**Supplementary Table 1**

| CC strain | Age | Initial weight | Treatment | Data | Weight loss (%) |
|---|---|---|---|---|---|
| 111A | 8.1 | 18.2 | Influenza infection | RNA-seq | 4.4 |
| 111A | 8.1 | 19.5 | Influenza infection | RNA-seq | 1.5 |
| 1488A | 9.6 | 23.4 | Influenza infection | RNA-seq | 4.7 |
| 1912A | 9.6 | 24 | Influenza infection | RNA-seq | 1.7 |
| 1912A | 9.6 | 22.7 | Influenza infection | RNA-seq | 2.6 |
| 2126A | 8.9 | 12.8 | Influenza infection | RNA-seq | 3.1 |
| 2126A | 8.9 | 14.9 | Influenza infection | RNA-seq | -0.3 |
| 21B | 8.1 | 15.5 | Influenza infection | RNA-seq | 1.9 |
| 21B | 8.1 | 15.1 | Influenza infection | RNA-seq | 0.0 |
| 2513A | 9.6 | 15.2 | Influenza infection | RNA-seq | 3.9 |
| 2513A | 9.6 | 18 | Influenza infection | RNA-seq | 1.7 |
| 2750A | 9.1 | 24.2 | Influenza infection | RNA-seq | 0.8 |
| 2750A | 9.1 | 20.1 | Influenza infection | RNA-seq | 2.5 |
| 3348A | 9.1 | 18.6 | Influenza infection | RNA-seq | 11.3 |
| 3348A | 9.1 | 18.5 | Influenza infection | RNA-seq | 6.5 |
| 3912A | 8.1 | 19 | Influenza infection | RNA-seq | 4.2 |
| 3912A | 8.1 | 20.5 | Influenza infection | RNA-seq | 4.4 |
| 4438A | 8.4 | 24.5 | Influenza infection | RNA-seq | 4.1 |
| 4438A | 8.4 | 20.6 | Influenza infection | RNA-seq | 4.4 |
| 5000A | 9.3 | 23.9 | Influenza infection | RNA-seq | 6.3 |
| 5000A | 9.3 | 22.3 | Influenza infection | RNA-seq | 9.0 |
| 5001A | 8.4 | 23.2 | Influenza infection | RNA-seq | 0.9 |
| 5003A | 9.7 | 18 | Influenza infection | RNA-seq | 1.1 |
| 5003A | 9.7 | 17.2 | Influenza infection | RNA-seq | -0.9 |
| 5004A | 8.7 | 19.6 | Influenza infection | RNA-seq | 2.0 |
| 5004A | 8.7 | 18 | Influenza infection | RNA-seq | -0.6 |
| 5010A | 8.3 | 19.7 | Influenza infection | RNA-seq | 8.6 |
| 5010A | 8.3 | 19.8 | Influenza infection | RNA-seq | 4.5 |
| 5021A | 8.9 | 19.6 | Influenza infection | RNA-seq | 5.1 |
| 5021A | 8.9 | 18.5 | Influenza infection | RNA-seq | 6.5 |
| 5022A | 8.7 | 19.7 | Influenza infection | RNA-seq | 2.0 |
| 5022A | 8.7 | 20.2 | Influenza infection | RNA-seq | 4.0 |
| 5023A | 9.4 | 19.6 | Influenza infection | RNA-seq | 5.6 |
| 5023A | 9.4 | 17 | Influenza infection | RNA-seq | 1.8 |
| 57B | 8.7 | 23.6 | Influenza infection | RNA-seq | 4.7 |
| 57B | 8.7 | 23.8 | Influenza infection | RNA-seq | 1.7 |
| 72A | 9.0 | 20.1 | Influenza infection | RNA-seq | 6.0 |
| 72A | 9.0 | 21.4 | Influenza infection | RNA-seq | 5.6 |
| 111A | 8.1 | 17.1 | PBS | RNA-seq | 0.6 |
| 1912A | 9.7 | 22.1 | PBS | RNA-seq | -0.7 |
| 2126A | 8.9 | 14.3 | PBS | RNA-seq | 2.1 |
| 2126A | 8.9 | 15 | PBS | RNA-seq | 0.7 |
| 21B | 8.9 | 18.8 | PBS | RNA-seq | 1.6 |
| 21B | 8.9 | 18.2 | PBS | RNA-seq | -0.5 |
| 21B | 8.1 | 22.1 | PBS | RNA-seq | -0.9 |
| 2513A | 9.6 | 16.7 | PBS | RNA-seq | 0.6 |
| 2513A | 9.6 | 16.2 | PBS | RNA-seq | -0.9 |
| 2750A | 9.0 | 19.9 | PBS | RNA-seq | 3.0 |
| 2750A | 9.0 | 24.7 | PBS | RNA-seq | 6.1 |
| 3912A | 8.1 | 19.5 | PBS | RNA-seq | 1.5 |
| 3912A | 8.1 | 19.6 | PBS | RNA-seq | 3.6 |
| 4438A | 8.3 | 18.4 | PBS | RNA-seq | -0.1 |
| 4438A | 8.3 | 17.5 | PBS | RNA-seq | -0.1 |
| 5000A | 9.3 | 19.9 | PBS | RNA-seq | 8.5 |
| 5000A | 9.3 | 22.5 | PBS | RNA-seq | 5.8 |
| 5001A | 8.4 | 19.2 | PBS | RNA-seq | 1.0 |
| 5001A | 8.4 | 19.5 | PBS | RNA-seq | 1.0 |
| 5003A | 9.7 | 18.7 | PBS | RNA-seq | -0.1 |
| 5003A | 9.7 | 18.2 | PBS | RNA-seq | 3.3 |
| 3348A | 9.1 | 19.3 | PBS | RNA-seq | 2.1 |
| 3348A | 9.1 | 18.4 | PBS | RNA-seq | 0.5 |

