Universal Cell Embeddings: A Foundation Model for Cell Biology

3	Yanay Rosen ^{1,*} , Yusuf Roohani ^{2,*} , Ayush Agarwal ¹ , Leon Samotorčan ¹ ,
4	Tabula Sapiens Consortium ³ , Stephen R. Quake ^{4,5,6,†} , Jure Leskovec ^{1,†}
5	¹ Department of Computer Science, Stanford University, Stanford, CA, USA
6	² Department of Biomedical Data Science, Stanford University, Stanford, CA, USA
7	³ Chan Zuckerberg BioHub, San Francisco, CA, USA
8	⁴ Department of Bioengineering, Stanford University, Stanford, CA, USA
9	⁵ Department of Applied Physics, Stanford University, Stanford, CA, USA
10	⁶ Chan Zuckerberg Initiative, Redwood City, CA, USA
11	[†] Corresponding author. Email: jure@cs.stanford.edu, quake@stanford.edu
12	*These authors contributed equally

13 Abstract

1

2

Developing a universal representation of cells which encompasses the tremendous molecular 14 diversity of cell types within the human body and more generally, across species, would be 15 transformative for cell biology. Recent work using single-cell transcriptomic approaches to 16 create molecular definitions of cell types in the form of cell atlases has provided the necessary 17 data for such an endeavor. Here, we present the Universal Cell Embedding (UCE) founda-18 tion model. UCE was trained on a corpus of cell atlas data from human and other species 19 in a completely self-supervised way without any data annotations. UCE offers a unified bio-20 logical latent space that can represent any cell, regardless of tissue or species. This universal 21 cell embedding captures important biological variation despite the presence of experimental 22 noise across diverse datasets. An important aspect of UCE's universality is that any new cell 23 from any organism can be mapped to this embedding space with no additional data labeling, 24

model training or fine-tuning. We applied UCE to create the Integrated Mega-scale Atlas, 25 embedding 36 million cells, with more than 1,000 uniquely named cell types, from hundreds 26 of experiments, dozens of tissues and eight species. We uncovered new insights about the or-27 ganization of cell types and tissues within this universal cell embedding space, and leveraged 28 it to infer function of newly discovered cell types. UCE's embedding space exhibits emergent 29 behavior, uncovering new biology that it was never explicitly trained for, such as identifying 30 developmental lineages and embedding data from novel species not included in the train-31 ing set. Overall, by enabling a universal representation for every cell state and type, UCE 32 provides a valuable tool for analysis, annotation and hypothesis generation as the scale and 33 diversity of single cell datasets continues to grow. 34

35 Introduction

Cells are the fundamental unit of life and biologists have long conceptualized cells as members of different universal landscapes [1–4]. A notable example of this is the Waddington landscape, which presents a theoretical framework for the developmental lineages of cells as they transition from pluripotent stages such as stem cells to more terminally differentiated end points [5]. Broadly, the field of cell biology has sought to map the range of phenotypes that cells might exhibit, their interrelationships, and the shifts between these states during development and disease [6–10].

The substantial growth in the size of single-cell RNA sequencing (scRNA-seq) datasets 42 presents a fresh opportunity to revisit these questions. Detailed transcriptomic snapshots of cells 43 are now widely available from a range of timepoints, tissues, donors, and species [11–13]. These 44 rich, high-dimensional states are typically distilled into low-dimensional vectors or embeddings 45 to facilitate computational analysis [14, 15]. However, existing computational approaches strug-46 gle to jointly analyze these diverse datasets. The unified representations they produce are often 47 unable to extend to new datasets due to species-specific constraints in their construction or the 48 presence of dataset-specific artifacts (or batch effects) which can obscure the underlying biologi-49 cal signal [16, 17]. 50

Some computational methods for scRNA-seq data have managed to overcome these limitations, but at the cost of requiring model tuning for each new dataset, thus rendering the representations non-universal [15, 18, 19]. As a result, whenever a new experiment is performed and new data is collected, it requires dedicated, resource-intensive data labeling and model training to perform even the most standard analyses, such as clustering or annotation. This process is both time consuming and inefficient, and results in sub-optimal analyses based on small, limited and private datasets.

Recent advances in the field of artificial intelligence have enabled general-purpose founda-58 tion models (such as ChatGPT [20, 21], PaLM [22], Llama [23] and SAM [24]) that can learn 59 universal representations that are then applied to diverse downstream tasks and analyses. These 60 foundation models are not specifically trained for these downstream tasks, thus presenting clear 61 instances of emergent capabilities [25]. This foundation model strategy has also found valu-62 able applications in biological contexts such as learning representations of protein and DNA se-63 quences [26,27]. While some recent work has applied foundation model architectures to single-cell 64 genomics data, the unique characteristics of these datasets necessitate a specialized modeling ap-65 proach to fully realize their potential [28, 29]. Directly modeling gene expression as text in the 66

⁶⁷ form of a sequence of genes is both inefficient from a learning perspective and often relies on
 ⁶⁸ inaccurate biological assumptions.

Here, we present Universal Cell Embedding (UCE), a foundation model for single-cell gene 69 expression that is designed to address questions in cell and molecular biology. UCE is uniquely 70 able to generate representations of new single-cell gene expression datasets with no model fine-71 tuning or retraining while still remaining robust to dataset and batch-specific artifacts. Moreover, 72 it does so while requiring no cell type annotation and no input dataset preprocessing, such as 73 gene selection. UCE can be applied to any set of genes from any species, even if they aren't 74 homologs of genes seen during training. UCE learns a universal representation of cell biology that 75 is intrinsically meaningful and can extend insights beyond the data that has been experimentally 76 observed. The representations learned by UCE display an emergent organization of cell types that 77 is consistent with known biology. These cell embeddings can be used to accurately predict cell 78 types with no additional model retraining, showing improved performance in dataset integration 79 against existing atlas-scale integration methods. 80

UCE presents a novel approach to analyzing cell states. It enables the mapping of new 81 data into a universal embedding space, already populated with annotated reference states. This 82 strategy addresses issues such as noisy measurements that limit data alignment across different 83 experiments, and reduces reliance on small sets of marker genes to translate insights across studies 84 [30]. UCE empowers researchers to utilize existing models on new data without needing data 85 labeling or model retraining. This can foster novel cross-dataset discoveries and overcome the 86 limitations currently faced when working with small, isolated datasets. For instance, a cell type 87 classifier trained to predict specific immune cell types can be seamlessly applied to a completely 88 new dataset. Thus, UCE offers a versatile, efficient, and broadly applicable framework for the 89 analysis of cell states. 90

91 **Results**

A biologically-informed foundation model for single cell gene expression.

