1 Surface protein imputation from single cell transcriptomes

2 by deep neural networks

- 3 Zilu Zhou^{1,2}, Chengzhong Ye³, Jingshu Wang², Nancy R. Zhang^{2*}
- 4 1) Graduate Group in Genomics and Computational Biology, University of Pennsylvania,
- 5 Philadelphia, PA
- 6 2) Department of Statistics, University of Pennsylvania, Philadelphia, PA
- 7 3) School of Medicine, Tsinghua University, Beijing, China
- 8 * Correspondence:
- 9
- 10 Nancy R. Zhang
- 11 <u>nzh@wharton.upenn.edu</u>
- 12 (215) 898-8007
- 13 Department of Statistics
- 14 The Wharton School
- 15 University of Pennsylvania
- 16
- 17 While single cell RNA sequencing (scRNA-seq) is invaluable for studying cell
- 18 populations, cell-surface proteins are often integral markers of cellular function and
- 19 serve as primary targets for therapeutic intervention. Here we propose a transfer learning
- 20 framework, single <u>cell Transcriptome to Protein prediction with deep neural network</u>
- 21 (cTP-net), to impute surface protein abundances from scRNA-seq data by learning from
- 22 existing single-cell multi-omic resources.
- 23 Keywords: genomics, immunophenotypes, single cell sequencing, deep learning, prediction,
- 24 multi-omics

25 Introduction

- 26 Recent technological advances allow the simultaneous profiling, across many cells in parallel, of
- 27 multiple omics features in the same cell ¹⁻⁵. In particular, high throughput quantification of the
- transcriptome and a selected panel of cell surface proteins in the same cell is now feasible

through the REAP-seq and CITE-seq protocols ^{2, 3}. Cell surface proteins can serve as integral 29 30 markers of specific cellular functions and primary targets for therapeutic intervention. Immunophenotyping by cell surface proteins has been an indispensable tool in hematopoiesis, 31 32 immunology and cancer research during the past 30 years. Yet, due to technological barriers 33 and cost considerations, most single cell studies, including Human Cell Atlas project ⁶, quantify 34 the transcriptome only and do not have cell-matched measurements of relevant surface proteins 35 ^{7,8}. Sometimes, which cell types and corresponding surface proteins are essential become apparent only after exploration by scRNA-seq. This motivates our inquiry of whether protein 36 37 abundances in individual cells can be accurately imputed by the cell's transcriptome. We propose cTP-net (single cell Transcriptome to Protein prediction with deep neural network), 38 39 a transfer learning approach based on deep neural networks that imputes surface protein 40 abundances for scRNA-seq data. Through comprehensive benchmark evaluations and 41 applications to Human Cell Atlas and acute myeloid leukemia data sets, we show that cTP-net 42 outperform existing methods and can transfer information from training data to accurately impute 24 immunophenotype markers, which achieve a more detailed characterization of 43 44 cellular state and cellular phenotypes than transcriptome measurements alone. cTP-net relies, 45 for model training, on accumulating public data of cells with paired transcriptome and surface 46 protein measurements.

47

48 Results

49 Method overview

50 An overview of cTP-net is shown in Figure 1a. Studies based on both CITE-seq and REAP-seq 51 have shown that the relative abundance of most surface proteins, at the level of individual cells,

is only weakly correlated with the relative abundance of the RNA of its corresponding gene ^{2, 3, 9}. 52 53 This is due to technical factors such as RNA and protein measurement error ¹⁰, as well as inherent stochasticity in RNA processing, translation and protein transport ¹¹⁻¹⁵. To accurately 54 impute surface protein abundance from scRNA-seq data, cTP-net employs two steps: (1) 55 56 denoising of the scRNA-seq count matrix and (2) imputation based on the denoised data through a transcriptome-protein mapping (Figure 1a). The initial denoising, by SAVER-X¹⁶. 57 produces more accurate estimates of the RNA transcript relative abundances for each cell. 58 59 Compared to the raw counts, the denoised relative expression values have significantly 60 improved correlation with their corresponding protein measurement (Figure 1b, S3a, S4ab). Yet, 61 for some surface proteins, such as CD45RA, this correlation for denoised expression is still extremely low. 62

63 The production of a surface protein from its corresponding RNA transcript is a complicated 64 process involving post-transcriptional modifications and transport ¹¹, translation ¹², posttranslational modifications ¹³ and protein trafficking ¹⁴. These processes depend on the state of 65 the cell and the activities of other genes ^{9, 15}. To learn the mapping from a cell's transcriptome to 66 the relative abundance of a given set of surface proteins, cTP-net employs a multiple branch 67 deep neural network (MB-DNN, Figure S1). Deep neural networks have recently shown success 68 in modeling complex biological systems ^{17, 18}, and more importantly, allow good generalization 69 across data sets^{16, 19}. Generalization performance is an important aspect of cTP-net, as we 70 would like to perform imputation on tissues that do not exactly match the training data in cell 71 72 type composition. Details of the cTP-net model and training procedure, as well as of alternative models and procedures that we have tried, are in Methods and Supplementary Note. 73

74

75 Imputation accuracy assessment and transfer learning

76 Evaluation via random holdout

77 To examine imputation accuracy, we first consider the ideal case where imputation is conducted on cells of types that exactly match those in training data. For benchmarking, we used 78 79 peripheral blood mononuclear cells (PBMCs) and cord blood mononuclear cells (CBMCs) processed by CITE-seq and REAP-seq^{2,3}, described in Table S1. We employed holdout 80 method, where the cells in each data set were randomly partitioned into two sets: a training set 81 82 with 90% of the cells and a holdout set with the remaining 10% of the cells for validation 83 (Methods, Figure S2a). Each cell type is well represented in both the training and validation 84 sets. Figure 1b and S3a show that, for all proteins examined in the CITE-seq PBMC data, cTPnet imputed abundances have much higher correlation to the measured protein levels, as 85 compared with the denoised and raw RNA counts of the corresponding genes. We obtained 86 87 similar results for the CITE-seq CBMC and REAP-seq PBMC data sets (Figure S4ab).

