

1 **CausalCell: applying causal discovery to single-cell analyses**

2

3 Yujian Wen^{1#}, Jielong Huang^{1#}, Hai Zhang^{4#}, Shuhui Guo¹, Yehezqel Elyahu² Alon Monsonego²,
4 Yanqing Ding^{*3}, Hao Zhu^{1,5*}

5

6 1 Bioinformatics Section, School of Basic Medical Sciences, Southern Medical University,
7 Guangzhou, 510515, China

8 2 The Shraga Segal Department of Microbiology, Immunology and Genetics, Faculty of
9 Health Sciences, Ben-Gurion University of the Negev, Beer-Sheva 8410501, Israel

10 3 Department of Pathology, School of Basic Medical Sciences, Southern Medical University,
11 Guangzhou, 510515, China

12 4 Network Center, Southern Medical University, Guangzhou, 510515, China

13 5 Guangdong-Hong Kong-Macao Greater Bay Area Center for Brain Science and
14 Brain-Inspired Intelligence, Southern Medical University, Guangzhou, 510515, China

15

16 # These authors contributed equally to the work.

17 * Corresponding authors. Email: dyqgz@126.com (Y.D.), zhuhao@smu.edu.cn (H.Z.)

18

19 **ABSTRACT**

20 Correlation between objects does not answer many scientific questions because of the lack
21 of causal but the excess of spurious information and is prone to happen by coincidence.
22 Causal discovery infers causal relationships from data upon conditional independence test
23 between objects without prior assumptions (e.g., variables have linear relationships and
24 data follow the Gaussian distribution). Causal interactions within and between cells provide
25 valuable information for investigating gene regulation, identifying diagnostic and
26 therapeutic targets, and designing experimental and clinical studies. The rapid increase of
27 single-cell data permits inferring causal interactions in many cell types. However, because no
28 algorithms have been designed for handling abundant variables and few algorithms have
29 been evaluated using real data, how to apply causal discovery to single-cell data remains a
30 challenge. We report a pipeline and web server

31 (<http://www.gaemons.net/causalcell/causalDiscovery/>) for accurately and conveniently
32 performing causal discovery. The pipeline has been developed upon the benchmarking of 18
33 algorithms and the analyses of multiple datasets. Our applications indicate that only
34 complicated algorithms can generate satisfactorily reliable results. Critical issues are
35 discussed, and tips for best practices are provided.

36 **Keywords:** Single-cell, scRNA-seq, feature selection, causal discovery, causal network, causal
37 analysis

38

39 INTRODUCTION

40 The cell-specific regulation of gene expression and protein interaction generate various
41 emergent signalling pathways which indicate that most interactions between genes and their
42 products are causal. Causation determines widely observed and varied correlation. Some
43 causal interactions are annotated in the "canonical" pathways (e.g., the KEGG pathways), but
44 most remain unannotated, especially those in cells during development and in diseases and
45 in small cell populations. On statistical data analysis, Judea Pearl wrote "*statistics alone*
46 *cannot tell which is the cause and which is the effect*" ([Pearl and Mackenzie, 2019](#)); however,
47 uncovering causation is more difficult than uncovering correlation. Causal discovery is a
48 science which infers causal interactions from data observations upon testing conditional
49 independence (CI) between variables ([Glymour et al., 2019](#)). Mathematically, CI is at the
50 heart of causal discovery and $CI \neq \text{unconditional independence} \neq \text{uncorrelation}$.

51

52 Researchers have used RNA-seq to detect gene expression in a lump of cells for years, but
53 causal interactions in such mixed cells are blurred. Also, the sizes of such samples (lumped
54 cells) are adequate for inferring only correlation but not causation between genes. Many
55 methods (e.g., weighted gene co-expression network analysis, WGCNA) have been developed
56 to construct networks of correlated genes in lumped cells upon RNA-seq data ([Joehanes,](#)
57 [2018](#)). Recently, scRNA-seq has been widely used to detect gene expression in single cells. In
58 many situations (especially scRNA-seq using 10X Genomics), numbers of many cell types
59 allow for inferring causal interactions between genes in each cell type.

60

61 Many different CI tests have been developed, from the quite fast Gauss CI test to the highly
62 time-consuming kernel-based CI tests ([Verbyla, 2018](#); [Zhang et al., 2011](#)). Gauss CI test is
63 based upon partial correlations between variables; kernel-based CI tests estimate the
64 dependence between variables upon their observations without assuming any relationship
65 between variables or data distribution. CI tests critically characterize causal discovery
66 algorithms and differentiate causal discovery from other network inference methods,

67 including regulatory network inference (Nguyen et al., 2021; Pratapa et al., 2020), causal
68 network inference (Lu et al., 2021), network inference (Deshpande et al., 2019), and gene
69 network inference (Marbach et al., 2012).

70

71 The PC algorithm (named after its developers Peter Spirtes and Clark Glymour) is a
72 state-of-the-art causal discovery algorithm and can work with different CI tests (Glymour et
73 al., 2019). The time consumption of CI tests (especially kernel-based ones) makes it
74 infeasible to apply the PC algorithm to all genes in a scRNA-seq dataset. On the other hand,
75 what CI tests best suit scRNA-seq data and how to make proper trade-offs between time
76 consumption and network size or accuracy remain unclear. Thus, benchmarking the PC
77 algorithm and CI tests using single-cell data is essential before developing causal discovery
78 pipelines and applying causal discovery to single-cell analysis.

79

80 This *Tools and Resources* article presents a solution to single-cell causal discovery by
81 combining feature selection algorithms and causal discovery algorithms. Upon
82 benchmarking 9 feature selection algorithms and 9 CI tests using simulated and real
83 scRNA-seq data, we developed a pipeline and web server (called CausalCell) to perform
84 causal discovery. Some measures are developed and imbeded into the pipeline to ensure
85 reliability of causal discovery. The analysis of multiple datasets were performed, with the
86 results indicating that complicated (time-consuming) CI tests are crucial for generating
87 reliable results. The inferred causal interactions provide informative clues for experimental
88 and clinical studies.

