

# HYPED: A MULTIMODAL HYBRID PERTURBATION GENE EXPRESSION AND IMAGING DATASET

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Integrating multimodal, high-resolution biological data is a useful way to characterize biological processes, such as how cells respond to perturbations. Cell perturbation prediction is a major experimental challenge and has motivated substantial research in machine learning for biology. In this work, we generated a multimodal benchmark dataset that captures the dynamic response of human fibroblasts to transient transcription factor perturbations. We performed time-series live cell imaging with fluorescent cell cycle reporters over 72 hours and collected long-read single-cell RNA sequencing data from the same population of cells. We release the processed dataset, preprocessing pipelines and benchmarking code along with the evaluation of existing models using our data as ground truth. This work supports the development and evaluation of machine learning methods for modeling dynamical systems from multimodal datasets. HYPED consists of RNA sequencing data from approximately 20,000 cells and 203 imaging timepoints across four experimental conditions, totaling 2030 imaging frames. HYPED makes the cell perturbation problem accessible to machine learning researchers with state-of-the-art experimental data.

## 1 INTRODUCTION

Direct cell reprogramming, or the process of converting one cell type into another without passing through an intermediate pluripotent state, holds promise for personalized and regenerative medicine (65; 53). Reprogramming requires the introduction of targeted perturbations capable of altering cell state, most often through the introduction of cell-type specific transcription factors (TFs) (22; 61). Current methods for delivering TFs, such as lentiviral vectors or CRISPR-based systems, permanently alter the cell’s DNA. This irreversible modification raises concerns about uncharacterized and potentially harmful effects, limiting their utility for a wide range of research and clinical applications. To address this limitation, TF delivery via transient RNA-based methods have emerged as safer alternatives (57; 23).

Despite recent progress in the field of direct reprogramming, key challenges remain: conversion efficiency remains low and identifying new TF combinations is difficult due to limited knowledge of their combinatorial effects (31; 61). Evaluating the effects of perturbations is essential for optimizing existing cell reprogramming protocols. However, most existing datasets are limited to a single data modality and do not use transient TF delivery methods (16; 49; 70; 20; 13; 4).

This lack of adequate data poses a challenge for evaluating and benchmarking machine learning models for cell perturbation prediction. In this work, we collected high-resolution, long-read single-cell transcriptomic data and live-cell fluorescent microscopy data from a state-of-the-art TF reprogramming protocol on human fibroblasts using a combination of transient forced overexpression (*MYOD1*) and gene suppression (*PRRX1*). Our dataset is a valuable resource for the machine-learning evaluation of cell perturbations.

**Biological Background.** The cell reprogramming problem originated in the 1980s with H. Weintraub’s use of the TF MyoD to convert fibroblasts into skeletal muscle (66). In the early 2000s, S. Yamanaka reprogrammed the first induced pluripotent stem cells (iPSC) (53). Both approaches typically use viral vectors to permanently integrate modified genes into the host cell genome. Many contemporary perturbation screens use a CRISPR-based approach (e.g., Perturb-seq) to modulate

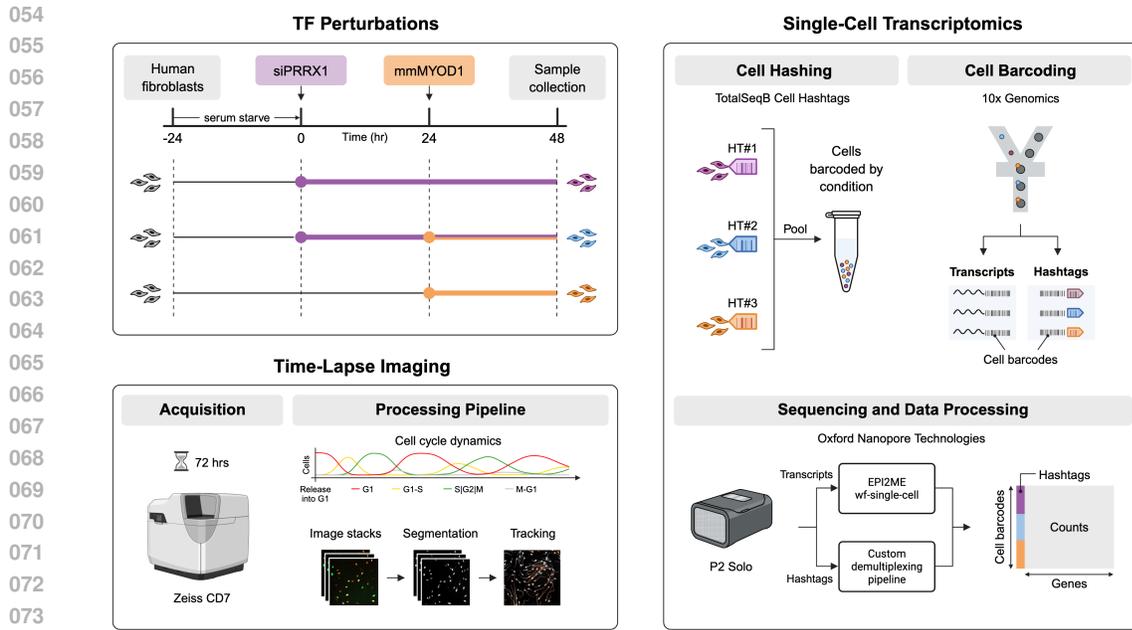


Figure 1: **Experimental Dataset.** The HYPED dataset includes time-series imaging and long-read sequencing to characterize transient cell perturbations. **Top left:** Human fibroblasts are perturbed using combinations of *MYOD1* overexpression and *PRRX1* attenuation, delivered via modified mRNA and siRNA, respectively. **Bottom left:** Live-cell fluorescent imaging with the Zeiss Celldiscoverer 7 (CD7) captures the temporal dynamics of the cell cycle under TF perturbation. **Right:** Long-read single-cell sequencing captures the final transcriptional state of the perturbed cells.

gene expression, but also irreversibly modify the host cell genome (10). These methods have produced useful datasets for perturbation prediction using machine learning. However, genome modification itself can result in off-target effects, making non-integrating approaches desirable from a safety and quality standpoint.

Non-integrating methods including modified mRNA (mmRNA) and small interfering RNA (siRNA) have been successfully used for direct reprogramming (62; 60; 56). mmRNA enables transient, high-level expression of reprogramming TFs without genomic integration, while siRNA facilitates the attenuation of endogenous mRNAs that maintain the original cell identity. Together, these approaches promote efficient and controlled lineage conversion without permanent genetic alteration (50; 34). To our knowledge, we present the first cell perturbation dataset based on state-of-the-art experimental methods for machine learning.

**Data Modalities.** HYPED is a multimodal dataset that combines high-throughput sequencing with live-cell imaging (fig. 1). Recent advancements in long-read omics technologies, such as single-cell RNA sequencing (scRNA-seq), have expanded our ability to investigate complex intracellular mechanisms across multiple regulatory levels (26; 58; 25). In particular, scRNA-seq with Oxford Nanopore Technologies enables the direct sequencing of full-length RNA transcripts, providing insights into transcriptional regulation and cell state dynamics. The long-read capability of Nanopore sequencing allows for the identification of alternative splicing events, quantification of isoform diversity, and detection of RNA modifications, which are often inaccessible features through short-read platforms (24; 26). This enables a more accurate and comprehensive characterization of the transcriptome, facilitating the study of gene expression heterogeneity at the single-cell level.

