

1 **Improved detection of tumor suppressor events in**
2 **single-cell RNA-Seq data**

3
4 Andrew E. Teschendorff^{1,2,*} and Ning Wang¹

5
6
7 1. CAS Key Laboratory of Computational Biology, CAS-MPG Partner Institute for Computational
8 Biology, Shanghai Institute of Nutrition and Health, Shanghai Institutes for Biological Sciences,
9 University of Chinese Academy of Sciences, Chinese Academy of Sciences, 320 Yue Yang Road,
10 Shanghai 200031, China.

11 2. UCL Cancer Institute, Paul O’Gorman Building, University College London, 72 Huntley Street,
12 London WC1E 6BT, United Kingdom.

13
14 *Corresponding author: Andrew E. Teschendorff- andrew@picb.ac.cn , a.teschendorff@ucl.ac.uk

15
16 **Keywords: single cell RNA-Seq; transcription factor; regulatory network**

Abstract

Tissue-specific transcription factors are frequently inactivated in cancer. To fully dissect the heterogeneity of such tumor suppressor events requires single-cell resolution, yet this is challenging because of the high dropout rate. Here we propose a simple yet effective computational strategy called SCIRA to infer regulatory activity of tissue-specific transcription factors at single-cell resolution and use this tool to identify tumor suppressor events in single-cell RNA-Seq cancer studies. We demonstrate that tissue-specific transcription factors are preferentially inactivated in the corresponding cancer cells, suggesting that these are driver events. For many known or suspected tumor suppressors, SCIRA predicts inactivation in single cancer cells where differential expression does not, indicating that SCIRA improves the sensitivity to detect changes in regulatory activity. We identify NKX2-1 and TBX4 inactivation as early tumor suppressor events in normal non-ciliated lung epithelial cells from smokers. In summary, SCIRA can help chart the heterogeneity of tumor suppressor events at single-cell resolution.

Keywords: single-cell RNA-Seq; cancer; transcription factor; tumor suppressor; regulatory activity; cell-type heterogeneity

Introduction

Tissue-specific transcription factors are required for the differentiated state of cells in a given tissue¹. They are often inactivated in cancer, which is associated with a lack of differentiation, a well-known cancer hallmark²⁻⁶. Many of these tissue-specific transcription factors (TFs) encode tumor suppressors and their inactivation may constitute driver events that are thought to occur in the earliest stages of carcinogenesis⁷⁻⁹. Estimating regulatory activity of such tissue-specific transcription factors (TFs) in both normal and cancer tissue is therefore a critically important task, as this can reveal which normal tissues are at risk of neoplastic transformation¹⁰. There are two main reasons why this task should be performed at single-cell resolution¹¹⁻¹³. First, TFs control cell-identity^{1,14}, and thus, estimation of regulatory activity in bulk tissue is subject to confounding by cell-type heterogeneity. Second, to fully characterize cancer heterogeneity requires identifying putative tumor suppressor events at the most fundamental scale, i.e. the single-cell¹⁵⁻¹⁸.

However, estimating regulatory activity of TFs at single-cell resolution is hard, because of the typically high dropout rate and low genomic coverage of single-cell assays¹⁹⁻²¹. In the context of single-cell RNA-Seq assays, one could in principle use TF expression as a surrogate marker of TF-activity (i.e. regulatory activity reflecting the effect of the TF on downstream expression of direct and indirect targets), and while this strategy works well on

