
A benchmark for prediction of transcriptomic responses to chemical perturbations across cell types

Artur Szalata^{1*}, Andrew Benz^{2*}, Robrecht Cannoodt^{3,4,5}, Mauricio Cortes², Jason Fong², Sunil Kuppasani², Richard Lieberman², Tianyu Liu⁶, Javier A. Mas-Rosario², Rico Meinel⁷, Jalil Nourisa⁸, Jared Tumiel⁷, Tin M. Tunjic⁹, Mengbo Wang¹⁰, Noah Weber¹¹, Hongyu Zhao⁶, Benedict Anchang¹², Fabian J. Theis^{1†}, Malte D. Luecken^{1†}, Daniel B. Burkhardt^{2†}

¹Helmholtz Munich, ²Cellarity, ³Data Intuitive, ⁴VIB Center for Inflammation Research, ⁵Ghent University, ⁶Yale University, ⁷Retro Biosciences, ⁸Helmholtz Center Hereon, ⁹TU Vienna, ¹⁰Purdue University, ¹¹Olden Labs, ¹²NIH, *,†Equal contribution
{artur.szalata, fabian.theis, malte.luecken}@helmholtz-munich.de;
{abenz, dburkhardt}@cellarity.com

Abstract

1 Single-cell transcriptomics has revolutionized our understanding of cellular hetero-
2 geneity and drug perturbation effects. However, its high cost and the vast chemical
3 space of potential drugs present barriers to experimentally characterizing the effect
4 of chemical perturbations in all the myriad cell types of the human body. To
5 overcome these limitations, several groups have proposed using machine learning
6 methods to directly predict the effect of chemical perturbations either across cell
7 contexts or chemical space. However, advances in this field have been hindered
8 by a lack of well-designed evaluation datasets and benchmarks. To drive innova-
9 tion in perturbation modeling, the Open Problems Perturbation Prediction (OP3)
10 benchmark introduces a framework for predicting the effects of small molecule per-
11 turbations on cell type-specific gene expression. OP3 leverages the Open Problems
12 in Single-cell Analysis benchmarking infrastructure and is enabled by a new single-
13 cell perturbation dataset, encompassing 146 compounds tested on human blood
14 cells. The benchmark includes diverse data representations, evaluation metrics,
15 and winning methods from our “Single-cell perturbation prediction: generaliz-
16 ing experimental interventions to unseen contexts” competition at NeurIPS 2023.
17 We envision that the OP3 benchmark and competition will drive innovation in
18 single-cell perturbation prediction by improving the accessibility, visibility, and
19 feasibility of this challenge, thereby promoting the impact of machine learning in
20 drug discovery.

21 1 Introduction

22 Examining gene expression in individual cells via single-cell RNA sequencing (scRNA-seq) provides
23 high-resolution insights into cellular behavior within healthy and diseased tissue. One emerging
24 application of single-cell technology is to profile cells under basal and perturbed states to characterize
25 the changes in cellular states associated with chemical treatments and to associate these changes with
26 healthy or pathological tissue phenotypes [1–5]. These technologies have the potential to transform

27 how drugs are discovered and bring new therapies to patients with unmet clinical needs [6–8]. Instead
28 of focusing on single molecular targets for drug discovery, it is possible to analyze how compounds
29 influence gene expression to shift cells from diseased to healthy states. This approach holds promise
30 for treating complex diseases where single-target methods have been less effective, as it addresses the
31 interplay of multiple genes and pathways within the cell.

32 However, associating small molecules with changes in cell state is challenging. One approach is to
33 brute-force screen compounds and measure the associated changes in gene expression, as has been
34 done to discover drug candidates for heart valve disorders [9]. However, chemical space is vast. There
35 are an estimated 10^{60} drug-like molecules [10]. Compounds can also have diverse impacts on gene
36 expression across different tissues, cell types, and individuals. Moreover, scRNA-seq experiments are
37 expensive and require highly-trained technicians to run. Hence, accurate prediction of the changes in
38 gene expression induced by compounds across different chemical structures and biological contexts
39 could provide immense time and cost savings.

40 Recently, machine learning methods to predict the impact on gene expression of small molecule
41 perturbations directly from chemical structures have been proposed [11–14]. However, understanding
42 such models’ effectiveness is difficult due to a lack of independent evaluations and limited availability
43 of benchmarking datasets [15]. Indeed, most existing datasets include only a single perturbation [16],
44 a single donor, or are limited to homogeneous cancer cell lines [1, 17]. Although these studies
45 represent important contributions to the field, a rigorous, standardized benchmark is needed to assess
46 their performance in diverse cell types across a wide range of chemical perturbations.

47 Here, we introduce the Open Problems Perturbation Prediction (OP3) benchmark, which is the first
48 standardized benchmark for predicting chemical perturbation effects across cell types. It includes a
49 formalized task, an open-source benchmarking platform, and a new dataset profiling 146 chemical
50 perturbations in human peripheral blood mononuclear cells (PBMCs) from three donors. We hosted a
51 NeurIPS 2023 Competition using this benchmark, and used the learnings and proposed methods to
52 improve the benchmark. OP3 provides a continuously updated, extensible benchmark for perturbation
53 prediction, promoting translation of these methods to applied science.

54 **2 Related work**

55 This work builds on previous efforts to generate single-cell chemical perturbation datasets and
56 evaluations performed alongside method development for perturbation prediction algorithms.

57 **Chemical perturbation datasets** Recently, several large-scale datasets with drug perturbations
58 have been published. The popular sci-Plex [17] dataset profiles 188 compounds in three cancer cell
59 lines, and its recent sequel, the sci-Plex-GxE [18] dataset, profiled 22 drugs combinatorially in three
60 cancer cell lines. While these datasets feature a large number of compounds, their use of cancer cell
61 lines limits their applicability, as cancer cell lines have a number of significant deviations from human
62 tissue. These datasets also use nuclei sequencing technologies which are less sensitive and have
63 higher noise compared to whole-cell sequencing used in our study [19]. In addition, a recent pre-print
64 introduced a scRNA-seq dataset of drug-perturbed human PBMCs [20], but its lack of replicates
65 makes it difficult to disentangle technical and biological noise from the drug perturbation signal.
66 Finally, a harmonized collection of public single-cell perturbation datasets was recently published,
67 but most datasets contain only a single cell type and few perturbations with overlap across datasets,
68 making them unsuitable for our benchmarking task [15].

69 **Perturbation prediction evaluation** The task of predicting the transcriptomic effects of small
70 molecule perturbations in single-cell data has been tackled by a few machine learning models [13,
71 14, 12, 21]. However, the evaluations of these models did not include drug perturbations on primary
72 tissue, used evaluation methods that are biased toward natural transcriptional variation [22], and
73 lacked assessments of stability across replicates and batches. No independent method evaluations
74 exist to our knowledge, which is essential to fairly compare algorithm performance [23].

75 3 A living benchmark for perturbation prediction

76 To drive innovation in algorithm development for single-cell perturbation analysis, we set up the
77 OP3 benchmark, including a formalized task definition, a fit-for-purpose benchmarking dataset,
78 and computational infrastructure to support continuously-updated, community-driven benchmarking
79 (Figure 1a). We outline these features below.

80 3.1 Task overview

81 Chemical perturbations induce cell type-specific gene expression changes by interacting with target
82 proteins and altering cellular processes. For example, tamoxifen, a breast cancer drug, binds the
83 estrogen receptor and inhibits cell growth, thereby acting selectively on cells expressing the estrogen
84 receptor [24]. However, the lack of knowledge about mechanisms of action for most compounds
85 hinders predicting their effects on specific cell types.

86 The goal of this task is to leverage data about chemical perturbations in some cell types to infer their
87 impact on gene expression in other cell types. The data is a tensor with three axes: compounds,
88 cell types, and genes. Each value in this tensor is a measurement of the impact on gene expression
89 observed in a specific cell type under a specific chemical perturbation (Section 3.3). Models are
90 provided with the changes in gene expression for all cell types for a subset of compounds. The
91 remaining compounds comprise the test set. These compounds have their differential expression
92 values masked for all genes for a subset of the cell types. The target of this task is to predict these
93 masked differential expression values (Figure 1b).

94 3.2 Generating a single-cell perturbation benchmarking dataset

95 **Considerations for data set generation** We identified the following properties of an ideal dataset
96 for benchmarking small molecule perturbation prediction:

- 97 1. **Disease-relevance:** To reflect the downstream application to drug discovery, an ideal dataset
98 ought to focus on a disease-relevant biological system.
- 99 2. **Balanced cellular heterogeneity:** Cell types must exhibit distinct perturbation responses
100 but be similar enough that translating compounds' effects is tractable.
- 101 3. **Diverse perturbations:** The compounds should perturb a range of biochemical pathways.
- 102 4. **Replicates across multiple donors:** Capturing perturbation effects across multiple donors
103 enables identifying effects that are preserved across diverse donors.
- 104 5. **Positive and negative controls:** Because of the high degree of technical and biological
105 variability in gene expression measurements, positive and negative controls are essential to
106 accurately estimate the variation attributable to perturbation effects.
- 107 6. **Open access & informed consent:** To ensure open access to benchmarking data collected
108 from human donors, samples must be collected under IRB supervision. This ensures donors
109 give informed consent for public sharing of any derived data.

110 **Dataset overview** We generated a novel scRNA-seq dataset profiling 146 compounds in PBMCs to
111 provide a high-quality reference benchmark dataset for single-cell perturbation prediction (Figure 1c).
112 We also included multiome single-nucleus RNA and chromatin accessibility measurements at base-
113 line to facilitate gene regulatory network inference. This effort represents, to date, the largest drug
114 perturbation dataset on primary human tissue with donor replicates [15], and was specifically de-
115 signed to satisfy all the criteria above. First, PBMCs comprise an important subset of the human
116 immune system and play a key role in various pathologies, including cancer, autoimmune diseases,
117 immunodeficiencies, and allergies. PBMCs also contain discrete cell types (including T-cells, B-cells,
118 myeloid cells, and NK cells) that perform distinct biological functions while sharing key biological
119 pathways, making perturbation prediction in PBMCs difficult yet tractable. The compounds in this
120 dataset were selected to span a wide range of mechanisms of action. Additionally, two positive control

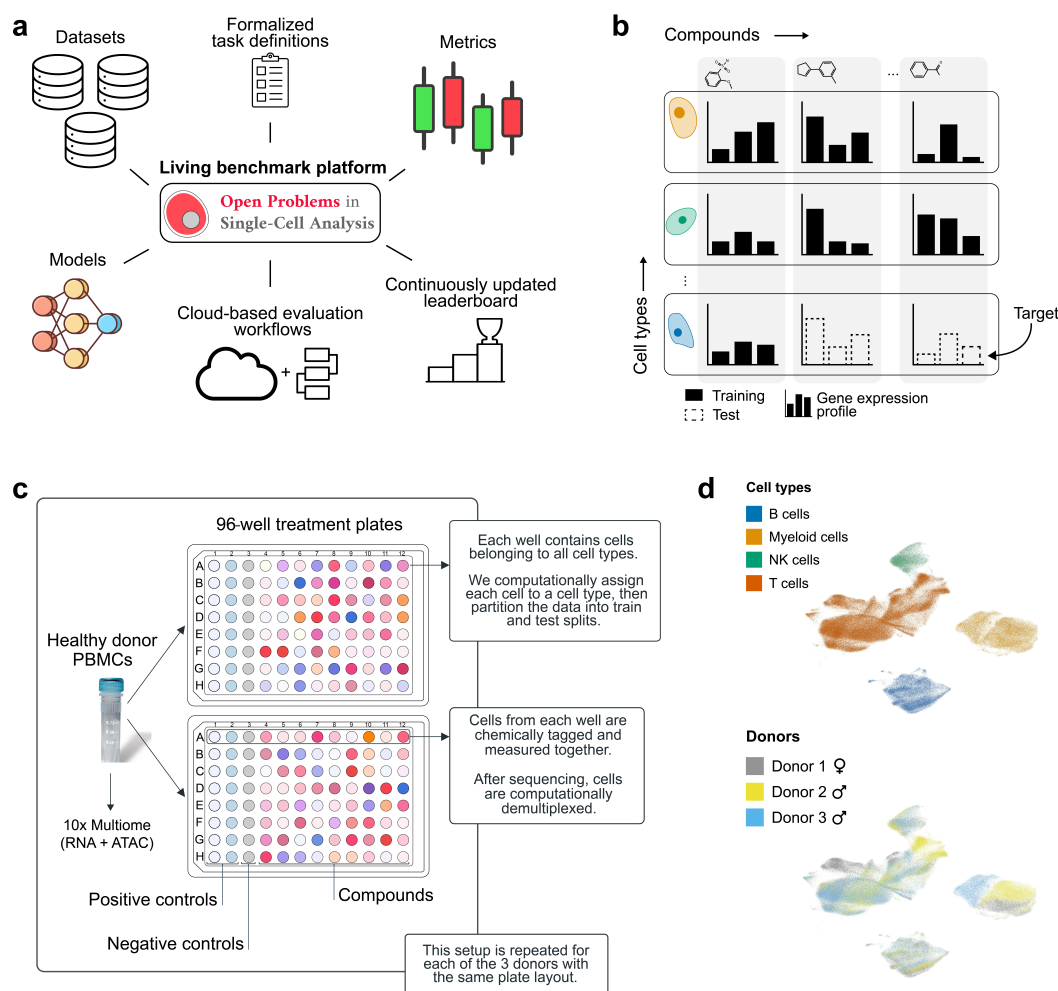


Figure 1: **Overview of the dataset.** (a) A overview of the Open Problems living benchmarking framework. (b) A graphical description of the perturbation prediction task. (c) The experimental setup for our benchmarking dataset. (d) UMAP representations of the resulting single-cell profiles colored by cell type (top) and donor (bottom).

121 compounds that were known to induce a strong transcriptional signature in PBMCs were included.
 122 Every perturbation was repeated in three healthy human donors, two male and one female. Finally,
 123 we performed this experiment using PBMCs that were commercially available with pre-obtained
 124 consent for public release.

125 **Data generation, processing and cell type annotation** PBMCs were cultured in six separate
 126 96-well plates, two for each donor (Figure 1c). After the cells were treated with compounds for 24
 127 hours, samples were collected, pooled to reduce batch effects and increase throughput, and sequenced.
 128 Sequencing reads were processed using the Cell Ranger pipeline [25], and a best-practice pipeline
 129 was followed to QC, normalize, reduce, and cluster the data [26]. We assigned each cluster to one
 130 of four cell type labels (B cell, T cell, NK cell, or myeloid cell) using established marker genes.
 131 Figure 1d shows the UMAP [27] visualization of the dataset with cell type and donor annotations.

132 The baseline multiome data (joint snRNA-seq/scATAC-seq) was processed by filtering out low-
 133 quality cells, along with both genes and chromatin accessibility features with low counts. Cells in
 134 this multiome data were then annotated based on marker gene expression in the same manner as the

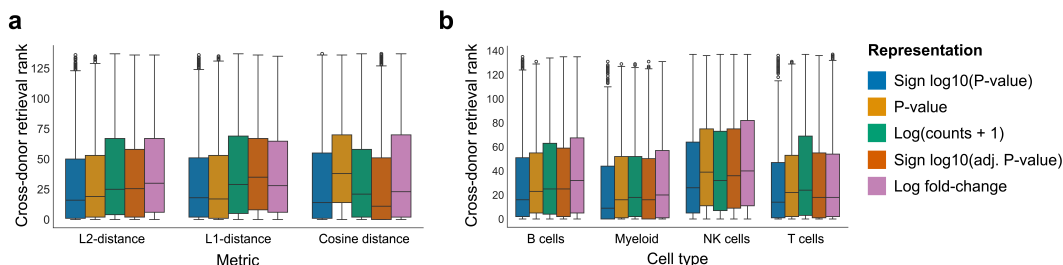


Figure 2: **Cross-donor retrieval analysis.** (a) For each pair of donors, for each compound in each cell type, the cross-donor retrieval rank was calculated using various distance metrics. The y-axis shows the retrieval rank (i.e., the rank of the same compound and cell type measurement in a different donor). The x-axis separates different retrieval distance metrics. Note that L1 distance is effectively a rescaled MAE, and L2 distance is effectively a rescaled RMSE. The hue differentiates box plots for different data representations according to the legend on the right. (b) We further examined the cross-donor retrieval rank per cell type using the L2-distance metric to ensure the results were consistent across cell types.

135 perturbational scRNA-seq data. For a detailed description of the experiment and analysis for both
 136 perturbational and baseline multiome data, please refer to **Appendix A**.

137 3.3 Representation of perturbation effects

138 In genomics, differential expression (DE) analysis is commonly used to identify how compounds
 139 affect gene activity in different cell types [26]. DE methods estimate perturbation effects by fitting
 140 generalized linear models to observed count data, explicitly accounting for biological and technical
 141 covariates. In this study, we performed DE analysis using the limma-voom framework [28], which
 142 provides estimates of effect size (e.g., log-fold change) and statistical significance while adjusting for
 143 variability associated with technical covariates.

144 Although using estimates of effect size or significance is standard in the genomics community, it is
 145 more common in machine learning benchmarks to directly predict a conditional distribution, such as
 146 the gene expression counts. To test whether the effect size (log-fold change), significance (p -values),
 147 or conditional counts are more suitable for benchmarking, we evaluated each of these representations
 148 using the replicates across donors in our dataset. We determined that an optimal representation would
 149 minimize the distance between observations of the same compound across donors, with lower median
 150 distance ranks indicating better identifiability of compounds across donors. We call this heuristic
 151 *cross-donor retrieval* (**Appendix C.1**).

152 We found that the measures of effect significance had better cross-donor retrieval (**Figure 2a** and
 153 **Appendix Figure 5**) than effect size or counts data, and this effect was consistent across cell types
 154 (**Figure 2b**). Based on these results, we decided on the following representation as a target for our
 155 benchmark: for a given compound c , cell type t , and gene g , let $p_{c,t,g}$ and $L_{c,t,g}$ be the p -value and
 156 log-fold change computed by `limma`, respectively. Then

$$\text{pert}_{c,t,g} = -\log_{10}(p_{c,t,g}) \times \text{sign}(L_{c,t,g}). \quad (1)$$

157 This representation captures both the direction and statistical significance of the perturbation effect on
 158 each gene. We do not claim that this representation is universally optimal for all tasks and analyses
 159 and note there are several challenges associated with DE analysis generally (**Appendix D**).

160 3.4 Evaluation metrics

161 We considered three metrics for evaluating model performance: mean row-wise root mean squared
 162 error (MRRMSE), mean absolute error (MAE), and cosine similarity. Mean row-wise indicates that
 163 we take a mean across predictions for compound-cell type pairs. Each of these metrics is related
 164 to the distance metrics used in the cross-donor retrieval task, e.g. MAE is effectively a rescaled

165 L1 distance, and MRRMSE is effectively a rescaled L2 distance. Using these relationships, we
166 concluded that cosine similarity had the best stability across donors, followed by MRRMSE and
167 MAE (**Figure 2a**). However, not all perturbations are expected to cause a change in gene expression,
168 and cosine similarity would not penalize models that incorrectly predict low p -values in such cases,
169 unlike MRRMSE and MAE. Hence, we primarily rely on MRRMSE for model evaluation, defined
170 as:

$$\text{MRRMSE} = \frac{1}{R} \sum_{i=1}^R \left(\frac{1}{n} \sum_{j=1}^n (y_{ij} - \hat{y}_{ij})^2 \right)^{1/2} \quad (2)$$

171 Where R is the number of (cell type, compound) tuples, and y_{ij} and \hat{y}_{ij} are the actual and predicted
172 values, respectively, and n is the number of genes.

