



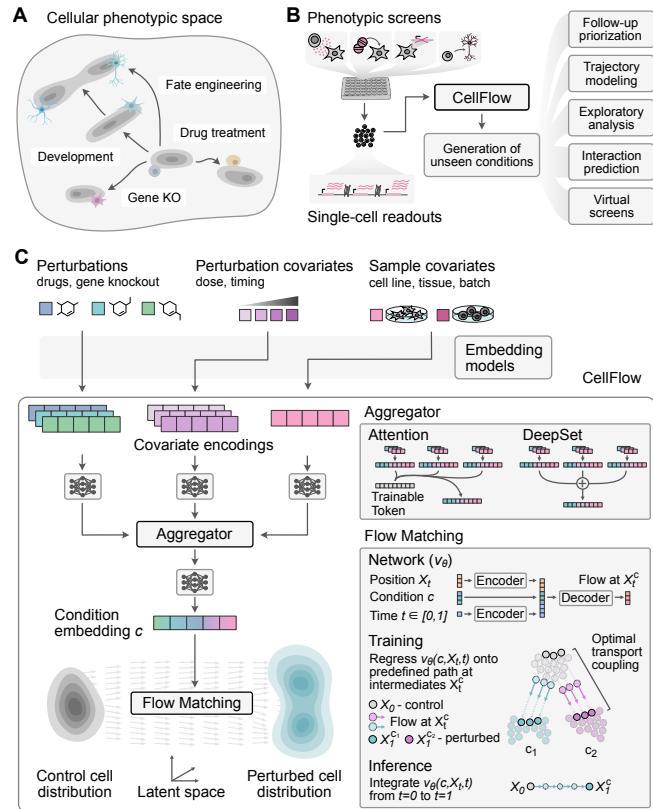
quality samples in computer vision<sup>20,21</sup>, video<sup>22</sup>, and molecular design<sup>23,24</sup>, in part through modeling complex distributions. We use optimal transport<sup>25–29</sup> to pair unperturbed and perturbed cells, enabling the distinction between inherent cellular heterogeneity and perturbation-induced distributional changes. To represent arbitrary types of perturbations and enable predictions for conditions out-of-distribution, CellFlow incorporates powerful pre-trained embeddings of biological entities<sup>30</sup>. To model arbitrary numbers of perturbations in a permutation-invariant manner, we employ set aggregation strategies including multihead attention, a key factor to foster the success of large language models<sup>31</sup>.

We demonstrate the capabilities and flexibility of CellFlow across a wide range of phenotypic screening applications. We show that it accurately models donor-specific cellular responses to cytokine treatment on a large perturbation dataset of almost ten million Peripheral Blood Mononuclear Cells (PBMCs)<sup>8</sup>. To demonstrate CellFlow's ability to model highly complex cellular distributions, we predict single-cell expression profiles across entire zebrafish embryos perturbed with different gene knockouts at various developmental stages. We show state-of-the-art performance on established perturbation prediction tasks, including prediction of gene knockout and drug treatment effects. Finally, we show CellFlow's ability to predict heterogeneous cell populations resulting from neuron fate engineering and during organoid development. As a proof-of-principle, we perform a virtual organoid protocol screen, which identifies previously untested treatment regimens with strong effects on organoid development. Together, our results show that CellFlow can accelerate discovery from phenotypic screens by extrapolating to unseen conditions, thereby informing the prioritization of follow up experiments and enabling efficient experimental design approaches (Figure 1B).

## Results

**CellFlow is a general framework to model perturbed single-cell phenotypes** CellFlow aims to predict single-cell phenotypes under diverse perturbations by conditionally mapping a source distribution (e.g. control cells) to a perturbed population of cells. The model first encodes experimental variables and aggregates combinatorial treatments into a common condition embedding, which is then injected into the flow matching module to guide the flow from source to perturbed distributions (Figure 1C, Methods).

To accommodate a wide range of phenotypic screening scenarios, we define a generic experimental setup with three types of variables: perturbations, perturbation covariates, and sample covariates. Perturbations represent observed experimental interventions, such as drug treatments or genetic modifications. To allow predic-

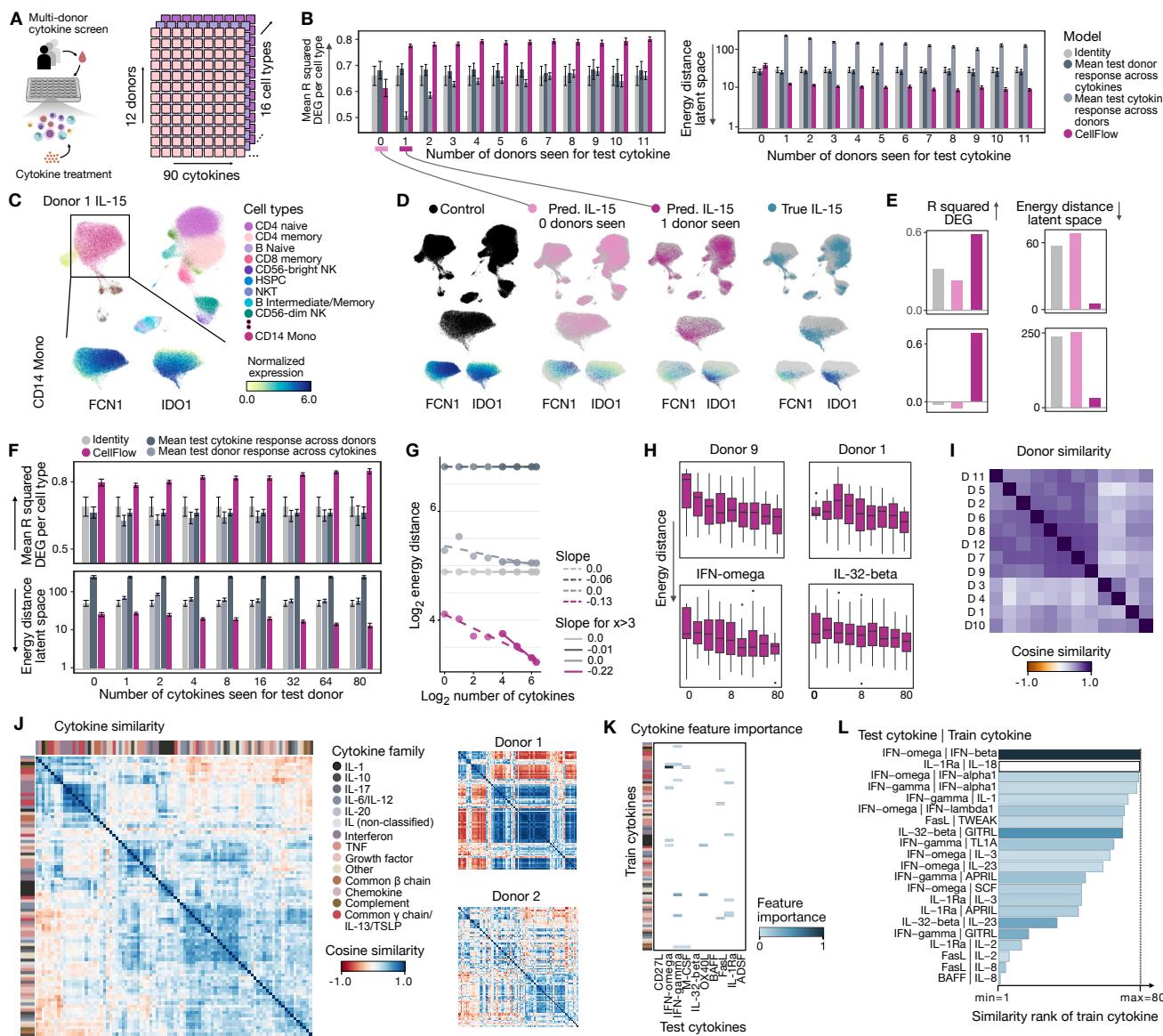


**Figure 1. CellFlow: a tool to explore cellular phenotypic space (A)** Cells can change their phenotype in response to various stimuli, such as drug treatment, gene knockout or developmental cues. **(B)** Phenotypic screens help to explore the relationships between perturbation and phenotype. CellFlow can learn from such screens and generate phenotypes of unseen conditions to facilitate various downstream tasks. **(C)** CellFlow takes into account various variables of a generic experimental setup and uses aggregations strategies such as attention mechanisms and DeepSets to obtain an aggregated condition embedding. We use flow matching, which regresses predicted flows onto a predefined probability path, to learn a conditional vector field transforming a source distribution onto the perturbed distribution.

tions for interventions that were unseen during training, we employ encodings, such as molecular fingerprints<sup>32</sup> for small molecules and ESM2<sup>30</sup> embeddings for the protein products of gene knockouts and proteins. Perturbation covariates represent additional modulators of treatments such as dosage or timing whereas sample covariates represent cell descriptors, independent of their treatment, such as the cell line or tissue. We refer to a combination of perturbations, perturbation covariates, and sample covariate as a single condition.

To model multiple interventions at the same time, it is necessary to encode combinations of treatments in a permutation-invariant manner. Therefore, CellFlow implements two set aggregation schemes, multi-head attention and deep sets, followed by an encoder consisting of a feed-forward network. This allows CellFlow to represent an arbitrary combination of perturbation variables in a single representation.

The condition embedding vector then serves as input into the flow matching module, which translates a source distribution to the perturbed cell population.



**Figure 2. CellFlow exhibits interpretable and scalable training behaviour on 10 million PBMCs.** (A) Peripheral Blood Mononuclear Cells (PBMCs) from twelve different donors are treated with 90 different cytokines, resulting in a dense experimental design capturing almost 10m cells. (B) Performance metrics for predicting donor-specific cytokine responses for varying numbers of donors for which the cytokine has been seen. The  $R^2$  of the top 50 differentially expressed genes (DEGs) per cell type is computed based on label transfer from real measurements to generated cells. One reported data point represents the mean of the  $R^2$  across cell types for a single condition (Methods). Plotted are the mean and the standard error across different training data sets with the same number of donors seen for the test cytokine (left). Analogously, the energy distance in latent space is reported for different numbers of conditions in the training data set (right, Methods). (C) UMAP of IL-15 treated PBMCs of donor 1, together with an inset of CD14 Mono cells colored by normalized gene expression of the top upregulated gene FCN1 and top down-regulated gene IDO1 (left). (D) Joint UMAP of control cells of donor 1, true IL-15 treated cells of donor 1, as well as generated cells of donor 1 treated with IL-15 from two different models. The first model was trained without IL-15 treated cells in the training dataset, while the other model included IL-15 treated cells of donor 8 in the training data set. Each column shows one of the four cell populations highlighted (top), together with an inset of CD14 Mono cells (middle), and the normalized gene expression of FCN1 and IDO1 (bottom). (E) Quantification of the visual results of the identity model (i.e. using the control distribution of donor 1 as predicted IL-15 treated cells), CellFlow trained with no IL-15 cells seen, and CellFlow trained with IL-15 treated cells from donor 8. (F) Performance of CellFlow and baselines for predicting the cytokine response of a new donor based on a varying number of cytokines measurements for the new donor. Mean and standard error are computed analogously to (b). (G) Linear fits of the mean energy distances displayed in (g) in log-log space. Slopes are computed from models trained on at least one cytokine for the test donor (dashed line) and trained on at least 16 cytokines for the test donor (solid line). (H) CellFlow's performance metrics filtered for predicted populations specific to donor 9 (top left) and donor 1 (top right), as well as specific to IFN-omega (bottom left) or IL-32 beta (bottom right). (I) Donor similarities computed from responses across all cytokines (Methods). (J) Cytokine similarities computed from responses across all donors (left), and from responses of donor 1 (top right) and donor 2 (bottom right, Methods). (K) Importance of presence of cytokine in training data set for a good performance of CellFlow based on the coefficients of a regularized linear model (Methods). (L) Ordered enumeration of non-zero importance scores displayed in (K).

To learn a meaningful mapping on the population and single-cell level, we assume cells follow the least-action principle in response to a perturbation, and thus lever-

age optimal transport to pair samples batch-wise. We train our models end-to-end, thus learning a space of condition representations, which guide the generation

of perturbed cell states from the source population. Training of flows is facilitated when the signal-to-noise ratio is high. As single-cell data can have sparse feature detection and technical variation, we leverage lower-dimensional cellular representations obtained from principal component analysis (PCA) or variational autoencoders (VAEs)<sup>33</sup>.

**Scalable prediction of cytokine responses in a large phenotypic screen** To assess the predictive capabilities and scaling properties of CellFlow on high-throughput phenotypic screens, we applied it to a recent cytokine perturbation screen on almost ten million Peripheral Blood Mononuclear Cells (PBMCs)<sup>8</sup>. This dataset comprises scRNA-seq samples obtained from twelve donors, each treated with each of 90 cytokines (Figure 2A). Due to their central role in the immune system, cytokines are a promising therapeutic target for numerous diseases like cancer and autoimmune disorders<sup>34</sup>.

We first investigated CellFlow's predictive performance in relation to the number of donors for which a cytokine treatment has been observed. For a given test combination of a donor and a cytokine, we systematically varied the number of donors treated with that cytokine included in the training data. For each number of donor-specific cytokine treatments included in the training data, we held out different sets of conditions during training (Methods). We evaluated the generative performance using the energy distance between the true and the predicted cell population in PCA space (Methods). To evaluate predicted gene expression, we transferred cell type labels of generated cells from true ones using the one-nearest neighbor classifier, and computed the R squared between the means of the true and the predicted gene expression for the top 50 cell type- and condition-specific differentially expressed genes (DEGs, Methods). To assess the ability to learn donor-specific cytokine responses beyond simple assumptions, we compared the performance against three baseline models: (1) an identity model assuming no effect (identity), (2) a model assuming a cytokine treatment has consistent effects across donors (donor mean), and (3) a model assuming all cytokine treatments have the same effect within a donor (cytokine mean). We found that when no samples of a test cytokine treatment were included in the training set, CellFlow was not able to predict a cytokine's donor-specific response more accurately than baseline methods (Figure 2B, Figure S1A). However, as soon as the cytokine response of one donor was included in the training data, the performance drastically improved (mean DEG R-squared of 0.76 vs. 0.60), outperforming all baselines (Figure 2B). This performance increase can also be illustrated in a UMAP embedding<sup>35</sup> (Figure 2C, Methods). While the donor-specific cell population generated without having seen IL-15 for any other donor resembled the control cells, having seen IL-15 for only one

other donor shifts predictions towards the true distribution, both on a population level and in terms of generated expression of the differentially expressed genes (Figure 2D,E). These results suggest that the ESM2 representation of the cytokine is not sufficient for CellFlow to extract the functional effect on gene expression level but it was able to calibrate to donor-specific effects as soon as the cytokine treatment effect has been measured for one donor. Moreover, measuring the cytokine treatment for multiple donors still increased the performance of CellFlow, but with decreasing marginal effects. Thus, given a sparse budget to explore the donor-specific responses of the remaining cytokine treatments, the most promising way to obtain reliable predictions is prioritizing the number of different cytokine treatments measured rather than the number of donors for a single treatment.

Next, we assessed the prediction of cytokine effects for new donors. To obtain a representation of each donor, we computed the donor-specific mean gene expression vector of the control population. This naive representation was sufficient for CellFlow to predict cytokine-specific donor responses more accurately than any baseline model (Figure 2F). In practice, it can be feasible to measure a limited number of cytokine treatments for a new donor. Therefore, we evaluated the performance of all models across different numbers of cytokine treatments measured for a new donor, which showed that adding more cytokine treatments further increased performance resulting in an almost 10-fold improvement over baselines (Figure 2G and Figure S3A). We further tested whether CellFlow could learn donor-specific responses. Analysis revealed positive correlations between predicted and true donor similarities based on test cytokine responses, although CellFlow tended to underestimate the magnitude of these differences between donors (Figure S2C-E).

