

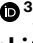










# Mapping disease regulatory circuits at cell-type resolution from single-cell multiomics data

Received: 16 November 2022

Accepted: 6 June 2023

Published online: 25 July 2023

 Check for updates

Xi Chen<sup>1,2</sup>, Yuan Wang <sup>3</sup>, Antonio Cappuccio<sup>4</sup>, Wan-Sze Cheng<sup>4</sup>, Frederique Ruf Zamojski <sup>4</sup>, Venugopalan D. Nair<sup>4</sup>, Clare M. Miller<sup>4</sup>, Aliza B. Rubenstein<sup>4</sup>, German Nudelman<sup>4</sup>, Alicja Tadych<sup>2</sup>, Chandra L. Theesfeld <sup>2</sup>, Alexandria Vornholt<sup>4</sup>, Mary-Catherine George<sup>4</sup>, Felicia Ruffin<sup>5</sup>, Michael Dagher<sup>5</sup>, Daniel G. Chawla<sup>6</sup>, Alessandra Soares-Schanoski<sup>4</sup>, Rachel R. Spurbeck<sup>7</sup>, Lishomwa C. Ndhlovu<sup>8</sup>, Robert Sebra<sup>9</sup>, Steven H. Kleinstei <sup>6,10</sup>, Andrew G. Letizia <sup>11</sup>, Irene Ramos <sup>4,12</sup>, Vance G. Fowler Jr<sup>5</sup>, Christopher W. Woods<sup>5</sup>, Elena Zaslavsky <sup>4,13</sup> , Olga G. Troyanskaya<sup>1,2,3,13</sup>  & Stuart C. Sealfon <sup>4,13</sup> 

Resolving chromatin-remodeling-linked gene expression changes at cell-type resolution is important for understanding disease states. Here we describe MAGICAL (Multiome Accessibility Gene Integration Calling and Looping), a hierarchical Bayesian approach that leverages paired single-cell RNA sequencing and single-cell transposase-accessible chromatin sequencing from different conditions to map disease-associated transcription factors, chromatin sites, and genes as regulatory circuits. By simultaneously modeling signal variation across cells and conditions in both omics data types, MAGICAL achieved high accuracy on circuit inference. We applied MAGICAL to study *Staphylococcus aureus* sepsis from peripheral blood mononuclear single-cell data that we generated from subjects with bloodstream infection and uninfected controls. MAGICAL identified sepsis-associated regulatory circuits predominantly in CD14 monocytes, known to be activated by bacterial sepsis. We addressed the challenging problem of distinguishing host regulatory circuit responses to methicillin-resistant and methicillin-susceptible *S. aureus* infections. Although differential expression analysis failed to show predictive value, MAGICAL identified epigenetic circuit biomarkers that distinguished methicillin-resistant from methicillin-susceptible *S. aureus* infections.

Gene expression can be modulated through the interplay of proximal and distal regulatory domains brought together in 3D space<sup>1</sup>. Chromatin regulatory domains, transcription factors (TFs), and downstream target genes form regulatory circuits<sup>2</sup>. Within circuits, the binding of TFs to chromatin regions and the 3D looping between these regions

and gene promoters represent the mechanisms governing how TFs transform regulatory signals into changes in RNA transcription<sup>3,4</sup>. In disease, these circuits could be dysregulated in a cell-type-specific manner and may not be observed from bulk samples<sup>5</sup>. Therefore, identifying the impact of disease on regulatory circuits requires a framework for

A full list of affiliations appears at the end of the paper. ✉ e-mail: [elena.zaslavsky@mssm.edu](mailto:elena.zaslavsky@mssm.edu); [ogt@genomics.princeton.edu](mailto:ogt@genomics.princeton.edu); [stuart.sealfon@mssm.edu](mailto:stuart.sealfon@mssm.edu)

mapping regulatory domains with chromatin accessibility changes to altered gene expression in the context of cell-type resolution<sup>6</sup>. Single-cell RNA sequencing (scRNA-seq) and single-cell transposase-accessible chromatin sequencing (scATAC-seq) characterizing disease states have improved the identification of differential chromatin sites and/or differentially expressed genes within individual cell types<sup>5,7,8</sup>.

Yet, advances in single-cell assay technology have outpaced the development of methods to maximize the value of multiomics datasets for studying disease-associated regulation, especially for the regulatory interactions that are not directly measured by the omics data. Recent computational approaches<sup>9–12</sup> to support the multiomics data analysis demonstrate the promise of this area but still lack the capacity to resolve regulation changes within individual cell types, which precludes elucidating regulatory circuits affected by the disease or showing different responses in varying disease states.

To address these, we developed MAGICAL, a method that models coordinated chromatin accessibility and gene expression variation to identify circuits (both the units and their interactions) that differ between conditions. MAGICAL analyzes scRNA-seq and scATAC-seq data using a hierarchical Bayesian framework. To accurately detect differences in regulatory circuit activity between conditions, MAGICAL introduces hidden variables for explicitly modeling the transcriptomic and epigenetic signal variations between conditions and optimization against the noise in both scRNA-seq and scATAC-seq datasets. Because regulatory circuits are cell-type specific<sup>13</sup>, MAGICAL reconstructs them at cell-type resolution. Systematic benchmarking against multiple public datasets supported the accuracy of MAGICAL-identified regulatory circuits.

*S. aureus*, a bacterium often resistant to common antibiotics, is a major cause of severe infection and mortality<sup>14,15</sup>. Using single-cell multiomics data generated from peripheral blood mononuclear cell (PBMC) samples of *S. aureus*-infected subjects and healthy controls, MAGICAL identified host response regulatory circuits that are modulated during *S. aureus* bloodstream infection, and circuits that discriminate the responses to methicillin-resistant *S. aureus* (MRSA) and methicillin-susceptible *S. aureus* (MSSA). Genes in the host circuits accurately predicted *S. aureus* infection in multiple validation datasets. Moreover, in contrast to conventional differential analysis that failed to identify specific genes for robust antibiotic-sensitivity prediction, MAGICAL-identified circuit genes can differentiate MRSA from MSSA. Therefore, MAGICAL can be used for multiomics-based gene signature development, providing a bioinformatic solution that can improve disease diagnosis.

## Results

### MAGICAL framework

MAGICAL identifies disease-associated regulatory circuits by comparing single-cell multiomics data (scRNA-seq and scATAC-seq) from disease and control samples (Fig. 1a). The framework incorporates TF motifs and chromatin topologically associated domain (TAD) boundaries as prior information to infer regulatory circuits comprising chromatin regulatory sites, modulatory TFs, and downstream target genes for each cell type. In brief, to build candidate disease-modulated circuits, differentially accessible sites (DAS) within each cell type are first associated with TFs by motif sequence matching and then linked to differentially expressed genes (DEG) in that cell type by genomic localization within the same TAD. Next, MAGICAL uses a Bayesian framework to iteratively model chromatin accessibility and gene expression variation across cells and samples in each cell type and to estimate the confidence of TF–peak and peak–gene linkages for each candidate circuit (Fig. 1b).

To accurately identify varying circuits between different conditions, MAGICAL explicitly models signal and noise in chromatin accessibility and gene expression data (see Methods section ‘MAGICAL’). A TF–peak binding variable and a hidden TF activity variable are jointly

estimated to fit to the chromatin accessibility variation across cells from the conditions being compared. These two variables are then used together with a peak–gene looping variable to fit the gene expression variation. Using Gibbs sampling, MAGICAL iteratively estimates variable values and optimizes the states of circuit TF–peak–gene linkages. Finally, high-confidence circuits fitting the signal variation in both data types are selected.

TF activity represents the regulatory capacity (protein level) of a particular TF protein<sup>16,17</sup>, which is distinct from TF expression. For each TF, we assume its hidden TF activities following an identical distribution across cells in the same cell type and the same sample, regardless of whether the cells are from the scATAC-seq assay, the scRNA-seq assay, or both. MAGICAL iteratively learns the activity distribution for each TF and estimates the specific activities of all TFs in each cell (Supplementary Fig. 1). This procedure eliminates the requirement of cell-level pairing of RNA-seq and ATAC-seq data. Thus, MAGICAL is a general tool that can analyze single-cell true multiome or sample-paired multiomics datasets.

We validated MAGICAL in multiple ways, demonstrating that it infers regulatory circuits accurately (Fig. 1c). MAGICAL-inferred linkages between chromatin sites and genes were validated using experimental 3D chromatin interactions. The resulting circuit genes, sites, and their regulatory TFs were evaluated in multiple independent studies. And finally, as one example of utility, we showed that the circuit genes can be used as features to classify disease states, providing a bioinformatics solution to challenging diagnostic problems.

### Comparative analysis of performance

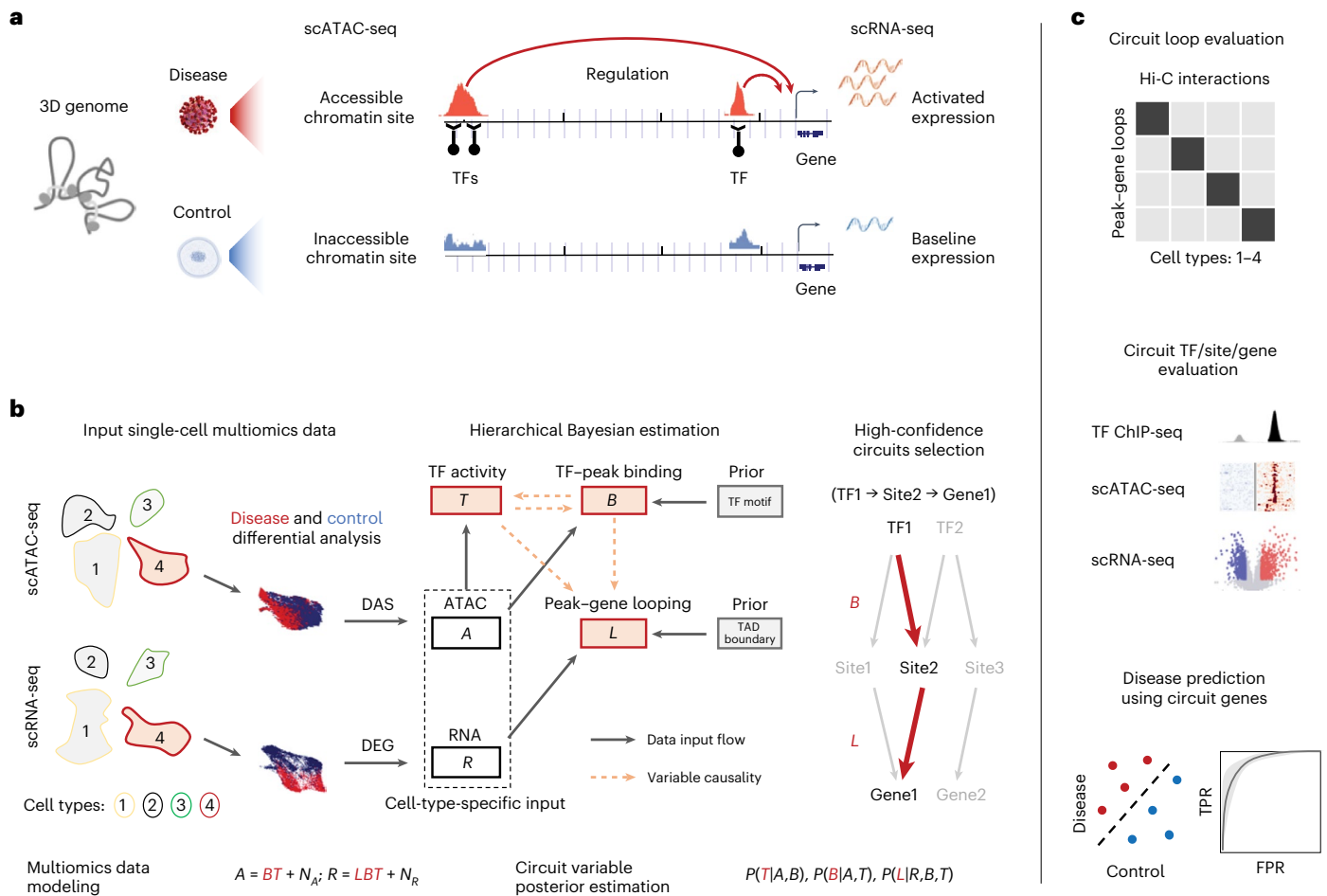
MAGICAL is a scalable framework. It can infer regulatory circuits of TFs, chromatin sites, and genes with differential activities between contrast conditions or infer regulatory circuits with active chromatin sites and genes in a single condition. Because existing integrative methods<sup>11,12,18</sup> can only be applied to single-condition data, to provide a comparative assessment of the performance of MAGICAL, we restricted MAGICAL to the single-condition data analysis possible with existing methods.

For peak–gene looping inference, we compared MAGICAL to the TRIPOD<sup>11</sup> and FigR<sup>18</sup> methods, using the same benchmark single-cell multiome datasets as used by the authors reporting these methods. In the comparison of MAGICAL with TRIPOD using a 10x multiome single-cell dataset, MAGICAL-inferred peak–gene loops showed significantly higher enrichment of experimentally observed chromatin interactions in blood cells in the 4DGenome database<sup>19</sup> ( $P < 0.0001$ , two-sided Fisher’s exact test, Supplementary Fig. 2a), the same validation data used by TRIPOD developers. MAGICAL also significantly outperformed FigR on the application to a GMI2878 SHARE-seq dataset<sup>10</sup>. In that case, the peak–gene loops in MAGICAL-selected circuits had significantly higher enrichment of H3K27ac-centric chromatin interactions<sup>20</sup> than did FigR ( $P < 0.0001$ , two-sided Fisher’s exact test, Supplementary Fig. 2b).

Because the MAGICAL framework, unlike TRIPOD and FigR, used chromatin TAD boundaries as prior information, we evaluated whether the improvement in performance resulted solely from this additional information. To investigate this, we eliminated the use of TAD boundaries and modified MAGICAL for this test by assigning candidate linkages between peaks and genes within 500 kb (a naive distance prior). As shown in Supplementary Fig. 2a,b, even without the prior TAD information, MAGICAL still outperformed the competing methods ( $P < 0.001$ , two-sided Fisher’s exact test). Overall, these results suggest that in addition to the benefit of priors, explicit modeling of signal and noise in both chromatin accessibility and gene expression data increased the accuracy of peak–gene looping identification.