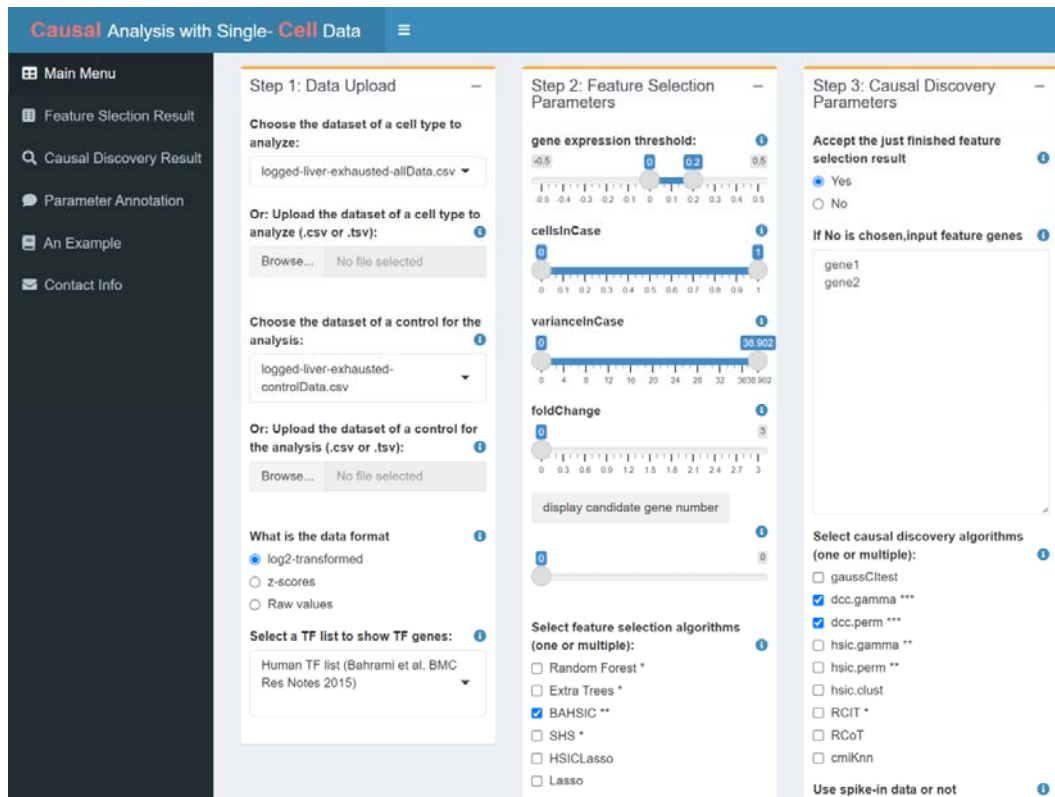
89

90 **METHOD DESCRIPTION**

91 **1. Software implementation**

92 The CausalCell pipeline consists mainly of feature selection and causal discovery. A parallel
93 version of the PC algorithm (Le et al., 2019), together with the Docker techniques, is used to
94 realize the parallel multi-task causal discovery, which is supported by a cluster of computers.
95 The user interface is implemented using the Shiny language (Figure 1). Annotations of
96 functions and parameters and a detailed description of an example are available online.

97



98

99 Figure 1. The user interface of CausalCell. Many functions are implemented to facilitate
100 performing feature selection and causal discovery.

101

102 2. Data input and display

103 scRNA-seq data generated by multiple protocols (e.g., 10X Genomics, smart-seq2) and
104 proteomics data (e.g., CyTOF) generated by mass cytometry can be analyzed ([Supplementary
105 Note 1](#)). Data can be in the log2-transformed or z-score normalized format, and online
106 transformation and normalization are available. A dataset can have or not have a control
107 dataset. If a control dataset is uploaded, the fold change of gene expression is computed
108 using the *FindMarkers* function in the *Seurat* package. Genes have multiple attributes (e.g.,
109 expression value, the percent of cells in which they are expressed, variance, and fold change);
110 all of these attributes can be used to order genes to reveal gene expression features and to
111 filter genes for performing feature selection (researchers often try to identify and analyse
112 highly differentially expressed genes or genes having high variance).

113

114 3. Feature selection

115 Combining feature selection and causal discovery enables causal discovery to be applied to a
116 arbitrary set of genes (feature genes). After genes are filtered upon conditions (i.e.,

117 expression threshold, the expressed cells, variance, and fold change) which generate
118 candidate genes for feature selection, one or multiple genes of primary interest are used as
119 the target genes (aka response variables) to select feature genes (aka features) from the
120 candidate genes (aka candidates). Upon the evaluation of the accuracy, time consumption,
121 and scalability of the 9 feature selection algorithms ([Supplementary Note 2](#)), BAHSIC is the
122 most recommended feature selection algorithm. We also recommend the joint use of
123 multiple algorithms (e.g., Random Forest + BAHSIC) to ensure reliability. Usually, feature
124 genes should be 50-70 (depending on what causal discovery algorithms are chosen). Genes
125 can be manually added into or removed from the feature gene list, to make feature genes
126 better reflect a biological question. Also, all feature genes can be manually selected without
127 performing feature selection, by which the user can examine any gene set.

128

129 **4. Causal discovery**

130 We implemented 9 causal discovery algorithms by combining the parallel version of the PC
131 algorithm with 9 CI tests ([Le et al., 2019](#)). We evaluated the accuracy, time consumption,
132 sample requirement, and stability of the 9 CI tests ([Figure 2](#); [Supplementary Note 3](#)). The
133 DCC algorithms are both most accurate and most time-consuming, suitable for small-scale
134 network inference; RCIT is reasonably accurate and relatively fast, suitable for large-scale
135 network inference. Multiple algorithms can be chosen in one run for a feature gene set, and a
136 consensus network can be constructed upon the networks inferred by some or all selected
137 algorithms. The consensus network is statistically more reliable. Edges in causal networks
138 have arrows that indicate activation or inhibition and show thickness that indicate CI test's
139 statistical significance.

140

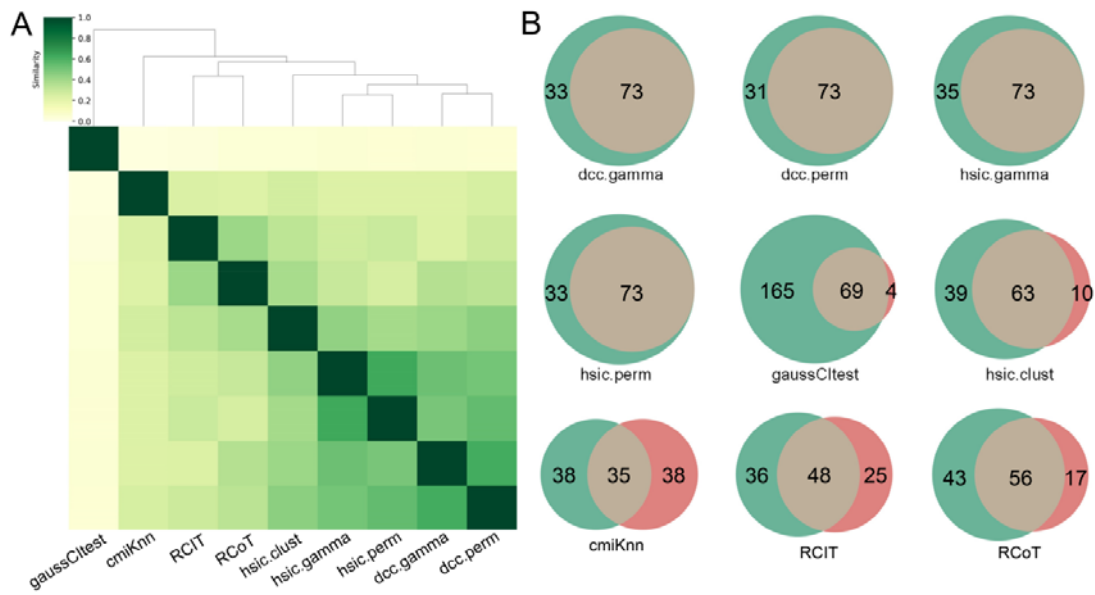
141 If the scRNA-seq dataset is too large, a subset of it should be sampled. Typically, for
142 Smart-seq2 data, 300 cells are enough, and for 10X Genomics data, 600 cells are enough.
143 Also, HSIC.perm and DCC.perm use permutations when performing the CI test. The random
144 sampling and permutation make causal networks inferred each time not identical. Our
145 benchmarking and data analyses reveal that interactions inferred by DCC algorithms are
146 highly stable ([Figure 3](#)).