We were able to observe a clear scaling relationship between CellFlow's performance and the number of seen conditions. This scaling law followed a linear fit in log-log space with an intercept of -0.22 and an R squared value of 0.999 for 16, 32, 64, and 80 cytokines included in the training data (Figure 2G, Methods). Scaling behavior differed across individual donors and cytokines (Figure 2H). For example, CellFlow's performance rapidly increased for donor 9. In contrast, the performance for donor 1 only increased when including a large number of treatments (Figure 2H). This may be explained by the high similarity of donor 9 to other donors, while donor 1 was highly dissimilar to all other donors (Figure 2I and Figure S2A). Analogously, we found a similar distinction between the scaling curves for IFN-omega and IL-32-beta (Figure 2I), suggesting the presence of a similar cytokine as IFN-omega and the lack of a cytokine similar to IL-32-beta. Indeed, IFN-omega had a highly similar effect as IFN-beta (cosine similarity=0.97), while the most similar cytokine to IL-

32-beta only had a cosine similarity of 0.55 (Figure 2J, Figure S2B, Methods). The dissimilarity led us to ask whether CellFlow was able to generate such a novel cell state. Indeed, we found IL-32-beta generations to be closer to the true set of cells than any population in the training set for 9 out of 12 donors (Figure S3B). An analogous analysis revealed that predicted populations are highly donor-specific (Figure S4B). To further investigate this scaling behavior, we assessed the relationship between predictive performance and the presence of certain cytokines in the training dataset. For this, we computed feature importance scores for each test cytokine with respect to all other cytokines using a regularized linear model (Methods). We observed that prediction accuracy for IFN-omega strongly depends on the presence of IFN-beta in the training data, which is consistent with the similar cellular responses elicited by these two cytokines. Similarly, the majority of identified train cytokines which positively influence CellFlow's performance could be explained by a high similarity to the test cytokine (Figure 2K). Together, this shows that CellFlow exhibits scaling behaviour with respect to the number of measured cytokines, which can be explained by similarity relationships between cytokines and donors.

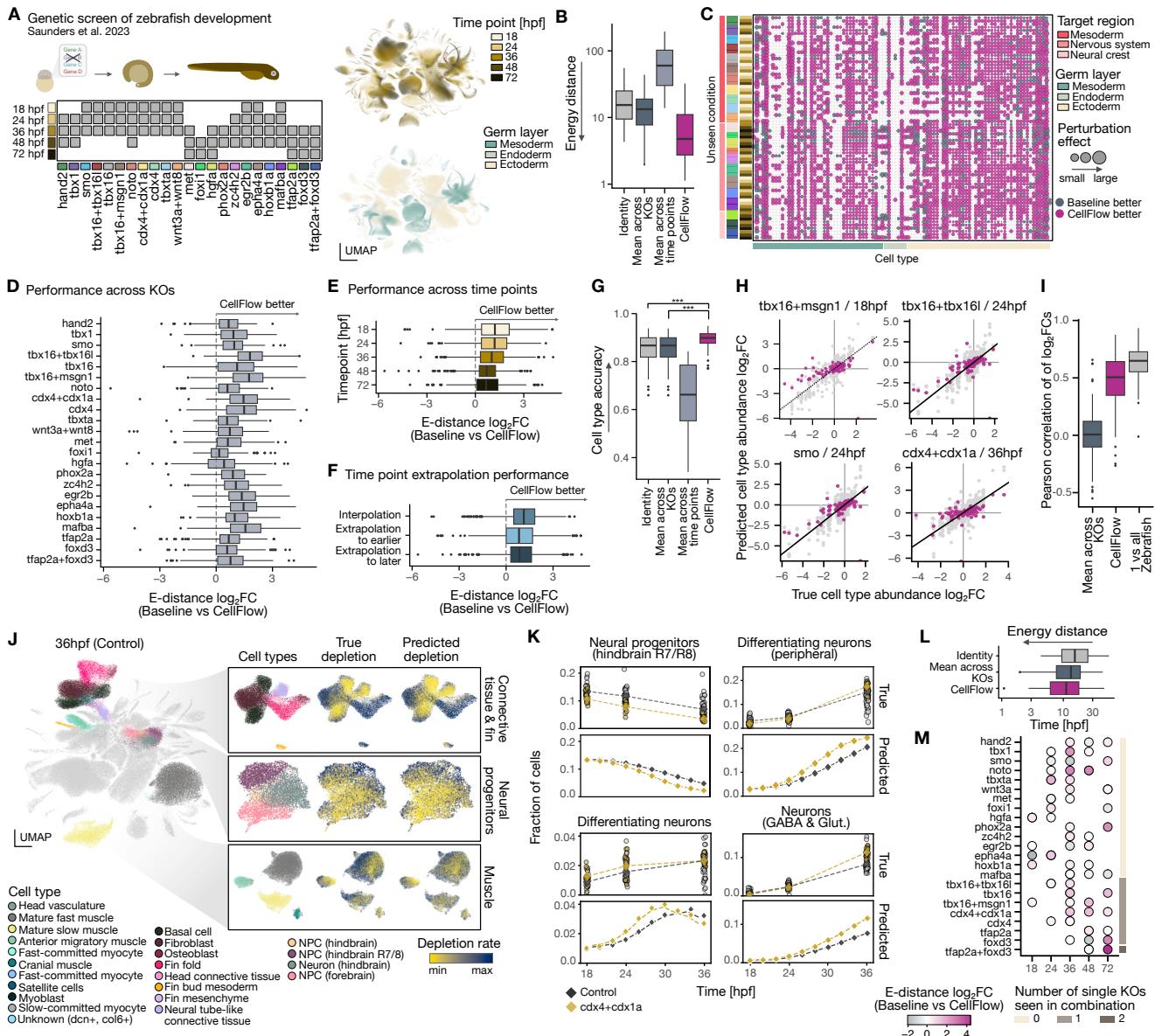
In summary, having trained 792 CellFlow models on up to ten million cells, resulting in more than 50 million generated cells (Methods), we found that CellFlow is stable to train and capable of accurate predictions exceeding baselines in almost all settings. We also observed a consistent and explainable training behaviour, satisfying scaling laws in the number of seen conditions.

**Modeling the perturbed development of zebrafish embryos** To understand CellFlow's ability to model perturbations of developing multicellular systems, we investigated the predictive performance on a phenotypic screen from an entire developing organism. We used ZSCAPE<sup>10</sup> (zebrafish single-cell atlas of perturbed embryos), which measured the effects of 23 gene knockouts induced through F0 genome editing measured on scRNA-seq profiles across five points during development (Figure 3A). Here, each knockout was only profiled at a subset of the five considered time points. Thus, we tested whether CellFlow could accurately predict the cellular phenotypes of perturbed developing embryos at unobserved developmental stages. We simulated this setting by predicting the cellular states of mutants, resulting in 71 trained models, each corresponding to one held-out condition defined by a combination of knockout and time point. We considered three baseline models: (1) An identity model assuming a gene knockout has no effect (identity); (2) the mean effect of all other perturbations observed at the same time point; and (3) the mean effect of the same gene knockout at the remaining observed time points (Methods). CellFlow outperformed all baselines across differ-

ent metrics (2x improvement in energy distance; Figure 3B, Figure S4A). To assess predictions on the cell type level, we transferred cell type annotations from measured cells to generated cells (Methods). We found that the performance was consistently above baseline across gene knockouts, developmental stage of the embryo, perturbation strength as well as time point interpolation and extrapolation scenarios (Figure 3C-F, Figure S4B). However, the extent of improvement decreased with more mature developmental stages of the embryo, possibly due to fewer phenotypes being observed at later time points and the increasing heterogeneity of the embryos. These evaluations show that CellFlow is able to learn meaningful gene knockout-specific developmental phenotypes.

One major phenotypic impact of the knockouts is relative changes in cell type proportions in the organism. Thus, we quantified CellFlow's ability to accurately predict changes in cell type proportions. We found that CellFlow accurately predicted cell type proportions based on transferred annotations (mean accuracy=0.89 vs. 0.85 of best baseline; Methods), outperforming all baselines (Figure 3G). To understand to what degree changes in cell type abundances were predicted, we compared true and predicted fold changes. We found predictions were generally correlated with true proportion changes (mean correlation=0.43) and accurately captured whether cell type abundance was increased, decreased or not affected (Figure 3H,I). However, strong effects were often underestimated (Figure 3H). This overly conservative behavior may in part be due to the variability of individual mutants (Figure 3H,I) and the high variability of the one-nearest neighbor classifier used for cell type annotation transfer (Figure S4B). To further assess CellFlow's predictions on the single-cell level, we computed depletion rates of single cells induced by the combined knockout of cdx4 and cdx1a, as previously investigated<sup>10</sup> (Methods). We used our model to extrapolate the effect of the genetic perturbation observed at 18hpf and 24hpf to time point 36hpf and found cell-specific true and predicted depletion rates to be highly correlated across affected organs of the zebrafish, even if the magnitude was often underestimated (Figure 3J, Figure S4C-E). For example, given observations in earlier time points, CellFlow was able to project the impact of the perturbation in hindbrain neural progenitor cells to 36h. This analysis demonstrates that CellFlow can accurately capture both organism-wide and tissue-specific cell type abundance changes, enabling the prediction of perturbed cell proportions for unobserved time points.

We next leveraged CellFlow to model the continuous, perturbed development of the central nervous system (Methods). We held out the cdx4/cdx1a mutant at time point 24, allowing us to compare the predicted cell type composition of the central nervous system with the true, unseen one. CellFlow's predictions of cell type fractions



**Figure 3. Organism-scale prediction of perturbed zebrafish development** **(A)** ZSCAPE7 (zebrafish single-cell atlas of perturbed embryos) captures the developing zebrafish at developmental stages ranging from 18 hpf (hours post fertilization) to 72 hpf with different genes knocked out. **(B)** Energy distance in latent space for CellFlow, as well as three baseline models for the task of predicting the cellular phenotype of a perturbed zebrafish at a certain time point: The identity assumes no perturbation effect, while the mean models assume either a constant effect across gene knockouts or a constant effect across time points. **(C)** Dotplot showing CellFlow's performance relative with the knockout mean baseline per cell type and time point. Cell types are categorized according to germ layer (bottom strip), and ordered by similarity (Methods), while gene knockouts are ordered by target region, gene knockout and developmental time (left strips). The size of the dot indicates the effect of the perturbation on the specific cell type, measured with the energy distance. **(D)** Comparison of CellFlow's performance (measured by energy distance per cell type) with the best baseline model aggregated across cell type and developmental time. **(E)** Performance gain by CellFlow with respect to the best baseline model aggregated across cell type and gene knockout. **(F)** Performance gain of CellFlow with respect to the best baseline model aggregated to the relative temporal position of the left perturbed developmental stage, categorized into interpolation (there are measurements for the same genetic knockout in both and earlier and a later time point), extrapolation to earlier time point (there is no measurement of the same genetic knockout in an earlier developmental stage), and extrapolation to a later time point (no measurement of the same genetic knockout in a later developmental stage). **(G)** Accuracy of predicted cell type proportions after perturbation of CellFlow and the three baseline models. \* $<10^{-2}$ , \*\* $<10^{-3}$ , \*\*\* $<10^{-4}$ , unpaired t-test. **(H)** True and predicted log fold changes of cell types for a certain perturbation and developmental time point of the zebrafish. Gray dots correspond to estimates of single perturbed zebrafish (as opposed to the union of cells across zebrafish of same perturbation and developmental stage, which is taken as the ground truth), providing a notion of noisiness of the data. **(I)** Pearson correlation between true and predicted log-fold changes across all 71 predicted perturbed states of the developing zebrafish, together with a reference obtained from log-fold changes of a single perturbed zebrafish with respect to the union of all perturbed zebrafish. **(J)** UMAP of control cells of a zebrafish at developmental stage 36hpf, with all cell types colored which depletion rates are calculated for (left, Methods). Insets visualize true and predicted depletion rates, which are computed from true cdx4/cdx1a-perturbed cells and predicted cdx4/cdx1a-perturbed cells at time point 36hpf. **(K)** Cell type proportions of selected cell types of the Central Nervous System at different developmental stages. The top boxes show the cell type fractions of single control zebrafish (gray dots) and single perturbed zebrafish (yellow dots), as well as the fraction of the union of all control zebrafish (dark gray diamond) and all perturbed zebrafish (yellow diamond). The dotted line depicts a linear interpolation between the time points. The bottom boxes show predicted cell type proportions for control and perturbed zebrafish, interpolated to densely sampled time points which have not been measured in the dataset. The model has been trained with the cdx4/cdx1a perturbed population only present at 18hpf and 36hpf, i.e. 24hpf was not included in the training data. **(L)** Performance of CellFlow and baseline models for the task of predicting perturbed zebrafish without having seen the same genetic perturbation at any time point. **(M)** Improvement of CellFlow with respect to the best baseline model, which models the perturbation as the mean displacement of all other perturbations at the same time point.

accurately matched the mean of the ground truth distribution for most cell types despite the high variance between individual zebrafish (Figure 3K, Figure S4F). This allowed us to further generate cells corresponding to a control or cdx4/cdx1a perturbed zebrafish for densely sampled time points, resulting in an interpolation of cell type abundance changes during development.

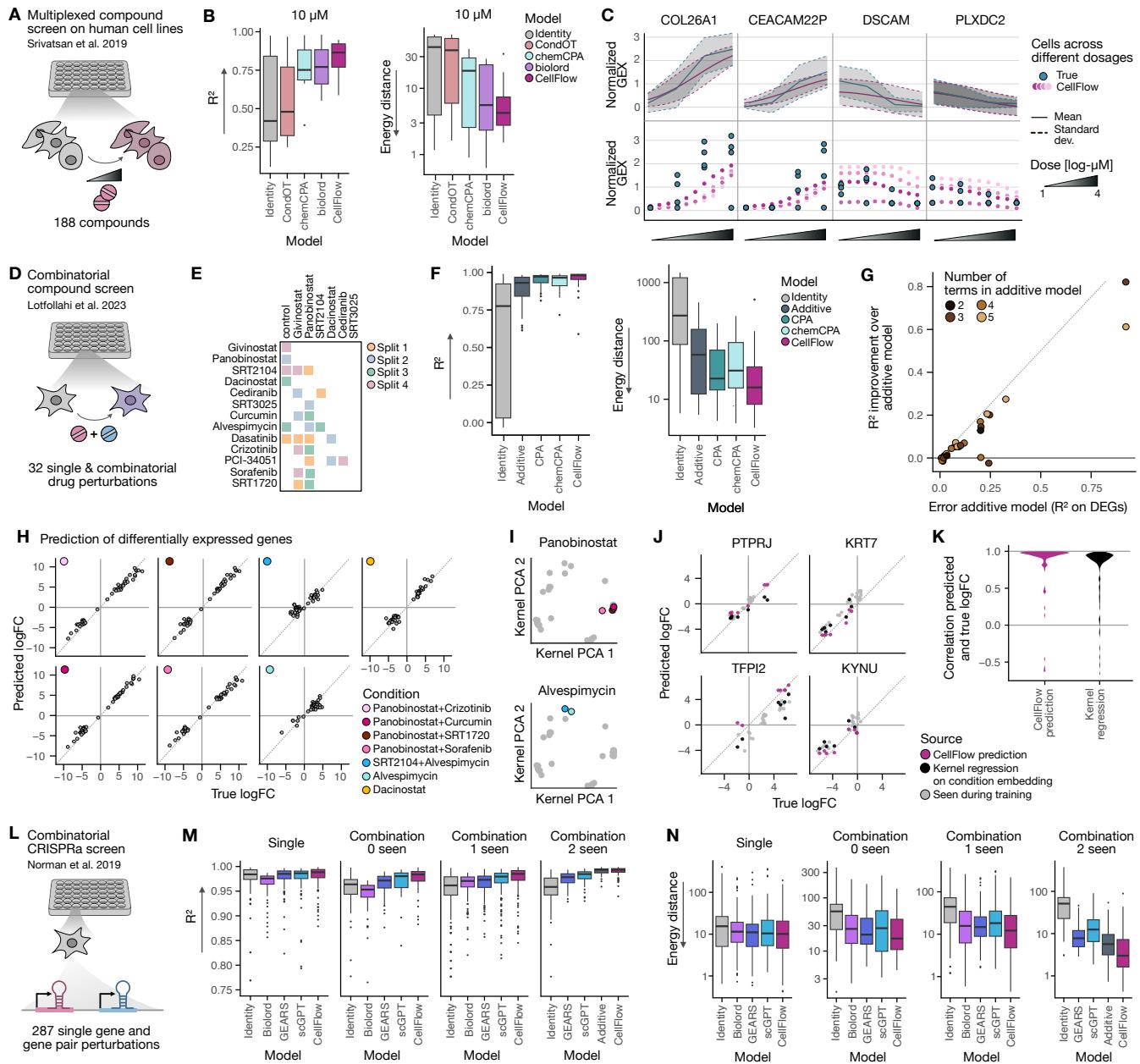
A more challenging task than predicting the effect of gene knockouts which have been seen at other developmental stages is modeling the evolution of an embryo with a completely unseen genetic perturbation. This requires CellFlow to fully rely on the discriminative power of the gene knockout embedding (here ESM2). While CellFlow performed above baseline on average, it did so less consistently than in the previous task (Figure 3L,M, Figure S4G). Analogously to the prediction of effects of new cytokines on PBMCs (Figure 2B), this suggests that the low number of observed gene knockouts in the training data makes it difficult to learn distinctive patterns in the data.