| CC strain | Age | Initial weight | Treatment | Data | Weight loss (%) |
|---|---|---|---|---|---|
| 5004A | 8.7 | 20.5 | PBS | RNA-seq | 0.5 |
| 5010A | 8.3 | 19.1 | PBS | RNA-seq | 3.1 |
| 5021A | 8.9 | 18.8 | PBS | RNA-seq | 1.1 |
| 5022A | 9.3 | 21.9 | PBS | RNA-seq | 1.8 |
| 5022A | 9.3 | 19.7 | PBS | RNA-seq | 1.5 |
| 5022A | 8.7 | 16.7 | PBS | RNA-seq | 2.4 |
| 5023A | 9.3 | 18.7 | PBS | RNA-seq | 1.1 |
| 57B | 8.1 | 22.6 | PBS | RNA-seq | 0.0 |
| 57B | 8.1 | 19 | PBS | RNA-seq | 0.0 |
| 72A | 9.0 | 22.1 | PBS | RNA-seq | 6.8 |
| 72A | 9.0 | 19.7 | PBS | RNA-seq | 3.0 |
| 2126A | 8.4 | 17.7 | Untreated | RNA-seq | Irrelevant |
| 111A | 9.1 | 17.6 | Untreated | RNA-seq | Irrelevant |
| 1912A | 8.7 | 21 | Untreated | RNA-seq | Irrelevant |
| 2750A | 8.9 | 19.6 | Untreated | RNA-seq | Irrelevant |
| 5023A | 8.1 | 14.5 | Untreated | RNA-seq | Irrelevant |
| C57BL/6J* | 7.6 | 18.4 | Untreated | RNA-seq | Irrelevant |
| 72A | 7.9 | 18.6 | Influenza infection | FACS | 3.8 |
| 72A | 9.4 | 19.5 | Influenza infection | FACS | 11.8 |
| 72A | 9.4 | 20.7 | Influenza infection | FACS | 9.7 |
| 3912A | 8.9 | 22.7 | Influenza infection | FACS | -0.9 |
| 3912A | 8.9 | 20.2 | Influenza infection | FACS | 3.5 |
| 4438A | 9.0 | 20.5 | Influenza infection | FACS | 8.8 |
| 4438A | 9.9 | 23.5 | Influenza infection | FACS | -0.9 |
| 5001A | 9.1 | 18 | Influenza infection | FACS | 9.4 |
| 5001A | 9.3 | 19.7 | Influenza infection | FACS | 5.6 |
| 5023A | 8.6 | 15 | Influenza infection | FACS | 5.3 |
| 5023A | 8.6 | 14.8 | Influenza infection | FACS | 4.1 |
| 72A | 7.9 | 19.5 | PBS | FACS | 0.5 |
| 3912A | 8.9 | 20.4 | PBS | FACS | 0.0 |
| 4438A | 9.1 | 26.9 | PBS | FACS | 5.2 |
| 5001A | 9.1 | 20 | PBS | FACS | 1.5 |
| 5023A | 8.6 | 14.3 | PBS | FACS | 3.5 |