88

89 Generalization to unseen cell types

Next, we considered the generalization accuracy of cTP-net, testing whether it produces accurate imputations for cell types that are not present in the training set. For each of the highlevel cell types in each data set in Table S2, all cells of the given type are held out during training, and cTP-net, trained on the rest of the cells, was then used to impute protein abundances for the held out cells (Methods, Figure S2b). We did this for each cell type and generated an "out-of-cell-type" prediction for every cells.

Across all benchmarking data sets and all cell types, these out-of-cell-type predictions still
improve significantly upon the corresponding RNA counts while slightly inferior in accuracy to
the traditional holdout validation predictions above (Figure 2a, S4a). This indicates that cTP-net

99 provides informative predictions on cell types not present during training, vastly improving upon100 using the corresponding mRNA transcript abundance as proxy for the protein level.

101

102 Generalization across tissue and lab protocol

103 To further examine the case where cell types in the training and test data are not perfectly 104 aligned, we considered a scenario where the model is applied to perform imputation on a tissue that differs from the training data. We trained cTP-net on PBMCs and then applied it to perform 105 imputation on CBMCs, and vice versa, using the data from Stoeckius et al.³ (Methods), Cord 106 107 blood is expected to be enriched for stem cells and cells undergoing differentiation, whereas peripheral blood contains well-differentiated cell types, and thus the two populations are 108 109 composed of different but related cell types. Figure 2a and S3b shows the result on training on 110 CBMC and then imputing on PBMC. Imputing across tissue markedly improves the correlation 111 to the measured protein level, as compared to the denoised RNA of the corresponding gene. 112 but is worse than imputation produced by model trained on the same population. For practical 113 use, we have trained a network using the all cell populations combined, which indeed achieves better accuracy than a network trained on each separately (Methods, Figure S3b, S4ac). The 114 115 weights for this network are publicly available at https://github.com/zhouzilu/cTPnet. 116 We then tested whether cTP-net's predictions are sensitive to the laboratory protocol, and in

seq's standard, and vice versa. Using a benchmarking design similar to above, we found that, in
general, cTP-net maintains good generalization power across these two protocols (Figure 2a,

particular, whether networks trained using CITE-seg data yields good predictions by REAP-

120 S3b).

121

117

122 Comparison to Seurat v3

Seurat v3 anchor transfer ²⁰ is a recent approach that uses cell alignment between data sets to 123 124 impute features for single cell data. For comparison, we applied Seurat v3 anchor transfer to the holdout validation and out-of-cell-type benchmarking scenarios above (Methods). In the 125 validation scenario, we found the performance of cTP-net and Seurat v3 to be comparable, with 126 127 cTP-net slightly better, as both methods can estimate protein abundance by utilizing marker 128 genes to identify the cell types. cTP-net, however, vastly improves upon Seurat in the out-of-129 cell-type scenario (Figure 2a, S5a). This is because cTP-net's neural network, trained across a 130 diversity of cell types, learns a direct transcriptome-protein mapping that can more flexibly 131 generalize to unseen cell types, while Seurat v3 depends on a nearest neighbor method that can only sample from the training dataset. As shown by the cross-population and out-of-cell-132 type benchmarking above, cTP-net does not require direct congruence of cell types across 133 134 training and test sets.

135 In addition to predictions on unseen cell type, cTP-net also improves upon the existing state-ofthe-art in capturing within cell-type variation in protein abundance. As expected, within cell-type 136 137 variation is harder to predict, but cTP-net's imputations nevertheless achieve high correlations with measured protein abundance for a subset of proteins and cell types (Figure S3c, S4d). 138 139 Compared to Seurat v3, cTP-net's imputations align more accurately with measured protein 140 levels when zoomed into cells of the same type (Figure 2b, S5b); see Figure 2c, for example, CD11c in CD14-CD16+ monocytes, CD2 in CD8 T cells, and CD16 in dendritic cells. All of these 141 142 surface proteins have important biological function in the corresponding cell types, as CD11c 143 helps trigger respiratory burst in monocyte ²¹, CD2 co-stimulates molecule on T cells ²² and CD16 differentiate DC subpopulation ²³. The learning of such within-type heterogeneity gives 144 145 cTP-net the potential to attain higher resolution in the discovery and labeling of cell states.

146

147 Network interpretation and feature importance

148 What types of features are being used by cTP-net to form its imputation? To interpret the 149 network, we conducted a permutation-based interpolation analysis, which calculates a 150 permutation feature importance for each protein-gene pair (Methods, Figure S6a). Interpolation can be done using all cells, or cells of a specific type, the latter allowing us to probe 151 152 relationships that may be specific to a given cell type. Applying this analysis to cTP-net trained 153 on PBMC, we found that, at the level of the general population that includes all cell types, the 154 most important genes for the prediction of each protein are those that exhibit the highest celltype specificity in expression (Table S3). This is because most of these surface proteins are 155 156 cell type markers, and thus when cells of all types are pooled together, "cell type" is the key latent variable that underlies their heterogeneity. Within cell type interpolation, on the other 157 hand, reveals genes related to RNA processing, RNA binding, protein localization and 158 159 biosynthetic processes, in addition to immune-related genes that differentiate the immune cell 160 sub-types (Table S4). This analysis shows that cTP-net combines different types of features, 161 both cell type markers and genes involved in RNA to protein conversion and transport, to 162 achieve multiscale imputation accuracy.

In addition, we analyzed the bottleneck layer with 128 nodes before the network branched out to the protein-specific layers. We performed dimension reduction (UMAP) directly on the bottleneck layer intermediate output of 7000 PBMCs from CITE-seq. Figure S6b shows that the cells are cleanly separated into different clusters, representing cell types as well as gradients in surface protein abundance. This confirms that the bottleneck layer captures the essential information on cell stages and transitions, and that each subsequent individual branch then predicts its corresponding protein's abundance.