**Predicting perturbation effects for diverse and complex experimental designs** We next asked how CellFlow would perform on established perturbation prediction benchmarks across diverse experimental designs. We first evaluated its ability to predict single-cell expression responses to drug treatments in cancer cell lines<sup>4</sup> (sciPlex3 dataset; Figure 4A), comparing against established perturbation prediction methods chemCPA<sup>16</sup>, biolord<sup>19</sup>, and CondOT<sup>17</sup> (Figure 4B, Figure S6A). CellFlow performed competitively across all metrics, with particularly strong performance in predicting high-dose drug effects. For lower dosages, which typically induce more subtle expression changes, all methods showed comparable performance, with none significantly outperforming a simple identity baseline (Figure S6A). To illustrate CellFlow's ability to model continuous dose-response relationships, we examined predictions for MCF7 cells treated with Dacinostat, a drug which had not been seen during training. The model correctly captured monotonic changes of differentially expressed genes in response to increasing drug concentration while maintaining realistic cell states close to the phenotypic manifold (Figure 4C, Figure S5A,B). Moreover, CellFlow's inherent optimal transport modeling approach allowed to trace the response of single cells across different dosages, overcoming the limitation of the destructive nature of sequencing technologies (Figure 4C, Figure S5C,D). We further leveraged the sciPlex3 dataset to compare PCA and VAEs as possible methods for encoding and reconstruction of cells. We computed upper bounds of CellFlow's performance by encoding and subsequently decoding held-out conditions independently of CellFlow. While we found PCA encoders to be more powerful with respect to distributional metrics, the VAE was more powerful at fitting the mean gene expression (Figure S5E, Methods).

We next assessed CellFlow's ability to predict the effect of drug combinations using the combosciplex dataset, which contains scRNA-seq data for 31 single or combinatorial drug treatments on the A549 lung cancer cell line<sup>6</sup>. We split the data such that each of the individual drugs in the test data has been observed in a different combination in the training data (Figure 4E). CellFlow outperformed other established methods as well as an expression-space additive model (Figure 4F-G, Figure S6B, Methods). As recent studies have highlighted the difficulty of beating simple additive baselines<sup>6,36</sup>, we focused on the comparison between CellFlow and the additive baseline. As not all drugs appearing in combinations were measured individually, the additive model sometimes required more than two terms (Methods). We found that CellFlow largely outperformed the additive model, independent of the number of additive terms, and reduced the error of the additive model by 66% on average as measured by R squared values for DEGs (Figure 3G). CellFlow effectively recovered log-fold changes of perturbation-specific DEGs, showing particular strength in predicting large-magnitude changes while struggling more with subtle expression shifts (Figure 4H), which is consistent with our observations in the sciPlex3 dataset.

Following recent interests in learning representation of drugs with foundation models<sup>37</sup>, we analyzed whether CellFlow learnt a meaningful functional embedding of the combinations of drugs. Visually, we found similar combinations of drugs to be close in embedding space (Figure 4I). We quantified the biological meaningfulness of the learnt latent space by predicting the log-fold change of gene expression from the learnt embedding. In particular, we modelled the response of TFPI2 and PTPRJ, whose overexpression is associated with better clinical outcome in lung cancer patients<sup>38,39</sup>, and KRT7 and KYNU whose overexpression results in poor prognosis<sup>40,41</sup>. We found that a linear predictor trained on the learnt embedding space produced meaningful predictions, but did not match the accuracy of predictions directly computed from generated cell profiles (Figure 4J,K).

To investigate CellFlow's ability to predict the effect of genetic perturbations, we evaluated its performance on a dataset capturing the response to gene overexpression<sup>9</sup> (Figure 4L). Following a previously established evaluation holdout strategy for this dataset<sup>18</sup>, we assessed predictions on three categories of held-out conditions: previously unseen single genes, combinations where all genes were unseen during training, combinations containing one unseen gene, and novel combinations of genes that had been individually observed. We benchmarked CellFlow against GEARS, biolord, and scGPT<sup>42</sup>, as well as the identity baseline. CellFlow achieved the highest median R-squared and the lowest median energy distance across all evaluation splits (Figure 4M,N, Figure S6C, Methods). In line with drug



**Figure 4. CellFlow outperforms other methods on diverse perturbation prediction tasks.** (A) The sciPlex3 dataset contains measurements of three cancer cell lines treated with various drugs at different dosages. (B) Performance of CellFlow and other established methods for the task of predicting the effect of unseen drugs of the highest dosage on all three cell lines, measured by  $R^2$  in normalized gene expression space as well as with the energy distance in latent space (Methods). (C) True and predicted gene expression of differentially expressed genes of MCF7 cells treated with Dacinostat, a drug which has not been seen during training, across different dosages. CellFlow enables tracing the gene expression of single cells (colored by different shades of magenta) across different interpolated dosages (bottom). In contrast, gene expression of measured cells is only available for four different dosages and cells are unpaired across different dosages. (D) The combiosciplex dataset captures gene expression responses of the A549 cell line treated with different drugs. (E) Combinations of drugs present in the combiosciplex dataset colored by test split. (F)  $R^2$ -squared of normalized gene expression aggregated across all four splits for CellFlow and other established methods (Methods), as well as the energy distance in latent space (Methods). (G) Comparison of the performance of CellFlow and the additive model with respect to the  $R^2$ -squared of treatment-specific differentially expressed genes. (H) True and predicted log-fold changes of treatment-specific differentially expressed genes for all drug perturbations in split 3. (I) Learned embedding space by CellFlow with all conditions in gray which have been seen during training, and all test-conditions containing Panobinostat (top) or Alvespimycin (bottom) colored. (J) Predicted log-fold changes induced by treatments for lung cancer suppressing genes PTPRJ, KRT7, TFPI2, and KYNU. Predictions are either obtained from generated gene expression CellFlow (magenta) or from a linear model trained on CellFlow's learnt embedding space (dark gray). (K) Correlations between true and predicted log-fold changes across all conditions and all four cancer suppressors PTPRJ, KRT7, TFPI2, and KYNU for predictions obtained from generated samples, and predictions obtained from a linear model on CellFlow's learnt embedding space. (L) The genetic perturbation dataset captures cellular responses to genetic perturbations yielding overexpression of genes on the K562 cancer cell line. (M) The genetic perturbation dataset captures cellular responses to genetic perturbations yielding overexpression of genes on the K562 cancer cell line. (N) Comparison of performance of different methods based on  $R^2$  of normalized gene expression and energy distance in latent space (Methods).

perturbation modeling as outlined above, we defined an additive model for the group of conditions comprising unseen combinations of seen genes. This additive model performed comparably to CellFlow when evaluating mean gene expression (0.99 mean R-squared for both models), but in terms of energy distance, CellFlow improved upon the additive model by 2-fold (Figure S6C).

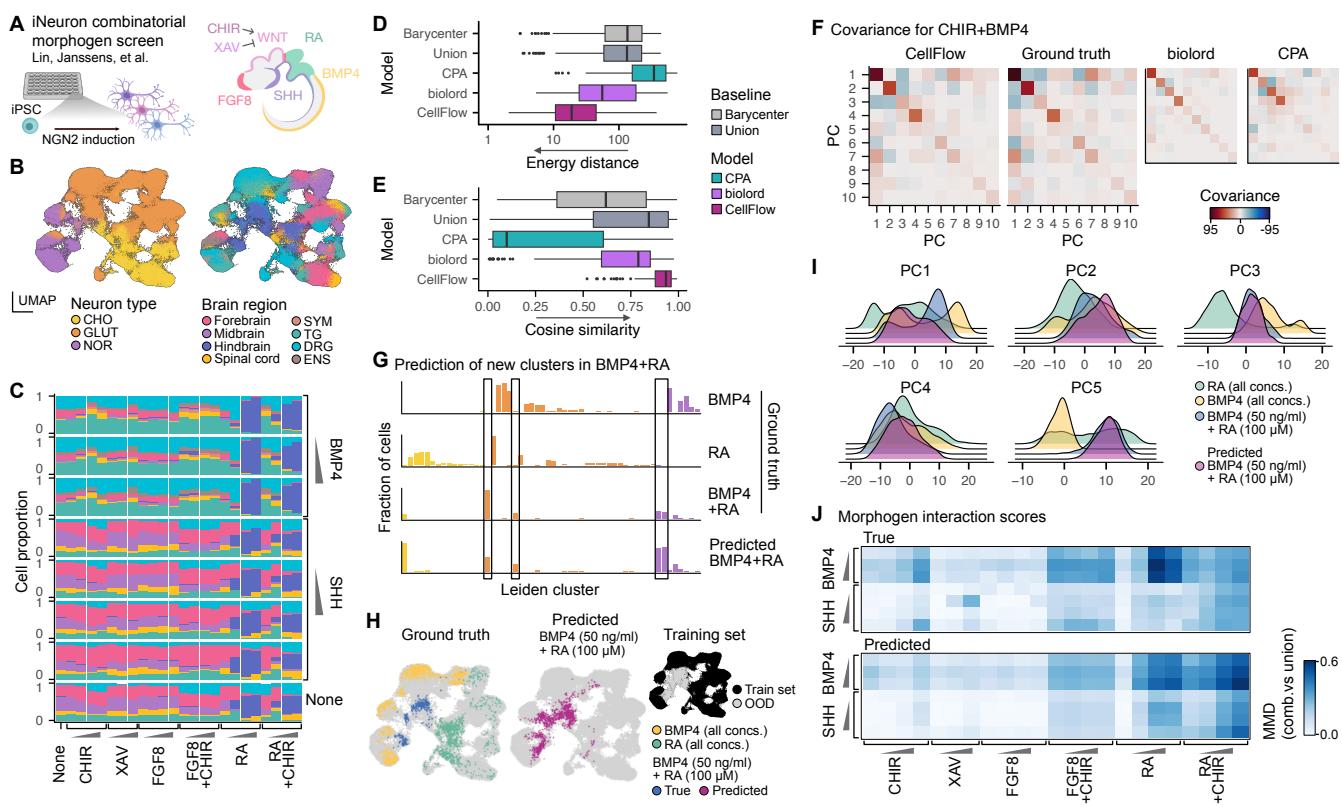
To further showcase CellFlow's ability to address complex experimental designs, we used a recent study of genetic perturbations on six different cancer cell lines stimulated with five different cytokines<sup>7</sup>. In this dataset, all cell lines were treated with each of the cytokines, but the each cytokine treatment was combined with a different set of gene knockdowns (Figure S7A). As with the PBMC and zebrafish examples, other methods were insufficient to capture this experimental design. We therefore compared CellFlow against two baselines: an identity model assuming no perturbation effect, and a model predicting the mean effect of all other perturbations in the same cell line/pathway combination (Methods). CellFlow is competitive with the additive model with respect to the performance on DEGs across perturbation strengths (Figure S7B, Methods). We then tested whether CellFlow could predict perturbation effects across cell lines (Figure S7C, Methods). Again, CellFlow clearly outperformed baseline models for strong effects, while for subtle expression changes, the identity model often yielded better predictions (Figure S7D). Based on previously described conserved and context-specific perturbation programs<sup>7</sup>, we evaluated CellFlow's predicted effect on curated sets of target genes on the interferon-gamma (IFNG) treated BXPC3 cell line, a condition manifesting strong and heterogeneous perturbation effects<sup>7</sup>. CellFlow's predictions were more similar to the ground truth than the responses of any other cell line or the mean thereof (Figure S7E). As expected, effects of cell line-specific perturbation responses were harder to predict than those that are conserved across cell lines (Figure S7E). Finally, we assessed whether CellFlow could predict gene knockdowns to a completely unseen cell line, a challenging task as the performance relies on CellFlow extracting patterns from only five cell line embeddings (Figure S7F). While predictions were better than the ones obtained from baseline models for strong perturbation effects, CellFlow was not able to outperform the identity model for smaller effects (Figure S7G).

The flexibility of CellFlow does not only apply to the choice of experimental designs, but also to the type of readout. As a proof of concept, we demonstrated its capability for predicting the effect of chemical perturbations on cancer cell lines measured by the 4i technology<sup>17,43</sup> (Figure S8). Altogether, we demonstrated CellFlow's performance on a large variety of perturbation experiments, allowing for more flexibility than previous methods which are mostly tailored towards spe-

cific experimental setups.

**Modeling neuron programming through combinatorial morphogen treatment** We next sought to use CellFlow to predict the outcome of cell fate programming experiments where conditions can generate unexpected and novel cell states. For this task, we leveraged a dataset based on a combinatorial morphogen screen on NGN2-induced human neurons (iNeurons)<sup>13</sup>. Such *in vitro* systems are challenging to model, as each condition can contain heterogeneous cell fates in response to complex interactions of intrinsic and extrinsic signalling modulators. In this screen, NGN2 expression was induced in iPSCs, followed by treatment with morphogen pathway modulators during neuronal differentiation (Figure 5A). The treatment conditions comprised combinations of modulators of anterior-posterior (AP) patterning (RA, CHIR99021, XAV-939, FGF8) with modulators of dorso-ventral (DV) patterning (BMP4, SHH), each applied in multiple concentrations. scRNA-seq readouts of patterned NGN2-iNs showed that treatments resulted in diverse distributions of neurons within each condition, including three distinct neuron classes and various brain region identities such as forebrain, midbrain, hindbrain and spinal cord as well as peripheral sympathetic and sensory neurons (Figure 5B,C).

To assess the ability of our model to predict neuron distributions induced by combinations of pathway modulators, we evaluated prediction performance on held-out combinatorial treatments. Specifically, in each training run, we withheld all experimental conditions containing certain combinations of AP and DV modulators (Methods), then assessed the ability to predict cellular states under these unseen treatment combinations. We compared the performance of CellFlow with CPA and biolord as well as two baseline models: the union of cellular distributions of individual treatments and their distributional mean (Wasserstein barycenter<sup>44</sup>). Intuitively, these baselines represent cases where combinatorial treatments generate a combination of states observed under individual treatments and an intermediate cellular state, respectively. We found that CellFlow strongly outperformed both existing methods and baselines in terms of energy distance and other distributional metrics ( $\geq 2.5x$  mean improvement in energy distance; Figure 5D and Figure S9A-C). We additionally sought to understand whether less accurate predictions default to predicting familiar training examples or generate states outside of the training distribution. While predictions were generally close to training conditions, they did not move closer as prediction error increased (Figure S9D, left). Instead, higher prediction errors correlated with greater divergence from the average training distribution (Figure S9D, right), indicating that CellFlow's errors tend to occur when attempting to model new cellular states rather than from incorrectly defaulting to known states.



**Figure 5. Neuron fate prediction from combinatorial morphogen treatment to enable condition prioritization.** (A) Schematic of the experimental setup for the iNeuron combinatorial morphogen screen. iNeurons were treated with morphogen pathway modulators during neural differentiation from pluripotency. (B) UMAP embedding of scRNA-seq data from all screening conditions colored by neuron type (left) and brain region identity (right). (C) Bar plot showing the proportion of brain regions in individual screen conditions. (D) Boxplot showing energy distance between predicted and true cell distributions for baselines, established perturbation prediction methods and CellFlow. (E) Boxplot showing cosine similarity of true and predicted leiden clusters. (F) Heatmap showing covariance between the first ten principal components for the true CHIR+BMP4 condition, as well as CellFlow, biolord and CPA predictions. (G) Bar plot showing the leiden cluster proportions for the true BMP4, RA and BMP4+RA condition as well as the BMP4+RA condition predicted by CellFlow. (H) UMAP embedding showing the cells belonging to the true BMP4, RA and BMP4+RA conditions (left), BMP4+RA cells predicted by CellFlow projected onto the UMAP embedding (middle) and UMAP embedding showing the cells included in the training set (right). (I) Density plots showing marginal distributions over the first five principal components for the true BMP4, RA and BMP4+RA condition as well as the BMP4+RA condition predicted by CellFlow. (J) Heatmap showing the maximum mean discrepancy between cell distributions arising from combinatorial morphogen treatment and the union of distributions from individual morphogen treatments.

To better evaluate the accuracy of predicted cell identity composition, we compared true and predicted Leiden cluster abundances, obtained by nearest neighbor-based label transfer<sup>45</sup> (Methods). This revealed that predictions achieved consistently high cosine similarity scores (mean=0.91) with true cluster distributions, exceeding the performance of other methods and baselines ( $\geq 70\%$  mean improvement; Figure 5D,E). CellFlow also accurately captured the covariance structure of cell state distributions, demonstrating its ability to model realistic cell populations (Figure 5F and Figure S9E). These results demonstrate that previous methods for perturbation prediction struggle to capture heterogeneous phenotype distributions that arise during fate programming experiments and often do not surpass simple baselines, highlighting the power of CellFlow to model cellular state transitions.

