLPs), and learning rate schedules. However, iterative evaluations revealed diminishing returns from broad exploration, leading to focused refinements in a narrower, critical subset of hyperparameters that showed higher sensitivity and impact on model performance. Specifically, hyperparameters such as `weight_decay`, which regulates the extent of L2 regularization thus affecting model complexity, and `kl_warmup_epochs`, influencing the gradual introduction of KL divergence regularization in the Variational Autoencoder (VAE), `use_layer_norm` vs `use_batch_norm`, and `dropout_rate` were found one of the most impactful ones. Additional sensitivity was observed for hyperparameters defining the strength of individual loss components, which directly modulate the trade-off between reconstruction fidelity and covariate-specific disentanglement. Subsequent to automated grid searches, selective *manual* refinements were performed *occasionally* when configurations surfaced by wandb-tracked experiments demonstrated promising but suboptimal performance, warranting fine-grained adjustments of individual component weights within the auxiliary loss. These adjustments ensured optimal disentanglement without compromising overall model robustness or reconstruction quality.

The auxiliary loss weight λ_C particularly governs the strength of the disentanglement constraints. As seen in [Figures 2](#) and [S2](#), increasing λ_C generally improves separation in the reserved latents and boosts clustering metrics like ASW. We did not observe notable overfitting or under-penalization for ranges of λ_C tested on *Afriat*, possibly because that dataset is sufficiently large and the disentanglement constraints are distributed across multiple latent subspaces. We recommend modest grid searches around the default λ_C when working with new datasets of distinctly different sizes or complexity, since an excessively large λ_C could in principle lead to diminished reconstruction if the data are too sparse or limited.

In many single-cell studies, exhaustive hyperparameter tuning is recommended but seldom practiced by the users, complicating reproducibility. By contrast, our aim was to mirror typical real-world usage: researchers often employ a single, default hyperparameter set rather than extensive re-optimizations per dataset. For *TarDis*, once the primary search on *Afriat* was complete, we retained these optimized settings as the default for all subsequent experiments, unless stated otherwise. Notably, this design choice

intentionally opts for a more conservative approach; it is likely that further data-specific tuning could yield even better performance on larger or differently structured datasets like *Braun*, *Suo* and *Norman* datasets.

Model and auxiliary loss hyperparameters:

Table: Hyperparameters for model configuration. This table lists the hyperparameters used in the model configuration, including their descriptions and assigned values.

Parameter	Description	Value
n_input	Number of input features.	
n_batch	Number of batches. If 0, no batch correction is performed.	0
n_labels	Number of labels.	0
n_hidden	Number of nodes per hidden layer. Passed into Encoder and Decoder.	512
n_latent	Dimensionality of the latent space.	$24 + 8 * N_K$
n_layers	Number of hidden layers. Passed into Encoder and Decoder.	3
n_continuous_cov	Number of continuous covariates.	0
n_cats_per_cov	A list of integers containing the number of categories for each categorical covariate.	None
dropout_rate	Dropout rate. Passed into Encoder but not Decoder.	0.25
dispersion	Flexibility of the dispersion parameter, which can be “gene”, “gene-batch”, “gene-label”, or “gene-cell”, when gene_likelihood is either nb or zinb.	“gene”
log_variational	If True, use torch.log1p on input data before encoding for numerical stability (not normalization).	True
gene_likelihood	Distribution to use for reconstruction in the generative process. (“zinb”, “nb”, “poisson”)	“nb”
latent_distribution	Distribution for the latent space. (“normal”, “ln”)	“normal”
encode_covariates	If True, covariates are concatenated to gene expression prior to passing through the encoder(s).	False
deeply_inject_covariates	If True and n_layers > 1, covariates are concatenated to the outputs of hidden layers in the encoder(s) and the decoder.	True
batch_representation	Method for encoding batch information. (“one-hot”, “embedding”)	“one-hot”
use_batch_norm	Specifies where to use torch.nn.BatchNorm1d in the model. (“encoder”, “decoder”, “none”, “both”)	None
use_layer_norm	Specifies where to use torch.nn.LayerNorm in the model. (“encoder”, “decoder”, “none”, “both”)	“both”
use_size_factor_key	If True, use the anndata.AnnData.obs column as defined by the size_factor_key parameter in the model’s setup_anndata method as the scaling factor in the mean of the conditional distribution.	False
use_observed_lib_size	If True, use the observed library size for RNA as the scaling factor in the mean of the conditional distribution.	True
library_log_means	Vector of shape (1, n_batch) of means of the log library sizes that parameterize the prior on library size.	None
library_log_vars	Vector of shape (1, n_batch) of variances of the log library sizes that parameterize the prior on library size.	None
var_activation	Callable used to ensure positivity of the variance of the variational distribution. Passed into Encoder. The default is the exponential function.	None
deeply_inject_disentangled_latents	If True, deeply inject disentangled latents.	True
include_auxillary_loss	If True, include auxiliary loss.	True
beta_kl_weight	Weight for the KL divergence term in the loss function.	0.5

Table: Hyperparameters used for optimization. It provides a comprehensive overview of the configurations necessary to monitor and enhance model performance throughout the training

Parameter	Description	Value
max_epochs	Maximum number of training epochs.	600
train_size	Proportion of data used for training.	0.8
batch_size	Number of samples per batch.	128
Parameter	Description	Value
check_val_every_n_epoch	Frequency of validation checks in epochs.	10
limit_train_batches	Fraction of training batches to use.	1.0
limit_val_batches	Fraction of validation batches to use.	1.0
learning_rate_monitor	Monitor learning rate during training.	True
early_stopping	Enable early stopping.	False
early_stopping_patience	Number of epochs with no improvement after which training will be stopped.	150
early_stopping_monitor	Metric to monitor for early stopping.	“elbo_train”
n_epochs_kl_warmup	Number of epochs for KL divergence warmup.	600
lr	Learning rate.	1e-4
weight_decay	Weight decay (L_2 penalty).	1e-4
optimizer	Optimizer to use.	“AdamW”
reduce_lr_on_plateau	Reduce learning rate when a metric has stopped improving.	True
lr_patience	Number of epochs with no improvement after which learning rate will be reduced.	100
lr_scheduler_metric	Metric to monitor for learning rate scheduler.	“elbo_train”

Table: Summary of \mathcal{L}_c configuration designed for covariates, namely status control, time, and zone in TarDis_{multiple} model trained on Afriat dataset. It provides insights into how each covariate contributes to the overall model loss.

Configuration				Auxiliary Losses			
Covariate	Res Dim	Target Type	Loss Type	Latent Group	Weight	Count Type	Opt Type
status	8	categorical	MSE	reserved	100	–	max
					10	+	min
				completely unreserved	10	–	min
					100	+	max
time	8	categorical	MSE	reserved	100	–	max
					completely unreserved	10	+
				completely unreserved		10	–
					completely unreserved	100	+
zone	8	categorical	MSE	reserved		100	–
					completely unreserved	10	+
				completely unreserved		10	–
					completely unreserved	100	+

Theoretical assumptions

- Gene Dependency: The model implicitly assumes that the expression of genes can be considered independently (conditional on the latent space and covariates) when calculating losses. However, genes often exhibit co-expression or are co-regulated, which the model might not account for without specific modifications.
- Homogeneity of Cell Populations: It’s implicitly assumed that cell populations are homogeneous within groups defined by covariates, which might not be the case in heterogeneous biological conditions such as tumors or developing tissues.
- Distribution of Gene Expression Counts: The model assumes that gene expression counts can be modeled effectively using a Negative Binomial distribution. This assumption is common but might not always capture the real variability and distribution in different types of datasets.