170

171 Application to Human Cell Atlas

172 Having benchmarked cTP-net's generalization accuracy across immune cell types, tissues, and technologies, we then applied the network trained on the combined CITE-seq dataset of 173 PBMCs,CBMCs and bone marrow mononuclear cells (BMMCs) ^{3, 24} to perform imputation for the 174 175 Human Cell Atlas CBMC and BMMC data sets (Table S1). Figure 3 shows the raw RNA count 176 and predicted surface protein abundance for 24 markers across 6023 BMMCs from sample 177 MantonBM1 and 4176 CBMCs from sample MantonCB1. (Similar plots for the other 7 BMMC 178 and 7 CBMC samples are shown in Figure S8, S9). Similar to what was observed for actual measured protein abundances in the CITE-seq and REAP-seq studies, the imputed protein 179 levels differ markedly from the RNA expression of its corresponding gene, displaying higher 180 181 contrast across cell types and higher uniformity within cell type. Thus, the imputed protein levels 182 can serve as interpretable intermediate features for the identification and labelling of cell states. 183 For example, imputed CD4 and CD8 levels separate CD4+ T cells from CD8+ T cells with high 184 confidence. Further separation of naïve T cells to memory T cells can be achieved through imputed CD45RA/CD45RO abundance, as CD45RA is a naïve antigen and CD45RO is a 185 186 memory antigen. Consistent with flow cytometry data, the large majority of CB T cells are naïve, 187 whereas the BM T cell population is more diverse ²⁵. Also, for BM B cells that have high imputed 188 CD19 levels, cTP-net allows us to confidently distinguish the Pre.B (CD38+, CD127+), immature B (CD38+, CD79b+), memory B (CD27+) and naïve B cells (CD27-), whose immunophenotypes 189 have been well characterized ²⁶. 190

In addition, consider natural killer cells, in which the proteins CD56 and CD16 serve as
indicators for immunostimulatory effector functions, including an efficient cytotoxic capacity ^{27, 28}.
We observe an opposing gradient of imputed CD56 and CD16 levels within transcriptomically
derived natural killer (NK) cell clusters that reveal CD56^{bright} and CD56^{dim} subsets, coherent with
previous studies³ (Figure 2f, Figure S10, F-test: p-value = 1.667e-15). This pattern is not found
in RNA abundances due to low expression (F-test: p-value = 0.9377). Between CD56^{bright} and

197 CD56^{dim} subsets, 7 out of 10 of previously studied differentially expressed genes are significant 198 in the single cell analysis (Fisher test: p-value = 1.07e-04) ^{3, 29, 30}. This gradient in CD56 and 199 CD16, where decrease in CD56 is accompanied by increase in CD16, is replicated across the 8 200 CBMC and 8 BMMC samples in HCA (Figure S8, S9, S10).

Consider also the case of CD57, which is a marker for terminally differentiated "senescent" cells
in the T and NK cell types. The imputed level of CD57 is lower in CBMCs (fetus's blood), and
rises in BMMCs (95% quantile: bootstrap p-value<1e-6). This is consistent with expectation
since CD57+ NK cell and T cell populations grow after birth and with ageing ³¹⁻³³ (Figure S8,
S9).

These results demonstrate how cTP-net, trained on a combination of PBMCs, CBMCs and BMMCs, can impute cell type, cell stage, and tissue-specific protein signatures in new data without explicitly being given the tissue of origin.

209

210 Application to Acute Myeloid Leukemia

We further apply cTP-net to an acute myeloid leukemia (AML) data set from Galen et al. ³⁴. AML 211 212 is a heterogeneous disease where the diversity of malignant cell types partially recapitulates the 213 stages of myeloid development. Mapping the malignant cells in AML to the differentiation stage 214 of their cell of origin strongly impacts tumor prognosis and treatment, as malignant cells that 215 originate from earlier stage progenitors have higher risk of relapse ^{35, 36}. In the original paper, 216 the authors sequenced 7698 cells from 5 healthy donors to build a reference map of cell types 217 during myeloid development, and then mapped 30712 cells from 16 AML patients across multiple time points to this reference to identify the differentiation stage of the malignant cells. 218 219 Here, by imputing 24 immunophenotype markers with cTP-net, we can directly characterize the 220 differentiation stage of cell-of-origin for the malignant cells.

221 Figure 4a is a UMAP plot based on imputed surface protein abundance of 5 normal BMs and 12 222 Day 0 samples from AML patients. The majority of the malignant cells as identified in the original paper reside on the right half of the plot, which recapitulate the myeloid differentiation 223 224 trajectory as revealed by the imputed values of canonical protein markers (Figure 4b): From 225 CD34+ progenitors to CD38+CD123+ cells in transition to CD11c+ and CD14+ mature monocytes ³⁷. All of the malignant cells have imputed protein values that place them along this 226 227 monocyte lineage. Using the transcriptome for visualization, on the other hand, reveals large 228 batch effects across samples, due to both technical batch and biological differences (Figure 229 S11). Thus, unlike the imputed protein data, the transcriptomic data cannot be directly 230 combined without alignment.

231 Based on the trajectory revealed by the imputed protein levels, we can determine the 232 differentiation cell stage(s) for the malignant cells of each tumor, according to which the 12 AML patients can be divided into three categories: (1) AMLs of single differentiation stage (AML420B, 233 AML556, AML707B and AML916; Figure 4c), (2) AMLs of two differentiation stages (AML210A, 234 235 AML328, AML419A and AML475; Figure 4e) and (3) AMLs of many differentiation stages (AML1012, AML329, AML870 and AML921A; Figure 4f). This stage assignment is consistent 236 237 with the original study ³⁴. For example, AML419A harbors two malignant cell types at opposite 238 ends of the monocyte differentiation axis, distinguished by imputed CD34 and CD11c levels as 239 CD34+CD11c- indicates progenitor-like and CD34-CD11c+ indicates differentiated monocyte-240 like cells (Figure 4d, 4e). AML707B, which carries a RUNX1/RUNX1T1 fusion, consists of cells 241 of a specific cell stage that is distinct from the normal myeloid trajectory (Figure 4c). Such 242 unique cell cluster was due to hyper CD38 level in surface protein prediction (Figure 4d). Such hyper-CD38 levels have been reported in AMLs with RUNX1/RUNX1T1 fusion³⁸⁻⁴⁰ and recent 243 studies have also shown that CD38 can be a potential target for adult AML^{41, 42}. 244

In this example, the imputed protein levels served as useful features for trajectory visualization.
This analysis also indicates that even though cTP-net is currently trained only on normal
immune cells, it can reveal disease-specific signatures in malignant cells and the imputed
protein levels are useful for characterizing tumor phenotypes.