RNA sequencing remains one of the most widely used data modalities to assess cell responses to perturbation (43; 55; 45). However, sequencing requires cell lysis, resulting in destruction of the cell during sample preparation. In contrast, advances in live-cell imaging technologies have significantly improved spatial resolution and throughput, offering a non-destructive, time-resolved method for

Table 1: **Overview of Related Single Cell Datasets.** Overview of benchmark datasets used in machine learning for cell imaging, sequencing, and perturbation.

Dataset ML Task		Imaging	Sequencing	Perturbation
<i>Only Imaging</i>				
LIVECell	(11)	✓		
NeurIPS Competition 2022	(33)	✓		
<i>Only Sequencing</i>				
CellXGene	(41)		✓	
NeurIPS Competition 2021	(32)		✓	
<i>Multi-Modal Perturbation (Perturbation + Sequencing)</i>				
scPerturb	(18)		✓	✓
PerturbBase	(64)		✓	✓
sc-pert	(27)		✓	✓
NeurIPS Competition 2023	(7)		✓	✓
<i>Multi-Modal (Imaging + Perturbation + Sequencing)</i>				
<b>HYPED (ours)</b>		✓	✓	✓

studying dynamic cell processes. These systems enable high-resolution monitoring of features like cell morphology, proliferation, migration, and apoptosis over time. When combined with fluorescent reporters such as the Fluorescent Ubiquitination-based Cell Cycle Indicator (FUCCI) system (17), live-cell imaging allows real-time tracking of specific protein expression to explore intracellular drivers of phenotypic changes (48). FUCCI distinguishes cell cycle phases by leveraging fluorescently-tagged proteins whose degradation is regulated by the cell cycle. Specifically, the Incucyte Cell Cycle Reporter (Sartorius 4779), uses the FUCCI system to capture cell cycle dependent changes in the expression patterns of Geminin and Cdt1 by linking fluorescent proteins TagGFP2 (Green) and mKate2 (Red) that can then be monitored with live-cell imaging.

This approach provides a dynamic view of cell cycle progression in live cells (40) that are not accounted for in existing machine learning perturbation models. In the HYPED dataset, time-series imaging offers high temporal resolution to capture perturbation dynamics, complementing sequencing-based approaches. This addition introduces a data modality for machine learning studies of cell perturbations, enabling multimodal and temporal analysis of the perturbation prediction problem (see table 1).

**Machine Learning Approaches.** Given the significance of the cell reprogramming problem, a wide range of machine learning approaches have been developed. Early efforts include graph- and network-based methods (42; 12), as well as studies that frame reprogramming through system identification, trajectory optimization, and control theory (43; 6). More recent work leverages transformers and other modern architectures to represent and embed cell states measured with sequencing (55; 38; 9). Although these approaches have achieved successes at the machine learning problem, there remain substantial gaps in the success of experimental perturbations, which motivates the development and application of new models to state-of-the-art experimental datasets.

**Contributions.** Imaging and transcriptomics each offer distinct yet complementary insights into cell biology; imaging captures spatial and morphological dynamics, while transcriptomics reveals the molecular and regulatory state of the cell. However, these modalities are typically collected in isolation, providing only fragmented views of genetic perturbations. Integrating these experimental modalities promises a more comprehensive understanding of the cell as a dynamical system, revealing how gene expression and physical behavior co-evolve over time. This integration is technically challenging, as it involves harmonizing data with vastly different structures, resolutions, and temporal characteristics. Emerging machine learning approaches, particularly multimodal representation learning and generative modeling, offer a powerful bridge, enabling alignment, translation, and predictive modeling across these disparate data types. Leveraging these computational frameworks

162 is key to unlocking unified models of cell behavior that incorporate both molecular and phenotypic  
163 dimensions.

164 To address this challenge, we designed and performed a perturbation experiment that captured  
165 live-cell imaging and scRNA-seq from the same initial population of cells that were treated under  
166 identical conditions and cell cycle synchronized, providing a state-of-the-art ground truth dataset  
167 for a variety of challenges, including improving perturbation prediction outcomes. We overcome  
168 several limitations of existing datasets, which lack synchronized, time-resolved measurements across  
169 modalities. Our contributions are a valuable resource for benchmarking multimodal models as well  
170 as a framework for studying dynamic cell processes through an integrative lens.  
171

## 172 2 RELATED WORK

### 173 2.1 CELL REPROGRAMMING

174 Direct cell reprogramming is a complex process that involves activating lineage-specific target net-  
175 works while simultaneously silencing the original cell identity (36). Most reprogramming strategies  
176 rely on viral vectors or plasmids to deliver lineage-specifying genes. Pioneer factors, such as *MYOD1*  
177 in muscle (66) or *OCT4* in pluripotency (53), can initiate reprogramming but often require cooper-  
178 ative TFs to fully activate and maintain the target gene regulatory network (19). Without the  
179 proper TFs, most cells fail to convert fully, instead stalling in unstable intermediate states (8; 37).  
180 Researchers have mapped regulatory networks that govern cell identity and have developed libraries  
181 of transcription factors capable of inducing cell transitions (47; 69; 67; 28).  
182

183 In parallel, RNA-based approaches have gained traction as non-integrative tools for cell repro-  
184 gramming (63; 62; 60). Synthetic mmRNA delivery allows for transient and tunable expression of  
185 reprogramming factors, minimizing risk of genomic integration and long-term mutagenesis (50; 6; 5).  
186 Similarly, siRNA can be used to transiently suppress endogenous factors that stabilize native cell  
187 identity (56; 30). While RNA-based systems offer less fine-tuned control compared to viral or  
188 plasmid-based delivery methods, they have significantly greater clinical promise (23). Their transient  
189 nature is well suited to control the cell cycle and reduces the risk of mutagenesis and tumorigenesis,  
190 making them attractive for therapeutic applications (3). As RNA delivery technologies and engineer-  
191 ing strategies mature, RNA-based reprogramming offers a safe path to clinical translation without  
192 compromising efficacy (31).  
193

### 194 2.2 DATASETS AND MODELS FOR GENE PERTURBATIONS

195 Gene perturbations are a core laboratory technique for cell reprogramming and have become a focus  
196 of machine learning in biology. RNA-seq is the primary technique for measuring cell responses  
197 to perturbations, with scRNA-seq enabling high-resolution analysis of development, disease, and  
198 treatment effects. This has led to the development of large single cell atlases (21; 52; 14; 54) and  
199 foundation models (55; 15; 51). Together, these datasets and models provide a rich resource for  
200 developing machine learning strategies to infer cell identity and predict perturbation outcomes.  
201

202 Perturb-seq, which combines CRISPR perturbations with scRNA-seq, enables high-throughput  
203 measurement of gene function and has led to the creation of several public datasets and resources  
204 (10). The Gene Perturbation Atlas (GPA) compiles single-gene perturbation data across diverse cell  
205 types to systematically assess how individual genes influence cell identity (68). The Perturbation Cell  
206 and Tissue Atlas (PCTA) extends this effort by integrating genetic and environmental perturbations  
207 with molecular and imaging readouts to support causal inference in cell and tissue biology (44).  
208

209 Early computational frameworks for direct reprogramming combine gene expression, regulatory  
210 networks, and epigenetic information to predict optimal TF combinations (29; 43; 42; 12). However,  
211 the low efficiency of current reprogramming protocols has positioned cell perturbation prediction  
212 as a standard machine learning challenge (45; 55; 27; 72). While many models focus on forward  
213 prediction—forecasting the cell response to a given perturbation—few address how to identify the  
214 optimal perturbation strategy. Moreover, existing methods often overlook important biological priors  
215 such as cell dynamics, cell cycle phase, and other factors known to influence reprogramming success  
(71).