- **Linearity and Gaussianity of Latent Space:** The auxiliary loss assumes a Gaussian distribution for the latent vectors \mathbf{z}_{nk} . This implies assumptions about linearity and normality in the latent space, which may not hold in more complex or non-linear biological data structures. This assumption is critical for the model's simplicity and tractability:

$$\mathbf{z}_{nk} \sim \mathcal{N}(\boldsymbol{\mu}_{nk}, \boldsymbol{\Sigma}_{nk}) \quad (\text{Equation 9})$$

- **Static Covariate Definition:** The model assumes static and well-defined positive or negative sample definitions in terms of covariate values. This is critical for the stability of the training process: $\mathbf{s}_{nk}^{(k)+}$ and $\mathbf{s}_{nk}^{(k)-}$ are fixed and consistent throughout the dataset.
- **Consistency and Availability of Covariate Labels:** Consistent and accurate labeling of covariates across all cells is required. Incomplete or inaccurate labels can undermine the model's effectiveness:

$$p(\mathbf{s}_{nk} = \mathbf{s}') = 1 \quad \forall n \in N_C \quad (\text{Equation 10})$$

- **Smoothness of Latent Space:** The auxiliary loss assumes the latent space is smooth and continuous, allowing for meaningful interpolation and extrapolation:

$$\forall \mathbf{z}_{nk}, \exists \text{ continuous function } g \text{ such that } g(\mathbf{z}_{nk}) = \mathbf{x}_n \quad (\text{Equation 11})$$

- **Sensitivity to Outliers:** The model does not explicitly account for outliers, which can skew learned representations. It's assumed that:

$$p(\mathbf{x}_n \text{ is outlier}) = 0 \quad (\text{Equation 12})$$

- **Assumption of Sufficient Sample Size:** The effectiveness of the model in disentangling and accurately representing biological phenomena is contingent upon having a sufficiently large number of samples to cover the variability and complexity of the data. Small sample sizes could lead to overfitting and poor generalization to new data:

$$\min_k \left(\sum_{n \in N_C} \mathbb{I}(\mathbf{s}_{nk} = \mathbf{s}') \right) \geq \text{threshold} \quad (\text{Equation 13})$$

- **Data Sparsity:** The model assumes it can handle sparsity in single-cell genomic data without additional modifications.
- **Consistency of Environmental and Experimental Conditions:** It's assumed that all cells are subject to similar environmental and experimental conditions, aside from the controlled variations represented by covariates. Variability in these conditions could introduce unmodeled noise and bias.

Dataset insights

Please be aware that this section contains embedded hyperlinks, which are essential for accessing the referenced datasets and additional resources. For optimal functionality and ease of navigation, it is highly recommended to consult the PDF version of this document. The PDF format ensures that all hyperlinks are active and can be directly accessed, facilitating seamless retrieval of the associated data and supplementary information.

Afriat Dataset: The *Afriat* dataset, named after the first author of the study, provides high-resolution single-cell RNA sequencing and single-molecule transcript imaging data of host and parasite gene expression during the liver stage of the rodent malaria parasite *Plasmodium berghei* ANKA. It highlights spatial differences in gene expression across hepatocyte lobule zones, revealing insights into the molecular interactions between host and parasite.⁵⁶

Number of Samples: The dataset comprises 19,053 individual cells.

Number of Features: It encompasses expression profiles across 8,203 genes.

Source: The data is publicly accessible. The raw dataset can be found under GEO accession number GSE181725. Processed data are available as a Seurat object⁷⁰ at Zenodo. The AnnData⁷¹ format, utilized in this study, was downloaded from Figshare, as prepared by Biolord study.³⁵ No preprocessing or subsetting was performed on our part.

Suo Dataset: Named after a co-author of the originating study, the *Suo* dataset offers a multi-organ, single-cell transcriptomic perspective, capturing dynamic immune system developments across nine prenatal human tissues during embryonic stages. This comprehensive dataset details the temporal and spatial maturation of immune cells, highlighting embryonic developmental timing and the interaction between different organ systems in shaping the immune landscape.⁵⁷

Number of Samples: From an initial count of 908,178 individual cells, 841,922 cells met quality control standards set by established single-cell best practices.²¹

Number of Features: The dataset, which initially profiled 33,538 genes, has been refined to focus on 8,192 highly variable genes (HVGs), following established single-cell sequencing best practices.²¹

Source: The raw dataset can be found under ArrayExpress accession number E-MTAB-11343. Processed data are available in AnnData format, accessible at CellAtlas portal. Additional metadata with more detailed annotation is available through the cellxgene server.⁷² The metadata was then refined and corrected for errors by the authors.

Braun Dataset: Named for the first author, the Braun dataset provides a comprehensive single-cell transcriptomic analysis of the human brain during the crucial first trimester. Spanning 5 to 14 postconceptional weeks across 26 brain specimens, the dataset includes over 1.66 million cells dissected into 111 distinct biological samples. This extensive dataset captures the early spatial and transcriptional blueprint of brain development, with detailed insights into neuronal and glial differentiation trajectories.⁵⁸

Number of Samples: From an initial count of 1,665,937 individual cells, 1,661,498 cells met quality control standards set by established single-cell best practices.²¹

Number of Features: The dataset, which initially profiled 59,459 genes, has been refined to focus on 8,192 highly variable genes (HVGs), following established single-cell sequencing best practices.²¹

Source: Raw sequencing data are available from the European Genome Phenome Archive under the accession number EGAS00001004107). The data can be browsed interactively at SciLifeLab Portal and cellxgene server. The metadata was then refined and corrected for errors by the authors.

Miller Dataset: The Miller dataset, named after the first author of the paper, provides a detailed single-cell mRNA sequencing atlas of human lung development from 11.5 to 21 weeks, integrated with studies on homogeneous human bud tip organoid cultures. This dataset specifically investigates the role of SMAD signaling in the differentiation of bud tip progenitors into airway lineages, showcasing how in vitro conditions mirror in vivo airway structures and function. This comprehensive atlas underscores critical insights into the cellular mechanisms guiding human airway differentiation.⁵⁹

Number of Samples: From an initial count of 8443 individual cells, 7405 cells met quality control standards set by established single-cell best practices.²¹

Number of Features: The dataset, which initially profiled 36,601 genes, has been refined to focus on 8,192 highly variable genes (HVGs), following established single-cell sequencing best practices.²¹

Source: The raw scRNA-seq data associated with this study are available in the EMBL-EBI ArrayExpress database under accession number E-MTAB-8221. The metadata was then refined and corrected for errors by the authors.

Sciplex Dataset: The Sciplex dataset, derived from the sci-Plex technology using nuclear hashing, quantifies transcriptional responses to chemical perturbations at single-cell resolution. Applied to three cancer cell lines and exposing them to 188 distinct compounds, it evaluates dose-dependent effects and different drug responses. This high-throughput chemical screen profiles approximately 650,000 single-cell transcriptomes across about 5000 samples in a single experiment, revealing cellular heterogeneity in drug response, commonalities within compound families, and nuanced differences within compound types, particularly histone deacetylase inhibitors.²⁷

Number of Samples: The dataset comprises 14,811 individual cells.

Number of Features: It encompasses expression profiles across 4999 genes.

Source: Both processed and raw data are accessible via NCBI GEO under accession number GSE139944. The dataset used, in its preprocessed and subsetted format, aligns with the methodology described in the CPA paper,⁶⁸ provided courtesy of the authors of CPA. No further preprocessing or subsetting was conducted by our team.

Norman Dataset: Named for the first author, the Norman dataset leverages high-content Perturb-seq (single-cell RNA-sequencing pooled CRISPR screens) to explore cellular and organismal complexity through combinatorial gene expression. The dataset features transcriptional responses from 284 different single or double gene knockouts, allowing for the exploration of genetic interactions at scale. This includes the mapping of regulatory pathways, classification of genetic interactions such as suppressors, and the mechanistic study of synergistic effects, notably between CBL and CNN1 in erythroid differentiation.⁵⁴

Number of Samples: The dataset comprises 108,497 individual cells.

Number of Features: It encompasses expression profiles across 5000 genes.

Source: Raw data is accessible via NCBI GEO under accession number GSE133344. The dataset used, in its preprocessed and subsetted format, aligns with the methodology described in the CPA paper,⁶⁸ provided courtesy of the authors of CPA. No further preprocessing or subsetting was conducted by our team.

Disentangled Latent Space: Gene Selection Discussion

The approach employed by *TarDis* offers distinct advantages over traditional differential expression (DE) analysis by enabling a deeper, more flexible exploration of gene expression data while addressing technical and biological variability. While DE analysis remains a robust method for identifying gene-covariate associations, it is inherently limited by its reliance on predefined class labels and its susceptibility to information loss when adjusting for confounders. In contrast, *TarDis* leverages covariate-specific latent spaces to uncover nuanced biological patterns, supports causal hypothesis generation, and preserves granularity in continuous covariates. The key advantages include the following:

- Exploratory Analysis within Covariate-Specific Latent Spaces: While differential expression (DE) analysis offers robust insights, *TarDis* enables an exploratory analysis within covariate-specific latent spaces that DE analysis cannot capture.