249

250 Discussion

251 Taken together, our results demonstrate that cTP-net can leverage existing CITE-seq and 252 REAP-seq datasets to predict surface protein relative abundances for new scRNA-seq data 253 sets, and that the predictions generalize to cell types that are absent from, but related to those in the training data. cTP-net was benchmarked on PBMC and CBMC immune cells, showing 254 255 good generalization across tissues and technical protocols. On Human Atlas Data, we show 256 that the imputed surface protein levels allow easy assignment of cells to known cell types, as 257 well as the revealing of intra-cell type gradients. We then demonstrate that, even though cTPnet used only immune cells from healthy individuals for training, it is able to impute 258 259 immunophenotypes for malignant cells from acute myeloid leukemia, and that these 260 immunophenotypes allow placement of the cells along the myeloid differentiation trajectory. 261 Furthermore, we show that cTP-net is able to impute protein signatures in the malignant cells that are disease specific and that are not easily detectable from the transcriptomic counts. 262

SAVER-X serves an important role in the training procedure of cTP-net. As shown in Table S5, without SAVER-X denoising, the cTP-net prediction performance retracts by 0.02 in correlation, more significant than any other parameter tweaks. This discrepancy in performance is due to: (1) SAVER-X makes use of the noise model to obtain estimates of the true RNA counts. This helps cTP-net learn the underlining relationship between true RNA counts and protein level, rather than the noisy raw counts and protein levels, which varies more across data sets and thus does not generalize well. (2) By denoising the scRNA-seq, the input for learning the RNAprotein relationship is less sparse. Manifold learning on a more continuous input space usually
works better^{43, 44}. (3) Comparing to other autoencoder based denoising method, SAVER-X
performs Bayesian shrinkage on top of autoencoder framework to prevent over-imputation
(over-smoothing) ^{16, 45}.

Despite these promising results, cTP-net has limitations. (1) cTP-net can only apply to count based expression input (UMI-based). CITE-seq data with TPM and RPKM expression metric is not available for testing. Thus, the prediction accuracy is unknown. (2) The generalization ability of cTP-net to unrelated cell types has limitations. Even though the final cTP-net model, trained on immune cells, has good results on immune cells from diverse settings, we have not tried to perform imputation of these immune-related markers on cells that are not of the hematopoietic lineage.

With the accumulation of publicly available CITE-seq and REAP-seq data across diverse proteins, cell types and conditions, cTP-net can be retrained to accommodate more protein targets and improve in generalization accuracy. The possibility of such cross-omic transfer learning underscores the need for more diverse multi-omic cell atlases, and demonstrate how such resources can be used to enhance future studies. The cTP-net package is available both in Python and R at <u>https://github.com/zhouzilu/cTPnet</u>.

287

288 Acknowledgements

We would like to thank the National Institutes of Health for the award 5R01-HG006137 (for Z.Z., and N.Z.), award 1U2CCA233285-01 (for N.Z.), the National Science Foundation for the

award DMS-1562665 (to J.W., N.Z.), and the Wharton Dean's Fund for Post-doctoral

292 Research (to J.W.).

293 Author Contributions

- Z.Z. and N.Z. conceptualized the study and planned the case studies. Z.Z. designed the model,
- developed the algorithm, implemented the cTP-net software and led the data analysis. C.Y.
- helped in CITE-seq and REAP-seq data denoising and cell type labeling. J.W. helped with model
- design and Human Cell Atlas data analysis. Z.Z. and N.Z. wrote the paper with feedback from
- 298 C.Y. and J.W.

299 **Competing Financial Interests Statement**

- 300 The authors declare no competing interests
- 301

302 **References**

- 1. Stuart, T. & Satija, R. Integrative single-cell analysis. *Nat Rev Genet* **20**, 257-272 (2019).
- Peterson, V.M. et al. Multiplexed quantification of proteins and transcripts in single cells.
 Nat Biotechnol 35, 936-939 (2017).
- 306 3. Stoeckius, M. et al. Simultaneous epitope and transcriptome measurement in single cells. *Nat Methods* **14**, 865-868 (2017).
- Macaulay, I.C. et al. G&T-seq: parallel sequencing of single-cell genomes and transcriptomes. *Nat Methods* 12, 519-522 (2015).
- Wang, X. et al. Three-dimensional intact-tissue sequencing of single-cell transcriptional
 states. *Science* 361 (2018).
- 312 6. Regev, A. et al. The Human Cell Atlas. *Elife* **6** (2017).
- 7. Tirosh, I. et al. Dissecting the multicellular ecosystem of metastatic melanoma by singlecell RNA-seq. *Science* **352**, 189-196 (2016).
- 8. Villani, A.C. et al. Single-cell RNA-seq reveals new types of human blood dendritic cells,
 monocytes, and progenitors. *Science* 356 (2017).
- Liu, Y., Beyer, A. & Aebersold, R. On the Dependency of Cellular Protein Levels on
 mRNA Abundance. *Cell* 165, 535-550 (2016).
- Svensson, V. et al. Power analysis of single-cell RNA-sequencing experiments. *Nat Methods* 14, 381-387 (2017).
- 11. Zhao, B.S., Roundtree, I.A. & He, C. Post-transcriptional gene regulation by mRNA
 modifications. *Nat Rev Mol Cell Biol* 18, 31-42 (2017).