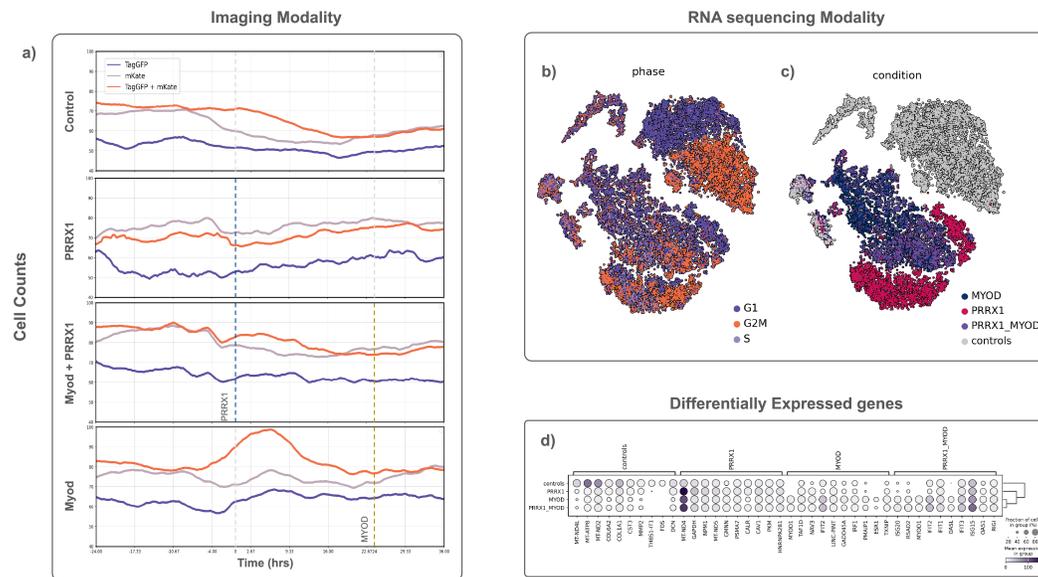


Figure 2: **HYPED Dataset Overview.** (a) Shows the Distribution of cells across different cell cycle phases, quantified from Incucyte imaging data across the different fluorescence channels. The (b) t-SNE projection of cells colored by cell cycle phase shows (c) t-SNE projection of cells colored by experimental condition. (d) Top 10 differentially expressed genes across experimental conditions.

### 3 DATASET

We present a multimodal dataset that captures the dynamic response of human fibroblasts to transcription factor (TF) perturbations. The dataset includes long-read scRNA-seq and time-lapse imaging data collected under four experimental conditions: *MYOD1* activation via modified mRNA (mmMYOD1), *PRRX1* suppression via siRNA (siPRRX1), sequential dual perturbation (siPRRX1 followed by mmMYOD1), and an unperturbed control. The scRNA-seq dataset comprises approximately 20,000 cells (~11,000 perturbed, ~9,000 controls), enabling isoform-level transcriptomic profiling. Time-lapse imaging data were collected at 20 minute intervals over 72 hours, capturing high-resolution trajectories of cell cycle dynamics and morphological changes. Sequencing data contains of 3 replicates in separate wells, that were pooled together for single-cell barcoding and sequencing. For the Imaging data, the negative controls and the MYOD + PRRX1 condition each contain three replicates, while the MYOD and PRRX1 conditions contain two replicates each. All replicates are included within the released dataset.

#### 3.1 EXPERIMENTAL SETUP

Human neonatal foreskin fibroblasts (BJ, ATCC CRL-2522) carrying the Incucyte Cell Cycle Reporter (Sartorius 4779) were used for perturbation experiments. Cells were cultured at 37°C in 5% CO<sub>2</sub> on standard cultureware in full media (FM; Dulbecco's Modified Eagle's Medium (DMEM, Gibco 11965-092) + 10% Fetal Bovine Serum (FBS, Corning 35-015-CV) + 1% Penicillin-Streptomycin (P/S, Gibco 15140-122)). Cells were seeded at 1.03 - 1.33 x 10<sup>4</sup> cells/cm<sup>2</sup> in 6- and 48-well plates, followed by 24 hours of serum starvation (0.2% FBS) for G0/G1 phase synchronization. At t=0h, cells were released from serum starvation, and respective cells were transfected with PRRX1 siRNA (25 nM; Dharmacon) in FM using the TransIT-X2 Dynamic Delivery System (Mirus MIR6000). At t=24h, siRNA was washed out with FM. Respective cells were then transfected with MYOD1 modified mRNA (1 ng/μl; Pseudouridine + Silica purification, Trilink WOTL39876) in FM supplemented with 2 μM 4OH-Tamoxifen (Sigma-Aldrich H7904) using the Lipofectamine MessengerMAX Transfection Reagent (Invitrogen 100026485). At t=48h, cells were harvested with TrypLE Express (Gibco 12604-013) for scRNA-seq sample preparation. For imaging plates, cells were treated the

Table 2: **Imaging Data Structure.** These dimensions are used to operate on the imaging files, named according to CZI specification.

Symbol	Description
H	Phase
S	Scene (plate/well)
T	Time
C	Channel
Z	Depth dimension
M	Mosaic (determines sub-tile)
Y	Vertical axis
X	Horizontal axis

same as described above with two differences: (1) Cells were cultured in imaging media (FluoroBrite DMEM (Gibco A18967-01) + 10% FBS + 1% P/S) instead of standard FM; (2) Cells were stained with 250 nM SiR-DNA (Cytoskeleton, Inc. CY-SC007) at t=-24h to enable tracking of cell nuclei.

### 3.2 SINGLE CELL RNA SEQUENCING WITH LONG READS

Cells from siPRRX1, siPRRX1/mmMYOD1, and mmMYOD1 were individually labeled with TotalSeq-B human cell hashing antibodies (BioLegend, Cat# 394631, #394633, and #394635, respectively). 20,000 cells from each condition were mixed together (viability >98%) and single cell barcoded on the 10x Genomics Chromium Controller (X) using the Next GEM Single Cell 3' Kit V4. Barcoded cDNA amplicons from transcript and hashing libraries were prepared according to the Oxford Nanopore Technologies (ONT) 3' cDNA protocol (SQK-LSK114, SST\_9198\_v114\_revJ\_13Nov2024). Libraries were assessed for quality following ONT recommendations. Prepared libraries were sequenced on the ONT PromethION Solo 2 (P2) sequencer. Raw reads were base-called with Dorado v0.9.1 (2) using the High Accuracy basecalling model and stored as fastq files.

For unperturbed controls, fibroblasts were sorted by fluorescence activated cell sorting (FACS) into their respective cell cycle phases after staining with 16  $\mu$ M Hoechst 33342 for 50 minutes. Cells in G1, S, and G2/M were individually labeled with TotalSeq-B hashing antibodies (BioLegend, Cat# 394631, #394603, #394605, respectively), mixed, and single cell barcoded (viability >97%). Transcript and hashing libraries were processed and sequenced as described above.