For instance, within a covariate-specific latent, data points can be further clustered into distinct groups, each characterized by unique DE genes. This level of granularity allows for the discovery of subgroups within covariate classes that traditional DE methods, constrained by pre-defined class labels, might overlook. Hence, *TarDis* enables a more nuanced exploration of gene expression dynamics, supporting sophisticated hypotheses generation beyond standard DE analysis. For example, within a certain developmental time point, *TarDis* might identify multiple Leiden clusters within the latent space reserved for *time* for a given cell type. Each of these clusters could represent a subpopulation of cells that, although originating from a similar stage of development, are on divergent paths towards differentiating into distinct cell types. Such subclusters could be crucial for identifying transitional states that are not apparent in traditional DE analysis. This ability allows for the characterization of unique DE genes that define each cluster, providing insights into the molecular mechanisms driving differentiation. In practice, this means that *TarDis* could help researchers uncover previously unrecognized cellular transitions within a homogeneous population, such as identifying progenitor cells in early developmental stages that eventually differentiate into different neuronal subtypes. Each Leiden cluster identified by *TarDis* could potentially highlight a unique pathway of development, characterized by distinct gene expression profiles, thereby offering a more detailed and dynamic view of cellular differentiation.

- Counterfactual Inference: *TarDis* extends beyond identifying associations by enabling counterfactual reasoning within the latent space. By manipulating specific latent dimensions while holding others constant, it is possible to observe the hypothetical outcomes on gene expression, offering a powerful tool for causal inference and understanding the impact of varying covariates in isolation.
- Preservation of Continuous Covariates Beyond Discretization: Traditional DE approaches often require the discretization of continuous covariates (e.g., age, dosage), potentially leading to loss of granularity and introducing biases. In contrast, *TarDis* maintains the integrity of continuous variables through its latent representations, employing a distance-weighted loss function that accurately captures subtle biological shifts, thus providing a richer and more precise characterization of gene expression changes over continuous covariates.
- Minimizing Artifact Propagation: Standard DE analyses can inadvertently adjust for batch effects in a manner that either strips away genuine biological signals or over-corrects, leading to skewed gene expression profiles. *TarDis* addresses this by explicitly modeling technical artifacts within distinct latent dimensions, thereby preserving true biological variability in the primary latent space reserved for invariant features. This targeted segregation helps in maintaining the purity of biological insights derived from the data.
- Scalability and Applicability to Multi-Omics Data: Considering the increasing size of single-cell datasets, *TarDis*'s architecture is designed to scale linearly with the number of cells and genes, facilitating efficient processing of large-scale data. Additionally, the flexibility of the underlying variational autoencoder framework makes *TarDis* adaptable to various omics data types, providing a unified approach for complex analyses that would be cumbersome with traditional DE strategies.

QUANTIFICATION AND STATISTICAL ANALYSIS

Loss functions

Without loss of generality, various choices for the loss function are investigated, focusing on elucidating the loss incurred between the anchor point \mathbf{x}_{nk} and the positive sample $(\mathbf{x}_{nk}^{(k)})^+$. The loss between the anchor point and the negative sample $(\mathbf{x}_{nk}^{(k)})^-$ can be derived similarly, with appropriate adjustments to maximize this loss.

Mean Squared Error (MSE)

The MSE between the latent representation of the anchor \mathbf{z}_{nk} and its positive counterpart $(\mathbf{z}_{nk}^{(k)})^+$ for the k th covariate is given by:

$$(\mathcal{L}_C^{(k)})_i^+(\varphi; \mathbf{x}_n, \mathbf{s}_n) = \text{MSE}(\mathbf{z}_{nk}, (\mathbf{z}_{nk}^{(k)})^+) = \frac{1}{|\mathbf{z}_{nk}|} \sum_{j=1}^{|\mathbf{z}_{nk}|} (\mathbf{z}_{nkj} - (\mathbf{z}_{nkj}^{(k)})^+)^2 \quad (\text{Equation 14})$$

However, minimizing the L_2 distance between normal vectors from distinct multivariate normal distributions with unique diagonal covariance matrices does not inherently ensure the convergence of their distributions. While this minimization may align distribution means, it disregards differences in variances and higher-order moments essential for comprehensive distributional characterization.

Mathematically speaking, if $\mathbf{z}_{nk} \sim \mathcal{N}(\boldsymbol{\mu}_{nk}, \boldsymbol{\Sigma}_{nk})$ and $(\mathbf{z}_{nk}^{(k)})^+ \sim \mathcal{N}((\boldsymbol{\mu}_{nk}^{(k)})^+, (\boldsymbol{\Sigma}_{nk}^{(k)})^+)$, by using linearity of expectation and properties of the transpose, the expected squared L_2 distance between \mathbf{z}_{nk} and $(\mathbf{z}_{nk}^{(k)})^+$ can be simplified to:

$$\mathbb{E}[\|\mathbf{z}_{nk} - (\mathbf{z}_{nk}^{(k)})^+\|_2^2] = \mathbb{E}[\mathbf{z}_{nk}^T \mathbf{z}_{nk}] - \mathbb{E}[(\mathbf{z}_{nk})^T (\mathbf{z}_{nk}^{(k)})^+] - \mathbb{E}[\mathbb{E}[(\mathbf{z}_{nk}^{(k)})^+]^T \mathbf{z}_{nk}] + \mathbb{E}[\mathbb{E}[(\mathbf{z}_{nk}^{(k)})^+]^T (\mathbf{z}_{nk}^{(k)})^+]] \quad (\text{Equation 15})$$

For any vector \mathbf{z}_{nk} with mean $\boldsymbol{\mu}_{nk}$ and covariance $\boldsymbol{\Sigma}_{nk}$, the following identity holds:

$$\mathbb{E}[\mathbf{z}_{nk}^T \mathbf{z}_{nk}] = \text{tr}(\boldsymbol{\Sigma}_{nk}) + \boldsymbol{\mu}_{nk}^T \boldsymbol{\mu}_{nk} \quad (\text{Equation 16})$$

Applying this to $(\mathbf{z}_{nk}^{(k)})^+$ and also knowing \mathbf{z}_{nk} and $(\mathbf{z}_{nk}^{(k)})^+$ are independent, we have:

$$\mathbb{E}\left[\left((\mathbf{z}_{nk}^{(k)})^+\right)^T (\mathbf{z}_{nk}^{(k)})^+\right] = \text{tr}((\Sigma_{nk}^{(k)})^+) + ((\boldsymbol{\mu}_{nk}^{(k)})^+)^T (\boldsymbol{\mu}_{nk}^{(k)})^+ \quad (\text{Equation 17})$$

$$\mathbb{E}\left[\mathbf{z}_{nk}^T (\mathbf{z}_{nk}^{(k)})^+\right] = \boldsymbol{\mu}_{nk}^T (\boldsymbol{\mu}_{nk}^{(k)})^+ \quad (\text{Equation 18})$$

$$\mathbb{E}\left[\left((\mathbf{z}_{nk}^{(k)})^+\right)^T \mathbf{z}_{nk}\right] = ((\boldsymbol{\mu}_{nk}^{(k)})^+)^T \boldsymbol{\mu}_{nk} \quad (\text{Equation 19})$$

where $\text{tr}(\cdot)$ denotes the trace of a matrix. Substituting back, we find:

$$\mathbb{E}\left[\|\mathbf{z}_{nk} - (\mathbf{z}_{nk}^{(k)})^+\|_2^2\right] = \text{tr}(\Sigma_{nk}) + \boldsymbol{\mu}_{nk}^T \boldsymbol{\mu}_{nk} - 2\boldsymbol{\mu}_{nk}^T (\boldsymbol{\mu}_{nk}^{(k)})^+ ((\boldsymbol{\mu}_{nk}^{(k)})^+)^T (\boldsymbol{\mu}_{nk}^{(k)})^+ \quad (\text{Equation 20})$$

To simplify further, recognizing the vector identity $\|\Delta\|_2^2 = \Delta^T \Delta$ for squared terms where $\Delta = (\boldsymbol{\mu}_{nk}^{(k)})^+ - (\boldsymbol{\mu}_{nk}^{(k)})^+$:

$$\mathbb{E}\left[\|\mathbf{z}_{nk} - (\mathbf{z}_{nk}^{(k)})^+\|_2^2\right] = \text{tr}(\Sigma_{nk}) + \text{tr}((\Sigma_{nk}^{(k)})^+) + \|\Delta\|_2^2 \quad (\text{Equation 21})$$

This expression reveals that the expected squared L_2 distance depends on both the aggregate covariances and the squared difference between the means. Minimizing this distance reduces the mean disparity term $\|\Delta\|_2^2$, but does not necessarily minimize the covariance term $\text{tr}(\Sigma_{nk} + (\Sigma_{nk}^{(k)})^+)$, which reflects distributional variability. However, it is crucial to ensure the convergence of our latent representations of similar pairs across their entire characteristics. As Tong and Kobayashi⁷³ demonstrated, differences in the diagonal covariances of multivariate normal distributions can significantly influence the optimal transport cost and Wasserstein distance, even when the means are aligned. This highlights the importance of considering both mean and covariance differences for accurate distribution comparison. Consequently, we redirect our focus towards statistical metrics like KL divergence, which encompass the entire distribution and provide a more comprehensive assessment of distributional convergence.