323	12.	Jackson, R.J., Hellen, C.U. & Pestova, T.V. The mechanism of eukaryotic translation
324		initiation and principles of its regulation. <i>Nat Rev Mol Cell Biol</i> 11 , 113-127 (2010).
325	13.	Mowen, K.A. & David, M. Unconventional post-translational modifications in
326		immunological signaling. Nat Immunol 15, 512-520 (2014).
327	14.	Schwartz, A.L. Cell biology of intracellular protein trafficking. Annu Rev Immunol 8, 195-
328		229 (1990).
329	15.	Roux, P.P. & Topisirovic, I. Signaling Pathways Involved in the Regulation of mRNA
330	10.	Translation. <i>Mol Cell Biol</i> 38 (2018).
	10	
331	16.	Wang, J. et al. Data denoising with transfer learning in single-cell transcriptomics. <i>Nat</i>
332	47	Methods 16 , 875-878 (2019).
333	17.	Webb, S. Deep learning for biology. Nature 554, 555-557 (2018).
334	18.	Tang, B., Pan, Z., Yin, K. & Khateeb, A. Recent Advances of Deep Learning in
335		Bioinformatics and Computational Biology. Front Genet 10, 214 (2019).
336	19.	Lopez, R., Regier, J., Cole, M.B., Jordan, M.I. & Yosef, N. Deep generative modeling for
337		single-cell transcriptomics. Nat Methods 15, 1053-1058 (2018).
338	20.	Stuart, T. et al. Comprehensive Integration of Single-Cell Data. Cell (2019).
339	21.	Martins, P.S. et al. Expression of cell surface receptors and oxidative metabolism
340		modulation in the clinical continuum of sepsis. <i>Crit Care</i> 12 , R25 (2008).
341	22.	Chen, L. & Flies, D.B. Molecular mechanisms of T cell co-stimulation and co-inhibition.
342	~~.	Nat Rev Immunol 13 , 227-242 (2013).
343	23.	Fromm, P. et al. CD16+Dendritic Cells Are a Unique Myeloid Antigen Presenting Cell
	23.	
344	0.4	Population. <i>Blood</i> 128 (2016).
345	24.	Stuart, T. et al. Comprehensive Integration of Single-Cell Data. <i>Cell</i> 177 , 1888-+ (2019).
346	25.	D'Arena, G. et al. Flow cytometric characterization of human umbilical cord blood
347		lymphocytes: immunophenotypic features. Haematologica 83, 197-203 (1998).
348	26.	Clavarino, G. et al. Novel Strategy for Phenotypic Characterization of Human B
349		Lymphocytes from Precursors to Effector Cells by Flow Cytometry. <i>Plos One</i> 11 (2016).
350	27.	Van Acker, H.H., Capsomidis, A., Smits, E.L. & Van Tendeloo, V.F. CD56 in the Immune
351		System: More Than a Marker for Cytotoxicity? Front Immunol 8, 892 (2017).
352	28.	Tsukerman, P. et al. Expansion of CD16 positive and negative human NK cells in
353		response to tumor stimulation. Eur J Immunol 44, 1517-1525 (2014).
354	29.	Poli, A. et al. CD56(bright) natural killer (NK) cells: an important NK cell subset.
355	20.	Immunology 126 , 458-465 (2009).
356	30.	Wendt, K. et al. Gene and protein characteristics reflect functional diversity of CD56(dim)
357	50.	and CD56(bright) NK cells. J Leukocyte Biol 80 , 1529-1541 (2006).
	21	
358	31.	d'Angeac, A.D. et al. CD57+ T lymphocytes are derived from CD57- precursors by
359	00	differentiation occurring in late immune responses. Eur J Immunol 24, 1503-1511 (1994).
360	32.	Musha, N. et al. Expansion of CD56+ NK T and gamma delta T cells from cord blood of
361		human neonates. Clin Exp Immunol 113, 220-228 (1998).
362	33.	Dalle, J.H. et al. Characterization of cord blood natural killer cells: implications for
363		transplantation and neonatal infections. <i>Pediatr Res</i> 57, 649-655 (2005).
364	34.	van Galen, P. et al. Single-Cell RNA-Seq Reveals AML Hierarchies Relevant to Disease
365		Progression and Immunity. Cell 176, 1265-+ (2019).
366	35.	Pollyea, D.A. & Jordan, C.T. Therapeutic targeting of acute myeloid leukemia stem cells.
367	-	<i>Blood</i> 129 , 1627-1635 (2017).
368	36.	McKenzie, M.D. et al. Interconversion between Tumorigenic and Differentiated States in
369	00.	Acute Myeloid Leukemia. <i>Cell Stem Cell</i> 25 , 258-+ (2019).
370	37.	Geissmann, F. et al. Development of Monocytes, Macrophages, and Dendritic Cells.
370 371	57.	Science 327, 656-661 (2010).
571		

372	38.	Jang, J.H. et al. Acute myeloid leukemia with del(X)(p21) and cryptic RUNX1/RUNX1T1
373		from ins(8;21)(q22;q22q22) revealed by atypical FISH signals. Ann Clin Lab Sci 40, 80-
374		84 (2010).
375	39.	Moroi, K. & Sato, T. Comparison between procaine and isocarboxazid metabolism in
376		vitro by a liver microsomal amidase-esterase. Biochem Pharmacol 24, 1517-1521
377		(1975).
378	40.	Shang, L. et al. The immunophenotypic characteristics and flow cytometric scoring
379		system of acute myeloid leukemia with t(8;21) (q22;q22); RUNX1-RUNX1T1. Int J Lab
380		<i>Hematol</i> 41 , 23-31 (2019).
381	41.	Naik, J. et al. CD38 as a therapeutic target for adult acute myeloid leukemia and T-cell
382		acute lymphoblastic leukemia. <i>Haematologica</i> 104 , E100-E103 (2019).
383	42.	Eveillard, M. et al. CD38 Expression in B-Lineage Acute Lymphoblastic Leukemia, a
384		Possible Target for Immunotherapy. <i>Blood</i> 128 (2016).
385	43.	An, G.Z. The effects of adding noise during backpropagation training on a generalization
386		performance. Neural Comput 8, 643-674 (1996).
387	44.	Reed, R. & MarksII, R.J. Neural smithing: supervised learning in feedforward artificial
388		neural networks. (Mit Press, 1999).
389	45.	Andrews, T.S. & Hemberg, M. False signals induced by single-cell imputation.
390		<i>F1000Res</i> 7 , 1740 (2018).
391	46.	LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. Nature 521, 436-444 (2015).
392	47.	Kingma, D. & Ba, J. Adam: a method for stochastic optimization (2014). arXiv preprint
393		arXiv:1412.6980 15 (2015).
394	48.	Subramanian, A. et al. Gene set enrichment analysis: a knowledge-based approach for
395		interpreting genome-wide expression profiles. Proc Natl Acad Sci U S A 102, 15545-
396		15550 (2005).