### MAGICAL analysis of COVID-19 single-cell multiomics data

To demonstrate the accuracy of the primary application of MAGICAL on contrast-condition data to infer disease-modulated circuits, we applied MAGICAL to sample-paired PBMC scRNA-seq and scATAC-seq



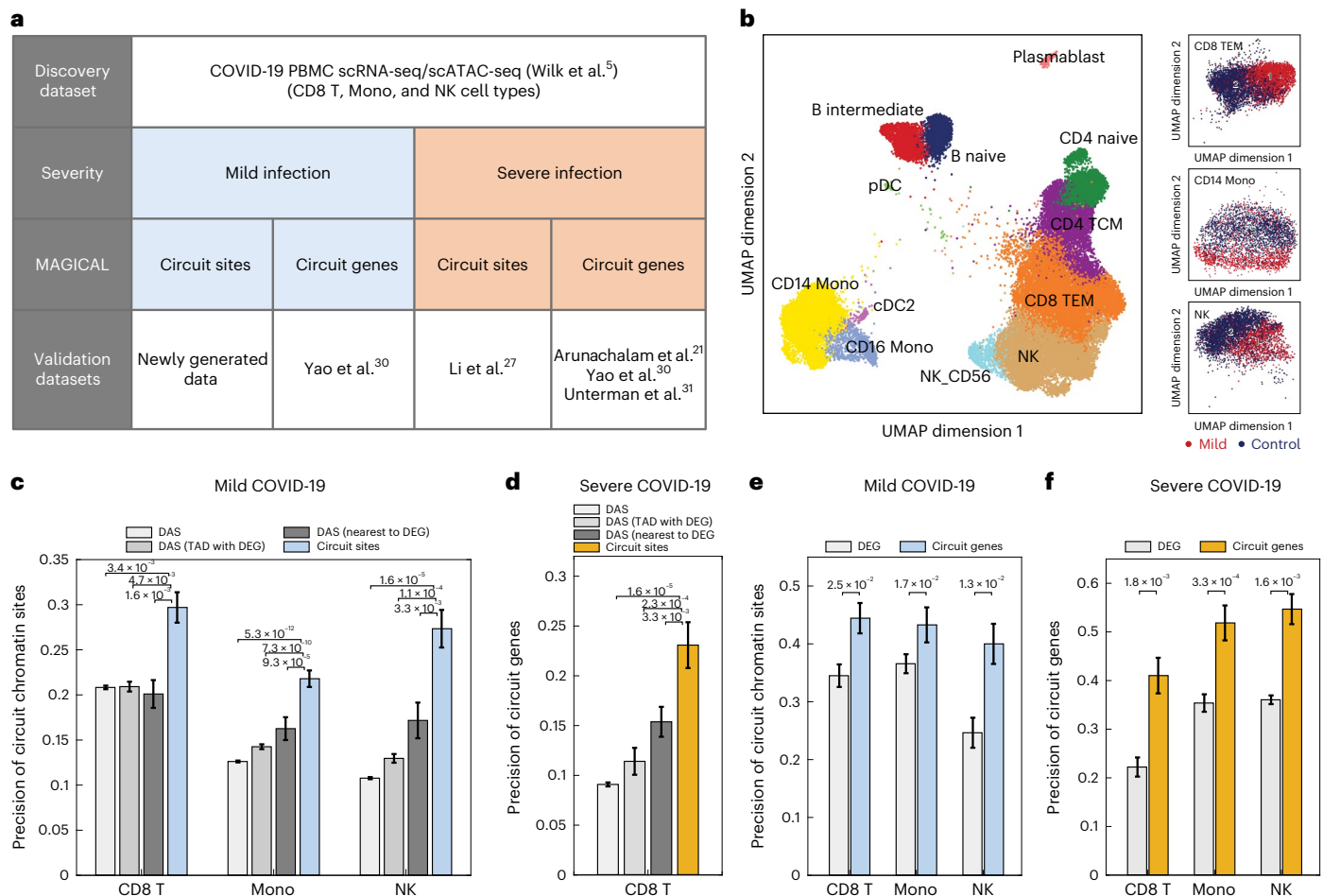
**Fig. 1 | Overview of MAGICAL for mapping disease-associated regulatory circuits from scRNA-seq and scATAC-seq data.** **a**, Disease-modulated regulatory circuitis. In the 3D genome, the altered gene expression in cells between disease and control conditions can be attributed to the chromatin accessibility changes of proximal and distal chromatin sites regulated by TFs. **b**, MAGICAL framework. To identify disease-associated regulatory circuits in a selected cell type (including ATAC assay cells and RNA assay cells from samples being compared), MAGICAL selects DAS as candidate chromatin sites (peaks) and DEG as candidate genes. Then, the filtered ATAC data and RNA data of DAS and DEG are used as input to a hierarchical Bayesian framework pre-embedded with the prior TF motifs and TAD boundaries. The chromatin activity  $A$  is modeled as a linear combination of TF–peak binding confidence  $B$  and the hidden TF activity  $T$ , with data noise contamination  $N_A$ . The gene expression  $R$  is modeled

as a linear combination of  $B$ ,  $T$ , and peak–gene looping confidence  $L$ , with data noise contamination  $N_R$ . MAGICAL estimates the posterior probabilities  $P(B|A, T)$ ,  $P(T|A, B)$ , and  $P(L|R, B, T)$  by iteratively sampling variables  $B$ ,  $T$ , and  $L$  to optimize against the data noise  $N_A$  and  $N_R$  in both modalities. Finally, regulatory circuits with high posterior probabilities of  $B$  and  $L$  (for example, a high confidence circuit with inferred interactions between TF1, Site2, and Gene1) are selected. **c**, Results validation. We evaluate the accuracy and cell-type specificity of the inferred peak–gene looping interactions by checking their enrichment with cell-type-matched chromatin interactions in Hi-C experiments. For the identified TFs, chromatin sites, and genes in circuits, we checked the accuracy of each using independent ChIP-seq, scATAC-seq, and scRNA-seq data. Finally, as a demonstration of the utility of MAGICAL, we used the circuit target genes as features to predict disease states.

data from individuals infected with SARS-CoV-2 and healthy controls<sup>5</sup>. Because immune responses in patients with COVID-19 differ according to disease severity<sup>21,22</sup>, MAGICAL inferred the regulatory circuits for mild and severe clinical groups separately. The chromatin sites and genes in the identified circuits were validated using newly generated and publicly available independent COVID-19 single-cell datasets (Fig. 2a). We primarily focused on three cell types that have been found to show widespread gene expression and chromatin accessibility changes in response to SARS-CoV-2 infection<sup>23,24</sup>, including CD8 effector memory T (TEM) cells, CD14 monocytes (Mono), and natural killer (NK) cells. In total, MAGICAL identified 1,489 high-confidence circuits (1,404 sites and 391 genes) in these cell types for mild and severe clinical groups (Supplementary Table 1; Methods section ‘MAGICAL analysis of COVID-19 single-cell multiomics data’).

To confirm the circuit chromatin sites selected by MAGICAL for mild COVID-19, we generated an independent PBMC scATAC-seq

dataset from six people infected with SARS-CoV-2 with mild symptoms and three uninfected (polymerase chain reaction (PCR)-negative) controls (Fig. 2b; Supplementary Table 2). Approximately 25,000 quality cells were selected after quality-control (QC) analysis. These cells were integrated, clustered, and annotated using ArchR<sup>25</sup> (Supplementary Fig. 3; Supplementary Table 3). Peaks were called from each cell type using MACS2<sup>26</sup>. In total, 284,909 peaks were identified (Supplementary Table 4). For the three selected cell types, differential analysis between COVID-19 and control groups returned 3,061 sites for CD8 TEM, 1,301 sites for CD14 Mono, and 1,778 sites for NK (Supplementary Table 5; Methods section ‘COVID-19 PBMC scATACseq data analysis’). This produced three validation peak sets for mild COVID-19 infection. For severe COVID-19, an existing study focused on T cells identified specific chromatin activity changes with severe COVID-19 in CD8 T cells<sup>27</sup>. We used their reported chromatin sites to validate the circuit chromatin sites identified in CD8 T cells. In all four validation sets, the precision



**Fig. 2 | Validation of COVID-19-associated circuit chromatin sites and genes.**

**a**, We applied MAGICAL to a COVID-19 PBMC single-cell multiomics dataset and identified circuits for the clinical mild and severe groups. We validated the MAGICAL-selected circuit sites and genes using newly generated and independent COVID-19 single-cell datasets. **b**, UMAPs of a newly generated independent scATAC-seq dataset including 16,000 cells from six people with COVID-19 and 9,000 cells from three controls showed chromatin accessibility changes in CD8 TEM, CD14 Mono, and NK cell types. **c,d**, The precision of MAGICAL-selected circuit sites is significantly higher than that of the original DAS, the nearest DAS to DEG, or all DAS in the same TAD with DEG. **e,f**, The precision of circuit genes are significantly higher than that of DEG. **c,e**, For mild

COVID-19, MAGICAL identified 645 sites in CD8 TEM, 599 sites in CD14 Mono, and 148 sites in NK, regulating 153 genes, 183 genes, and 60 genes, respectively. **d,f**, For severe COVID-19, MAGICAL identified 78 sites, 202 sites, and 62 sites in the three cell types, regulating 25 genes, 81 genes, and 26 genes, respectively. **c-f**, Precision is defined as the proportion of the selected sites and genes to be differentially accessible and differentially expressed in the same cell type between infection and control conditions in independent datasets. Results are presented as bar plots where the heights represent the precision and the error bars represent the 95% confidence interval. Significance is evaluated using a two-sided Fisher's exact test and *P* values between bars are shown.

(proportion of sites that are differential in the validation data) of the MAGICAL-selected chromatin sites was significantly higher than the original DAS (*P* < 0.001, two-sided Fisher's exact test, Fig. 2c,d).

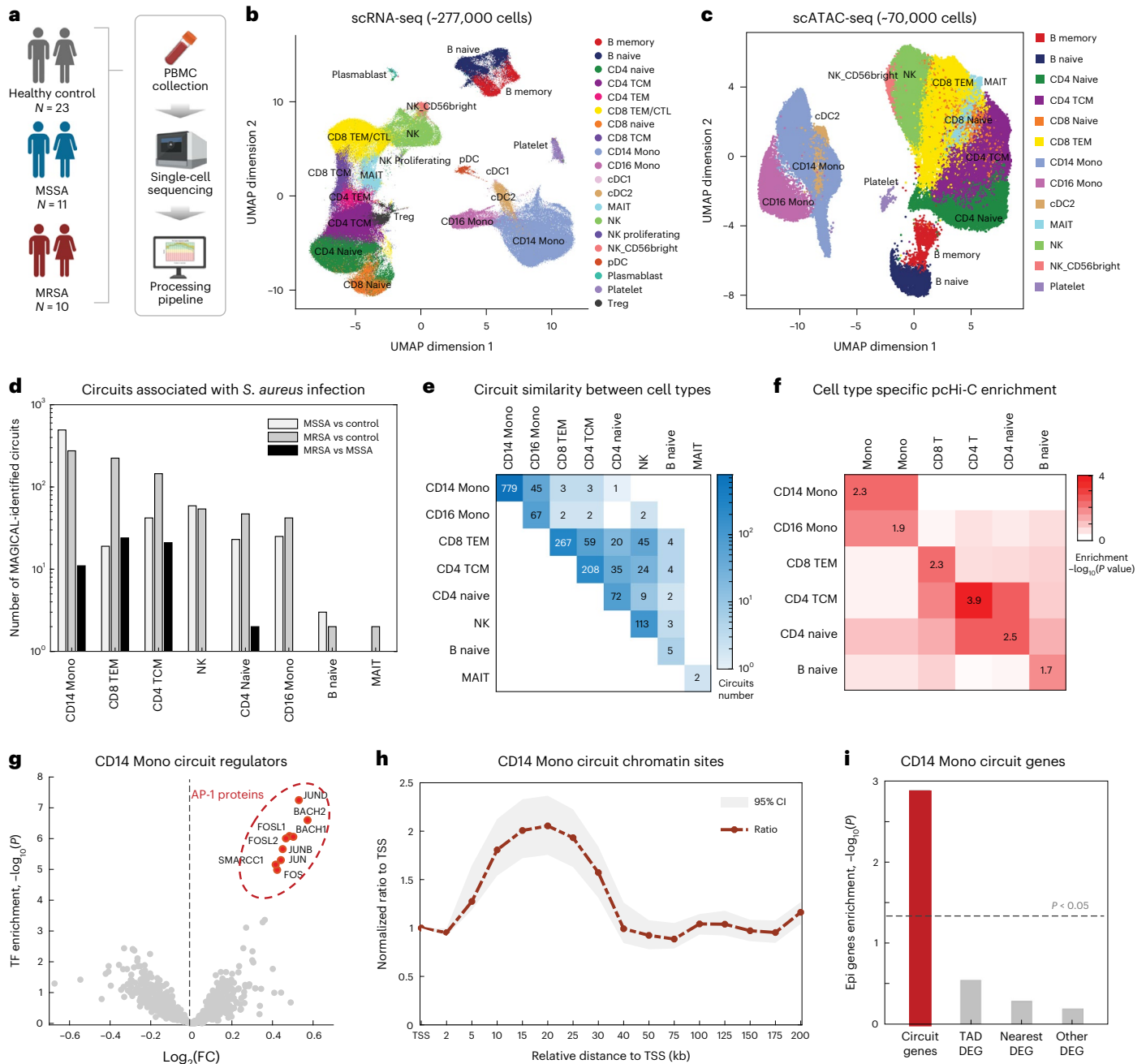
When multiple potential chromatin regulatory loci are identified in the vicinity of a specific gene, it is commonly assumed that the locus closest to the transcriptional starting site (TSS) is likely to be the most important regulatory site. Challenging this assumption, however, are experimental studies that show genes may not be regulated by the nearest region<sup>28,29</sup>. Supporting the importance of more distal regulatory loci, MAGICAL-selected chromatin sites significantly outperformed the nearest DAS to the TSS of DEG or all DAS within the same TAD with DEG, and the improvement is substantial (precision is approximately 50% better with MAGICAL, *P* < 0.05, two-sided Fisher's exact test, Fig. 2c,d).

To validate the circuit genes modulated by mild or severe COVID-19, we used genes reported by external COVID-19 single-cell studies<sup>21,30,31</sup>. In total, we collected six validation gene sets (three cell types for mild COVID-19 and three cell types for severe COVID-19). The precision of MAGICAL-selected circuit genes is significantly higher

than that of original DEG in all validations (precision is approximately 30% better with MAGICAL, *P* < 0.05, two-sided Fisher's exact test, Fig. 2e,f). These results confirmed the increased accuracy of disease association for both chromatin sites and genes in MAGICAL-identified regulatory circuits.