**Table S1**. **Description of the animals under study**. Shown are the Collaborative Cross (CC) strain identifiers (Column A), mouse age (weeks, Column B) and weight (in grams; Column C) of each mouse before the treatment. The experimental procedure—treatment applied and subsequent measurements—are indicated in Columns D and E, respectively. The clinical outcome, namely the percentage of body weight loss, is indicated in Columns F. * The C57BL/6J strain is one of the CC founder strains.

# SUPPLEMENTARY NOTE 1

## The contribution of the CPM algorithm

CPM, the algorithm developed in this study, utilizes the SVR deconvolution approach, since (i) SVR accuracy has been previously demonstrated within the Cibersort framework[1], and (ii) SVR can be applied on both absolute and relative bulk genomics data (unlike alternative approaches, such as the DCQ framework[2]). A substantial challenge in the application of deconvolution methods to single-cell reference data is in maintaining robustness with a large number of reference profiles. Although deconvolution methods are generally robust and practical with a relatively small number of reference profiles, scaling to a large reference collection is problematic: most methods can handle only a few dozen of reference profiles, while DCQ, which is scalable to 200 profiles, was designed specifically to relative bulk data[3]. Another challenge is the biased representation of cell states within the reference data. Here we make two contributions that make the CPM methodology accurate and scalable: (1) To solve the scaling issue, we applied SVR multiple times, each time on a different subset of the reference profiles, and then aggregated the inferred values into a final abundance prediction for each reference profile. This avoids the need to utilize many reference profiles simultaneously and to improve robustness with bootstrapping. (2) To avoid biases due to imbalance of reference single cells over the cell-state space, we ensured that the cell subsets used in the SVR are uniformly distributed over the cell-state space.

## Synthetic data generation and the accuracy score

A single 'synthetic data collection' consisted of 100 input bulk profiles, where each profile is generated by mimicking the heterogeneity of cells within a biological complex tissue. Gene expression of a gene $j$ in a synthetic bulk profile $k$, denoted $z_{jk}$, includes a mix of isolated single cells:

$$z_{jk} = \sum_{i=1..L} f_{ik} \cdot b_{ij} + \varepsilon_{ik} \quad [1]$$

where $b_{ij}$ denotes the gene expression value of gene $j$ in the reference single cell $i$, $L$ is the number of single cells in the reference data, and $f_{ik}$ is the fraction of reference cell $i$ in sample $k$, formalized as:

$$f_{ik} = \frac{c_{ik}}{\sum_{i=1..L} c_{ik}} + \eta_{ik} \quad [2]$$

where $c_{ik}$ is the quantity of cell $i$ in sample $k$. We further introduce noise level $\varepsilon_{ik}$ and $\eta_{ik}$ to the values of $z_{jk}$ and $f_{ik}$, respectively, by sampling from a normal distribution with a zero mean and a variance that is proportional (by a factor $\gamma_g$) to the expected value (namely, the average of $z_{jk}$ or $f_{ik}$ across the 100 profiles). The level of the '*expression noise*' in our synthetic data is therefore determined by the proportion factor $\gamma_g$. To avoid biases due to imbalance of cell densities, all synthetic bulk profiles were generated as a mixture of pre-selected

*L*=378 single cells (derived from the infected mouse at 2 days p.i.) that were uniformly distributed over the cell-state space.