399 Figure legends

- 400 Figure 1. cTP-net analysis pipeline and imputation of example proteins.
- 401 (a) Overview of cTP-net analysis pipeline, which learns a mapping from the denoised scRNA-
- seq data to the relative abundance of surface proteins, capturing multi-gene features that reflect
- 403 the cellular environment and related processes. (b) For three example proteins, cross-cell
- 404 scatter and correlation of CITE-seq measured abundances vs. (1) raw RNA count, (2) SAVER-X
- 405 denoised RNA level, and (3) cTP-net predicted protein abundance.
- 406 Figure 2. Benchmark evaluation on CITE-seq PBMC data

407 (a) Benchmark evaluation of cTP-net on CITE-seq PBMC data, with comparisons to Seurat v3,

408 in validation, across cell type, across tissue and across technology scenarios. The table on the

409 left shows the training scheme of each test, the heatmap shows correlations with actual

410 measured protein abundances. (b) Within cell type correlations between imputed and measured

411 protein abundance on the CITE-seq PBMC data, Seurat v3 versus cTP-net. Each point (color

and shape pair) indicates a cell type and surface protein pair, where the x-axis is correlation

413 between actual measured abundance and Seurat imputation and y-axis is the correlation

between actual measured abundance and cTP-net imputation. (c) Scatter of imputed versus

415 measured abundance for the three (surface protein, cell type) pairs marked by arrows in (b):

416 CD11c in CD14-CD16+ monocytes, CD2 in CD8 T cells, and CD19 in dendritic cells.

417 Figure 3. Imputation results analysis on Human Cell Atlas data sets.

(a) Left panel: UMAP visualization of MantonBM1 BMMCs T cell subpopulation based on RNA
expression, colored by cell type. CD4 T: mature CD4+ T cells; mature CD8 T: CD8+ T cells;
naïve CD4 T: naïve CD4+ T cells; naïve CD8 T: naïve CD8+ T cells; CD8 senescent T: CD8+
senescent T cells. Right panel: Related imputed protein abundance and RNA expression of its

422 corresponding gene. (b) UMAP visualization of MantonBM1 BMMCs based on RNA expression. 423 colored by cell type. B: B cells; CD4 T: CD4+ T cells; CD8 T: CD8+ T cells; cMono: classical monocyte; ncMono: non-classical monocyte; NK: natural killer cells; Pre.: precursors; Plasma: 424 425 plasma cells. (c) Left panel: UMAP visualization of MantonBM1 BMMCs B cell subpopulation 426 based on RNA expression, colored by cell type. Pre.B: B cell precursors; immature B: immature 427 B cells; memory B: memory B cells; naïve B: naïve B cells. Right panel: Related imputed protein 428 abundance and RNA expression of its corresponding gene. (d) UMAP visualization of 429 MantonCB2 CBMCs based on RNA expression, colored by cell type, (e) cTP-net imputed 430 protein abundance and RNA read count of its corresponding gene for 24 surface proteins. (f) UMAP visualization of MantonCB2 CBMCs NK cell subpopulation colored by CD56 and CD16 431 432 imputed protein abundance and RNA read count. Reverse gradient is observed in cTP-net 433 prediction but not in the read count for its corresponding RNA. (f) Contour plot of cells based on 434 imputed CD56 and CD16 abundance in NK cell populations. Strong negative correlation with 435 two subpopulation observed.

436 Figure 4. Imputation results analysis on Acute Myeloid Leukemia data sets.

(a) UMAP visualization of normal cells and malignant cells from 12 AML samples at Day0 based
on imputed protein abundance (red: malignant cells; grey: normal cells). (b) UMAP visualization
of the myeloid trajectory. cTP-net imputed protein abundance of markers that perfectly
recapitulate the myeloid development. (c, e, f) UMAP visualization of the myeloid trajectory with
corresponding malignant cells from AML sample highlighted. (d) Scatter plot of normal and AML
malignant cells based on imputed protein expression.

444 **Online Methods**

445 Data sets and pre-processing

Table S1 summarizes the five data sets analyzed in this study: CITE-PBMC, CITE-CBMC, 446 447 REAP-PBMC, HCA-CBMC and HCA-BMMC. Among these, CITE-PBMC, CITE-CBMC and 448 REAP-PBMC have paired scRNA-seq and surface protein counts, while HCA-CBMC and HCA-BMMC have only scRNA-seq counts. For all scRNA-seq data sets, low quality gene (< 10 449 450 counts across cells) and low-quality cells (less than 200 genes detected) are removed, and the 451 count matrix (C) for all remaining cells and genes is used as input for denoising. scRNA data denoising was performed with SAVER-X using default parameters. Denoised counts (Λ) were 452 453 further transformed with Seurat default LogNormalize function,

454
$$X_{ij} = \log\left(\frac{\Lambda_{ij} * 10,000}{m_j}\right)$$

where Λ_{ij} is the denoised molecule count of gene *i* in cell *j*, and *m_j* is the sum of all molecule counts of cell *j*. The normalized denoised count matrix *X* is the training input for the subsequent multiple branch neural network. For the surface protein counts, we adopted the relative abundance transformation from Stoeckius et al.³. For each cell *c*,

459
$$y_c = \left[\ln\left(\frac{p_{1c}}{g(p_c)}\right), \ln\left(\frac{p_{2c}}{g(p_c)}\right) \dots \ln\left(\frac{p_{dc}}{g(p_c)}\right) \right]$$

where p_c is vector of antibody-derived tags (ADT) counts, and $g(p_c)$) is the geometric mean of p_c . The network trained using this transformed relative protein abundance as the response vector yields better prediction accuracy than the network trained using raw protein barcode counts.