### MAGICAL analysis of *S. aureus* single-cell multiomics data

We applied MAGICAL to the clinically important challenge of distinguishing MRSA and MSSA infections<sup>32-34</sup>. We profiled sample-paired scRNA-seq and scATAC-seq data using human PBMCs from adults whose blood cultures were positive for *S. aureus* (10 MRSA and 11 MSSA), and from 23 uninfected control subjects (Fig. 3a; Supplementary Table 6). To integrate scRNA-seq data from all samples, we implemented a Seurat-based<sup>35</sup> batch correction and cell type annotation pipeline (Methods section 'S. aureus scRNA-seq data analysis'). In total, 276,200 quality cells were selected and labeled (Fig. 3b; Supplementary Fig. 4; Supplementary Table 7). For scATAC-seq data, we integrated the fragment files from quality samples using ArchR<sup>25</sup> and selected and annotated 70,174

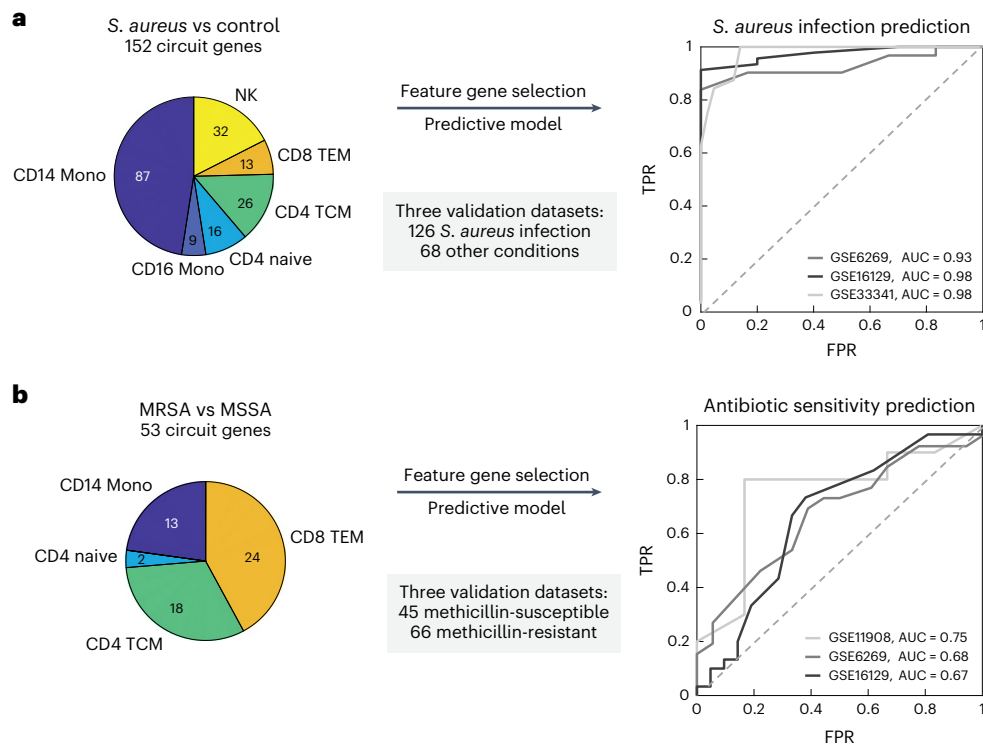


**Fig. 3 | MAGICAL accurately identified distal regulatory chromatin sites and epi-driven genes associated with *S. aureus* infection.** **a**, We collected PBMC samples from 10 subjects infected with MRSA, 11 with MSSA, and 23 uninfected control subjects and generated sample-paired scRNA-seq and scATAC-seq data using separate assays. **b**, UMAP of integrated scRNA-seq data with 18 PBMC cell subtypes. **c**, UMAP of integrated scATAC-seq data with 13 PBMC cell subtypes. Under-represented subtypes including cDC1, CD4 TEM, CD8 CTL, pDC, and Plasmablast (representing less than 5% of cells in the scRNA-seq data in total), were not recovered from the scATAC-seq data. **d**, The number of MAGICAL-identified regulatory circuits in contrast analysis for each cell type. **e**, The number of shared and specific circuits between cell types. **f**, Enrichment of circuit peak–gene interactions in each cell type with cell-type-specific pHi-C interactions. **g–i**, We specifically analyzed MAGICAL-identified regulatory

circuits for CD14 Mono. **g**, TF motif enrichment analysis in circuit sites showed that AP-1 proteins are mostly significantly enriched at chromatin regions with increased accessibility in the infection condition. The  $\log_2(\text{FC})$  is calculated for each TF by dividing the number of binding sites with increased chromatin activity in the infection condition by the number of sites with decreased activity. **h**, In total, 633 circuit sites were identified by MAGICAL. Compared with all accessible chromatin sites, an increased proportion of circuit sites were in the range of 15 kb to 25 kb relative to gene TSS. In the curve, the center points represent the FC between the proportions of circuit sites and background sites at each location. The upper and lower points represent the 95% confidence interval. **i**, The circuit genes were significantly enriched with experimentally confirmed epigenetically driven genes (epi-genes) in monocytes. All significance was assessed using adjusted *P* values from a one-sided hypergeometric test.

quality cells (Fig. 3c; Supplementary Fig. 5; Supplementary Table 8). In total, 388,860 peaks were identified (Supplementary Fig. 5b; Supplementary Table 9; Methods section ‘*S. aureus* scATAC-seq data analysis’).

In total, 13 major cell types that surpassed the 200-cell threshold in both scRNA-seq and scATAC-seq data were selected for subsequent analysis (Supplementary Fig. 6). Differential analysis for three contrasts



**Fig. 4 | MAGICAL-identified circuit genes robustly predict *S. aureus* infection and bacteria antibody sensitivity. a**, Circuit genes in common to MRSA and MSSA infections achieved a near-perfect classification of *S. aureus* infected and uninfected samples in multiple independent datasets (one adult dataset and

two pediatric datasets). **b**, Circuit genes that differed between MRSA and MSSA showed predictive value of antibiotic sensitivity in independent patient samples (three pediatric datasets).

(MRSA versus control, MSSA versus control, and MRSA versus MSSA) in each cell type returned a total of 1,477 DEG and 23,434 DAS (Supplementary Fig. 7; Supplementary Tables 10 and 11).

MAGICAL identified 1,513 high-confidence regulatory circuits (1,179 sites and 371 genes) within cell types for three contrasts (MRSA versus control, MSSA versus control, and MRSA versus MSSA) (Supplementary Table 12; Methods: MAGICAL analysis of *S. aureus* single-cell multiomics data). It has been reported that activation of CD14 Mono plays a principal role in response to *S. aureus* infection<sup>36,37</sup>. In MAGICAL analysis, CD14 Mono showed the highest number of regulatory circuits (Fig. 3d). Comparing circuits between cell types we found that these disease-associated circuits are cell-type-specific (Fig. 3e). For example, circuits rarely overlapped between very distinct cell types like monocytes and T cells. Between relevant cell types like CD14 Mono and CD16 Mono, or between subtypes of T cells, most circuits are still specific for one cell subtype. These circuits were further validated using cell type-specific chromatin interactions reported in a reference promoter capture (pc) Hi-C dataset<sup>13</sup>. In all the cell types for which the cell-type-specific pcHi-C data was available (B cells, CD4 T cells, CD8 T cells, CD14 Mono), the circuit peak–gene interactions showed significant enrichment of pcHi-C interactions in matched cell types (Fig. 3f;  $P < 0.01$ , one-sided hypergeometric test). For comparison, we also performed the peak–gene interaction enrichment analysis between different cell types, finding significantly lower enrichment levels. These results indicate the cell-type specificity of MAGICAL-identified circuits.

In CD14 Mono, MAGICAL identified AP-1 complex proteins as the most important regulators, especially at chromatin sites with increased activity in cells exposed to infections (Fig. 3g). This finding is consistent with the importance of this protein complex in gene regulation in response to a variety of infections<sup>5,38,39</sup>. Supporting the accuracy of the identified TFs, we compared circuit chromatin sites with ChIP-seq

peaks from the Cistrome database<sup>40</sup>. The most similar TF ChIP-seq profiles were from AP-1 complex proteins JUN and FOS in blood or bone marrow samples (Supplementary Fig. 8). Moreover, functional enrichment analysis<sup>41</sup> of the circuit genes showed that cytokine signaling, a known pathway mediated by AP-1 complex and associated with the inflammatory responses in macrophages<sup>42,43</sup>, was the most enriched (adjusted  $P = 2.4 \times 10^{-11}$ , one-sided hypergeometric test).

MAGICAL modeled regulatory effects of both proximal and distal sites on genes. We examined the chromatin site location relative to the target gene TSS, for the identified circuits in CD14 Mono. Compared to all ATAC peaks called around the circuit genes, a substantially increased proportion of circuit chromatin sites were located 15Kb to 25Kb away from the TSS (Fig. 3h). This pattern is consistent with the 24Kb median enhancer distance found by CRISPR-based perturbation in a blood cell line<sup>44</sup>. In addition, nearly 50% of circuit chromatin sites were overlapping with enhancer-like regions in the ENCODE database<sup>45</sup>, further emphasizing that MAGICAL circuits are enriched in distal regulatory loci. We also found that these circuit chromatin sites were significantly enriched in inflammatory-associated genomic loci reported in the genome-wide association studies (GWAS) catalog database<sup>46</sup>, suggesting active host epigenetic responses to infectious diseases (Supplementary Fig. 9;  $P < 0.005$  when compared to control diseases, two-sided Wilcoxon rank sum test). Notably, one distal chromatin site (hg38 chr6: 32,484,007–32,484,507) looping to gene HLA-DRB1 is within the most significant GWAS region (hg38 chr6: 32,431,410–32,576,834) previously reported to associate with *S. aureus* infection<sup>47</sup>.

We finally compared circuit genes to existing epi-genes whose transcriptions were significantly driven by epigenetic perturbations in CD14 Mono<sup>48</sup>. MAGICAL-identified circuit genes were significantly enriched with epi-genes (Fig. 3i; adjusted  $P < 0.005$ , one-sided hypergeometric test) while the remaining DEG not selected by MAGICAL, or DEG mappable with DAS either within the same topological domains

or closest to each other showed no evidence of being epigenetically driven. These results suggest that MAGICAL accurately identified regulatory circuits activated in response to *S. aureus* infection.

### ***S. aureus* infection prediction**

Early diagnosis of *S. aureus* infection and the strain's antibiotic sensitivity is critical to choosing appropriate treatment for this life-threatening condition. We first evaluated whether the MAGIC-identified circuit genes that are common to MRSA and MSSA infections could provide a robust signature for predicting the diagnosis of *S. aureus* infection in general. Within each cell type, we selected circuit genes common to both the MRSA versus control and MSSA versus control analyses, resulting in 152 genes (Fig. 4a; Supplementary Table 12). To evaluate the prediction accuracy of these molecular features on *S. aureus* infection, we collected external, public expression data of *S. aureus* infected subjects. In total, we found one adult whole-blood<sup>49</sup> and two pediatric PBMC bulk microarray datasets<sup>50,51</sup> that comprised a total of 126 subjects infected with *S. aureus* and 68 uninfected controls. The use of pediatric validation data has the advantage of providing a much more rigorous test of the robustness of MAGIC-identified circuit genes for classifying disease samples in this very different cohort.

To allow validation using public bulk transcriptome datasets, we refined the 152 circuit genes set by selecting those with robust performance in our dataset at pseudobulk level. We calculated an area under the receiver operating characteristic curve (AUROC) for each circuit gene by classifying *S. aureus* infected and control subjects using pseudobulk gene expression (aggregated from the discovery scRNA-seq data). In total, 117 circuit genes with AUROCs greater than 0.7 were selected (Supplementary Table 13; Supplementary Fig. 10a). Functional gene enrichment analysis showed that interleukin (IL)-17 signaling was significantly enriched (adjusted  $P = 2.4 \times 10^{-4}$ , one-sided hypergeometric test), including genes from the AP-1, Hsp90, and S100 families. IL-17 has been found to be essential for the host defense against cutaneous *S. aureus* infection in mouse models<sup>52</sup>. We trained a support vector machine (SVM) model using the selected circuit genes as features and the discovery pseudobulk gene expression data as input. We then applied the trained SVM model to each of the three validation datasets. The model achieved high prediction performance on all datasets, with AUROCs from 0.93 to 0.98 (Fig. 4a).

This generalizability of circuit genes for predicting infection in different cohorts suggested that MAGIC identifies regulatory processes that are fundamental to the host response to *S. aureus* sepsis. We further evaluated this by comparing the 117 circuit genes to 366 filtered DEG (with per gene AUROC > 0.7 in the discovery pseudobulk gene expression data). We examined the differential expression  $\pi$  value<sup>33</sup> (a statistic score that combines both fold change (FC) and  $P$ -values) of genes in the validation datasets and found significantly higher  $\pi$  values for the circuit genes (Supplementary Fig. 10b;  $P = 9.0 \times 10^{-3}$ , one-sided Wilcoxon rank sum test).

### ***S. aureus* antibiotic sensitivity prediction**

We then addressed the more challenging problem of predicting strain antibiotic sensitivity among *S. aureus* infected subjects. When we tested the predictive models trained with DEG for the contrast of MRSA and MSSA on three pediatric PBMC microarray datasets<sup>50,51,54</sup> (comprising a total of 66 MRSA and 45 MSSA samples), we did not find predictive value (median of prediction areas under the curve (AUCs) close to 0.5; Supplementary Fig. 10d–f). And in all tests, the statistical difference of DEG-based prediction scores between MRSA and MSSA samples in the validation datasets was never significant. These results suggest that using host scRNA-seq data alone fails to identify robust molecular features for predicting the antibiotic sensitivity of the infected strain. Our observation echoes previous studies showing that in some challenging cases, differential expression analysis using RNA-seq data had

limited power to identify robust features for disease-control sample classification<sup>55</sup>.

With MAGIC we identified 53 circuit genes from the comparative multiomics data analysis between MRSA and MSSA (Supplementary Table 14). A model trained using 32 circuit genes that were robustly differential in the discovery pseudobulk data (per gene discovery AUROC > 0.7, Supplementary Fig. 10c) best distinguished antibiotic-resistant and antibiotic-sensitive samples in all three validation datasets, with AUROCs from 0.67 to 0.75 (Fig. 4b). And the statistical difference between prediction scores of MRSA and MSSA samples was significant ( $P = 9.2 \times 10^{-3}$ , two-sided Wilcoxon rank sum test). The success of the circuit-gene-based model demonstrated that MAGIC captured generalizable regulatory differences in the host immune response to these closely related bacterial infections.