147

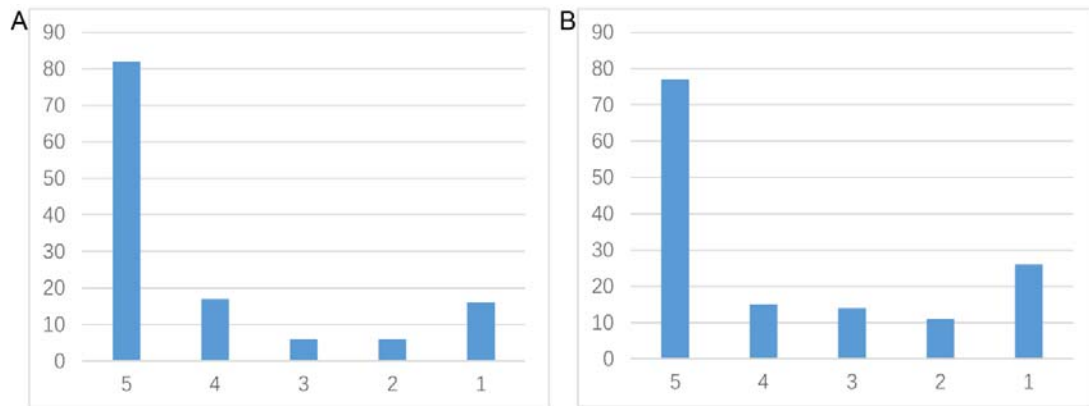
148 The following three parameters greatly influence causal discovery. "Set the alpha level"
149 determines the statistical significance cutoff of CI test; a large alpha level causes more causal
150 interactions to be inferred. "Select the number of cells" controls sample size; selecting more
151 cells for causal discovery makes the inference more reliable but more time-consuming.
152 "Select how a subset of cells is sampled" determines the way of a subset of cells is sampled. If
153 a subset is sampled randomly, the inferred causal network is not exactly reproducible, but

154 by running multiple times the inferred causal networks are highly consistent (Figure 3).
 155 Since each causal discovery task takes at least hours, providing an email address is necessary
 156 to make the result sent to the user automatically when it completes.

157
 158
 159



160
 161 Figure 2. The accuracy of the 9 CI algorithms (based on 9 CI tests). (A) The cluster map
 162 measures the consistency between causal networks generated by the 9 algorithms. Darker
 163 colors indicate higher similarity, and the networks of DCC.gamma, DCC.perm, HSIC.gamma,
 164 and HSIC.perm have the highest similarity values. A consensus network built upon the four
 165 DCC and HSIC networks was used as the reference to evaluate algorithms. (B) For each
 166 algorithm's network (green circled area), interactions overlapping the interactions in the
 167 consensus network (pink circled area) were examined. There are 73 overlapping
 168 interactions between DCC.gamma's network and the consensus network; thus, the true
 169 positive rate of the DCC.gamma network (TPR)=73/(73+33)=68.9%. The TPR of DCC.gamma,
 170 DCC.perm, HSIC.gamma, HSIC.perm, gaussCItest, HSIC.clust, cmiKnn, RCIT and RCoT are
 171 68.9%, 70.2%, 67.6%, 68.9%, 29.5%, 61.8%, 47.9%, 57.1%, and 56.6%.
 172



173

174 **Figure 3.** The shared and distinct interactions inferred by DCC.gamma (A) and DCC.perm (B)
175 by running the algorithm 5 times using the dataset of lung cancer cell line H2228. 78% and
176 64.3% of interactions occurred stably in ≥ 4 networks and many distinct interactions
177 occurred in just one network, indicating that the networks inferred by the two algorithms
178 are stable.

179

180 **5. Evaluating and ensuring the reliability**

181 A challenge for all kinds of network inferences is to verify or validate inferred networks.
182 Inspired by using RNA spike-in to measure RNA sequencing quality, we developed a method
183 to evaluate and ensure the reliability of causal discovery. This method includes three steps:
184 extracting the data of some well-known genes and their interactions from some datasets as
185 the "spike-in" data, integrating the spike-in data into the primary dataset, and applying
186 causal discovery to the integrated dataset. In the first step, the user can pick up a spike-in
187 data stored in the web server or design and upload a specific one; the following two steps
188 are performed automatically. In the inferred causal network, if genes and their interactions
189 in the spike-in data are clearly separated from genes and interactions in the primary dataset,
190 the causal discovery should be pretty reliable ([Supplementary Note 4](#)).

191

192 **6. Key features of different algorithms**

193 Upon one or several response variables (i.e., genes of interest), feature selection chooses a
194 subset of features (i.e., variables, genes) from the whole dataset by removing features
195 unrelated or less related to response variables. A feature selection algorithm combines a
196 search technique and an evaluation measure. After obtaining a measure between the
197 response variable(s) and each feature, a subset of features most related to the response
198 variable(s) is extracted. Constraint-based causal discovery algorithms identify causal
199 relationships in a set of features in two steps: skeleton estimation (determining the skeleton
200 of the causal network) and orientation (determining the direction of edges in the causal

201 network). Algorithms are different in that they use different CI tests to perform the first step
 202 (the most time-consuming step). We combine the PC algorithm with 9 CI tests to form 9
 203 causal discovery algorithms. Table 1 and Table 2 briefly describe the features and
 204 advantages/disadvantages of these feature selection and causal discovery algorithms. “+++”
 205 and “+” in the tables indicate the most and least recommended ones.

206

207

Table 1 Performance of the 9 feature selection algorithms

Algorithm	Category	Time consumption	Accuracy	Scalability	Stability	Advantage /disadvantage
RandomForest	An ensemble learning-based method (e.g., random forest) uses many trees of a random forest to calculate the importance of features, then performs regression based on the response variable(s) to identify the most relevant features.	+	++	++	+	This kind of algorithms is indeterministic (the same input may generate somewhat different outputs). Both ExtraTrees and RandomForest are good, the accuracy of XGBoost is unsatisfactory.
ExtraTrees		+	++	++	+	
XGBoost		++	+	+	++	
BAHSIC	Hilbert-Schmidt independence criterion (HSIC) is used as the measure of dependency between the response variable and features.	+	+++	+	++	BAHSIC and SHS are the best and second best, fast and accurate.
SHS		+	+++	+	++	
HSIC Lasso	BAHSIC, SHS, and HSIC Lasso are three HSIC-based algorithms.	++	++	++	++	Inferior to BAHSIC and SHS.
Lasso	Lasso is a regression analysis method that performs both variable selection and regularization. Regularization adds additional constraints or penalty to a regression model. Features which have non-zero regression coefficients are 'selected' by Lasso algorithms. Lasso, RidgeRegression, and ElasticNet are	+++	+	+++	++	Inferior to BAHSIC and SHS. Accuracy is not high and scalability is poor.
RidgeRegression		+++	+	+++	++	
ElasticNet		+++	+	+++	++	

	three regulation terms.					
--	-------------------------	--	--	--	--	--

208

209

Table 2 Performance of the 9 causal discovery algorithms

Algorithm	Category	Time consumption	Accuracy	Sample size	Stability #1	Advantage/dis advantage
GaussCItest	Gauss CI test examines CI using partial correlation, assuming that all variables are multivariate Gaussian. This assumption impairs the performance of GussCItest, especially when data are complex.	+++	+	+++	+++	Fast but inaccurate
CMIknn	Conditional mutual information (CMI) is a measure based on mutual information, which can be used to measure mutual dependence between two variables.	+++	++	+	+	Fast but inaccurate
RCIT	The Kernel Conditional Independence Test (KCIT) is a powerful but time-consuming CI test. RCIT and RCoT are two approximation methods of KCIT.	++	++	++	++	Moderately accurate, fast, recommended for large-scale networks
RCoT		++	++	++	++	
HSIC.clust	HSIC is a measure of dependency between two variables; HSIC(X, Y) = 0 if X and Y are unconditionally independent. HSIC.gamma and HSIC.perm employ gamma test and permutation test to estimate a <i>p</i> value.	+	++	++	++	Slow yet accurate, do not need large samples, recommended for small networks
HSIC.gamma		+	+++	++	++	
HSIC.perm		+	+++	++	+	
DCC.gamma	Distance covariance is an alternative to HSIC for measuring independence. DCC.gamma and DCC.perm employ Gamma test and permutation test to estimate a <i>p</i> value.	+	+++	+++	++	Slow yet most accurate, do not need large samples, recommended for small networks
DCC.perm		+	+++	+++	+	

210

#1 see Supplementary Table 3.