173 3.5 Control methods

174 Including control methods in each benchmarking task is one of the basic quality controls required by
175 Open Problems not only to verify the integrity of the benchmarking workflow but to also normalize
176 the metric outputs. In this benchmark, we implemented six control methods, where each returns
177 either a solution derived from the ground truth data (positive control), a naive baseline prediction, or
178 a randomly sampled prediction (negative control). The positive and negative control methods define
179 an upper and lower bound for the performance metrics, which is used to normalize metric outputs.
180 Full descriptions of the control methods can be found in **Appendix E.2.4**.

181 4 The Single-cell Perturbation Prediction Competition at NeurIPS 2023

182 To identify the state-of-the-art for perturbation prediction in unseen cell types, we hosted a Kaggle
183 competition as part of the NeurIPS 2023 Competitions track called *Single-cell perturbation prediction:
184 generalizing experimental interventions to unseen contexts*. This competition ran from September 12,
185 2023 through November 30, 2023 and used an earlier version of the dataset and benchmark before
186 it was updated based on learnings from the competition (**Appendix B**). We ran the competition in
187 two tracks. The Leaderboard Track followed the traditional data science competition setup with a
188 public and private leaderboard tracking a single metric on public and private test sets (**Appendix B**).
189 We also ran a Judges’ Prize track where participants were judged based on a write-up addressing
190 specific questions about perturbation prediction and the specific challenges of using our dataset to
191 tackle this task. \$50,000 in prizes were awarded for each track. The competition web page with the
192 final leaderboard, code submissions, and discussions is available at: [https://www.kaggle.com/
193 competitions/open-problems-single-cell-perturbations](https://www.kaggle.com/competitions/open-problems-single-cell-perturbations)

194 4.1 Leaderboard Track

195 In the leaderboard competition, competitors trained models on the training set and submitted CSV
196 files with predictions for the public and private test tests. During the development phase (3 months),
197 only the results from the public test set were used to calculate leaderboard rankings. During the final
198 phase (5 days), competitors selected their top submission. Final scores were judged on the private
199 test set only visible after the final submission deadline. Due to the limitations of the Kaggle platform,
200 we ran the competition with a single metric, MRRMSE, decided on in collaboration with Kaggle data
201 scientists. The participants were encouraged to use any publicly available external data.

202 Over the competition, 1,318 participants from 84 countries, forming 1,097 teams, submitted 25,529
203 solutions to our Leaderboard Track. This makes our competition one of the largest single-cell data
204 science competitions to date. Although participants were only required to submit CSV predictions,
205 the Kaggle platform has a strong culture of solution sharing. As such, we were able to read through
206 reported submission code and identified trends among the best performing methods. We found
207 that the top-scoring methods relied on diverse deep learning approaches, including transformer,
208 LSTM, GRU, CNN and MLP architectures. The models used diverse loss functions, such as

209 mean squared error, mean absolute error, LogCosh ($L(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N \log(\cosh(y_i - \hat{y}_i))$), binary
210 cross-entropy, MRRMSE, and Huber loss [29]. Despite several reported attempts, only the first of
211 the three top-performing models relied on data other than the training set. The winning method
212 used ChemBERTa [30], a pre-trained transformer, to encode the small molecule structure SMILES
213 representation. According to the competitors’ reports, data preprocessing proved to be very impactful.
214 In particular, multiple competitors reported that target encoding and singular value decomposition
215 of the high-dimensional input data were effective. One method used pseudolabels [31] for model
216 training. All of the top three methods relied on model ensembles. We provide detailed descriptions of
217 these methods in **Appendix E.1**.

218 4.2 Judges’ Prize

219 In the Judges’ Prize, participants were asked to address how biological priors or alternative model
220 architectures influence leaderboard performance, to describe technical challenges that make perturba-
221 tion prediction difficult, to characterize how data noise or downsampling affect model robustness,
222 and to present well-documented and packaged model code. To identify winners, the write-ups were
223 scored by a panel of single-cell experts. 17 teams submitted write-ups for a judges’ prize, all of
224 whom also participated in the leaderboard prize.

225 Many of the submissions provided valuable insights and were exceptionally detailed—the top-
226 scoring team wrote a 33-page report. For example, several participants mentioned their efforts
227 on integrating gene regulatory networks (GRN) inferred from ATAC and RNA data as an extra
228 modality for prediction task [32, 33]. Although distinct patterns among cell types were observed
229 from the provided ATAC-seq data, attempts at incorporating inferred GRNs in model predictions,
230 even only for expression-enriched regulators, resulted in performance decreases in their models.
231 Other groups attempted to use molecular interactions as an additional modality for model design. For
232 example, GSEA-MsigDB [34] provides valuable information about pathways activated in various
233 cell types. From these, a correlation network can be constructed based on shared pathways or
234 shared regulation target genes. However, the models overall did not benefit from these efforts, which
235 suggests that further filtering over inferred regulation/correlation relationships might be necessary.
236 Finally, many submissions also investigated challenges associated with data representations and data
237 pre-processing, which are described in the following section. We provide detailed descriptions of the
238 Judges’ Prize-winning methods in **Appendix E.2**.

239 4.3 Lessons learned

240 Here, we list several key learnings and opportunities to improve our benchmarking setup.

241 **False positives for unexpressed/lowly expressed genes:** DE analysis is sensitive to low-count genes,
242 which can lead to overestimation of relative expression changes. This is especially problematic for
243 compounds with subtle gene effects. To mitigate this, we employed a stricter gene filtering strategy
244 per cell type [35], resulting in a reduced 5,317 genes (originally 18,211).

245 **Inconsistent annotations:** Proportions of T cell subtypes were inconsistent across donors (**Appendix**
246 **Table 3, Appendix Figure 7**). These subtypes had low cell counts and subtle differences in expression
247 that suggested misannotation, which may have been caused by perturbation impacts on marker gene
248 expression. To resolve this, we grouped all of the T cells together in the final annotations **Figure 1b**.

249 **Outlier samples:** Certain samples had very few cells, which may be caused by perturbation-associated
250 toxicity and was correlated with a high fraction of low p -values. To address this, we removed samples
251 with < 10 cells or inconsistent cell type proportions across donors. We also removed three compounds
252 for which we could not confidently annotate cell types (**Appendix F.2**), likely due to toxicity.

253 **Design matrix:** Due to a high number of factors and collinearity, the design matrix used in the
254 competition (**Appendix Figure 6**) was not full-rank, potentially leading to parameter estimation
255 issues. We updated the linear model to $f(g_j) = x_1 cc_i + x_2 p_i$, where g_j is a gene, cc_i is (cell type,
256 compound) tuple, and p_i is the plate. The resulting design matrix is full rank.

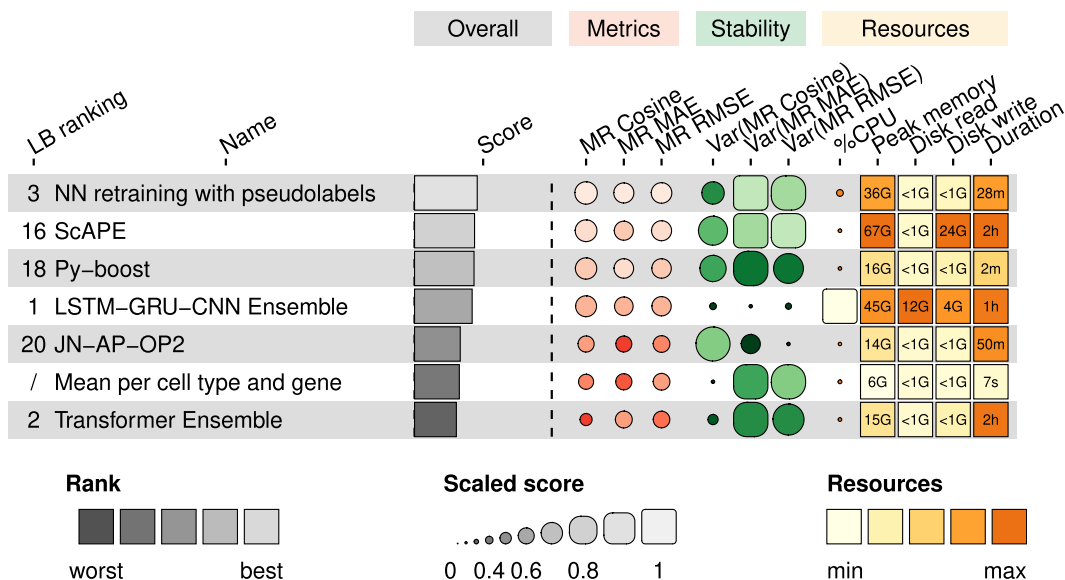


Figure 3: **An overview of the benchmarking results** of the six selected methods and one control method. Methods are ordered by the arithmetic mean of the three metrics. The MR Cosine, MR MAE, and MR RMSE were computed by comparing a method’s predictions to the ground-truth data. Each of these metric values were min-max scaled between the positive control and random sample. The resources column group shows the resource usage of the various methods throughout their execution.

257 **Outlier p -values:** Our dataset contained some very low p -values ($1e-180$). As we do not want to
 258 penalize models for not differentiating between very small p -values, we clipped p -values in the
 259 dataset at $1e-4$.

260 **Submit algorithms, not predictions:** Even though the competition participants submitted methods
 261 implementations, we were unable to exactly reproduce all of the results. We recommend requiring
 262 competitors to submit algorithms instead of predictions to promote the development of reusable
 263 tools. In addition, it allows algorithms to be more easily adapted, ultimately accelerating scientific
 264 discovery.

265 4.4 Updating the living benchmark

266 A central challenge in machine learning competitions is translating state-of-the-art methods according
 267 to competition leaderboards to impact applied science. A review of 10 years of machine learning
 268 competitions in dementia [36] found that no competition winners had been applied in clinical
 269 settings, suggested that winning methods may be overfitted to the competition dataset and metric,
 270 and suggested making methods available for testing in other settings. To enable further testing and
 271 evaluation of top-performing methods from our competition, we implemented and retrained the
 272 top 3 methods according to the leaderboard and the top 3 according to judges’ scores in our Open
 273 Problems Perturbation Prediction living benchmark. This final benchmark includes the changes
 274 listed in the preceding section. Additionally, the public test set is now part of the training set.
 275 The results are shown in **Figure 3** and the latest results of the living benchmark are available at
 276 https://openproblems.bio/results/perturbation_prediction.

277 Examining model performance across compounds, we observed that for all 6 methods, the error
 278 residuals correlated with the number of DE genes. This indicates that the methods are better at
 279 predicting no change in gene expression than a significant change. Indeed, the top performing method,
 280 NN retraining with pseudolabels, predicts high p -values more often than they occur in the dataset
 281 (Appendix Figure **Figure 8**).

282 5 Discussion

283 In this study, we presented a living benchmark for single-cell perturbation prediction. The Open
284 Problems Perturbation Prediction (OP3) benchmark features a newly generated fit-for-purpose dataset
285 that is the largest of its kind, optimized data representations and metrics, positive and negative
286 baseline methods that define performance ranges, and a cloud-based infrastructure that enables users
287 to add new methods, metrics, and datasets to the benchmark. Using this benchmarking setup we ran
288 the Single-cell Perturbation Prediction competition at NeurIPS 2023, in which over 1,300 participants
289 contributed over 25,000 method solutions to address the challenge of predicting perturbation responses
290 across drugs and cell types. This competition successfully made the topic of single-cell perturbation
291 prediction accessible to a non-specialist community (more than half of the surveyed participants never
292 worked with single-cell data **Appendix F.1**), while leveraging the expertise from this community to
293 improve upon current state-of-the-art methods (via Leaderboard Track winners) and provide feedback
294 on the task definition and implementation (via Judges’ Prize winners). To promote the translation of
295 competition outputs to domain impact in perturbation prediction, we used this competitor feedback to
296 update the OP3 benchmark and populated it with the top-performing solutions from the competition.
297 This enables methods to be further scrutinized by the community on the generalizability of their
298 performance across data contexts and metrics.

299 To power our single-cell perturbation prediction competition and benchmark, we generated the largest
300 multi-donor single-cell drug perturbation dataset on primary human tissue. However, despite profiling
301 146 drug perturbations in over 300,000 cells, the training data size is still limited from the perspective
302 of building models that generalize across drugs, donors, and cell types. There are over 16,600
303 clinical-stage drugs [37], which typically elicit heterogeneous responses across cell types [38] and
304 individuals [39]. Predicting the cell-type-specific response of a drug on an unseen individual will
305 likely require data generation efforts that are not feasible by individual groups, but rather coordinated
306 across consortia. Such efforts would also be needed to ensure aspects such as differing drug efficacy
307 across genetic backgrounds [40, 41] can be taken into account, which is not feasible with existing
308 perturbation datasets that often only profile cells from a single genotype [15]. In this context, our
309 OP3 benchmark and dataset represent a first step towards this larger goal.

310 A further limitation of our competition, and indeed most other Kaggle competitions, derives from
311 the use of a single performance metric, which is a limitation of the Kaggle platform. Goodhart’s
312 law suggests that when a performance metric becomes the optimization target, the metric ceases
313 to be a good metric [42, 43]. This phenomenon is especially challenging when the chosen metric
314 represents a proxy for good performance that is easy to evaluate during a model development loop
315 (i.e. is differentiable and quickly calculable). In our case, perturbation prediction would ideally
316 assess how well an unseen candidate drug treats a disease of unknown pathology in a particular
317 patient. To make this tractable, we instead evaluate the transcriptome response in an unseen hold-out
318 donor, cell type, and drug combination. A mitigation strategy for overfitting to this setting is to
319 define additional relevant tasks related to perturbation prediction to evaluate method performance on
320 different criteria. To promote innovation towards the overall goal of improving perturbation prediction,
321 we specifically enable such a multi-task evaluation setup via the OP3 living benchmark and the design
322 of our dataset. To promote generalizability of developed solutions [44], future competitions in this
323 direction may further include orthogonal readouts, such as cell type proportions, rates of cell death,
324 or inflammation [45].

325 Taken together, the OP3 benchmark and corresponding competition represent the first community-
326 extensible standard for predicting perturbation responses from single-cell transcriptomic data. While
327 several algorithms existed for this task also prior to our competition, the competition has been
328 successful in greatly expanding the set of possible solutions available, which can be further scrutinized
329 via the OP3 living benchmark. Indeed, the combination of a large-scale competition and a cloud-
330 based living benchmark represents a promising approach to promoting innovation towards critical
331 domain-specific challenges. We envision that the OP3 benchmark will lay the groundwork for further
332 method development for this question, which is of critical importance to realize the promise of
333 personalized medicine and optimized drug discovery.

334 Acknowledgments and Disclosure of Funding

335 **Acknowledgements** We would like to thank Lijun Zhao, Roman Montez, Nicole Robichaud, Nina
336 Colon, Sakina Saif, Laura Isacco, and Cameron Reilly at Cellarity for their contributions in generating
337 the single-cell RNA sequencing data used in the publication. We also thank Yuge Ji for proofreading
338 the manuscript. Saturn Cloud donated compute to support analysis of the winning methods.

339 **Funding** This work was supported by funds from the Chan Zuckerberg Initiative, Cellarity Inc., the
340 Helmholtz Association and Helmholtz Munich. This work was co-funded by the European Union
341 (ERC, DeepCell -101054957, to A.S. and F.J.T.). Views and opinions expressed are however those of
342 the authors only and do not necessarily reflect those of the European Union or the European Research
343 Council. Neither the European Union nor the granting authority can be held responsible for them.

344 **Competing interests** A.B., M.C., J.F., S.K., R.L., J.M-R., D.B. are paid employees of and have
345 equity interest in Cellarity Inc. N.W. is a paid employee of and has equity interest in Olden Labs PBC.
346 B.A. works for the US government. A.S. consults for Exvivo Labs Inc. R.M., J.T. are paid employees
347 of and have equity interests in Retro Biosciences. F.J.T. consults for Immunai Inc., Singularity Bio
348 B.V., CytoReason Ltd, Cellarity, and has ownership interest in Dermagnostix GmbH and Cellarity.
349 M.D.L. contracted for the Chan Zuckerberg Initiative and received speaker fees from Pfizer and
350 Janssen Pharmaceuticals.

351 References

- 352 [1] Thomas M. Norman, Max A. Horlbeck, Joseph M. Replogle, Alex Y. Ge, Albert Xu, Marco
353 Jost, Luke A. Gilbert, and Jonathan S. Weissman. Exploring genetic interaction manifolds
354 constructed from rich single-cell phenotypes. *Science*, 365(6455):786–793, August 2019.
355 doi: 10.1126/science.aax4438. URL <https://www.science.org/doi/10.1126/science.aax4438>. Publisher: American Association for the Advancement of Science.
- 357 [2] Jonas Simon Fleck, Sophie Martina Johanna Jansen, Damian Wollny, Fides Zenk, Makiko
358 Seimiya, Akanksha Jain, Ryoko Okamoto, Malgorzata Santel, Zhisong He, J. Gray Camp, and
359 Barbara Treutlein. Inferring and perturbing cell fate regulomes in human brain organoids. *Nature*,
360 621(7978):365–372, September 2023. ISSN 1476-4687. doi: 10.1038/s41586-022-05279-8.
361 URL <https://www.nature.com/articles/s41586-022-05279-8>. Publisher: Nature
362 Publishing Group.
- 363 [3] Marco Jost, Amy N Jacobson, Jeffrey A Hussmann, Giana Cirolia, Michael A Fischbach,
364 and Jonathan S Weissman. CRISPR-based functional genomics in human dendritic cells.
365 *eLife*, 10:e65856, April 2021. ISSN 2050-084X. doi: 10.7554/eLife.65856. URL <https://doi.org/10.7554/eLife.65856>. Publisher: eLife Sciences Publications, Ltd.
- 367 [4] Joseph M. Replogle, Reuben A. Saunders, Angela N. Pogson, Jeffrey A. Hussmann, Alexander
368 Lenail, Alina Guna, Lauren Mascibroda, Eric J. Wagner, Karen Adelman, Gila Lithwick-
369 Yanai, Nika Iremadze, Florian Oberstrass, Doron Lipson, Jessica L. Bonnar, Marco Jost,
370 Thomas M. Norman, and Jonathan S. Weissman. Mapping information-rich genotype-phenotype
371 landscapes with genome-scale Perturb-seq. *Cell*, 185(14):2559–2575.e28, July 2022. ISSN
372 0092-8674, 1097-4172. doi: 10.1016/j.cell.2022.05.013. URL [https://www.cell.com/cell/abstract/S0092-8674\(22\)00597-9](https://www.cell.com/cell/abstract/S0092-8674(22)00597-9). Publisher: Elsevier.
- 374 [5] Atray Dixit, Oren Parnas, Biyu Li, Jenny Chen, Charles P. Fulco, Livnat Jerby-Aron, Ne-
375 manja D. Marjanovic, Danielle Dionne, Tyler Burks, Raktima Raychowdhury, Britt Adam-
376 son, Thomas M. Norman, Eric S. Lander, Jonathan S. Weissman, Nir Friedman, and Aviv
377 Regev. Perturb-Seq: Dissecting Molecular Circuits with Scalable Single-Cell RNA Profiling of
378 Pooled Genetic Screens. *Cell*, 167(7):1853–1866.e17, December 2016. ISSN 0092-8674. doi:
379 10.1016/j.cell.2016.11.038. URL <https://www.sciencedirect.com/science/article/pii/S0092867416316105>.