## **Discussion**

MAGIC addressed the previously unmet need of identifying differential regulatory circuits based on single-cell multiomics data from contrast conditions. Critically, it identifies regulatory circuits involving distal chromatin sites. The previously difficult-to-predict distal regulatory sites are increasingly recognized as key for understanding gene regulatory mechanisms. As MAGIC uses DAS and DEG called from a pre-selected cell type, for less distinct cell types or conditions, it is harder for MAGIC to infer circuits at cell-type resolution as there will be few candidate peaks and genes. Also, MAGIC analyzes each cell type separately, and cell-type specificity is not directly modeled for disease circuit identification. Incorporating an approach to directly identify cell-type-specific circuits regulated in disease conditions would be valuable. In future work, we hope to extend the MAGIC framework to improve circuit identification when cell types are poorly defined and to model cell-type specificity.

## **Methods**

### **Human participants**

The COVID-19 study protocol was approved by the Naval Medical Research Center institutional review board (protocol number NMRC.2020.0006) in compliance with all applicable federal regulations governing the protection of human subjects. The *S. aureus* sepsis study protocol was reviewed and approved by the Duke Medical School institutional review board (protocol number Pro00102421). Subjects provided written informed consent prior to participation.

### **Statistics and reproducibility**

No statistical methods were used to pre-determine sample sizes. No data were excluded from the analyses. The experiments were not randomized. The investigators were not blinded to allocation during experiments and outcome assessment.

### ***S. aureus* patient and control samples selection**

Patients with culture-confirmed *S. aureus* bloodstream infection transferred to Duke University Medical Center were eligible if pathogen speciation and antibiotic susceptibilities were confirmed by the Duke Clinical Microbiology Laboratory. DNA and RNA samples, PBMCs, clinical data, and the bacterial isolate from the subject were cataloged using an institutional review board-approved notification of decedent research. We excluded samples with prior enrollment of the patient in this investigation (to ensure statistical independence of observations) or they were polymicrobial (that is, more than one organism in blood or urine culture). In total, 21 adult patients were selected (10 MRSA and 11 MSSA). None of them received any antibiotics in the 24 h before the bloodstream infection. Control samples were obtained from uninfected healthy adults matching the sample number and age range of the patient group. In total, 23 samples were collected from two cohorts: 14 controls (provided by Weill Cornell Medicine, New York, NY), and 9 controls (provided by the Battelle Memorial Institute,

Columbus, OH). Meta information of the selected subjects are provided in Supplementary Table 6.

### PBMC thawing

Frozen PBMC vials were thawed in a water bath at 37 °C for 1–2 minutes and placed on ice. Roswell Park Memorial Institute (RPMI) medium with 20% fetal bovine serum (FBS) (500 µl) was added dropwise to the thawed vial, the content was aspirated and added dropwise to 9 ml of RPMI/20% FBS. The tube was gently inverted to mix, before being centrifuged at 300g for 5 min. After removal of the supernatant, the pellet was resuspended in 1–5 ml of RPMI/10% FBS depending on the size of the pellet. Cell count and viability were assessed with Trypan Blue on a Countess II cell counter (Invitrogen).

### *S. aureus* scRNA-seq data generation

scRNA-seq was performed (10x Genomics, Pleasanton, CA), following the Single Cell 3' Reagents Kits V3.1 User Guidelines. Cells were filtered, counted on a Countess instrument, and resuspended at a concentration of 1,000 cells µl<sup>-1</sup>. The number of cells loaded on the chip was determined based on the 10x Genomics protocol. The 10x chip (Chromium Single Cell 3' Chip kit G PN-200177) was loaded to target 5,000–10,000 final cells. Reverse transcription was performed in the emulsion and complementary DNA was amplified following the Chromium protocol. Quality control and quantification of the amplified cDNA were assessed on a Bioanalyzer (High-Sensitivity DNA Bioanalyzer kit) and the library was constructed. Each library was tagged with a different index for multiplexing (Chromium i7 Multiplex Single Index Plate T Set A, PN-2000240) and quality controlled using a Bioanalyzer prior to sequencing.

### *S. aureus* scRNA-seq data analysis

Reads of scRNA-seq experiments were aligned to human reference genome (hg38) using 10x Genomics Cell Ranger software (version 1.2). The filtered feature-by-barcode count matrices were then processed using Seurat<sup>35</sup>. Quality cells were selected as those with more than 400 features (transcripts), fewer than 5,000 features, and less than 10% of mitochondrial content (Supplementary Fig. 4; Supplementary Table 7). Cell cycle phase scores were calculated using the canonical markers for G2M and S phases embedded in the Seurat package. Finally, the effects of mitochondrial reads and cell cycle heterogeneity were regressed out using SCTransform.

To integrate cells from heterogeneous disease samples, we first built a reference by integrating and annotating cells from the uninfected control samples using a Seurat-based pipeline. For batch correction, we identified the intrinsic batch variants and used Seurat to integrate cells together with the inferred batch labels. All control samples were integrated into one harmonized query matrix. Each cell was assigned a cell-type label by referring to a reference PBMC single cell dataset. The cell-type label of each cell cluster was determined by most cell labels in each. Canonical markers were used to refine the cell-type label assignment. This integrated control object was used as reference to map the infected samples.

To avoid artificially removing the biological variance between each infected sample during batch correction, we computationally predicted and manually refined cell types for each sample. All infection samples were projected onto the UMAP (Uniform Manifold Approximation and Projection) of the control object for visualization purpose. In total, 276,200 high-quality cells and 19 cell types with at least 200 cells in each were selected for the subsequent analysis. Within each cell type, DEG between contrast conditions were first called using the FindMarkers function of the Seurat V4 package<sup>35</sup> with default parameters. DEG with Wilcoxon test false discovery rate (FDR) < 0.05,  $|\log_2(\text{FC})| > 0.1$  and actively expressed in at least 10% of cells (pct > 0.1) from either condition were selected. To correct potential bias caused by the different sequencing depth between samples, we ran DESeq2<sup>56</sup> on the aggregated

pseudobulk gene expression data. Refined DEG passing pseudobulk differential statistics  $P < 0.05$  and  $|\log_2(\text{FC})| > 0.3$  were selected as the final DEG (Supplementary Table 10).

### Nuclei isolation for scATAC-seq

Thawed PBMCs were washed using phosphate-buffered saline (PBS) with 0.04% bovine serum albumin (BSA). Cells were counted and 100,000–1,000,000 cells were added to a 2 ml microcentrifuge tube. Cells were centrifuged at 300g for 5 min at 4 °C. The supernatant carefully completely removed, and 0.1X lysis buffer (1x: 10 mM Tris-HCl pH 7.5, 10 mM NaCl, 3 mM MgCl<sub>2</sub>, nuclease-free H<sub>2</sub>O, 0.1% v/v NP-40, 0.1% v/v Tween-20, 0.01% v/v digitonin) was added. After 3 min incubation on ice, 1 ml of chilled wash buffer was added. The nuclei were pelleted at 500g for 5 min at 4 °C and resuspended in a chilled diluted nuclei buffer (10x Genomics) for scATAC-seq. Nuclei were counted and the concentration was adjusted to run the assay.

### *S. aureus* scATAC-seq data generation

scATAC-seq was performed immediately after nuclei isolation and following the Chromium Single Cell ATAC Reagent Kits V1.1 User Guide (10x Genomics, Pleasanton, CA). Transposition was performed in 10 µl at 37 °C for 60 min on at least 1,000 nuclei, before loading of the Chromium Chip H (PN-2000180). Barcoding was performed in the Gel Bead-in-Emulsion (GEMs) (12 cycles) following the Chromium protocol. After post-GEM cleanup, libraries were prepared following the protocol and were indexed for multiplexing (Chromium i7 Sample Index N, Set A kit PN-3000427). Each library was assessed on a Bioanalyzer (High-Sensitivity DNA Bioanalyzer kit).

### *S. aureus* scATAC-seq data analysis

Reads of scATAC-seq experiments were aligned to human reference genome (hg38) using 10x Genomics Cell Ranger software (version 1.2). The resulting fragment files were processed using ArchR<sup>25</sup>. Quality cells were selected as those with TSS enrichment >12, the number of fragments >3,000 and <30,000, and nucleosome ratio <2 (Supplementary Fig. 5a; Supplementary Table 8). The likelihood of doublet cells was computationally assessed using the addDoubletScores function and cells were filtered using the filterDoublets function with default settings. Cells passing quality and doublet filters from each sample were combined into a linear dimensionality reduction using the addIterativeLSI function with the input of the tile matrix (read counts in binned 500 bp across the whole genome) with iterations = 2 and varFeatures = 20,000. This dimensionality reduction was then corrected for batch effect using the Harmony method<sup>57</sup>, via the addHarmony function. The cells were then clustered based on the batch-corrected dimensions using the addClusters function. We annotated scATAC-seq cells using the addGeneIntegrationMatrix function, referring to a labeled multimodal PBMC single cell dataset. Doublet clusters containing a mixture of many cell types were manually identified and removed. In total, 70,174 high-quality cells and 13 cell types with at least 200 cells in each were selected.

Peaks were called for each cell type using the addReproduciblePeakSet function with the MACS2 peak caller<sup>26</sup> (Supplementary Fig. 5b). In total, 388,859 peaks were identified (Supplementary Table 9). Within each cell type, differentially accessible chromatin sites (DAS) between contrast conditions (MRSA versus control, MSSA versus control or MRSA versus MSSA) were called from the single cell chromatin accessibility count data using the getMarkerFeatures function<sup>25</sup>, with parameter settings as testMethod = wilcoxon, bias = log10(nFragments), normBy = ReadsInPeaks, and maxCells = 15,000. Peaks with single cell differential statistics FDR < 0.05,  $|\log_2(\text{FC})| > 0.1$ , and actively accessible in at least 10% of cells (pct > 0.1) from either condition were selected as DAS. Owing to the high false positive rate in single-cell-based differential analysis<sup>58</sup>, we further refined the DAS by fitting a linear model to the aggregated and normalized pseudobulk



chromatin accessibility data and tested DAS individually about their covariance with sample conditions<sup>56</sup>. Refined DAS passing pseudobulk differential statistics  $P < 0.05$  and  $|\log_2(\text{FC})| > 0.3$  between the contrast conditions were selected as the final DAS (Supplementary Table 11).

**MAGICAL**

To build candidate regulatory circuits, TFs were mapped to the candidate chromatin sites by searching for human TF motifs from the chromVARmotifs library<sup>59</sup> using the addMotifAnnotations function (ArchR). The TF binding sites were then linked with the candidate genes by requiring them in the same TAD within boundaries. Then, a candidate circuit is constructed with a chromatin site and a gene in the same domain, with at least one TF binding at the site.

For each cell type (that is, the  $i$ -th cell type), MAGICAL inferred the confidence of TF–peak binding and peak–gene looping in each candidate circuit using a hierarchical Bayesian framework with two models: a model of TF–peak binding confidence ( $B$ ) and hidden TF activity ( $T$ ) to fit chromatin accessibility ( $A$ ) for MTFs and  $P$  chromatin sites in  $K_{A,S,i}$  cells with scATAC-seq measures from  $S$  samples; a second model of peak–gene interaction ( $L$ ) and the refined (noise removed) regulatory region activity ( $BT$ ) to fit gene expression ( $R$ ) of  $G$  genes in  $K_{R,S,i}$  cells with scRNA-seq measures from the same  $S$  samples.

$$A_{P \times K_{A,S,i}} = B_{P \times M,i} T_{M \times K_{A,S,i}} + N_{P \times K_{A,S,i}}, \tag{1}$$

$$R_{G \times K_{R,S,i}} = L_{G \times P,i} B_{P \times M,i} T_{M \times K_{R,S,i}} + N_{G \times K_{R,S,i}}, \tag{2}$$

$A_{P \times K_{A,S,i}}$  was a  $P$  by  $K_{A,S,i}$  matrix with each element  $a_{p,k_{A,S,i}}$  representing the ATAC read count of  $p$ -th chromatin site (ATAC peak) in the  $k_{A,S,i}$ -th cell in the  $s$ -th sample.

$R_{G \times K_{R,S,i}}$  was a  $G$  by  $K_{R,S,i}$  matrix with each element  $r_{g,k_{R,S,i}}$  representing the RNA read count of  $g$ -th gene in the  $k_{R,S,i}$ -th cell of the  $s$ -th sample.

$N_{P \times K_{A,S,i}}$  and  $N_{G \times K_{R,S,i}}$  represented data noise corresponding to  $A_{P \times K_{A,S,i}}$  and  $R_{G \times K_{R,S,i}}$ .

$B_{P \times M,i}$  was a  $P$  by  $M$  matrix with each element  $b_{p,m,i}$  representing the binding confidence of the  $m$ -th TF on the  $p$ -th candidate chromatin site.

$L_{G \times P,i}$  was a  $G$  by  $P$  matrix with each element  $l_{p,g,i}$  representing the interaction between the  $p$ -th chromatin site and the  $g$ -th gene.

$T_{M \times K_{A,S,i}}$  was an  $M$  by  $K_{A,S,i}$  matrix with each element  $t_{m,k_{A,S,i}}$  representing the hidden TF activity of the  $m$ -th TF in the  $k_{A,S,i}$ -th ATAC cell of  $s$ -th sample.

$T_{M \times K_{R,S,i}}$  was an  $M$  by  $K_{R,S,i}$  matrix with each element  $t_{m,k_{R,S,i}}$  representing the hidden TF activity of the  $m$ -th TF in the  $k_{R,S,i}$ -th RNA cell of  $s$ -th sample.

$T_{M \times K_{A,S,i}}$  and  $T_{M \times K_{R,S,i}}$  were both extended from the same  $T_{M \times S,i}$  (with elements  $t_{m,s,i}$ ) by assuming that in the  $i$ -th cell type and the  $s$ -th sample, the  $m$ -th TF's regulatory activities in all ATAC cells and all RNA cells followed an identical distribution of a single variable  $t_{m,s,i}$ . Therefore,  $K_{A,S,i}$  and  $K_{R,S,i}$  can be different numbers and MAGICAL will only estimate the matrix  $T_{M \times S,i}$ .