Integrating single-cell RNA sequencing (scRNA-seq) datasets is challenging for two primary reasons: scRNA-seq data does not always contain the same genes, or features, and those features are plagued by dataset-specific experimental artifacts or batch effects, which means models have to be built separately for each dataset. UCE overcomes these challenges by abstracting cells as ⁹⁷ 'bags of RNA' [31]. UCE (Fig. 1a) converts the RNA gene expression of a single cell into an ⁹⁸ expression weighted sample of its corresponding genes. Next, UCE represents the sample's genes ⁹⁹ by their protein products, using a large protein language model. This allows UCE to meaningfully ¹⁰⁰ represent any gene, from any species, regardless of whether the species had appeared in the training ¹⁰¹ data. Finally, after incorporating additional metadata about genes' chromosomal locations, this ¹⁰² representation is fed into a large transformer model [32]. UCE is able to map any cell, from any ¹⁰³ tissue, or any species, into one shared universal space, with no additional training.

In particular, UCE takes as an input (1) scRNA-seq count data and (2) the corresponding 104 protein embeddings, generated by a large protein language model, ESM2 [33], for the genes in the 105 dataset. The ESM2 protein language model takes amino acid sequences as an input and produces 106 a numerical representation called a protein embedding. Given the expression count data for a cell, 107 UCE takes a weighted and normalized sample, with replacement, of the cell's genes. This sample 108 can only contain genes which had non-zero expression, and can contain multiple copies of each 109 gene. These genes are then tokenized by converting them to the protein embedding representation 110 of the protein that they code for [34]. Genes belonging to the same chromosome are grouped 111 together by placing them in between special tokens and are then sorted by genomic location. A 112 special token representing the entire cell, the 'CLS' token, is appended to the beginning of the cell 113 representation [35]. This combined representation is passed into a transformer neural network. The 114 embedding of a cell is taken as the embedding of the CLS token at the final layer of the transformer 115 (Fig. 1a). 116

¹¹⁷ UCE is trained in a completely self-supervised manner, and thus does not make use of any ¹¹⁸ cell type or dataset-based annotations. In particular, during training, a random subset (20%) of ¹¹⁹ genes that were expressed are masked before sampling. These expressed genes are combined ¹²⁰ with a random subset of genes which had zero expression (non-expressed genes) to form a set of ¹²¹ query genes. Each of these query genes' protein embedding tokens is combined with the UCE ¹²² embedding of the cell they were generated from, and this joint embedding is passed into a fully ¹²³ connected neural network that predicts if that gene was expressed.

UCE is a 33 layer model consisting of over 650 million parameters. UCE was trained across more than 300 datasets that are largely collected from the CellXGene corpus [36] consisting of over 36 million cells, for 40 days across 24 A100 80GB GPUs (Methods, Extended Data Table

5

¹²⁷ 2, Supplementary Table 2). The model's weights and implementation are freely available and the
 ¹²⁸ model will be hosted as an openly available resource for the research community to run inference
 ¹²⁹ on new datasets.

130 UCE creates an Integrated Mega-scale Atlas (IMA) of 36 million cells.

We apply UCE to generate an Integrated Mega-scale Atlas (IMA) of 36 million cells sampled from diverse biological conditions, demonstrating the emergent organization of UCE cell representations (Fig. 1b). We find that cells within the UCE space naturally cluster by biological conditions like cell type, while mixing among experimental conditions like batch (Fig. 1b). Since UCE is trained in a self-supervised manner, this organization represents an emergent behavior of the model. The IMA contains numerous cell type alignments, across tissues and species.

To investigate the emergent organization of the IMA, we inspect how tissue residency can influence the state of cell types. Although macrophages found in different tissues are characterized by diverse transcriptional identities [37], they align closely in the UCE space (Extended Data Table 1). For the purpose of our analysis below, we first determine the central location of each cell type and tissue combination in the IMA space, by averaging the UCE embeddings of the cells from that combination, creating a tissue and cell type 'centroid'.

Cells in the IMA have been pre-labeled by their cell type. As these labels were never used for training the UCE model, we use them to validate the quality of the learned representation. For example, in the IMA, human macrophages are found in 73 different tissues and among these tissues, 72% (53) of tissue-specific macrophage centroids were embedded closest to a macrophage centroid from another tissue. Considering the 3 nearest centroids increases this percentage to 93% (Extended Data Table 1). Similar cross-tissue homogeneity can also be identified in other prolific cell types, like endothelial cells or neurons. This demonstrates that UCE, without any explicit
training or labels, identifies that macrophages have a unique cellular identity that is shared across
tissues. More broadly, it is an example of UCE's emergent organization that is consistent with
known biology even though not explicitly trained for.

153 UCE embeds new datasets without additional model training.

We evaluated the universality of UCE representations by directly mapping new datasets 154 which were not part of the training set into the embedding, without any additional training or 155 refinement of the UCE model. This is referred to as a 'zero-shot' capability, since the model was 156 never trained on any samples from the new dataset (Fig. 2a). While a variety of deep learning 157 models have been proposed for this task, we choose to compare the performance of UCE to other 158 self-supervised transformer-based methods. This is because they do not rely on cell type annota-159 tion, are trained on large datasets, have high model capacity, and can be run in a zero-shot setting. 160 In particular, we compare against Geneformer [28] and scGPT [29], both of which represent cells 161 using ordered lists of gene tokens. 162

We assess the performance of these methods on a completely new and yet unreleased dataset 163 (as of the publication of this manuscript), Tabula Sapiens v2, which contains diverse human data 164 from 581,430 cells, 27 tissues, 167 batches and 162 unique cell types. We use established metrics 165 for embedding quality that measure the conservation of cell type information and the correction 166 of batch effects (Methods). We compared several methods and found that UCE substantially out-167 performs the next best method Geneformer by 9.0% on overall score, 10.6% on biological conser-168 vation score, and 7.4% on batch correction score (Supplementary Table 1). To comprehensively 169 assess the value of these zero-shot embeddings, we also compare UCE to fine-tuned methods that 170

are conventionally used for this task. Notably, UCE even performs slightly better than non-zeroshot methods that require dataset-specific training: scVI [15] and scArches [18].

We also investigate the Tabula Sapiens v1 [11] (which was part of the training set) and v2 173 embeddings of each model visually by creating UMAP embeddings (Fig. 2b). UCE embeddings 174 distinctly separate cell types more effectively than other methods tested in zero-shot. Even though 175 UCE is not trained on the Tabula Sapiens v2 dataset, its embeddings more closely resemble those of 176 fine-tuned methods, which are directly trained on it. Moreover, cell types align correctly regardless 177 of whether the data was drawn from new donors or previously seen ones (Supplementary Fig. 1). 178 For all cell types in Tabula Sapiens v2, we calculate the silhouette width score of each zero-179 shot embedding method. For 67% of cell types, UCE has the highest silhouette score of any 180 method. UCE outperforms Geneformer on 80% of cell types, tGPT on 73% of cell types, and 181 scGPT on 83% of cell types. Notably, UCE accurately embeds B cells, while Geneformer and 182 scGPT fail to do so (Supplementary Fig. 2a). In Tabula Sapiens v2, the silhouette width score of 183 B cells is 93% higher in UCE versus scGPT and 25% higher versus Geneformer. Additionally, B 184 cells within the UCE embedding space can be accurately mapped to an existing reference. We train 185 a simple logistic classifier on the UCE embeddings of the Immune Cell Atlas [38], and then apply 186 the classifier to B cell embeddings from Tabula Sapiens v2. This classifier accurately classifies the 187 Tabula Sapiens v2 cells as memory and naive B cells (Supplementary Fig. 2b), which is confirmed 188 with marker gene analysis (Supplementary Fig. 2c). Overall, these results illustrate that UCE has 189 the unique capability to meaningfully integrate new, previously unseen datasets into a universal 190 cell representation space with no additional model training. 19

¹⁹² UCE embeds diverse cell types from organisms that were not part of the training data.