- 381 [6] Bence Szalai and Dániel V. Veres. Application of perturbation gene expression profiles in drug
382 discovery—From mechanism of action to quantitative modelling. *Frontiers in Systems Biology*,
383 3, February 2023. ISSN 2674-0702. doi: 10.3389/fsysb.2023.1126044. URL [https://www.
384 frontiersin.org/articles/10.3389/fsysb.2023.1126044](https://www.frontiersin.org/articles/10.3389/fsysb.2023.1126044). Publisher: Frontiers.
- 385 [7] Bram Van de Sande, Joon Sang Lee, Euphemia Mutasa-Gottgens, Bart Naughton, Wendi Bacon,
386 Jonathan Manning, Yong Wang, Jack Pollard, Melissa Mendez, Jon Hill, Namit Kumar, Xiao-
387 hong Cao, Xiao Chen, Mugdha Khaladkar, Ji Wen, Andrew Leach, and Edgardo Ferran. Applica-
388 tions of single-cell RNA sequencing in drug discovery and development. *Nature Reviews Drug
389 Discovery*, 22(6):496–520, June 2023. ISSN 1474-1784. doi: 10.1038/s41573-023-00688-4.
390 URL <https://www.nature.com/articles/s41573-023-00688-4>. Publisher: Nature
391 Publishing Group.
- 392 [8] Ido Yofe, Rony Dahan, and Ido Amit. Single-cell genomic approaches for developing the
393 next generation of immunotherapies. *Nature Medicine*, 26(2):171–177, February 2020. ISSN
394 1546-170X. doi: 10.1038/s41591-019-0736-4. URL [https://www.nature.com/articles/
395 s41591-019-0736-4](https://www.nature.com/articles/s41591-019-0736-4). Publisher: Nature Publishing Group.
- 396 [9] Christina V. Theodoris, Ping Zhou, Lei Liu, Yu Zhang, Tomohiro Nishino, Yu Huang, Aleksan-
397 dra Kostina, Sanjeev S. Ranade, Casey A. Gifford, Vladimir Uspensky, Anna Malaschicheva,
398 Sheng Ding, and Deepak Srivastava. Network-based screen in iPSC-derived cells reveals
399 therapeutic candidate for heart valve disease. *Science (New York, N.Y.)*, 371(6530), February
400 2021. doi: 10.1126/science.abd0724.
- 401 [10] Jean-Louis Reymond. The Chemical Space Project. *Acc. Chem. Res.*, 48(3):722–730, March
402 2015. ISSN 0001-4842. doi: 10.1021/ar500432k.
- 403 [11] Mohammad Lotfollahi, F. Alexander Wolf, and Fabian J. Theis. scGen predicts single-
404 cell perturbation responses. *Nature Methods*, 16(8):715–721, August 2019. ISSN 1548-
405 7105. doi: 10.1038/s41592-019-0494-8. URL [https://www.nature.com/articles/
406 s41592-019-0494-8](https://www.nature.com/articles/s41592-019-0494-8). Publisher: Nature Publishing Group.
- 407 [12] Mohammad Lotfollahi, Anna Klimovskaia Susmelj, Carlo De Donno, Leon Hetzel, Yuge Ji,
408 Ignacio L Ibarra, Sanjay R Srivatsan, Mohsen Naghipourfar, Riza M Daza, Beth Martin, Jay
409 Shendure, Jose L McFaline-Figueroa, Pierre Boyeau, F Alexander Wolf, Nafissa Yakubova,
410 Stephan Günemann, Cole Trapnell, David Lopez-Paz, and Fabian J Theis. Predicting cellular
411 responses to complex perturbations in high-throughput screens. *Molecular Systems Biology*,
412 19(6):e11517, June 2023. ISSN 1744-4292. doi: 10.15252/msb.202211517. URL [https://
413 www.embopress.org/doi/full/10.15252/msb.202211517](https://www.embopress.org/doi/full/10.15252/msb.202211517). Publisher: John Wiley &
414 Sons, Ltd.
- 415 [13] Leon Hetzel, Simon Boehm, Niki Kilbertus, Stephan Günemann, Mohammad Lotfollahi, and
416 Fabian Theis. Predicting Cellular Responses to Novel Drug Perturbations at a Single-Cell
417 Resolution. *Advances in Neural Information Processing Systems*, 35:26711–26722, Decem-
418 ber 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/hash/
419 aa933b5abc1be30baece1d230ec575a7-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/aa933b5abc1be30baece1d230ec575a7-Abstract-Conference.html).
- 420 [14] Zoe Piran, Niv Cohen, Yedid Hoshen, and Mor Nitzan. Disentanglement of single-cell data
421 with biolord. *Nature Biotechnology*, pages 1–6, January 2024. ISSN 1546-1696. doi: 10.1038/
422 s41587-023-02079-x. URL <https://www.nature.com/articles/s41587-023-02079-x>.
423 Publisher: Nature Publishing Group.
- 424 [15] Stefan Peidli, Tessa D. Green, Ciyue Shen, Torsten Gross, Joseph Min, Samuele Garda, Bo Yuan,
425 Linus J. Schumacher, Jake P. Taylor-King, Debora S. Marks, Augustin Luna, Nils Blüthgen,
426 and Chris Sander. scPerturb: harmonized single-cell perturbation data. *Nat. Methods*, 21(3):
427 531–540, March 2024. ISSN 1548-7105. doi: 10.1038/s41592-023-02144-y.

- 428 [16] Hyun Min Kang, Meena Subramaniam, Sasha Targ, Michelle Nguyen, Lenka Maliskova,
429 Elizabeth McCarthy, Eunice Wan, Simon Wong, Lauren Byrnes, Cristina M. Lanata, Rachel E.
430 Gate, Sara Mostafavi, Alexander Marson, Noah Zaitlen, Lindsey A. Criswell, and Chun Jimmie
431 Ye. Multiplexed droplet single-cell RNA-sequencing using natural genetic variation. *Nature*
432 *Biotechnology*, 36(1):89–94, January 2018. ISSN 1546-1696. doi: 10.1038/nbt.4042. URL
433 <https://www.nature.com/articles/nbt.4042>. Publisher: Nature Publishing Group.
- 434 [17] Sanjay R. Srivatsan, José L. McFaline-Figueroa, Vijay Ramani, Lauren Saunders, Junyue
435 Cao, Jonathan Packer, Hannah A. Pliner, Dana L. Jackson, Riza M. Daza, Lena Chris-
436 tiansen, Fan Zhang, Frank Steemers, Jay Shendure, and Cole Trapnell. Massively multiplex
437 chemical transcriptomics at single-cell resolution. *Science*, 367(6473):45–51, January 2020.
438 doi: 10.1126/science.aax6234. URL <https://www.science.org/doi/10.1126/science.aax6234>.
439 Publisher: American Association for the Advancement of Science.
- 440 [18] José L. McFaline-Figueroa, Sanjay Srivatsan, Andrew J. Hill, Molly Gasperini, Dana L. Jack-
441 son, Lauren Saunders, Silvia Domcke, Samuel G. Regalado, Paul Lazarchuck, Sarai Alvarez,
442 Raymond J. Monnat, Jay Shendure, and Cole Trapnell. Multiplex single-cell chemical ge-
443 nomics reveals the kinase dependence of the response to targeted therapy. *Cell Genomics*,
444 4(2), February 2024. ISSN 2666-979X. doi: 10.1016/j.xgen.2023.100487. URL [https://www.cell.com/cell-genomics/abstract/S2666-979X\(23\)00339-7](https://www.cell.com/cell-genomics/abstract/S2666-979X(23)00339-7).
445 Publisher: Elsevier.
446
- 447 [19] Jiarui Ding, Xian Adiconis, Sean K. Simmons, Monika S. Kowalczyk, Cynthia C. Hession,
448 Nemanja D. Marjanovic, Travis K. Hughes, Marc H. Wadsworth, Tyler Burks, Lan T. Nguyen,
449 John Y. H. Kwon, Boaz Barak, William Ge, Amanda J. Kedaigle, Shaina Carroll, Shuqiang
450 Li, Nir Hacohen, Orit Rozenblatt-Rosen, Alex K. Shalek, Alexandra-Chloé Villani, Aviv
451 Regev, and Joshua Z. Levin. Systematic comparison of single-cell and single-nucleus RNA-
452 sequencing methods. *Nat. Biotechnol.*, 38:737–746, June 2020. ISSN 1546-1696. doi: 10.1038/
453 s41587-020-0465-8.
- 454 [20] Jialong Jiang, Sisi Chen, and Matt Thomson. D-SPIN constructs gene regulatory network models
455 from multiplexed scRNA-seq data revealing organizing principles of cellular perturbation
456 response, May 2023. URL <https://data.caltech.edu/records/2cjsj-wgh69>.
- 457 [21] Charlotte Bunne, Andreas Krause, and Marco Cuturi. Supervised Training of Conditional
458 Monge Maps. *Advances in Neural Information Processing Systems*, 35:6859–6872, Decem-
459 ber 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/hash/2d880acd7b31e25d45097455c8e8257f-Abstract-Conference.html.
- 461 [22] Haotian Cui, Chloe Wang, Hassaan Maan, Kuan Pang, Fengning Luo, Nan Duan, and Bo Wang.
462 scGPT: toward building a foundation model for single-cell multi-omics using generative AI. *Nature*
463 *Methods*, pages 1–11, February 2024. ISSN 1548-7105. doi: 10.1038/s41592-024-02201-0.
464 URL <https://www.nature.com/articles/s41592-024-02201-0>. Publisher: Nature
465 Publishing Group.
- 466 [23] Kevin Musgrave, Serge Belongie, and Ser-Nam Lim. A Metric Learning Reality Check. *arXiv*,
467 March 2020. doi: 10.48550/arXiv.2003.08505.
- 468 [24] PubChem. Tamoxifen, June 2024. URL <https://pubchem.ncbi.nlm.nih.gov/compound/Tamoxifen>. [Online; accessed 1. Jun. 2024].
- 470 [25] Grace X. Y. Zheng, Jessica M. Terry, Phillip Belgrader, Paul Ryvkin, Zachary W. Bent, Ryan
471 Wilson, Solongo B. Ziraldo, Tobias D. Wheeler, Geoff P. McDermott, Junjie Zhu, Mark T.
472 Gregory, Joe Shuga, Luz Montesclaros, Jason G. Underwood, Donald A. Masquelier, Stefanie Y.
473 Nishimura, Michael Schnall-Levin, Paul W. Wyatt, Christopher M. Hindson, Rajiv Bharadwaj,
474 Alexander Wong, Kevin D. Ness, Lan W. Beppu, H. Joachim Deeg, Christopher McFarland,
475 Keith R. Loeb, William J. Valente, Nolan G. Ericson, Emily A. Stevens, Jerald P. Radich,

- 476 Tarjei S. Mikkelsen, Benjamin J. Hindson, and Jason H. Bielas. Massively parallel digital
477 transcriptional profiling of single cells. *Nat. Commun.*, 8(14049):1–12, January 2017. ISSN
478 2041-1723. doi: 10.1038/ncomms14049.
- 479 [26] Lukas Heumos, Anna C. Schaar, Christopher Lance, Anastasia Litinetskaya, Felix Drost, Luke
480 Zappia, Malte D. Lücken, Daniel C. Strobl, Juan Henao, Fabiola Curion, Herbert B. Schiller,
481 and Fabian J. Theis. Best practices for single-cell analysis across modalities. *Nat. Rev. Genet.*,
482 24:550–572, August 2023. ISSN 1471-0064. doi: 10.1038/s41576-023-00586-w.
- 483 [27] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation
484 and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- 485 [28] Charity W. Law, Monther Alhamdoosh, Shian Su, Xueyi Dong, Luyi Tian, Gordon K. Smyth,
486 and Matthew E. Ritchie. RNA-seq analysis is easy as 1-2-3 with limma, Glimma and edgeR.
487 *F1000Research*, 5, 2016. doi: 10.12688/f1000research.9005.3.
- 488 [29] Peter J. Huber. Robust Estimation of a Location Parameter. In *Breakthroughs in Statistics*, pages
489 492–518. Springer, New York, NY, New York, NY, USA, 1992. ISBN 978-1-4612-4380-9. doi:
490 10.1007/978-1-4612-4380-9_35.
- 491 [30] Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. ChemBERTa: Large-Scale
492 Self-Supervised Pretraining for Molecular Property Prediction, October 2020. URL [http://](http://arxiv.org/abs/2010.09885)
493 arxiv.org/abs/2010.09885. arXiv:2010.09885 [physics, q-bio].
- 494 [31] Dong-Hyun Lee. Pseudo-label: The simple and efficient semi-supervised learning method
495 for deep neural networks. In *Workshop on challenges in representation learning, ICML*,
496 volume 3, page 896. Atlanta, 2013. URL [https://www.kaggle.com/blobs/download/](https://www.kaggle.com/blobs/download/forum-message-attachment-files/746/pseudo_label_final.pdf)
497 [forum-message-attachment-files/746/pseudo_label_final.pdf](https://www.kaggle.com/blobs/download/forum-message-attachment-files/746/pseudo_label_final.pdf). Issue: 2.
- 498 [32] Carmen Bravo González-Blas, Seppe De Winter, Gert Hulselmans, Nikolai Hecker, Irina
499 Matetovici, Valerie Christiaens, Suresh Poovathingal, Jasper Wouters, Sara Aibar, and Stein
500 Aerts. Scenic+: single-cell multiomic inference of enhancers and gene regulatory networks.
501 *Nature methods*, 20(9):1355–1367, 2023.
- 502 [33] Kenji Kamimoto, Christy M Hoffmann, and Samantha A Morris. Celloracle: Dissecting cell
503 identity via network inference and in silico gene perturbation. *BioRxiv*, pages 2020–02, 2020.
- 504 [34] Anthony S Castanza, Jill M Recla, David Eby, Helga Thorvaldsdóttir, Carol J Bult, and Jill P
505 Mesirov. Extending support for mouse data in the molecular signatures database (msigdb).
506 *Nature methods*, 20(11):1619–1620, 2023.
- 507 [35] Yunshun Chen, Aaron T. L. Lun, and Gordon K. Smyth. From reads to genes to pathways:
508 differential expression analysis of RNA-Seq experiments using Rsubread and the edgeR quasi-
509 likelihood pipeline, August 2016. URL <https://f1000research.com/articles/5-1438>.
- 510 [36] Esther E. Bron, Stefan Klein, Annika Reinke, Janne M. Pappas, Lena Maier-Hein, Daniel C.
511 Alexander, and Neil P. Oxtoby. Ten years of image analysis and machine learning competitions
512 in dementia. *Neuroimage*, 253:119083, June 2022. ISSN 1053-8119. doi: 10.1016/j.neuroimage.
513 2022.119083.
- 514 [37] Craig Knox, Mike Wilson, Christen M. Klinger, Mark Franklin, Eponine Oler, Alex Wilson,
515 Allison Pon, Jordan Cox, Na Eun (Lucy) Chin, Seth A. Strawbridge, Marysol Garcia-Patino,
516 Ray Kruger, Aadhavya Sivakumaran, Selena Sanford, Rahil Doshi, Nitya Khetarpal, Omolola
517 Fatokun, Daphnee Doucet, Ashley Zubkowski, Dorsa Yahya Rayat, Hayley Jackson, Karxena
518 Harford, Afia Anjum, Mahi Zakir, Fei Wang, Siyang Tian, Brian Lee, Jaanus Liigand, Harrison
519 Peters, Ruo Qi (Rachel) Wang, Tue Nguyen, Denise So, Matthew Sharp, Rodolfo da Silva,
520 Cyrella Gabriel, Joshua Scantlebury, Marissa Jasinski, David Ackerman, Timothy Jewison, Tan-
521 vir Sajed, Vasuk Gautam, and David S. Wishart. DrugBank 6.0: the DrugBank Knowledgebase

- 522 for 2024. *Nucleic Acids Res.*, 52(D1):D1265–D1275, January 2024. ISSN 0305-1048. doi:
523 10.1093/nar/gkad976.
- 524 [38] Mohammad Lotfollahi, F. Alexander Wolf, and Fabian J. Theis. scGen predicts single-cell
525 perturbation responses. *Nat. Methods*, 16:715–721, August 2019. ISSN 1548-7105. doi:
526 10.1038/s41592-019-0494-8.
- 527 [39] James F. Wilson, Michael E. Weale, Alice C. Smith, Fiona Gratrix, Benjamin Fletcher, Mark G.
528 Thomas, Neil Bradman, and David B. Goldstein. Population genetic structure of variable drug
529 response. *Nat. Genet.*, 29:265–269, November 2001. ISSN 1546-1718. doi: 10.1038/ng761.
- 530 [40] Wenndy Hernandez, Keith Danahey, Xun Pei, Kiang-Teck J. Yeo, Edward Leung, Samuel L.
531 Volchenbom, Mark J. Ratain, David O. Meltzer, Barbara E. Stranger, Minoli A. Perera, and
532 Peter H. O’Donnell. Pharmacogenomic genotypes define genetic ancestry in patients and enable
533 population-specific genomic implementation. *Pharmacogenomics J.*, 20:126–135, February
534 2020. ISSN 1473-1150. doi: 10.1038/s41397-019-0095-z.
- 535 [41] Polina V. Rusina, Maria J. Falaguera, Juan Maria R. Romero, Ellen M. McDonagh, Ian Dunham,
536 and David Ochoa. Genetic support for FDA-approved drugs over the past decade. *Nat. Rev.
537 Drug Discovery*, 22:864, October 2023. doi: 10.1038/d41573-023-00158-x.
- 538 [42] David Manheim and Scott Garrabrant. Categorizing Variants of Goodhart’s Law. *arXiv*, March
539 2018. doi: 10.48550/arXiv.1803.04585.
- 540 [43] Damien Teney, Ehsan Abbasnejad, Kushal Kaffle, Robik Shrestha, Christopher Kanan,
541 and Anton van den Hengel. On the Value of Out-of-Distribution Testing: An Ex-
542 ample of Goodhart’s Law. *Advances in Neural Information Processing Systems*,
543 33:407–417, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/
544 045117b0e0a11a242b9765e79cbf113f-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/045117b0e0a11a242b9765e79cbf113f-Abstract.html).
- 545 [44] Jack W. Scannell, James Bosley, John A. Hickman, Gerard R. Dawson, Hubert Truebel, Guil-
546 herme S. Ferreira, Duncan Richards, and J. Mark Treherne. Predictive validity in drug discovery:
547 what it is, why it matters and how to improve it. *Nat. Rev. Drug Discovery*, 21:915–931, De-
548 cember 2022. ISSN 1474-1784. doi: 10.1038/s41573-022-00552-x.
- 549 [45] Anika Liu, Srijit Seal, Hongbin Yang, and Andreas Bender. Using chemical and biological
550 data to predict drug toxicity. *SLAS Discovery*, 28(3):53–64, April 2023. ISSN 2472-5552. doi:
551 10.1016/j.slasd.2022.12.003.
- 552 [46] F. Alexander Wolf, Philipp Angerer, and Fabian J. Theis. SCANPY: large-scale single-cell gene
553 expression data analysis. *Genome Biol.*, 19(1):1–5, December 2018. ISSN 1474-760X. doi:
554 10.1186/s13059-017-1382-0.
- 555 [47] Terms and Definitions – ENCODE, June 2024. URL [https://www.encodeproject.org/
556 data-standards/terms](https://www.encodeproject.org/data-standards/terms). [Online; accessed 10. Jun. 2024].
- 557 [48] Chuan Xu, Martin Prete, Simone Webb, Laura Jardine, Benjamin J. Stewart, Regina Hoo,
558 Peng He, Kerstin B. Meyer, and Sarah A. Teichmann. Automatic cell-type harmonization and
559 integration across Human Cell Atlas datasets. *Cell*, 186(26):5876–5891.e20, December 2023.
560 ISSN 0092-8674. doi: 10.1016/j.cell.2023.11.026.
- 561 [49] Danila Bredikhin, Ilia Kats, and Oliver Stegle. MUON: multimodal omics analysis frame-
562 work. *Genome Biol.*, 23(1):1–12, December 2022. ISSN 1474-760X. doi: 10.1186/
563 s13059-021-02577-8.
- 564 [50] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9
565 (8):1735–1780, November 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL
566 <https://doi.org/10.1162/neco.1997.9.8.1735>.