The combination of BMP4 and RA resulted in new cell states that were not observed in either individual treatment<sup>13</sup>. To assess whether such interactions between pathway modulators could be captured by CellFlow we compared the distribution of true and predicted cells

across Leiden clusters, which revealed that CellFlow accurately predicted emergent populations of cells (Figure 5G). We further visualized these predictions on a UMAP projection (Methods), illustrating that CellFlow can predict new cell populations that were not observed during training (Figure 5H). This is also reflected in the marginal distributions along individual principal components, which show that CellFlow correctly learns an intermediate distribution which deviates from distributional averaging or addition of individual treatment effects (Figure 5I and Figure S9F).

Given CellFlow's ability to accurately predict combinatorial treatment outcomes, we explored its utility in identifying morphogen interactions resulting in emergent new phenotypes, thereby enabling the prioritization of experimental follow-ups. To measure the degree to which morphogen combinations induce new cell states, we devised an interaction score based on the maximum mean discrepancy between combinatorial treatment effects and the union of distributions from individual treatments (Methods). When comparing these interaction scores between predicted and experimen-

tal data, we found that CellFlow successfully captured the overall pattern of interactions, particularly identifying combinations with high interaction scores such as BMP4+RA and BMP4+FGF8+CHIR. However, we observed discrepancies in certain combinations, such as SHH+XAV and SHH+CHIR, where interaction strength was highly concentration-dependent (Figure 5J). Despite these limitations, these results indicate CellFlow's potential as a tool for guiding experimental design by predicting which combinations of pathway modulators are likely to yield novel and unexpected cellular states.

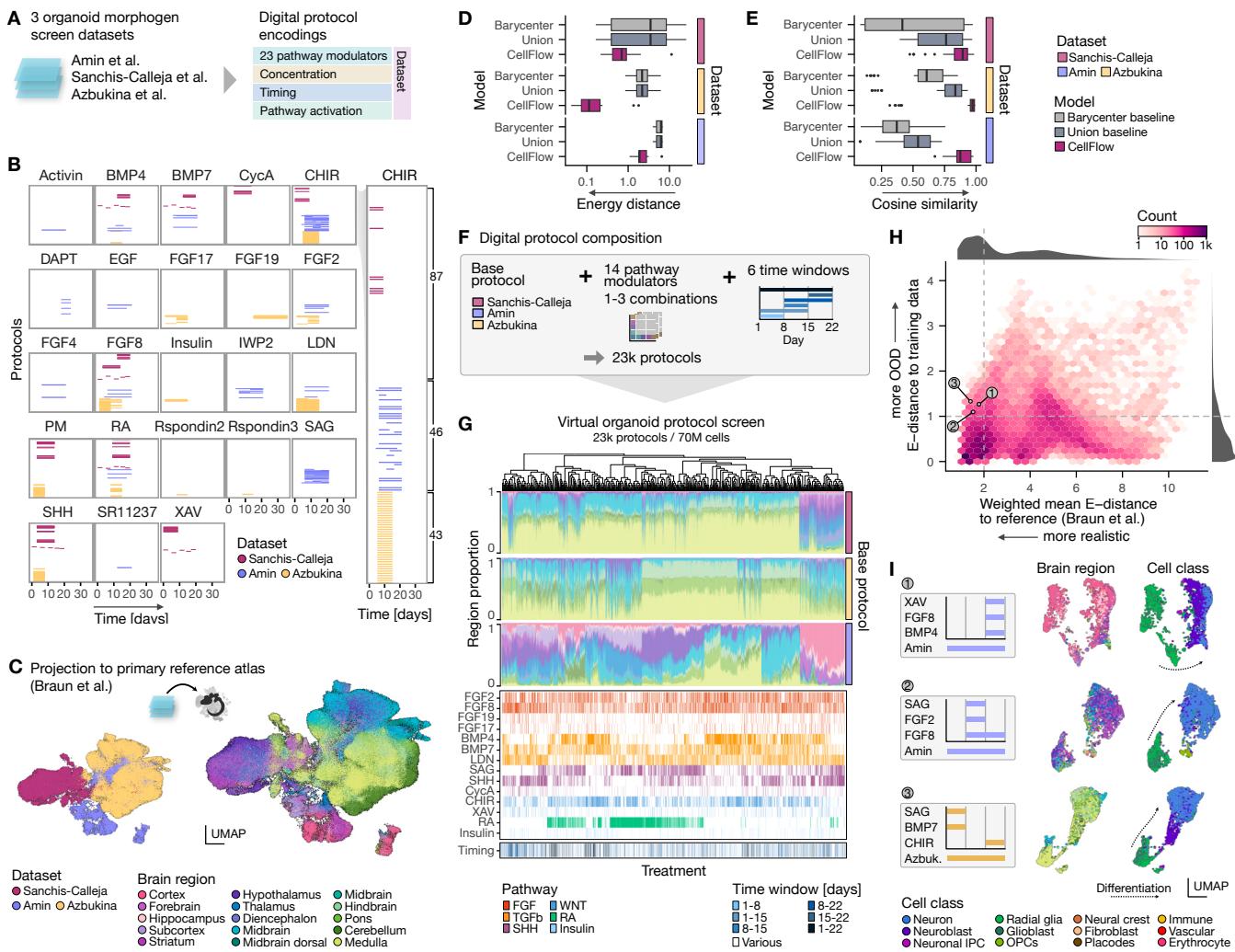
**A virtual organoid protocol screen** *In vitro* organoid models of brain development are complex tissues that provide inroads to study the mechanisms underlying human patterning. We applied CellFlow to model the cell type composition of human brain organoids induced by multi-step organoid protocols<sup>46</sup>. To cover diverse organoid protocols, we assembled three scRNA-seq datasets of human brain organoid morphogen screens<sup>47-49</sup> for training and evaluation. In these screens, organoids were grown from pluripotency and exposed to different morphogen pathway modulators at various timings during development. Across the three datasets, 23 pathway modulators were applied in various combinations at defined time windows between day 1 and day 36 of development, comprising a total of 176 conditions. To train the model jointly on all datasets, we harmonized protocol configurations into a common protocol encoding, capturing the pathway modulators, their respective concentrations, timings and information about the modulated pathway (Figure 6A,B, Methods). This representation also incorporated a one-hot encoding of the dataset label, to account for other dataset-specific experimental parameters (base protocol), such as media composition or usage of matrigel. To integrate scRNA-seq data of the three datasets and allow for comparison of organoid cell states with primary counterparts, we projected scRNA-seq data from all datasets to a single-cell transcriptomic atlas of the developing human brain<sup>50</sup> using scANVI<sup>45,51</sup> (Methods). This joint latent space allowed us to annotate brain region and cell types across datasets through nearest neighbor-based label transfer, revealing the broad distribution of brain regions generated by the diverse brain organoid culture conditions (Figure 6C and Figure S10A).

To assess the ability of CellFlow to learn from these systematic screens to predict cell type distributions for new protocol configurations, we first evaluated its performance on held out combinatorial treatments. We withheld all conditions containing certain combinations of pathway modulators during training, then assessed predictions for these unseen treatment combinations (Methods). Analogously to the evaluation on iNeurons, we compared the performance of CellFlow against distributional baselines (Wasserstein barycenter, union) based on energy distance and cosine similarity of pre-

dicted leiden clusters (Figure 6D,E). CellFlow achieved low energy distances and high cosine similarity scores (median=0.94), largely outperforming baselines (5.7x mean improvement in energy distance over union, Figure S10B,C). Corroborating our results from iNeurons, we found that the model was able to capture the interaction between BMP7 and CHIR, where combinatorial treatment leads to the emergence of new cell populations (Figure S10D,E).

To further assess to what degree information learned in one dataset could help make predictions on other datasets, we next focused on pathway modulators that were used in multiple datasets and systematically withheld all conditions including them from one dataset at a time during training. We compared predictions against a baseline distribution comprising all training conditions from the respective dataset, to evaluate the prediction of effects beyond the protocol differences between datasets. CellFlow largely outperformed this baseline in two datasets (1.7x mean improvement in energy distance) but showed limited predictive power in the Azbukina et al. dataset (Figure S10F,G). We hypothesize that this stems from the exclusive use of higher-order modulator combinations in this dataset (Figure S10H,I), making it challenging to extrapolate from the simpler treatment conditions present in the other datasets. Generally, we conclude that CellFlow enables prediction of cellular compositions arising from new combinatorial treatments in organoid protocols, and shows promising cross-dataset transfer ability for lower-order morphogen combinations.

To systematically explore the space of untested organoid protocols, we used CellFlow to perform a proof-of-principle virtual organoid protocol screen. Based on our previous observations of reduced prediction accuracy for higher-order combinations, we restricted the analysis to single, double, and triple modulator treatments. We focused on 14 pathway modulators that were well-represented across datasets ( $\geq 4$  conditions) at their highest observed concentration. We then composed protocol conditions by applying each modulator in one of 6 time windows during development (Figure 6F). We trained CellFlow on all available conditions from the three datasets to obtain predictions for all composed protocols. Altogether, our systematic *in silico* screen generated cellular compositions for more than 23,000 protocols, comprising over 70 million cells. We annotated all generated cells through transfer of region, cell type and cluster labels from the primary reference atlas. This revealed that generated organoids were predominantly composed of progenitor and neuron populations with regional identities spanning all major brain regions (Figure 6G and Figure S10J). Overall, the base protocol had a strong impact on brain region and cell type composition. We note that generated organoids corresponding to one base protocol (Sanchis-Calleja et al.) also contained neural crest populations, consistent



**Figure 6. Integrated learning from systematic morphogen experiments enables virtual neural organoid protocol screening.** (A) Schematic of integrating three organoid morphogen screen datasets to obtain consistent digital protocol encodings. (B) UMAP embedding obtained by projecting scRNA-seq data of all three datasets onto a primary reference atlas<sup>45</sup> using scANVI<sup>46</sup>. (C) UMAP embedding obtained by projecting scRNA-seq data of all three datasets onto a primary reference atlas<sup>45</sup> using scANVI<sup>46</sup>. (D and E) Box plots showing the prediction performance on held out morphogen combinations in terms of energy distance of distributions (D) and cosine similarity of predicted leiden clusters (E). (F) Schematic of digital protocol composition by combining a base protocol with combinatorial pathway modulator treatments during defined time windows, resulting in 23k virtual protocols. (G) Predictions of the virtual protocol screen. Bar plot showing brain region composition of all predictions (top) as well as morphogen usage (middle) and timings (bottom). (H) Density scatter plot showing weighted mean energy distance to the reference atlas (measure of realism) versus energy distance to the closest training condition (measure of novelty) for all predictions. (I) UMAP embeddings showing examples for three predictions resulting in different brain region-specific progenitor and neuron populations. Colors indicate brain region (left) and cell class (right).

with observations from the original experimental study where specific morphogen combinations induced neural crest<sup>47</sup>. These differences between base protocols are expected, as protocol-specific experimental parameters can strongly influence organoid development<sup>45</sup>. To assess the impact of combinatorial modulator effects across base protocols, we analyzed their influence on brain region composition using a regularized linear model (Methods). We found that both individual modulators and specific combinations showed consistent relationships with specific brain regions across protocols (Figure S10K). For example, retinoic acid (RA) treatment was positively associated with pons and medulla cell identities, while the combination of BMP4 and RA promoted cerebellar fates. This combination in isolation was not included in any of the training datasets, but

the predicted effect is consistent with known dorsalizing effects of BMP4 and role of RA in hindbrain patterning<sup>52–55</sup>.

To further characterize all predicted protocols, we computed two complementary metrics. To assess the biological plausibility of predicted cell populations, we quantified the energy distance between each predicted cell cluster with its counterpart in the reference atlas, then computed a weighted mean across clusters for each condition. We quantified the novelty of predictions with respect to the training data with the energy distance to the closest condition in the training datasets (Figure 6H). Most protocols generated cell populations that closely resembled both primary reference data and existing training conditions, indicating that CellFlow's predictions largely remained within the bounds of plausi-

ble cellular states. A subset of protocols exhibited both high realism (by reference similarity) and novelty compared to the training data. These predictions spanned a range of predicted cellular compositions and naturally clustered into three major categories based on regional identity: mid/hindbrain, ventral forebrain, and dorsal forebrain (Figure S10L). We examined an example from each of these categories as UMAP embeddings, showing trajectories from progenitor to neuron states with distinct brain regional identities that are specific to the protocol composition (Figure 6I). Together, these results demonstrate that CellFlow can generate cell type distributions that resemble organoid and primary cell states and make predictions for untested combinations of morphogen treatments and timing regimens. As a result, CellFlow can help explore and iterate organoid protocol design.

## Discussion

Computational approaches that can learn from high-content phenotypic screens have the potential to accelerate scientific discovery. Our work demonstrates that CellFlow can effectively learn from existing phenotypic screens to predict cellular responses to uncharacterized conditions. CellFlow builds upon flow matching and learns a functional embedding space from experimental variables to encode sophisticated experimental designs. We showed that CellFlow accurately predicts complex phenotypic distributions across diverse biological contexts. We observed a clear scaling law with respect to the number of training conditions, demonstrated by the relationship between prediction accuracy and the number of observed cytokine treatments in the PBMC analysis. This scalable performance indicates that as more perturbation data becomes available, such models will become increasingly accurate and generalizable<sup>5</sup>.

A key strength of CellFlow is its ability to model heterogeneous cellular distributions that arise from perturbations in developmentally dynamic systems. Other perturbation modeling approaches often struggle to capture the heterogeneous outcomes of fate programming experiments, where interventions can generate diverse mixtures of cell states rather than uniform shifts. Our results demonstrate CellFlow's capacity to model developmental effects of genetic perturbations across entire zebrafish embryos and predict the emergence of new cell populations in neuron fate engineering as well as organoid cell type composition from various protocol configurations. These applications highlight the importance of modeling distributional perturbation effects, especially in all contexts where cellular heterogeneity is integral to the studied biological system.

Our analyses highlight that out-of-distribution prediction is still a very challenging task. In particular, the effectiveness of predicting responses to entirely unseen perturbations strongly depends on the quality of con-

dition embeddings and their ability to capture functional properties. While we leveraged existing embeddings like ESM2 for proteins and molecular fingerprints for small molecules, these representations may not optimally encode the functional roles relevant to cellular responses. This limitation was evident in our results for completely unseen cytokines or gene knockouts. Finding more powerful representations of perturbed molecular entities that better capture their functional roles could significantly improve extrapolation to unseen contexts. In this context, CellFlow's trained condition embedding space offers a promising avenue to learn mappings between structure and function.