Next, using our basic simulation setting (Eqs. 1-2), we generated three types of simulations. All three simulation types rely on the partition of single cells into nine cell types, and mainly differ in their cell quantity ($c_{ik}$) values. In the first simulation type—'cell-type simulation'—cells of different types attained different quantities, whereas all cells within each cell type attained the same cell quantity value. The two other simulations, in contrast, are focused on intra-cell-type heterogeneity: either an intra-cell-type gradual change in the quantity of cells throughout a certain trajectory of cell states (the 'gradual-change simulation'), or alternatively, an intra-cell-type changes in quantities of a selected cell subpopulation (the 'cell-subtype simulation').

Specifically, in the cell-type simulations, the quantity of each cell *i* within selected cell types was set to $c_{ik}$=1+ $e_c$ ($e_c$ is denoted the '*effect size*' of the cell type simulation), whereas the quantity of the remaining cells was set to $c_{ik}$=1. In the case of the gradual-change simulation, for each cell type, the position of single cells along the activation-state trajectory were used as the cell-state space (see details in the 'reference single-cell data' section). Cell quantities were generated based on this cell-state space in four steps. In the first step, we select the *Kc* cell types within which we simulate gradual change in cell quantities. Among these selected cell types, we further randomly selected *Kr* cell types whose activation trajectory was reversed. The second step introduced a normally-distributed noise in the positions of all cells with a variance that is proportional to the standard deviation of cell positions, using a proportion factor $\gamma_p$. The level of the '*cell space noise*' in our synthetic data is therefore determined by $\gamma_p$. In the third step, we standardized the positions from the second step to lie in [0,1]. In the fourth step, cell quantities ($c_{ik}$) for the selected cell types were calculated using an exponential function $f(x_g) = b^{x_g}$ in which *b* is the positions from the third step and the exponent $x_g$ is the 'effect size' of the gradual-change simulation. The value of $c_{ik}$ in all remaining (unchanged) cells in all samples was the mean standardized score of their trajectory.

For the cell-subtype simulation, for each cell type, we used the first two principle components as the cell-state space, and further introduced a cell space noise to the positions of cells within this space, as described above. Cell quantities within each given cell type were generated as follows: we first selected a single reference cell and then used its neighboring cells, assuming an Euclidian cell space, as the selected cell subtype (using a certain 'cell subpopulation size' parameter). Next, the quantity of each cell *i* within the selected subset was set to $c_{ik}$=1+ $e_s$ where $e_s$ is the magnitude of cell quantity changes, denoted the 'effect size' of the cell-subtype simulation. The default value of $c_{ik}$ in all remaining (unchanged) cells in all samples was set to 1.

Whereas the data generation above refers to absolute expression values, we also generated relative-expression synthetic data collections. This was done by generating additional 100 control synthetic profiles for each synthetic collection (in which $c_{ik}$ is always the mean trajectory standardized score or 1 for either gradual-change simulation and cell-subtype simulation, respectively), and then calculating the differential expression values of each synthetic data profile.