54 expression data derived from bulk tissue (see e.g. ¹), it is unclear how well this works for
55 scRNA-Seq assays ^{22,23}. Thus, it is also unclear how best to infer regulatory activity in the
56 majority of scRNA-Seq cancer studies that are performed in solid epithelial tissues.
57 Here we present a novel strategy called SCIRA (SCalable Inference of Regulatory Activity in
58 single cells), which applies an existing regulatory inference method ⁸ to a suitably powered
59 bulk multi-tissue RNA-Seq dataset to identify tissue-specific TFs and their regulons (i.e. their
60 direct and indirect targets), from which regulatory activity in single cells can then be
61 estimated. We comprehensively validate SCIRA and demonstrate through a power calculation
62 and application to real scRNA-Seq data, that SCIRA can estimate regulatory activity even for
63 TFs that are highly expressed only in relatively minor fractions (~5%) of cells within a bulk
64 tissue. We subsequently apply SCIRA to several scRNA-Seq datasets containing both normal
65 and cancer cells, where it reveals preferential inactivation of tissue-specific TFs in
66 corresponding single cancer cells, an observation strongly consistent with analogous results
67 obtained in bulk tissue ⁵, whilst also revealing novel tumor suppressor events at single-cell
68 resolution. We further showcase an important application of SCIRA to identify tumor
69 suppressor events in single normal cells (lung epithelial cells) exposed to a cancer risk factor
70 (smoking). Our results underscore the critical need for a method like SCIRA, since ordinary
71 differential expression fails to reveal the same insights, even after imputation of dropouts.

72

73 **Results**

74 **Inferring regulatory activity with SCIRA: rationale**

75 SCIRA identifies tissue-specific TFs, builds regulons for these TFs, and uses these regulons
76 to estimate regulatory activity of the TFs in scRNA-Seq data (**Methods**). SCIRA adapts the
77 SEPIRA algorithm (previously published by us ⁸) to infer tissue-specific TFs and regulons
78 from the large GTEx multi-tissue bulk RNA-Seq dataset (8555 samples, 30 tissue-types) ²⁴
79 (**Methods, Fig.1A**). We note that the tissue-specific TFs are derived by adjusting for cell-type
80 (stromal) heterogeneity, which can otherwise strongly confound differential expression
81 analyses (**Methods**) ²⁵. To justify inferring TFs and their regulons from bulk tissue data, we
82 performed a careful power calculation, which revealed that SCIRA has reasonable sensitivity
83 to detect tissue-specific TFs that are highly expressed even if only in a relatively
84 underrepresented cell-type within the tissue (**Methods, Fig.1B**). For instance, using
85 reasonable values for the average fold-change (**SI fig.S1**), we estimated that for tissues like
86 lung, pancreas and liver, for which there are more than 100 samples in GTEx (total number
87 of samples is 8555), sensitivity to detect TFs expressed in only 5% of cells within the tissue
88 (i.e. a minor cell fraction MCF=0.05) were generally still over 50% (**Fig.1B, SI fig.S2**). The
89 inferred TF-regulons can subsequently be applied to suitably-matched scRNA-Seq data in a
90 linear regression framework ²⁶ (**Methods**) to estimate regulatory activity for each single cell.
91 By using the actual regulon of the TF, this inference should be robust to dropouts, i.e. even if

92 the TF itself is not detected across most if not all of the cells in the study (**Fig.1C**). Finally,
93 one can construct regulatory activity maps across the relevant cells within the tissue (**Fig.1C**),
94 which can reveal deregulated TFs at single-cell resolution.

95

96 **Validation of SCIRA in normal tissue**

97 As a proof of principle we applied SCIRA to four tissue-types (lung, liver, kidney and
98 pancreas) using the GTEX dataset to infer corresponding tissue-specific TFs and regulons.
99 We identified on average about 30 tissue-specific TFs for each of the 4 tissue-types and on
100 average about 40 to 50 regulon genes per TF (**SI tables S1-S4, Supplementary File 1**). The
101 TF lists contained well-known tissue-specific factors: e.g. for liver, the list included the
102 well-known hepatocyte factors *HNF1A*, *HNF4A* and *FOXA1* (*HNF3A*); for lung, the list
103 included well-known lung alveolar differentiation factors *TBX2* and *FOXA2*²⁷⁻²⁹, and *FOXJ1*,
104 a factor required for ciliogenesis³⁰. In order to test the reliability of the TFs and regulons, we
105 performed four separate validation analyses.