KL Divergence

Unlike the L_2 distance, which primarily measures central tendency, the KL divergence accounts for both dispersion and correlation structure. Specifically, KL divergence is sensitive to differences in the means and covariance matrices of the distributions, offering a comprehensive measure of how well one distribution approximates another, beyond merely the distance between their centers.

To frame our problem contextually, assume we have determined the representation of a positive data point in a lower-dimensional space, i.e., $(\mathbf{z}_{nk}^{(k)})^+$ is fixed. With this in mind, we aim to represent the anchor point to reflect its partial similarity in its corresponding latent representation \mathbf{z}_{nk} . Therefore, we utilize the encoder distribution of the positive sample, $q_\varphi((\mathbf{z}_{nk}^{(k)})^+ | (\mathbf{x}_{nk}^{(k)})^+, (\mathbf{s}_{nk}^{(k)})^+) = \mathcal{N}((\boldsymbol{\mu}_{nk}^{(k)})^+, (\Sigma_{nk}^{(k)})^+)$ as the target for the current point's distribution, $q_\varphi(\mathbf{z}_{nk} | \mathbf{x}_{nk}, \mathbf{s}_{nk}) = \mathcal{N}(\boldsymbol{\mu}_{nk}, \Sigma_{nk})$ given that the gradients for the forward pass of the positive sample are not computed.

Based on the KL divergence between these two multivariate Gaussian distributions, the positive pair loss $(\mathcal{L}_C^{(k)})^+(\varphi; \mathbf{x}_n, \mathbf{s}_n) = -D_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}_{nk}, \Sigma_{nk}) || \mathcal{N}((\boldsymbol{\mu}_{nk}^{(k)})^+, (\Sigma_{nk}^{(k)})^+))$ can be calculated using a straightforward and efficient formula:

$$(\mathcal{L}_C^{(k)})^+(\varphi; \mathbf{x}_n, \mathbf{s}_n) = \frac{1}{2} \left[\text{tr}(\text{inv}((\Sigma_{nk}^{(k)})^+) \Sigma_{nk}) + ((\boldsymbol{\mu}_{nk}^{(k)})^+ - \boldsymbol{\mu}_{nk})^T \text{inv}((\Sigma_{nk}^{(k)})^+) ((\boldsymbol{\mu}_{nk}^{(k)})^+ - \boldsymbol{\mu}_{nk}) - |\mathbf{z}_{nk}| + \log \frac{|\Sigma_{nk}^{(k)}|}{|\Sigma_{nk}^{(k)}|} \right] \quad (\text{Equation 22})$$

Here, $\text{inv}(\cdot)$ stands for the inverse of a matrix, $|\cdot|$ represents the determinant of a matrix, $|\mathbf{z}_{nk}|$ is the dimensionality of the distributions, $\Sigma_{nk} = \text{diag}((\sigma_{nk1})^2, \dots, (\sigma_{nk|\mathbf{z}_{nk}})^2)$ and $(\Sigma_{nk}^{(k)})^+ = \text{diag}(((\sigma_{nk1}^{(k)})^+)^2, \dots, ((\sigma_{nk|\mathbf{z}_{nk}}^{(k)})^+)^2)$. Furthermore, the determination of the determinant for such matrices is simplified, requiring only the multiplication of their diagonal elements. Therefore, equation 22 becomes:

$$(\mathcal{L}_C^{(k)})^+(\varphi; \mathbf{x}_n, \mathbf{s}_n) = \frac{1}{2} \sum_{j=1}^{|\mathbf{z}_{nk}|} \left[\frac{(\sigma_{nkj})^2}{((\sigma_{nkj}^{(k)})^+)^2} + \frac{((\mu_{nkj}^{(k)})^+ - \mu_{nkj})^2}{((\sigma_{nkj}^{(k)})^+)^2} - 1 + 2 \log \left((\sigma_{nkj}^{(k)})^+ - 2 \log \sigma_{nkj} \right) \right] \quad (\text{Equation 23})$$

We propose summing the KL divergence over all covariates k , analogous to the total correlation (TC) in the objective function of the Relevance Factor VAE (RF-VAE).⁷⁴ This approach is designed to promote independence among latent variables. Consequently, we apply this method to the KL loss term by calculating the KL divergence between each latent representation and the standard normal distribution individually, and then summing the results.

Additionally, instead of assigning a weight to each positive pair loss function with respect to covariate k and the KL divergence between its latent representation and the prior distribution (standard normal distribution), we introduce relevance indicators, $\mathbf{r}^{(k)}$ and $\mathbf{r}_j^{(0)}$ respectively. These indicators can be learned via a variational approach. They are parameterized and updated during the training process.

$$\begin{aligned} \mathbf{r}_j^{(0)} &= \mathbf{W}_j^{(0)} \cdot \mathbf{z}_{nj} + \mathbf{b}_j^{(0)} \quad \forall j \in \{0\} \cup \mathbf{J}_k \\ \mathbf{r}^{(k)} &= \mathbf{W}^{(k)} \cdot \mathbf{z}_{nk} + \mathbf{b}^{(k)} \quad \forall k \in \mathbf{J}_k \end{aligned} \quad (\text{Equation 24})$$

Hence the primary objective function to maximize for becomes:

$$\begin{aligned} (\mathcal{L}_C)^+(\phi; \mathbf{x}_n, \mathbf{s}_n) &= \frac{1}{|\mathbf{J}_k|} \sum_{k \in \mathbf{J}_k} \left[-\mathbf{r}^{(k)} D_{\text{KL}} \left(\mathcal{N}(\boldsymbol{\mu}_{nk}, \boldsymbol{\Sigma}_{nk}) \parallel \mathcal{N}((\boldsymbol{\mu}_{nk}^{(k)})^+, (\boldsymbol{\Sigma}_{nk}^{(k)})^+) \right) \right] \\ &\quad + \frac{1}{|\mathbf{J}_k|+1} \sum_{\forall j \in \{0\} \cup \mathbf{J}_k} \left[-\mathbf{r}_j^{(0)} D_{\text{KL}} \left(\mathcal{N}(\boldsymbol{\mu}_{nj}, \boldsymbol{\Sigma}_{nj}) \parallel \mathcal{N}(\mathbf{0}, \mathbf{I}) \right) \right] \end{aligned}$$

Bhattacharyya Loss

When comparing the Bhattacharyya Loss (D_B) to the KL divergence, several key distinctions arise. KL divergence can be less effective in handling outliers and noise compared to D_B , which provides a more robust measure in noisy environments.⁷⁵ Studies have demonstrated that in high-dimensional data scenarios, D_B can outperform KL divergence in both clustering accuracy and robustness to data anomalies.⁷⁶

Incorporating D_B as a loss function offers several additional advantages. First, it has shown superior performance in distinguishing between different distributions, which is essential for effective novelty detection⁷⁷ and a key aspect of disentanglement. Disentangling different factors of variation in the data often requires a measure that can accurately differentiate between various underlying distributions. Thus, the superior performance of D_B in this regard directly supports its use in disentanglement tasks. In the domain of single-cell RNA sequencing (scRNA-seq), D_B has been successfully applied to detect fear-memory-related genes from neuronal data, demonstrating its ability to handle the high heterogeneity and dropout noise inherent in such datasets.⁷⁸ Furthermore, it has been integrated into k-means clustering, enhancing the efficiency and memory-saving capabilities for large-scale scRNA-seq data analysis.⁷⁹ D_B is also robust to outliers and noise, ensuring more reliable and consistent results, which is crucial for noisy datasets.⁸⁰ Disentangling factors of variation in noisy datasets requires a measure that can reliably handle outliers and noisy data points without compromising the integrity of the disentangled components. D_B 's robustness makes it a suitable choice for such tasks. Additionally, its symmetry and comprehensive capture of distributional differences enhance the accuracy of various analytical models.⁸¹ For disentanglement, accurately capturing and separating the underlying factors of variation in the data is essential. D_B 's mathematical properties ensure that it can provide a more precise and reliable measure of these differences, facilitating better disentanglement.