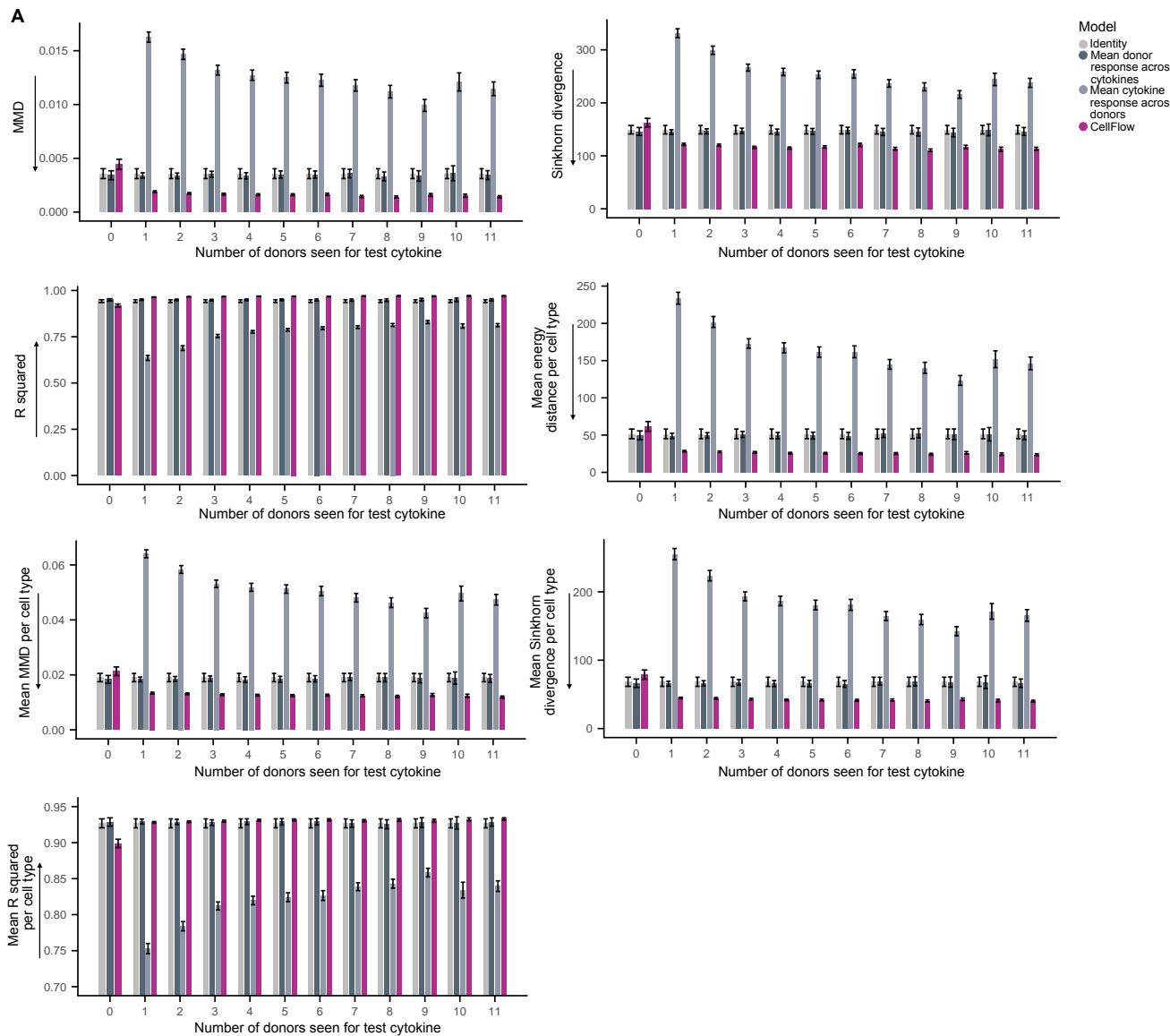
The modular architecture of CellFlow provides flexibility for extension and improvement across multiple aspects. For cell representation, we explored both PCA and VAE embeddings and found that they present a trade-off between reconstruction capability and flexibility. Future work could explore alternative cell embedding approaches, such as those derived from single-cell foundation models<sup>42,56</sup> or disentangled representations<sup>57</sup>. Moreover, modifications to the flow matching module provide promising avenues for improving CellFlow, for example by directly generating count data<sup>58</sup>. Stochastic extensions<sup>59,60</sup> could enable modeling of the inherent randomness in cellular responses, allowing for better uncertainty estimation on single-cell level, while stochastic parameterization of the encoder module would allow for uncertainty estimation on a distributional level. Furthermore, fine tuning approaches, which are an established part of the training pipeline in computer vision<sup>61</sup>, could enhance the ability to capture subtle transcriptomic differences. All of these technical improvements could be readily integrated into CellFlow's modular framework, enhancing its capabilities without requiring major modifications of the overall architecture.

Altogether, CellFlow presents a significant advance in computational modeling of perturbed single-cell phenotypes across diverse biological contexts. By addressing the limitations of existing approaches and providing a modular architecture for ongoing improvement, CellFlow has the potential to accelerate scientific discovery through systematic learning from phenotypic screens. One particularly exciting perspective is the integration of computational modeling with experimental design. Screens can be tailored to exploit the strengths of CellFlow in predicting combinatorial perturbations through sparse experimental designs, and predictions can then point to informative follow-up experiments. This lab-in-the-loop approach could be particularly valuable for exploring complex biological systems, where the space of possible experiments is large but resources are limited. As bigger and more diverse perturbation data sets become available, we anticipate that CellFlow will be a valuable tool to explore the vast space of cellular phenotypes.

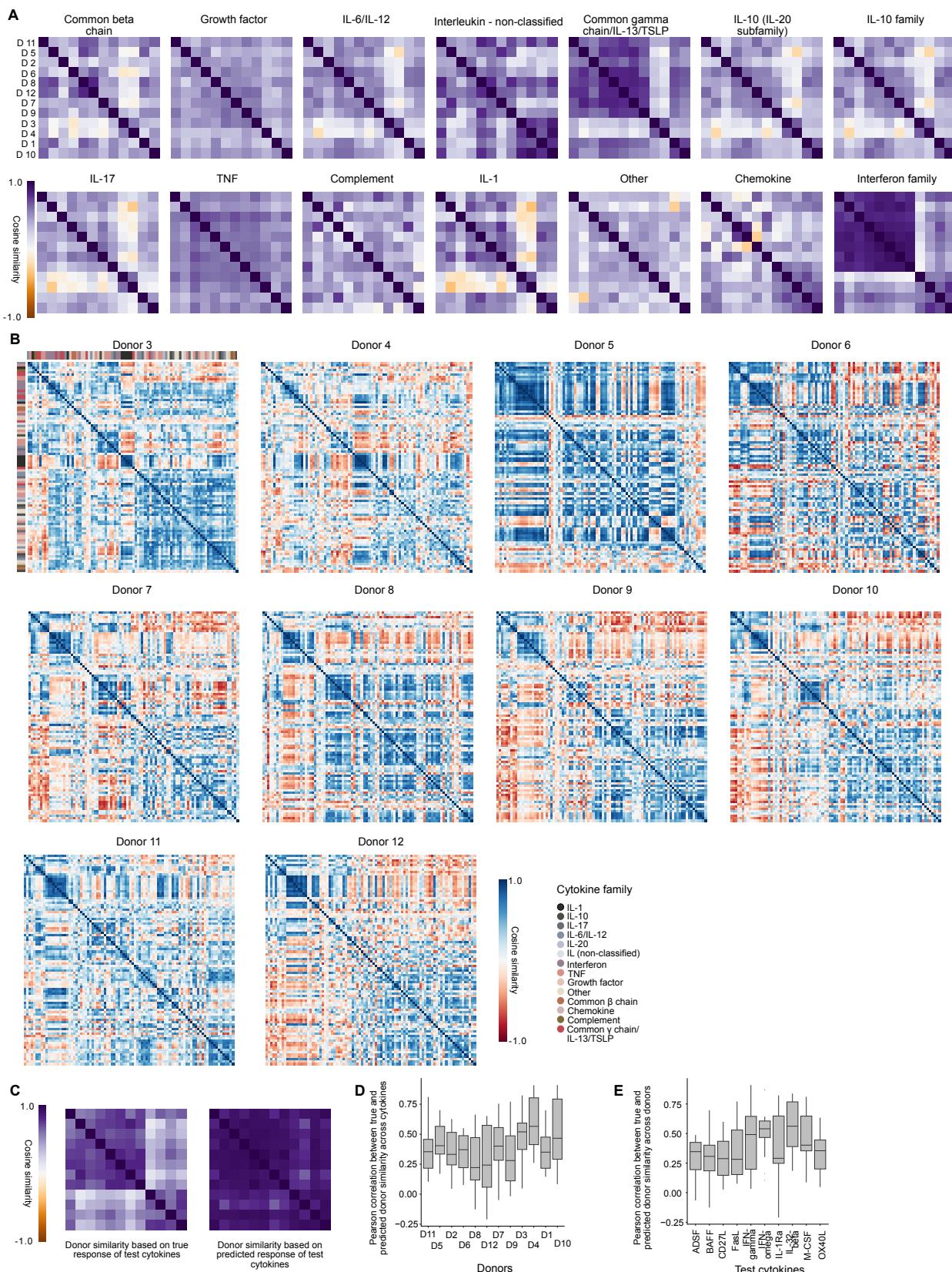




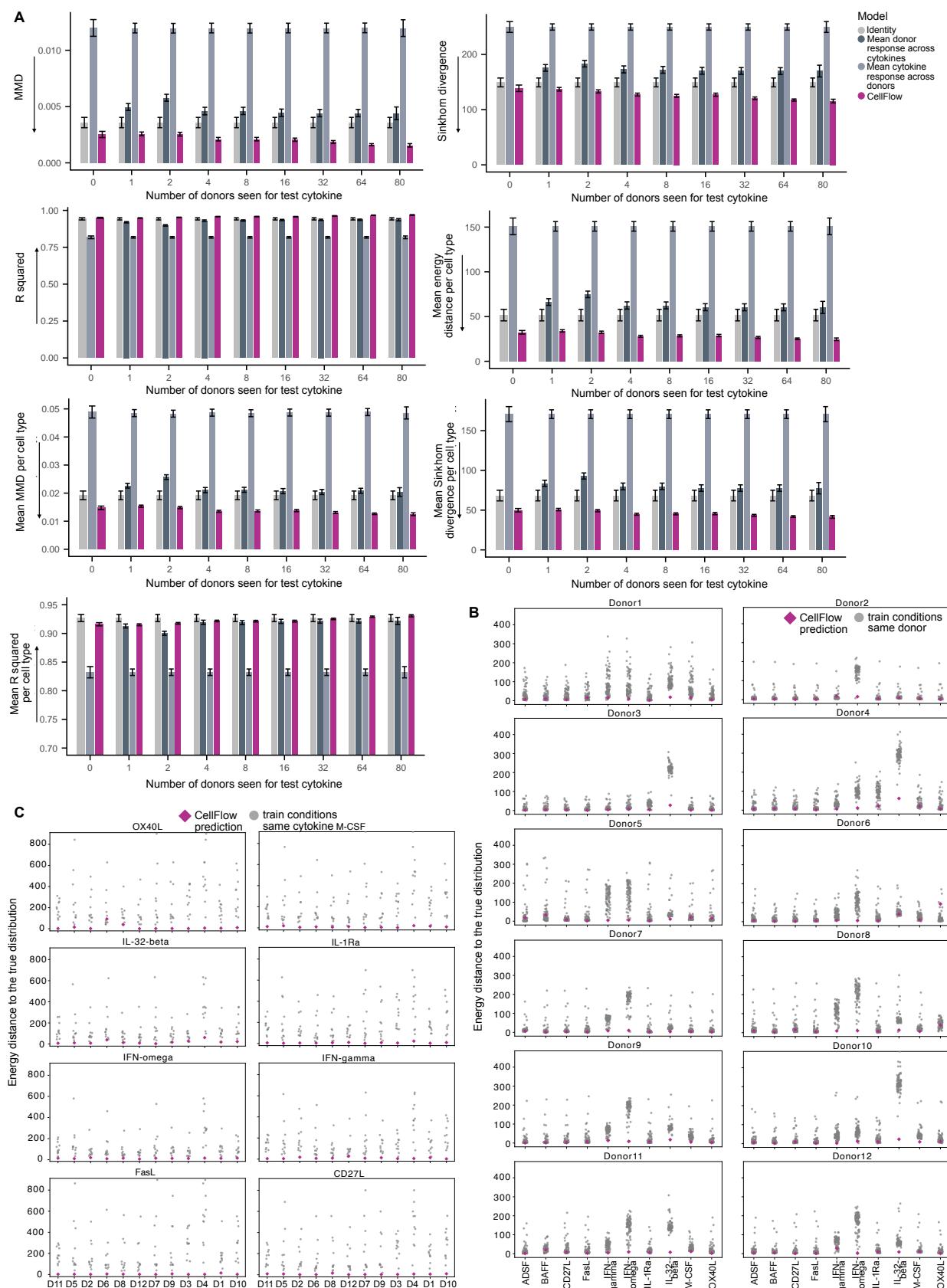
## Supplementary Figures



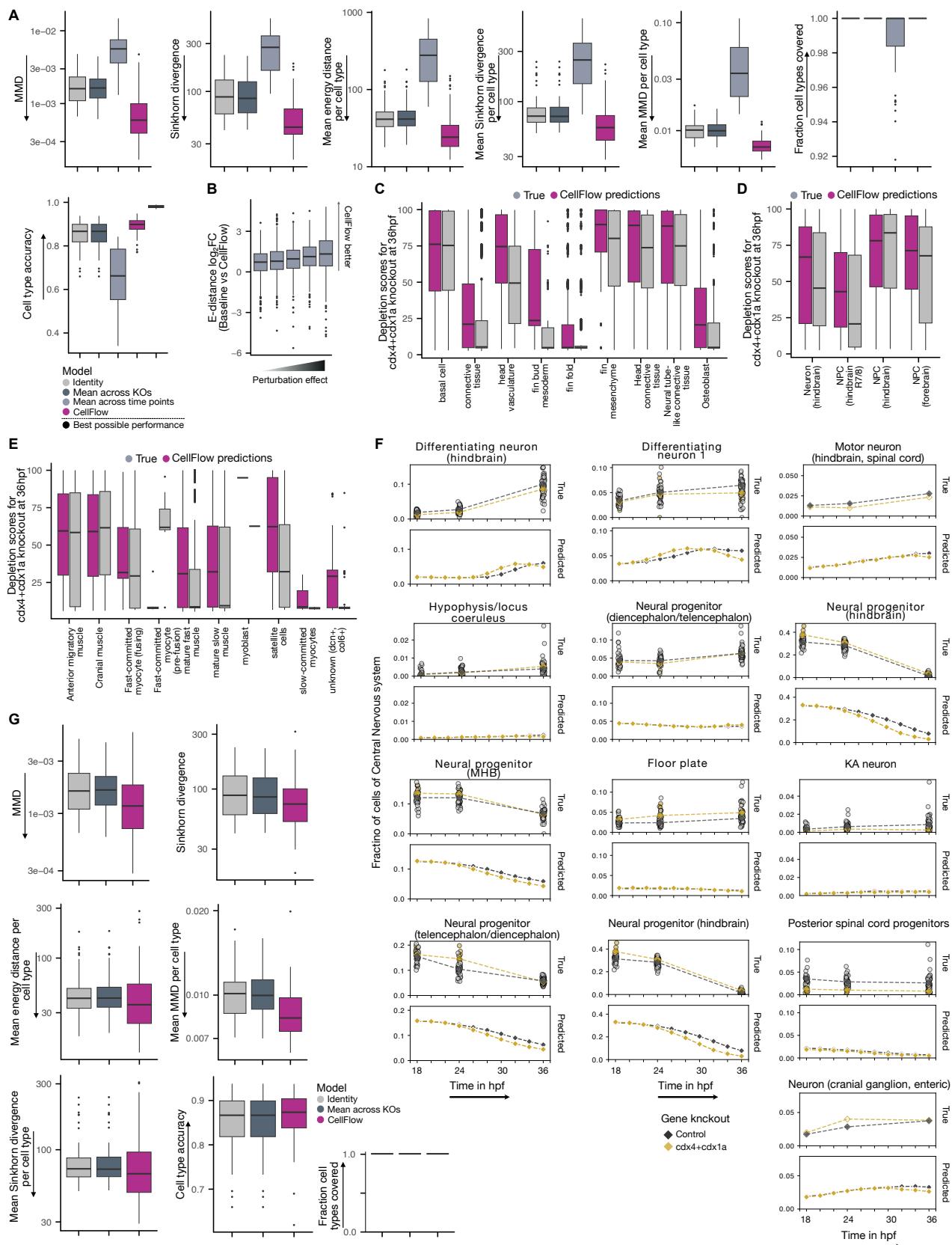
**Figure S1. CellFlow is superior to baseline models with respect to different metrics for predicting the effect of a new cytokine. (A)** Performance metrics MMD (Maximum mean discrepancy) and Sinkhorn divergence in latent space, and R squared of true and predicted mean normalized gene expression for CellFlow and baseline methods reported as mean and standard error across different test cytokines and different sets of donors of size  $k = 0, \dots, 11$  which the test cytokine has been seen for. Analogously, performance metrics are reported per cell type, where cell type labels are transferred using one-nearest neighbor (Methods).