¹⁹³ UCE is also able to align datasets from novel species without additional model training. This ¹⁹⁴ is due to the fact that UCE is not dependent on any particular genome—each gene of interest is ¹⁹⁵ translated to a corresponding protein sequence, which is then embedded in a universal protein ¹⁹⁶ space. The representation in this space is independent of species and importantly does not require ¹⁹⁷ any judgment about whether particular pairs of genes are homologs or not. Since UCE can analyze ¹⁹⁸ cell atlas data from distinct species that were not part of the training set, the extent to which it ¹⁹⁹ succeeds in this task is a stringent test of whether UCE displays emergent behavior.

UCE's training data is composed of datasets from eight species: human, mouse, mouse lemur, zebrafish, pig, rhesus macaque, crab eating macaque and western clawed frog. We apply UCE to embed datasets from three novel species that were not included in the training set. For each species, we generate a zero-shot embedding and then determine the nearest cell type centroid from the IMA for each of the dataset's existing annotated cell types. For all three species we observed very high agreement between independent annotations of the novel species' data and the nearest cell type centroids in the IMA.

Within a dataset of green monkey lymph node and lung cells [39], for 13 of the 17 cell type centroids, the closest centroid from another species corresponds to the same cell type in the green monkey. This match extends to all 17 centroids when considering the three nearest centroids (Extended Data Table 1, Fig. 2c, 2d). Moreover, a population of lymph node cells that were originally labelled as B cells, form a distinct cluster in UCE space (Supplementary Fig. 3b). Differential expression analysis revealed that this cluster predominantly expresses a T cell marker, *Cd3d* (Supplementary Fig. 3a, 3c).

214

In the case of naked mole rat spleen and circulating immune cells [40], for 17 out of 24 cell

types, the nearest cross species centroid matches the naked mole rat cell type (Extended Data Table 215 1, Supplementary Fig. 4b). In the case of chicken, we embed two distinct chicken datasets, chick 216 retina [41] and developing chick heart [42] (Supplementary Fig. 5a, 5b). Different eye-specific 217 neurons within the chick retina map to mouse lemur neurons, such as chick oligodendrocytes, 218 which are closest to mouse lemur oligodendrocytes (Extended Data Table 1). In chicken heart, 12 219 of 15 cell type centroids are matched within the nearest two cross species centroids (Extended Data 220 Table 1). No bird species were included when training UCE. Altogether, these results highlight that 22 UCE can be directly applied to investigate new and diverse datasets from previously unobserved 222 species. 223

²²⁴ UCE learns a meaningful organization of cell types in previously unseen data.

Moving beyond metrics focused on individual cell type clusters, we also examined the structure of the universal embedding space as a whole, through the relative positioning of different cells within it. A meaningful arrangement of cell types emerges upon embedding all the cells from the Tabula Sapiens v2 dataset from the lung tissue (Fig. 3a). Not only do distinct cell types like T cells, monocytes and endothelial cells cluster together, but higher-level categories, such as immune cells and epithelial cells, are also clearly distinguished.

To systematically assess this organization of cells within the embedding, we compared distances between pairs of cell types across all tissues in the embedding space to their distances in the Cell Ontology tree [43] (Fig. 3b). We hypothesized that cells that are known to be similar based on the cell ontology would likely also be closer together in the embedding space, and that the degree of closeness would be correlated with ontological similarity. The results validate this relationship: at each additional unit of separation between cell types in the cell ontology tree, there is

a significant increase in the embedding distance in UCE between those cell types. We consistently
observed this trend up to a distance of 5 hops in the ontology tree (Fig. 3b). However, beyond that,
the effect levels off (Supplementary Fig. 6). This is expected due to the curse of dimensionality
in high-dimensional spaces and the variability in the level of ontological refinement in different
branches of the ontology (Supplementary Note 3).

We also noted significant colocalization among cells originating from the same developmen-242 tal lineages, in particular from the mesoderm, endoderm, and ectoderm germ layers. For Tabula 243 Sapiens v2, 90 out of 97 of the centroids for mesoderm-derived cell types had other mesoderm-244 derived cell type centroids as their closest neighbors. A similar pattern was observed for 46 of the 245 56 endoderm-derived cell types and 22 of the 30 ectoderm-derived cell types (Supplementary Fig. 246 7a). A neural network classifier trained to predict the germ layer of origin for individual held-out 247 cell types using their universal embeddings showed an accuracy of over 80% (Supplementary Fig. 248 7b). 249

The accuracy of cell type organization in the Tabula Sapiens v2 lung dataset was evaluated by 250 comparing it with other lung datasets in the IMA (Fig. 3c, Supplementary Fig. 8). Four different 25 endothelial cell subtypes are observed to map correctly to their corresponding counterparts in the 252 IMA. Similarly, lung ciliated cells correctly map to their counterpart in the larger corpus despite 253 the presence of four different ciliated cell subtypes (Fig. 3c). Further analysis of the alignment of 254 cell type centroids between Tabula Sapiens v2 and the IMA across all tissues showed an average 255 correct alignment of 56% for each tissue, as detailed in the Methods section. This alignment, based 256 on the three nearest neighbor cell type centroids, is 60% more accurate compared to that measured 257 in the original gene expression space (Fig. 3d). When focusing on the single nearest centroid, the 258

alignment accuracy improves by 93%. These results demonstrate that UCE can effectively learn a
universal representation of cell biology that not only enables discrimination between individual cell
types but also captures their relative similarities across scales with the potential to reveal deeper
insights into development and function.

²⁶³ A workflow for decoding the function of newly discovered cell types.

UCE's zero-shot embedding capabilities unlock novel computational analyses of scRNA-264 seq data and aid in hypothesis generation. Beyond identifying novel cell type clusters, UCE differs 265 from other methods in that the same cell type can also be easily compared against all previously 266 assayed cells across tissues, disease states and species. Moreover, UCE is not biased in this process 267 by existing annotations, opening the door for discovery of novel function (Fig. 4a). With existing 268 fine tuning based methods, every searched dataset would need to be integrated, requiring repeated 269 model retraining. Thus, UCE enables a new workflow for scRNA-seq data analysis that performs 270 an unbiased search across the universe of cell biology. 27

We present an example of this analysis by using the recently identified kidney Norn cell as a case study. The kidney Norn cell is the long-sought erythropoietin (*Epo*) producing cell in the kidney, and is characterized as fibroblast-like. We perform a zero-shot embedding of mouse renal cells from [44], which produces a cluster of cells corresponding to Norn cells (Fig. 4b).