Therefore, we can write the positive pair loss utilizing D_B as follows:

$$\begin{aligned} (\mathcal{L}_C^{(k)})^+(\phi; \mathbf{x}_n, \mathbf{s}_n) &= DB(\mathbf{z}_{nk}, (\mathbf{z}_{nk}^{(k)})^+) = \left[\frac{1}{8} \left((\boldsymbol{\mu}_{nk}^{(k)})^+ - \boldsymbol{\mu}_{nk} \right)^T \left(\frac{\boldsymbol{\Sigma}_{nk} + (\boldsymbol{\Sigma}_{nk}^{(k)})^+}{2} \right)^{-1} \left((\boldsymbol{\mu}_{nk}^{(k)})^+ - \boldsymbol{\mu}_{nk} \right) \right. \\ &\quad \left. + \frac{1}{2} \ln \left(\frac{|\boldsymbol{\Sigma}_{nk} + (\boldsymbol{\Sigma}_{nk}^{(k)})^+|}{2} \right) - \frac{1}{2} \ln \left(\sqrt{|\boldsymbol{\Sigma}_{nk}|} \left| (\boldsymbol{\Sigma}_{nk}^{(k)})^+ \right| \right) \right] \\ &= \frac{1}{4} \sum_{j=1}^{|\mathbf{z}_{nk}|} \frac{\left((\boldsymbol{\mu}_{nkj}^{(k)})^+ - \boldsymbol{\mu}_{nkj} \right)^2}{(\boldsymbol{\sigma}_{nkj}^{(k)})^2 + \left((\boldsymbol{\sigma}_{nkj}^{(k)})^+ \right)^2} + \frac{1}{2} \sum_{j=1}^{|\mathbf{z}_{nk}|} \ln \left(\frac{(\boldsymbol{\sigma}_{nkj}^{(k)})^2 + \left((\boldsymbol{\sigma}_{nkj}^{(k)})^+ \right)^2}{2 \cdot \boldsymbol{\sigma}_{nkj} (\boldsymbol{\sigma}_{nkj}^{(k)})^+} \right) \end{aligned} \quad (\text{Equation 25})$$

Mahalanobis Loss

Mahalanobis Loss (D_M) is a robust metric for quantifying the distance-like measure between a point and a distribution, or between two points within a distribution-defined space. Unlike KL divergence and D_B , D_M measures the deviation of a point from the mean of a distribution and can be extended to compare the central tendencies of two distributions.

The innovative use of D_M significantly enhances data interpretation and clustering accuracy. The DR-A model, combining a VAE with a generative adversarial network (GAN) leverages D_M for dimensionality reduction, achieving superior clustering and more precise low-dimensional representations of scRNA-seq data.⁸² This precision is crucial for accurately representing covariates in lower-dimensional spaces.

The scDREAMER framework integrates D_M within an adversarial VAE to tackle skewed cell types and nested batch effects, improving batch correction and preserving biological variability across heterogeneous datasets.⁸³ Table 1 highlights that while our model excels in batch correction, there is room for improvement in biological conservation. Therefore, we can adopt D_M to measure the dissimilarity between the latent representation of the anchor point \mathbf{z}_{nk} and the respective posterior distributions $\mathbf{q}_\phi((\mathbf{z}_{nk}^{(k)})^+ | (\mathbf{x}_{nk}^{(k)})^+, (\mathbf{s}_{nk}^{(k)})^+)$ as follows:

$$(\mathcal{L}_C^{(k)})^+(\mathbf{x}_n, \mathbf{s}_n) = D_M(\mathbf{z}_{nk}, (\mathbf{z}_{nk}^{(k)})^+) = \left(\sqrt{(\mathbf{z}_{nk} - (\boldsymbol{\mu}_{nk}^{(k)})^+)^T \left((\boldsymbol{\Sigma}_{nk}^{(k)})^+ \right)^{-1} (\mathbf{z}_{nk} - (\boldsymbol{\mu}_{nk}^{(k)})^+)} \right)^2 \quad (\text{Equation 26})$$

The inverse covariance matrix computation simplifies to the reciprocal of each diagonal element, resulting in:

$$(\mathcal{L}_C^{(k)})^+ (\varphi; \mathbf{x}_n, \mathbf{s}_n) = \sum_{j=1}^{|\mathbf{z}_{nk}|} \frac{(\mathbf{z}_{nkj} - (\mu_{nkj}^{(k)})^+)^2}{((\sigma_{nkj}^{(k)})^+)^2} \quad (\text{Equation 27})$$

Minimizing D_M encourages \mathbf{z}_n and $(\mathbf{z}_{nk}^{(k)})^+$ to be located within high-probability regions of the latent space, as defined by the Gaussian distribution. The latent representation of the positive example $(\mathbf{z}_{nk}^{(k)})^+$ serves as a reference, with all adjustments made relative to the current anchor point \mathbf{z}_{nk} .

Fisher Information

Fisher information can be used to measure the amount of information that a random variable $(\mathbf{z}_{nk}^{(k)})^+$ carries about the unknown parameters μ_{nk} and Σ_{nk} of a probability distribution modeling $(\mathbf{z}_{nk}^{(k)})^+$. This measurement allows for a more precise identification of the most informative latent factors, leading to more interpretable representations. Because Fisher information is grounded in information theory, the resulting disentangled factors are often more meaningful and easier to understand, which is beneficial for tasks requiring human interpretability of covariates.⁸⁴ Representations derived using Fisher information have been shown to improve performance in downstream tasks such as classification, clustering, and anomaly detection,⁸⁵ which is the ultimate goal of learning latent representations of single-cell RNA-seq data. Therefore, in the context of VAEs, Fisher information aids in analyzing information loss during the encoding process:

$$I_{\mu_{nkj}}(\mu_{nk}, \Sigma_{nk}) = \mathbb{E}_{q_\varphi((\mathbf{z}_{nk}^{(k)})^+ | \mathbf{x}_n, \mathbf{s}_n)} \left[\left(\frac{\partial}{\partial \mu_{nkj}} \log q_\varphi((\mathbf{z}_{nk}^{(k)})^+ | \mathbf{x}_n, \mathbf{s}_n) \right)^2 \right] \quad (\text{Equation 28})$$

$$I_{\sigma_{nkj}}(\mu_{nk}, \Sigma_{nk}) = \mathbb{E}_{q_\varphi((\mathbf{z}_{nk}^{(k)})^+ | \mathbf{x}_n, \mathbf{s}_n)} \left[\left(\frac{\partial}{\partial \sigma_{nkj}} \log q_\varphi((\mathbf{z}_{nk}^{(k)})^+ | \mathbf{x}_n, \mathbf{s}_n) \right)^2 \right] \quad (\text{Equation 29})$$

$$(\mathcal{L}_C^{(k)})^+ (\varphi; \mathbf{x}_n, \mathbf{s}_n) = \sum_{j=1}^{|\mathbf{z}_{nk}|} \left[I_{\mu_{nkj}}(\mu_{nk}, \Sigma_{nk}) + I_{\sigma_{nkj}}(\mu_{nk}, \Sigma_{nk}) \right] \quad (\text{Equation 30})$$

In our case, the log-likelihood function for a single observation \mathbf{x}_n is given by:

$$\log \left(q_\varphi((\mathbf{z}_{nk}^{(k)})^+ | \mathbf{x}_n, \mathbf{s}_n) \right) = -\frac{1}{2} \left[|\mathbf{z}_{nk}| \log(2\pi) + \sum_{j=1}^{|\mathbf{z}_{nk}|} \frac{(\mathbf{z}_{nkj}^{(k)})^+ - \mu_{nkj}}{\sigma_{nkj}^2} \right] \log \left[\right] \quad (\text{Equation 31})$$

For the mean parameter μ_{nkj} :

$$\begin{aligned} I_{\mu_{nkj}}(\mu_{nk}, \Sigma_{nk}) &= \mathbb{E} \left[\frac{\partial}{\partial \mu_{nkj}} \frac{((\mathbf{z}_{nkj}^{(k)})^+ - \mu_{nkj})^2}{\sigma_{nkj}^2} \right] \\ &= \frac{2}{\sigma_{nkj}^2} \cdot ((\mathbf{z}_{nkj}^{(k)})^+ - \mu_{nkj}) \end{aligned} \quad (\text{Equation 32})$$

For the variance parameter σ_{nkj}^2 :

$$\begin{aligned} I_{\sigma_{nkj}^2}(\mu_{nk}, \Sigma_{nk}) &= \mathbb{E} \left[\frac{\partial}{\partial \sigma_{nkj}^2} \log(q_\varphi((\mathbf{z}_{nk}^{(k)})^+ | \mathbf{x}_n, \mathbf{s}_n)) \right] \\ &= 2\sigma_{nkj} - 4 \cdot \frac{((\mathbf{z}_{nkj}^{(k)})^+ - \mu_{nkj})^2}{\sigma_{nkj}^2} \end{aligned} \quad (\text{Equation 33})$$

Evaluation metrics

Average Silhouette Width

The average silhouette width (ASW)⁸⁶ evaluates clustering quality by measuring the relationship between within-cluster and between-cluster distances. ASW values range from -1 to 1, where -1 indicates misclassification, 0 indicates overlapping clusters, and 1 indicates well-separated clusters.