106 First, although there is no logical requirement for regulon genes to be direct targets³¹, some
107 enrichment for direct binding targets is expected. Approximately 65% of our TF-regulons
108 exhibited statistically significant enrichment for corresponding ChIP-Seq TF-binding targets
109 (**SI fig.S3-S4**), as determined using data from the ChIP-Seq Atlas³² (**Methods**). For instance,
110 in the case of liver we could find ChIP-Seq data for 12 of the 22 liver-specific TFs, and for
111 9/12 we observed statistically significant enrichment (**SI fig.S3D-E**). In many instances,
112 proportions of regulon genes that were direct TF binding targets were considerable. For
113 example, for the liver-specific TF *HNF4G*, 57% of its 37 regulon genes (i.e. 21 genes) were
114 bound by *HNF4G* within +/- 5kb of the gene's transcription start site (TSS) (**SI fig.S3D**). For
115 *FOXA1*, 8 of its 10 regulon genes were bound by *FOXA1* within +/-1kb of the TSS (**SI**
116 **fig.S3D**). Statistical significance estimates were independent of the choice of threshold on
117 binding intensity values (**Methods**), and also robust to parameter choices in SCIRA (**SI**
118 **fig.S5, Methods**). Second, we were able to validate the tissue-specificity of the regulons and
119 derived regulatory activity estimates in independent multi-tissue bulk RNA-Seq (ProteinAtlas
120³³) and microarray data from Roth et al³⁴ (**SI fig.S6-S9**). Given these successful validations,
121 we estimated on average only 10% of TF regulon-gene associations to be false positives (**SI**
122 **fig.S10**). Third, we collated and analysed scRNA-Seq datasets representing differentiation
123 timecourses into mature epithelial cell-types present within the given tissues, encompassing
124 two species (human & mouse) and 3 different single-cell technologies (Fluidigm C1,
125 DropSeq & Smart-Seq2) (**SI table S5, Methods**)³⁵⁻³⁸. We reasoned that most of our
126 tissue-specific TFs would exhibit higher regulatory activity in the corresponding mature
127 differentiated cells compared to the immature progenitors, a hypothesis that we were able to
128 strongly validate in each of the four tissue-types (**SI fig.S11-S14**). These results could not
129 have arisen by random chance and were not seen if we used tissue-specific TFs from other
130 unrelated (non-epithelial) tissues like skin or brain (**SI fig.S15**). We further observed that,

131 owing to the high dropout rate, SCIRA's regulatory activity estimates were much more
132 sensitive than expression itself (**SI fig.S11-S14, Fig.2A**). As a concrete example, SCIRA's
133 regulatory activity estimates for lung alveolar differentiation factors *TBX2* and *FOXA2*²⁷⁻²⁹
134 were higher in the mature alveolar cell-types compared to the immature progenitors, as
135 required, whilst expression levels could not detect an increase (**SI fig.S11**). SCIRA displayed
136 improved sensitivity and prevision (i.e. lower false discovery rate) over differential
137 expression (DE) even after application of imputation methods (scImpute³⁹, MAGIC⁴⁰,
138 Scrabble⁴¹), or even when compared to other regulatory activity estimation methods like
139 SCENIC/GENIE3⁴² (**Fig.2A-C, Methods**). SCIRA also displayed improved sensitivity over
140 the combined use of VIPER^{43,44} and the dorothea TF-regulon database^{45,46} ("VIPER-D"), as
141 well as lower FDRs (**Fig.2A-C, Methods**). This is noteworthy given that the TF-regulons
142 from dorothea are not tissue-specific. Fourth, we validated the power calculation underlying
143 SCIRA by applying it to a differentiation timecourse scRNA-Seq dataset in liver³⁶, which
144 revealed the expected bifurcation of hepatoblasts into hepatocytes and cholangiocytes, as well
145 as identifying cholangiocyte specific factors, despite their very low frequency (5-10%) in
146 liver tissue (**SI fig.S16B-E, SI fig.S17**). We note that the bifurcation and dynamic expression
147 patterns were not revealed when analyzing TF expression levels (**SI fig.S18**), further
148 supporting the view that SCIRA can improve the sensitivity to detect correct patterns of
149 TF-activity.