Using a simple logistic classifier trained on the embedding of mouse renal cells, we identify Norn cell clusters in many kidney datasets. Since this classifier takes universal cell embeddings as an input, we can directly apply it to all 36 million cells in the IMA, in a manner unbiased by cell type annotations ascribed by previous studies. We also confirm these cell's Norn identity using marker gene analysis. Cells classified as Norn cells in the top 13 kidney datasets by Norn

abundancy demonstrate preferential expression of the Norn markers Dcn, Lpar1, Colla1, Cxcl12, 28 and Cfh (Extended Data Table 3). Notably, Epo transcripts, which are often missing from datasets 282 and lowly expressed, are not typically differentially expressed in these cells. Cxcl14, another 283 marker of Norn cells, displays mixed expression patterns in these predicted Norn cells (Fig. 4c). 284 The same pattern of marker gene expression is also found in cells from other tissues, including lung 285 and heart datasets (Fig. 4c). Additionally, these predicted cells also share a common set of genes 286 that are lowly expressed in mouse renal Norn cells (Supplementary Fig. 9). The tissues with the 287 highest number of predicted Norn cells were gonad, heart and lung. While Epo expression has been 288 previously observed in the heart and lung tissue, the mechanisms and cell types associated with 289 this expression, and their relation to kidney Norn cells have not been previously determined [45]. 290 Overall, this demonstrates that UCE can serve as an unbiased tool for predicting the existence of 291 novel cell types. 292

²⁹³ UCE helps interrogate alternate lung disease outcomes.

Lastly, we apply UCE and our simple Norn cell classifier to investigate Norn cells in lung diseases. We generate an embedding of lung cells sampled from patients with idiopathic pulmonary fibrosis (IPF), chronic obstructive pulmonary disease (COPD), or patients from a control group [46]. We identify Norn-like lung cells that preferentially express Norn markers in all three groups (Fig. 4d).

For these Norn-like lung cells, we identify differences across disease groups (Fig. 4e). COPD and IPF are both associated with elevated bloodstream *Epo*, but COPD has levels higher than IPF. Additionally, in patients with IPF, secondary erythrocytosis is absent or reduced compared to patients with COPD [47–49]. Given the identification of Norn-like cells in the lung, and Norn cell's production of *Epo*, it is possible that this difference in disease prognosis could be related to
 disease associated differences in Norn-like cells.

In COPD predicted Norn cells, there is a significantly greater ratio of *Epas1* : *Egln1* transcripts (p=0.035) than in IPF predicted Norn cells (Fig. 4e). *Epas1* is a master regulator of *Epo* transcription, which is degraded by the oxygen sensing enzyme encoded by *Egln1* [44]. Control and COPD predicted Norn cells express genes (*Bgn*, *Crispld2*) involved in glycosaminoglycan pathways at different levels than IPF predicted Norn cells [50,51]. IPF cells also have significantly lower expression of *Il6st* than cells in control or COPD groups.

Taken as a whole, these results indicate that Norn-like cells may be found in other tissues in the body, and may play a previously unidentified role in disease. UCE greatly facilitates an analysis of this scale and diversity because it is a universal model.

314 Discussion

³¹⁵ UCE is a single-cell foundation model that is built from the ground up to represent cell biology ³¹⁶ across the wide array of single-cell datasets. We envision UCE as an embedding approach that ³¹⁷ enables researchers to map any new data, including entire atlases, into an accurate, meaningful and ³¹⁸ universal space. The embedding space that emerges from UCE is highly structured and diverse and ³¹⁹ aligns cell types across tissues and species. Additionally, these cell types organize themselves in a ³²⁰ pattern that reflects existing biological knowledge.

The UCE model has broad implications for the creation of large foundation models for single cell biology. For large foundation models to be truly useful for scientific discovery, they must have unique qualities that distinguish them from existing methods. Zero-shot embeddings are one such important capability because it enables an intrinsically meaningful representation that can extend

insights beyond the data that has already been observed and annotated experimentally. Our results
 demonstrate that UCE can achieve such a generalizable representation across different datasets
 while maintaining accuracy on individual datasets, comparable to methods that require retraining
 for each specific dataset.

By building UCE, we enable new and novel analyses of scRNA-seq data. However, these 329 analyses and corresponding benchmarks are still far from perfect, as they are generally limited by a 330 focus on coarse cell type labels. To better understand single-cell foundation models, and especially 331 how they scale, new analyses and benchmarks that surpass this resolution limit must be developed. 332 For a precise representation of biology, models must incorporate core biological motivation. To 333 this end, we recognize that current scRNA-seq foundation models, including UCE, do not account 334 for any information contained in the raw RNA transcripts. By aligning these transcripts to the 335 reference genome, vital data on genetic variation and crucial RNA-splicing processes are discarded 336 [52]. Future single cell foundation models should seek to include this genomic precision at the 337 transcript level. As these models adopt more biologically-relevant features, they will increasingly 338 be able to simulate the biological processes of cells, leading to the creation of "Virtual Cells". 339

In 2002, Nobel laureate Sydney Brenner identified many of the core motivations for the creation of cell atlases and virtual cells. Virtual cells should be the goal of biological foundation modeling, because cells are the "real units of function and structure in an organism" [53]. Brenner also identified the need for such models to be computationally efficient, predictive, and able to generate new cell types. We believe that UCE represents a significant advancement in the progress towards a virtual cell. Through learning a universal representation of every cell state and type, we expect that UCE will be a valuable tool for analysis, annotation and hypothesis generation as the

³⁴⁷ scale and diversity of single-cell datasets continues to grow.

348 Methods

349 Overview of UCE.

UCE (Universal Cell Embedding) is a machine learning model for mapping single-cell gene expression profiles into a universal embedding space, denoted as \mathcal{U} . In this space, each cell c_i is represented as a d_{emb} -dimensional vector, where $d_{emb} = 1280$.

The model takes as input a dataset \mathcal{D} with N cells $\{\mathbf{c}_i\}_{i=1}^N$. Cells in \mathcal{D} can be drawn from one or more distinct scRNA-seq experiments. Each cell c_i in \mathcal{D} is described by a gene expression vector $\mathbf{x}^i \in \mathbf{N}^{K_i}$, where K_i is the number of genes measured in c_i and can differ across \mathcal{D} . The gene expression vectors $\mathbf{x}^i \in \mathbb{N}^{K_i}$ are not subset to those with high variance. UCE defines a function $f_u: \{\mathbb{N}^{K_i} \to \mathbb{R}^{d_{emb}}\}_{i=1}^N$ that maps each gene expression vector \mathbf{x}^i to its cell embedding vector \mathbf{h}^i .