- 567 [51] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical Evaluation
568 of Gated Recurrent Neural Networks on Sequence Modeling, December 2014. URL <http://arxiv.org/abs/1412.3555>. arXiv:1412.3555 [cs].
569
- 570 [52] Yann LeCun, Bernhard Boser, John Denker, Donnie Henderson, R. Howard, Wayne
571 Hubbard, and Lawrence Jackel. Handwritten Digit Recognition with a Back-
572 Propagation Network. In *Advances in Neural Information Processing Systems*, volume 2.
573 Morgan-Kaufmann, 1989. URL <https://proceedings.neurips.cc/paper/1989/hash/53c3bce66e43be4f209556518c2fcb54-Abstract.html>.
574
- 575 [53] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N
576 Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is All you Need.
577 In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://papers.nips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html.
578
579
- 580 [54] Leonid Iosipoi and Anton Vakhrushev. SketchBoost: Fast Gradient Boosted Decision Tree
581 for Multioutput Problems. *Advances in Neural Information Processing Systems*, 35:25422–
582 25435, December 2022. URL https://papers.nips.cc/paper_files/paper/2022/hash/a36c3dbe676fa8445715a31a90c66ab3-Abstract-Conference.html.
583
- 584 [55] Yan Wu, Esther Wershof, Sebastian M Schmon, Marcel Nassar, Błażej Osiński, Ridvan Eksi,
585 Kun Zhang, and Thore Graepel. Perturbench: Benchmarking machine learning models for
586 cellular perturbation analysis, 2024. URL <https://arxiv.org/abs/2408.10609>.
- 587 [56] Isaac Virshup, Sergei Rybakov, Fabian Theis, Philipp Angerer, and Fabian Alexander Wolf.
588 anndata: Access and store annotated data matrices, February 2024. URL <https://doi.org/10.5281/zenodo.10705762>.
589
- 590 [57] Robrecht Cannoodt, Hendrik Cannoodt, Dries Schaumont, Kai Waldrant, Eric Van de Kerckhove,
591 Andy Boschmans, Dries De Maeyer, and Toni Verbeiren. Viash: A meta-framework for
592 building reusable workflow modules. *Journal of Open Source Software*, 9(93):6089, 2024. doi:
593 10.21105/joss.06089. URL <https://doi.org/10.21105/joss.06089>.

594 **Checklist**

- 595 1. For all authors...
- 596 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
597 contributions and scope? [Yes] Our three main claims are that we (a) introduce a
598 valuable benchmarking dataset, (b) we used it to run a NeurIPS 2023 competition
599 “Single-cell perturbation prediction: generalizing experimental interventions to unseen
600 contexts“, analyzed the results, and (c) implemented the lessons learned in a new
601 benchmark. The sections that describe these contributions are correspondingly (a) 3.2,
602 (b) 4, (c) 3,4.3, 4.4.
- 603 (b) Did you describe the limitations of your work? [Yes] See Section 5.
- 604 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We
605 believe our work does not have potential for negative societal impact.
- 606 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
607 them? [Yes]
- 608 2. If you are including theoretical results...
- 609 (a) Did you state the full set of assumptions of all theoretical results? [N/A] No theoretical
610 results included.
- 611 (b) Did you include complete proofs of all theoretical results? [N/A] No theoretical results
612 included.
- 613 3. If you ran experiments (e.g. for benchmarks)...
- 614 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
615 mental results (either in the supplemental material or as a URL)? [Yes] We included the
616 links to code (Appendix I, H) and data (Appendix ??) in the supplementary material.
- 617 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
618 were chosen)? [Yes] See Appendix B for data splits. The hyperparameters of the
619 models in the benchmark are specified in the attached code (Appendix I). As they were
620 not developed by us, we do not provide an explanation for their choice.
- 621 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
622 ments multiple times)? [Yes] See Appendix Figure 8 for the results of the experiments
623 under dataset bootstrapping, with the corresponding standard deviation error bars.
- 624 (d) Did you include the total amount of compute and the type of resources used (e.g., type
625 of GPUs, internal cluster, or cloud provider)? [Yes] Described in Appendix I.5.
- 626 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 627 (a) If your work uses existing assets, did you cite the creators? [Yes] We implemented
628 models that were submitted to the competition, see Appendix E and the code linked in
629 Appendix I for credits.
- 630 (b) Did you mention the license of the assets? [Yes] See Appendix E and the repository
631 linked in I for information on the implemented methods and code license information
- 632 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
633 Yes, see the repositories linked in Appendix I and H for the code of the benchmark and
634 analysis, respectively. See Appendix A for the description of the introduced dataset.
- 635 (d) Did you discuss whether and how consent was obtained from people whose data you’re
636 using/curating? [Yes] See Appendix J.
- 637 (e) Did you discuss whether the data you are using/curating contains personally identifiable
638 information or offensive content? [Yes] See Appendix J.
- 639 5. If you used crowdsourcing or conducted research with human subjects...
- 640 (a) Did you include the full text of instructions given to participants and screenshots, if
641 applicable? [N/A]

- 642 (b) Did you describe any potential participant risks, with links to Institutional Review
643 Board (IRB) approvals, if applicable? [Yes] See Appendix J.
- 644 (c) Did you include the estimated hourly wage paid to participants and the total amount
645 spent on participant compensation? [N/A] We did not directly hire participants; we
646 only used data acquired from human subjects as described in Appendix J.

647 A Detailed dataset description

648 A.1 Overview

649 We measure the impact of 146 compounds in human PBMCs with three replicates, one per each donor.
650 The dataset was generated in a 96-well plate format with sample multiplexing such that each of the
651 12 wells in each row of the plate were pooled in a single lane of the 10x Chromium chip. We included
652 controls in 3 out of 12 wells in each row of the plate such that each resulting multiplexed library
653 contains a negative control DMSO well treatment and two positive controls of either belinostat or
654 dabrafenib. The remaining 9 wells per row each contained a different treatment condition. The result
655 is 576 unique scRNA samples (Appendix A.5). The dose of belinostat is $0.1\mu\text{M}$, DMSO $14.1\mu\text{M}$,
656 and the rest of the compounds $1\mu\text{M}$. After sample demultiplexing, preprocessing, and quality control
657 filtering, we retained 301,785 single cells and 21265 genes for further analysis. Further filtering and
658 processing were performed to better tailor the dataset for the perturbation prediction task. Differential
659 expression was computed with `limma` to create the representation of perturbation effects used in the
660 benchmark (Appendix A.8).

661 A.2 Data Availability

662 As is standard in the computational biology field, processed counts data is publicly available through
663 the Gene Expression Omnibus (GEO) with accession GSE279945 and raw sequencing data is available
664 through the Sequencing Read Archive (SRA) with accession PRJNA1149320.

665 **Maintenance plan** The dataset will be stored on GEO and SRA indefinitely. Any updates to the
666 dataset will be made available on these platforms. The source code of the components and workflows
667 used in this study are stored on GitHub at github.com/openproblems-bio/task_perturbation_prediction.
668 At the time of publication, the project was published on GitHub and Zenodo as release 1.0.0.
669 Each component is backed by a Docker container published at ghcr.io/openproblems-bio, also
670 using tag 1.0.0. Any feedback or found errors can be reported through GitHub issues at
671 github.com/openproblems-bio/task_perturbation_prediction.

672 **Responsibility** We, the authors, bear all responsibility to withdraw our paper and data in case
673 of violation of licensing or patient privacy rights. The dataset will be distributed under a Creative
674 Commons license (CC BY 4.0).

675 A.3 Culture of PBMCs

676 Human PBMCs from three donors were purchased from AllCells (www.allcells.com). Donor
677 information is described in Table 1, and the informed donor consent process is described in Appendix J.
678 PBMCs from one female and two males were used and were selected due to similarities in age and
679 BMI, the absence of reported use of medications, and sufficient cell inventory for data generation.

Table 1: PBMC Donor Information

Donor Name	Donor ID	Lot #	Age	Sex	BMI	Blood Type	Race	Smoker	CMV+
Donor 1	110044355	3097601	45	F	25.4	O+	White	-	Neg
Donor 2	110044590	3096819	52	M	37.2	O-	White	No	Pos
Donor 3	888676709	3094710	45	M	24.7	A+	White	No	Neg

680 PBMCs were thawed in RPMI (Gibco cat # 11875-093) supplemented with 10% heat inactivated
681 fetal bovine serum (HI-FBS, Gibco cat # 10082-147) and centrifuged for 8 minutes at $300 \times G$.
682 The cell pellet was resuspended in RPMI supplemented with 10% HI-FBS, counted on a Luna
683 fluorescent cytometer (Logos Biosystems) using AO/PI stain (Logos Biosystems, cat # F23001) per
684 the manufacturer’s instructions, and centrifuged to wash cells. The cell pellet was then resuspended
685 to a concentration of 1,000,000 cells/mL in RPMI supplemented with 10% HI-FBS. For perturbation

686 experiments, cells were plated at 200,000 cells/well in 96-well V-bottom plates (Thermo Scientific
687 cat # 277143) in 200 μ L media and were cultured for a total of 48 hours prior to collection. For
688 multiome profiling experiments, PBMCs were seeded into a T75 flask (Corning cat # 430641U) and
689 were cultured 24 hours before collection.

690 **A.4 Characterization of PBMCs Across Donors**

691 Flow cytometry was used to characterize the major cell populations in the PBMC samples from the
692 three donors after thaw (0 hours) and at 48 hours of culture in 96-well V-bottom plates. This was
693 performed to confirm the relative consistency of cell types across donors and to ensure that the time
694 in culture and media conditions did not bias the survival of specific cell types. 200,000 PBMCs per
695 well were seeded in a 96-well V-bottom plate and centrifuged for 8 minutes at 300 x G. The cell pellet
696 was resuspended in antibodies against established cell lineage markers, which were diluted in Cell
697 Staining Buffer (Biolegend cat # 420201) and incubated at 4C in the dark for 25 minutes. PBMCs
698 incubated in a Cell Staining Buffer without adding antibodies were used as unstained controls. Cells
699 were washed by centrifugation for 8 minutes at 300 x G and resuspended in a Cell Staining Buffer.
700 Events were captured on a Novocyte Quanteon (config. V8B7Y6R4) with an average of 56,500
701 PBMCs per well acquired for analysis.

702 Prior to quantification, the spectral overlap of our conjugated antibodies was adjusted for using
703 UltraComp eBeads™ Plus Compensation Beads (Invitrogen cat # 01-3333-42), per the manufacturer's
704 instructions. Briefly, two conditions were used to compensate for spectral overlap: 1) unstained beads,
705 and 2) single-colored controls with each antibody applied individually to the beads. Antibodies were
706 incubated together with beads for 15 minutes, washed, and resuspended in a Cell Staining Buffer,
707 following which events were captured on the cytometer. The compensation matrix was generated on
708 FlowJo 10.8.1 and applied before the quantification of cell populations within PBMCs. The gating
709 strategy used to quantify CD3+ T-cells, CD14+ myeloid cells, CD19+ B-cells, and CD56+ NK cells
710 is described in Appendix Figure 4.

711 Overall, we observed that the four major cell populations measured were relatively consistent across
712 all donors at each time point, with CD3+ T-cells comprising most cells within the sample. We
713 noticed a reduction in the CD14+ myeloid compartment following the culture of the cells, which
714 was consistent across all donors. We speculate this could be due to the myeloid population tending
715 to differentiate and adhere following time in culture. We also acknowledge that the broad cell type
716 markers used for characterization via flow cytometry do not permit quantification of more specific
717 cell clusters (e.g., CD4+ vs CD8+ T-cells, monocytes vs. dendritic cells) that can be performed
718 using gene markers in the sequencing data, making a direct comparison of cell numbers across the
719 modalities more challenging. In sum, we selected PBMCs from three donors that contain relatively
720 consistent numbers of cell types within each sample and perform similarly after 48 hours of culture.

721 **A.5 Compound information and treatment of PBMCs**

722 146 compounds were applied to PBMCs from three healthy donors 24 hours after thawing and seeding
723 into 96 well V-bottom plates. Compounds were selected based on two criteria:

- 724 1. Inclusion in Library of Integrated Network-Based Cellular Signatures (LINCS) Connectivity
725 Map dataset, and
- 726 2. Robust transcriptional response observed in CD34+ hematopoietic stem cells (data not
727 released).

728 These compounds also span a diverse range of mechanisms of action.

729 Compounds were resuspended in DMSO to 1 mM and arrayed onto a 96-well PCR plate. Each of the
730 first three columns on the plate contained, respectively:

- 731 1. Belinostat, an HDAC inhibitor that we previously observed to induce a large transcriptional
732 response in PBMCs (positive control).

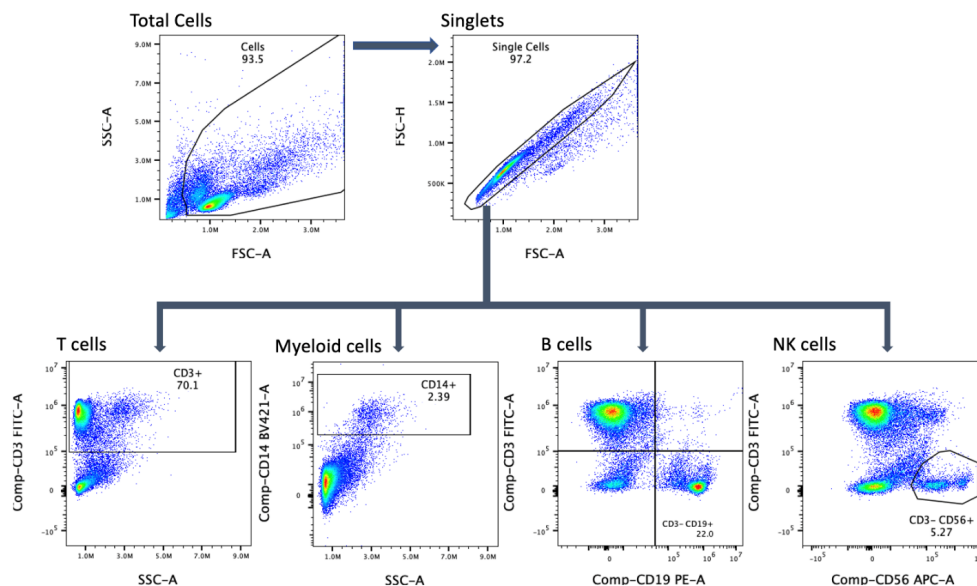


Figure 4: **Gating strategy for quantification of cell types within PBMCs.** FlowJo software was used to quantify the population of cells within PBMCs using the antibodies described in the methods section. First, a gate was generated by visualizing the forward scatter area and side scatter area to define the total live cells and to remove non-viable cells and debris from the analysis. Next, a gate was generated visualizing forward scatter height and area to define single cells and exclude doublets. Lastly, two-parameter density plots were used to assess the percentage of T cells (CD3+), myeloid cells (CD14+), B cells (CD3-, CD19+), and natural killer cells (CD3-, CD56+). Using this gating strategy, the percentage of cells within each population was quantified in PBMCs from the three donors used prior to full-scale data generation to ensure consistency in the samples and that time in cluster did not skew populations in a donor-specific manner.

- 733 2. Dabrafenib, a BRAF-inhibitor that we previously observed to induce a moderately-strong
 734 transcriptional response in PBMCs (positive control).
 735 3. No compound treatment except for DMSO (negative control).

736 Each well in columns 4-12 of the 96-well PCR plate (72 wells) contained a unique treatment
 737 compound. On the day of treatment, compounds were diluted and mixed directly in the PCR to
 738 100 μ M in RPMI (Gibco cat # 11875-093) using an Integra Viaflo 384 automated liquid handler.
 739 2 μ L of compound in solution was then transferred using the Integra Viaflo 384 automated liquid
 740 handler from the PCR plate and applied to PBMCs cultured in a 96-well V-bottom plate, described
 741 above. The use of an automated liquid handler enabled simultaneous application of 96 compounds
 742 and limited errors in transferring. The final concentration of compound applied to the cells was 1 μ M.
 743 Cells were treated 24 hours before collection.

744 A.6 Single-cell sequencing of perturbed PBMCs

745 48 hours after thaw, PBMCs cultured and treated in 96-well V-bottom plates were mixed with an
 746 Integra Viaflo 384, and a sample of cell suspension was transferred into a Thermo Scientific U-
 747 bottom plate (cat # 163320) containing CountBright Plus Absolute Counting Beads (Invitrogen cat #
 748 C36995) and SYTOX AADvanced Ready Flow (Invitrogen cat # R37173) diluted in DPBS, per the
 749 manufacturer's instructions. Total viability and live cell number per well were quantified via flow
 750 cytometry using a BD FACSCelesta Cell Analyzer (BD Biosciences).

751 The remaining treated PBMCs were centrifuged in the 96 well V-bottom plate at 300 x G for 8
 752 minutes. Culture media was aspirated using semi-automated liquid handling to not disturb the cell
 753 pellet and washed once with Cell Staining Buffer. Cells were centrifuged and resuspended in 12

754 unique Cell Multiplexing Oligonucleotides (3' CellPlex Kit Set A, 10X Genomics, cat # 1000261)
755 applied in columns 1 through 12, incubated for 5 minutes at room temperature, and then quenched
756 and washed three times in DPBS supplemented with 4% human serum albumin (Grifols cat # NDC
757 68516-5216-1). This allowed for consolidation of a 96-well plate into 8 pools, each containing cells
758 from a well labeled with a unique Cell Multiplexing Oligonucleotide. To ensure equal sequencing
759 representation from each compound treatment, 100,000 cells per well (calculated from the initial
760 cell count) were pooled by row into a 5 mL conical tube, resulting in a total of 8 pools, 1 conical
761 tube collected per row, using an Integra Assist Plus and associated D-ONE module automated liquid
762 handling instrument (Integra Biosciences).

763 Final pools therefore contained a DMSO negative control, 2 transcriptionally active positive controls,
764 and 9 experimental compounds. Cells were pelleted by centrifugation and resuspended to 1.2 x 10⁶
765 cells/mL in Cell Staining Buffer for downstream preparation for single-cell sequencing. Single-cell
766 libraries were prepared from each pool using Chromium Next GEM Single-cell 3' Kit v3.1 (10X
767 Genomics, cat # 1000268) following the manufacturer's recommended protocol (10X Genomics,
768 CG000388 Rev C).

769 Briefly, a total of 12,000 cells (1,200 per multiplexing oligo) were loaded into a single channel of a
770 Chromium Next GEM Chip G (10X Genomics, Cat # 1000120) and partitioned into droplets with
771 gel beads using a Chromium controller (10X Genomics, cat # 1000204). After emulsion droplets
772 were formed and collected, reverse transcription reactions were incubated at 53C for 45 minutes.
773 Barcoded transcripts were purified, amplified and size fractionated to create separate libraries for
774 the transcriptome and feature barcode fractions from each sample. Transcriptome libraries were
775 fragmented and ligated to indexed sequencing adapters according to the manufacturer's recommended
776 protocol. Feature barcode libraries were prepared using 3' Feature Barcode Kit (10X Genomics,
777 cat # 1000262) following the manufacturer's recommended protocol (10X Genomics, CG000388
778 Rev C). Transcriptome libraries were sequenced on an Illumina NovaSeq6000 with paired end reads
779 as follows: Read 1 = 28 cycles, i7 Index = 10 cycles, i5 = 10 cycles, Read 2 = 89 cycles. Feature
780 barcode libraries were sequenced on an Illumina NovaSeq6000 with paired end reads as follows:
781 Read 1 = 28 cycles, i7 Index = 10 cycles, i5 = 10 cycles, Read 2 = 35 cycles. Cell Ranger (v5.0.1)
782 mkfastq was used to generate all demultiplexed FASTQ files from the raw sequencing data.