Raw sequencing files obtained from scRNA-seq of the perturbed cells and controls were processed using the EPI2ME wf-single-cell pipeline v3.1.0 (1) with the human GRCh38 reference genome. The processed gene expression matrices were combined into a single AnnData object containing cells as rows (obs) and genes as columns (var) (59).

Respective perturbation conditions (siPRRX1, siPRRX1/mmMYOD1, mmMYOD1) were assigned to individual cell barcodes using a custom demultiplexing pipeline. We add the *assigned\_condition* column to AnnData.obs.

### 3.3 TIME SERIES MICROSCOPY

The Zeiss Celldiscoverer 7 (CD7) live-cell imaging system was used to capture time-lapse images over the course of perturbation experiments. Oblique contrast and fluorescence microscopy was performed with a Plan-Apochromat 20x/0.7 objective and 0.5x tube lens. Images were taken using an AxioCam 506 with 14 bit resolution. Cells were imaged at 37°C in 5% CO<sub>2</sub> in imaging media. Images were captured every 20 minutes over 72 hours. Raw CZI files were exported from the Zen Blue 3.0 software for image processing and downstream analyses.

### 3.4 DATA STORAGE

The HYPED dataset is made available on kaggle along with the software used to process the raw experimental data from the sequencers. Each modality of the dataset is structured separately. The raw CZI files exported from the microscope contain data with dimensions HSTCZMYX (table 2). The four imaging channels correspond to markers that can be used to delineate the cell cycle phase:

- Cy5 - cell nuclei
- MKate - G2/M phase
- tagGFP - G1 phase
- Oblique - oblique contrast

Each scene is associated with one of the four experimental conditions. We extract only the relevant scenes across time frames where each frame is split into mosaics of  $6 \times 5$ . These mosaic tiles are stitched together to get the full frame. Each frame was then converted into a four-mode numpy array (T, C, X, Y) with dimensions (202, 4, 2826, 3245). To make the dataset standardized and accessible, tensor data of each frame was embedded into a HDF5 (Hierarchical Data Format). Each HDF5 file consists of a single frame with the four channels at a particular image capture point of the experiment. The dataset is designed and to load frames across any range of time points into compute memory

The long read sequencing gene expression data is stored in the HDF5 AnnData format *.h5ad*, a commonly-used HDF5-based format with extensive support in Python and R. This format stores the measured gene expression in each single cell in a cell by gene matrix, with the rows and columns annotated to correspond to the genes and experimental conditions of the cells.

## 4 BENCHMARKING

**Sequencing.** We compare the ground truth results obtained from the experiment and the outputs from predictive models. We use two models to benchmark: the foundation model Geneformer (55) and the perturbation prediction model GEARS (46) and report on cosine similarity, Mean Squared Error (MSE), Mean Absolute Error (MAE) and Pearson correlation coefficient between the predicted and observed values (see table 3).

Geneformer is a Transformer-based Deep Learning model which utilizes unsupervised self-attention mechanism to generate Transcriptomic representations. This generates embeddings in a lower-dimensional space that can be fine-tuned for a variety of tasks.

We tokenize the transcriptome of each cell in the data set with rank value encoding. For the in-silico perturbation, we set the expression level of *MYOD1* to a value greater than the highest expressed gene, which places it at the top of the ranked value encoding, simulating the activation of *MYOD1* within the cell. Similarly, we simulate suppression of *PRRX1* by setting the expression level to zero before rank value encoding. The tokenized data is then run through Geneformer generating embeddings. This embedded data is compared with those of experimentally perturbed cells, and the metrics are reported.

GEARS is a Deep Learning perturbation prediction model which can predict gene expression vectors under certain perturbation conditions. The GEARS model takes in a combinations of perturbed genes as input and predicts output gene-expression vectors. The model architecture supports splitting the dataset based on different perturbation combinations. The model is first fine tuned with existing perturbation data (39) which was split into

- Train - 139 combinations
- Test - 107 combinations
- Validation - 31 combinations

This pretrained model is further fine tuned with the control group data from our experiment as 3 new conditions. This is the suggested method from the authors of GEARS to integrate new datasets into the model.

We then performed in silico prediction for all three conditions using GEARS. This generates the predicted gene expression vector for the top 500 genes which we compare with the average expression of our held out conditional cells, and calculate different metrics.

---

378 **Algorithm 1** Geneformer Validation Experiment

---

379 **Require:** perturbed\_group, control\_group, perturbation\_genes

380 **Ensure:** cos\_sim

381 1: // Initialize lists to store control and perturbed group embeddings

382 2:  $E_c, E_p \leftarrow [], []$

383 3:

384 4: // Compute Geneformer embeddings of perturbed data

385 5: **for** each cell  $x$  in perturbed\_group **do**

386 6:    $e \leftarrow \text{Geneformer}(x)$

387 7:    $E_p \leftarrow [E_p \ e]$

388 8: **end for**

389 9:

390 10: // Compute Geneformer embeddings with in silico perturbations

391 11: **for** each cell  $x$  in control\_group **do**

392 12:   **if** MYOD  $\in$  perturbation\_genes **then**

393 12:      $x[\text{MYOD}] \leftarrow \max(x)$

394 13:   **end if**

395 14:   **if** PRRX1  $\in$  perturbation\_genes **then**

396 14:      $x[\text{PRRX1}] \leftarrow 0$

397 15:   **end if**

398 16:    $e \leftarrow \text{Geneformer}(x)$

399 17:    $E_c \leftarrow [E_c \ e]$

400 18: **end for**

401 19:

402 20: // Determine the average embedding of each experimental group

403 21:  $e_p \leftarrow \text{mean}(E_p)$

404 22:  $e_c \leftarrow \text{mean}(E_c)$

405 23:

406 24: **Return** CosineSimilarity( $u, v$ ) =0

---

407 **Table 3: Benchmarking Perturbation Models Against Experimental Data.** Single-cell control  
 408 data were computationally perturbed using Geneformer and GEARS, and the resulting profiles were  
 409 compared to experimentally perturbed counterparts.

Model	Perturbation	Cell Cycle	Cosine Sim.	MSE	MAE	Pearson
Geneformer	+ MYOD1	G1	0.9808	0.0148	0.0902	0.9807
		S	0.9776	0.0172	0.0988	0.9775
		G2/M	0.9715	0.0221	0.1119	0.9714
	- PRRX1	G1	0.9913	0.0068	0.0605	0.9913
		S	0.9900	0.0078	0.0660	0.9900
		G2/M	0.9900	0.0078	0.0680	0.9900
	- PRRX1 + MYOD1	G1	0.9902	0.0076	0.0645	0.9902
		S	0.9887	0.0086	0.0697	0.9887
		G2/M	0.9851	0.0115	0.0810	0.9851
GEARS	- PRRX1	All Phases	0.9518	0.0502	0.1079	0.9358
	+ MYOD1	All Phases	0.9741	0.0129	0.0622	0.9625
	- PRRX1 + MYOD1	All Phases	0.9772	0.0116	0.0598	0.9676

427

428

429 The high cosine similarities observed for both the GEARS and Geneformer predicted gene expression  
 430 indicate that the trained machine learning models reflect experimental behaviour. HYPED dataset  
 431 provides a strong empirical ground truth for training and evaluating future foundation models that  
 can leverage both modalities.