For each simulation setting we generated a single synthetic data collection and analyzed the performance of each method on this collection. In particular, for each synthetic profile collection, an 'accuracy' metric was calculated as the Pearson correlation between the actual fractions $f_{ik}$ and predicted fractions across all $L$ reference single cells and 100 bulk profiles. Overall we generated 1671 synthetic data collection, for three simulations settings (cell-type, cell-subtype and gradual-change simulations), seven different levels of expression noise ($\gamma_g$ ranges from 0.001 to 0.9), seven levels of cell space noise ($\gamma_p$ ranges from 0.1 to 3), seven cell subtypes in the cell-subtype simulations (fractions ranging between 0.1 to 0.7 of cells) assuming either absolute or relative bulk data, nine combinations of $Kc$ and $Kr$, and 5 effect size levels in the gradual-change simulation. The effect sizes of the cell type and cell subtype simulations ($e_c$ and $e_s$) were fixed to 0.5 as the range of effects parallels the abovementioned changes in noise factors. In the cell-type simulation, the effect was added to cells within six arbitrarily-selected cell types. Our default set of parameter is: cell space noise = 0.5, expression noise = 0.1, cell subpopulation size = 0.3, $Kc$ = 6 and $Kr$=0, assuming relative bulk data, and effect size $x_g$= 1 (corresponding to a linear function). In all cases, we report only the particular parameters that are different from this default setting.

## Generating reference data for the compared methods

Since the compared deconvolution methods rely on a relatively small number of input reference profiles, the reference data was constructed using the input scRNA-seq profiles. In particular, we implemented the methodology that was previously applied in scRNA-seq-based deconvolution studies[4,5]: single cells of each cell type were partitioned into cell groups using DBSCAN clustering (using ε=3 and minPts=15; as previously described[6]). Then, the center of each such group was used as a reference profile; the 'mean center' (namely, the averaged cell profiles) of each cell group was used as its center. We refer to this approach as the 'DBSCAN + mean center method'. To generalize the reference-construction approach to a pre-selected size of reference dataset, we clustered the cells within each cell type using $K$-means clustering (instead of DBSCAN) and then identified the center of each cell group (the clustering relies on the cell-state space). $K$ therefore specifies the number of single cell groups within each cell type and is referred to as the 'level of granularity'. We applied three alternative center-identification methods: (1) the 'mean-center', as detailed above; (2) the 'median-center', which is the median vector across cells; and (3) the 'harmonic-center', defined as the in verse of the mean of the in verse of the values. In all settings, gene marker selection was done by selecting a set of genes that minimizes the condition number. We refer to these approaches as the '$K$-means + mean-/median-/harmonic-center' methods. All reported results, except from **Supplementary Fig. 4E**, use $K$-means clustering followed by the mean-center statistics. Specifically, each of the compared methods was analyzed using a variety of granularity ($K$) values.

## Synthetic data analysis demonstrates the tradeoff between complexity and scalability

The lower performance of the alternative deconvolution methods (compared to CPM; e.g., **Figure 2CD**) could be due to either a lack of scalability to a large number of reference profiles, or simply due to the cell-state complexity within the bulk cell population. Examination of the performance across varying levels of granularity indicated that the existing methods compromise between scaling issues and cell-state complexity: as the number of reference profiles increases, the ability to accurately estimate the cell-state complexity

increases; however, further increase in the number of reference profiles leads to decreased accuracy due to scaling issues (**Supplementary Fig. 2F**).

We further observed that the high performance of CPM depended not only on its 'scalability', but also on its ability to handle a high cell-state complexity: analysis of an "enrichment scheme", which analyzes each reference profile independently and thereby is scalable to a large reference collection, resulted in substantially lower accuracy compared to CPM (**Supplementary Fig. 2G**). In particular, the enrichment scheme was implemented as previously descried[7]. In brief, we used the single-sample GSEA (ssGSEA[8]) method to determine the enrichment of gene markers associated with each reference profile within the top-ranked genes of the complex tissue. The set of gene markers associated with each reference profile was defined as the $Ng$ top ANOVA-scored genes. Importantly, the comparison of CPM to ssGSEA was applied on the cell-subtype simulations, but using only a single cell type in each synthetic data collection (a 'single-cell-type design'). Current enrichment-based methods differ in their post-processing step that typically compensates between cell types but does not compensate between cell states of the same cell type[7]; our single-cell-type design therefore provides a broad comparison to the different enrichment-based approaches regardless their particular post-processing compensation strategy.