150

151

152 **SCIRA predicts inactivation of tissue-specific TFs in corresponding tumor epithelial** 153 **cells**

154 Next, we applied SCIRA to a recent lung cancer scRNA-Seq study (Lambrecht et al)⁴⁷ which
155 profiled a total of 52,698 cells (10X Chromium) derived from 5 lung cancer patients (2 lung
156 adenoma carcinomas – LUAD, 2 lung squamous cell carcinomas – LUSC and 1 non-small
157 cell lung cancer -NSCLC). We hypothesized that many of our previously identified
158 lung-specific TFs would be inactivated in lung epithelial tumor cells^{5,8}, since lack of
159 differentiation is a well-known cancer hallmark⁶. We used the same dimensional reduction
160 and tSNE-approach as in Lambrecht et al⁴⁷, to first categorize specific clusters of cells as
161 normal alveolar epithelial (n=1709) and tumor epithelial (n=7450) (**Fig.3A**). We verified that
162 the alveolar cells expressed relatively high levels of an alveolar marker (*CLDN18*) (**Fig.3B**),
163 whilst both alveolar and tumor epithelial cells expressed relatively high levels of *EPCAM*, a
164 well-known epithelial marker (**Fig.3C**). As noted by Lambrecht et al, the great majority of
165 alveolar cells were from non-malignant specimens representing normal (squamous)
166 epithelium and clustered together irrespective of patient-ID⁴⁷, whilst cancer cells clustered
167 according to patient (**Fig.3A**)⁴⁷. Next, we used SCIRA to estimate regulatory activity for all
168 38 lung-specific TFs in each of the (1709+7450) cells, and computed t-statistics of
169 differential activity between alveolar and tumor epithelial cells. Remarkably, 35 out of the 38

170 TFs exhibited a statistically significant (Bonferroni adjusted $P < 0.05$) reduction in regulatory
171 activity in tumor cells (**Fig.3D**, Wilcox-test $P < 1e-8$). Using 1000 Monte-Carlo
172 randomizations of the regulons, we verified that this number of inactivated TFs could not
173 have arisen by chance (**Fig.3D**, Monte Carlo $P < 0.001$). Among the most significantly
174 inactivated TFs, we observed *FOXA2*, a TF required for alveolarization and which regulates
175 airway epithelial cell differentiation^{28,29} (**Fig.3E**), and *NKX2-1*, a master TF of early lung
176 development⁴⁸ (**SI fig.S19**). Other inactivated TFs included *SOX13*, which has been broadly
177 implicated in lung morphogenesis⁴⁹, *HIF3A*, which has been shown to be highly expressed in
178 alveolar epithelial cells and thought to be protective of hypoxia-induced damage⁵⁰, and the
179 aryl hydrocarbon receptor (*AHR*), which is a regulator of mucosal barrier function and
180 activation of which enhances CD4+ T-cell responses to viral infections^{51,52} (**SI fig.S19**).
181 Importantly, these findings would not have been obtained had we performed DE or VIPER-D
182 analysis on the 38 TFs (**Fig.3D & 3F**). Indeed, according to a Wilcoxon rank sum test, 21 TFs
183 were differentially expressed between alveolar and tumor epithelial cells, but with no clear
184 trend towards underexpression in tumor cells (**Fig.3D**). For instance, according to single-cell
185 DE analysis, TFs such as *TBX4* and *FOXJ1*, both with important roles in lung tissue
186 development, were not underexpressed in tumor cells, yet found to be inactivated according
187 to SCIRA (**Fig.3F**). Given that *TBX4* and *FOXJ1* have been found to be
188 inactivated/underexpressed in bulk lung cancer tissue⁸, this further supports the view that
189 SCIRA improves sensitivity over ordinary DE analysis. To explore this further we compared
190 the differential activity and differential expression patterns between normal and cancer cells
191 to the differential expression patterns in the two TCGA lung cancer studies^{53,54}. A stronger
192 agreement with the bulk RNA-Seq data of both TCGA cohorts was observed for SCIRA's
193 differential activity profiles compared to differential expression or when using VIPER-D to
194 infer differential activity (**Fig.3F-G**). Indeed, approximately 30 of the 38 TFs exhibited
195 differential activity patterns at the single-cell level that were consistent with differential
196 expression in bulk, whilst for differential expression and VIPER-D this number was only
197 around 10 (**Fig.3H**).