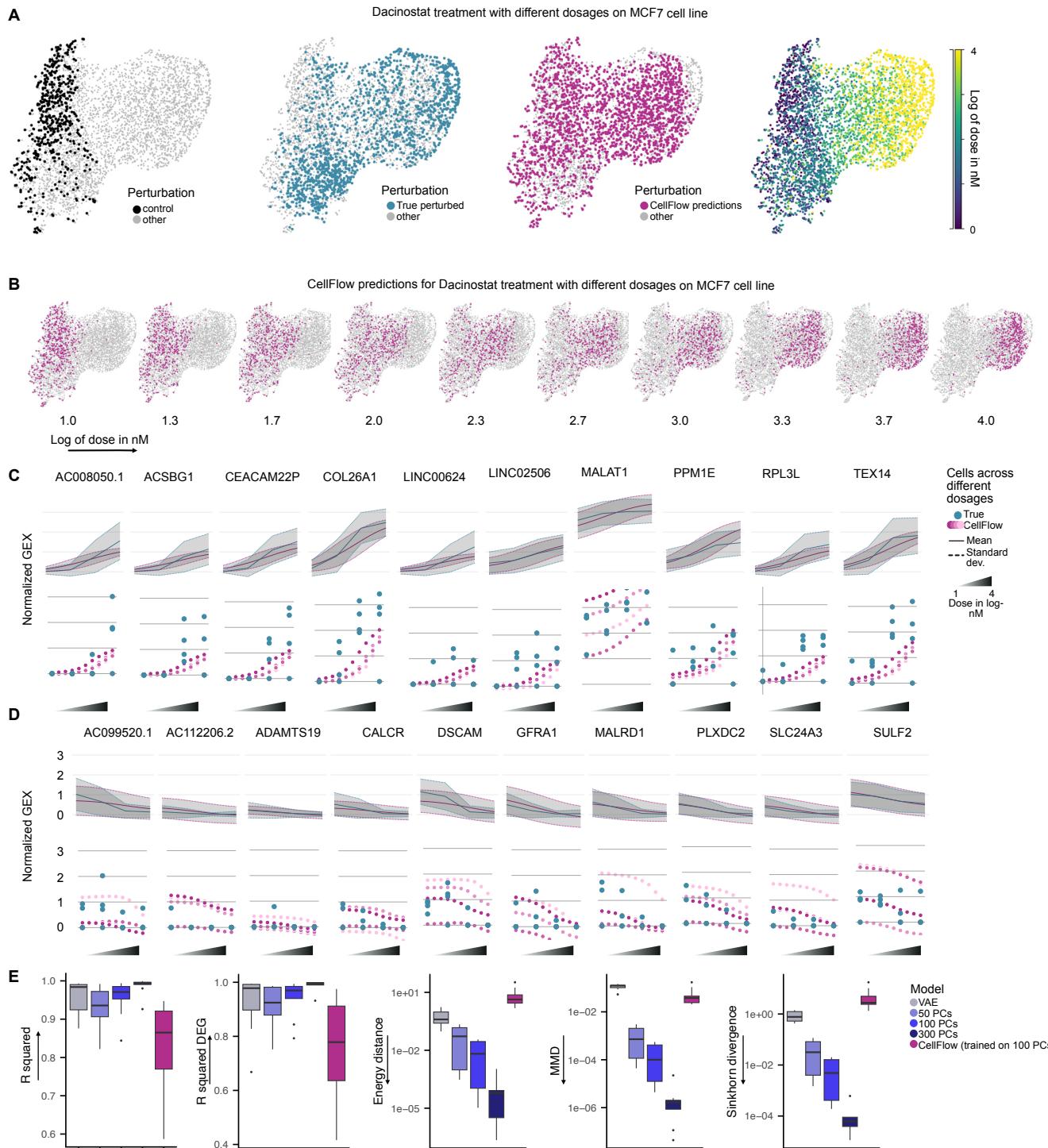
**Figure S2. Donor and cytokine similarities of cytokine-treated PBMCs** **(A)** Donor similarities computed from treatment responses of cytokines belonging to one cytokine family (Methods). **(B)** Cytokine similarities computed from treatment responses per donor. **(C)** Donor similarities computed from true cytokine responses based on the ten test cytokines (left), and donor similarities computed from predicted cytokine responses based on the test test cytokines. Predictions were computed from the models including 80 training cytokines for the test donor. **(D)** Pearson correlation between true and predicted donor similarities separated by donor. One data point in a box plot corresponds to one Pearson correlation coefficient between true and predicted donor similarity based on the response of one test cytokine. **(E)** Analogous to (D), but with one data point of the box plot denoting the Pearson correlation coefficient between the true and the predicted donor similarities for one test cytokine.



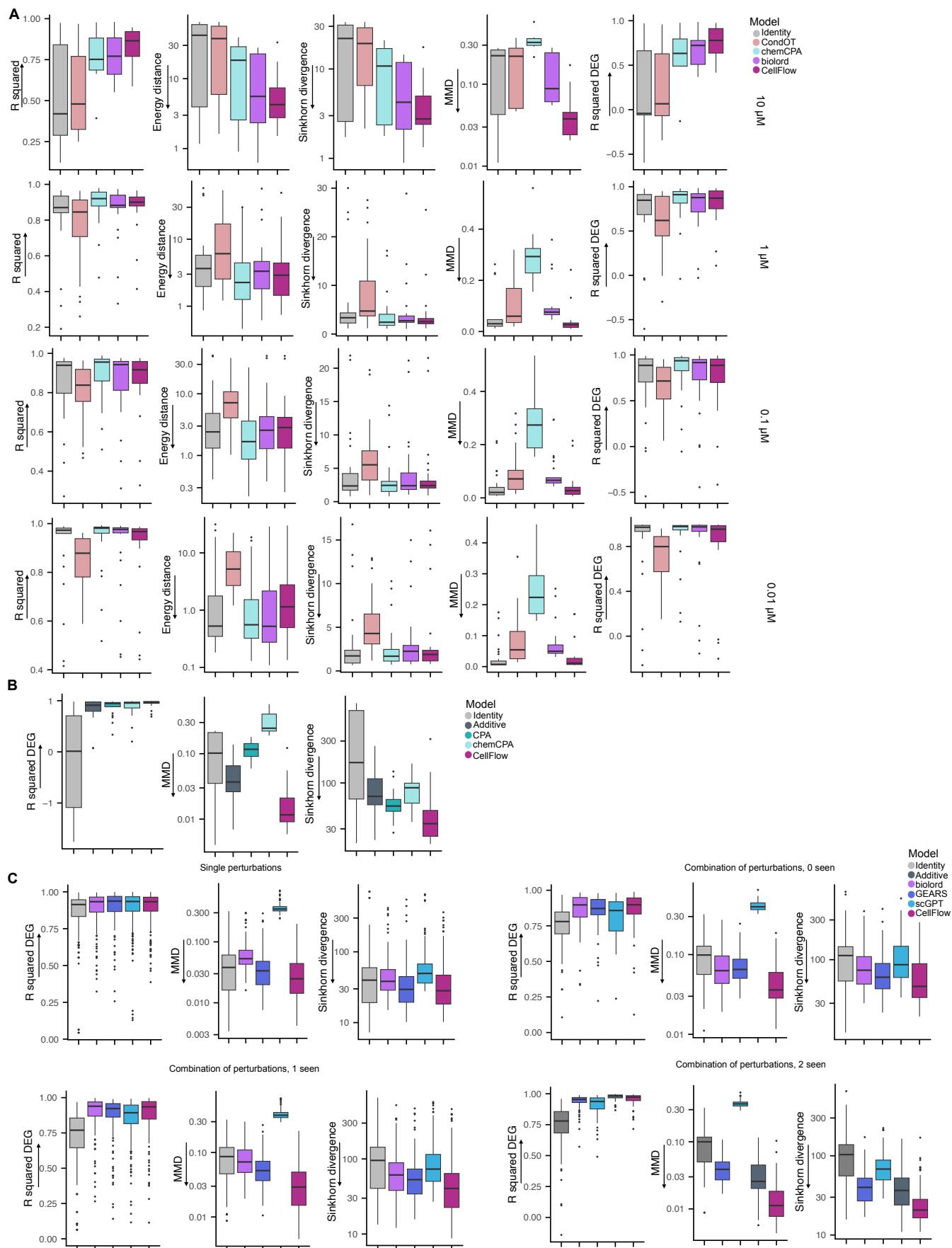
**Figure S3. CellFlow out-performs baseline models for predicting the response of a new donor across metrics** **(A)** Performance metrics MMD (Maximum mean discrepancy) and Sinkhorn divergence in latent space, and R squared of true and predicted mean normalized gene expression for CellFlow and baseline methods reported as mean and standard error across different test donors and different sets of cytokines of size  $k$  which have been seen for the test donor. Analogously, performance metrics are reported per cell type, where cell type labels are transferred using one-nearest neighbor (Methods). **(B)** Comparisons of the distance of CellFlow's prediction to the ground truth cell distribution (magenta diamond) with distances of training populations of the same donor to the ground truth. CellFlow's predictions are based on model trainings including 80 cytokines for the test donor. **(C)** Comparisons of the distance of CellFlow's prediction to the ground truth cell distribution with distances of training populations of the same cytokine to the ground truth.



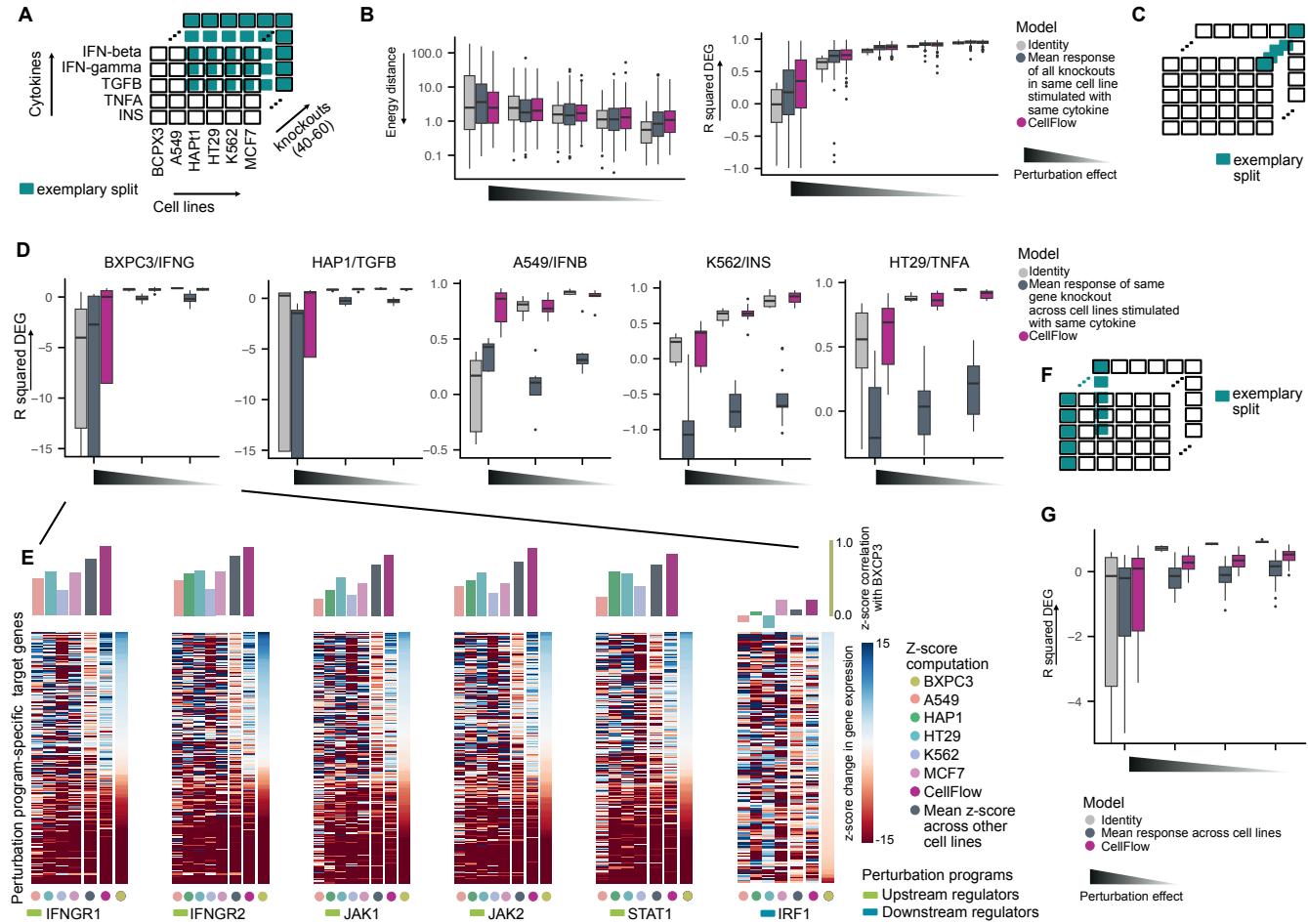
**Figure S4. Developmental modeling of perturbed zebrafish embryos.** (A) Additional metrics for the task of predicting the cellular phenotype of a mutant at an unseen combination of gene knockout and time point (Methods). (B) Comparison of CellFlow with the best base line method stratified with respect to perturbation strength (Methods). (C) True and predicted depletion scores aggregated to cell type level for the fin. (D) Analogous to (c) for neural progenitors. (E) Analogous to (C) for cell types of the muscle and connective tissue. (F) True and predicted cell type proportions of cdx4+cdx1a mutants at different developmental stages. The training data contains observations of cdx4+cdx1a mutants at time points 18hpf and 36hpf, but not at 24hpf. Observations for other time points are not given in the dataset, and hence only CellFlow's predictions can be reported. (G) Additional metrics for the task of predicting the cellular phenotype of a completely unseen mutation across all time points (Methods).



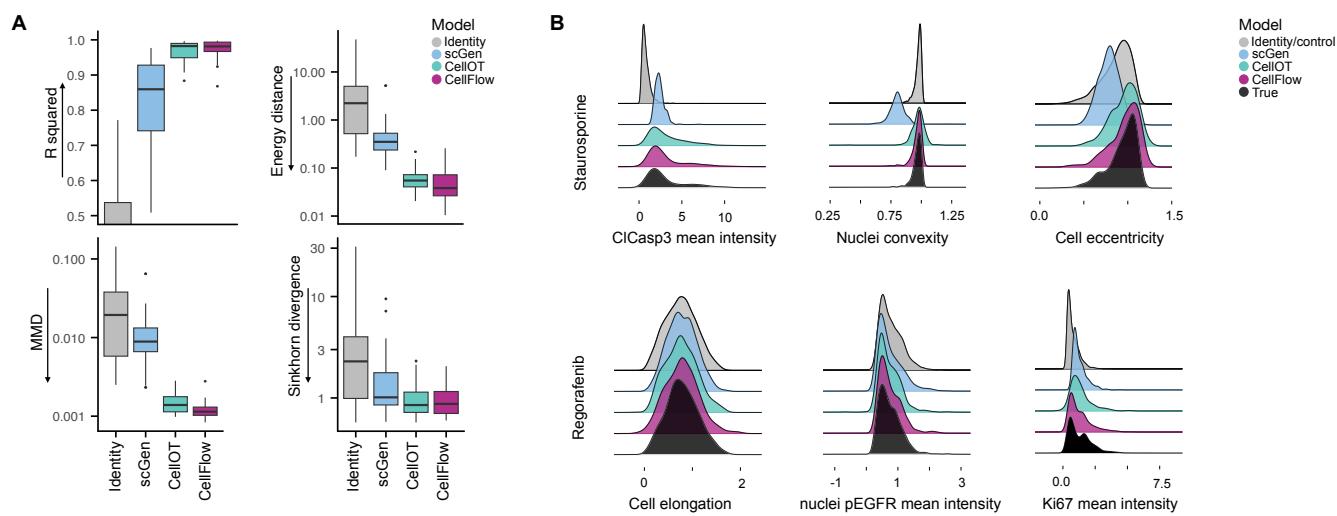
**Figure S5. Interpolation of effects across drug dosages.** (A) UMAP of control cells of MCF7, along with true and predicted Dacinostat-treated cells for dosages 10, 100, 1000, and 10000nM, colored by whether cells are true control, true perturbed, or predicted perturbed cells (from left to right), as well as altogether colored by dose. (B) Generated cells for Dacinostat treatment at interpolated dosages. (C) Normalized expression of top upregulated genes for Dacinostat treatment of MCF7 with 10000nM. (D) Normalized gene expression of top downregulated genes for Dacinostat treatment of MCF7 with 10000nM. (E) Upper bound of CellFlow's performance by encoding-decoding method (Methods) computed by encoding and decoding the true perturbed population, compared with actually achieved performance of CellFlow (trained on 100 PCs).