358 Model input: Gene representation.

The expression of gene g in cell c_i is denoted by x_g^i , where g represents any protein-coding gene. The corresponding token embedding p_g is a pretrained embedding for the protein(s) encoded by the gene g. These embeddings are derived from a pretrained protein language model that takes an amino acid sequence as input and returns a d_p -dimensional embedding vector as output. To create p_g , we take the average of all proteins coded by gene g. In the context of UCE, we can formulate this as a dictionary that maps each gene g to a d_p -dimensional protein embedding vector. Specifically, we employ the ESM2 model, which yields embeddings of size $d_p = 5120$ [33, 34].

366 Model input: Cell representation.

For each cell c_i in the input dataset \mathcal{D} , we identify two distinct sets of protein-coding genes: the expressed genes G_i^+ and the non-expressed genes G_i^- . These sets are defined as follows:

$$G_i^+ = \{ g \mid x_a^i > 0 \}$$
(1)

$$G_i^- = \{g \mid x_g^i = 0\}$$
(2)

For producing the cell embedding, a multi-set of 1024 non-unique genes G_i^s are sampled from the expressed genes G_i^+ , with replacement. The probability of sampling a gene $g \in G_i^+$ is weighted by the log normalized expression of that gene, which can be formulated as:

$$P(g \mid c_i) = \frac{\log(x_g^i)}{\sum_{g \in G_i^+} \log(x_g^i)}$$
(3)

where x_g^i is the expression count of gene g in cell c_i , and the sum in the denominator is over all genes in G_i^+ .

Once the multi-set G_i^s is compiled for each cell c_i , we arrange the genes within each chromosome according to their genomic positions. Different chromosomes are specified using special chromosome start and end tokens. Start tokens are unique to each chromosome and species. Every chromosome group is combined into a single sequence, with chromosome order randomly determined. A cell-level *CLS* token is appended to the start of the sequence. It is designed to capture the cell-level embedding upon training the model. The final sequence of genes ordered by genomic location and separated by chromosome is referred to as the cell sentence S_i for cell c_i .

381 Transformer Architecture.

Each cell sentence S_i is fed into a transformer that consists of n_{lay} layers. Each layer contains a multi-head self-attention mechanism with n_{head} attention heads and a feedforward network operating over a hidden space of dimensionality d_{hid} . We also initialize sinusoidally-varying po-

sitional embeddings. Gene token embeddings are compressed using a single layer MLP to d_{emb} dimensional vectors before passing through the transformer.

387 Model output: Cell embedding.

The final output from the model is the cell embedding vector $\mathbf{h}_{cell}^i \in \mathcal{U}$ which corresponds to the d_{emb} -dimensional embedding of the CLS token in the final layer of the model following decoding with an additional MLP.

391 Model training: Cell representation.

At the time of training, we generate a set $G_i^{M+} \subset G_i^+$ by randomly selecting a certain percentage (r_{mask}) of genes from G_i^+ , without replacement. This set is used for computing the loss during training, and is masked from the cell representation.

The probability of sampling a gene $g \in G_i^+ \setminus G_i^{M+}$ (Equation 3) is then updated to be:

$$P(g \mid c_i) = \frac{\log(x_k^i)}{\sum_{g \in G_i^+ \setminus G_i^{M+}} \log(x_j^i)},\tag{4}$$

We also establish two additional gene sets to be used for loss computation: $G_i^{L+} \in G$ and $G_i^{L-} \in G$. G_i^{L+} and G_i^{L-} are randomly selected from the masked set of expressed genes G_i^{M+} and the set of unexpressed genes G_i^- respectively. Both G_i^{L+} and G_i^{L-} are of equal size, specifically $N_{loss}/2$. In the case of G_i^{L-} , the sampling is done without replacement unless $|G_i^-| < N_{loss}/2$. Similarly G_i^{L-} , is also sampled without replacement unless $|G_i^{M+}| < N_{loss}/2$. In this case, G_i^{M+} is used as-is alongwith additional samples drawn with replacement from the full set of expressed genes G_i^+ .

Model training: Loss Function.

To calculate the loss function for a given cell c_i , the cell embedding vector \mathbf{h}_{emb}^i is individu-

ally concatenated with every gene g within both G_i^{L+} and G_i^{L-} . These concatenated vectors then serve as input to a feedforward multilayer perceptron (MLP), which computes the probability that gene g is expressed within cell c_i .

⁴⁰⁸ \mathbf{h}_{cell}^{i} represents the embedding vector for cell c_i and p_g represents the token embedding for ⁴⁰⁹ gene g. Then the concatenated vector \mathbf{z}_{q}^{i} that serves as input to the MLP for cell c_i and gene g is:

$$p'_{q} = \mathsf{MLP}(p_{g}) \tag{5}$$

$$\mathbf{z}_{g}^{i} = [\mathbf{h}_{cell}^{i} || p_{g'}] \tag{6}$$

where || denotes the concatenation operation and p'_g is the compressed protein embedding. The MLP then processes this concatenated input to produce the predicted probability that gene g is expressed:

$$p(y_a^i) = \mathsf{MLP}(\mathbf{z}_a^i) \tag{7}$$

This probability is then used in the binary cross-entropy loss function. The true classification labels for each gene's expression status in cell c_i are represented by the vector \mathbf{y}^i . UCE is trained to accurately predict the expression of genes in G_i^{L+} and the lack of expression in G_i^{L-} . The model is trained using a binary cross-entropy loss, which is averaged across all N_{loss} genes and all Ncells in the minibatch as follows:

$$L = -\frac{1}{N} \sum_{i=1}^{N} \frac{1}{N_{loss}} \sum_{j=1}^{N_{loss}} \left[y_j^i \log(p(y_j^i)) + (1 - y_j^i) \log(1 - p(y_j^i)) \right]$$
(8)

418

For further details on hyperparmeter choices please see Supplementary Table 2.

419 Creating the IMA and dataset preprocessing.

The Integrated Mega-scale Atlas (IMA) used to train UCE was created by combining scRNAseq datasets from multiple publicly available sources. The majority of IMA data (33.9 million cells and 285 datasets) is human and mouse data downloaded from CZI Cell X Gene (CxG) Census [36] version "2023-07-10" (July 10th, 2023). Duplicate cells were removed by selecting primary cells only. The remainder of the IMA is composed of 2.3 million cells from 28 datasets, from eight different species: human, mouse, zebrafish, rhesus macaque, crab-eating macaque, mouse lemur, frog, and pig.

For datasets from the CxG Census, preprocessing only involved filtering cells by minimum gene counts (200) and genes by a minimum cells count of 10. No highly variable gene selection was applied. For datasets collected from other sources, preprocessing was not uniform.