465 **cTP-net neural network structure and training parameters**

466 Figure S1 shows the structure of cTP-net. Here, we have a normalized expression matrix X of N cells and D genes, and a normalized protein abundance matrix Y of the same N cells and d 467 surface proteins. Let's denote cTP-net as a function *F* that maps from \mathbb{R}^{D} to \mathbb{R}^{d} . Starting from 468 the input layer, with dimension equals to number of genes D, the first internal layer has 469 dimension 1000, followed by a second internal layer with dimension 128. These two layers are 470 471 designed to learn and encode features that are shared across proteins, such as features that 472 are informative for cell type, cell state and common processes such as cell cycle. The remaining layers are protein specific, with 64 nodes for each protein that feed into a one node output layer 473 giving the imputed value. All layers except the last layer are fully connected (FC) with rectified 474 475 linear unit (ReLU) activation function ⁴⁶, while the last layer is a fully connected layer with 476 identity activation function for output. The objective function here is,

477
$$\operatorname{argmin}_{F} |Y - F(X)|_{1}$$

where the loss is L1 norm. The objective function was optimized stochastically with Adam ⁴⁷ with
learning rate set to 10e-5 for 139 epochs (cross-validation). Other variations of cTP-net, which
we found to have inferior performance, are illustrated in more details in Supplementary Note.

481

485

482 Benchmarking procedure

Validation set testing procedure. Figure S2a shows the validation set testing procedure. Given
limited amount of data, we keep only 10% of the cells as the testing set, and use the other 90%

486 *Out-of-cell type prediction procedure.* We perform the out-of-cell type prediction based on

of the cells for training. The optimal model was selected based on the testing error.

487 Figure S2b. This procedure mimics cross-validation, except that, instead of selecting the test set

488 cells randomly, we partition the cells by their cell types. Iteratively, we designate all cells of a 489 given cell type for testing and use the remaining cells for training. We then perform prediction on the hold-out cell type using the model trained on all other cell types. In the end, every cell has 490 been tested once and has the corresponding predictions. In the benchmark against the 491 492 validation set testing procedure, we limit comparisons to the same cells that were in the 493 validation set in the holdout scheme to account for variations between subsets. 494 Cell population and technology transfer learning procedure. To apply the models we trained in 495 validation set testing procedure to different cell populations and technologies, the inputs have to 496 be in the same feature space. Even though all data sets considered are from human cells, the

497 list of genes differs between experiments and technologies. Genes that are in the training data

but not in the testing data are filled with zeros. Because cTP-net utilizes overrepresented

499 number of genes to predict the surface proteins level, having a small number of genes missing

500 has little effect on the performance. After prediction, we selected only the shared proteins

501 between two data sets for comparison.

502

503 **cTP-net interpolation**

To better interpret the relationships that the neural network is learning, we developed a 504 permutation-based interpolation scheme that can calculate an influence score epi for each gene 505 506 in the imputation of each protein (Figure S6). The idea is to assess how much changing the 507 expression value of certain genes in the training data affects the training errors for a given 508 model F. In each epoch, we interpolate all of the genes in a stochastic manner. Let's denote X as the expression matrix (N by G matrix, where N is the number of cells and G is number of 509 genes), Y as protein abundance matrix and L as the loss function. The algorithm goes as follow 510 511 (Figure S6):

512 (1) Estimate the original model error $e^{orig} = L(Y, F(X))$.

513 (2) Sampling batch of genes denote by qs. Generate expression matrix X^{perm} by permuting 514 genes in *gs* in the data *X*. This breaks the association between *gs* and protein abundance Y, i.e. the cell order within gs does not coordinate with protein abundance Y. 515 (3) Estimate error $e^{perm} = L(Y, F(X^{perm}))$ based on the predictions of the permuted data. 516 (4) Calculate permutation feature importance $\Delta_{gs} = |\epsilon^{orig} - \epsilon^{perm}|$ of gene set *gs* to this 517 518 model F. 519 We set batch size as 100 with 500 epochs. Furthermore, by picking different cells to interpolate, 520 we could identify gene influence score in different cell types. For example, if matrix X belongs to

a given cell type, the cell type specific genes are consistent across cells of the given cell type,

and thus, the permutation will not influence these genes. Genes that influence the surface

523 protein abundance within the cell type, such as cell cycle genes and protein synthesis genes,

tend to be rewarded with high influence scores in such a cell-type specific interpolation analysis.

525 For the top 100 highest influence scored genes from the following scenarios in CITE-PBMC: (1)

526 CD45RA in CD14-CD16+ monocytes, (2) CD11c in CD14-CD16+ monocytes, (3) CD45RA in

527 CD8 T cells, (4) CD45RA in CD4 T cells, (5) CD11c in CD14+CD16+ monocytes, (6) CD45RA

528 in dendritic cells, and (7) CD11c in dendritic cells, we employed a Gene Ontology analysis ⁴⁸

which identify top 10 pathways based on GO gene sets with FDR q-value < 0.05 as significant(Table S4).