198 To test the generality of our observations, we next considered a scRNA-Seq study profiling
199 normal colon epithelial cells and tumor colon epithelial cells⁵⁵. We first used SCIRA to
200 derive a colon-specific regulatory network from GTEX, resulting in 56 colon-specific TFs
201 and associated regulons (**SI table S6, Supplementary File 1**). This list included many well
202 known intestinal factors such as the enterocyte differentiation factors *CDX1/CDX2*⁵⁶, the
203 crypt epithelial factor *KLF5*⁵⁷ and the intestinal master regulator *ATOH1*^{58,59}. Next, we
204 obtained TF-activity (TFA) estimates for all 56 colon-TFs across a total of 432 single cells
205 (160 normal epithelial + 272 cancer epithelial, C1 Fluidigm) from 11 different colon-cancer
206 patients. Hierarchical clustering over this TFA-matrix revealed clear segregation of single
207 cells by normal/cancer status and not by patient (**Fig.4A**). Of the 56 TFs, 23 exhibited
208 differential activity (Bonferroni $P < 0.05$) with the great majority (87%, 20/23) exhibiting

209 inactivation, indicating a strong statistical tendency for inactivation in cancer cells (Binomial
210 test, $P=3e-5$, **Fig.4B**). Once again, had we relied on TF-expression itself, no segregation of
211 single-cells by normal/cancer status was evident (**Fig.4A**), and only 13 TFs were
212 differentially expressed (Bonferroni $P<0.05$) with no obvious trend towards underexpression
213 in cancer (Binomial test, $P=0.13$, **Fig.4B**). Of note, while *CDX1* and *CDX2* were found to be
214 both inactivated and underexpressed, several TFs like *KLF5* or *ATOH1* with established
215 tumor suppressor roles in colorectal cancer^{60,61}, were only found inactivated via SCIRA
216 (**Fig.4C**). Interestingly, using VIPER-D there was only moderate correlation with SCIRA's
217 predictions, with VIPER-D not predicting preferential inactivation and failing to predict
218 inactivity of established tumor suppressors like *KLF5* and *CDX1* (**Fig.4B**). Performing the
219 analysis on a per-patient level and focusing on the 3 patients with the largest numbers of both
220 normal and tumor epithelial cells, revealed a similar skew towards inactivation with 8, 15 and
221 21 TFs exhibiting significantly lower activity across cancer cells (**Fig.4D**), and with
222 effectively no TF exhibiting increased activity. For several TFs and for each of the 3 patients,
223 inactivation events were seen across most if not all cancer cells (**Fig.4D**): for instance, this
224 was the case for *ATOH1*, or the autophagy inducer *TRIM31*⁶², thus implicating disruption of
225 this novel and specific autophagy pathway in colon cancer⁶³. Using the 5 patients with both
226 normal and cancer cells profiled, we estimated the frequency of inactivation of all 56
227 colon-specific TFs across the 5 patients, which revealed that *CDX2* and *TRIM31* were
228 inactivated in 80% of the patients, whilst *ATOH1*, *HNF4A*, *CDX1* and *TBX10* were
229 inactivated in 60% (**SI fig.S20**).

230

231 **Tissue-specificity of TF inactivation in cancer**

232 The observed frequent inactivation of tissue-specific TFs in corresponding single cancer cells
233 suggests that these could be driver events. To obtain supporting evidence for this, we
234 reasoned that TFs specific for other unrelated tissue-types would exhibit much lower
235 frequencies of inactivation. We thus compared the lung and colon-specific TFs to additional
236 TFs specific to skin and brain, two non-epithelial tissue types, as well as to stomach-specific
237 TFs which should bear more resemblance to colon-TFs. Consistent with our expectation, in
238 the case of lung cancer cells, the TFs specific to colon, stomach, brain and skin exhibited
239 much lower frequencies of inactivation compared to lung-TFs (**SI fig.21A**). In the case of
240 colon cancer cells, colon and stomach-specific TFs exhibited the highest inactivation
241 frequencies, and were about two-fold higher than for skin and brain-specific TFs (**SI**
242 **fig.S21B**).