783 Cell Ranger count was used to align transcriptome reads to the human GRCh38 genome reference,
784 identify corresponding feature barcode reads according to a csv reference file containing all the
785 relevant information needed for downstream processing, and quantify gene and UMI counts.

786 **A.7 Multiome ATAC + gene expression profiling of PBMCs at baseline**

787 24 hours after thaw, PBMCs cultured in T75 flasks were collected in a 50 mL conical tube, centrifuged
788 at 300 x G for 8 minutes, and washed once with DPBS. Viability and total live cells/mL were quantified
789 using AO/PI stain on a Luna fluorescent cytometer. Nuclei were isolated from cells using Chromium
790 Next GEM Single-cell Multiome ATAC + Gene Expression Reagent Bundle (10X Genomics, cat #
791 1000283) and Chromium Nuclei Isolation Kit with RNase Inhibitor (10X Genomics, cat # 1000494)
792 following the manufacturer's recommended protocol (10X Genomics, CG000365, Rev C). Briefly,
793 1.2 million cells were pelleted and resuspended in 100 µL of lysis buffer and incubated on ice for
794 5 minutes. Multiple rounds of washing were followed by resuspension in 150 µL of diluted nuclei
795 buffer and filtered through a 40 µm Flowmi Cell Strainer (Fisher Scientific, cat # 14100150). Nuclei
796 were counted on a Nexcelom Cellometer K2.

797 Multiome ATAC + Gene Expression libraries were prepared using Chromium Next GEM Single-cell
798 Multiome ATAC + Gene Expression Reagent Bundle (10X Genomics, cat # 1000283) following
799 the manufacturer's recommended protocol (10X Genomics, CG000338, Rev F). Briefly, a total of
800 8,000 nuclei were targeted for loading into a transposition reaction, which incubated at 37 C for
801 60 minutes. The output was then loaded into a single channel of a Chromium Next GEM Chip J
802 (10X Genomics, cat # 1000234) and partitioned into droplets with gel beads using a Chromium
803 controller (10X Genomics, cat # 1000204). After emulsion droplets were formed and collected,

804 reverse transcription and transposed DNA fragment barcoding reactions were incubated at 37C for
 805 45 minutes. Both products were purified, amplified and size fractionated to create the ATAC and
 806 transcriptome fractions from each sample. ATAC libraries had indexed sequencing adapters added.
 807 Transcriptome libraries were fragmented and ligated to indexed sequencing adapters. ATAC libraries
 808 were sequenced on an Illumina NovaSeq6000 with paired end reads as follows: Read 1 = 50 cycles,
 809 i7 Index = 8 cycles, i5 = 24 cycles, Read 2 = 49 cycles. Transcriptome libraries were sequenced
 810 on an Illumina NovaSeq6000 with paired end reads as follows: Read 1 = 28 cycles, i7 Index = 10
 811 cycles, i5 = 10 cycles, Read 2 = 89 cycles. Cell Ranger ARC (v2.0) mkfastq was used to generate
 812 all demultiplexed FASTQ files from the raw sequencing data. Cell Ranger ARC count was used to
 813 align transcriptome reads to the human GRCh38 genome reference provided by 10X Genomics and
 814 generate all downstream count matrices.

815 A.8 Processing of scRNA-seq perturbation data

816 Starting with the counts matrix output from Cell Ranger, cells with total numbers of transcripts that
 817 fell below or above specific thresholds were filtered out of the dataset. These transcript thresholds
 818 were determined per sequencing pool. All cells that had a mitochondrial counts fraction higher than
 819 0.2 were also removed. The Python package `scrublet` was then used to label cells with a probability
 820 of being doublets. These probabilities were smoothed over a k -nearest neighbor graph constructed
 821 from the cells, and cells with a smoothed doublet probability of greater than 0.2 were filtered out.

822 During the pooling process (Appendix A.6), cells from each of the twelve wells in a plate row were
 823 tagged with distinct cell multiplexing oligonucleotides to increase throughput and decrease batch
 824 effects across wells. This multiplexing procedure necessitated a demultiplexing step in the processing
 825 pipeline, whereby a multivariate Gaussian-mixture model was used to identify the well from which
 826 each cell most likely originated. Cells that could not be conclusively labeled as belonging to any
 827 particular well were dropped from the dataset.

828 Prior to cell-type annotation, counts were normalized to sum to 1000 in each cell and then transformed
 829 with the mapping $x \mapsto \ln(x + 1)$. Cell-type annotation was performed by first running Leiden
 830 clustering with `resolution = 1` on a k -nearest neighbor graph built from the 2000 most highly-
 831 variable genes, then manually assigning a cell type label (T-cells, B-cells, myeloid cells, or NK cells)
 832 to each cluster based on expression of the marker genes in Table 2. In addition, we filtered samples
 833 of certain compounds as described in Appendix F.2.

Table 2: PBMC Marker Genes

Cell Type	Marker Genes
T-cells	CCR6, CCR7, CD2, CD27, CD3D, CD3E, CD3G, CD4, CD6, CD8A, CTLA4, FOXP3, GZMB, IL2RA, PTPRC, TRDV2, TRGC1, TRGV9
B-cells	CD19, CD24, CD24, CD27, CD38, CD38, CD38, CD79A, CD79B, DERL3, FKBP11, HLA-DQA1, HLA-DQB1, IGLL5, IGLL5, IGLL5, JCHAIN, MS4A1, PAX5, PRDX4, PTPRC, SEC11C, SSR4, TCL1A, VPREB3
Myeloid cells	CD14, CD163, CD1C, CD68, CD83, ITGAX
NK cells	CD2, CD69, COX6A2, FCGR3A, GNLY, GZMA, GZMB, GZMM, KIR2DL4, KLRB1, NCAM1, NCR1, NKG7, NKG7, ZMAT4

834 Next, differential expression (DE) was performed to produce a representation of the perturba-
 835 tion effects of each drug. To ensure our DE computation was as robust as possible, we used the
 836 `filterByExpr` function from the EdgeR package to filter down to 5317 genes that were consistently
 837 expressed across all cell types. Counts from these 5317 genes are then summed across the cells
 838 of each type in every well to produce what is known as a *pseudobulked* expression object. The

839 pseudobulked counts are then fed into the `limma/voom` pipeline to compute moderated p -values and
840 log-fold change statistics per gene for each condition. The linear model fit by `limma` included an
841 additional covariate that captured the plate-to-plate batch effects. This covariate also reflected the
842 variability in perturbation effect across donors, as each plate contained samples from only one donor.

843 All of the processing steps described in this section, unless explicitly stated otherwise, were performed
844 using the `scanpy` library [46].

845 **A.9 Processing of baseline Multiome snRNA-seq/scATAC-seq data**

846 For processing the joint snRNA-seq/scATAC-seq measurements, we start with the outputs provided
847 by the Cell Ranger pipeline, namely:

- 848 1. the fragments file, which lists both the region of the chromosome and the cell barcode for
849 each detected ATAC-seq fragment, and
- 850 2. the filtered feature-barcode matrix, which contains both the detected genes (snRNA-seq)
851 and called peaks (scATAC-seq) assigned to each cell barcode.

852 We first split up the feature-barcode matrix into a cell-by-gene snRNA-seq matrix and a cell-by-peaks
853 scATAC-seq matrix. The QC steps for the snRNA-seq measurements were nearly identical to the
854 process described in Appendix A.8 for the scRNA-seq data, albeit with a slight hand-tuning of the
855 filtering thresholds. Namely, cells with low counts (below 500 transcripts), high mitochondrial counts
856 percentage (above 0.2), or high probability of being doublets (above 0.2) were removed, and genes
857 that were expressed in fewer than 100 cells were also filtered out, resulting in 17438 distinct genes.
858 Following this, counts were normalized to sum to 1000 in each cell and then rescaled using the
859 mapping $x \mapsto \ln(x + 1)$.

860 Cells were further filtered using the scATAC-seq measurements. Specifically, cells that met any of
861 the following criteria were removed:

- 862 1. fewer than 1000 fragments,
- 863 2. fewer than 750 peaks,
- 864 3. *transcription start site (TSS) enrichment score* below 0.8, or
- 865 4. *nucleosome signal* below 2.0.

866 Both the TSS enrichment score and nucleosome signal are common metrics for evaluating the quality
867 of chromatin accessibility measurements. The TSS enrichment score is calculated by taking windows
868 of 2000bp around either side of TSSs, then computing the average ratio of read depth at 100bps on
869 either side of these windows to the read depth at the respective TSS in the center of the window [47].
870 For the sake of computational efficiency, we estimate the TSS enrichment score by computing this
871 average ratio over a random subset of 3000 TSSs rather than every TSS.

872 The nucleosome signal is the ratio of the number of single-nucleosome fragments (between 147bp
873 and 294bp) to the number of nucleosome-free fragments (shorter than 147bp). Again for the sake
874 of computational efficiency, we estimate the nucleosome signal using a subset of the ATAC-seq
875 fragments.

876 Specific peaks that were observed in fewer than 20 cells were also dropped.

877 After filtering, cell type annotation was performed per-donor by running Leiden clustering, then
878 assigning all the cells in each cluster a cell type label using `celltypist` [48]. These annotations
879 were then validated based on the expression of the marker genes listed in Table 2. If a cluster could
880 not be conclusively labeled with a specific cell type, the cells from that cluster were filtered out. All
881 of these preprocessing steps were performed with either `scanpy` (for snRNA-seq) and `muon` (for
882 scATAC-seq) [46, 49].

883 **A.10 Datasheet for datasets**

884 **Motivation**

885 **For what purpose was the dataset created?** Was there a specific task in mind? Was there a specific
886 gap that needed to be filled? Please provide a description.

887 OP3 dataset was created to enable research on predicting cell-type specific transcriptomic response to
888 drugs. The dataset was created intentionally with that task in mind, providing donor replicates to
889 account for the variability of outcomes.

890 **Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g.,
891 company, institution, organization)?**

892 The authors of this paper, along with the additional scientists at Cellarity listed in the acknowledgment
893 section, namely Lijun Zhao, Roman Montez, Nicole Robichaud, Nina Colon, Sakina Saif, Laura
894 Isacco, and Cameron Reilly.

895 **Who funded the creation of the dataset?** If there is an associated grant, please provide the name of
896 the grantor and the grant name and number.

897 This work was supported by funds from the Chan Zuckerberg Initiative, Cellarity Inc., the Helmholtz
898 Association and Helmholtz Munich. This work was co-funded by the European Union (ERC,
899 DeepCell -101054957).

900 **Any other comments?**

901 None.

902 **Composition**

903 **What do the instances that comprise the dataset represent (e.g., documents, photos, people,
904 countries)?** Are there multiple types of instances (e.g., movies, users, and ratings; people and
905 interactions between them; nodes and edges)? Please provide a description.

906 OP3 contains scRNA-seq data of PBMCs across three donors. Cells are either control (DMSO) or
907 were exposed to one of 146 drugs. It also provides multimodal (joint snRNA-seq and scATAC-seq)
908 data for the same three donors at baseline. Processed data contains p -values and log-fold change per
909 cell type and gene for each drug, which indicate the significance and magnitude of gene expression
910 change induced by a given compound in a given cell type.

911 **How many instances are there in total (of each type, if appropriate)?**

912 There are 449650 cells collected across 576 samples in the raw scRNA-seq dataset. After filtering for
913 the perturbation prediction task, this becomes 298087 cells across 567 samples.

914 Meanwhile, the raw multiome snRNA-seq/scATAC-seq data contains 53086 cells, which are filtered
915 down to 22591 during processing.

916 **Does the dataset contain all possible instances or is it a sample (not necessarily random) of
917 instances from a larger set?** If the dataset is a sample, then what is the larger set? Is the
918 sample representative of the larger set (e.g., geographic coverage)? If so, please describe how
919 this representativeness was validated/verified. If it is not representative of the larger set, please
920 describe why not (e.g., to cover a more diverse range of instances, because instances were withheld
921 or unavailable).

922 While individual cells and samples were removed from the raw data for failing to pass quality-control,
923 the raw data is available to download and represents all the samples that were collected in this
924 experiment.

925 **What data does each instance consist of? “Raw” data (e.g., unprocessed text or images) or
926 features?** In either case, please provide a description.

927 The most raw form of the data is a collection of .bc1 files from the Illumina sequencer (not released).
928 These are then processed into .fastq files, which we have made available on the Sequencing Read
929 Archive (SRA), as is standard practice for the computational biology field. .fastq files are then
930 converted into raw counts matrices through standard Cell Ranger bioinformatics pipelines. For

931 the scRNA-seq data, the combined raw counts matrix has 449650 rows (cells) and 58676 columns
932 (genes). The majority of columns contain either all zeros or very few measurements. For multimodal
933 snRNA-seq/scRNA-seq data, the raw counts matrix has 53086 rows (cells) and 172019 columns. Of
934 these columns, 36601 are gene expression measurements, while the other 135418 measure chromatin
935 accessibility. Similar to the scRNA-seq data, this matrix is extremely sparse.

936 **Is there a label or target associated with each instance?** If so, please provide a description.

937 The only information that is known about any given cell with absolute certainty is which sequencing
938 library, plate, and donor the cell originated from. However, if a cell can be assigned to a given well
939 in the demultiplexing process (Appendix A.8), then well-level metadata, which includes compound
940 treatment, can be attached to the cell. Moreover, marker gene expression can be used to label the
941 majority of cells with high-confidence cell type annotations.

942 The processed dataset (after DGE analysis) contains a $-\log_{10}(p\text{-value}) \times \text{sign}(\log\text{-fold change})$
943 statistic for each cell type, compound, and gene, which indicates the significance and direction of a
944 gene expression change.

945 **Is any information missing from individual instances?** If so, please provide a description, ex-
946 plaining why this information is missing (e.g., because it was unavailable). This does not include
947 intentionally removed information, but might include, e.g., redacted text.

948 Single-cell RNA-seq data is sparse, meaning that counts for the majority of genes are missing from
949 each individual cell. This sparsity is caused by a number of different factors, including stochasticity of
950 gene expression and constraints on read depth per cell. In addition, the wells with certain compounds
951 had few or no viable cells to sequence, which might have been a result of compound toxicity or
952 experimental conditions in a given well.

953 **Are relationships between individual instances made explicit (e.g., users' movie ratings, social
954 network links)?** If so, please describe how these relationships are made explicit.

955 Which cells belong to the same donor or were cultured on the same plate can be determined directly
956 from the raw data. Among the cells that can be successfully demultiplexed (Appendix A.8), it can
957 be further determined which cells came from the same well and which were treated with the same
958 compound.

959 **Are there recommended data splits (e.g., training, development/validation, testing)?** If so, please
960 provide a description of these splits, explaining the rationale behind them.

961 See Appendix B.

962 **Are there any errors, sources of noise, or redundancies in the dataset?** If so, please provide a
963 description.

964 See Appendix I.2.

965 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,
966 websites, tweets, other datasets)?** If it links to or relies on external resources, a) are there guarantees
967 that they will exist, and remain constant, over time; b) are there official archival versions of the
968 complete dataset (i.e., including the external resources as they existed at the time the dataset was
969 created); c) are there any restrictions (e.g., licenses, fees) associated with any of the external resources
970 that might apply to a future user? Please provide descriptions of all external resources and any
971 restrictions associated with them, as well as links or other access points, as appropriate.

972 The dataset is entirely self-contained.

973 **Does the dataset contain data that might be considered confidential (e.g., data that is pro-
974 tected by legal privilege or by doctor-patient confidentiality, data that includes the content of
975 individuals non-public communications)?** If so, please provide a description.

976 The dataset contains human samples that were obtained with the consent of the subjects. See
977 Appendix J.

978 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening,
979 or might otherwise cause anxiety?** If so, please describe why.

980 No.

981 **Does the dataset relate to people?** If not, you may skip the remaining questions in this section.
982 Yes. The data was derived from human blood samples.

983 **Does the dataset identify any subpopulations (e.g., by age, gender)?** If so, please describe how
984 these subpopulations are identified and provide a description of their respective distributions within
985 the dataset.

986 Yes, we included the age, sex, BMI, race, smoker status, and CMV+ status of the donors. The data
987 was collected through the general health interview described in Appendix J.

988 **Is it possible to identify individuals (i.e., one or more natural persons), either directly or
989 indirectly (i.e., in combination with other data) from the dataset?** If so, please describe how.

990 See Appendix J.

991 **Does the dataset contain data that might be considered sensitive in any way (e.g., data that
992 reveals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or
993 union memberships, or locations; financial or health data; biometric or genetic data; forms
994 of government identification, such as social security numbers; criminal history)?** If so, please
995 provide a description.

996 The data contains the racial origin and health information, including BMI, smoker status, and CMV+
997 status of the donors that were collected through the general health interview described in Appendix J.
998 In theory, unique gene expression patterns could be used to identify donors.

999 **Any other comments?**

1000 None.

Collection Process

1001

1002 **How was the data associated with each instance acquired?** Was the data directly observable (e.g.,
1003 raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived
1004 from other data (e.g., part-of-speech tags, model-based guesses for age or language)? If data was
1005 reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If
1006 so, please describe how.

1007 We performed the scRNA-seq and multimodal snRNA-seq/scATAC-seq assays to study the effects of
1008 the drugs on the gene expression, as described in Appendix A.1.

1009 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or
1010 sensor, manual human curation, software program, software API)?** How were these mechanisms
1011 or procedures validated?

1012 The experiments are described in detail in Appendix A.

1013 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,
1014 probabilistic with specific sampling probabilities)?**

1015 The raw data is available to download and represents all the samples that were collected in this
1016 experiment.

1017 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and
1018 how were they compensated (e.g., how much were crowdworkers paid)?**

1019 The information on sample collection is available in Appendix J.

1020 **Over what timeframe was the data collected? Does this timeframe match the creation timeframe
1021 of the data associated with the instances (e.g., recent crawl of old news articles)?** If not, please
1022 describe the timeframe in which the data associated with the instances was created.

1023 Cells were collected from patients in 2021-2022, while the perturbation experiments were performed
1024 at Cellarity in June and July of 2023.

1025 **Were any ethical review processes conducted (e.g., by an institutional review board)?** If so,
1026 please provide a description of these review processes, including the outcomes, as well as a link or
1027 other access point to any supporting documentation.

1028 See Appendix J.

1029 **Does the dataset relate to people?** If not, you may skip the remaining questions in this section.

1030 Yes.

1031 **Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (e.g., websites)?**

1032 We purchased commercially available human tissue samples from AllCells, Inc.

1033 **Were the individuals in question notified about the data collection?** If so, please describe (or show with screenshots or other information) how notice was provided, and provide a link or other access point to, or otherwise reproduce, the exact language of the notification itself.

1034 Yes, see Appendix J.

1035 **Did the individuals in question consent to the collection and use of their data?** If so, please describe (or show with screenshots or other information) how consent was requested and provided, and provide a link or other access point to, or otherwise reproduce, the exact language to which the individuals consented.

1036 Yes, see Appendix J.

1037 **If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses?** If so, please provide a description, as well as a link or other access point to the mechanism (if appropriate).