531

532 Seurat anchor-transfer analysis

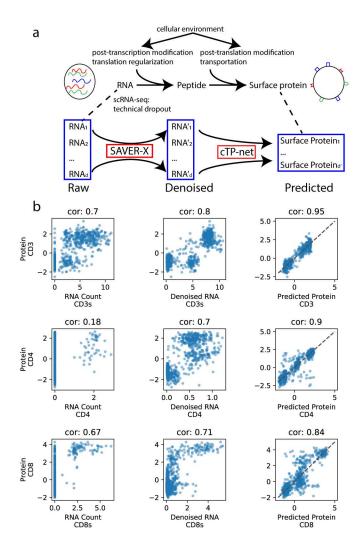
We compared cTP-net with an anchor-based transfer learning method developed in Seurat v3
 For Seurat v3, RNA count data are normalized by LogNormalization method, while surface
 protein counts are normalized by centered log-ratio (CLR) method. In validation test setting, we

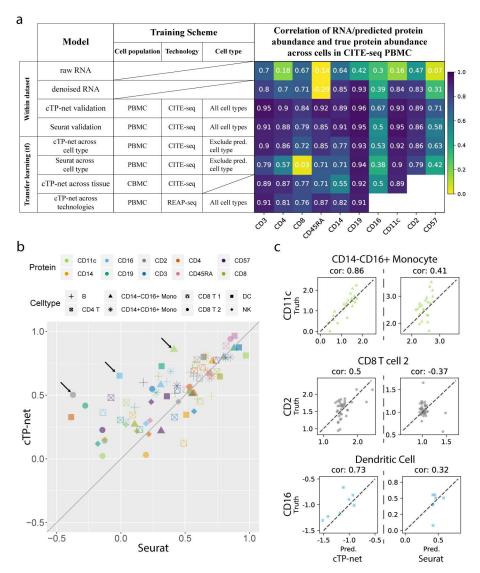
used the same cells for training and testing as in cTP-net so as to be directly comparable to
cTP-net. For out-of-cell type prediction, default parameters did not work for several cell types in
anchor-transfer step, because, for those cell types, there are few anchors shared between the
training and testing sets. To overcome this, we reduced the number of anchors iteratively until
the function ran successfully.

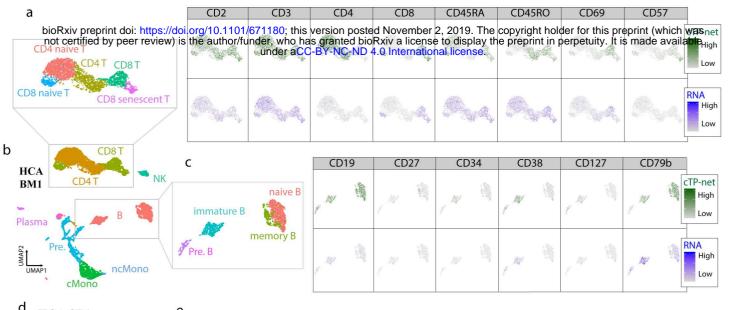
541

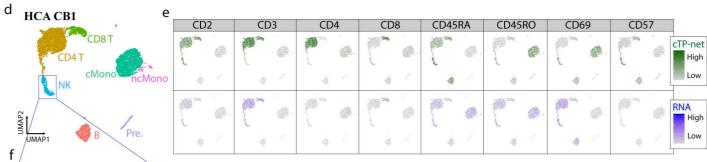
542 HCA data analysis

- 543 HCA RNA-seq transcriptome data analysis. HCA RNA-seq data sets are processed as
- discussed above, resulting in log-normalized denoised values. We applied default pipeline of
- 545 Seurat and generated t-SNE plot for both data sets (Figure S7). Cells are clearly clustered by
- 546 individuals, indicating strong batch effects. As a result, the following analysis was performed on
- 547 cells of each individual. Major cell types were determined by known gene markers.
- 548 Surface protein prediction by cTP-net. From the log-normalized denoised expression value, we
- 549 predict the surface protein abundance with cTP-net model trained jointly on CITE-seq PBMC
- and CBMC data sets. We embedded 12 surface protein abundance across 16 individuals on t-
- 551 SNE plot, showing consistent results with cell type information (Figure S8, S9).





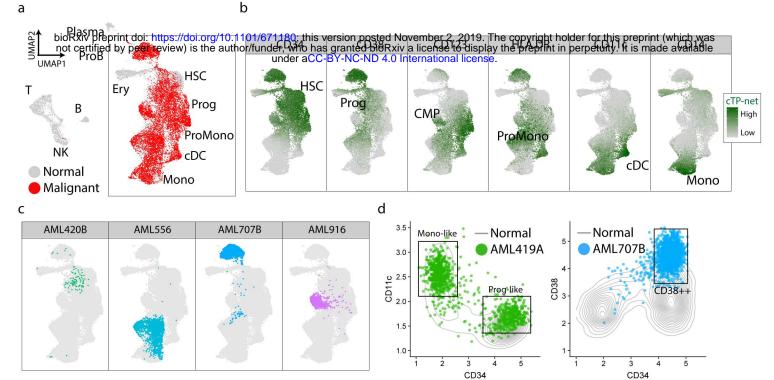




CD56	CD16	\square	CD5	6	CD	16	CD	19	CD2	.7	CD3	4	CD3	8	CD1	27	CD7	9b	
			T		-	ø.		Ø.	1	Ø.		Ø.		Ø	7	Ø.	E	Ø	cTP-net High
			, in the second se	ŀ	ė	1	٩	1	ġ	1	@	Ţ	â	1	ġ	1	۹	1	Low
				ø.		Ø.,	7	ø	2	0		Ø.	P	Ø.		Ø.		Ø.	RNA High
	Stor Star	1	-		-	2	ġ	I.	ģ	1	ģ	1	4	Ţ	ġ	ŀ,	4	1	Low

g	NK cells
3.5 -	Spearman cor. -0.47
3.0 -	High
9100 2.0 -	
U 2.0-	Low
1.5 -	
1.0 -	² ³ ⁴ CD56
	CD56

CD1	1a	CD1	1c	CD	14	CD1	23	CD2	8	CD16	51	CD2	78	HLA.	DR	
7	Ø.		Ø.		Ø.		Ø.		Ø.	P	Ø.		Ø.		Ø.	cTP-net High
ġ	Ž	-	7	0		Ø	Ζ	4	Z		Z	ė			7	Low
	Ø.	P.	Ø.		Ø		Ø.	C	Ø	P	Ø.		Ø.	2	Ø.	RNA High
4	2	4		à		à		á		à	2	\$				Low



1	٦	
r	-	
Ľ	-	

AML210A	AML328	AML419A	AML475	AML1012	AML329	AML870	AML921A

f