## Synthetic data analysis of different sequencing depths

To explore whether the quality of scRNA-seq data can derive improved performance, we further analyzed the impact of the number of single cells as well as the sequencing depth per cell. As the quality of scRNA-seq increased, the performance of CPM also increased; specifically, the impact of data quality on CPM appeared to be more substantial than that on alternative methods (**Supplementary Fig. 6AB**). In addition, the absence of an entire cell type from the reference dataset was also evaluated (**Supplementary Fig. 6C**), suggesting that CPM is more robust to missing cell types than the existing deconvolution methods.

## Pre-processing of bulk RNA-Seq data

Reads alignment and transcript quantification were performed as described earlier[2], with several modifications. In brief, reads were aligned using the HISAT aligner[9] to the mouse reference genome (NCBI 37, mm9). Reads that were mapped to multiple positions were excluded. Recorded are those reads that were mapped to mouse gene exons (using the UCSC transcript annotation). Expression levels were then calculated and normalized by the total number of mapped reads per experiment, using HOMMER[10]. The absolute bulk profile of an infected mouse are the HOMMER-normalized expression values of the relevant lung sample.

## Additional analyses support the gradual change predicted by the CPM method

To support the observed stepwise change in cell-to-phenotype correlations (**Figure 3C**), we performed several analyses. First, additional analyses showed that the same gradual changes also appeared using alternative techniques by which the trajectory was defined (**Supplementary Fig. 7C**), for absolute input profiles (**Supplementary Fig. 7D**), when analysing the average prediction of each genetic background (**Supplementary Fig. 7E**), and using reference data from another mouse (**Supplementary Fig. 7F**).

Second, we analyzed microarray gene expression data from a public repository (GEO accession number GSE30506), consisting of bulk profiles of the lung tissue at 4 days after influenza virus infection, across a collection of 44 pre-CC mice with extreme weight-loss phenotypes. In this analysis, CPM combined bulk and single-cell expression datasets that were generated by different labs and experimental settings: bulk pre-CC profiles at 4 days p.i.[11] and scRNA-seq profiles that were generated at 2 days p.i.[12]. Comparison with weight loss data at 4 days post influenza infection across the pre-CCs[11] confirmed the gradual change in cell-to-phenotype correlations (**Supplementary Fig. 7G**).

Third, to further confirm the gradual change in cell-to-phenotype coefficients over the trajectory, we designed and applied two statistical tests: a gradual-change test and a stepwise-change test. (i) *Gradual-change testing.* To assess gradual (ever-increasing) changes, we calculated the average cell-to-phenotype coefficients through a 50-cells sliding window along the antiviral trajectory. Each window was assigned a "+" sign if its average coefficient was the maximal compared to all its predecessors and a "-" sign otherwise. The test statistics was defined as the percentage of windows carrying a "+" sign. Statistical significance was evaluated by repeating this procedure 100 times, each time using permuted data that was generated by shuffling the position of cells along the trajectory; each reported P-value was then calculated based on the approximated distribution of permuted test statistics. (ii) *Stepwise-change testing.* We observe that the cell-to-phenotype coefficients are negative in low-antiviral state and positive in high-antiviral state, indicating a transition in the level of the cell-to-phenotypes coefficients during the progression of cells along the activation trajectory. Our null hypothesis is therefore a standard model of a one-step transition, and the alternative hypothesis is a stepwise transition. In accordance, the cell-state trajectory was divided into either two intervals (the null hypothesis) or ten equal consecutive intervals (the alternative hypothesis), and Gaussian parameters were fitted for the cell-to-phenotype coefficients of cells within each such interval. We then calculate a likelihood ratio (LR) score as the ratio between the goodness of fit of the two models - namely, the ratio between the maximal likelihood value when using a ten-Gaussian model versus the maximal likelihood among all 2-bins models that differ in their particular division cutoff. To estimate statistical significance, we repeated this procedure 100 times with randomly shuffled positions of cells along the activation trajectory, and then calculated a P-value for the observed LR score based on the approximated distribution of shuffling-based LR scores.