211

212

213 APPLICATIONS

214 1. The analysis of lung cancer cell lines and alveolar epithelial cells

215 Down-regulated MHC-II genes help cancer cells avoid being recognized by immune cells
216 (Rooney et al., 2015); thus, identifying genes and interactions related to the down-regulation
217 is important. To assess if causal discovery helps identify the related interactions, we
218 examined 5 lung cancer cell lines (A549, H1975, H2228, H838, and HCC827) and the normal
219 alveolar epithelial cells (Tian et al., 2019; Travaglini et al., 2020). For each of the six datasets,
220 we took the 5 MHC-II genes (*HLA-DPA1*, *HLA-DPB1*, *HLA-DRA*, *HLA-DRB1*, *HLA-DRB5*) as the
221 target genes and selected 50 feature genes (using BAHASIC, unless otherwise stated) from all
222 genes expressed in >50% cells. Then, we applied 9 causal discovery algorithms to the 50
223 genes in 300 cells sampled from each of the datasets. The two DCC algorithms performed the
224 best when processing the H2228 cells and lung alveolar epithelial cells (Figure 2;
225 Supplementary Note 5).

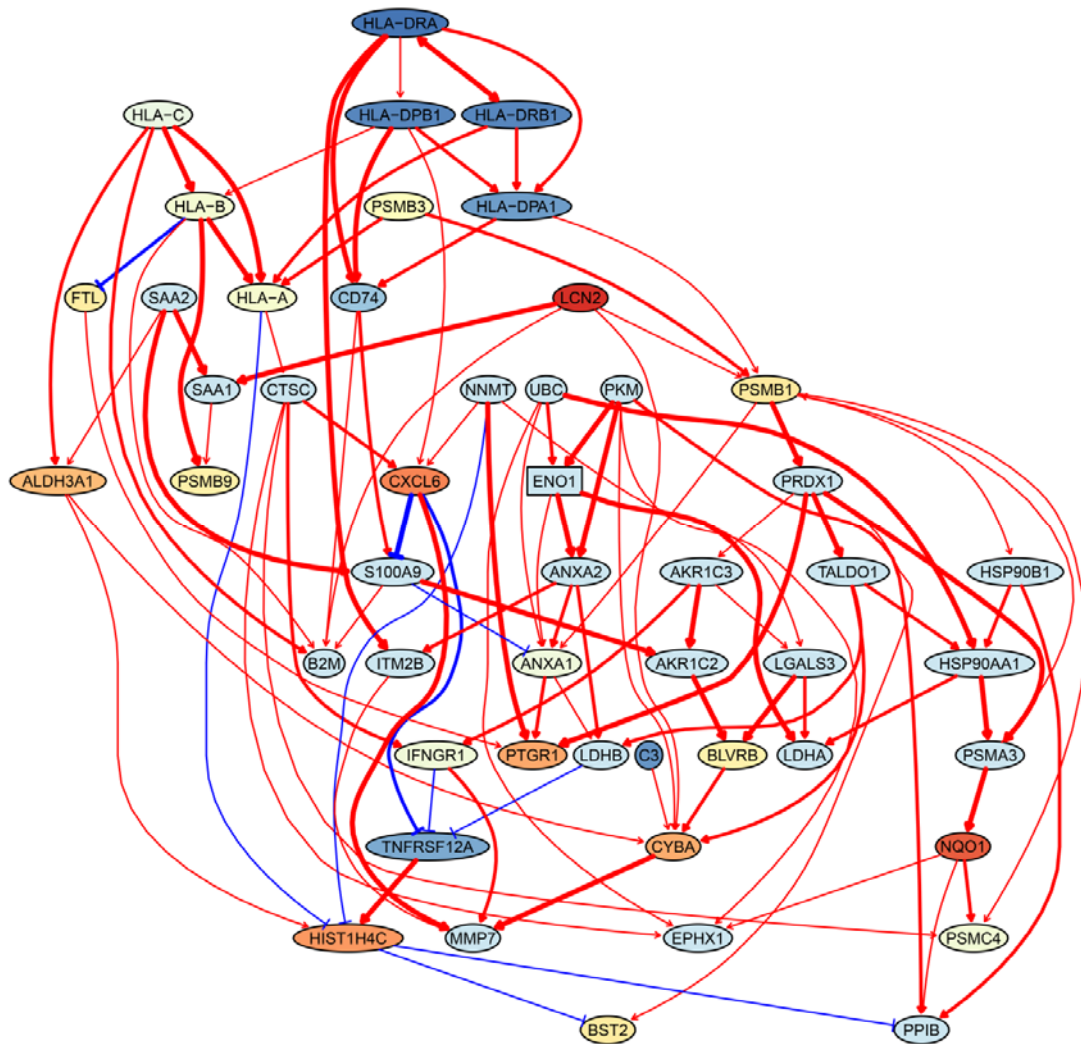
226

227 Inferred networks show that down-regulated genes weakly, but up-regulated genes strongly,
228 regulate downstream targets and that loss of activation (or inhibition) leads to down (or up)
229 regulation. These features are biologically reasonable. Many interactions, including those
230 among MHC-II genes and CD74, among CXCL genes, and among MHC-I genes and B2M, are
231 supported by the STRING database (<http://string-db.org>) and experimental findings (Figure
232 4; Supplementary Fig. 12) (Castro et al., 2019; Karakikes et al., 2012; Szklarczyk et al., 2021).
233 An interesting finding is the PRDX1→TALDO1→HSP90AA1→NQO1→PSMC4 cascade in
234 H2228 cells. Interactions between PRDX1/TALDO1/HSP90AA1 and NQO1 were reported
235 (Mathew et al., 2013; Yin et al., 2021), but between NQO1 and PSMC4 were not. Previous
236 findings on NQO1 include that it determines cellular sensitivity to the antitumor agent
237 Napabucasin in many cancer cell lines (Guo et al., 2020), is a potential poor prognostic
238 biomarker, and is a promising therapeutic target for patients with lung cancers (Cheng et al.,
239 2018; Siegel et al., 2012), and that mutations in *NQO1* are associated with susceptibility to
240 various forms of cancer. Previous findings on PSMC4 include that high levels of PSMC4 (and
241 other PSMC) transcripts were positively correlated with poor breast cancer survival (Kao et
242 al., 2021). Thus, the inferred NQO1→PSMC4 probably somewhat explains the mechanism
243 behind these experimental findings.

244

245

246



247

248 Figure 4. The network of the 50 genes inferred by DCC.gamma from the H2228 dataset (the
 249 alpha level for CI test was 0.1). Red → and blue -| arrows indicate activation and inhibition,
 250 and colors indicate fold changes of gene expression (from -2 to 2) compared with genes in
 251 the alveolar epithelial cells.