For visualization of the IMA (Fig. 1b), predicting green monkey cell types (Fig. 2d), match-430 ing new species centroids (Extended Data Table 1), and prediction of Norn-like cells (Fig. 4, 431 Supplementary Fig. 9) a representative sample of the IMA was used in place of the full 36 million 432 cells. This representative sample was used in order to speed up computationally intensive tasks like 433 UMAP calculation. The sample was created by randomly choosing 10,000 cells from each dataset, 434 without replacement. For datasets with fewer than 10,000 cells, the entire dataset was included. In 435 total, this representative sample has 2,969,114 cells. The average number of cells per dataset in the 436 sample is 9486. For visualization and centroid calculation, cell types in the sample were coarsened 437 by mapping them to a set of 51 coarse cell types. 438

439 Model Evaluation.

• Zero-shot embedding quality and clustering For evaluating the quality of embeddings, we

- used metrics from the single-cell integration benchmark [16].
- **Cell type organization** For each cell type dendrogram the Euclidean distance was used to perform hierarchical clustering across all cells.
- **Comparison to cell ontology** Here, we used the tree distance between any two cell types in Cell Ontology [43]. To determine the Euclidean distance distribution, we sampled 100,000 random pairs of cells from Tabula Sapiens v2.
- Zero-shot cell type alignment to IMA For each cell type θ , a centroid was identified separately for data from Tabula Sapiens v2 (TSv2) c_{θ}^{T} and from IMA c_{θ}^{I} . For each cell type that is present in both TSv2 and the IMA, the 3 nearest neighbor cell type centroids N_{θ}^{T} to the centroid in Tabula Sapiens c_{θ}^{T} were identified. These neighbors could be either from Tabula Sapiens or from the IMA.
- If this set of neighbors \mathbf{N}_{θ}^{T} to the anchor centroid from TSv2 data c_{θ}^{T} contains the centroid for the same cell type in IMA data c_{θ}^{I} , then this was counted as a correct match.

This analysis was performed per tissue, both in the UCE embedding space as well as in the original expression space (after log-normalization). In case of the original data representation, the set of 5704 shared genes across all human datasets were used to represent each cell.

458 Differential expression analysis of predicted Norn cells.

⁴⁵⁹ A logistic classifier was trained to predict cell types from UCE embeddings on mouse kidney ⁴⁶⁰ cells. This classifier was then applied to UCE embeddings from the representative sample of IMA ⁴⁶¹ datasets. Datasets were then split by tissue, and the datasets with the most predicted norn cells in

462 each tissue were used for differential expression analysis. The top 13 kidney datasets, top 6 lung
463 and top 6 heart datasets were chosen.

For each individual (full) dataset, RNA counts were log normalized, and then differential expression was run using default settings as implemented in Scanpy [54], comparing predicted Norn cells to all other cells in the dataset. The results of these differential expression tests were used to determine the log fold change of marker genes in predicted Norn cells (Fig. 4c, Supplementary Fig. 9).



Figure 1: The Universal Cell Embedding Model is a large foundation model for single cell biology (a) Overview of the Universal Cell Embedding (UCE) model. UCE has a unique, biologically motivated representation of cells (blue) and training scheme (purple). Given the gene expression for a single cell, UCE samples with replacement genes that were expressed, weighted by their level of expression. Each gene is represented using a 'token' corresponding to its protein product. Gene tokens are represented numerically by using ESM2 protein embeddings, a 15 billion parameter protein language model that takes amino acid sequences as an input. The gene tokens are sorted by genomic location and grouped by chromosome. Chromosome groups are delineated by specific chromosome start tokens and end tokens, joined, and then passed into a transformer neural network. The embedding of the cell is determined by taking the final layer output of a special CLS token that is appended before all the other tokens. To train the UCE model, a portion of genes that were expressed are masked. The model next combines the protein embeddings corresponding to each of these genes with the embedding of the cell, and passes this joint representation through a neural network that predicts if a given gene was expressed in the cell or not. This objective function is then used to update the weights of the model. (b) UMAP visualizations of the universal cell embedding space. We apply UCE to embed 36 million cells, with more than 1,000 uniquely named cell types, from hundreds of datasets, dozens of tissues and eight species, creating an Integrated Mega-scale Atlas (IMA) spanning the universe of cell biology.



Figure 2: Zero-shot cell embedding capabilities of UCE (a) Comparison of zero-shot and fine-tuned single-cell embedding models. A zero-shot embedding model maps new data directly to the to the representation space, with no additional model training. In contrast, fine-tuned models must first be retrained on a given dataset, and only then can be applied on that dataset, fundamentally altering the model's representation space. (b) UMAP embeddings of UCE and other methods for Tabula Sapiens v1 and v2, colored by cell type. UCE zero-shot embeddings closely resemble the embeddings of fine-tuned methods scVI and scArches, demonstrating clusters that correspond to cell types, in contrast to the other zero-shot methods Geneformer and scGPT. (c) UMAP of cells from a new species, green monkey colored by cell type. UCE is able to generate high-quality zero-shot embeddings of novel species that were never seen during training. The UCE embedding for green monkey mediastinal lymph node [39] recaptures cell type clusters. Notably, a population of B cells (blue) clusters nearby to T cells, potentially due to expression of *Cd3* (Supplementary Fig. 1). (d) Green monkey lymph node cells can be accurately annotated using the IMA. A logistic classifier is first trained to predict cell types based on UCE embeddings of human lymph node cells. The classifier is then directly applied on green monkey cells to predict the cell types. Predicted cell types have high agreement with the original cell type annotations, demonstrating that UCE can be used to transfer cell type annotations to novel species.



Figure 3: UCE learns meaningful organization of cell types (a) The UCE space generated for new, previously unseen data shows a meaningful arrangement of cell types. Lung data was used from new donors from the Tabula Sapiens Consortium. Dendrogram of hierarchical clustering of all annotated cell types in the UCE embedding space. Closely connected cell types in the dendrogram show meaningful biological relationships both at finer and coarser scale resolutions. (b) Evaluation of the organization of cell types in the embedding space when compared to Cell Ontology. The x-axis depicts the density of Euclidean distances between all pairs of cells across all tissues for these new donors from the Tabula Sapiens Consortium. The y-axis shows the corresponding tree distance between cell types as found in the Cell Ontology. Stars denote statistical significance, which was established using a one-sided t-test. (c) Mapping data from new donors to the Integrated Mega-scale Atlas (IMA) across multiple lung datasets. Red labels correspond to data from new donors, grey are from IMA datasets. All cell type labels from multiple datasets are displayed as-is, with no modifications or reformatting of text. Accurate alignment between the new dataset and IMA is observed at finer resolution. Four different subtypes of endothelial cells are shown to correctly map to their corresponding counterparts in the complete mega-scale atlas. In the case of lung ciliated cells, they map more closely to their matching counterpart as compared to all other ciliated cell subtypes also present in the IMA. (d) Quantification of cell type alignment between new dataset and IMA. Accuracy in 3-nearest centroid matches between new dataset and IMA cell types at the finest level of original annotation. Results are measured across all 27 tissues in Tabula Sapiens v2 for both the UCE space and the original gene expression space. Tissues are ordered by accuracy in the UCE space.