Overall, the data supported gradual changes in the levels of cell-to-phenotype correlations over the activation process ($p < 10^{-5}$, gradual-change test; **Supplementary Fig. 7H, top**) and a better fit of the stepwise model compared to the one-step model ($p < 0.1$, stepwise-change test; **Supplementary Fig. 7H, bottom**).

Finally, we asked whether the gradual-change trend is also revealed by alternative deconvolution methods or by using an unrelated reference dataset. **Supplementary Fig. 7I** shows that using each of the alternative deconvolution methods, the increasing trend in cell-to-phenotype correlation cannot be determine due to the lack of consistency among different levels of granularity (e.g., Cibersort captured the trend with a granularity of 10 but have missed the trend with a granularity of 4 and 20). Similarly, the lack of a trend when using an uninfected-reference dataset (**Supplementary Fig. 7J**) exemplifies the importance of using reference and bulk data derived from a similar experimental setting. These findings therefore highlight the advantage of a CPM

model that is based on a continuous space of cell states, and which is constructed using reference data of a similar protocol.

## Analysis of deconvolution using naive mice

We used lung tissues from naive (untreated) mice to demonstrate the robustness of deconvolution on a very different dataset: whereas the infected lung tissue harbors substantial cell-activation heterogeneity within each cell type, the naive lungs typically harbor discrete cell types with a limited cell-activation heterogeneity within each cell type[12]. To test deconvolution in this case, we generated RNA-seq data of the lung tissue from five naive CC mice (**Supplementary Table 1**), and applied CPM and alternative deconvolution methods on this bulk data. As a reference data we used scRNA-seq profiles of an uninfected mouse (2075 cells that were partitioned into nine cell types[12]; data from GEO accession number GSE107947). Since we focused on inter-cell-type heterogeneity, granularity=1 (a single group for each cell type) was used to construct the reference data for the alternative methods. Using comparison with known cell-type fractions in naive mice (**Supplementary Fig. 1B**), we find strong support for the accuracy of all compared methods in predicting the quantities of discrete cell types.

## Cell population maps, inferred by CPM, are a valuable resource for future investigations

CPM can be used not only to predict the cell composition by their lineage association and state, but also to utilize those predictions for future investigations. For instance, these cell-state-specific quantities may be used to calculate cell-state-specific expression within a complex tissue, as previously described for the case of cell types[13]. As another example, **Supplementary Figure 8** demonstrates how temporal dynamics can be inferred from cell-state-specific quantities. In particular, our simplified cell-state transition model suggested relationships between cell-state transition rates and physiological outcomes. Further studies are needed in order to combine additional key features of this dynamic process (such as a potential 'lag' in onset time), which may contribute to inter-individual variation in influenza-infection outcomes. While this study demonstrates a simplified stochastic model, the interpretation of CPM-reconstructed cell population maps should be further enhanced by advanced utilization of complex neighborhood structures (such as cyclic and bifurcating trajectories), as well as by continuous stochastic modeling, as suggested in previous studies[14].

Furthermore, extensive single-cell resource catalogues such as the Human Cell Atlas[15] are currently being accumulated, suggesting that CPM may soon become applicable to analyze both archived and newly generated bulk profiles without requiring expertise in single-cell technologies. For instance, application of CPM to the growing number of genomic datasets, such as TCGA[16] and GTeX[17], should provide a resource for studying the genetic basis of cellular heterogeneity and its relationships with disease outcome. This may motivate the development of a new generation of personalized predictive tools and risk factor identification methodologies[18] that are especially designed to exploit complex patterns of inter-individual variation in the cellular population structure.