For each data point \mathbf{x}_n , the silhouette coefficient $s(\mathbf{x}_n)$ is calculated as:

$$s(\mathbf{x}_n) = \frac{d_{\text{inter}}(\mathbf{x}_n) - d_{\text{intra}}(\mathbf{x}_n)}{\max(d_{\text{intra}}(\mathbf{x}_n), d_{\text{inter}}(\mathbf{x}_n))} \quad (\text{Equation 34})$$

where $d_{\text{intra}}(\mathbf{x}_n)$ is the average distance from point \mathbf{x}_n to all other points within the same cluster (intra-cluster distance) and $d_{\text{inter}}(\mathbf{x}_n)$ is the minimum average distance from point \mathbf{x}_n to points in any other cluster (nearest-cluster distance). The overall ASW is the mean of the silhouette coefficients for all points in the dataset:

$$\text{ASW} = \frac{1}{N_C} \sum_{n=1}^{N_C} s(\mathbf{x}_n) \quad (\text{Equation 35})$$

where N_C is the total number of data points. ASW is particularly relevant in single-cell genomics for assessing how well cells cluster based on their gene expression profiles.⁸⁶ This metric provides an intuitive measure of clustering quality and batch mixing, crucial for understanding both biological conservation and batch effect removal. It is particularly useful in clustering-based analyses but may be sensitive to noise and outliers.

Cell Type Average Silhouette Width

Cell type average silhouette width (Cell type ASW)⁴¹ evaluates cell clustering quality in single-cell transcriptomics by measuring how well cells are grouped based on type labels. The silhouette coefficient for each cell is computed similarly to general ASW. To scale the ASW values between 0 and 1, the following transformation is applied:

$$\text{celltypeASW} = \frac{\text{ASW}_c + 1}{2} \quad (\text{Equation 36})$$

where ASW_c is the ASW computed over all cell type labels c .

Batch Average Silhouette Width

Batch average silhouette width (Batch ASW)⁴¹ assesses the quality of batch mixing in integrated datasets, which is essential in single-cell transcriptomics to ensure that technical variations do not obscure biological signals. The silhouette coefficient for each cell, based on batch labels, is computed similarly to general ASW.

To obtain a Batch ASW score between 0 and 1, the following transformation is applied for each batch label j :

$$\text{batchASW}_j = \frac{1}{|C_j|} \sum_{\mathbf{x}_n \in C_j} (1 - s_{\text{batch}}(\mathbf{x}_n)) \quad (\text{Equation 37})$$

where C_j is the set of cells with batch label j , $|C_j|$ is the size of this set, and $s_{\text{batch}}(n)$ is the silhouette coefficient for each cell n based on batch labels. The final Batch ASW score is calculated by averaging the batch ASW values across all batch labels:

$$\text{batchASW} = \frac{1}{|B|} \sum_{j \in B} \text{batchASW}_j \quad (\text{Equation 38})$$

where B is the set of unique batch labels. A Batch ASW score closer to 0 indicates good batch mixing, meaning batch effects have been effectively corrected.⁸⁷

Isolated Label F1 Score

Precision, also known as positive predictive value, gauges the proportion of correctly predicted positive instances among the total predicted positives. It's calculated by considering True Positives (TP) against False Positives (FP), following the formula:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (\text{Equation 39})$$

In contrast, Recall, also called sensitivity or true positive rate, measures how well the model identifies actual positive instances, crucial when false negatives are costly. Its calculation focuses on TP relative to FN, given by:

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (\text{Equation 40})$$

The F1 score, a harmonic mean of precision and recall, offers a single metric balancing both aspects, with high values indicating a well-balanced model. It is calculated as:

$$\text{F1} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (\text{Equation 41})$$

Isolated Label Scores are used to evaluate the clustering and separation of cell identity labels shared by a few batches. Specifically, the isolated label F1 score, also known as the class-wise F1 score, evaluates the F1 score for individual classes and is optimized to achieve the best clustering of these isolated labels, ensuring effective integration of rare cell types. This metric is particularly valuable for handling imbalanced datasets, such as those in single-cell genomics, where it assesses the accuracy of identifying rare cell types.^{41,88} The original scIB package typically employs a cluster-based F1 scoring method by default. However, for the sake of

speed and simplicity, we are opting to use the ASW instead as implemented in scib-metrics package.⁸⁹ The isolated label ASW measures the separation quality of these labels. These scores address the challenge of integrating rare cell types, ensuring that integration methods can effectively manage rare cell populations. However, the performance of these scores is heavily influenced by the quality of initial annotations.

Mutual Information

Mutual information (MI) quantifies the reduction in uncertainty about one variable given knowledge of another between variables in complex systems, making it a valuable measure in both theoretical analyses and practical applications.^{90,91} It measures the amount of information shared between two random variables \mathbf{z}_n^+ and \mathbf{z}_n^- as follows:

$$I(\mathbf{z}_n^+, \mathbf{z}_n^-) = \rho(\mathbf{z}_n^+, \mathbf{z}_n^-) \log \left(\frac{\rho(\mathbf{z}_n^+, \mathbf{z}_n^-)}{\rho(\mathbf{z}_n^+) \rho(\mathbf{z}_n^-)} \right) \quad (\text{Equation 42})$$

where $\rho(\mathbf{z}_n^+, \mathbf{z}_n^-)$ is the joint probability distribution of \mathbf{z}_n^+ and \mathbf{z}_n^- , and $\rho(\mathbf{z}_n^+)$ and $\rho(\mathbf{z}_n^-)$ are their marginal distributions.

The value of MI is non-negative, $I(\mathbf{z}_n^+, \mathbf{z}_n^-) \geq 0$, and measures the reduction in uncertainty of \mathbf{z}_n^+ given \mathbf{z}_n^- and vice versa. When $I(\mathbf{z}_n^+, \mathbf{z}_n^-) = 0$, the variables are statistically independent, meaning that knowing \mathbf{z}_n^+ does not provide any information about \mathbf{z}_n^- . A higher value of MI indicates a greater level of dependency between the variables.

Normalized Mutual Information

MI is influenced by dataset size and cluster entropy, complicating comparisons across datasets. Normalization techniques, which adjust MI to a standard range, typically $[0, 1]$, enable more equitable comparisons.

$$\text{NMI}(\mathbf{z}_n^+, \mathbf{z}_n^-) = \frac{I(\mathbf{z}_n^+, \mathbf{z}_n^-)}{\sqrt{H(\mathbf{z}_n^+)H(\mathbf{z}_n^-)}} \quad (\text{Equation 43})$$

where $H(\mathbf{z}_n^+)$ and $H(\mathbf{z}_n^-)$ are the entropies of \mathbf{z}_n^+ and \mathbf{z}_n^- . The higher values indicate superior clustering quality.⁹² In the context of single-cell genomics, the normalized mutual information (NMI) is critical for evaluating how well clusters correspond to known cell types.⁴¹ This metric evaluates how well cell-type labels are preserved post-integration. It is often used in scenarios requiring validation of clustering results against known labels. While it provides an intuitive measure, it may not distinguish well between near-perfect and perfect clustering.