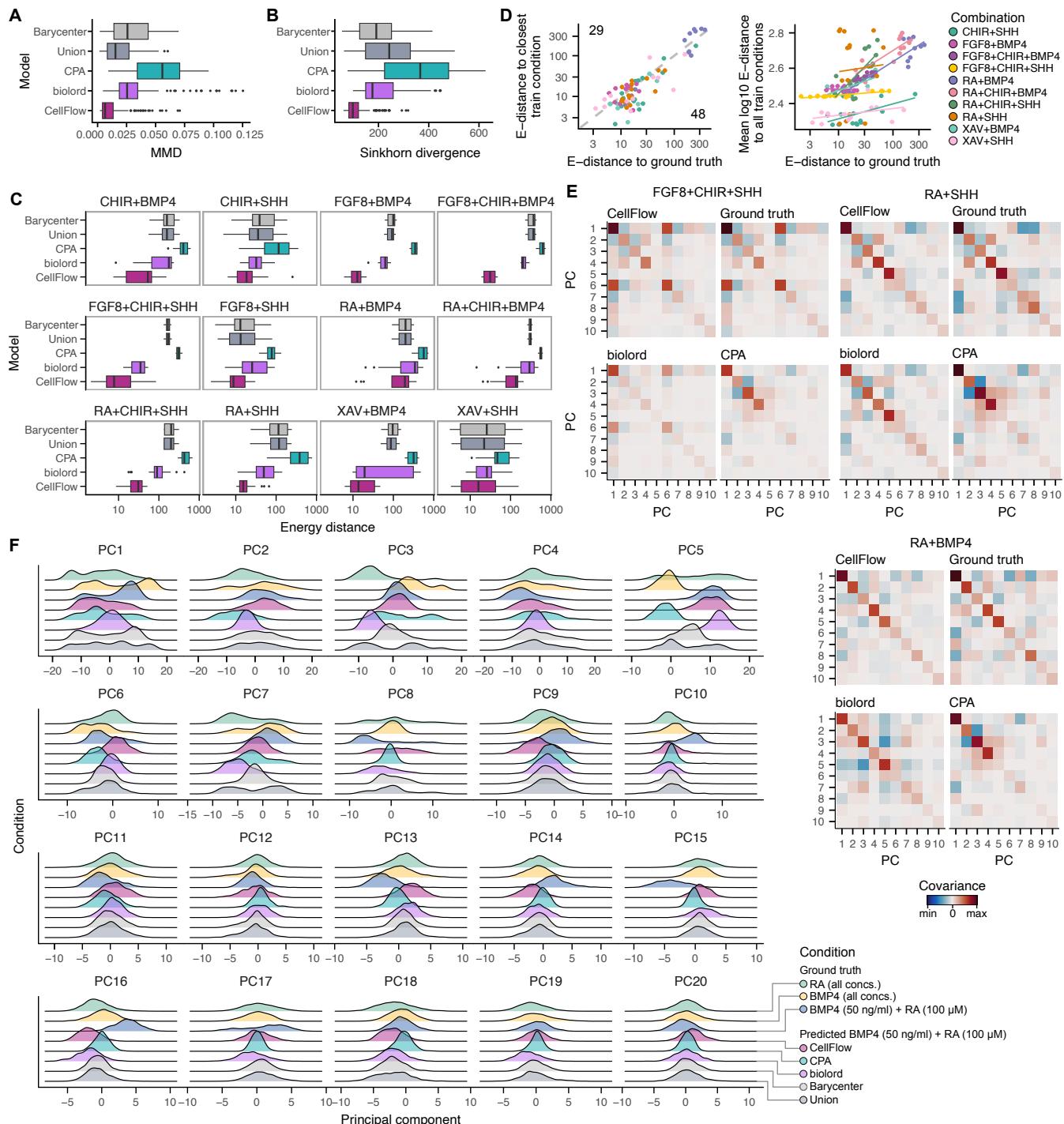
**Figure S6. Performance metrics on various perturbation effect prediction tasks.** Additional metrics for the task of **(A)** predicting the effect of unseen drugs on the sciPlex3 dataset for different dosages **(B)** predicting the effect of combinations of drugs on the combosciplex dataset and **(C)** predicting the effect of genetic perturbations.



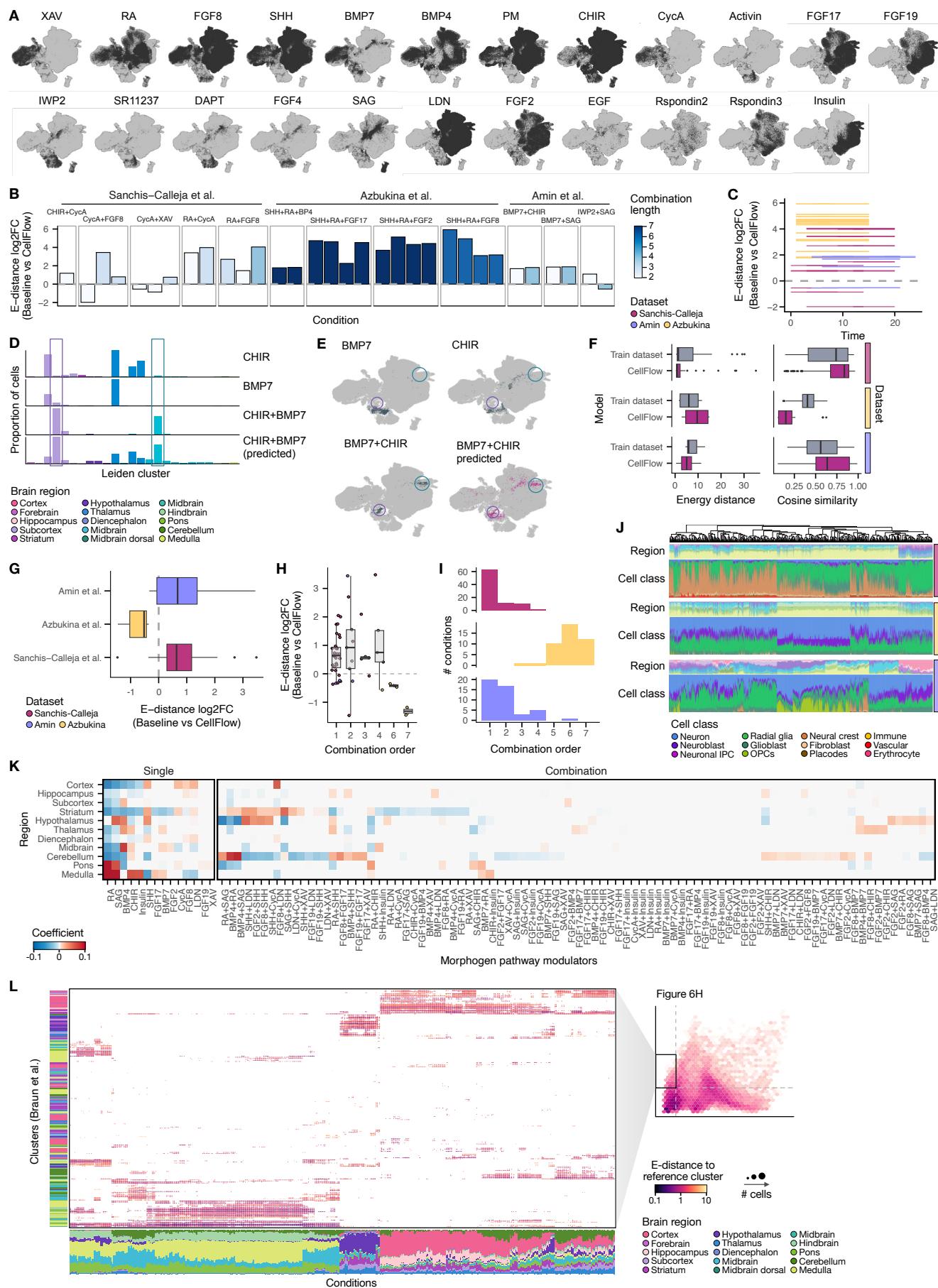
**Figure S7. Gene knockout effect prediction across cytokine-stimulated cell lines** **(A)** The dataset contains measurements of gene knockouts which are specific to the cytokine treatment of cell lines<sup>7</sup>. **(B)** Quantification of the performance of CellFlow for predicting the effect of unseen gene knockouts in cytokine-stimulated cell lines (Methods). Reported are the energy distance in latent space and the  $R^2$  of the condition-specific differentially expressed genes (DEGs) (Methods). Baseline models are the identity, which predicts the perturbation to have no effect, and the mean model, which predicts the perturbation to have the mean effect of the remaining perturbations in the same cytokine-treated cell line. Perturbations are binned into five equally sized groups based on the strength of the perturbation (measured by  $R^2$  of DEGs). **(C)** Sketch of the experimental setup of predicting gene knockouts in a new combination of cytokines and cell line, while the cell line and the cytokine have been seen in other combinations in the training data (Methods). **(D)**  $R^2$  squared of normalized gene expression of condition-specific DEGs for each combination of cell line and cytokine which has been held out from the training set. Gene knockouts are evenly split into three groups according to the perturbation strength (measured by  $R^2$  of DEGs). **(E)** Z-scores of up- and downregulation of selected genes specific to perturbation programs<sup>7</sup>. **(F)** Sketch of the experimental setup of holding out a complete cell line from the training data. **(G)**  $R^2$  squared on DEGs for the task of predicting the effect of all gene knockouts on the BXPC3 cell line for all cytokine treatments for CellFlow and baseline models. For all plots,  $*p < 10^{-2}$ ,  $**p < 10^{-3}$ ,  $***p < 10^{-4}$  (unpaired t-test).



**Figure S8. CellFlow predicts the effect of drug perturbations on 4i data. (A)** Performance metrics for predicting the effect of seen drugs on unseen cells of CellFlow, scGen<sup>15</sup>, CellOT<sup>17</sup>, CellFlow, and the identity model, which assumes the drugs to have no effect. The 4i dataset contains readouts for 35 drug treatments on two melanoma cell lines<sup>17</sup> (Methods). **(B)** True and predicted marginal distributions of a selected set of molecular and morphological features for two different drug treatments.



**Figure S9. Performance of neuron fate prediction from combinatorial morphogen treatment. (A and B)** Boxplot showing maximum mean discrepancy (MMD) (**A**) Sinkhorn divergence (**B**) and between predicted and true cell distributions for baselines, other established methods and CellFlow. (**C**) Boxplots showing Energy distance between predicted and true cell distributions for individual morphogen combinations. (**D**) Scatterplots showing the relationship between model performance and energy distance to the closest training condition (left) and mean  $\log_{10}$ (energy distance) across all training conditions (right). Lines represent the linear fits per combination. (**E**) Heatmap showing covariance between the first ten principal components of different conditions for the ground truth, as well as CellFlow, biolord and CPA predictions. (**F**) Density plots showing marginal distributions over all principal components for the true BMP4, RA and BMP4+RA condition predicted by CellFlow, CPA, biolord models and the barycenter and union baselines.



**Figure S10. Prediction of organoid cell type composition through digital protocol encodings.** **(A)** UMAP embeddings of three combined morphogen screen datasets colored by usage of pathway modulators in screen conditions. All screen conditions using the respective morphogen are colors in dark grey. **(B)** Bar plot showing performance of organoid composition prediction from combinatorial morphogen treatment for all hold-out splits and conditions. Energy distance log2 fold change was computed between CellFlow predictions and the union baseline. **(C)** Treatment time windows and performance for all evaluation conditions. **(D)** Bar plot showing the leiden cluster proportions for the true BMP7, CHIR and BMP7+CHIR condition as well as the BMP7+CHIR condition predicted by CellFlow. **(E)** UMAP embedding showing the cells belonging to the true BMP7, CHIR and BMP7+CHIR conditions (left), BMP4+RA cells predicted by CellFlow projected onto the UMAP embedding (middle). **(F-H)** Performance of predicting the effects of pathway modulators that were held out in one dataset at a time. The baseline comprises all training conditions from the respective dataset to evaluate prediction performance beyond the dataset label. **(F)** Box plots showing the performance in terms of energy distance and cosine similarity of predicted leiden clusters. **(G and H)** Box plot showing log2 fold change of energy distances between CellFlow predictions and baseline split by dataset **(G)** and combination length **(H)**. **(I)** Histogram showing the distribution of combination lengths (how many morphogens were combined) across all conditions in the three datasets. **(J)** Bar pot showing brain region and cell class composition of all predictions from the virtual protocol screen. **(K)** Heatmap showing the coefficients of a L1-regularized linear model to measure the association between pathway modulator treatment and brain region composition based on the virtual protocol screen. **(L)** Dot plot showing the energy distance between predicted cell clusters and corresponding primary reference clusters<sup>50</sup> for selected protocols from the virtual protocol screen (energy distance to training data > 1, mean energy distance to reference < 2).

# CellFlow methods

## Contents

<b>1 CellFlow</b>	<b>2</b>
1.1 One formulation to unify the setup of a perturbation experiment . . . . .	2
1.2 Learning a single representation from all perturbation factors . . . . .	2
1.3 Optimal Transport for pairing cells . . . . .	3
1.4 Flow Matching for generating perturbed cell populations . . . . .	5
1.5 The CellFlow algorithm . . . . .	7
1.6 Parameters in CellFlow . . . . .	8
1.7 Implementation of CellFlow . . . . .	10
<b>2 Metrics</b>	<b>10</b>
2.1 R squared . . . . .	10
2.2 Maximum Mean Discrepancy . . . . .	11
2.3 Energy distance . . . . .	12
2.4 Sinkhorn divergence . . . . .	12
2.5 Evaluation of metrics . . . . .	13
<b>3 Applications</b>	<b>13</b>
3.1 PBMCs treated with cytokine . . . . .	13
3.2 Developing zebrafish . . . . .	15
3.3 Sciplex data . . . . .	16
3.4 Combosciplex data . . . . .	17
3.5 Perturb-seq for gene overexpression . . . . .	18
3.6 Perturb-seq with pathway activations . . . . .	19
3.7 4i data . . . . .	20
3.8 NGN2-induced neuron morphogen screen . . . . .	20
3.9 Brain organoid morphogen screens . . . . .	21

## 1 CellFlow

Given a population of control cells and an experimental intervention, CellFlow aims to generate a perturbed population of cells. We represent the phenotypes of both control (source) and perturbed (target) cells as vectors in a multidimensional embedding space created from principal component analysis (PCA) or variational autoencoders (VAEs). To generate a prediction, CellFlow takes into account all variables defining an experimental intervention and embeds them to a single representation vector, which then guides the transformation of control cells into perturbed cells.

This transformation is governed by a neural ordinary differential equation <sup>1</sup> (neural ODE), which is parameterized by a time-dependent neural vector field. The neural vector field is a feedforward neural network that takes as inputs encoded experimental conditions, time points at which the vector field should be evaluated, and interpolated cell phenotype vectors at those time points. Here, time serves as an abstract way to represent a transition from the initial ( $t = 0$ ) to the perturbed ( $t = 1$ ) cellular states. The network outputs velocity vectors defining the flow between these states, and the predicted phenotypes of perturbed cells are obtained by integrating the vector field from  $t = 0$ , where cell phenotypes are those of the given source cells, to  $t = 1$  using a numerical ODE solver.

We train the neural vector field using flow matching and optimal transport. Flow matching <sup>2</sup> is a simulation-free way of training neural vector fields, enabling fast and stable optimization. Optimal Transport <sup>3,4</sup> (OT) aligns control and perturbed cells by minimizing the total displacement cost, and has been shown to accelerate and improve flow matching training <sup>5-7</sup>. At each training iteration, a batch of source cells, an experimental condition (e.g. drug or gene knockout), and a batch of target cells corresponding to this condition are randomly sampled. The neural vector field is evaluated at time points which are sampled uniformly from  $[0, 1]$ . Subsequently, source and target cells from the sampled batches are paired according to the optimal transport solution for straightening paths of the neural vector field which connects control and perturbed cells. The encoder for experimental conditions and the neural vector field are then optimized end-to-end, thus implicitly learning a functional representation of the condition.

In the following, we elaborate on the most relevant concepts underlying CellFlow.

### 1.1 One formulation to unify the setup of a perturbation experiment

To cover a wide range of perturbation experiments, we formulated a generic setup of phenotypic screens (Figure 1c). We partition the set of variables defining a perturbation experiment into three groups, namely perturbations  $\{p_i\}_{i \in \mathcal{I}_p}$ , perturbation covariates  $\{q_i\}_{i \in \mathcal{I}_q}$ , and sample covariates  $\{s_i\}_{i \in \mathcal{I}_s}$ , with  $\mathcal{I}_p, \mathcal{I}_q$ , and  $\mathcal{I}_s$  denoting index sets. Perturbations include all interventions like drug treatments, gene knockouts, and cytokine treatments, while perturbation covariates add complementary information about the perturbation, e.g. the dosages of drugs or the timing of the intervention. Perturbations (and thus perturbation covariates) can be combined together, for instance, multiple drugs can be given. In contrast, we define sample covariates to capture the cellular state independent of the perturbation. Examples are categorizations of cells into cell lines, tissues, donors, or batches. Sample covariates are not combinable; for instance, a cell cannot belong to multiple batches at the same time.

Thus, a condition  $\mathbf{c}$  a 3-tuple defining the intervention and cellular state  $\mathbf{c} = (\{p_{\mathbf{c}_j}\}_{j \in n_c}, \{q_{\mathbf{c}_j}\}_{j \in n_c}, s_{\mathbf{c}})$ . Here,  $n_c = \{1, \dots, n\}, n \in \mathbb{N}$  denotes the number of perturbations (and thus perturbation covariates) of condition  $\mathbf{c}$ . Moreover, condition  $\mathbf{c} \in \mathcal{C}$  is uniquely tied to a pair of source and target distribution  $(\mu_{\mathbf{c}}, \nu_{\mathbf{c}})$ .