Table 4: **Image Benchmarking.** Image Quality Metrics Averaged Across Timepoints.

Model	MSE	PSNR (dB)	SSIM	Calculated SNR	Denoised SNR
FM2S	0.000146	86.480553	0.877678	14.082471	15.051044

**Imaging.** To evaluate the performance of noise removal techniques, we performed a benchmarking task on our imaging dataset. The extraction of regions of interest was performed through a standard pipeline using the following steps:

1. **Clipping and normalization:** For each channel, intensities above 98% were clipped and scaled to [0,1] range.
2. **Channel stacking:** mKate and TagGFP channels were stacked on each other to get the G1-S phase transitions.
3. **Gaussian filtering and Otsu thresholding:** Used to separate foreground cells from background.
4. **Watershed Algorithm:** Was used to differentiate individual cell boundaries.
5. **Segment mapping:** Segmented regions were converted to binary masks
6. **Tracking:** Bayesian tracking was used to assign object identities over time.

These cells are then labeled as foreground, and everything else is considered noise and the Signal-to-Noise Ratio is calculated. We then run the images through a machine learning algorithm and PSNR and SSIM were calculated. FM2S (35) is a Deep Neural Network based Fluorescent Microscopy Imaging denoiser. We Run this model on different conditions of our dataset and report the average metrics (see table 4).

## 5 DISCUSSION

Here, we introduced the HYPED benchmark dataset for cell reprogramming, designed using state-of-the-art perturbation techniques. Unlike widely used lentiviral, CRISPR, and Perturb-seq datasets which can have unforeseen off-target effects, HYPED offers the first multimodal perturbation dataset generated using transient RNA application for improved learning and prediction of perturbation dynamics.

Moreover, HYPED’s live cell imaging provides one of the first perturbation datasets to capture cell cycle dynamics with perturbations. It has long been recognized that cell perturbations are transient processes, yet many current models including GEARS, Geneformer, scGPT, and others do not account for or model the dynamics. This dataset provides an improved opportunity to train models that account for perturbation and biological dynamics at a higher temporal resolution and under current experimental conditions than previously possible.

To advance cell reprogramming and related biological problems, it is essential that machine learning models are aligned with the capabilities and limitations of contemporary experiments by training on modern experimental data modalities. Many existing approaches are trained on idealized or outdated datasets that fail to reflect the transient, dynamic nature of real biological systems. By employing modified mmRNA and siRNA and including live-cell imaging with long-read single-cell transcriptomics, the HYPED dataset serves as a resource for researchers to develop models that are not only more predictive, but also more actionable in laboratory and clinical settings.

## 6 LIMITATIONS

We identified two main limitations in our study. First, the experiment conducted includes the effects of three unique transcription factor perturbations captured on skin fibroblasts. Although this describes the activity of a small subset, the vast combinatorial space of transcription factors and cell types remains largely unexplored. Second, we considered a limited number of tasks for each modality, but there are a variety of challenges that can be explored with our dataset.

486 7 ETHICAL CONSIDERATIONS  
487

488 All resources provided as part of this cell reprogramming study are strictly for research purposes only  
489 and should not be used in clinical settings and diagnostic procedures. No sensitive information is  
490 included in the dataset. With the aforementioned restrictions, we have not identified any potential  
491 adverse impacts from the HYPED dataset.

492  
493 ACKNOWLEDGMENTS  
494

495 We would like to thank all Rajapakse lab members for helpful and inspiring discussions. This  
496 work was supported by the Defense Advanced Research Projects Agency (DARPA) award number  
497 [HR00112490472 to I.R.], the Air Force Office of Scientific Research (AFOSR) award number  
498 [FA9550-22-1-0215 to I.R.], support from NVIDIA [to I.R.], and support from the National Institute  
499 of General Medical Sciences (NIGMS) award number [GM150581 to J.P.].

500  
501 REFERENCES  
502

- 503 [1] epi2me-labs/wf-single-cell. original-date: 2022-08-19T09:19:03Z.  
504  
505 [2] nanoporetech/dorado. original-date: 2022-05-17T23:12:13Z.  
506  
507 [3] Aurelio Balsalobre and Jacques Drouin. Pioneer factors as master regulators of the epigenome  
508 and cell fate. *Nature Reviews Molecular Cell Biology*, 23(7):449–464, 2022.  
509  
510 [4] Brent A. Bidy, Wenjun Kong, Kenji Kamimoto, Chuner Guo, Sarah E. Waye, Tao Sun, and  
511 Samantha A. Morris. Single-cell mapping of lineage and identity in direct reprogramming.  
512 *Nature*, 564(7735):219–224, December 2018. Publisher: Nature Publishing Group.  
513  
514 [5] Anna K Blakney, Paul F McKay, and Robin J Shattock. Delivery of mrna-based vaccines and  
515 therapeutics with lipid nanoparticles. *Biochemical Society Transactions*, 47(4):1209–1218,  
516 2019.  
517  
518 [6] Simone Bruno, Thorsten M Schlaeger, and Domitilla Del Vecchio. Epigenetic oct4 regulatory  
519 network: stochastic analysis of cellular reprogramming. *npj Systems Biology and Applications*,  
520 10(1):3, 2024.  
521  
522 [7] Daniel Burkhardt, Andrew Benz, Robrecht Cannoodt, Mauricio Cortes, Scott Gigante, Christo-  
523 pher Lance, Richard Lieberman, Malte Luecken, and Angela Pisco. Single-cell pertur-  
524 bation prediction: generalizing experimental interventions to unseen contexts. <https://neurips.cc/virtual/2023/competition/66586>, 2023. NeurIPS 2023 Compe-  
525 tition Track.  
526  
527 [8] Patrick Cahan, Hu Li, Samantha A. Morris, Edroaldo Lummertz da Rocha, George Q. Daley,  
528 and James J. Collins. CellNet: Network Biology Applied to Stem Cell Engineering. *Cell*,  
529 158(4):903–915, August 2014.  
530  
531 [9] Haotian Cui, Chloe Wang, Hassaan Maan, Kuan Pang, Fengning Luo, Nan Duan, and Bo Wang.  
532 scgpt: toward building a foundation model for single-cell multi-omics using generative ai.  
533 *Nature Methods*, 21(8):1470–1480, 2024.  
534  
535 [10] Atray Dixit, Oren Parnas, Biyu Li, Jenny Chen, Charles P Fulco, Livnat Jerby-Arnon, Ne-  
536 manja D Marjanovic, Danielle Dionne, Tyler Burks, Raktima Raychowdhury, et al. Perturb-seq:  
537 dissecting molecular circuits with scalable single-cell rna profiling of pooled genetic screens.  
538 *cell*, 167(7):1853–1866, 2016.  
539  
540 [11] Christoffer Edlund, Timothy R Jackson, Nabeel Khalid, Nicola Bevan, Timothy Dale, Andreas  
541 Dengel, Sheraz Ahmed, Johan Trygg, and Rickard Sjögren. Livecell—a large-scale dataset for  
542 label-free live cell segmentation. *Nature methods*, 18(9):1038–1045, 2021.  
543  
544 [12] Ryohei Eguchi, Momoko Hamano, Michio Iwata, Toru Nakamura, Shinya Oki, and Yoshihiro  
545 Yamanishi. Transdire: data-driven direct reprogramming by a pioneer factor-guided trans-omics  
546 approach. *Bioinformatics*, 38(10):2839–2846, 2022.