1038 Yes, see Appendix J.

1039 **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data protection impact analysis) been conducted?** If so, please provide a description of this analysis, including the outcomes, as well as a link or other access point to any supporting documentation.

1040 See Appendix J.

1041 **Any other comments?**

1042 None.

1053 **Preprocessing/cleaning/labeling**

1054 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)?** If so, please provide a description. If not, you may skip the remainder of the questions in this section.

1055 We provide the raw version of the dataset, processed, and the code used for data processing. Data processing is described in Appendix A.8 and A.9.

1056 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)?** If so, please provide a link or other access point to the “raw” data.

1057 Raw data for both the perturbational scRNA-seq and baseline snRNA-seq/scATAC-seq data are currently available through the Sequencing Read Archive (SRA) with accession PRJNA1149320.

1058 **Is the software used to preprocess/clean/label the instances available?** If so, please provide a link or other access point.

1059 Yes, see github.com/theislab/task-dge-perturbation-prediction-analysis and github.com/openproblems-bio/task_perturbation_prediction for the code used for data processing. In addition, other steps not included in the code are outlined in Appendix A.

1060 **Any other comments?**

1061 None.

1071 **Uses**

1072 **Has the dataset been used for any tasks already?** If so, please provide a description.

1073 The dataset has been used for the Kaggle competition as part of the NeurIPS 2023 Competitions
1074 track called *Single-cell perturbation prediction: generalizing experimental interventions to unseen*
1075 *contexts*. It was also used to develop the benchmark described in this paper, see Section 3.

1076 **Is there a repository that links to any or all papers or systems that use the dataset?** If so, please
1077 provide a link or other access point.

1078 The dataset will be officially released with this publication. Hence, no other papers used this dataset.

1079 **What (other) tasks could the dataset be used for?**

1080 Aside from the use outlined in this study, the dataset enables myriad other inquiries, including but not
1081 limited to: the impact of different compound mechanisms of action on gene expression, the variance
1082 in compound effects across donors, pathway-based analyses of perturbation effects, etc.

1083 **Is there anything about the composition of the dataset or the way it was collected and prepro-**
1084 **cessed/cleaned/labeled that might impact future uses?** For example, is there anything that a future
1085 user might need to know to avoid uses that could result in unfair treatment of individuals or groups
1086 (e.g., stereotyping, quality of service issues) or other undesirable harms (e.g., financial harms, legal
1087 risks) If so, please provide a description. Is there anything a future user could do to mitigate these
1088 undesirable harms?

1089 Some medical information in the dataset could be used for deanonymization. However, given the
1090 limited scope of the provided data, it is highly unlikely that particular individuals or groups would be
1091 unfairly treated as a result of using this dataset.

1092 **Are there tasks for which the dataset should not be used?** If so, please provide a description.

1093 Given the limited scope of this dataset, it should not be used to influence immediate medical decision-
1094 making.

1095 **Any other comments?**

1096 None.

1097 **Distribution**

1098 **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution,**
1099 **organization) on behalf of which the dataset was created?** If so, please provide a description.

1100 Yes, the dataset will be publicly available on the internet.

1101 **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)** Does the dataset
1102 have a digital object identifier (DOI)?

1103 As is standard in the computational biology field, processed counts data is publicly available through
1104 the Gene Expression Omnibus (GEO) with accession GSE279945 and raw sequencing data is available
1105 through the Sequencing Read Archive (SRA) with accession PRJNA1149320.

1106 **When will the dataset be distributed?**

1107 If this paper is accepted into the Datasets and Benchmarks track, the dataset will be distributed
1108 publicly with the submission of the camera-ready version of the paper, at the latest. However, we will
1109 likely release the dataset sooner due to interest in the single-cell research community.

1110 **Will the dataset be distributed under a copyright or other intellectual property (IP) license,**
1111 **and/or under applicable terms of use (ToU)?** If so, please describe this license and/or ToU, and
1112 provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU,
1113 as well as any fees associated with these restrictions.

1114 The dataset will be distributed under a Creative Commons license (CC BY 4.0).

1115 **Have any third parties imposed IP-based or other restrictions on the data associated with**
1116 **the instances?** If so, please describe these restrictions, and provide a link or other access point
1117 to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these
1118 restrictions.

1119 No.

1120 **Do any export controls or other regulatory restrictions apply to the dataset or to individual**
1121 **instances?** If so, please describe these restrictions, and provide a link or other access point to, or
1122 otherwise reproduce, any supporting documentation.

1123 No.

1124 **Any other comments?**

1125 None.

Maintenance

1126 **Who will be supporting/hosting/maintaining the dataset?**

1128 The authors of this paper. The dataset will be hosted on the GEO platform indefinitely.

1129 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

1130 Contact the last author of this paper (dburkhardt@cellarity.com).

1131 **Is there an erratum?** If so, please provide a link or other access point.

1132 No.

1133 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?**

1134 If so, please describe how often, by whom, and how updates will be communicated to users (e.g.,
1135 mailing list, GitHub)?

1136 If any correction is needed, such as adjustments to metadata or refiltering of cells, we will upload a
1137 new version of the dataset to GEO. This will be noted on the OP3 benchmark GitHub page.

1138 **If the dataset relates to people, are there applicable limits on the retention of the data associated**
1139 **with the instances (e.g., were individuals in question told that their data would be retained for a**
1140 **fixed period of time and then deleted)?** If so, please describe these limits and explain how they will
1141 be enforced.

1142 There is no such limit. See Appendix J.

1143 **Will older versions of the dataset continue to be supported/hosted/maintained?** If so, please
1144 describe how. If not, please describe how its obsolescence will be communicated to users.

1145 Older versions will be available to download on GEO.

1146 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for**
1147 **them to do so?** If so, please provide a description. Will these contributions be validated/verified? If
1148 so, please describe how. If not, why not? Is there a process for communicating/distributing these
1149 contributions to other users? If so, please provide a description.

1150 Changes to data postprocessing can be proposed in GitHub issues and Pull Requests at
1151 github.com/openproblems-bio/task_perturbation_prediction. For all other changes,
1152 contact the authors of the paper.

1153 **Any other comments?**

1154 None.

1155 **B Data splits**

1156 To derive the competition training and test splits, the compounds were divided into three groups,
1157 public train, public test, and private test, at a ratio of 1:3:5. This lopsided train-test split was chosen to
1158 determine whether we could model perturbation signatures in unseen cell types while only measuring
1159 roughly 10% of the compounds in those cell types. Differential expression values were provided to
1160 competitors for all cell types for compounds in the train set but masked in B and myeloid cells for
1161 test perturbations, although they could evaluate their models against the public test set an unlimited
1162 number of times. The final score was computed only on the private test set.

1163 To avoid data leakage from the test set, we fit the training and test set DE models separately. To
 1164 generate the training data, we fit the DE model on only the samples from the training set. To generate
 1165 the private and public test data, we fit the DE model to all samples in the experiment. This kept the
 1166 test data private and ensured the test data was the most accurate.

1167 For the benchmark, we use only two splits, train and test, where the train split contains public train
 1168 and public test data, and the test split contains only private test data.

1169 C Benchmarking representations of perturbation effect

1170 C.1 Cross-donor retrieval

1171 As mentioned in Section 3.3, we developed the cross-donor retrieval heuristic to compare different
 1172 representations of perturbation effects. This heuristic is calculated as follows. Let:

- 1173 1. $C = \{c_1, \dots, c_{140}\}$ be the list of compounds,
- 1174 2. $G = \{g_1, \dots, g_{5317}\}$ be the list of genes,
- 1175 3. $T = \{t_1, t_2, t_3, t_4\}$ be the list of cell types, and
- 1176 4. $D = \{d_1, d_2, d_3\}$ be the list of donors.

1177 First, we compute differential expression (DE) across all genes for each donor-compound-cell type
 1178 combination (d, c, t) . Note that this is slightly different from the approach we take in computing DE
 1179 for the task data. In that context, we include data from all donors in our model and then add a donor
 1180 covariate to regress out donor-specific effects. For computing cross-donor retrieval, we compute DE
 1181 for each donor separately¹.

We let $\text{pert}_{d,c,t,g} \in \mathbb{R}$ denote the representation for gene g of the perturbation (d, c, t) , and let

$$\text{pert}_{d,c,t} \in \mathbb{R}^{|G|}$$

be the vector of representations for all genes. The for a fixed donor pair (d_i, d_j) and cell type t_k we
 compute the pairwise distance matrix

$$\begin{bmatrix} \left\| \text{pert}_{d_i, c_1, t_k} - \text{pert}_{d_j, c_1, t_k} \right\| & \cdots & \left\| \text{pert}_{d_i, c_1, t_k} - \text{pert}_{d_j, c_{140}, t_k} \right\| \\ \vdots & \ddots & \vdots \\ \left\| \text{pert}_{d_i, c_{140}, t_k} - \text{pert}_{d_j, c_1, t_k} \right\| & \cdots & \left\| \text{pert}_{d_i, c_{140}, t_k} - \text{pert}_{d_j, c_{140}, t_k} \right\| \end{bmatrix}.$$

1182 Now we replace each value in this pairwise distance matrix with its (ascending) rank among the
 1183 values in the same row. After computing the ranked distance matrix for all three pairs of donors, we
 1184 extract the diagonals of these matrices. This distribution of values for various representations and
 1185 metrics can be seen in Figure 2 of the main paper.

1186 C.2 Perturbation effect representation

1187 In Figure 2, we compared the following representations:

- 1188 1. $\log(\text{counts} + 1)$: natural log of raw counts per gene, with an additional pseudocount to
 1189 prevent taking the logarithm of 0.
- 1190 2. log-fold change: base-2 logarithm of change in normalized gene expression under the effect
 1191 of a perturbation, taken directly from the `logFC` output from `limma`.
- 1192 3. p -value: significance of change of gene expression, taken directly from the `P.Value` output
 1193 from `limma`.

¹While there is only one well for all treatment compounds per donor (besides the positive controls), there are 16 negative control wells for each donor. Hence, we can obtain estimates for the statistical significance of perturbation effects by comparing gene expression in the 1 treatment well against the 16 negative control wells.



Figure 5: **Cross-donor retrieval for binarized significance.** "DE X" stands for the binarized representation, with sign indicating direction, and X indicates the threshold of significance.

- 1194 4. $-\log_{10}(p\text{-value}) \times \text{sign}(\log\text{-fold change})$: the magnitude of this value correlates with the
 1195 statistical significance of the change in gene expression, while the sign corresponds to the
 1196 direction of the change.
- 1197 5. $-\log_{10}(\text{FDR-adjusted } p\text{-value}) \times \text{sign}(\log\text{-fold change})$: the FDR-adjusted p -value is the
 1198 adj .P.Val output from limma, which is computed using the Benjamini-Hochberg proce-
 1199 dure from the original p -values.

1200 In addition, we considered multiple strategies of binarizing the significance of change in gene
 1201 expression to cast the task as a classification problem. We found that $-\log_{10}(p\text{-value}) \times$
 1202 $\text{sign}(\log\text{-fold change})$ performs better as a perturbation representation according to the cross-donor
 1203 retrieval (Figure 5).

1204 D Limitations of Differential Expression Analysis

1205 Predicting transcriptional differential effects using standard tools in scRNA-seq data from hetero-
 1206 geneous cell populations, such as PBMCs treated with targeted drugs, presents several challenges.
 1207 Statistically, these tools must contend with batch effects, which can arise from processing times,
 1208 reagents, or sequencing runs. Although adjustments for batch effects can be incorporated into the
 1209 analysis design, the confounding of batch and treatment effects can still obscure true biological signals.
 1210 Small sample sizes or high biological variability can further hinder accurate dispersion estimates of
 1211 parametric methods like negative binomial models, thereby reducing statistical power. This limitation
 1212 is especially pronounced in low-abundance cell types, where variability is high, and transcript detec-
 1213 tion is low. Additionally, high-quality, consistent data across all samples is recommended, which is
 1214 challenging in practice. Insufficient sequencing depth and biological variability between donors can
 1215 obscure true differential effects. Biologically, the complexity of PBMCs introduces further limitations.
 1216 These cells engage in intricate interactions and signaling pathways that influence transcriptional

1217 responses indirectly, complicating the identification of direct drug effects. Heterogeneity within
1218 PBMC populations and baseline variability among donors can obscure drug-induced transcriptional
1219 changes. To address these issues, experiments should use matched samples from the same donors and
1220 apply robust normalization methods. Additionally, differential gene expression analysis may miss
1221 regulatory effects at other levels, such as protein activity and epigenetic modifications. Complement-
1222 ing scRNA-seq data with other omics data, such as proteomics or epigenomics, and integrating these
1223 datasets can provide a more comprehensive view of drug effects.

1224 **E Summary of methods**

1225 Below, we describe 6 methods submitted by challenge participants and the control methods. Note that
1226 the methods needed to be updated to generalize to different datasets, which might have impacted their
1227 performance. Despite contacting the authors and our efforts, we suspect that the implementations
1228 of LSTM-GRU-CNN Ensemble and Transformer ensemble might have worsened their predictions.
1229 All of the methods were released under MIT license [https://www.kaggle.com/competitions/
1230 open-problems-single-cell-perturbations/rules](https://www.kaggle.com/competitions/open-problems-single-cell-perturbations/rules).

1231 **E.1 Leaderboard winners**

1232 **E.1.1 LSTM-GRU-CNN Ensemble**

1233 Kaggle user jeannkouagou had the highest score on the private test with a method that integrated
1234 additional biological knowledge into the feature space. Towards that, they utilized ChemBERTa [30]
1235 embeddings for SMILES encodings of small molecules which resulted in notable improvements in
1236 predictive performance. Furthermore, a 5-fold cross-validation setting was utilized, incorporating
1237 three model architectures (LSTM [50], GRU [51], and 1d-CNN [52]) with multiple loss functions
1238 and three distinct input feature representations (initial, light, and heavy) to optimize model accuracy.
1239 The method also included additional data augmentation techniques, such as randomly replacing input
1240 features with zeros to simulate biological noise.

1241 **E.1.2 Transformer ensemble**

1242 Kaggle user Elior Kalfon proposed a method based on a transformer [53] ensemble and scored 2nd
1243 place on the leaderboard. This method employed an ensemble of four transformer models, each
1244 with different weights and trained on slightly varying feature sets. Their method considered both the
1245 strategies for both feature normalization and data sampling. The feature engineering process involved
1246 one-hot encoding of categorical labels, target encoding using mean and standard deviation, and
1247 enriching the feature set with the standard deviation of target variables. Their method also considered
1248 to normalize data based on both mean value and standard deviation (std), or only mean value. A
1249 sophisticated sampling strategy based on K-Means clustering was employed to partition the data into
1250 training and validation sets, ensuring a representative distribution. The model architecture leveraged
1251 sparse and dense feature encoding, along with a transformer for effective learning.

1252 **E.1.3 NN retraining with pseudolabels**

1253 Kaggle user Okon2000 scored 3rd place in the competition leaderboard using their multi-stage
1254 MLP approach. Both stages use an ensemble of MLPs that underwent individual hyperparameter
1255 optimization to select model dimensions, learning rate and dropout. The first round trains an ensemble
1256 of 7 MLPs to predict pseudolabels [31] for the entire test set. These pseudo labels are added to
1257 the training dataset and used in the second round, where an ensemble of 20 MLPs to predict the
1258 output. 4-fold cross-validation, averaged over 2 repeats per fold, was used to avoid overfitting. The
1259 submission finds benefit to replacing one-hot encoding with an embedding layer, but did not find
1260 improvements with various dataset denoising and label normalization schemes. The robustness of the
1261 model to increasing dataset size, noisy labels, and noisy inputs is examined, demonstrating small
1262 benefits to adding noise to training labels.

1263 E.2 Judge prize winners

1264 E.2.1 JN-AP-OP2

1265 The solution by Antoine Passemiers and Jalil Nourisa earned the 1st judge prize. They employed a
1266 deep neural network architecture for perturbation modeling. Initially, the training data was encoded
1267 using a leave-one-out encoder based on unique pairs of compounds and cell types, converting the
1268 data into a format of (n_samples, n_genes, n_encode), referred to as X, where n_encode is 2. Then,
1269 the encoded data, X, was fed into the first multi-layer perceptron (MLP1). MLP1 processed X in a
1270 sample-wise manner and utilized fully connected layers to learn inter-gene relationships by sharing
1271 the encoded data across genes. Next, the output of MLP1 was concatenated with the original encoded
1272 data X to form a new representation of (n_samples, n_genes, 2*n_encode), which merged the learned
1273 encoding with the original encoding. This combined data was then inputted into a second multi-layer
1274 perceptron (MLP2) in a gene-wise manner, resulting in a final representation of n_samples * n_genes.

1275 E.2.2 ScAPE

1276 Kaggle user Los Rodríguez proposed their method named ScAPE, which won 2nd place for the
1277 Judge’s award in the competition. With a similar design of chemCPA [13], the core of ScAPE is an
1278 auto-encoder that utilizes drug and cell features and outputs signed log(p-values). Specifically, it
1279 has separate encoders to learn the latent representations of cells and drugs, respectively, with noise
1280 introduced. ScAPE computes the features as the median of signed log(p-values) from differential
1281 expression analysis results calculated on single-cell level. In addition, it computes differential
1282 expression on pseudobulk level to get mean log(fold-changes) as extra information. The method
1283 uses cell features both in the encoding part and the decoding part of the neural network, which is
1284 non-probabilistic, as the authors didn’t observe further advantages, either with respect to accuracy or
1285 generalization ability, with additional variational inference. Using cell latent features during decoding
1286 gives the method better scores in the leaderboard, though there’s not much improvement observed
1287 during training. The model also employs a leave-one-drug-out cross-validation strategy to assess
1288 generalization to unseen drugs, which ensures robust predictions by leveraging both raw fold changes
1289 and the most variable genes, thus it results in a competitive performance. Besides, the authors also
1290 proposed several other designs of methods and benchmarked the performances. Their exploration
1291 of both the problem and methodology are well documented which could provide useful insights for
1292 further studies.

1293 E.2.3 Py-boost

1294 This solution earned the third judge prize. Kaggle user AmbrosM implemented a gradient-boosted
1295 decision tree model using the py-boost framework [54]. The data is preprocessed in two ways before
1296 model training. First, $-\log_{10}(\text{pvalue})\text{sign}(\text{lfc})$ values are converted to t-statistic. This mapping is
1297 continuous and bijective, so there is no loss of information. Second, the training data is compressed
1298 down to 50 dimensions with PCA. After training, model outputs are mapped through the (pseudo-
1299)inverse of this PCA transform, then converted back into $-\log_{10}(\text{P-value})\text{sign}(\text{lfc})$.

1300 E.2.4 Control methods

1301 We implemented six control methods described below:

- 1302 1. **Ground truth** (id: ground_truth): Return the test set as output.
- 1303 2. **Constant zero** (id: zeros): Predict no differential expression for any of the samples.
- 1304 3. **Random sample** (id: sample): Randomly sample counts from the training set per gene.
- 1305 4. **Mean outcome**: We used three average-based baselines. One that averages over all of the
1306 compounds and cell types $\hat{y}_{ij} = \sum_{i=1}^R y_{\text{train}_{ij}}$ (id: mean_outcome), one that averages across
1307 all of the cell types for a given compound (id: mean_across_compound), and one that
1308 averages across all of the compounds for a given cell type (id: mean_across_celltypes).

Table 3: Comparison of coefficient of variation across Kaggle competition cell types. We observe high variation in T cells CD8+ and T regulatory cells in control compounds.