252

253 2. The analysis of macrophages isolated from glioblastoma

254 Macrophages critically influence glioma formation, maintenance, and progression (Gutmann,
 255 2020), and CD74 is the master regulator of macrophage functions in glioblastoma (Alban et al.,
 256 2020; Quail and Joyce, 2017; Zeiner et al., 2015). To examine the function of CD74 in
 257 macrophages in gliomas, we used CD74 as the target gene and selected 50 genes from genes
 258 expressed in >50% macrophages isolated from glioblastoma patients (Nefel et al., 2019). In
 259 the networks of DCC algorithms (Supplementary Note 6), CD74 regulates MHC-II genes,
 260 agreeing with the finding that CD74 is an MHC-II chaperone and plays a role in the

261 intracellular sorting of MHC class II molecules. In the network, there are interactions between
262 C1QA/B/C, agreeing that they form the complement C1q complex. The identified
263 TYROBP→TREM2→A2M→APOE→APOC1 cascade is supported by the reports that TREM2 is
264 expressed in tumor macrophages in over 200 human cancer cases ([Molgora et al., 2020](#)) and that
265 there are interactions between TREM2/A2M, TREM2/APOE, A2M/APOE, and APOE/APOC1
266 ([Krasemann et al., 2017](#)).

267

268 **3. The analysis of tumor-infiltrating exhausted CD8 T cells**

269 Tumor-infiltrating exhausted CD8 T cells are highly heterogeneous yet share common
270 differentially expressed genes ([McLane et al., 2019](#); [Zhang et al., 2018](#)), suggesting that CD8
271 T cells undergo different processes to reach exhaustion. We analyzed three exhausted CD8 T
272 datasets isolated from human liver, colorectal, and lung cancers ([Supplementary Note 7](#))
273 ([Guo et al., 2018](#); [Zhang et al., 2018](#); [Zheng et al., 2017](#)). A key feature of CD8 T cell
274 exhaustion identified in mice is PDCD1 upregulation by TOX ([Khan et al., 2019](#); [Scott et al.,](#)
275 [2019](#); [Seo et al., 2019](#)). Using TOX and PDCD1 as the target gene, we selected 50 genes
276 expressed in >50% exhausted CD8 T cells and 50 genes expressed in >50% non-exhausted
277 CD8 T cells, respectively. Transcriptional regulation of PDCD1 by TOX was observed in
278 LVMV-infected mice without mentioning any role of CXCL13 ([Khan et al., 2019](#)). Here
279 indirect TOX→PDCD1 (via genes such as CXCL13) was inferred in exhausted CD8 cells, and
280 direct TOX→PDCD1 was inferred in non-exhausted CD8 T cells (although the expression of
281 TOX and PDCD1 is low in these cells) ([Supplementary Figure 17](#)). Recently, CXCL13 was
282 found to play a critical role in T cells for effective responses to anti-PD-L1 therapies ([Zhang](#)
283 [et al., 2021](#)). The causal discovery results help reveal differences in CD8 T cell exhaustion
284 between species and under different pathological conditions. The PDCD1→TOX inferred in
285 exhausted and non-exhausted CD8 T cells may indicate some feedback between TOX and
286 PDCD1; on the proteome level, a related report is that the binding of PD1 to TOX in the
287 cytoplasm facilitates the endocytic recycling of PD1 ([Wang et al., 2019](#)).

288

289 **4. Identifying genes and inferring interactions that signify CD4 T cell age**

290 How immune cells age and whether some senescence signatures reflect the aging of all cells
291 draw wide attention ([Gorgoulis et al., 2019](#)). We analyzed gene expression in naive, TEM,
292 rTreg, naive_Isg15, cytotoxic, and exhausted CD4 T cells from young (2-3 months, n=4) and
293 old (22-24 months, n=4) mice ([Supplementary Note 8](#)) ([Elyahu et al., 2019](#)). For each cell
294 type, we compared the combined data from all four young mice with the data of each old
295 mouse to identify differentially expressed genes. If genes were expressed in >25% cells and
296 consistently up/down-regulated ($|\text{fold change}|>0$) in most of the 24 comparisons, we
297 assumed them as aging-related ([Supplementary Table 4](#)). Some of these identified genes

298 play important roles in the aging of T cells or other cells, such as the mitochondrial genes
299 encoding cytochrome C oxidases and the gene *Sub1* in the mTOR pathway (Bektas et al.,
300 2019; Gorgoulis et al., 2019; Goronzy and Weyand, 2019; Walters and Cox, 2021). We
301 directly used these genes, plus one CD4-specific biomarker (Cd28) and two reported aging
302 biomarkers (Cdkn1b, Cdkn2d) (Gorgoulis et al., 2019; Larbi and Fulop, 2014), as feature
303 genes to infer their interactions in different CD4 T cells in young and old mice. The causal
304 networks unveil multiple findings (Supplementary Figure 18). First, B2m→H2-Q7 (a mouse
305 MHC class I gene), Gm9843→Rps27rt (Gm9846), and the interactions between the five
306 mitochondrial genes (MT-ATP6, MT-CO1/2/3, MT-Nd1) were inferred in nearly all CD4 T
307 cells. Second, many interactions are supported by the STRING database (Supplementary
308 Figure 13). Third, some interactions agree with experimental findings, including
309 Sub1-|Lamtor2 (Chen et al., 2021) and the regulation of these mitochondrial genes by
310 Lamtor2 (Morita et al., 2017). Fourth, Gm9843→Rps27rt→Junb were inferred in multiple
311 CD4 T cells, and both Gm9843 and Rps27rt are mouse-specific. Since JUNB belongs to the
312 AP-1 family transcription factors that are increased in all immune cells during human aging
313 (Zheng et al., 2020), Gm9843→Rps27rt→Junb could highlight a counterpart regulation of
314 JUNB in human immune cells.

315

316 **DISCUSSION**

317 Various methods have been developed to infer interactions between variables from data. As
318 surveyed recently (Nguyen et al., 2021; Pratapa et al., 2020), most methods assume linear
319 relationships between variables and the Gaussian distribution of data. The assumptions
320 enable these methods to run fast, capable of handling many genes or performing
321 genome-wide predictions. Our results indicate that networks inferred by such fast methods
322 deserve serious concern. Instead, based on kernel-based CI tests, causal discovery performs
323 inference directly upon data observations without assuming any relationship between
324 variables and the distribution of data (Glymour et al., 2019; Imbens and Rubin, 2015). The
325 cost in time consumption pays off in terms of accuracy. Interacting genes and molecules
326 within and between cells may have varied quantitative relationships, so causal discovery
327 employing kernel-based CI tests best satisfies inferring causal interactions in varied single
328 cells.

329

330 Several conclusions can be drawn from the benchmarking and applications. First, although
331 kernel-based CI tests are time-consuming (Shah and Peters, 2020), applying causal discovery
332 to a set of genes can be reasonably performed. Of note, the most time-consuming CI tests
333 generate the most reliable results. Second, dropouts and noises in scRNA-seq data, which
334 concern researchers and trouble correlation computation (Hou et al., 2020; Mohan and Pearl,

335 [2018; Tu et al., 2019](#)), can be well tolerated by kernel-based CI tests if the dataset is large
336 enough to provide sufficient observations. Third, latent and unobserved variables influence
337 causal discovery (just as they influence any network inference), and a solution to this
338 problem is to evaluate whether the inference is reliable by using the "spike-in" data. Fourth,
339 it is difficult to judge inferred interactions if without relevant information (e.g., related
340 findings and domain knowledge).

341

342 Here are three examples showing the help of relevant information for judging inferred causal
343 interactions. First, upon the report TOX activating PDCD1 in mice ([Khan et al., 2019](#)),
344 whether CXCL13 is involved (or even required in humans) in the TOX-PDCD1 interaction in
345 exhausted CD8 T cells is unclear until CXCL13 was reported to play critical roles in T cells for
346 effective responses to anti-PD-L1 therapies ([Zhang et al., 2021](#)). Second, upon data from
347 different cancers, inferred networks in exhausted CD8 T cells are quite different, and a recent
348 study reports that exhausted CD8 T cells show high heterogeneity and exhaustion can follow
349 different paths ([Zheng et al., 2021](#)). Third, it was difficult to explain the multiple genes
350 encoding ribosomal proteins in the inferred networks in CD4 cells from old mice; a new
351 study reports that aging impairs the ability of ribosomes to synthesize proteins efficiently
352 ([Stein et al., 2022](#)).