- 540 [13] Mirko Francesconi, Bruno Di Stefano, Clara Berenguer, Luisa de Andrés-Aguayo, Marcos  
541 Plana-Carmona, Maria Mendez-Lago, Amy Guillaumet-Adkins, Gustavo Rodriguez-Esteban,  
542 Marta Gut, Ivo G Gut, Holger Heyn, Ben Lehner, and Thomas Graf. Single cell RNA-seq  
543 identifies the origins of heterogeneity in efficient cell transdifferentiation and reprogramming.  
544 *eLife*, 8:e41627, March 2019. Publisher: eLife Sciences Publications, Ltd.
- 545 [14] Oscar Franzén, Li-Ming Gan, and Johan L M Björkegren. PanglaoDB: a web server for  
546 exploration of mouse and human single-cell RNA sequencing data. *Database*, 2019:baz046,  
547 January 2019.
- 548 [15] Xi Fu, Shentong Mo, Alejandro Buendia, Anouchka P. Laurent, Anqi Shao, Maria del Mar  
549 Alvarez-Torres, Tianji Yu, Jimin Tan, Jiayu Su, Romella Sagatelian, Adolfo A. Ferrando, Alberto  
550 Ciccia, Yanyan Lan, David M. Owens, Teresa Palomero, Eric P. Xing, and Raul Rabadan. A  
551 foundation model of transcription across human cell types. *Nature*, 637(8047):965–973, January  
552 2025. Publisher: Nature Publishing Group.
- 553 [16] Andreia M. Gomes, Ilia Kurochkin, Betty Chang, Michael Daniel, Kenneth Law, Namita  
554 Satija, Alexander Lachmann, Zichen Wang, Lino Ferreira, Avi Ma’ayan, Benjamin K. Chen,  
555 Dmitri Papatsenko, Ihor R. Lemischka, Kateri A. Moore, and Carlos-Filipe Pereira. Cooper-  
556 ative Transcription Factor Induction Mediates Hemogenic Reprogramming. *Cell Reports*,  
557 25(10):2821–2835.e7, December 2018.
- 558 [17] Gavin D Grant, Katarzyna M Kedziora, Juanita C Limas, Jeanette Gowen Cook, and Jeremy E  
559 Purvis. Accurate delineation of cell cycle phase transitions in living cells with pip-fucci. *Cell*  
560 *cycle*, 17(21-22):2496–2516, 2018.
- 561 [18] Tessa Durakis Green, Stefan Peidli, Ciyue Shen, Torsten Gross, Joseph Min, Samuele Garda,  
562 Jake P Taylor-King, Debora Susan Marks, Augustin Luna, Nils Blüthgen, et al. scperturb:  
563 Information resource for harmonized single-cell perturbation data. In *NeurIPS 2022 Workshop*  
564 *on Learning Meaningful Representations of Life*, 2022.
- 565 [19] Chuner Guo and Samantha A Morris. Engineering cell identity: establishing new gene regulatory  
566 and chromatin landscapes. *Curr. Opin. Genet. Dev.*, 46:50–57, October 2017.
- 567 [20] Lin Guo, Lihui Lin, Xiaoshan Wang, Mingwei Gao, Shangtao Cao, Yuanbang Mai, Fang Wu,  
568 Junqi Kuang, He Liu, Jiaqi Yang, Shilong Chu, Hong Song, Dongwei Li, Yujian Liu, Kaixin Wu,  
569 Jiadong Liu, Jinyong Wang, Guangjin Pan, Andrew P. Hutchins, Jing Liu, Duanqing Pei, and  
570 Jiekai Chen. Resolving Cell Fate Decisions during Somatic Cell Reprogramming by Single-Cell  
571 RNA-Seq. *Molecular Cell*, 73(4):815–829.e7, February 2019. Publisher: Elsevier.
- 572 [21] Zhisong He, Leander Dony, Jonas Simon Fleck, Artur Szałata, Katelyn X. Li, Irena Slišković,  
573 Hsiu-Chuan Lin, Malgorzata Santel, Alexander Atamian, Giorgia Quadrato, Jieran Sun, Sergiu P.  
574 Paşca, J. Gray Camp, Fabian J. Theis, and Barbara Treutlein. An integrated transcriptomic  
575 cell atlas of human neural organoids. *Nature*, 635(8039):690–698, November 2024. Publisher:  
576 Nature Publishing Group.
- 577 [22] Kenichi Horisawa and Atsushi Suzuki. Direct cell-fate conversion of somatic cells: Toward  
578 regenerative medicine and industries. *Proceedings of the Japan Academy, Series B*, 96(4):131–  
579 158, 2020.
- 580 [23] Masahito Inagaki. Cell reprogramming and differentiation utilizing messenger rna for regenera-  
581 tive medicine. *Journal of Developmental Biology*, 12(1):1, 2023.
- 582 [24] Miten Jain, Robin Abu-Shumays, Hugh E Olsen, and Mark Akeson. Advances in nanopore  
583 direct rna sequencing. *Nature methods*, 19(10):1160–1164, 2022.
- 584 [25] Miten Jain, Sergey Koren, Karen H Miga, Josh Quick, Arthur C Rand, Thomas A Sasani, John R  
585 Tyson, Andrew D Beggs, Alexander T Dilthey, Ian T Fiddes, et al. Nanopore sequencing and  
586 assembly of a human genome with ultra-long reads. *Nature biotechnology*, 36(4):338–345,  
587 2018.
- 588 [26] Miten Jain, Hugh E Olsen, Benedict Paten, and Mark Akeson. The oxford nanopore minion:  
589 delivery of nanopore sequencing to the genomics community. *Genome biology*, 17:1–11, 2016.