In response to concerns about the relatively low ARI values (below 0.3) in Table 1, especially when contrasted with the higher NMI scores (above 0.6), which may raise questions about the models' overall performance, we offer the following clarifications:

- **Biological Relevance Despite Lower ARI:** While ARI values are influenced by cluster granularity, the key biological groups (e.g., cell types or states) remain discernible in the integrated data. For example, in our case studies (Section 3.2–3.3), marker gene analysis and UMAP visualizations confirm that major cell populations are well-separated. The lower ARI reflects technical discrepancies in cluster counts rather than a failure to capture biological signals. Thus, the integration quality is sufficient for downstream tasks like differential expression or trajectory inference, which rely on meaningful biological variation rather than strict alignment with predefined annotations.
- **Default scIB Clustering Settings:** In our workflow, we adhere to the default scIB parameters (of scib_metrics implementation⁴²) and do not manually optimize the number of clusters for each approach. Our focus is to provide a fair, consensus-oriented comparison of different data integration methods, rather than fine-tuning each model to achieve the best possible clustering outcome. Tuning these parameters individually could raise concerns of overfitting and is also *not* the primary aim of this study, which is to investigate TarDis's efficacy in disentangling covariates and preserving biological signals. Importantly, in scIB, the interest lies predominantly in the relative ranking across different models rather than absolute metric values. Consequently, our key takeaway is the comparative improvement in capturing meaningful cell-group distinctions rather than the absolute level of cluster-label alignment.
- **Clustering Metrics and Biological Interpretability:** Although ARI values near 0.3 may appear lower in comparison to NMI above 0.6, we find that such an ARI range is frequently observed when the number of predicted clusters differs significantly from the number of annotated categories in single-cell data. ARI heavily penalizes discrepancies in the cardinality of clusters, potentially magnifying even small misalignments between predicted clusters and ground-truth labels. We note that default settings in the scIB metrics pipeline often lead to higher cluster counts, which can naturally deflate ARI values.
- **Biological Validation Beyond ARI:** While ARI is a recognized external clustering metric, our analyses also involve UMAP visualizations and additional single-cell integration assessments (such as silhouette width, LISI score, and more). These complementary metrics more holistically capture how effectively a method corrects batch effects while retaining crucial biological variation.

Maximum Mutual Information Gap

The maximum mutual information gap (maxMIG) is a metric designed to evaluate the disentanglement of latent variables in complex datasets where the number of covariates exceeds two, a complexity that only particular methods are equipped to manage^{39,93–96} due to its ability to generalize and be unbiased.^{39,68,97} This measure quantifies the MI between latent representations and observed covariates, focusing on how effectively these latent variables independently capture the informative characteristics of each covariate.

The maxMIG is defined for a set of latent variables $\{\mathbf{z}_k\}_{k=1}^{N_k}$ and corresponding covariates $\{\mathbf{s}_k\}_{k=1}^{N_k}$ as:

$$\text{maxMIG}(\mathbf{z}_1, \dots, \mathbf{z}_{N_k}; \mathbf{s}_1, \dots, \mathbf{s}_{N_k}) = \frac{1}{N_k} \sum_{k=1}^{N_k} \frac{1}{H(\mathbf{s}_k)} \max_{j \neq k} [\text{MI}(\mathbf{z}_k, \mathbf{s}_k) - \text{MI}(\mathbf{z}_k, \mathbf{s}_j)] \quad (\text{Equation 44})$$

The maxMIG score is computed by averaging the normalized differences between the mutual information of each latent variable with its corresponding covariate and the highest mutual information with any other covariate. This focus on maximizing the information gap helps evaluate the specificity and relevance of each latent variable to its respective covariate. Higher maxMIG values suggest better disentanglement, indicating that each latent variable is more uniquely aligned with a specific covariate, thus enhancing the model's interpretability and generalizability.

Rand Index

The Rand index (*RI*) serves as a pivotal metric for evaluating the concordance between two clustering outcomes. It quantifies the degree of similarity by scrutinizing the allocation of data points into clusters across two distinct clustering results. Computed as the ratio of the sum of agreements to the total number of data point pairs, *RI* encapsulates both intra-cluster cohesion and inter-cluster separation. The formula for calculating the Rand Index is as follows:

$$\text{RI} = \frac{\text{TP} + \text{TN}}{\binom{N}{2}} \quad (\text{Equation 45})$$

where $N = \text{TP} + \text{TN} + \text{FP} + \text{FN}$. While the Rand Index offers valuable insights into clustering performance, it may have limitations when dealing with varying cluster sizes or datasets with an uncertain number of clusters.

Adjusted Rand Index

The *RI* quantifies the proportion of agreements between the two clusterings out of all possible pairings of elements. However, because the *RI* does not adjust for the chance grouping of elements, the Adjusted Rand Index (*ARI*)^{41,98} is often preferred, which is defined as:

$$\text{ARI} = \frac{\text{RI} - \text{Expected RI}}{\text{Max RI} - \text{Expected RI}} \quad (\text{Equation 46})$$

where the Expected *RI* is the expected value of the *RI* for random clusterings and the Max *RI* is the maximum possible value of the *RI*. Mathematically, the *ARI* can be expressed as:

$$\text{ARI} = \frac{\sum_{ij} \binom{n_{ij}}{2} - \frac{\left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right]}{\binom{n}{2}}}{\frac{1}{2} \left[\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right] - \frac{\left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right]}{\binom{n}{2}}} \quad (\text{Equation 47})$$

where n_{ij} is the number of elements in the intersection of cluster i in X and cluster j in Y , a_i is the number of elements in cluster i of X , b_j is the number of elements in cluster j of Y , and $\binom{n}{2}$ denotes the binomial coefficient. This adjustment provides a corrected-for-chance measure, making the *ARI* a more reliable metric for clustering comparison.

Values of *ARI* above zero indicate better-than-random agreement, with a value of 1 representing perfect agreement.⁹⁸ In single-cell data analysis, *ARI* is useful for validating the consistency of cell type assignments across different clustering methods. This metric is key for evaluating clustering performance in the presence of noise and is commonly used to validate clustering results in datasets with known ground truth. However, it can be less intuitive to interpret compared to simpler metrics.

k-nearest neighbor Batch Effect Test

The k-nearest neighbor batch effect test (*kBET*)^{41,99} assesses batch effects in high-dimensional datasets by testing the homogeneity of batch labels within the k-nearest neighbors of each data point. It evaluates whether the neighbors of a cell are more likely to come from the same batch than expected under random mixing. *kBET* is a robust method designed to quantify batch effects in single-cell RNA sequencing (scRNA-seq) data. To implement *kBET*, one first constructs a k-nearest-neighbor (*kNN*) graph for each cell in the dataset, using an appropriate distance metric such as Euclidean distance in a principal component analysis (*PCA*)-reduced space. For each cell n , the algorithm identifies its k nearest neighbors and calculates the proportion of cells from each batch within this neighborhood, denoted as p'_n , where j indexes the batches. Under the null hypothesis of no batch effect, the expected proportion of cells from each batch should reflect the overall batch composition in the dataset, represented as q_j . The *kBET* then compares the observed batch proportions p'_n with the expected proportions q_j using a statistical test, such as the Chi-square test or a permutation-based test. The test statistic for each cell n is computed as

$$\chi_n^2 = \sum_{j=1}^{|B|} \frac{(p_n^j - q_j)^2}{q_j}$$

where $|B|$ is the number of batches. The p-value associated with the Chi-square statistic indicates the likelihood that the observed batch composition within the neighborhood of cell n is consistent with the global batch composition. These p-values are aggregated across all cells to assess the overall presence of batch effects in the dataset. The kBET statistic is:

$$\text{kBET} = \frac{1}{N} \sum_{n=1}^N \mathbf{1}_{(p_n < \alpha)} \quad (\text{Equation 48})$$

where N is the number of neighborhoods tested, p_n is the p-value from a chi-squared test, and α is the significance threshold.

This method was evaluated using peripheral blood mononuclear cells (PBMCs) from healthy donors, effectively distinguishing cell-type-specific inter-individual variability from changes in relative proportions of cell populations. kBET is crucial for evaluating the effectiveness of batch effect correction methods in single-cell transcriptomics. The kBET tool and its detailed implementation are available on the kBET GitHub repository.

Graph Connectivity

Graph connectivity evaluates whether the kNN graph of integrated data effectively connects all cells with the same identity. For each cell identity label, a subset kNN graph is created. The graph connectivity score is then computed as the average size of the largest connected component relative to the number of nodes with that cell identity.⁴¹ This metric ensures that cells of the same type remain connected post-integration, a critical aspect for evaluating graph-based methods. Despite its importance, calculating graph connectivity can be computationally intensive for large datasets.

In single-cell genomics, graph connectivity assesses the robustness of cell interaction networks. The formula for graph connectivity is:

$$\text{Graph Connectivity} = \frac{1}{|C|} \sum_{c \in C} \frac{|LCC(G(N_c, E_c))|}{|N_c|} \quad (\text{Equation 49})$$

where C is the set of cell identity labels, $LCC(G(N_c, E_c))$ is the largest connected component of the graph for cells with label c , and $|N_c|$ is the number of nodes with cell identity c .