To select high-confidence regulatory circuits, MAGICAL estimated the confidence (probability) of TF–peak binding  $B_{P \times M,i}$  and peak–gene interaction  $L_{G \times P,i}$  together with the hidden variable  $T_{M \times S,i}$  in a Bayesian framework.

$$P(B, T, L|A, R) \propto P(R|L, B, T)P(A|B, T)P(L)P(B)P(T). \tag{3}$$

Based on the regulatory relationship among chromatin sites, upstream TFs, and downstream genes (as illustrated in Fig. 1), the posterior probability of each variable can be approximated as:

$$P(T|A, B) \propto P(A|B, T)P(T), \tag{4}$$

$$P(B|A, T) \propto P(A|B, T)P(B), \tag{5}$$

$$P(L|R, B, T) \propto P(R|L, B, T)P(L). \tag{6}$$

Although the prior states of  $b_{p,m,i}$  and  $l_{p,g,i}$  were obtained from the prior information of TF motif–peak mapping and topological-domain-based peak–gene pairing, their values were unknown. We assumed zero-mean Gaussian priors for  $B, L$  and the hidden variable  $T$  by assuming that positive regulation and negative regulation would have the same priors, which is likely to be true given the fact that there were usually similar numbers of upregulated and downregulated peaks and genes after the differential analysis. We set a high variance (non-informative) in each prior distribution to allow the algorithm to learn the distributions from the input data.

$$b_{p,m,i} \sim \text{normal}(\mu_B, \sigma_B^2), \tag{7}$$

$$t_{m,s,i} \sim \text{normal}(\mu_T, \sigma_T^2), \tag{8}$$

$$l_{p,g,i} \sim \text{normal}(\mu_L, \sigma_L^2). \tag{9}$$

where  $(\mu_B, \sigma_B^2)$ ,  $(\mu_T, \sigma_T^2)$ , and  $(\mu_L, \sigma_L^2)$  are hyperparameters representing the prior mean and variance of TF–peak binding, TF activity, and peak–gene looping variables.

The likelihood functions  $P(A|B, T)$  and  $P(R|L, B, T)$  represent the fitting performance of the estimated variables to the input data. These two conditional probabilities are equal to the probabilities of the fitting residues  $N_{P \times K_{A,S,i}}$  and  $N_{G \times K_{R,S,i}}$ , for which we assumed zero-mean Gaussian distributions.

$$A|B, T \sim \text{normal}(\mu_{N_A}, \sigma_{N_A}^2), \sigma_{N_A}^2 \sim \text{inverse gamma}(\alpha_{N_A}, \beta_{N_A}), \tag{10}$$

$$R|L, B, T \sim \text{normal}(\mu_{N_R}, \sigma_{N_R}^2), \sigma_{N_R}^2 \sim \text{inverse gamma}(\alpha_{N_R}, \beta_{N_R}), \tag{11}$$

where  $(\mu_{N_A}, \sigma_{N_A}^2)$  and  $(\mu_{N_R}, \sigma_{N_R}^2)$  are hyperparameters representing the prior mean and variance of data noise in the ATAC and RNA measures. Here, the variance of the signal noise is modeled using inverse gamma distributions, with hyperparameters  $(\alpha_{N_A}, \beta_{N_A})$  and  $(\alpha_{N_R}, \beta_{N_R})$  to control the variance of fitting residues (very low probabilities on large variances).

Then, the posterior probability of each variable defined in equations (4)–(6) was still a Gaussian distribution with poster mean  $\hat{\mu}$  and variance  $\hat{\sigma}$  as shown below:

$$\hat{b}_{p,m,i} \sim \text{normal}(\hat{\mu}_{B,m,i}, \hat{\sigma}_{B,m,i}^2), \tag{12}$$

$$\hat{t}_{m,s,i} \sim \text{normal}(\hat{\mu}_{T,m,s,i}, \hat{\sigma}_{T,m,s,i}^2), \tag{13}$$

$$\hat{l}_{p,g,i} \sim \text{normal}(\hat{\mu}_{L,i}, \hat{\sigma}_{L,i}^2). \tag{14}$$

Gibbs sampling was used to iteratively learn the posterior distribution mean and variance of each set of variables and draw samples of their values accordingly.

For the TF–peak binding events, the posterior mean  $\hat{\mu}_{B,m,i}$  and variance  $\hat{\sigma}_{B,m,i}^2$  were estimated specifically for  $m$ -th TF since the number of binding sites and the positive or negative regulatory effects between TFs could be very different.

$$\hat{\mu}_{B,m,i} = \frac{\sum_s \sum_k t_{m,s,i} (a_{p,k_{A,S,i}} - \sum_m b_{p,m,i} t_{m,s,i}) \sigma_B^2 + \mu_{B,i} \sigma_{N_A}^2}{\sum_s \sum_k t_{m,s,i} \sigma_B^2 + \sigma_{N_A}^2} \text{ and} \tag{15}$$

$$\hat{\sigma}_{B,m,i}^2 = \frac{\sigma_{N_A}^2 \sigma_B^2}{\sum_s \sum_k t_{m,s,i} \sigma_B^2 + \sigma_{N_A}^2}.$$

For TF activities, the posterior mean  $\hat{\mu}_{T,m,s,i}$  and variance  $\hat{\sigma}_{T,m,s,i}^2$  were estimated specifically for the  $m$ -th TF and  $s$ -th sample using chromatin accessibility data as follows:

$$\hat{\mu}_{T,m,s,i} = \frac{\sum_p K_{A,s} b_{p,m} (a_{p,k,s,i} - \sum_m b_{p,m} t_{m,s}) \sigma_T^2 + \mu_T \sigma_{N_A}^2}{\sum_p K_{A,s} b_{p,m}^2 \sigma_T^2 + \sigma_{N_A}^2} \text{ and} \tag{16}$$

$$\hat{\sigma}_{T,m,s,i}^2 = \frac{\sigma_{N_A}^2 \sigma_T^2}{\sum_p K_{A,s} b_{p,m}^2 \sigma_T^2 + \sigma_{N_A}^2}.$$

Then, based on the estimated distribution parameters of  $\hat{\mu}_{T,m,s,i}$  and  $\hat{\sigma}_{T,m,s,i}^2$  of  $\hat{t}_{m,s,i}$  for the  $k_{R,s}$ -th RNA cell in the same  $s$ -th sample we draw a TF regulatory activity sample as  $t_{m,k_{R,s},s,i}$ . For  $p$ -th peak, we were able to reconstruct its chromatin activity in the RNA cell as  $\hat{a}_{p,k_{R,s},s,i} = \sum_m \hat{b}_{p,m,i} \hat{t}_{m,k_{R,s},s,i}$  and for  $g$ -th gene, we further estimated the interaction confidence  $\hat{l}_{p,g,i}$  between  $p$ -th peak and  $g$ -th gene. The peak–gene interaction distribution parameters  $\hat{\mu}_{L,i}$  and  $\hat{\sigma}_{L,i}^2$  were estimated as follows:

$$\hat{\mu}_{L,i} = \frac{\sum_s \sum_k \hat{a}_{p,k_{R,s},s,i} (t_{g,k_{R,s},s,i} - \sum_p \hat{t}_{p,g,i} \hat{a}_{p,k_{R,s},s,i}) \sigma_L^2 + \mu_L \sigma_{N_R}^2}{\sum_s \sum_k (\hat{a}_{p,k_{R,s},s,i})^2 \sigma_L^2 + \sigma_{N_R}^2} \text{ and} \tag{17}$$

$$\hat{\sigma}_{L,i}^2 = \frac{\sigma_{N_R}^2 \sigma_L^2}{\sum_s \sum_k (\hat{a}_{p,k_{R,s},s,i})^2 \sigma_L^2 + \sigma_{N_R}^2}.$$

In  $n$ -th round of Gibbs estimation, after learning all distributions, we estimated the confidence of each linkage by linearly mapping the sampled values of  $\hat{b}_{p,m,i}$  and  $\hat{l}_{p,g,i}$  in the range of  $(-\infty, \infty)$  to probabilities in  $(0,1)$  as follows:

$$P(\text{state}(b_{p,m,i}|n) = 1) = \frac{\exp\{(\hat{b}_{p,m,i} - \hat{\mu}_{B,m,i})/2\hat{\sigma}_{B,m,i}^2\}}{\exp\{(\hat{b}_{p,m,i} - \hat{\mu}_{B,m,i})/2\hat{\sigma}_{B,m,i}^2\} + \exp\{(0 - \hat{\mu}_{B,m,i})/2\hat{\sigma}_{B,m,i}^2\}}. \tag{18}$$

$$P(\text{state}(l_{p,g,i}|n) = 1) = \frac{\exp\{(\hat{l}_{p,g,i} - \hat{\mu}_{L,i})/2\hat{\sigma}_{L,i}^2\}}{\exp\{(\hat{l}_{p,g,i} - \hat{\mu}_{L,i})/2\hat{\sigma}_{L,i}^2\} + \exp\{(0 - \hat{\mu}_{L,i})/2\hat{\sigma}_{L,i}^2\}}. \tag{19}$$

Binary state samples were then drawn based on the confidence of each linkage and were then used to initiate the next round of estimations. After running a long sampling process (in total,  $N$  rounds) and accumulating enough samples on the binary states of TF–peak bindings and peak–gene interactions, we calculated the sampling frequency of each linkage as a posterior probability.

$$\begin{cases} P(\text{state}(b_{p,m,i}) = 1) = \frac{\sum_n \text{state}(b_{p,m,i}|n)}{N} \\ P(\text{state}(l_{p,g,i}) = 1) = \frac{\sum_n \text{state}(l_{p,g,i}|n)}{N} \end{cases} \tag{20}$$

**MAGICAL analysis of *S. aureus* single-cell multiomics data**

For each cell type, given DAS and DEG of contrast conditions (MRSA versus control, MSSA versus control or MRSA versus MSSA), MAGICAL was first initialized by mapping prior TF motifs from the chromVAR-motifs library to DAS using addMotifAnnotations (ArchR). Because there is no PBMC cell type Hi-C data publicly available, we are using TAD boundaries from a lymphoblastoid cell line, GM12878, which was originally generated by EBV transformation of PBMCs<sup>60</sup>. The TAD boundary structure is closely conserved between the lymphoblastoid cell lines and primary PBMC<sup>61</sup> and between cell types<sup>62,63</sup>. We called TAD boundaries from a GM12878 cell line Hi-C profile<sup>63</sup> using TopDom<sup>64</sup>. About 6,000 topological domains were identified. For each contrast, we built candidate circuits by pairing DAS with TF binding sites with DEG in the same domain. MAGICAL was run 10,000 times to ensure that the sampling process converged to stable states. This process was repeated for all cell types and the top 10% high confidence circuit predictions were selected from each cell type for validation analysis.

**MAGICAL analysis of COVID-19 single-cell multiomics data**

As a proof of concept for contrast condition, single-cell multiomics data analysis, MAGICAL was applied to a public PBMC COVID-19 single-cell multiomics dataset<sup>5</sup> with samples collected from patients with different severity and healthy controls. For each of the three selected cell subtypes (CD8 TEM, CD14 Mono, and NK), from the original publication we downloaded DEG for two contrasts: mild versus control and severe versus control. For each of the selected cell types, DAS were called respectively for mild versus control and severe versus control using ArchR functions and thresholds as introduced in the paper. MAGICAL was initialized by mapping TF motifs from the chromVARmotifs library to DAS using addMotifAnnotations (ArchR). As explained above, we used TAD boundary information of ~6,000 domains identified in the GM12878 cell line<sup>63</sup> as prior to pair DAS with TF binding sites and DEG. Then, the initial candidate regulatory circuits were constructed. Respectively for mild and severe COVID-19, MAGICAL was run 10,000 times to ensure that the sampling process converged to stable states. This process was repeated for the three selected cell types. The chromatin sites and genes in the top 10% predicted high confidence circuits in each cell type were selected as disease associated.

**COVID-19 PBMC samples of validation scATAC-seq data**

To validate chromatin sites associated with mild COVID-19, PBMC samples were obtained from the COVID-19 Health Action Response for Marines (CHARM) cohort study, which has been previously described<sup>65</sup>. The cohort is composed of Marine recruits who arrived at Marine Corps Recruit Depot–Parris Island for basic training between May and November 2020, after undergoing two quarantine periods (first a home quarantine, and next a supervised quarantine starting at enrollment in the CHARM study) to reduce the possibility of SARS-CoV-2 infection at arrival. Participants were regularly screened for SARS-CoV-2 infection during basic training by PCR, serum samples were obtained using serum separator tubes at all visits, and a follow-up symptom questionnaire was administered. At selected visits, blood was collected in BD Vacutainer CPT Tube with Sodium Heparin and PBMC were isolated following the manufacturer’s recommendations. We used PBMC samples from six participants (five males and one female) who had a positive COVID-19 PCR test and had mild symptoms (sampled 3–11 days after the first PCR positive test), and from three control participants (three males) who had a PCR negative test at the time of sample collection and were seronegative for SARS-CoV-2 immunoglobulin G. New scATAC-seq data were generated following the same protocol as described in “*S. aureus* scATAC-seq data generation” (Supplementary Table 2).

**COVID-19 PBMC scATACseq data analysis**

Reads of scATAC-seq experiments were aligned to human reference genome (hg38) using 10x Genomics Cell Ranger software (version 1.2). The resulting fragment files were processed using ArchR<sup>25</sup>. Quality cells were selected as those with TSS enrichment >12, a number of fragments >3,000 and <30,000, and a nucleosome ratio <2. The likelihood of doublet cells was computationally assessed using the addDoubletScores function and cells were filtered using the filterDoublets function with default settings. A total of 15,836 high-quality cells in the infection group and 9,125 cells in the control group were selected after QC analysis (Supplementary Fig. 3; Supplementary Table 3). These cells were combined into a linear dimensionality reduction using the addIterativeLSI function with the input of the tile matrix (read counts in binned 500 bp across the whole genome) with iterations = 2 and varFeatures = 20,000. The cells were then clustered using the addClusters function. We annotated scATAC-seq cells using the addGeneIntegrationMatrix function, referring to a labeled multimodal PBMC single cell dataset. Doublet clusters containing a mixture of many cell types were manually identified and removed.

Peaks were called for each cell type using the `addReproduciblePeakSet` function with peak caller MACS2<sup>26</sup> (Supplementary Fig. 3). In total, 284,525 peaks were identified (Supplementary Table 4). For each of the three selected cell types (CD8 TEM, CD14 Mono, and NK), chromatin sites with single cell differential statistics  $FDR < 0.05$  and  $|\log_2(FC)| > 0.1$  between COVID-19 and control conditions and actively accessible in at least 10% of cells ( $pct > 0.1$ ) from either condition were selected. Refined peaks passing pseudobulk differential statistics  $P < 0.05$  and  $|\log_2(FC)| > 0.3$  between the contrast conditions were finally selected as the validation peak set (Supplementary Table 5).