- 594 [27] Yuge Ji, Mohammad Lotfollahi, F Alexander Wolf, and Fabian J Theis. Machine learning for  
595 perturbational single-cell omics. *Cell Systems*, 12(6):522–537, 2021.  
596
- 597 [28] Julia Joung, Sai Ma, Tristan Tay, Kathryn R. Geiger-Schuller, Paul C. Kirchgatterer, Vanessa K.  
598 Verdine, Baolin Guo, Mario A. Arias-Garcia, William E. Allen, Ankita Singh, Olena Kuksenko,  
599 Omar O. Abudayyeh, Jonathan S. Gootenberg, Zhanyan Fu, Rhiannon K. Macrae, Jason D.  
600 Buenrostro, Aviv Regev, and Feng Zhang. A transcription factor atlas of directed differentiation.  
601 *Cell*, 186(1):209–229.e26, January 2023.
- 602 [29] Uma S Kamaraj et al. Computational methods for direct cell conversion. *Cell Cycle*,  
603 15(24):3343–3354, 2016.  
604
- 605 [30] Rosemary Kanasty, J. Robert Dorkin, Arturo Vegas, and Daniel Anderson. Delivery materials  
606 for sirna therapeutics. *Nature Materials*, 12(11):967–977, 2013.
- 607 [31] Riya Keshri, Damien Detraux, Ashish Phal, Clara McCurdy, Samriddhi Jhajharia, Tung Ching  
608 Chan, Julie Mathieu, and Hannele Ruohola-Baker. Next-generation direct reprogramming.  
609 *Frontiers in Cell and Developmental Biology*, 12:1343106, 2024.
- 610 [32] Christopher Lance, Malte D Luecken, Daniel B Burkhardt, Robrecht Cannoodt, Pia Rauten-  
611 strauch, Anna Laddach, Aidyn Ubingazhibov, Zhi-Jie Cao, Kaiwen Deng, Sumeer Khan, et al.  
612 Multimodal single cell data integration challenge: results and lessons learned. *BioRxiv*, pages  
613 2022–04, 2022.  
614
- 615 [33] Kwanyoung Lee, Hyungjo Byun, and Hyunjung Shim. Cell segmentation in multi-modality  
616 high-resolution microscopy images with cellpose. In *Competitions in Neural Information  
617 Processing Systems*, pages 1–11. PMLR, 2023.
- 618 [34] Hanqin Li, Houbo Jiang, Xinzhen Yin, Jonathan E Bard, Baorong Zhang, and Jian Feng.  
619 Attenuation of *prrx2* and *hey2* enables efficient conversion of adult human skin fibroblasts to  
620 neurons. *Biochemical and biophysical research communications*, 516(3):765–769, 2019.  
621
- 622 [35] Jizhihui Liu, Qixun Teng, Qing Ma, and Junjun Jiang. Fm2s: Towards spatially-correlated noise  
623 modeling in zero-shot fluorescence microscopy image denoising, 2025.
- 624 [36] Samantha A Morris. Direct lineage reprogramming via pioneer factors; a detour through  
625 developmental gene regulatory networks. *Development*, 143(15):2696–2705, August 2016.  
626
- 627 [37] Samantha A. Morris, Patrick Cahan, Hu Li, Anna M. Zhao, Adrianna K. San Roman, Ramesh A.  
628 Shivdasani, James J. Collins, and George Q. Daley. Dissecting Engineered Cell Types and  
629 Enhancing Cell Fate Conversion via CellNet. *Cell*, 158(4):889–902, August 2014.
- 630 [38] Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Michael Wornow, Callum Birch-  
631 Sykes, Stefano Massaroli, Aman Patel, Clayton Rabideau, Yoshua Bengio, et al. Hyenadna:  
632 Long-range genomic sequence modeling at single nucleotide resolution. *Advances in neural  
633 information processing systems*, 36:43177–43201, 2023.  
634
- 635 [39] Thomas M Norman, Max A Horlbeck, Joseph M Replogle, Alex Y Ge, Albert Xu, Marco Jost,  
636 Luke A Gilbert, and Jonathan S Weissman. Exploring genetic interaction manifolds constructed  
637 from rich single-cell phenotypes. *Science*, 365(6455):786–793, 2019.
- 638 [40] Paul Nurse, Yoshio Masui, and Leland Hartwell. Understanding the cell cycle. *Nature medicine*,  
639 4(10):1103–1106, 1998.  
640
- 641 [41] CZI Cell Science Program, Shibli Abdulla, Brian Aevertmann, Pedro Assis, Seve Badajoz,  
642 Sidney M Bell, Emanuele Bezzi, Batuhan Cakir, Jim Chaffer, Signe Chambers, et al. Cz  
643 cellxgene discover: a single-cell data platform for scalable exploration, analysis and modeling  
644 of aggregated data. *Nucleic Acids Research*, 53(D1):D886–D900, 2025.
- 645 [42] Owen JL Rackham, Jaber Firas, Hai Fang, Matt E Oates, Melissa L Holmes, Anja S Knaupp,  
646 FANTOM Consortium, Harukazu Suzuki, Christian M Nefzger, Carsten O Daub, et al. A  
647 predictive computational framework for direct reprogramming between human cell types.  
*Nature genetics*, 48(3):331–335, 2016.

- 648 [43] Scott Ronquist, Geoff Patterson, Lindsey A Muir, Stephen Lindsly, Haiming Chen, Markus  
649 Brown, Max S Wicha, Anthony Bloch, Roger Brockett, and Indika Rajapakse. Algorithm for  
650 cellular reprogramming. *Proceedings of the National Academy of Sciences*, 114(45):11832–  
651 11837, 2017.
- 652 [44] Jennifer E Rood, Anna Hupalowska, and Aviv Regev. Toward a foundation model of causal cell  
653 and tissue biology with a perturbation cell and tissue atlas. *Cell*, 187(17):4520–4545, 2024.
- 654 [45] Yusuf Roohani, Kexin Huang, and Jure Leskovec. Predicting transcriptional outcomes of novel  
655 multigene perturbations with gears. *Nature Biotechnology*, 42(6):927–935, 2024.
- 656 [46] Yusuf Roohani, Kexin Huang, and Jure Leskovec. Predicting transcriptional outcomes of  
657 novel multigene perturbations with GEARS. *Nature Biotechnology*, 42(6):927–935, June 2024.  
658 Publisher: Nature Publishing Group.
- 659 [47] Rita Silvério-Alves, Ilia Kurochkin, Anna Rydström, Camila Vazquez Echegaray, Jakob Haider,  
660 Matthew Nicholls, Christina Rode, Louise Thelaus, Aida Yifter Lindgren, Alexandra Gabriela  
661 Ferreira, Rafael Brandão, Jonas Larsson, Marella F. T. R. de Bruijn, Javier Martin-Gonzalez, and  
662 Carlos-Filipe Pereira. GATA2 mitotic bookmarking is required for definitive haematopoiesis.  
663 *Nature Communications*, 14(1):4645, August 2023. Publisher: Nature Publishing Group.
- 664 [48] Amar M Singh, Robert Trost, Benjamin Boward, and Stephen Dalton. Utilizing fucci reporters  
665 to understand pluripotent stem cell biology. *Methods*, 101:4–10, 2016.
- 666 [49] Cooper Stansbury, Jillian Cwycyshyn, Joshua Pickard, Walter Meixner, Indika Rajapakse, and  
667 Lindsey A. Muir. Data-guided direct reprogramming of human fibroblasts into the hematopoi-  
668 etic lineage, August 2024. `tex.ids= stansburyDataguidedDirectReprogramming2024a` pages:  
669 2024.08.26.609589 section: New Results.
- 670 [50] Clara Steichen, Eléonor Luce, Jérôme Maluenda, Lucie Tosca, Inmaculada Moreno-Gimeno,  
671 Christophe Desterke, Noushin Dianat, Sylvie Goulinet-Mainot, Sarah Awan-Toor, Deborah  
672 Burks, et al. Messenger rna-versus retrovirus-based induced pluripotent stem cell reprogram-  
673 ming strategies: analysis of genomic integrity. *Stem Cells Translational Medicine*, 3(6):686–691,  
674 2014.
- 675 [51] Artur Szalata, Karin Hrovatin, Sören Becker, Alejandro Tejada-Lapuerta, Haotian Cui, Bo Wang,  
676 and Fabian J. Theis. Transformers in single-cell omics: a review and new perspectives. *Nature*  
677 *Methods*, 21(8):1430–1443, August 2024. Publisher: Nature Publishing Group.
- 678 [52] John A. Tadross, Lukas Steuernagel, Georgina K. C. Dowsett, Katherine A. Kentistou, Sofia  
679 Lundh, Marta Porniece-Kumar, Paul Klemm, Kara Rainbow, Henning Hvid, Katarzyna Kania,  
680 Joseph Poley-Wolf, Lotte Bjerre-Knudsen, Charles Pyke, John R. B. Perry, Brian Y. H. Lam,  
681 Jens C. Brüning, and Giles S. H. Yeo. Human HYPOMAP: A comprehensive spatio-cellular  
682 map of the human hypothalamus, September 2023. Pages: 2023.09.15.557967 Section: New  
683 Results.
- 684 [53] Kazutoshi Takahashi and Shinya Yamanaka. Induction of pluripotent stem cells from mouse  
685 embryonic and adult fibroblast cultures by defined factors. *cell*, 126(4):663–676, 2006.
- 686 [54] THE TABULA SAPIENS CONSORTIUM. The Tabula Sapiens: A multiple-organ, single-cell  
687 transcriptomic atlas of humans. *Science*, 376(6594):eabl4896, May 2022. Publisher: American  
688 Association for the Advancement of Science.
- 689 [55] Christina V Theodoris, Ling Xiao, Anant Chopra, Mark D Chaffin, Zeina R Al Sayed, Matthew C  
690 Hill, Helene Mantineo, Elizabeth M Brydon, Zexian Zeng, X Shirley Liu, et al. Transfer learning  
691 enables predictions in network biology. *Nature*, 618(7965):616–624, 2023.
- 692 [56] Yasuhiro Tomaru, Ryota Hasegawa, Takahiro Suzuki, Taiji Sato, Atsutaka Kubosaki, Masanori  
693 Suzuki, Hideya Kawaji, Alistair RR Forrest, Yoshihide Hayashizaki, FANTOM Consortium,  
694 et al. A transient disruption of fibroblastic transcriptional regulatory network facilitates trans-  
695 differentiation. *Nucleic acids research*, 42(14):8905–8913, 2014.