Coefficient of determination in VAE

The R^2 Reconstruction metric, often referred to as the coefficient of determination, is a statistical measure used to evaluate the performance of VAEs in reconstructing input data. This metric quantifies how well the reconstructed outputs from a VAE approximate the original inputs, indicating the proportion of variance in the data that is captured by the model. R^2 Reconstruction is particularly useful in the evaluation of VAEs because it provides a clear metric to gauge the accuracy of data reconstructions, facilitates comparison between different VAE architectures or configurations on the same dataset, helps identify areas where the model might be lacking, guiding further refinements. This metric is critical for researchers and practitioners using VAEs to ensure that their models not only generate new data that is statistically similar to the input data but also effectively reconstruct specific instances of input data.^{69,100}

In the context of VAEs, the R^2 Reconstruction is defined as:

$$R^2 = 1 - \frac{\sum_{n=1}^{N_C} \|\mathbf{x}_n - \hat{\mathbf{x}}_n\|^2}{\sum_{n=1}^{N_C} \|\mathbf{x}_n - \bar{\mathbf{x}}\|^2} \quad (\text{Equation 50})$$

where \mathbf{x}_n represents the original input data, $\hat{\mathbf{x}}_n$ represents the reconstructed data produced by the VAE, and $\bar{\mathbf{x}}$ is the mean of the original input data.

The R^2 value ranges from 0 to 1, where a higher value indicates that the model has effectively captured more of the variance in the input data through its reconstructions. An R^2 value of 1 signifies perfect reconstruction, whereas a value close to 0 indicates that the model performs no better than a model that would simply predict the mean of the input data for all outputs.

Coefficient of determination for Differentially Expressed Genes in VAE

In computational biology, the evaluation of VAEs reconstruction often focuses on differentially expressed genes (DEG), which show significant changes in expression under different conditions, are critical for understanding biological processes and disease mechanisms. The R^2 Reconstruction metric is adapted in this context to specifically assess how well VAEs can reconstruct the expression patterns of these DEG. Refer to the description of R^2 reconstruction score for its details.^{69,100}

The R^2 Reconstruction for DEG is defined as:

$$R_{\text{DEG}}^2 = 1 - \frac{\sum_{n=1}^{N_C} \|\mathbf{x}_n - \hat{\mathbf{x}}_n\|^2}{\sum_{n=1}^{N_C} \|\mathbf{x}_n - \bar{\mathbf{x}}_{\text{DEG}}\|^2} \quad (\text{Equation 51})$$

where \mathbf{x}_n represents the expression levels of DEG in the original data, $\hat{\mathbf{x}}_n$ represents their reconstructed levels from the VAE, and $\bar{\mathbf{x}}_{\text{DEG}}$ is the mean expression level of DEG.

Focusing on DEG, the R^2 Reconstruction metric specifically evaluates how effectively the VAE captures the variability and regulatory patterns in gene expression that are most biologically relevant and likely to be impacted by experimental conditions. A high R^2 value indicates that the VAE has effectively learned to model the critical aspects of gene expression relevant to the study's goals.

Reconstructing differentially expressed genes is inherently more difficult yet more critical than reconstructing overall gene expression due to several factors:

- (i) *Biological Relevance* DEG often carry more biological significance than stably expressed genes, directly reflecting the cellular responses to biological stimuli or disease states.
- (ii) *High Variability* DEG typically exhibit high variability in expression levels, making accurate reconstruction a complex challenge that tests the model's sensitivity and precision.
- (iii) *Data Reduction* By concentrating on DEG, researchers can reduce the dimensionality of the data, focusing computational resources and analytical efforts on the most informative parts of the dataset.
- (iv) *Improved Sensitivity* Models tuned to capture changes in DEG can be more sensitive to subtle but biologically important changes that might be overlooked when considering all genes.

Evaluating VAE performance using the R^2 Reconstruction metric on DEG provides insights into the model's ability to handle the most critical and dynamic components of biological data, facilitating the development of more accurate and biologically informative models.

Principal Component Regression

The principal component regression (PCR) quantifies batch removal by calculating the variance contribution of the batch effect per principal component (PC).⁴¹ The variance contribution of the batch effect is computed as the product of the variance explained by each PC and the corresponding R^2 value from a linear regression of the batch variable onto each PC. Mathematically, it is expressed as:

$$\text{Var}(C|B) = \sum_{g=1}^G \text{Var}(C|PC_g) \times R^2(PC_g|B) \quad (\text{Equation 52})$$

where $\text{Var}(C|PC_g)$ is the variance of the data matrix C explained by the g th principal component and $R^2(PC_g|B)$ is the coefficient of determination for the batch variable B . This metric provides a quantitative measure of batch effects, allowing for direct comparison between methods, and is essential for assessing how well integration methods remove technical variability, particularly in large-scale multi-batch studies. However, it may not fully capture non-linear batch effects.

Local Inverse Simpson's Index

The graph local inverse Simpson's index (LISI) is a metric for evaluating batch mixing (iLISI) and cell-type separation (cLISI) in integrated single-cell datasets. It uses graph-based distances and the inverse Simpson's index to measure diversity within neighborhood compositions. Scores are rescaled from 1 to the total number of batches to a range of 0 to 1, where 0 indicates minimal integration or separation, and 1 indicates optimal mixing or segregation. This metric is especially useful for graph-based integration methods and allows for cross-method comparisons, although it requires careful parameter tuning and interpretation.^{41,101}

cLISI assesses the integration of diverse cell types within a combined dataset. For each cell, its kNN are identified, and the composition of cell types within this neighborhood is analyzed. The diversity is quantified using the Inverse Simpson's Index:

$$D_{\text{cLISI}} = \frac{1}{\sum_{n=1}^{N_C} p_n^2} \quad (\text{Equation 53})$$

where p_n is the proportion of the n -th cell type in the neighborhood, and N_C is the total number of distinct cell types. The average cLISI score across all cells indicates how well cell types are mixed, with high values showing effective mixing and low values indicating poor mixing.

iLISI measures dataset mixing within the local neighborhood of each cell, quantifying how well cells from different datasets are integrated. iLISI close to the number of datasets suggests good mixing, meaning datasets are well integrated where cLISI close to 1 indicates good preservation of cell types, meaning different cell types remain well separated.

Balancing iLISI and cLISI ensures datasets are integrated effectively while preserving distinct cell type identities. Graph LISI's unified measure for both batch mixing and cell-type separation makes it a valuable tool for single-cell data integration studies, providing a standardized framework for comparing integration methods and identifying optimal strategies.

CI calculation and DE gene selection procedure

Steps and rationale for preparing [Figure 3](#).

- Gene Expression Smoothing: To reduce noise and capture underlying trends, each gene's expression profile along the continuous latent variable (PC1) was convolved with a Gaussian kernel. The kernel's window size (256 samples) and standard deviation (256) were chosen to balance resolution and smoothness, ensuring local variations are preserved while suppressing high-frequency noise. This generates a moving average trajectory for each gene.

- Confidence Interval Construction: A symmetric confidence interval (CI) around zero was derived under the null hypothesis of no differential expression. Assuming Gaussian-distributed noise in the smoothed expression values, the CI width was determined by computing the cumulative distribution function (CDF) at ± 2.25 standard deviations. This corresponds to a $\sim 97.5\%$ two-tailed probability coverage (i.e., $CI = 2 \times \Phi(2.25) - 1 \approx 0.975$), where Φ is the standard normal CDF. Genes with trajectories remaining within ± 1.9 across the latent axis were classified as non-DE (within CI).
- Threshold-Based DE Detection: Genes were deemed differentially expressed (DE) if their smoothed trajectory exceeded the ± 1.9 threshold at any point along the latent axis. This threshold corresponds to the 97.5% quantile under the null distribution, controlling the false positive rate (FPR) at $\sim 2.5\%$ per tail. The threshold was applied post-smoothing to account for multiple testing across the latent axis while maintaining sensitivity to localized expression changes.
- Statistical Rationale: The Gaussian kernel induces spatial correlation in the smoothed trajectories, implicitly modeling local dependencies in gene expression along the latent axis. By thresholding the maximum absolute deviation of the smoothed trajectory, we identify genes with sustained expression changes exceeding random fluctuations. This approach approximates a family-wise error rate (FWER) control across the latent axis, as the Gaussian smoothing reduces independent multiple comparisons.
- Empirical Validation: The threshold (1.9) was calibrated to match the 97.5% CI derived from the Gaussian null model, ensuring theoretical FPR control. Significant genes were visually highlighted (Figure 3, right), while non-DE genes were aggregated with low opacity to depict population-level variability within the CI.