### COVID-19 circuit peaks and genes accuracy evaluation

The number of infection-associated peaks/genes reported by each COVID-19 study would be different owing to the difference in the number of recruited patients and collected cells. To overcome the issue caused by the imbalanced number between discovery and validation datasets or between differential peaks/genes and circuit sites/genes, in each comparison, the larger peak/gene set was randomly downsampled to match the smaller number of peaks/genes in the other set. The precision (site reproduction rate) is calculated to assess the accuracy of each peak/gene set.

### MAGICAL analysis of 10x PBMC single-cell true multiome data

For benchmarking, MAGICAL was applied to a 10x PBMC single-cell multiome dataset including 108,377 ATAC peaks, 36,601 genes, and 11,909 cells from 14 cell types. MAGICAL used the same candidate peaks and genes as selected by TRIPOD<sup>11</sup> for fair performance comparison. Two different priors were used to pair candidate peaks and genes: (1) the peaks and genes were within the same TAD from the GM12878 cell line; (2) the centers of peaks and the TSS of genes were within 500 kbp. MAGICAL inferred regulatory circuits with each prior and used the top 10% of predictions for accuracy assessment. High-confidence peak–gene interactions predicted by TRIPOD on the same data were directly downloaded from the supplementary tables of their publication<sup>11</sup>. Two baseline approaches of peak–gene pairing were included: pairing all peaks with a gene if they are in the same TAD or pairing only the nearest peak to a gene based on their genomic distance. To fairly assess the accuracy of MAGICAL weighted peak–gene interactions and the results (paired or non-paired) from TRIPOD or baseline approaches, we selected the top 10% of predictions by MAGICAL as the final peak–gene pairing. We overlapped these pairs with the curated 3D genome interactions in blood context from the 4DGenome database<sup>19</sup> and calculated the precision for each approach.

### MAGICAL analysis of GM12878 cell line SHARE-seq data

For benchmarking, MAGICAL was also applied to a GM12878 cell line SHARE-seq dataset<sup>10</sup>. For fair comparison, MAGICAL used the same candidate peaks and genes as selected by FigR<sup>18</sup>. MAGICAL was initialized with two different priors to pair candidate peaks and genes: (1) the peaks and genes were within the same prior TAD from the GM12878 cell line; (2) the centers of peaks and the TSS of genes were within 500 kbps. MAGICAL inferred regulatory circuits under each setting and used the top 10% predictions for accuracy assessment. High-confidence peak–gene interactions predicted by FigR were directly downloaded from the supplementary tables of the original publication<sup>10</sup>. Similarly, the top 10% predictions by MAGICAL and interactions paired by the two baseline approaches mentioned above were selected. We overlapped peak–gene interactions predicted by each approach with GM12878 H3K27ac HiChIP chromatin interactions<sup>20</sup> for precision evaluation.

### Validating predicted peak–gene interactions

To assess the precision of the predicted circuit peak–gene interactions, we assumed a correctly inferred peak–gene pair should be also connected by a chromatin interaction reported by Hi-C or similar experiments. To check this, each peak was extended to 2 kb long and then

checked for overlap with one end of a physical chromatin interaction. For genes, we checked whether the gene promoter (–2 kb to 500 b of TSS) overlapped the other end of the interaction. Precision was calculated as the proportion of overlapped chromatin interactions among the predicted peak–gene pairs. The significance of enrichment of overlapped chromatin interactions was assessed using hypergeometric  $P$  value, with all candidate peak–gene pairs as background.

### GWAS enrichment analysis

To assess the enrichment of GWAS loci of inflammatory diseases in circuit chromatin sites in each cell type, significant GWAS loci were downloaded from GWAS catalog<sup>46</sup> for inflammatory diseases and control diseases. GREGOR<sup>66</sup> was used to assess the enrichment of GWAS loci at which either the index single nucleotide polymorphism (SNP) or at least one of its LD proxies overlaps with a circuit chromatin site, using pre-calculated LD data from 1,000 G EUR samples. The enrichment  $P$  value of each disease GWAS was converted to a z-score. With each cell type, enrichment scores for traits with fewer than five overlapped GWAS SNPs with circuit sites were hold out. Also, as all reference data used by GREGOR is hg19-based, genome coordinates of testing regions were mapped from hg38 to hg19.

### Predicting *S. aureus* infection state

To refine circuit genes for predicting infection diagnosis in microarray gene expression data, the capability of each circuit gene on distinguishing infection and control samples, or MRSA and MSSA samples, was assessed using sample level pseudobulk gene expression data, aggregated from the discovery scRNA-seq datasets. The total number of reads of each sample was normalized to  $1 \times 10^7$ . The normalized RNA read counts across all samples were then log and z-score transformed. For each circuit gene, a discovery AUROC was calculated by comparing the scRNA-seq gene-expression-based sample ranking against the contrasted sample groups. Circuit genes were prioritized based on AUROCs. An SVM model was trained using the top-ranked circuit genes as features and their normalized pseudobulk expression data as input. The model was then tested on independent microarray datasets. The microarray gene expression data was also log and z-score transformed to ensure a similar distribution to the training data. For comparison, top DEG prioritized by discovery AUROC or by other approaches like the minimum redundancy maximum relevance algorithm or LASSO regression were also tested on the same microarray datasets.

### Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

### Data availability

The 10x PBMC single cell multiome dataset can be downloaded from [https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/1.0.0/pbmc\\_granulocyte\\_sorted\\_10k](https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/1.0.0/pbmc_granulocyte_sorted_10k). Users will need to provide their contact information to access the download webpage where the filtered feature barcode matrix (HDF5 format) can be downloaded. The reference multimodal PBMC single cell dataset (H5Seurat data file) can be downloaded from <https://atlas.fredhutch.org/nygc/multimodal-pbmc/>. The GWAS catalog database can be accessed at <https://www.ebi.ac.uk/gwas/docs/file-downloads>. SNPs associated with each disease used in this paper can be extracted from the downloadable file “All associations v1.0”. Home sapiens chromatin interactions data can be downloaded from <https://4dgenome.research.chop.edu/Download.html>. Home sapiens TF ChIP-seq profiles can be downloaded at <http://cistrome.org/db/>. Users can also provide their customized peaks in BED format to the server <http://dbtoolkit.cistrome.org/> and identify TFs that have a significant binding overlap. Home sapiens candidate enhancers annotated by ENCODE can be downloaded at <https://screen.encodeproject.org/>. The chromVAR motifs library is available at

<https://github.com/GreenleafLab/chromVARmotifs>. The source single cell data collected in this study is publicly accessible at the GEO repository (<https://www.ncbi.nlm.nih.gov/geo/>, accession no. GSE220190) and the Zenodo repository<sup>67</sup>. Source data for Figs. 2–4 is available with this manuscript.

## Code availability

The source code of MAGICAL is available on GitHub at <https://github.com/xichensf/magical> and the Zenodo repository<sup>68</sup>.

## References

- Schoenfelder, S. & Fraser, P. Long-range enhancer-promoter contacts in gene expression control. *Nat. Rev. Genet.* **20**, 437–455 (2019).
- Kim, H. D., Shay, T., O’Shea, E. K. & Regev, A. Transcriptional regulatory circuits: predicting numbers from alphabets. *Science* **325**, 429–432 (2009).
- modENCODE Consortium et al. Identification of functional elements and regulatory circuits by *Drosophila* modENCODE. *Science* **330**, 1787–1797 (2010).
- Marbach, D. et al. Tissue-specific regulatory circuits reveal variable modular perturbations across complex diseases. *Nat. Methods* **13**, 366–370 (2016).
- Wilk, A. J. et al. Multi-omic profiling reveals widespread dysregulation of innate immunity and hematopoiesis in COVID-19. *J. Exp. Med.* **218**, e20210582 (2021).
- Krijger, P. H. & de Laat, W. Regulation of disease-associated gene expression in the 3D genome. *Nat. Rev. Mol. Cell Biol.* **17**, 771–782 (2016).
- Cao, J. et al. Joint profiling of chromatin accessibility and gene expression in thousands of single cells. *Science* **361**, 1380–1385 (2018).
- Kreitmaier, P., Katsoula, G. & Zeggini, E. Insights from multi-omics integration in complex disease primary tissues. *Trends Genet.* **39**, 46–58 (2022).
- Stuart, T. et al. Comprehensive integration of single-cell data. *Cell* **177**, 1888–1902 (2019).
- Ma, S. et al. Chromatin potential identified by shared single-cell profiling of RNA and chromatin. *Cell* **183**, 1103–1116 (2020).
- Jiang, Y. et al. Nonparametric single-cell multiomic characterization of trio relationships between transcription factors, target genes, and cis-regulatory regions. *Cell Syst.* **13**, 737–751 (2022).
- Cao, Z. J. & Gao, G. Multi-omics single-cell data integration and regulatory inference with graph-linked embedding. *Nat. Biotechnol.* **40**, 1458–1466 (2022).
- Javierre, B. M. et al. Lineage-specific genome architecture links enhancers and non-coding disease variants to target gene promoters. *Cell* **167**, 1369–1384 (2016).
- Arnold, S. R. et al. Changing patterns of acute hematogenous osteomyelitis and septic arthritis: emergence of community-associated methicillin-resistant *Staphylococcus aureus*. *J. Pediatr. Orthop.* **26**, 703–708 (2006).
- Saavedra-Lozano, J. et al. Changing trends in acute osteomyelitis in children: impact of methicillin-resistant *Staphylococcus aureus* infections. *J. Pediatr. Orthop.* **28**, 569–575 (2008).
- Liao, J. C. et al. Network component analysis: reconstruction of regulatory signals in biological systems. *Proc. Natl Acad. Sci. USA* **100**, 15522–15527 (2003).
- Tran, L. M., Brynildsen, M. P., Kao, K. C., Suen, J. K. & Liao, J. C. gNCA: a framework for determining transcription factor activity based on transcriptome: identifiability and numerical implementation. *Metab. Eng.* **7**, 128–141 (2005).
- Kartha, V. K. et al. Functional inference of gene regulation using single-cell multi-omics. *Cell Genom.* **2**, 100166 (2022).
- Teng, L., He, B., Wang, J. & Tan, K. 4DGenome: a comprehensive database of chromatin interactions. *Bioinformatics* **31**, 2560–2564 (2015).
- Mumbach, M. R. et al. Enhancer connectome in primary human cells identifies target genes of disease-associated DNA elements. *Nat. Genet.* **49**, 1602–1612 (2017).
- Arunachalam, P. S. et al. Systems biological assessment of immunity to mild versus severe COVID-19 infection in humans. *Science* **369**, 1210–1220 (2020).
- Lucas, C. et al. Longitudinal analyses reveal immunological misfiring in severe COVID-19. *Nature* **584**, 463–469 (2020).
- Mathew, D. et al. Deep immune profiling of COVID-19 patients reveals distinct immunotypes with therapeutic implications. *Science* **369**, eabc8511 (2020).
- Schulte-Schrepping, J. et al. Severe COVID-19 is marked by a dysregulated myeloid cell compartment. *Cell* **182**, 1419–1440 (2020).
- Granja, J. M. et al. ArchR is a scalable software package for integrative single-cell chromatin accessibility analysis. *Nat. Genet.* **53**, 403–411 (2021).
- Feng, J., Liu, T., Qin, B., Zhang, Y. & Liu, X. S. Identifying ChIP-seq enrichment using MACS. *Nat. Protoc.* **7**, 1728–1740 (2012).
- Li, S. et al. Epigenetic landscapes of single-cell chromatin accessibility and transcriptomic immune profiles of T cells in COVID-19 patients. *Front Immunol.* **12**, 625881 (2021).
- Jung, I. et al. A compendium of promoter-centered long-range chromatin interactions in the human genome. *Nat. Genet.* **51**, 1442–1449 (2019).
- Chen, X. et al. Tissue-specific enhancer functional networks for associating distal regulatory regions to disease. *Cell Syst.* **12**, 353–362 (2021).
- Yao, C. et al. Cell-type-specific immune dysregulation in severely ill COVID-19 patients. *Cell Rep.* **34**, 108590 (2021).
- Unterman, A. et al. Single-cell multi-omics reveals dyssynchrony of the innate and adaptive immune system in progressive COVID-19. *Nat. Commun.* **13**, 440 (2022).
- Magill, S. S. et al. Changes in prevalence of health care-associated infections in U.S. hospitals. *N. Engl. J. Med.* **379**, 1732–1744 (2018).
- Tong, S. Y., Davis, J. S., Eichenberger, E., Holland, T. L. & Fowler, V. G. Jr. *Staphylococcus aureus* infections: epidemiology, pathophysiology, clinical manifestations, and management. *Clin. Microbiol. Rev.* **28**, 603–661 (2015).
- Marquez-Ortiz, R. A. et al. USA300-related methicillin-resistant *Staphylococcus aureus* clone is the predominant cause of community and hospital MRSA infections in Colombian children. *Int J. Infect. Dis.* **25**, 88–93 (2014).
- Hao, Y. et al. Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573–3587 (2021).
- Skjeflo, E. W., Christiansen, D., Espevik, T., Nielsen, E. W. & Mollnes, T. E. Combined inhibition of complement and CD14 efficiently attenuated the inflammatory response induced by *Staphylococcus aureus* in a human whole blood model. *J. Immunol.* **192**, 2857–2864 (2014).
- Kusunoki, T., Hailman, E., Juan, T. S., Lichenstein, H. S. & Wright, S. D. Molecules from *Staphylococcus aureus* that bind CD14 and stimulate innate immune responses. *J. Exp. Med.* **182**, 1673–1682 (1995).
- Ludwig, S. et al. Influenza virus-induced AP-1-dependent gene expression requires activation of the JNK signaling pathway. *J. Biol. Chem.* **276**, 10990–10998 (2001).
- Gjertsson, I., Hultgren, O. H., Collins, L. V., Pettersson, S. & Tarkowski, A. Impact of transcription factors AP-1 and NF- $\kappa$ B on the outcome of experimental *Staphylococcus aureus* arthritis and sepsis. *Microbes Infect.* **3**, 527–534 (2001).