- 702 [57] Connor A Tsuchida, Kevin M Wasko, Jennifer R Hamilton, and Jennifer A Doudna. Targeted  
703 nonviral delivery of genome editors in vivo. *Proceedings of the National Academy of Sciences*,  
704 121(11):e2307796121, 2024.
- 705 [58] Erwin L Van Dijk, Yan Jaszczyszyn, Delphine Naquin, and Claude Thermes. The third  
706 revolution in sequencing technology. *Trends in Genetics*, 34(9):666–681, 2018.
- 707 [59] Isaac Virshup, Sergei Rybakov, Fabian J. Theis, Philipp Angerer, and F. Alexander Wolf. ann-  
708 data: Access and store annotated data matrices. *Journal of Open Source Software*, 9(101):4371,  
709 2024.
- 710 [60] Aline Yen Ling Wang. Application of modified mrna in somatic reprogramming to pluripotency  
711 and directed conversion of cell fate. *International journal of molecular sciences*, 22(15):8148,  
712 2021.
- 713 [61] Haofei Wang, Yuchen Yang, Jiandong Liu, and Li Qian. Direct cell reprogramming: approaches,  
714 mechanisms and progress. *Nature Reviews Molecular Cell Biology*, 22(6):410–424, June 2021.  
715 Number: 6 Publisher: Nature Publishing Group.
- 716 [62] Luigi Warren and Cory Lin. mrna-based genetic reprogramming. *Molecular Therapy*, 27(4):729–  
717 734, 2019.
- 718 [63] Luigi Warren, Philip D Manos, Tim Ahfeldt, Yui-Han Loh, Hu Li, Frank Lau, Wataru Ebina,  
719 Pankaj K Mandal, Zachary D Smith, Alexander Meissner, et al. Highly efficient reprogramming  
720 to pluripotency and directed differentiation of human cells with synthetic modified mrna. *Cell*  
721 *stem cell*, 7(5):618–630, 2010.
- 722 [64] Zhiting Wei, Duanmiao Si, Bin Duan, Yicheng Gao, Qian Yu, Zhenbo Zhang, Ling Guo, and  
723 Qi Liu. Perturbbase: a comprehensive database for single-cell perturbation data analysis and  
724 visualization. *Nucleic Acids Research*, 53(D1):D1099–D1111, 2025.
- 725 [65] Harold Weintraub, Robert Davis, Stephen Tapscott, Matthew Thayer, Michael Krause, Robert  
726 Benezra, T Keith Blackwell, David Turner, Ralph Rupp, Stanley Hollenberg, et al. The myod  
727 gene family: nodal point during specification of the muscle cell lineage. *Science*, 251(4995):761–  
728 766, 1991.
- 729 [66] Harold Weintraub, Stephen J Tapscott, Robert L Davis, Mathew J Thayer, Mohammed A Adam,  
730 Andrew B Lassar, and A Dusty Miller. Activation of muscle-specific genes in pigment, nerve,  
731 fat, liver, and fibroblast cell lines by forced expression of myod. *Proceedings of the National*  
732 *Academy of Sciences*, 86(14):5434–5438, 1989.
- 733 [67] Thomas P. Wytock and Adilson E. Motter. Cell reprogramming design by transfer learning  
734 of functional transcriptional networks. *Proceedings of the National Academy of Sciences*,  
735 121(11):e2312942121, March 2024. `tex.ids= wytockCellReprogrammingDesign2024a` pub-  
736 lisher: Proceedings of the National Academy of Sciences.
- 737 [68] Yun Xiao, Yonghui Gong, Yanling Lv, Yujia Lan, Jing Hu, Feng Li, Jinyuan Xu, Jing Bai, Yulan  
738 Deng, Ling Liu, et al. Gene perturbation atlas (gpa): a single-gene perturbation repository  
739 for characterizing functional mechanisms of coding and non-coding genes. *Scientific reports*,  
740 5(1):10889, 2015.
- 741 [69] Zhiyuan Xie, Ilya Sokolov, Maria Osmala, Xue Yue, Grace Bower, J. Patrick Pett, Yinan Chen,  
742 Kai Wang, Ayse Derya Cavga, Alexander Popov, Sarah A. Teichmann, Ekaterina Morgunova,  
743 Evgeny Z. Kvon, Yimeng Yin, and Jussi Taipale. DNA-guided transcription factor interactions  
744 extend human gene regulatory code. *Nature*, pages 1–10, April 2025. Publisher: Nature  
745 Publishing Group.
- 746 [70] Qiao Rui Xing, Chadi El Farran, Pradeep Gautam, Yu Song Chuah, Tushar Warriar, Cheng-  
747 Xu Delon Toh, Nam-Young Kang, Shigeki Sugii, Young-Tae Chang, Jian Xu, James J. Collins,  
748 George Q. Daley, Hu Li, Li-Feng Zhang, and Yui-Han Loh. Diversification of reprogramming  
749 trajectories revealed by parallel single-cell transcriptome and chromatin accessibility sequencing.  
750 *Science Advances*, 6(37):eaba1190, September 2020. Publisher: American Association for the  
751 Advancement of Science.

756 [71] Shinya Yamanaka. Elite and stochastic models for induced pluripotent stem cell generation.  
757 *Nature*, 460(7251):49–52, 2009.  
758

759 [72] Bo Yuan, Ciyue Shen, Augustin Luna, Anil Korkut, Debora S Marks, John Ingraham, and Chris  
760 Sander. Cellbox: interpretable machine learning for perturbation biology with application to the  
761 design of cancer combination therapy. *Cell systems*, 12(2):128–140, 2021.  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809