Figure 4: Norn Cell Case Study: UCE unlocks new analyses of single cell datasets (a) Overview of a novel single cell analysis workflow that UCE facilitates. Analysis begins with (1) the identification of a novel cell type (circled) within the embedding space, using methods such as clustering and confirmation using marker gene analysis. (2) Next, the novel cell type can be easily identified in other datasets profiled from the same tissue (for example, kidney). A simple classifier, such as a logistic classifier, is trained to predict cell types from universal cell embeddings, and is then applied to embeddings from other datasets of the same tissue (kidney), to confirm the cell type's existence and improve its characterization. (3) Finally, the same simple classifier can be applied to the embeddings of cells from any other tissue, to find cell types with similar biological functions or patterns of gene expression. (b) Identification of novel Norn cells in mouse kidney. UMAP visualization of zero-shot embedding of mouse renal cells from Kragesteen et al. [44]. Norn cells form a distinct cluster within the embedding space (circled). (c) Identification of Norn cells and Norn-like cells across tissues. A logistic classifier is trained to predict Norn cells from universal cell embeddings, and is then applied to other kidney datasets (left) and datasets from lung and heart (right). The log fold change of known Norn marker genes between cells predicted to be Norn cells and the remaining cells within each dataset is visualized. Cells which are predicted to be Norn-like preferentially express Norn markers in kidney, as well as in lung and heart. Notably, Cxcl14 has a mixed pattern of expression among some datasets. (d) Cells predicted to be Norn cells within a lung disease dataset [46] express known Norn markers, as demonstrated by log fold change (LFC). e Differential gene expression in predicted Norn cells, grouped by disease status. There are significant differences in gene expression of important Norn markers and genes involved in the production of erythropoietin (Epo) between cells from IPF, COPD and control patients. Patients with IPF and COPD are known to have elevated levels of blood stream Epo, with COPD patients having greater bloodstream Epo levels than patients with IPF.

Data availability

The full list of datasets used to train UCE are in Extended Data Table 2. Most of these datasets are available to download from CellXGene [36]. Tabula Sapiens v2, used for model evaluation, will be made available upon publication.

Datasets analyzed in the paper are publicly available to download. The green monkey lung and lymph node dataset is available with accession code GSE156755. The naked mole rat dataset is available with accession code GSE132642. The chicken retina dataset is available with accession code GSE159107. The chicken heart dataset is available with accession code GSE149457. The mouse kidney dataset is available with accession code GSE193321. The human lung disease

dataset is available with accession code GSE136831.

479 Code availability

480 UCE was written in Python using the PyTorch library. The source code is available on Github at 481 https://github.com/snap-stanford/uce.

Acknowledgements

We thank Rok Sosič, Kexin Huang, Charlotte Bunne, Hanchen Wang, Michihiro Yasunaga, Michael 483 Moor, Minkai Xu, Mika Jain, George Crowley, Maria Brbić, Jonah Cool, Nicholas Sofroniew, 484 Andrew Tolopko, Ivana Jelic, Ana-Maria Istrate and Pablo Garcia-Nieto for discussions and for 485 providing feedback on our manuscript. We acknowledge support from Robert C. Jones for help 486 with accessing and analyzing the Tabula Sapiens v2 dataset. We acknowledge support from the 487 Chan Zuckerberg Initiative, including help with accessing and processing CxG datasets. We grate-488 fully acknowledge the support of DARPA under Nos. N660011924033 (MCS); NSF under Nos. 489 OAC-1835598 (CINES), CCF-1918940 (Expeditions), Stanford Data Science Initiative, Wu Tsai 490 Neurosciences Institute, Amazon, Genentech, GSK, Hitachi, Juniper Networks, and KDDI. Y. 491 RH. acknowledges funding support form GlaxoSmithKline. L.S. was supported by the American 492 Slovenia Education Foundation (ASEF). Icons created with BioRender.com. 493

494 Author information

Y.RS., Y.RH., S.Q. and J.L. conceived the study. Y.RS, Y.RH., S.Q. and J.L. performed research,
contributed new analytical tools, designed algorithmic frameworks, analyzed data and wrote the
manuscript. Y.RS. and Y.RH. performed experiments and developed the software. A.A. and L.S.
contributed to code and performed analyses. T.S. provided annotated data.

499 **References**

500

- ⁵⁰¹ 1. Vaishnav, E. D. *et al.* The evolution, evolvability and engineering of gene regulatory dna. ⁵⁰² *Nature* **603**, 455–463 (2022).
- Onuchic, J. N., Luthey-Schulten, Z. & Wolynes, P. G. Theory of protein folding: the energy landscape perspective. *Annual review of physical chemistry* 48, 545–600 (1997).
- Moris, N., Pina, C. & Arias, A. M. Transition states and cell fate decisions in epigenetic landscapes. *Nature Reviews Genetics* 17, 693–703 (2016).
- 4. Treutlein, B. *et al.* Reconstructing lineage hierarchies of the distal lung epithelium using single-cell rna-seq. *Nature* **509**, 371–375 (2014).
- 5. Waddington, C. H. *The strategy of the genes* (Routledge, 1957).
- 510 6. Schaum, N. *et al.* Single-cell transcriptomics of 20 mouse organs creates a tabula muris: The tabula muris consortium. *Nature* 562, 367 (2018).
- ⁵¹² 7. Regev, A. *et al.* The human cell atlas. *elife* **6**, e27041 (2017).
- 8. Rood, J. E., Maartens, A., Hupalowska, A., Teichmann, S. A. & Regev, A. Impact of the human cell atlas on medicine. *Nature medicine* 28, 2486–2496 (2022).
- 515 9. Tanay, A. & Regev, A. Scaling single-cell genomics from phenomenology to mechanism.
 516 *Nature* 541, 331–338 (2017).
- ⁵¹⁷ 10. Plass, M. *et al.* Cell type atlas and lineage tree of a whole complex animal by single-cell transcriptomics. *Science* **360**, eaaq1723 (2018).
- ⁵¹⁹ 11. Consortium^{*}, T. S. *et al.* The tabula sapiens: A multiple-organ, single-cell transcriptomic atlas ⁵²⁰ of humans. *Science* **376**, eabl4896 (2022).
- Siletti, K. *et al.* Transcriptomic diversity of cell types across the adult human brain. *Science* **382**, eadd7046 (2023).
- Li, H. *et al.* Fly cell atlas: A single-nucleus transcriptomic atlas of the adult fruit fly. *Science* 375, eabk2432 (2022).
- Eraslan, G., Avsec, Ž., Gagneur, J. & Theis, F. J. Deep learning: new computational modelling
 techniques for genomics. *Nature Reviews Genetics* 20, 389–403 (2019).
- 15. Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for single-cell transcriptomics. *Nature methods* **15**, 1053–1058 (2018).
- Luecken, M. D. *et al.* Benchmarking atlas-level data integration in single-cell genomics. *Nature methods* 19, 41–50 (2022).
- ⁵³¹ 17. Argelaguet, R., Cuomo, A. S., Stegle, O. & Marioni, J. C. Computational principles and ⁵³² challenges in single-cell data integration. *Nature biotechnology* **39**, 1202–1215 (2021).
- 18. Lotfollahi, M. *et al.* Mapping single-cell data to reference atlases by transfer learning. *Nature biotechnology* 40, 121–130 (2022).