353

354 **Limitations of the methods and study**

355 First, the time consumption of the most accurate causal discovery algorithms disables the
356 inference of large-scale networks. Inferring multiple networks with shared genes and
357 merging these networks into a big one is a way to circumvent this problem, but the
358 effectiveness of the strategy remains to be confirmed. Second, it deserves noting that
359 although time consumption pays off in accuracy, small networks could be biologically
360 inaccurate and unreliable due to potential lack of highly related genes. Third, to make the
361 trade-off properly between time consumption, network accuracy, and network size may
362 need multiple rounds of trials. Fourth, the current programming language support parallel
363 computing but does not support high-performance computing. The most time-consuming
364 parts of the codes are to be replaced using C codes.

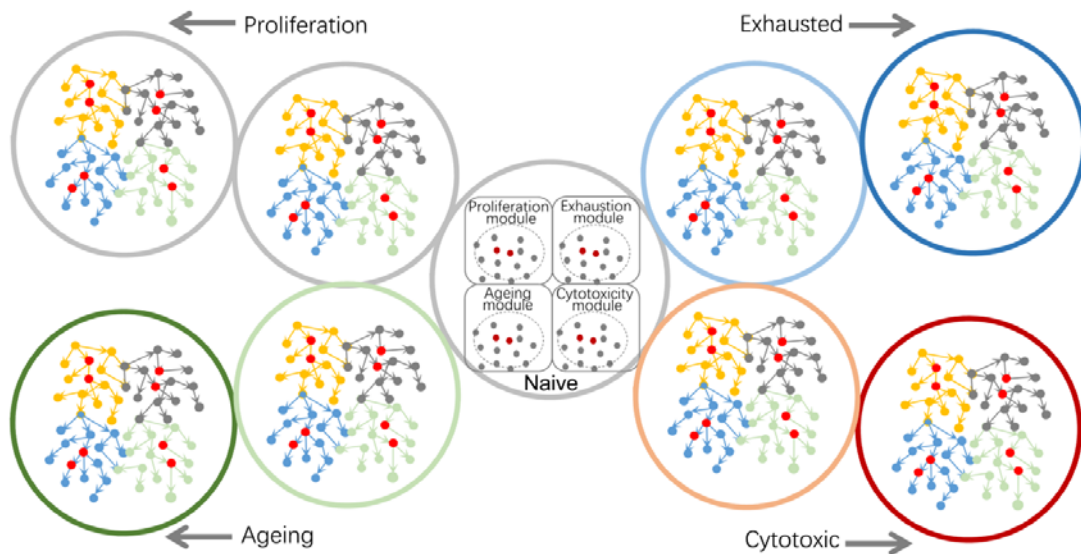
365

366 **TIPS FOR BEST PRACTICES**

367 First, exploring different modules or processes needs different target genes ([Figure 5](#)). When
368 it is unclear what gene is suitable or whether multiple genes can be co-selected, it is better to
369 examine one by one and inspect the shared feature genes. Second, BAHSIC and SHS are the
370 best feature selection algorithms. Third, selecting feature genes from too many candidate
371 genes may be unreliable. Usually, filtering out some genes is necessary upon conditions such

372 as genes are expressed in too few cells or have too low fold changes. Fourth, sometimes it is
373 advisable to apply causal discovery to a set of genes (e.g., differentially expressed genes)
374 without choosing a target gene and performing feature selection. Fifth, the two DCC
375 algorithms are most recommended; it is often sufficient just to use their results to build the
376 consensus network. Sixth, there are trade-offs between the scale, reliability, and accuracy for
377 causal discovery. To examine many genes, using RCIT is a proper trade-off. If the dataset is
378 large, choosing a subset of cells (e.g., 300) is a must. More cells are needed if feature genes
379 are expressed in a small portion (e.g., 25%) of cells or if scRNA-seq data are sparse. Seventh,
380 using a spike-in dataset and repeating causal discovery multiple rounds are two ways to
381 ensure and improve reliability. Eighth, carefully inspect the influence of cell heterogeneity on
382 causal discovery. Ninth, randomly sampling cells from the dataset and sampling cells with
383 more feature genes expressed suit large and small datasets, respectively. Tenth, causal
384 discovery identifies cell-specific causality when applied to homogeneous cells but identifies
385 more general causality when applied to heterogeneous cells (Figure 5); in the latter case,
386 caution is needed to interpret the results.

387
388



389
390
391
392
393
394
395
396

Figure 5. Using causal discovery to analyze different cells, cells at different stages, or different biological processes in cells. The red and grey dots within the four circles in the central cell indicate the four modules' core genes and related genes. When exploring different biological processes, core genes in different modules should be chosen as target genes.

397 **Declaration of Competing Interest**

398 The authors declare no competing interest.

399

400 **Additional information**

401 This manuscript has one supplementary file containing supplementary tables and figures.

402

403 **Author contributions**

404 H. Zhu designed the study and drafted the manuscript. Y.W. and J.H. performed algorithm
405 integration and benchmarking. Y.W., H. Zhang, S.G. developed the web server. H. Zhu, A.M.,
406 Y.E., and Y.D. analyzed data. A.M. revised the manuscript. All authors have read the
407 manuscript and consent to its publication.

408

409 **Acknowledgments**

410 This work was supported by the National Natural Science Foundation of China (31771456)
411 and the Department of Science and Technology of Guangdong Province
412 (2020A1515010803). We appreciate the help from Prof Ruichu Cai at the Guangdong
413 University of Technology.

414

415 **Data and code availability**

416 The web server is at <http://www.gaemons.net/causalcell/causalDiscovery/> (letters are
417 capital sensitive and “http” is without ‘s’).

418

419 **References**

- 420 1. Alban, T.J., Bayik, D., Otvos, B., Rabljenovic, A., Leng, L., Jia-Shiun, L., Roversi, G., Lauko, A.,
421 Momin, A.A., Mohammadi, A.M., *et al.* (2020). Glioblastoma Myeloid-Derived Suppressor
422 Cell Subsets Express Differential Macrophage Migration Inhibitory Factor Receptor
423 Profiles That Can Be Targeted to Reduce Immune Suppression. *Front Immunol* *11*, 1191.
- 424 2. Bektas, A., Schurman, S.H., Gonzalez-Freire, M., Dunn, C.A., Singh, A.K., Macian, F., Cuervo,
425 A.M., Sen, R., and Ferrucci, L. (2019). Age-associated changes in human CD4(+) T cells
426 point to mitochondrial dysfunction consequent to impaired autophagy. *Aging (Albany*
427 *NY)* *11*, 9234-9263.
- 428 3. Castro, A., Ozturk, K., Pyke, R.M., Xian, S., Zanetti, M., and Carter, H. (2019). Elevated
429 neoantigen levels in tumors with somatic mutations in the HLA-A, HLA-B, HLA-C and
430 B2M genes. *BMC Med Genomics* *12*, 107.
- 431 4. Chen, L., Liao, F., Wu, J., Wang, Z., Jiang, Z., Zhang, C., Luo, P., Ma, L., Gong, Q., Wang, Y., *et al.*
432 (2021). Acceleration of ageing via disturbing mTOR-regulated proteostasis by a new
433 ageing-associated gene PC4. *Aging Cell* *20*, e13370.