40. Liu, T. et al. Cistrome: an integrative platform for transcriptional regulation studies. *Genome Biol.* **12**, R83 (2011).
41. Gillespie, M. et al. The reactome pathway knowledgebase 2022. *Nucleic Acids Res.* **50**, D687–D692 (2022).
42. Kyriakis, J. M. Activation of the AP-1 transcription factor by inflammatory cytokines of the TNF family. *Gene Expr.* **7**, 217–231 (1999).
43. Hannemann, N. et al. The AP-1 transcription factor c-Jun promotes arthritis by regulating cyclooxygenase-2 and arginase-1 expression in macrophages. *J. Immunol.* **198**, 3605–3614 (2017).
44. Gasperini, M. et al. A genome-wide framework for mapping gene regulation via cellular genetic screens. *Cell* **176**, 377–390 (2019).
45. Consortium, E. P. et al. Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature* **583**, 699–710 (2020).
46. Buniello, A. et al. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* **47**, D1005–D1012 (2019).
47. DeLorenze, G. N. et al. Polymorphisms in HLA class II genes are associated with susceptibility to *Staphylococcus aureus* infection in a white population. *J. Infect. Dis.* **213**, 816–823 (2016).
48. Chen, L. et al. Genetic drivers of epigenetic and transcriptional variation in human immune cells. *Cell* **167**, 1398–1414 (2016).
49. Ahn, S. H. et al. Gene expression-based classifiers identify *Staphylococcus aureus* infection in mice and humans. *PLoS One* **8**, e48979 (2013).
50. Ramilo, O. et al. Gene expression patterns in blood leukocytes discriminate patients with acute infections. *Blood* **109**, 2066–2077 (2007).
51. Ardura, M. I. et al. Enhanced monocyte response and decreased central memory T cells in children with invasive *Staphylococcus aureus* infections. *PLoS One* **4**, e5446 (2009).
52. Cho, J. S. et al. IL-17 is essential for host defense against cutaneous *Staphylococcus aureus* infection in mice. *J. Clin. Invest.* **120**, 1762–1773 (2010).
53. Xiao, Y. et al. A novel significance score for gene selection and ranking. *Bioinformatics* **30**, 801–807 (2014).
54. Chaussabel, D. et al. A modular analysis framework for blood genomics studies: application to systemic lupus erythematosus. *Immunity* **29**, 150–164 (2008).
55. Wenric, S. & Shemirani, R. Using supervised learning methods for gene selection in RNA-Seq case-control studies. *Front. Genet.* **9**, 297 (2018).
56. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
57. Korsunsky, I. et al. Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat. Methods* **16**, 1289–1296 (2019).
58. Squair, J. W. et al. Confronting false discoveries in single-cell differential expression. *Nat. Commun.* **12**, 5692 (2021).
59. Schep, A. N., Wu, B., Buenrostro, J. D. & Greenleaf, W. J. chromVAR: inferring transcription-factor-associated accessibility from single-cell epigenomic data. *Nat. Methods* **14**, 975–978 (2017).
60. Anderson, M. A. & Gusella, J. F. Use of cyclosporin A in establishing Epstein-Barr virus-transformed human lymphoblastoid cell lines. *Vitro* **20**, 856–858 (1984).
61. Tan, L., Xing, D., Chang, C. H., Li, H. & Xie, X. S. Three-dimensional genome structures of single diploid human cells. *Science* **361**, 924–928 (2018).
62. McArthur, E. & Capra, J. A. Topologically associating domain boundaries that are stable across diverse cell types are evolutionarily constrained and enriched for heritability. *Am. J. Hum. Genet.* **108**, 269–283 (2021).
63. Rao, S. S. et al. A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell* **159**, 1665–1680 (2014).
64. Shin, H. et al. TopDom: an efficient and deterministic method for identifying topological domains in genomes. *Nucleic Acids Res.* **44**, e70 (2016).
65. Letizia, A. G. et al. SARS-CoV-2 seropositivity and subsequent infection risk in healthy young adults: a prospective cohort study. *Lancet Respir. Med.* **9**, 712–720 (2021).
66. Schmidt, E. M. et al. GREGOR: evaluating global enrichment of trait-associated variants in epigenomic features using a systematic, data-driven approach. *Bioinformatics* **31**, 2601–2606 (2015).
67. Chen, X. Source data for paper “Mapping disease regulatory circuits at cell-type resolution from single-cell multiomics data”. Zenodo <https://doi.org/10.5281/zenodo.7992711> (2023).
68. Chen, X. MAGICAL (v1.1). Zenodo <https://doi.org/10.5281/zenodo.7951577> (2023).

## Acknowledgements

We thank the Single-Cell and Spatial Technologies team at the Center for Advanced Genomics Technology, Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai for providing the experimental, computational, data resources, and staff expertise. S.C.S. was supported by the Defense Advanced Research Projects Agency under contract N6600119C4022. A.G.L. is supported by Defense Health Agency grant 9700130 through the Naval Medical Research Center. O.G.T. is supported by the National Institutes of Health under grant R01GM071966 and Simons Foundation under grant 395506. L.C.N. is supported by the National Institute of Mental Health, the National Institute of Neurological Disorders and Stroke, the National Institute of Diabetes and Digestive and Kidney Diseases, the National Heart, Lung, and Blood Institute, and the National Institute of Allergy and Infectious Diseases under grant UM1AI164559 and by the National Institute on Drug Abuse under grants U01 DA53625 and U01DA058527.

## Author contributions

S.C.S., O.G.T., and E.Z. conceived the study and supervised the research. X.C. designed and implemented the computational framework, conducted benchmarks and case studies with Y.W., wrote the code, and set up the web access with the help of A.T.; A.C. was involved in the *S. aureus* study. W.-S.C. managed and processed single-cell sequencing data with help from A.B.R., G.N., and A.V.; S.H.K. and D.G.C. conducted the public microarray data search. F.R.Z., V.D.N., M.-C.G., and R.S. generated the PBMC single-cell multiomics data for the *S. aureus* infected and control subjects. The *S. aureus* patient blood samples were provided by C.W.W., V.G.F., F.R., and M.D. The control samples were provided by R.R.S. and L.C.N.; I.R. and C.M.M. provided immunological interpretations of the results. I.R., A.G.L., and A.S.-S. provided the validation PBMC scATAC-seq data of patients with COVID-19 and uninfected controls. S.C.S., O.G.T., X.C., E.Z., A.C., and C.L.T. wrote the first draft of the manuscript. All authors proofread the submitted version.

## Competing interests

A.G.L. is a military service member. This work was prepared as part of his official duties. Title 17, US Code §105 provides that copyright protection under this title is not available for any work of the US Government. Title 17, US code §101 defines a US Government work as a work prepared by a military service member or employee of the US Government as part of that person's official duties. The views expressed in the article are those of the authors and do not necessarily express the official policy and position of the US Navy, the Department of Defense, the US Government, or the

institutions affiliated with the authors. V.G.F. reports personal fees from Novartis, Debiopharm, Genentech, Achaogen, Affinium, Medicines Co., MedImmune, Bayer, Basilea, Affinergy, Janssen, Contrafect, Regeneron, Destiny, Amphlphi Biosciences, Integrated Biotherapeutics; C3J, Armata, Valanbio; Akagera, Aridis, Roche, grants from NIH, MedImmune, Allergan, Pfizer, Advanced Liquid Logics, Theravance, Novartis, Merck; Medical Biosurfaces; Locus; Affinergy; Contrafect; Karius; Genentech, Regeneron, Deep Blue, Basilea, Janssen; Royalties from UpToDate, stock options from Valanbio and ArcBio, Honoraria from Infectious Diseases of America for his service as Associate Editor of Clinical Infectious Diseases, and a patent sepsis diagnostics pending. L.C.N. has received consulting fees from work as a scientific advisor for AbbVie, ViiV Healthcare, and Cytodyn and also serves on the Board of Directors of CytoDyn and has financial interests in Ledidi AS, all for work outside of the submitted work. S.C.S. is a founder of GNOMX Corp and serves as chief scientific officer. The remaining authors declare no competing interests.

### Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s43588-023-00476-5>.

**Correspondence and requests for materials** should be addressed to Elena Zaslavsky, Olga G. Troyanskaya or Stuart C. Sealfon.

**Peer review information** *Nature Computational Science* thanks the anonymous reviewers for their contribution to the peer review of this work. Primary Handling Editor Kaitlin McCardle, in collaboration with the *Nature Computational Science* team.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023, corrected publication 2023

<sup>1</sup>Center for Computational Biology, Flatiron Institute, New York, NY, USA. <sup>2</sup>Lewis-Sigler Institute of Integrative Genomics, Princeton University, Princeton, NJ, USA. <sup>3</sup>Department of Computer Science, Princeton University, Princeton, NJ, USA. <sup>4</sup>Department of Neurology, Icahn School of Medicine at Mount Sinai, New York, NY, USA. <sup>5</sup>Division of Infectious Diseases, Department of Medicine, Duke University School of Medicine, Durham, NC, USA. <sup>6</sup>Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT, USA. <sup>7</sup>Battelle Memorial Institute, Columbus, OH, USA. <sup>8</sup>Division of Infectious Diseases, Department of Medicine, Weill Cornell Medicine, New York, NY, USA. <sup>9</sup>Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA. <sup>10</sup>Department of Pathology and Department of Immunobiology, Yale School of Medicine, New Haven, CT, USA. <sup>11</sup>Naval Medical Research Center, Silver Spring, MD, USA. <sup>12</sup>Precision Immunology Institute, Icahn School of Medicine at Mount Sinai, New York, NY, USA. <sup>13</sup>These authors jointly supervised this work: Elena Zaslavsky, Olga G. Troyanskaya, Stuart C. Sealfon. ✉e-mail: [elena.zaslavsky@mssm.edu](mailto:elena.zaslavsky@mssm.edu); [ogt@genomics.princeton.edu](mailto:ogt@genomics.princeton.edu); [stuart.sealfon@mssm.edu](mailto:stuart.sealfon@mssm.edu)

## Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

### Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement
- A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided  
*Only common tests should be described solely by name; describe more complex techniques in the Methods section.*
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g.  $F$ ,  $t$ ,  $r$ ) with confidence intervals, effect sizes, degrees of freedom and  $P$  value noted  
*Give  $P$  values as exact values whenever suitable.*
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's  $d$ , Pearson's  $r$ ), indicating how they were calculated

*Our web collection on [statistics for biologists](#) contains articles on many of the points above.*

### Software and code

Policy information about [availability of computer code](#)

Data collection

Data analysis

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio [guidelines for submitting code & software](#) for further information.

### Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our [policy](#)

The 10X PBMC single cell multiome dataset can be downloaded from [https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/1.0.0/pbmc\\_granulocyte\\_sorted\\_10k](https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets/1.0.0/pbmc_granulocyte_sorted_10k). Users will need to provide their contact information to access the download webpage where the filtered feature barcode matrix (HDF5 format) can be downloaded. The reference multimodal PBMC single cell dataset (H5 Seurat data file) can be downloaded from <https://atlas.fredhutch.org/nygc/multimodal-pbmc/>. The GWAS catalog database can be accessed at <https://www.ebi.ac.uk/gwas/docs/file-downloads>. SNPs associated with each disease used

in this paper can be extracted from the downloadable file "All associations v1.0". Home sapiens chromatin interactions data can be downloaded from <https://4dgenome.research.chop.edu/Download.html>. Home sapiens transcription factor ChIP-seq profiles can be downloaded at <http://cistrome.org/db/>. Users can also provide their customized peaks in BED format to the server <http://dbtoolkit.cistrome.org/> and identify transcription factors that have a significant binding overlap. Home sapiens candidate enhancers annotated by ENCODE can be downloaded at <https://screen.encodeproject.org/>. The chromVARmotifs library is available at <https://github.com/GreenleafLab/chromVARmotifs>. The source single cell data collected in this study is publicly accessible at the GEO repository (<https://www.ncbi.nlm.nih.gov/geo/>, accession no. GSE220190) and the Zenodo repository. Source data for Figures 2-4 is available with this manuscript.

## Human research participants

Policy information about [studies involving human research participants and Sex and Gender in Research](#).

Reporting on sex and gender	We are studying the gene regulatory mechanisms in human immune system in response to <i>S.aureus</i> infection or COVID-19 infection. Therefore, we did not use sex information of the participants in data analysis.
Population characteristics	Meta information of human subjects involved in this study including sex, age, and infection state are provided in Supplementary Tables 2 (COVID-19) and 6 ( <i>S.aureus</i> ).
Recruitment	<p>In the <i>S. aureus</i> study, Patients with culture-confirmed <i>S. aureus</i> bloodstream infection transferred to DUMC are eligible if pathogen speciation and antibiotic susceptibilities are confirmed by the Duke Clinical Microbiology Laboratory. We excluded samples if prior enrollment of the patient in this investigation (to ensure statistical independence of observations) or they are polymicrobial (i.e., more than one organism in blood or urine culture). In total, 21 adult patients were selected with 10 MRSAs and 11 MSSAs. None of them received any antibiotics in the 24 h before the bloodstream infection. Control samples were obtained from uninfected healthy adults matching the sample number and age range of the patient group. In total, 23 samples were collected from two cohorts: 14 controls provided by from the Weill Cornell Medicine, New York, NY, and 9 controls (provided by the Battelle Memorial Institute, Columbus, OH).</p> <p>In the COVID-19 study, PBMC samples were obtained from the COVID-19 Health Action Response for Marines (CHARM) cohort study. The cohort is composed of Marine recruits that arrived at Marine Corps Recruit Depot—Parris Island (MCRDPI) for basic training between May and November 2020, after undergoing two quarantine periods (first a home-quarantine, and next a supervised quarantine starting at enrolment in the CHARM study) to reduce the possibility of SARS-CoV-2 infection at arrival. Participants were regularly screened for SARS-CoV-2 infection during basic training by PCR, serum samples were obtained using serum separator tubes (SST) at all visits, and a follow-up symptom questionnaire was administered. At selected visits, blood was collected in BD Vacutainer CPT Tube with Sodium Heparin and PBMC were isolated following the manufacturer's recommendations. We used PBMC samples from six participants (five males and one female) who had a COVID-19 PCR positive test and had mild symptoms (sampled 3-11 days after the first PCR positive test), and from three control participants (three males) that had a PCR negative test at the time of sample collection and were seronegative for SARS-CoV-2 IgG.</p> <p>Subjects provided written informed consent prior to participation.</p>
Ethics oversight	The staphylococcus sepsis protocol was reviewed and approved by the Duke Medical School institutional review board (protocol number Pro00102421). The Navy COVID-19 study protocol was approved by the Naval Medical Research Center institutional review board (protocol number NMRC.2020.0006) in compliance with all applicable Federal regulations governing the protection of human subjects.

Note that full information on the approval of the study protocol must also be provided in the manuscript.

## Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences       Behavioural & social sciences       Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://nature.com/documents/nr-reporting-summary-flat.pdf)

## Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	No sample size calculation was performed. In the <i>S. aureus</i> infection study, we collected PBMC samples from 10 subjects diagnosed with MRSA, from 11 subjects diagnosed with MSSA and from 23 control samples for contrast analysis. In the COVID-19 study, we collected PBMC samples from 6 subjects with a COVID-19 PCR positive test and mild symptoms (sampled 3-11 days after the first PCR positive test), and from 3 control participants (three males) that had a PCR negative test at the time of sample collection.
Data exclusions	No data was excluded in data analysis.
Replication	We assessed our findings on different datasets. All replication attempts were successful.



Randomization	In the <i>S. aureus</i> study, samples are grouped based on the diagnosed infection state (MRSA, MSSA, or Control). In the COVID-19 study, samples are grouped based on the symptoms severity (severe-COVID-19, mild-COVID-19 or PCR negative).
Blinding	In both <i>S. aureus</i> and COVID-19 studies, we study the varying regulatory circuits between different conditions. Therefore, disease-state-based grouping in the discovery dataset is important to train the model. When we validate our model on public data, the grouping information is blinded for testing the model prediction accuracy.

## Behavioural & social sciences study design

All studies must disclose on these points even when the disclosure is negative.

Study description	Briefly describe the study type including whether data are quantitative, qualitative, or mixed-methods (e.g. qualitative cross-sectional, quantitative experimental, mixed-methods case study).
Research sample	State the research sample (e.g. Harvard university undergraduates, villagers in rural India) and provide relevant demographic information (e.g. age, sex) and indicate whether the sample is representative. Provide a rationale for the study sample chosen. For studies involving existing datasets, please describe the dataset and source.
Sampling strategy	Describe the sampling procedure (e.g. random, snowball, stratified, convenience). Describe the statistical methods that were used to predetermine sample size OR if no sample-size calculation was performed, describe how sample sizes were chosen and provide a rationale for why these sample sizes are sufficient. For qualitative data, please indicate whether data saturation was considered, and what criteria were used to decide that no further sampling was needed.
Data collection	Provide details about the data collection procedure, including the instruments or devices used to record the data (e.g. pen and paper, computer, eye tracker, video or audio equipment) whether anyone was present besides the participant(s) and the researcher, and whether the researcher was blind to experimental condition and/or the study hypothesis during data collection.
Timing	Indicate the start and stop dates of data collection. If there is a gap between collection periods, state the dates for each sample cohort.
Data exclusions	If no data were excluded from the analyses, state so OR if data were excluded, provide the exact number of exclusions and the rationale behind them, indicating whether exclusion criteria were pre-established.
Non-participation	State how many participants dropped out/declined participation and the reason(s) given OR provide response rate OR state that no participants dropped out/declined participation.
Randomization	If participants were not allocated into experimental groups, state so OR describe how participants were allocated to groups, and if allocation was not random, describe how covariates were controlled.

## Ecological, evolutionary & environmental sciences study design

All studies must disclose on these points even when the disclosure is negative.

Study description	Briefly describe the study. For quantitative data include treatment factors and interactions, design structure (e.g. factorial, nested, hierarchical), nature and number of experimental units and replicates.
Research sample	Describe the research sample (e.g. a group of tagged <i>Passer domesticus</i> , all <i>Stenocereus thurberi</i> within Organ Pipe Cactus National Monument), and provide a rationale for the sample choice. When relevant, describe the organism taxa, source, sex, age range and any manipulations. State what population the sample is meant to represent when applicable. For studies involving existing datasets, describe the data and its source.
Sampling strategy	Note the sampling procedure. Describe the statistical methods that were used to predetermine sample size OR if no sample-size calculation was performed, describe how sample sizes were chosen and provide a rationale for why these sample sizes are sufficient.
Data collection	Describe the data collection procedure, including who recorded the data and how.
Timing and spatial scale	Indicate the start and stop dates of data collection, noting the frequency and periodicity of sampling and providing a rationale for these choices. If there is a gap between collection periods, state the dates for each sample cohort. Specify the spatial scale from which the data are taken
Data exclusions	If no data were excluded from the analyses, state so OR if data were excluded, describe the exclusions and the rationale behind them, indicating whether exclusion criteria were pre-established.
Reproducibility	Describe the measures taken to verify the reproducibility of experimental findings. For each experiment, note whether any attempts to repeat the experiment failed OR state that all attempts to repeat the experiment were successful.
Randomization	Describe how samples/organisms/participants were allocated into groups. If allocation was not random, describe how covariates were controlled. If this is not relevant to your study, explain why.

## Blinding

Describe the extent of blinding used during data acquisition and analysis. If blinding was not possible, describe why OR explain why blinding was not relevant to your study.

Did the study involve field work?  Yes  No

## Field work, collection and transport

## Field conditions

Describe the study conditions for field work, providing relevant parameters (e.g. temperature, rainfall).

## Location

State the location of the sampling or experiment, providing relevant parameters (e.g. latitude and longitude, elevation, water depth).

## Access &amp; import/export

Describe the efforts you have made to access habitats and to collect and import/export your samples in a responsible manner and in compliance with local, national and international laws, noting any permits that were obtained (give the name of the issuing authority, the date of issue, and any identifying information).

## Disturbance

Describe any disturbance caused by the study and how it was minimized.

## Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

### Materials & experimental systems

- |                                     |  |
|-------------------------------------|--|
| n/a                                 | Included in the study                                  |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Antibodies                    |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Eukaryotic cell lines         |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Palaeontology and archaeology |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Animals and other organisms   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Clinical data                 |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Dual use research of concern  |

### Methods

- |                                     |   |
|-------------------------------------|---|
| n/a                                 | Included in the study                           |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> ChIP-seq               |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Flow cytometry         |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> MRI-based neuroimaging |

## Antibodies

## Antibodies used

Describe all antibodies used in the study; as applicable, provide supplier name, catalog number, clone name, and lot number.

## Validation

Describe the validation of each primary antibody for the species and application, noting any validation statements on the manufacturer's website, relevant citations, antibody profiles in online databases, or data provided in the manuscript.

## Eukaryotic cell lines

Policy information about [cell lines and Sex and Gender in Research](#)

## Cell line source(s)

State the source of each cell line used and the sex of all primary cell lines and cells derived from human participants or vertebrate models.

## Authentication

Describe the authentication procedures for each cell line used OR declare that none of the cell lines used were authenticated.

## Mycoplasma contamination

Confirm that all cell lines tested negative for mycoplasma contamination OR describe the results of the testing for mycoplasma contamination OR declare that the cell lines were not tested for mycoplasma contamination.

Commonly misidentified lines  
(See [ICLAC](#) register)

Name any commonly misidentified cell lines used in the study and provide a rationale for their use.

## Palaeontology and Archaeology

## Specimen provenance

Provide provenance information for specimens and describe permits that were obtained for the work (including the name of the issuing authority, the date of issue, and any identifying information). Permits should encompass collection and, where applicable, export.

Specimen deposition

Dating methods

Tick this box to confirm that the raw and calibrated dates are available in the paper or in Supplementary Information.

Ethics oversight

Note that full information on the approval of the study protocol must also be provided in the manuscript.

## Animals and other research organisms

Policy information about [studies involving animals](#); [ARRIVE guidelines](#) recommended for reporting animal research, and [Sex and Gender in Research](#)

Laboratory animals

Wild animals

Reporting on sex

Field-collected samples

Ethics oversight

Note that full information on the approval of the study protocol must also be provided in the manuscript.

## Clinical data

Policy information about [clinical studies](#)

All manuscripts should comply with the ICMJE [guidelines for publication of clinical research](#) and a completed [CONSORT checklist](#) must be included with all submissions.

Clinical trial registration

Study protocol

Data collection

Outcomes

## Dual use research of concern

Policy information about [dual use research of concern](#)

### Hazards

Could the accidental, deliberate or reckless misuse of agents or technologies generated in the work, or the application of information presented in the manuscript, pose a threat to:

No	Yes	
<input type="checkbox"/>	<input type="checkbox"/>	Public health
<input type="checkbox"/>	<input type="checkbox"/>	National security
<input type="checkbox"/>	<input type="checkbox"/>	Crops and/or livestock
<input type="checkbox"/>	<input type="checkbox"/>	Ecosystems
<input type="checkbox"/>	<input type="checkbox"/>	Any other significant area

## Experiments of concern

Does the work involve any of these experiments of concern:

- | No                       | Yes                      |   |
|--------------------------|--------------------------|---|
| <input type="checkbox"/> | <input type="checkbox"/> | Demonstrate how to render a vaccine ineffective                             |
| <input type="checkbox"/> | <input type="checkbox"/> | Confer resistance to therapeutically useful antibiotics or antiviral agents |
| <input type="checkbox"/> | <input type="checkbox"/> | Enhance the virulence of a pathogen or render a nonpathogen virulent        |
| <input type="checkbox"/> | <input type="checkbox"/> | Increase transmissibility of a pathogen                                     |
| <input type="checkbox"/> | <input type="checkbox"/> | Alter the host range of a pathogen  |
| <input type="checkbox"/> | <input type="checkbox"/> | Enable evasion of diagnostic/detection modalities                           |
| <input type="checkbox"/> | <input type="checkbox"/> | Enable the weaponization of a biological agent or toxin                     |
| <input type="checkbox"/> | <input type="checkbox"/> | Any other potentially harmful combination of experiments and agents         |

## ChIP-seq

### Data deposition

- Confirm that both raw and final processed data have been deposited in a public database such as [GEO](#).
- Confirm that you have deposited or provided access to graph files (e.g. BED files) for the called peaks.

#### Data access links

May remain private before publication.

For "Initial submission" or "Revised version" documents, provide reviewer access links. For your "Final submission" document, provide a link to the deposited data.

#### Files in database submission

Provide a list of all files available in the database submission.

#### Genome browser session

(e.g. [UCSC](#))

Provide a link to an anonymized genome browser session for "Initial submission" and "Revised version" documents only, to enable peer review. Write "no longer applicable" for "Final submission" documents.

## Methodology

### Replicates

Describe the experimental replicates, specifying number, type and replicate agreement.

### Sequencing depth

Describe the sequencing depth for each experiment, providing the total number of reads, uniquely mapped reads, length of reads and whether they were paired- or single-end.

### Antibodies

Describe the antibodies used for the ChIP-seq experiments; as applicable, provide supplier name, catalog number, clone name, and lot number.

### Peak calling parameters

Specify the command line program and parameters used for read mapping and peak calling, including the ChIP, control and index files used.

### Data quality

Describe the methods used to ensure data quality in full detail, including how many peaks are at FDR 5% and above 5-fold enrichment.

### Software

Describe the software used to collect and analyze the ChIP-seq data. For custom code that has been deposited into a community repository, provide accession details.

## Flow Cytometry

### Plots

Confirm that:

- The axis labels state the marker and fluorochrome used (e.g. CD4-FITC).
- The axis scales are clearly visible. Include numbers along axes only for bottom left plot of group (a 'group' is an analysis of identical markers).
- All plots are contour plots with outliers or pseudocolor plots.
- A numerical value for number of cells or percentage (with statistics) is provided.

### Methodology

#### Sample preparation

Describe the sample preparation, detailing the biological source of the cells and any tissue processing steps used.

#### Instrument

Identify the instrument used for data collection, specifying make and model number.

Software

Describe the software used to collect and analyze the flow cytometry data. For custom code that has been deposited into a community repository, provide accession details.

Cell population abundance

Describe the abundance of the relevant cell populations within post-sort fractions, providing details on the purity of the samples and how it was determined.

Gating strategy

Describe the gating strategy used for all relevant experiments, specifying the preliminary FSC/SSC gates of the starting cell population, indicating where boundaries between "positive" and "negative" staining cell populations are defined.

Tick this box to confirm that a figure exemplifying the gating strategy is provided in the Supplementary Information.

## Magnetic resonance imaging

### Experimental design

Design type

Indicate task or resting state; event-related or block design.

Design specifications

Specify the number of blocks, trials or experimental units per session and/or subject, and specify the length of each trial or block (if trials are blocked) and interval between trials.

Behavioral performance measures

State number and/or type of variables recorded (e.g. correct button press, response time) and what statistics were used to establish that the subjects were performing the task as expected (e.g. mean, range, and/or standard deviation across subjects).

### Acquisition

Imaging type(s)

Specify: functional, structural, diffusion, perfusion.

Field strength

Specify in Tesla

Sequence &amp; imaging parameters

Specify the pulse sequence type (gradient echo, spin echo, etc.), imaging type (EPI, spiral, etc.), field of view, matrix size, slice thickness, orientation and TE/TR/flip angle.

Area of acquisition

State whether a whole brain scan was used OR define the area of acquisition, describing how the region was determined.

Diffusion MRI

 Used

 Not used

### Preprocessing

Preprocessing software

Provide detail on software version and revision number and on specific parameters (model/functions, brain extraction, segmentation, smoothing kernel size, etc.).

Normalization

If data were normalized/standardized, describe the approach(es): specify linear or non-linear and define image types used for transformation OR indicate that data were not normalized and explain rationale for lack of normalization.

Normalization template

Describe the template used for normalization/transformation, specifying subject space or group standardized space (e.g. original Talairach, MNI305, ICBM152) OR indicate that the data were not normalized.

Noise and artifact removal

Describe your procedure(s) for artifact and structured noise removal, specifying motion parameters, tissue signals and physiological signals (heart rate, respiration).

Volume censoring

Define your software and/or method and criteria for volume censoring, and state the extent of such censoring.

### Statistical modeling & inference

Model type and settings

Specify type (mass univariate, multivariate, RSA, predictive, etc.) and describe essential details of the model at the first and second levels (e.g. fixed, random or mixed effects; drift or auto-correlation).

Effect(s) tested

Define precise effect in terms of the task or stimulus conditions instead of psychological concepts and indicate whether ANOVA or factorial designs were used.

Specify type of analysis:  Whole brain  ROI-based  Both

Statistic type for inference  
(See [Eklund et al. 2016](#))

Specify voxel-wise or cluster-wise and report all relevant parameters for cluster-wise methods.

Correction

Describe the type of correction and how it is obtained for multiple comparisons (e.g. FWE, FDR, permutation or Monte Carlo).

## Models & analysis

- n/a | Involved in the study
- Functional and/or effective connectivity
- Graph analysis
- Multivariate modeling or predictive analysis

Functional and/or effective connectivity

*Report the measures of dependence used and the model details (e.g. Pearson correlation, partial correlation, mutual information).*

Graph analysis

*Report the dependent variable and connectivity measure, specifying weighted graph or binarized graph, subject- or group-level, and the global and/or node summaries used (e.g. clustering coefficient, efficiency, etc.).*

Multivariate modeling and predictive analysis

*Specify independent variables, features extraction and dimension reduction, model, training and evaluation metrics.*