Cell type	Dabrafenib	Belinostat	Dimethyl Sulfoxide
B cells	0.319051	0.338520	0.307461
Myeloid cells	0.184550	0.275649	0.185540
NK cells	0.240455	0.577534	0.222283
T cells CD4+	0.129801	0.162064	0.106406
T cells CD8+	0.488251	2.288442	0.498569
T regulatory cells	0.411224	1.894598	0.317219

1309 F Competition learnings

1310 F.1 Participant survey

1311 We surveyed 35 competitors to learn more about the participants’ backgrounds and their experience
 1312 of the competition. 57% of respondents haven’t worked with single-cell data before, and the same
 1313 number never participated in a Kaggle competition before. 91% have not participated in an Open
 1314 Problems competition before. The respondents come from 16 different countries. 31% work in
 1315 industry, and 54% in academia. Only 9% used other single-cell datasets, and 3% used external
 1316 references (e.g. KEGG or Gene Ontology) in their solutions.

1317 F.2 Outlier compounds

1318 One of the 20 clusters identified by the Leiden algorithm (Appendix A.8) could not be conclusively
 1319 labeled as belonging to any particular cell type. Over 96% of the cells in this cluster were from the
 1320 wells of three compounds (Delanzomib, Oprozomib, and MLN2238), all of which shared the same
 1321 mechanism of action, proteasome inhibition. To avoid biasing the perturbation prediction models
 1322 with low-confidence cell type labels, these three compounds were removed from the dataset. Due to
 1323 either low counts induced by toxicity or high variability in cell type proportions across replicates,
 1324 three other compounds were also dropped: CGP60474, BXU45ZH6LI, and Alvocidib.

1325 G Single-cell perturbation prediction evaluation

1326 Single-cell perturbation models can also be applied to our benchmark task. According to a recent
 1327 single-cell perturbation benchmark, PerturBench, a latent additive model performs best in this
 1328 category [55]. We used the parameters from the PerturBench run that performed best on the sci-Plex
 1329 dataset [17]. We trained the model on unnormalized counts. We then used the `limma` package for
 1330 differential expression analysis on the predicted counts, and the resulting outcomes were used as
 1331 model predictions. The latent additive model performed worse than our benchmark control methods
 1332 according to mean row-wise RMSE and mean row-wise MAE (Table 4).

1333 H Data analysis reproducibility

1334 The code for reproducing the figures and data analysis, including cell type annotation and filtering,
 1335 is available at github.com/theislab/task-dge-perturbation-prediction-analysis. The
 1336 code is provided under MIT license.

1337 I Benchmark details

1338 Benchmark code is available at github.com/openproblems-bio/task_perturbation_prediction,
 1339 DOI:10.5281/zenodo.11537124. The code of the benchmark is provided under MIT license.

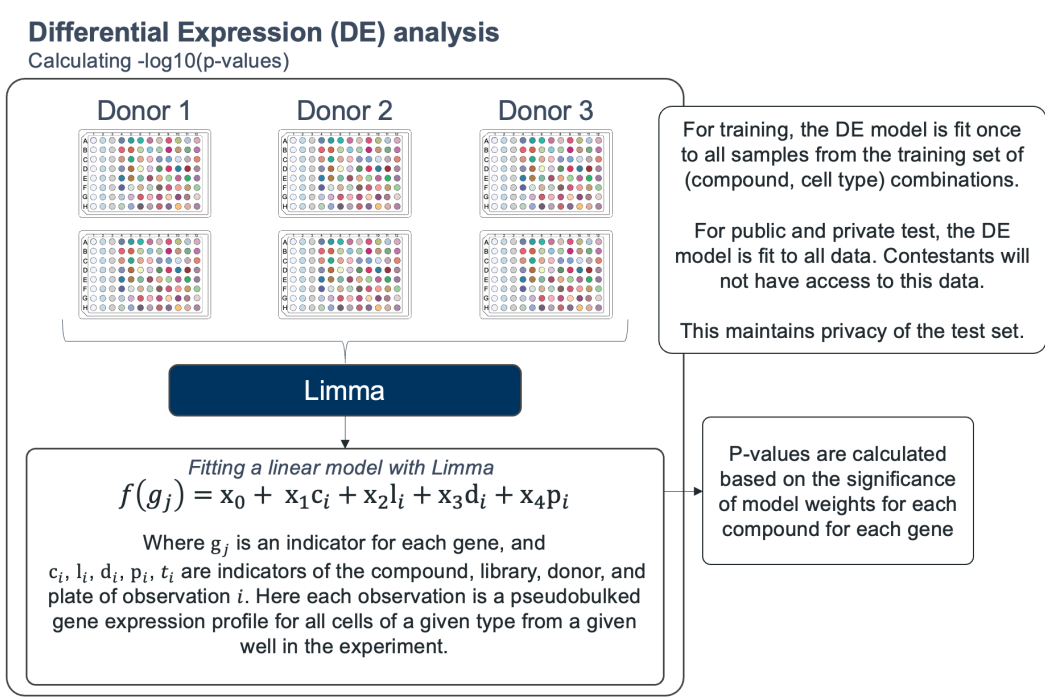


Figure 6: High-level overview of the Kaggle competition dataset DE computation, including the design matrix.

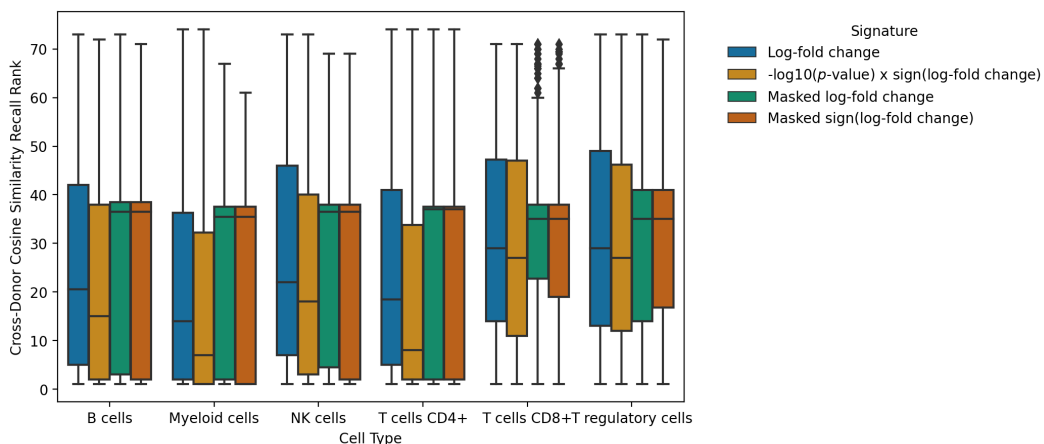


Figure 7: Cross-donor retrieval on the Kaggle competition dataset with cosine-similarity as a metric. The scores of T cells CD8+ and T regulatory cells stand out.

Table 4: Latent additive model comparison to OP3 benchmark models, sorted by mean row-wise RMSE.

Model	Mean rowwise RMSE	Mean rowwise MAE
Ground truth	0.0000	0.0000
NN retraining with pseudolabels	0.7562	0.5464
LSTM-GRU-CNN Ensemble	0.7921	0.5756
Py-boost	0.7957	0.5609
Mean per cell type and gene	0.8925	0.6437
JN-AP-OP2	0.8965	0.6518
Mean per gene	0.8992	0.6356
Zeros	0.9179	0.6351
Mean per compound and gene	0.9428	0.6979
Latent additive	1.162	0.8223

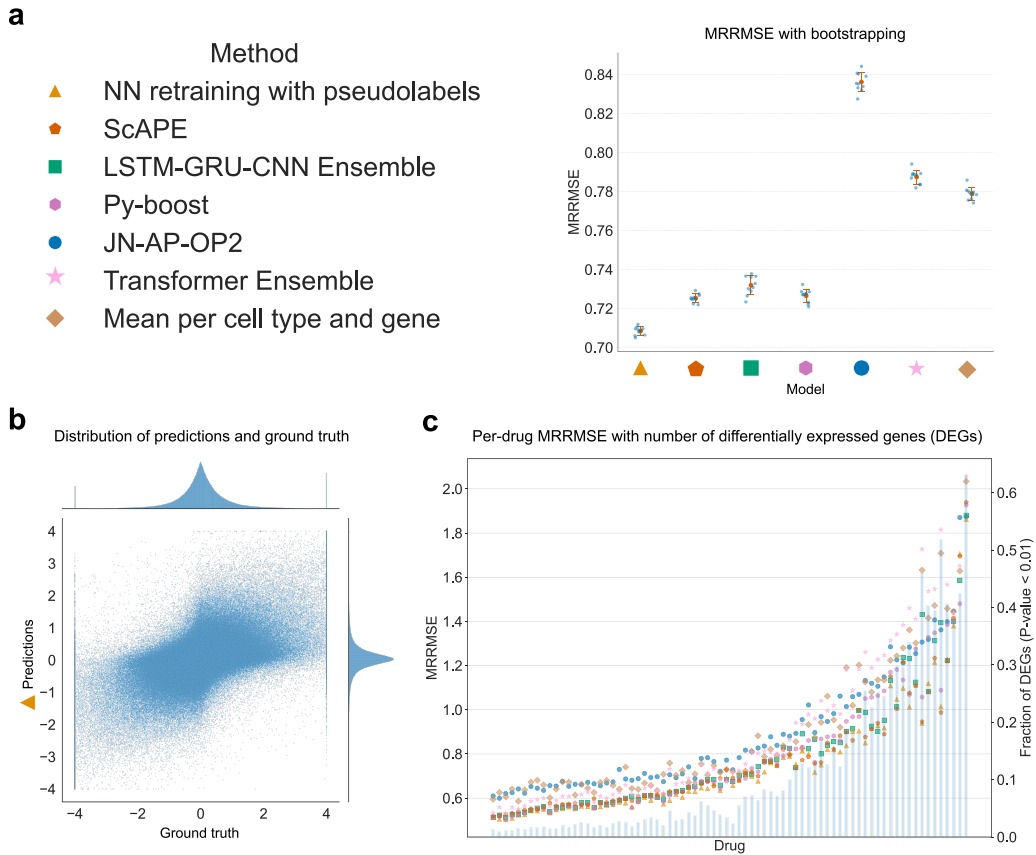


Figure 8: **Benchmark results.** (a) Results of rerunning the methods with dataset bootstrapping with 10 bootstraps. The error bars are standard deviation. Note that bootstrapping was performed by sampling cells in both the training and test sets. (b) Distribution of ground truth and the predictions of the top-performing method, NN retraining with pseudolabels. We note that the predictions are biased toward lower than true significance. (c) Per-drug MRRMSE and the fraction of genes for a given compound with a P-value lower than 0.01 (the latter shown with a bar chart). We note that the errors are larger for compounds with a high fraction of DEGs. The differences in errors across the methods and the baseline are smaller in samples with a low fraction of DEGs.

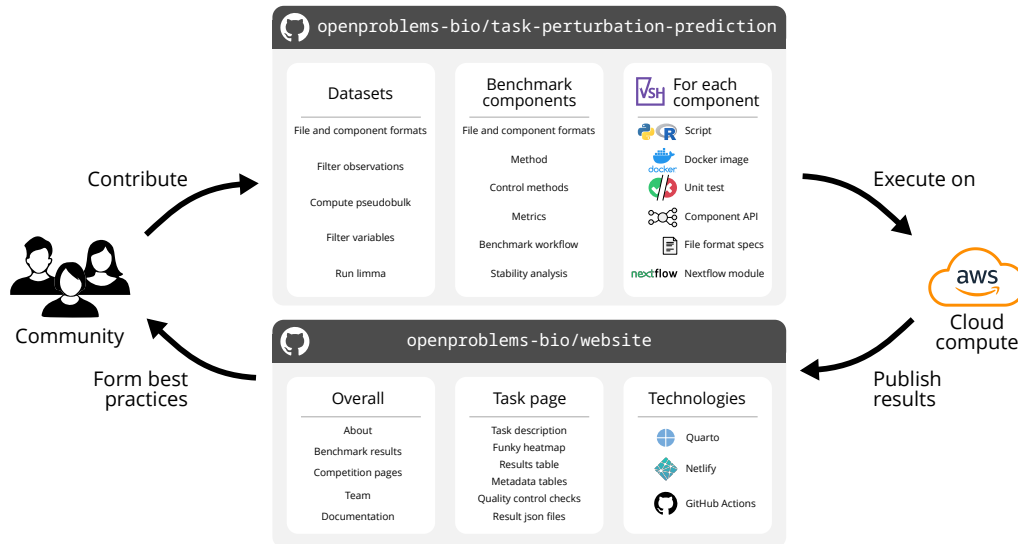


Figure 9: An overview of the technology stack of the perturbation prediction living benchmark within the OpenProblems ecosystem.

1340 **Data formats:** To ensure interoperability between components, the repository uses AnnData [56] as
 1341 the standard data format for both input and output files of components, and strict requirements are
 1342 imposed on the format of these files.

1343 **Components:** Workflows are comprised of Viash components and are themselves also Viash compo-
 1344 nents [57]. A Viash component is a small amount of metadata combined with a script implemented
 1345 in Python, R, Bash, or Nextflow. Viash can use this information to build a component-specific
 1346 Docker container, and turn the component into a Docker-backed Nextflow workflow. These Nextflow
 1347 workflows can be used as a standalone module, or as a submodule for another workflow.

1348 I.1 Workflows

1349 The repository consists of three main workflows: `process_dataset`, `run_benchmark`, and
 1350 `run_stability_analysis` (Figure 10).

1351 I.2 Workflow: Process dataset

1352 The data processing steps used to transform the single-cell RNA-seq expression matrix into the
 1353 Perturbation Differential Gene Expression (DGE) matrix for the perturbation prediction task (Figure
 1354 10 top). It consists of the following components:

- 1355 • **Filter obs:** Remove low-quality observations from the dataset. The conditions are designed
 1356 to exclude cells that could introduce bias or noise into the downstream analysis, such as
 1357 cells from certain donors, cells treated with certain molecules, or certain cell types.
- 1358 • **Compute pseudobulk:** Aggregate cell types into pseudobulks.
- 1359 • **Filter vars:** Subset the genes
- 1360 • **Limma on train:** Run limma on the train and control splits, per cell type and per small
 1361 molecule. The resulting information is stored as an AnnData object we call DE train.
- 1362 • **Limma on train and test:** Run limma on train, control and test split, per cell type and per
 1363 small molecule. The resulting information is stored as an AnnData object we call DE test.
- 1364 • **Extract ID map:** Extract a data frame containing a combination of the cell types and small
 1365 molecules which methods will need to predict. The resulting information is called ID map.

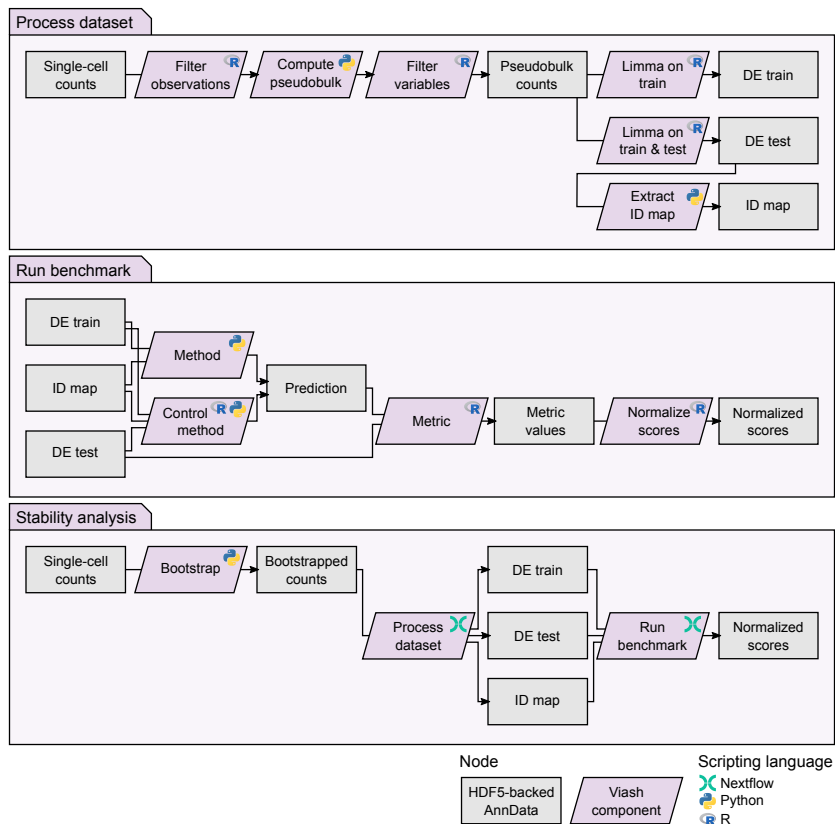


Figure 10: The different workflows used to perform the analyses in this study, `process_dataset`, `run_benchmark`, and `run_stability_analysis`. Each workflow uses HDF5-backed AnnData (h5ad) files (grey rectangle) as a common data format, and is comprised of Viash components (purple rhombus) implemented in Nextflow, Python, or R. Since each workflow is also a Viash component, it can in turn be used as a subworkflow of a larger workflow.

1366 **I.3 Workflow: Run benchmark**

1367 Evaluate the performance of methods and control methods using a set of metrics (Figure 10 middle).
 1368 This workflow accepts the DE train, DE test and ID map objects and inputs and runs the various
 1369 control methods and methods on it. Each prediction generated by the methods is evaluated using
 1370 each of the metrics. In the end, all output results is stored, alongside the dataset metadata, method
 1371 metadata, metric metadata, and runtime resource information. The workflow consists of the following
 1372 components:

- 1373 • **Method:** A method for predicting the perturbation response of small molecules on certain
 1374 cell types.
- 1375 • **Control method:** A control method to serve as a quality control for the perturbation
 1376 prediction benchmark.
- 1377 • **Metric:** A metric to compare a perturbation prediction to the ground truth.
- 1378 • **Normalize scores:** Normalise the metric values by min-max scaling the values between the
 1379 worst control method result and the best control method result.

1380 **I.4 Workflow: Stability analysis**

1381 This workflow is used to perform a stability analysis of the methods (Figure 10 bottom). It bootstraps
1382 the original single-cell counts matrix, and runs the Process dataset and Run benchmark workflows to
1383 perform a benchmark on each of the bootstrapped datasets. It consists of the following components:

- 1384 • **Bootstrap:** This component bootstraps the single-cell dataset by sampling the same number
1385 of cells with replacement from the dataset.
- 1386 • **Process dataset:** The process dataset workflow mentioned earlier.
- 1387 • **Run benchmark:** The run benchmark workflow mentioned earlier.

1388 **I.5 Execution environment**

1389 Workflows were executed on AWS Batch, where components could run completely in parallel
1390 depending on the topology of the workflow. Components were run on different instance types
1391 depending on the specific memory / CPU / GPU requirements of the component. The following is a
1392 list of suitable instance types depending on the requirements of the component:

- 1393 • GPU required: g4dn.8xlarge, 32 vCPUs, 128 GB memory, 1 Nvidia T4 GPU.
- 1394 • Low memory: m4.2xlarge, 8 vCPUs, 32 GB memory
- 1395 • Medium memory: m4.4xlarge, 16 vCPUs, 64 GB memory
- 1396 • High memory: m4.10xlarge, 40 vCPUs, 160 GB memory

1397 All method components required a GPU to run, whereas dataset processing components, control
1398 methods, and metrics did not require a GPU to run.