- 434 5. Cheng, X., Liu, F., Liu, H., Wang, G., and Hao, H. (2018). Enhanced glycometabolism as a
435 mechanism of NQO1 potentiated growth of NSCLC revealed by metabolomic profiling.
436 *Biochem Biophys Res Commun* 496, 31-36.
- 437 6. Deshpande, A., Chu, L.F., Stewart, R., and Gitter, A. (2019). Network inference with
438 Granger causality ensembles on single-cell transcriptomic data. *BioRxiv*, p.534834.
- 439 7. Elyahu, Y., Hekselman, I., Eizenberg-Magar, I., Berner, O., Strominger, I., Schiller, M., Mittal,
440 K., Nemirovsky, A., Eremenko, E., Vital, A., *et al.* (2019). Aging promotes reorganization
441 of the CD4 T cell landscape toward extreme regulatory and effector phenotypes. *Sci Adv*
442 5, eaaw8330.
- 443 8. Glymour, C., Zhang, K., and Spirtes, P. (2019). Review of Causal Discovery Methods Based
444 on Graphical Models. *Front Genet* 10, 524.
- 445 9. Gorgoulis, V., Adams, P.D., Alimonti, A., Bennett, D.C., Bischof, O., Bishop, C., Campisi, J.,
446 Collado, M., Evangelou, K., Ferbeyre, G., *et al.* (2019). Cellular Senescence: Defining a
447 Path Forward. *Cell* 179, 813-827.
- 448 10. Goronzy, J.J., and Weyand, C.M. (2019). Mechanisms underlying T cell ageing. *Nat Rev*
449 *Immunol* 19, 573-583.
- 450 11. Guo, G., Gao, Z., Tong, M., Zhan, D., Wang, G., Wang, Y., and Qin, J. (2020). NQO1 is a
451 determinant for cellular sensitivity to anti-tumor agent Napabucasin. *Am J Cancer Res*
452 10, 1442-1454.
- 453 12. Guo, X., Zhang, Y., Zheng, L., Zheng, C., Song, J., Zhang, Q., Kang, B., Liu, Z., Jin, L., Xing, R.,
454 *et al.* (2018). Global characterization of T cells in non-small-cell lung cancer by
455 single-cell sequencing. *Nat Med* 24, 978-985.
- 456 13. Gutmann, D.H. (2020). The Sociobiology of Brain Tumors. *Adv Exp Med Biol* 1225,
457 115-125.
- 458 14. Hou, W., Ji, Z., Ji, H., and Hicks, S.C. (2020). A systematic evaluation of single-cell
459 RNA-sequencing imputation methods. *Genome Biol* 21, 218.
- 460 15. Imbens, G.W., and Rubin, D.B. (2015). *Causal Inference for Statistics, Social and*
461 *Biomedical Sciences*. Cambridge University Press.
- 462 16. Joehanes, R. (2018). Network Analysis of Gene Expression. *Methods Mol Biol* 1783,
463 325-341.
- 464 17. Kao, T.J., Wu, C.C., Phan, N.N., Liu, Y.H., Ta, H.D.K., Anuraga, G., Wu, Y.F., Lee, K.H., Chuang,
465 J.Y., and Wang, C.Y. (2021). Prognoses and genomic analyses of proteasome 26S subunit,
466 ATPase (PSMC) family genes in clinical breast cancer. *Aging (Albany NY)* 13, 17970.
- 467 18. Karakikes, I., Morrison, I.E., O'Toole, P., Metodieva, G., Navarrete, C.V., Gomez, J.,
468 Miranda-Sayago, J.M., Cherry, R.J., Metodiev, M., and Fernandez, N. (2012). Interaction of
469 HLA-DR and CD74 at the cell surface of antigen-presenting cells by single particle image
470 analysis. *FASEB J* 26, 4886-4896.

- 471 19. Khan, O., Giles, J.R., McDonald, S., Manne, S., Ngiow, S.F., Patel, K.P., Werner, M.T., Huang,
472 A.C., Alexander, K.A., Wu, J.E., *et al.* (2019). TOX transcriptionally and epigenetically
473 programs CD8(+) T cell exhaustion. *Nature* 571, 211-218.
- 474 20. Krasemann, S., Madoze, C., Cialic, R., Baufeld, C., Calcagno, N., El Fatimy, R., Beckers, L.,
475 O'Loughlin, E., Xu, Y., Fanek, Z., *et al.* (2017). The TREM2-APOE Pathway Drives the
476 Transcriptional Phenotype of Dysfunctional Microglia in Neurodegenerative Diseases.
477 *Immunity* 47, 566-581 e569.
- 478 21. Larbi, A., and Fulop, T. (2014). From "truly naive" to "exhausted senescent" T cells: when
479 markers predict functionality. *Cytometry A* 85, 25-35.
- 480 22. Le, T.D., Hoang, T., Li, J., Liu, L., Liu, H., and Hu, S. (2019). A Fast PC Algorithm for High
481 Dimensional Causal Discovery with Multi-Core PCs. *IEEE/ACM Trans Comput Biol*
482 *Bioinform* 16, 1483-1495.
- 483 23. Lu, J., Dumitrescu, B., McDowell, I.C., Jo, B., Barrera, A., Hong, L.K., Leichter, S.M., Reddy,
484 T.E., and Engelhardt, B.E. (2021). Causal network inference from gene transcriptional
485 time-series response to glucocorticoids. *PLoS Comput Biol* 17, e1008223.
- 486 24. Marbach, D., Costello, J.C., Kuffner, R., Vega, N.M., Prill, R.J., Camacho, D.M., Allison, K.R.,
487 Consortium, D., Kellis, M., Collins, J.J., *et al.* (2012). Wisdom of crowds for robust gene
488 network inference. *Nat Methods* 9, 796-804.
- 489 25. Mathew, B., Jacobson, J.R., Siegler, J.H., Moitra, J., Blasco, M., Xie, L., Unzueta, C., Zhou, T.,
490 Evenoski, C., Al-Sakka, M., *et al.* (2013). Role of migratory inhibition factor in age-related
491 susceptibility to radiation lung injury via NF-E2-related factor-2 and antioxidant
492 regulation. *Am J Respir Cell Mol Biol* 49, 269-278.
- 493 26. McLane, L.M., Abdel-Hakeem, M.S., and Wherry, E.J. (2019). CD8 T Cell Exhaustion
494 During Chronic Viral Infection and Cancer. *Annu Rev Immunol* 37, 457-495.
- 495 27. Mohan, K., and Pearl, J. (2018). Graphical models for processing missing data.
496 arxiv.org/pdf/180103583.
- 497 28. Molgora, M., Esaulova, E., Vermi, W., Hou, J., Chen, Y., Luo, J., Brioschi, S., Bugatti, M.,
498 Omodei, A.S., Ricci, B., *et al.* (2020). TREM2 Modulation Remodels the Tumor Myeloid
499 Landscape Enhancing Anti-PD-1 Immunotherapy. *Cell* 182, 886-900 e817.
- 500 29. Morita, M., Prudent, J., Basu, K., Goyon, V., Katsumura, S., Hulea, L., Pearl, D., Siddiqui, N.,
501 Strack, S., McGuirk, S., *et al.* (2017). mTOR Controls Mitochondrial Dynamics and Cell
502 Survival via MTFP1. *Mol Cell* 67, 922-935 e925.
- 503 30. Neftel, C., Laffy, J., Filbin, M.G., Hara, T., Shore, M.E., Rahme, G.J., Richman, A.R., Silverbush,
504 D., Shaw, M.L., Hebert, C.M., *et al.* (2019). An Integrative Model of Cellular States,
505 Plasticity, and Genetics for Glioblastoma. *Cell* 178, 835-849 e821.
- 506 31. Nguyen, H., Tran, D., Tran, B., Pehlivan, B., and Nguyen, T. (2021). A comprehensive
507 survey of regulatory network inference methods using single cell RNA sequencing data.