### 1.2 Learning a single representation from all perturbation factors

To learn a representation of the above mentioned factors defining a perturbation experiment, all factors have to be made machine-readable, and subsequently aggregated in a permutation-invariant manner. While there are multiple ways to represent any perturbation, and we leave a systematic evaluation to future work, we employ embeddings from ESM2 <sup>8</sup> (Evolutionary Scale Modeling) to represent cytokines and genetic perturbations, while we leverage molecular fingerprints for drugs <sup>9,10</sup>. If a representation is not available for the applied treatments, they can be one-hot encoded, which only allows predictions of their unseen combinations but not of unseen individual perturbations. As perturbation covariates are typically scalar values (i.e. time of treatment, dosage of treatment), we transform the values to a near-linear scale and neural network-typical range (e.g. applying log-transformation to dosages). However, a user can provide any representation for encoding perturbation covariates. Analogously to perturbations, there is no single representation for all the possible sample covariates. We typically leverage embeddings of the Cancer Cell Line Encyclopedia <sup>11</sup> to embed cancer cell lines. For representing patients, we leverage statistics of donor-specific control populations. When learning across datasets, we encode those using one-hot encoding. We













**Encoder module** Each perturbation, covariate, and cell line embedding can be further processed by a feedforward neural network (`mlp`) or self-attention module (`self-attention`) before being concatenated, the entirety of which we refer to `layers_before_pool`. These concatenated vectors are then pooled using standard attention (`attention_token`), set attention (`attention_seed`), or mean-pooling/deep sets (`mean`). When the set of vectors to combine is of size one, attention falls back to self-attention, while pooling is equivalent to the identity. The resulting vector is then fed forward through another feedforward neural network (`layers_after_pool`), before it is projected to a vector of size `condition_embedding_dim`, the output of the encoder module. In summary, we have the following parameters together with default values in parentheses

- Pooling (`attention_token`)
- `cond_output_dropout` (0.9)
- `condition_embedding_dim` (256)
- `layers_before_pool` (specified for perturbation, perturbation covariate and sample covariate):
  - `layer_type` (`mlp`)
  - `dims` ([1024, 1024])
  - `dropout_rate` (0.0)
- `layers_after_pool`
  - `layer_type` `mlp`
  - `dims` ([1024, 1024])
  - `dropout_rate` (0.0)

**Flow Matching module** The flow matching module transforms the source (control) distribution to a target distribution (perturbed) given a condition. The `time_encoder_dims` is a feed-forward MLP encoding the sinusoidal embedding of dimension `time_freqs` of the time (not experimental time, but time of the neural differential equation). The `hidden_dims` specify the architecture of the MLP processing samples (cells) in the source (control) distribution. The `decoder_dims` define the layers of the MLP processing the concatenated vectors (linearly projected if `linear_projection_before_concatenation`, and layer-normalized if `layer_norm_before_concatenation`) of the condition embedding (from the condition encoder), the time embedding (output of `output_dims`), and the processed cell embedding (output of `hidden_dims`). The noise schedule can be chosen arbitrarily (see 1.4), we recommend adding constant Gaussian noise (`constant_noise`) with standard deviation `flow_noise`, or noise following the noise schedule of a Schroedinger bridge bridge.

- `time_freqs` (1024)
- `time_encoder_dims` ([2048, 2048, 2048])
- `time_encoder_dropout` (0.0)
- `hidden_dims` ([4096, 4096, 4096])
- `hidden_dropout` (0.0)
- `linear_projection_before_concatenation` (False)
- `layer_norm_before_concatenation` (False)
- `decoder_dims` ([4096, 4096, 4096])
- `decoder_dropout` (0.0)
- `flow_type` (`constant_noise`)
- `flow_noise` (0.1).

**Training** In contrast to most other deep learning methods, samples in one batch (of size `batch_size`) are not independent of each other due to the resampling step according to the OT plan. The number of iterations (`num_iterations`) is a hyperparameter which has to be optimised; due to the relatively costly inference, early stopping is not as cheaply available as in one-step models. `multi_steps` is one of the most relevant parameters. It denotes the number of conditions (not cells) seen between gradient updates. It inherits its name from its implementation, i.e., due to the data structure, we first sample a source distribution (sample covariate), and then corresponding perturbations and perturbation covariates. Given this configuration of conditions, we sample cells from the control and perturbed population, denoting one "step" of the training procedure. As this corresponds to only seeing one data point in the condition encoder module, we increase the number of steps to `multi_steps`, for the condition encoder module not to overfit. In summary, we have the following parameters together with default values

- `batch_size` (1024)
- `num_iterations` (500\_000)
- `multi_steps` (20)
- `optimizer` (Adam as implemented in optax<sup>31</sup>)
- `learning_rate` (0.00005).

## 1.7 Implementation of CellFlow

CellFlow is implemented in JAX<sup>32</sup>, and thus makes use of its deep learning ecosystem including `flax`, `optax`, and `diffraex`<sup>33</sup>. We make use of `ott-jax`<sup>24</sup> for pairing cells with optimal transport, and also use their implementations for parts of the flow matching module. For preparing and processing the data, we build upon the scverse ecosystem<sup>34</sup> including `anndata`<sup>35</sup> and `scanpy`<sup>36</sup> as well as its GPU-accelerated counterpart `rapids-singlecell`. For CellFlow's preprocessing module, we build upon `pertpy`<sup>37</sup> and HuggingFace.

## 2 Metrics

In this section, we introduce several distributional metrics to evaluate how well generative models capture the underlying distribution of perturbed cells. We note that not all of them are actual metrics, but for the sake of readability we refer to them as such.

### 2.1 R squared

R squared ( $R^2$ ) between the true and predicted mean gene expression is arguably the most common metric used in the perturbation modeling community<sup>38–41</sup>. Given true samples  $\{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^n \sim_{\text{iid}} \mu$  and predicted samples  $\{\mathbf{y}_i \in \mathbb{R}^d\}_{i=1}^n \sim_{\text{iid}} \nu$ , we first compute the empirical means

$$\tilde{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \in \mathbb{R}^d, \quad \tilde{\mathbf{y}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \in \mathbb{R}^d,$$

and then define the empirical  $R^2$  between these two empirical means as

$$R^2(\hat{\mu}_n, \hat{\nu}_n) = 1 - \frac{\sum_{j=1}^d (\tilde{\mathbf{y}}_j - \tilde{\mathbf{x}}_j)^2}{\sum_{j=1}^d (\tilde{\mathbf{x}}_j - \bar{\tilde{\mathbf{x}}})^2}.$$

with

$$\bar{\tilde{\mathbf{x}}} = \frac{1}{d} \sum_{j=1}^d \tilde{\mathbf{x}}_j$$

There are two major advantages of the  $R^2$  metric for evaluation of perturbation predictions. First, it is independent of any hyperparameter choices. Second, it does not suffer from the curse of dimensionality, i.e. it can be evaluated in very high dimensions without losing expressivity, which distinguished this metric from the remaining ones.

However,  $R^2$  values must be handled with care. In particular, it only takes into account the mean, i.e. the first moment of a distribution. This entails that it does not distinguish between dimensions of little and dimensions of much variation, weighing each axis equally. Thus, while it might be useful for evaluating predictions on "simple" distributions like cell lines, it is not suitable for evaluating heterogeneous populations.





















**CellFlow training** For training the model, we used 30-dimensional PCA space. The embedding used for training was always constructed using only cells from the training set. All of the morphogen signaling modulators were one-hot encoded and the one-hot encodings multiplied by the corresponding dosages were used as input perturbation representations.

**Baselines** As a baseline for any combination of treatment outcomes, we defined two baseline models:

- **Union.** Here, we computed the union of cells from the individual treatment conditions. This baseline assumes that a combination of treatments will result in the combination of states produced from individual treatments.
- **Barycenter.** Here, we computed the Wasserstein barycenter<sup>64</sup> between the distributions produced by individual treatments using the FreeWassersteinBarycenter implementation in ott-jax with a Sinkhorn solver and a fixed  $\epsilon$  of 0.1. This baseline assumes that a combination of treatments will result in intermediate states between states produced from individual treatments.

**Hyperparameter optimization** We performed hyperparameter tuning over a fixed set of hyperparameter values using the Bayesian optimization-based method Optuna<sup>65</sup>. For CellFlow, a sweep was optimized to minimize E-distance on a held-out test set comprising the combinatorial treatments CHIR+BMP4 and FGF8+SHH. These conditions were excluded from any subsequent model evaluations. The model with the best MMD was used for evaluation. Biolord and CPA were hyperparameter-tuned to minimize MMD on the same holdout split using a grid of hyperparameters of comparable size to that of CellFlow. Each run of Optuna included between 400 and 600 of trained models. All optimized hyperparameters are reported in Supplementary Table 8.

**Evaluation** For all of the models, a full gene expression matrix was reconstructed from the obtained predictions, with subsequent projection to the 20-dimensional PCA space of the whole ground truth dataset for evaluation. In addition to computing distributional metrics, predicted cells were assigned to the ground truth data clusters using weighted k-nearest neighbors, as described for the annotation of the human neural organoid cell atlas based on the primary atlas<sup>66</sup>. We then computed the cosine similarity between predicted ( $x$ ) and ground truth ( $y$ ) cluster compositions as

$$1 - \frac{x \cdot y}{\|x\|_2 \|y\|_2}$$

To visualize predictions, we projected them onto the UMAP embedding computed from the full dataset using UMAP.fit and UMAP.transform from cuML (<https://github.com/rapidsai/cuml>).

**Morphogen interaction score** To analyze interactions between treatments and to assess whether combinations result in new cell states, we computed the MMD between the cells from combinatorial condition, either predicted or ground truth, and the union of cells from the corresponding single-treatment conditions.

### 3.9 Brain organoid morphogen screens

To train and evaluate CellFlow for predicting organoid engineering outcomes, we obtained three scRNA-seq datasets of organoid morphogen screens. The Amin dataset<sup>67</sup> was obtained from the link provided in the publication, while the other two datasets<sup>68,69</sup> (Azbukina and Sanchis-Calleja) were provided by the authors. All datasets were consistently preprocessed using scanpy<sup>36</sup>.

**Integration and annotation** To integrate all three datasets, we used HNOCA-tools (<https://github.com/devsystemslab/HNOCA-tools>) and a scANVI model<sup>70</sup> that was previously pre-trained on an atlas of the developing human brain<sup>66,71</sup>. We performed query-to-reference mapping using the model on the organoid datasets using the datasets label as the batch covariate. We trained the model with

- `retrain = "partial"`
- `max_epochs = 200`
- `batch_size = 1024`
- `weight_decay = 0.0`

and otherwise default parameters. This shared latent space with the three organoid datasets as well as the primary cell atlas also allowed us to obtain consistent brain region and cell class annotations for all organoid datasets. For this, we used CellFlow's `compute_wknn` and `transfer_labels` functions to transfer labels as described previously<sup>66</sup>.

**Holdout splits** To evaluate different aspects of model performance, we constructed two types of holdout splits:

- **Combination prediction.** Here, we sought to test the ability of the model to predict the outcome of unseen combinatorial morphogen treatments. For this, we selected morphogen combinations where conditions were available for both individual and combinatorial treatments, resulting in 12 separate splits. During each training run, we included individual treatments but held out their combination and all conditions including this combination from training.
- **Information transfer.** Here, we tested to what degree a condition seen in one dataset could inform predictions on other datasets. For this, we selected morphogens that were used in all three datasets, but in no more than 30% of conditions for each dataset, resulting in 7 separate splits. During each training run, we held out all conditions including this morphogen in one dataset from training.

All splits for both tasks are listed in Supplementary Table 9. For each holdout split, we retrained the scANVI model on the training data only, to obtain the latent embedding for training CellFlow and for computing baseline models.

**CellFlow training and digital protocol encodings.** We trained CellFlow on all three combined datasets using the shared latent space of the scANVI model, which was recomputed for each split. To construct conditions for training, we obtained digital protocol encodings as follows: For each morphogen treatment, the morphogen was one-hot encoded and the one-hot encoding multiplied by the corresponding log-transformed dosages. We concatenated this representation of the morphogen with the start and end time of treatment (in days) as well as a one-hot encoding of the modulated pathway and the directions of modulation (-1 for inhibition, 1 for activation). Each condition was represented as a combinatorial set of such treatments and the dataset label was used as a sample covariate. Because not all datasets had controls that were not treated with any morphogen, we derived a common source distribution of 10000 cells by repeatedly computing the mean of 10 cells sampled from the primary atlas of the developing human brain. We trained CellFlow to generate predictions from this source distribution for all datasets based on the set of treatments and the dataset covariate.

**Baselines** As a baseline for combinatorial treatment prediction we used the Union and Barycenter baselines as described above. For the information transfer task, we used an additional baseline:

- **Train dataset.** Here, we used all conditions in the respective dataset excluding the held-out test conditions as a baseline. This baseline assesses whether the model learns relevant treatment-specific effects beyond the effects of the base protocol used in each dataset, i.e. the dataset covariate label.

**Hyperparameter optimization** We performed hyperparameter tuning over a fixed set of hyperparameter values using the Bayesian optimization-based method Optuna<sup>65</sup>. We optimized for 200 iterations to minimize MMD on a holdout split where all conditions including BMP4 in the Sanchis-Calleja dataset<sup>69</sup> were held out during training. All optimized hyperparameters are reported in Supplementary Table 10.

**Evaluation** For evaluation of CellFlow and baseline models, we computed distributional metrics directly in scANVI latent space as described above. In addition to computing distributional metrics, predicted cells were assigned to the ground truth data clusters using WKNN-based label transfer from the full dataset in scANVI latent space and cosine similarity was computed on these cluster compositions as described above. To visualize predictions, we projected them onto the UMAP embedding computed from the full dataset using `UMAP.fit` and `UMAP.transform` from cuML (<https://github.com/rapidsai/cuml>).

**Virtual neural organoid protocol screen** To perform a virtual organoid protocol screen, we first composed protocols based on single, double and triple combinations of 14 morphogen pathway modulators as well as their concentration and timing. For all morphogens we fixed the concentration to the highest concentration observed across all datasets. To comprehensively probe the space of possible protocols without ending up with intractably many configurations, digital protocols were computed in two parts:





- Cell Atlas Organoid Biological Network, J Gray Camp, Fabian J Theis, and Barbara Treutlein. An integrated transcriptomic cell atlas of human neural organoids. *Nature*, 635(8039):690–698, 2024.
- 67. Neal D. Amin, Kevin W. Kelley, Konstantin Kaganovsky, Massimo Onesto, Jin Hao, Yuki Miura, James P. McQueen, Noah Reis, Genta Narazaki, Tommy Li, Shravanti Kulkarni, Sergey Pavlov, and Sergiu P. Pasca. Generating human neural diversity with a multiplexed morphogen screen in organoids. *Cell Stem Cell*, 31(12):1831–1846.e9, 2024. ISSN 1934-5909. doi: <https://doi.org/10.1016/j.stem.2024.10.016>.
  - 68. Nadezhda Azbukina, Zhisong He, Hsiu-Chuan Lin, Małgorzata Santel, Bijan Kashanian, Ashley Maynard, Tivadar Török, Ryoko Okamoto, Marina Nikolova, Sabina Kanton, Valentin Brösamle, Rene Holtackers, J. Gray Camp, and Barbara Treutlein. Multi-omic human neural organoid cell atlas of the posterior brain. *bioRxiv*, 2025. doi: 10.1101/2025.03.20.644368.
  - 69. Fátima Sanchís-Calleja, Akanksha Jain, Zhisong He, Ryoko Okamoto, Charlotte Rusimbi, Pedro Rífes, Gaurav Singh Rathore, Małgorzata Santel, Jasper Janssens, Makiko Seimiya, Jonas Simon Fleck, Agneta Kirkeby, J. Gray Camp, and Barbara Treutlein. Decoding morphogen patterning of human neural organoids with a multiplexed single-cell transcriptomic screen. *bioRxiv*, 2024. doi: 10.1101/2024.02.08.579413.
  - 70. Chenling Xu, Romain Lopez, Edouard Mehlman, Jeffrey Regier, Michael I Jordan, and Nir Yosef. Probabilistic harmonization and annotation of single-cell transcriptomics data with deep generative models. *Molecular Systems Biology*, 17(1):e9620, 2021. doi: <https://doi.org/10.1525/msb.20209620>.
  - 71. Emelie Braun, Miri Danan-Gothold, Lars E. Borm, Ka Wai Lee, Elin Vinsland, Peter Lönnberg, Lijuan Hu, Xiaofei Li, Xiaoling He, Žaneta Andrusiová, Joakim Lundberg, Roger A. Barker, Ernest Arenas, Erik Sundström, and Sten Linnarsson. Comprehensive cell atlas of the first-trimester developing human brain. *Science*, 382(6667):eadf1226, 2023. doi: 10.1126/science.adf1226.
  - 72. Jerome Friedman, Robert Tibshirani, and Trevor Hastie. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1):1–22, 2010. doi: 10.18637/jss.v033.i01.