# References

1       Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue expression profiles. *Nature Methods* **12**, 453-457, doi:10.1038/nmeth.3337 (2015).

2       Altboum, Z. *et al.* Digital cell quantification identifies global immune cell dynamics during influenza infection. *Molecular systems biology* **10**, 720-720, doi:10.1002/msb.134947 (2014).

3       Frishberg, A., Brodt, A., Steuerman, Y. & Gat-Viks, I. ImmQuant: a user-friendly tool for inferring immune cell-type composition from gene-expression data. *Bioinformatics* **32**, 3842-3843, doi:10.1093/bioinformatics/btw535 (2016).

4       Puram, S. V. *et al.* Single-Cell Transcriptomic Analysis of Primary and Metastatic Tumor Ecosystems in Head and Neck Cancer. *Cell* **171**, 1611-1624.e1624, doi:10.1016/j.cell.2017.10.044 (2017).

5       Tirosh, I. *et al.* Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-seq. *Science* **352**, 189, doi:10.1126/science.aad0501 (2016).

6       Schubert, E., Sander, J., Ester, M., Kriegel, H. P. & Xu, X. DBSCAN Revisited, Revisited: Why and How You Should (Still) Use DBSCAN. *ACM Trans. Database Syst.* **42**, 19:11–19:21, doi:10.1145/3068335 (2017).

7       Aran, D., Hu, Z. & Butte, A. J. xCell: digitally portraying the tissue cellular heterogeneity landscape. *Genome Biol* **18**, 220, doi:10.1186/s13059-017-1349-1 (2017).

8       Hanzelmann, S., Castelo, R. & Guinney, J. GSVA: gene set variation analysis for microarray and RNA-seq data. *BMC Bioinformatics* **14**, 7, doi:10.1186/1471-2105-14-7 (2013).

9       Kim, D., Langmead, B. & Salzberg, S. HISAT: a fast spliced aligner with low memory requirements. *Nature Methods* **12**, 357-360, doi:10.1038/nmeth.3317 (2015).

10      Heinz, S. *et al.* Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. *Molecular Cell* **38**, 576-589, doi:https://doi.org/10.1016/j.molcel.2010.05.004 (2010).

11      Bottomly, D. *et al.* Expression quantitative trait Loci for extreme host response to influenza a in pre-collaborative cross mice. *G3 (Bethesda, Md.)* **2**, 213-221, doi:10.1534/g3.111.001800 (2012).

12      Steuerman, Y. *et al.* Dissection of influenza infection in vivo by single-cell RNA sequencing. *Cell Systems* **6**, 679-691, doi:10.1016/j.cels.2018.05.008 (2018).

13      Shen-Orr, S. S. *et al.* Cell type–specific gene expression differences in complex tissues. *Nature Methods* **7**, 287, doi:10.1038/nmeth.1439 (2010).

14      Tanay, A. & Regev, A. Scaling single-cell genomics from phenomenology to mechanism. *Nature* **541**, 331, doi:10.1038/nature21350 (2017).

15      Regev, A. *et al.* The Human Cell Atlas. *eLife* **6**, e27041, doi:10.7554/eLife.27041 (2017).

16      Hutter, C. & Zenklusen, J. C. The Cancer Genome Atlas: Creating Lasting Value beyond Its Data. *Cell* **173**, 283-285, doi:10.1016/j.cell.2018.03.042 (2018).

17      Kaul, R. Enhancing GTEx by bridging the gaps between genotype, gene expression, and disease. *Nature Genetics* **49**, 1664–1670, doi:10.1038/ng.3969 (2017).

18      West, M., Ginsburg, G. S., Huang, A. T. & Nevins, J. R. Embracing the complexity of genomic data for personalized medicine. *Genome Res* **16**, 559-566, doi:doi: 10.1101/gr.3851306 (2006).