- 508 Brief Bioinform 22.
- 509 32. Pearl, J., and Mackenzie, D. (2019). *The Book of Why - The New Science of Cause and*
510 *Effect*. Penguin.
- 511 33. Pratapa, A., Jalihal, A.P., Law, J.N., Bharadwaj, A., and Murali, T.M. (2020). Benchmarking
512 algorithms for gene regulatory network inference from single-cell transcriptomic data.
513 *Nat Methods* 17, 147-154.
- 514 34. Quail, D.F., and Joyce, J.A. (2017). The Microenvironmental Landscape of Brain Tumors.
515 *Cancer Cell* 31, 326-341.
- 516 35. Rooney, M.S., Shukla, S.A., Wu, C.J., Getz, G., and Hacohen, N. (2015). Molecular and
517 genetic properties of tumors associated with local immune cytolytic activity. *Cell* 160,
518 48-61.
- 519 36. Scott, A.C., Dundar, F., Zumbo, P., Chandran, S.S., Klebanoff, C.A., Shakiba, M., Trivedi, P.,
520 Menocal, L., Appleby, H., Camara, S., *et al.* (2019). TOX is a critical regulator of
521 tumour-specific T cell differentiation. *Nature* 571, 270-274.
- 522 37. Seo, H., Chen, J., Gonzalez-Avalos, E., Samaniego-Castruita, D., Das, A., Wang, Y.H.,
523 Lopez-Moyado, I.F., Georges, R.O., Zhang, W., Onodera, A., *et al.* (2019). TOX and TOX2
524 transcription factors cooperate with NR4A transcription factors to impose CD8(+) T cell
525 exhaustion. *Proc Natl Acad Sci U S A* 116, 12410-12415.
- 526 38. Shah, R.D., and Peters, J. (2020). The hardness of conditional independence testing and
527 the generalised covariance measure. *Ann Statist* 48, 1514-1538.
- 528 39. Siegel, D., Yan, C., and Ross, D. (2012). NAD(P)H:quinone oxidoreductase 1 (NQO1) in the
529 sensitivity and resistance to antitumor quinones. *Biochem Pharmacol* 83, 1033-1040.
- 530 40. Stein, K.C., Morales-Polanco, F., van der Lienden, J., Rainbolt, T.K., and Frydman, J. (2022).
531 Ageing exacerbates ribosome pausing to disrupt cotranslational proteostasis. *Nature*
532 601, 637-642.
- 533 41. Szklarczyk, D., Gable, A.L., Nastou, K.C., Lyon, D., Kirsch, R., Pyysalo, S., Doncheva, N.T.,
534 Legeay, M., Fang, T., Bork, P., *et al.* (2021). The STRING database in 2021: customizable
535 protein-protein networks, and functional characterization of user-uploaded
536 gene/measurement sets. *Nucleic Acids Res* 49, D605-D612.
- 537 42. Tian, L., Dong, X., Freytag, S., Le Cao, K.A., Su, S., JalalAbadi, A., Amann-Zalcenstein, D.,
538 Weber, T.S., Seidi, A., Jabbari, J.S., *et al.* (2019). Benchmarking single cell RNA-sequencing
539 analysis pipelines using mixture control experiments. *Nat Methods* 16, 479-487.
- 540 43. Travaglini, K.J., Nabhan, A.N., Penland, L., Sinha, R., Gillich, A., Sit, R.V., Chang, S., Conley,
541 S.D., Mori, Y., Seita, J., *et al.* (2020). A molecular cell atlas of the human lung from
542 single-cell RNA sequencing. *Nature* 587, 619-625.
- 543 44. Tu, R., Zhang, C., Ackermann, P., Mohan, K., Kjellstrom, H., and Zhang, K. (2019). Causal
544 discovery in the presence of missing data. *Proceedings of the 22nd International*

- 545 Conference on Artificial Intelligence and Statistics (AISTATS) 2019.
- 546 45. Verbyla, P. (2018). Network Inference Using Independence Criteria. (PhD dissertation).
547 Cambridge University.
- 548 46. Walters, H.E., and Cox, L.S. (2021). Intercellular Transfer of Mitochondria between
549 Senescent Cells through Cytoskeleton-Supported Intercellular Bridges Requires mTOR
550 and CDC42 Signalling. *Oxid Med Cell Longev* 2021, 6697861.
- 551 47. Wang, X., He, Q., Shen, H., Xia, A., Tian, W., Yu, W., and Sun, B. (2019). TOX promotes the
552 exhaustion of antitumor CD8(+) T cells by preventing PD1 degradation in hepatocellular
553 carcinoma. *J Hepatol* 71, 731-741.
- 554 48. Yin, H., Huang, Y.H., Best, S.A., Sutherland, K.D., Craik, D.J., and Wang, C.K. (2021). An
555 Integrated Molecular Grafting Approach for the Design of Keap1-Targeted Peptide
556 Inhibitors. *ACS Chem Biol* 16, 1276-1287.
- 557 49. Zeiner, P.S., Preusse, C., Blank, A.E., Zachskorn, C., Baumgarten, P., Caspary, L., Braczynski,
558 A.K., Weissenberger, J., Bratzke, H., Reiss, S., *et al.* (2015). MIF Receptor CD74 is
559 Restricted to Microglia/Macrophages, Associated with a M1-Polarized Immune Milieu
560 and Prolonged Patient Survival in Gliomas. *Brain Pathol* 25, 491-504.
- 561 50. Zhang, K., Peters, J., Janzing, D., and Schölkopf, B. (2011). Kernel-based conditional
562 independence test and application in causal discovery. In Proceedings of the
563 Twenty-Seventh Conference on Uncertainty in Artificial Intelligence (Barcelona, Spain:
564 AUAI Press), pp. 804–813.
- 565 51. Zhang, L., Yu, X., Zheng, L., Zhang, Y., Li, Y., Fang, Q., Gao, R., Kang, B., Zhang, Q., Huang, J.Y.,
566 *et al.* (2018). Lineage tracking reveals dynamic relationships of T cells in colorectal
567 cancer. *Nature* 564, 268-272.
- 568 52. Zhang, Y., Chen, H., Mo, H., Hu, X., Gao, R., Zhao, Y., Liu, B., Niu, L., Sun, X., Yu, X., *et al.*
569 (2021). Single-cell analyses reveal key immune cell subsets associated with response to
570 PD-L1 blockade in triple-negative breast cancer. *Cancer Cell* 39, 1578-1593 e1578.
- 571 53. Zheng, C., Zheng, L., Yoo, J.K., Guo, H., Zhang, Y., Guo, X., Kang, B., Hu, R., Huang, J.Y., Zhang,
572 Q., *et al.* (2017). Landscape of Infiltrating T Cells in Liver Cancer Revealed by Single-Cell
573 Sequencing. *Cell* 169, 1342-1356 e1316.
- 574 54. Zheng, L., Qin, S., Si, W., Wang, A., Xing, B., Gao, R., Ren, X., Wang, L., Wu, X., Zhang, J., *et al.*
575 (2021). Pan-cancer single-cell landscape of tumor-infiltrating T cells. *Science* 374,
576 abe6474.
- 577 55. Zheng, Y., Liu, X., Le, W., Xie, L., Li, H., Wen, W., Wang, S., Ma, S., Huang, Z., Ye, J., *et al.*
578 (2020). A human circulating immune cell landscape in aging and COVID-19. *Protein Cell*
579 11, 740-770.

580