1399 A run of the run_benchmark workflow requires:

- 1400 • 216 jobs on non-GPU and GPU instances
- 1401 • Wall time: 3h 2m
- 1402 • CPU time: 173 CPU hours
- 1403 • Total memory: 232 GB
- 1404 • Disk read: 24 GB
- 1405 • Disk write: 27 GB

1406 A run of the run_stability_analysis requires the following resources:

- 1407 • 1271 jobs on non-GPU and GPU instances
- 1408 • Wall time: 6h 52m
- 1409 • CPU time: 2162 CPU hours
- 1410 • Total memory: 3101 GB
- 1411 • Disk read: 321 GB
- 1412 • Disk write: 440 GB

1413 **J Informed consent for PBMC donors**

1414 For this study, we purchase commercially available human tissue samples from AllCells, Inc. AllCells
1415 is a tissue bank licensed by the State of California Department of Public Health, USA (Tissue Bank
1416 ID#: CTB 00080812). AllCells is responsible for maintaining IRB approval for all human subjects
1417 research. Below is one of the informed consent documents signed by one of the donors (name and
1418 signature redacted). More information is available from AllCells upon request.

1419

Discovery Life Sciences, LLC
DLS-BB018-V.11

RESEARCH SUBJECT INFORMATION, CONSENT, AND AUTHORIZATION FORM

TITLE:	Collection of Human Apheresis Specimens from Healthy Donors for Future Scientific and Medical Research
--------	--------------------------------------------------------------------------------------------------------

This consent form contains important information to help you decide whether to take part in a research study.

The study staff will explain this study to you. Ask questions about anything that is not clear at any time. You may take home an unsigned copy of this consent form to think about and discuss with family or friends.

- **Being in a study is voluntary – your choice.**
- **If you join this study, you can still stop at any time.**
- **Do not join this study unless all your questions are answered.**

After reading and discussing the information in this consent form you should know:

- Why this study is being done;
- What will happen during the study;
- Any possible benefits to you;
- The possible risks to you;
- Other options you could choose instead of being in this study;
- How your personal health information will be treated, used, and **disclosed** during the study and after the study is over;
- Whether being in this study could involve any cost to you; and
- What to do if you have problems or questions about this study.

Please read this consent form carefully.

RESEARCH SUBJECT INFORMATION, CONSENT, AND AUTHORIZATION FORM

TITLE: Collection of Human Apheresis Specimens from Healthy Donors for Future Scientific and Medical Research

PROTOCOL NO.: DLS-BB018-V.11
IRB Protocol #20130996

SPONSOR: Discovery Life Sciences, LLC

INVESTIGATOR: Timothy M. Howard, MD
800 Hudson Way
Huntsville, Al 35806
USA

SITE(S): Discovery Life Sciences, LLC 800 Hudson Way
Huntsville, Al 35806
USA
American Red Cross
100 Peartree Lane
Raleigh, NC 27600
USA

American Red Cross
1101 Washington St NW
Huntsville, Al 35801
USA
American Red Cross
2425 Park Road
Charlotte, NC 28203
USA

American Red Cross
700 Caldwell Trace
Birmingham, Al 35242
USA
American Red Cross
2751 Bull Street
Columbia, SC 29230
USA

American Red Cross
2179 Roswell Road
Marietta, GA 30062
USA
American Red Cross
100 Rustcraft Road
Dedham, MA 02026
USA

American Red Cross
337 Stoneridge Lane
Gahanna, OH 43230
American Red Cross
7539 Oswego Road
Liverpool, NY 13090

**STUDY-RELATED
PHONE NUMBER(S):** Discovery Life Sciences, LLC
Study Coordinator
256-327-9828 (24 Hours)

1421

SUMMARY

We invite you to take part in a research study (the “Study”). The purpose of this consent and authorization form (the “Consent”) is to help you decide if you want to be in the Study, and if you agree to have your health information used and disclosed for the Study. This Consent may contain words that you do not understand. Please ask the Study staff to explain any words or information that you do not clearly understand. You may have this Consent read to you.

Things to know before deciding to take part in a research study:

- The main goal of a research study is to learn things to help patients in the future.
- The main goal of regular medical care is to help each patient.
- Basic health information will be collected during the time you are taking part in this Study. This health information may be looked at and/or copied by Discovery, government agencies, and/or other groups associated with the Study.

If you take part in this research study, you will be given a signed copy of this Consent.

PURPOSE OF THE STUDY

You are being invited to take part in a Study because human apheresis samples (for example, white blood cells and plasma) are needed to support research. Research on samples and health information can help scientists discover more about what causes diseases, how to prevent them, and how to cure them. The Study specifications will be provided to the subject by the Study Doctor.

Apheresis is the process of collecting particular parts of the blood (e.g. white blood cells, plasma, platelets) by passing the blood through an apheresis machine. The machine separates the part of the blood, collects those that are necessary, then returns the remainder of the blood back to the donor.

DURATION OF THE STUDY

If you decide to take part in the Study, your participation is expected to last indefinitely or until you choose to no longer take part. There is no limit to the number of donors enrolled in the Study. The total number of donors expected to take part is unknown. Only adult donors will be included in this Study.

PROCEDURES

If you decide to take part in the Study, after you sign this Consent, you will be required to complete a general health interview and meet the specified inclusion criteria to be eligible for the apheresis procedure.

General Health Interview

Your vital signs (height, weight, blood pressure, pulse, and temperature) will be taken and recorded. Blood will be drawn and tested for blood counts (complete blood count and retic), metabolic function (comprehensive metabolic panel and hemoglobin A1C), lipid measurements (lipid panel), blood type (ABO/Rh), blood antibodies via direct and indirect antiglobulin tests, pregnancy, HLA typing, and the following list of diseases, as applicable: Covid-19, human immunodeficiency virus (HIV), hepatitis C virus (HCV), hepatitis B virus (HBV), hepatitis A

1422

virus (HAV), herpes simplex virus 1 (HSV-1), herpes simplex virus 2 (HSV-2), Varicella-zoster virus (VZV), Epstein-barr virus (EBV), cytomegalovirus (CMV), human herpesvirus 6 (variants A & B), human herpesvirus 7, Kaposi's sarcoma virus, West Nile virus (WNV), relevant cell-associated communicable disease agents and diseases (including human t-cell lymphotropic virus (HTLV)), human transmissible spongiform encephalopathy (including Creutzfeldt-Jakob disease), *Trypanosoma cruzi*, *Treponema pallidum*, syphilis, malaria, Zika Virus, and Parvovirus. Exclusive of the list above, any additional testing that is necessary to satisfy the project-specific inclusion and exclusion criteria may also be performed as long as the test does not require reporting to federal or state agencies. Women of childbearing potential will be tested for pregnancy. Pregnancy testing may be performed on blood, or a urine sample may be requested.

- If you test negative for all tested diseases and pregnancy, as applicable, and you meet the inclusion criteria as described below, you will be eligible for an apheresis procedure.
- If you have a positive test result for any tested disease, you will be removed from enrollment. If you have a positive pregnancy test, as applicable, and/or do not meet the inclusion criteria as described below, you will be enrolled in the Study but will not be eligible for an apheresis procedure at this time.
- If any of your disease tests are positive, the Sponsor will follow all applicable laws in the notification to appropriate agencies. The Sponsor will notify you of the positive result and request you seek follow up care with your general practitioner.

Apheresis Procedures

If you are eligible for an apheresis procedure you will be scheduled for the procedure within 3 weeks of the general health interview. Your vital signs (blood pressure, temperature, pulse) will be taken before the apheresis procedure begins. You will have an intravenous (IV) line (either a needle or catheter) placed in both arms. A nurse will monitor you and your vital signs will continue to be taken throughout the procedure. The apheresis procedure will take approximately four to five hours.

Once enrolled you are eligible to donate no more than one procedure of apheresis collection every sixty days (60) as long as you meet the inclusion/exclusion criteria described below. Mononuclear cells and other apheresis samples may be collected for a total product volume of no more than 550 mLs (a little less than 2 ½ cups).

The samples may be kept by the Sponsor in a bank and stored indefinitely.

In addition to the qualification and apheresis procedures, donor information about you (for example, age, race, and gender) and your pertinent medical information will be obtained from you or your donor record. This information will be linked to your samples. However, before the samples and information are released to any researcher, they will be given a special code without your name or private information on them that directly identifies you. The Site, Sponsor, Study Doctor, and Study staff may have access to the key that links this special code to your private information. However, no researchers will have access to your directly identifiable private information through this Study.

This Consent allows for more than one collection during your participation.

The following procedures may be performed in this Study:

- Apheresis collection – Blood component separation procedure in which whole blood is removed from your vein and passed through a device that separates the blood into components. Particular components are collected, and the remaining components are returned back to you. Up to four blood volumes can be collected once every sixty (60) days.
- Nasal swab(s) collection – A procedure in which a sample of nasal secretions is taken. This is usually performed by wiping the inner nostril with a cotton-tipped swab.
- Nasopharyngeal swab(s) collection – A procedure in which nasal secretions from the back of the nose and throat is taken. This is usually performed by inserting a cotton-tipped swab into the nostril and rotating over the surface of the posterior nasopharynx.
- Urine collection – A procedure in which urine is collected in a sterile, plastic container.
- Venipuncture – A procedure in which blood is removed from one of your veins using a needle

INCLUSION CRITERIA

- Age 18-70 years old (must be a legal adult in state of the Site)
- Weigh at least 110 lbs
- Baseline Blood Pressure: Systolic: 90 -180 mm Hg, Diastolic: 50-100 mm Hg
- Temperature: less than 99.5°F
- Pulse rate: 50-110 beats/minute and regular
- Negative for all tested diseases as listed in the General Health Interview section
- Hemoglobin:
 - Females: no less than 11.5 g/dcL
 - Males: no less than 12.2 g/dcL
- Hematocrit:
 - Females: No less than 35.2%
 - Males: No less than 38.2%

EXCLUSION CRITERIA

- Donors who do not give informed consent
- Donors who do not understand the informed consent
- Women who are pregnant or breastfeeding
- Donors with any history of heart, lung, liver, or kidney disease
- Donors with any history of blood or bleeding disorders, including sickle cell disease
- Donors with any history of neurologic disorders
- Donors with any history of cancer
- Donors with any history of diabetes
- Donors with a positive test result for any disease tested for as listed in the General Health Interview section
- Steroid use within two weeks of apheresis procedure

RISKS AND DISCOMFORTS

There are potential unforeseen risks with any procedure. The known potential risks are as follows:

- Apheresis – potential risks include:
 - Citrate toxicity: muscle cramping, numbness, chills, tingling sensations. Citrate toxicities are managed symptomatically using oral calcium supplements.
 - Bleeding, bruising, irritation, infiltration, inflammation at the venipuncture sites, or risk of arterial puncture
 - Allergic reaction
 - Vasovagal episode: lightheadedness, hot flashes, nausea, vomiting, decreased heart rate, and decreased blood pressure.
 - Syncope: fainting, risk of injury/fall
 - Hyperventilation
 - Infection at venipuncture site
 - Air embolus from machine malfunction: gas bubble enters the blood stream
 - Long term effects of donor apheresis are unknown
- Nasopharyngeal swab(s) collection – potential discomfort or pressure is associated with this procedure.
- Venipuncture - potential risks include pain, bruising, lightheadedness, or, on rare occasions, infection.
- There are no known risks associated with nasal swab and urine collections. However, there may be infectious pathogens that can be spread to others. Hands should be washed thoroughly with antibacterial soap after collection of these biospecimens. There may be minor bleeding, bruising, or discomfort from the nasal swab.
- Confidentiality – There is a possible loss of confidentiality of your health information, although all reasonable efforts will be made to protect your information as described in this Consent.

Due to scientific advances or human error, your identity and health information may become known. Since DNA (the chemical that makes up genes) information is unique to you, in the future this link could occur. For this link to occur, it would require someone to take another sample from you, analyze the DNA, and compare it with the data resulting from this research project.

COMPENSATION FOR INJURY

If you are injured, you should obtain treatment as you would for any other injury, or you may contact your Apheresis Nurse/Study Doctor who can refer you for treatment. There are no plans to compensate you for any injuries you suffer as a result of this Study.

USE OF SAMPLES AND INFORMATION

Samples may be used to explore possible links between different types of molecules (for example, DNA, RNA, proteins) and features of the people (for example, age, gender, family history of certain medical conditions). The medical conditions studied will be widespread including some that you may not have. None of the results will be linked directly to you. They will be linked only to the group of people. Researchers may perform a variety of tests including genetic tests, tests of the cells that make up your samples, DNA or RNA sequencing or gene editing and even future medical research that is currently unknown at this time.

Your samples may be stored in ways that allow the cells to grow and multiply. These multiplying cells may grow to what is called a cell line. Cell lines can be used for many future studies and these cells may be kept alive for many years. None of your donated samples will be infused into another human being.

Researchers may develop products based on things they learn from your samples. Any information obtained by the researchers as a result of testing your samples will not be provided to you, as applicable. Any applicable information provided to you will come from the Study Doctor. These researchers will use your samples as needed and destroy unused portions per government regulations. The tests done on your samples are for research purposes only.

NEW INFORMATION

You will be told about any new information that might change your decision to be in this Study. You may be asked to sign a new Consent if this occurs.

BENEFITS

If you agree to take part in this Study, there will be no direct medical benefit to you.

COSTS

There are no costs to you for taking part in the Study.

COMPENSATION FOR PARTICIPATION

You will be compensated for the time and effort you devote to this Study. The compensation for taking part is up to \$1000.

The site where the procedure is performed will be reimbursed in accordance with separately negotiated agreements between the Site and the Sponsor.

COMMERCIAL USES

Any samples you provide that are used in research may result in new products, tests, or discoveries. In some instances, these developments may have commercial value. There are no plans for you to share in any financial benefits from these products, tests, or discoveries.

1426

ALTERNATIVE TREATMENT

The Study is for research purposes only. The only alternative is to not take part in this Study.

VOLUNTARY PARTICIPATION AND WITHDRAWAL

Your participation in this study is voluntary. You may decide not to take part or you may leave the Study at any time. Your decision will not result in any penalty or loss of benefits to which you are otherwise entitled.

You may withdraw from taking part in the Study at any time. You do this by providing written or verbal notification to your Study Doctor or Study Staff. If you withdraw your permission, you will not be able to continue taking part in this Study. Upon withdrawal, information that has already been gathered and samples already distributed before the date of withdrawal may still be used to make the research reliable. Remaining samples collected during the period of time you had given consent may be used for research. Information that has already been gathered will be maintained to ensure the accuracy of the research.

Your participation in this Study may be stopped at any time by the Study Doctor or the Sponsor without your consent.

SOURCE OF FUNDING FOR THE STUDY

Funding for this Study is provided by Discovery Life Sciences, LLC, the Sponsor.

QUESTIONS

If you have questions, concerns, or complaints, or think this research has hurt you or made you sick, talk to the research team at the phone number listed above on the second page.

This research is being overseen by an Institutional Review Board (“IRB”). An IRB is a group of people who perform independent review of research studies. You may talk to them at 855-818-2289 or researchquestions@wcgirb.com if:

- You have questions, concerns, or complaints that are not being answered by the research team.
- You are not getting answers from the research team.
- You cannot reach the research team.
- You want to talk to someone else about the research.
- You have questions about your rights as a research subject

AUTHORIZATION TO USE AND DISCLOSE INFORMATION FOR RESEARCH PURPOSES

Federal regulations give you certain rights related to your health information. These include the right to know who will be able to get the limited information and why they may be able to get it. The Study doctor must get your authorization (permission) to use or give out any health information that might identify you.

What information may be used and given to others?

If you choose to be in this Study, the Study doctor will get limited personal information about you. This may include information that might identify you. The Study doctor may also get limited information about your health including:

- Past, present, and future medical records
- Research records
- Questionnaire information collected as part of the Study
- Records about your Study visits
- Disease registry information.

Who may use and give out information about you?

The limited information about your health may be used and given to others by the Study Doctor, Study staff, or the Sponsor. They might see the research information during and after the Study.

Who will get this information?

The Sponsor of this Study will have access to your limited personal and medical information. Sponsor means any persons or companies that are:

- working for or with the Sponsor, or
- owned by the Sponsor.

Researchers will receive certain limited information about you. This limited information will not directly identify you.

The limited information about you and your health, which might identify you, may be given to:

- The U.S. Food and Drug Administration (FDA),
- Department of Health and Human Services (DHHS) agencies,
- Governmental agencies in other countries, and
- Institutional Review Board (IRB).

Why will this information be used and/or given to others?

The Sponsor will analyze and evaluate the results of the Study. The Sponsor will be visiting the research site. They will follow how the Study is done, and they will be reviewing your limited information for this purpose.

The limited information about you may be given to researchers to carry out the Study, but your identity will not be disclosed.

The limited information about you may be given to the FDA. It may also be given to governmental agencies in other countries. The limited information may be used to meet the reporting requirements of governmental agencies.

The results of this research may be published in scientific journals or presented at medical meetings, but your identity will not be disclosed.

1428

What if I decide not to give permission to use and give out my limited health information?

Then you will not be able to be in this Study.

May I review or copy the limited information obtained from me or created about me?

Yes, but only after the Study is closed.

May I withdraw or revoke (cancel) my permission?

Yes, but this permission will not stop automatically.

You may withdraw or take away your permission to use and disclose your limited health information at any time. You do this by notifying the Study Doctor or Study staff in writing or verbally. If you withdraw your permission, you will not be able to stay in the Study.

When you withdraw your permission, no new health information identifying you will be gathered for the Study after that date. Once the Sponsor receives your withdrawal notice, it will not further disclose your limited information, but it may still use the limited information to make the Study reliable.

However, your withdrawal will not affect any action that has already been taken in reliance on your authorization. For example, if the Sponsor has already released your limited information to another researcher for future use, it may continue to be used and disclosed, and it will not be possible to get the limited information back.

Is my limited health information protected after it has been given to others?

The Sponsor has processes in place to protect your limited identifying information; for example, your name is replaced by a number and you are referenced only by that number with others who do not have the ability to tie that number back to your name. However, there is a risk that your limited information will be released to others who may not have the same legal obligation to protect that limited information.

When does my permission to use my limited information expire?

There is no current plan to end the Study. Your limited information may be held in a repository (or multiple repositories) indefinitely, and your permission to use this limited information will not expire.

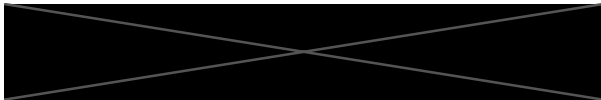
1429

CONSENT TO PARTICIPATE IN THE STUDY

I have read this Consent (or it has been read to me). All my questions about the Study and my part in it have been answered. I freely and voluntarily consent to take part in this Study.

By signing this Consent, I give permission for my samples and limited health information to be used and stored for current and future research of my medical diagnosis or other medical diagnoses.

By signing this consent form, I have not given up any of my legal rights.



Signature of Subject

5.2.22
Date

12:59 p
Time



Subject's Name (Printed)

PERSON CONDUCTING INFORMED CONSENT DISCUSSION:

I confirm that the Study was thoroughly explained to the subject, including but not limited to the risks and benefits of participation, and that it is voluntary. I reviewed the Consent with the subject and answered the subject's questions. The subject appeared to have understood the information with verbal recall about the Study upon my questioning.


Signature of Person Conducting the
Informed Consent Discussion

5/2/22
Date

12:59 PM
Time


Printed Name of Person Conducting the
Informed Consent Discussion

----- **Use this witness section only if applicable** -----

If this Consent is read to the donor because the donor is unable to read the Consent, an impartial witness not affiliated with the research or investigator must be present for the consent and sign the following statement:

I confirm that the information in the Consent and any other written information was accurately explained to, and apparently understood by, the donor. The donor freely consented to be in the Study.

Signature of Impartial Witness

Date

Time

Printed Name of Impartial Witness

Note: This signature block cannot be used for translations into another language. A translated Consent is necessary for enrolling donors who do not speak the language of this consent.