- ⁵³⁵ 19. Tarashansky, A. J. *et al.* Mapping single-cell atlases throughout metazoa unravels cell type ⁵³⁶ evolution. *Elife* **10**, e66747 (2021).
- ⁵³⁷ 20. Brown, T. *et al.* Language models are few-shot learners. *Advances in neural information* ⁵³⁸ *processing systems* **33**, 1877–1901 (2020).
- ⁵³⁹ 21. OpenAI. Gpt-4 technical report (2023). 2303.08774.
- 540 22. Anil, R. et al. Palm 2 technical report. arXiv preprint arXiv:2305.10403 (2023).
- ⁵⁴¹ 23. Touvron, H. *et al.* Llama: Open and efficient foundation language models. *arXiv preprint* ⁵⁴² *arXiv:2302.13971* (2023).
- ⁵⁴³ 24. Kirillov, A. *et al.* Segment anything. *arXiv preprint arXiv:2304.02643* (2023).
- ⁵⁴⁴ 25. Bommasani, R. *et al.* On the opportunities and risks of foundation models. *arXiv preprint* ⁵⁴⁵ *arXiv:2108.07258* (2021).
- ⁵⁴⁶ 26. Avsec, Ž. *et al.* Effective gene expression prediction from sequence by integrating long-range
 ⁵⁴⁷ interactions. *Nature methods* 18, 1196–1203 (2021).

Rives, A. *et al.* Biological structure and function emerge from scaling unsupervised learning
 to 250 million protein sequences. *Proceedings of the National Academy of Sciences* 118,
 e2016239118 (2021).

- ⁵⁵¹ 28. Theodoris, C. V. *et al.* Transfer learning enables predictions in network biology. *Nature* 1–9 (2023).
- ⁵⁵³ 29. Cui, H. *et al.* scgpt: Towards building a foundation model for single-cell multi-omics using generative ai. *bioRxiv* 2023–04 (2023).
- ⁵⁵⁵ 30. Stegle, O., Teichmann, S. A. & Marioni, J. C. Computational and analytical challenges in ⁵⁵⁶ single-cell transcriptomics. *Nature Reviews Genetics* **16**, 133–145 (2015).
- ⁵⁵⁷ 31. Quake, S. R. The cell as a bag of rna. *Trends in Genetics* **37**, 1064–1068 (2021).
- ⁵⁵⁸ 32. Vaswani, A. *et al.* Attention is all you need. *Advances in neural information processing* ⁵⁵⁹ *systems* **30** (2017).
- 33. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language
 model. *Science* 379, 1123–1130 (2023).
- ⁵⁶² 34. Rosen, Y. *et al.* Towards universal cell embeddings: Integrating single-cell rna-seq datasets
 ⁵⁶³ across species with saturn. *bioRxiv* (2023).
- ⁵⁶⁴ 35. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. Bert: Pre-training of deep bidirectional ⁵⁶⁵ transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- ⁵⁶⁶ 36. Biology, C. S.-C. *et al.* Cz cellxgene discover: A single-cell data platform for scalable explo-⁵⁶⁷ ration, analysis and modeling of aggregated data. *bioRxiv* 2023–10 (2023).
- ⁵⁶⁸ 37. Gordon, S., Plüddemann, A. & Martinez Estrada, F. Macrophage heterogeneity in tissues: ⁵⁶⁹ phenotypic diversity and functions. *Immunological reviews* **262**, 36–55 (2014).
- ⁵⁷⁰ 38. Conde, C. D. *et al.* Cross-tissue immune cell analysis reveals tissue-specific features in humans. *Science* **376**, eabl5197 (2022).

- ⁵⁷² 39. Speranza, E. *et al.* Single-cell rna sequencing reveals sars-cov-2 infection dynamics in lungs ⁵⁷³ of african green monkeys. *Science translational medicine* **13**, eabe8146 (2021).
- 40. Hilton, H. G. *et al.* Single-cell transcriptomics of the naked mole-rat reveals unexpected features of mammalian immunity. *PLoS Biology* **17**, e3000528 (2019).
- 41. Yamagata, M., Yan, W. & Sanes, J. R. A cell atlas of the chick retina based on single-cell transcriptomics. *Elife* **10**, e63907 (2021).
- 42. Mantri, M. *et al.* Spatiotemporal single-cell rna sequencing of developing chicken hearts iden tifies interplay between cellular differentiation and morphogenesis. *Nature communications* 12, 1771 (2021).
- 43. Bard, J., Rhee, S. Y. & Ashburner, M. An ontology for cell types. *Genome biology* **6**, 1–5 (2005).
- 44. Kragesteen, B. K. *et al.* The transcriptional and regulatory identity of erythropoietin producing
 cells. *Nature medicine* 1–10 (2023).
- 45. Haine, L. *et al.* Cytoprotective effects of erythropoietin: What about the lung? *Biomedicine & Pharmacotherapy* 139, 111547 (2021).
- 46. Adams, T. S. *et al.* Single-cell rna-seq reveals ectopic and aberrant lung-resident cell populations in idiopathic pulmonary fibrosis. *Science advances* **6**, eaba1983 (2020).
- 47. Tassiopoulos, S. *et al.* Erythropoietic response to hypoxaemia in diffuse idiopathic pulmonary
 fibrosis, as opposed to chronic obstructive pulmonary disease. *Respiratory Medicine* 95, 471–
 475 (2001).
- 48. Abdel-Aziz, C., Okaily, N. & Kasem, A. Erythropoietin: role in idiopathic pulmonary fibrosis
 revisited. *The Egyptian Journal of Chest Diseases and Tuberculosis* 69, 716 (2020).
- 49. Tsantes, A. E. *et al.* Red cell macrocytosis in hypoxemic patients with chronic obstructive pulmonary disease. *Respiratory medicine* **98**, 1117–1123 (2004).
- 596 50. Safran, M. *et al.* The GeneCards suite. In Abugessaisa, I. & Kasukawa, T. (eds.) *Practical* 597 *guide to life science databases*, 27–56 (Springer Singapore, Singapore, 2021).
- ⁵⁹⁸ 51. Stelzer, G. *et al.* The genecards suite: from gene data mining to disease genome sequence ⁵⁹⁹ analyses. *Current Protocols in Bioinformatics* **54**, 1.30.1–1.30.33 (2016).
- 52. Amaral, P. et al. The status of the human gene catalogue. Nature 622, 41–47 (2023).
- ⁶⁰¹ 53. Brenner, S. Nature's gift to science (nobel lecture). *Chembiochem* **4**, 683–687 (2003).
- ⁶⁰² 54. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data ⁶⁰³ analysis. *Genome Biology* **19**, 15 (2018).