243

244 **Inactivation of tumor suppressors in normal cells at risk of cancer**

245 An important application of SCIRA is to normal cells at risk of cancer, which could reveal
246 early inactivation of key tumor suppressor TFs. To demonstrate this, we applied SCIRA to a
247 scRNA-Seq dataset (CEL-Seq) encompassing 564 lung epithelial cells, obtained from

248 bronchial brushings of 6 healthy individuals (6 never-smokers, 6 current smokers) ⁶⁴
249 (**Methods**). We inferred regulatory activity profiles for our 38 lung-specific TFs in each of
250 the 564 lung epithelial cells, and subsequently used t-stochastic neighborhood embedding
251 (tSNE) ⁶⁵ for dimensional reduction and visualization, as well as DBSCAN ⁶⁶ for clustering
252 (**Methods**), which revealed two main clusters (**Fig.5A**). Overlaying the transcription factor
253 activity (TFA) profiles over the cells revealed that *FOXJ1* (a marker for ciliated cells) was
254 significantly more active in the smaller cluster, suggesting that this cluster defines ciliated
255 cells (**Fig.5A**). Confirming this, *FOXJ1* expression was also higher in this cluster, whilst
256 expression of basal (*KRT5*), club (*SCGB1A1*) and goblet (*MUC5AC*) markers were higher in
257 the larger cluster, suggesting that this larger cluster is composed of non-ciliated lung
258 epithelial cells (i.e. basal cells, goblets and club cells) (**Fig.5B**). Of note, *FOXJ1* was one of
259 the few transcription factors for which activity and expression were reasonably well
260 correlated. For instance, *TBX4* exhibited higher regulatory activity in non-ciliated cells
261 (**Fig.5A**), yet it exhibited a 100% dropout rate across all lung epithelial cells (**Fig.5C**). Other
262 key lung-specific TFs with very high expression in lung tissue, as assessed in our GTEX bulk
263 RNA-Seq data, but with 100% dropout rates included *GATA2* and *TBX2* (**Fig.5C**). Thus,
264 SCIRA is able to retrieve biologically relevant variation in regulatory activity of key TFs,
265 when expression alone can not.

266 Despite the tSNE diagram being derived from the regulatory activity profiles of only 38
267 lung-specific TFs, the larger cluster of non-ciliated cells revealed clear segregation of cells
268 according to whether they derived from current or never-smokers, suggesting that smoking
269 exposure has a dramatic effect on the regulatory activity of lung-specific TFs (**Fig.5D**). We
270 verified this by applying PCA to the activity profiles over the non-ciliated cells only (Wilcox
271 test $P=5e-32$, **Fig.5D**). We identified a total of 6 TFs exhibiting significantly lower and 6
272 exhibiting significantly higher regulatory activity in the cells of smokers (**Fig.5E**).
273 Interestingly, among the 6 TFs exhibiting lower activation in cells from smokers, all 6 were
274 also seen to be inactivated in single lung cancer cells, whilst 2 of the 6 exhibiting activation
275 in exposed cells also exhibited increased activity in lung cancer (**Fig.5F**). Among the 6 TFs
276 exhibiting lower activity in both lung epithelial cells of smokers and cancer patients, it is
277 worth noting *NKX2-1*, a putative tumor suppressor for lung cancer as noted recently ⁴⁸, and
278 *TBX4*, another putative tumor suppressor for non-small cell lung cancer ^{67,68}. Among the TFs
279 exhibiting increased regulatory activity in smokers we observed *EHF* (**Figs.5E**), a
280 transcription factor which has been implicated in goblet cell hyperplasia ⁶⁹. Consistent with
281 this, goblet hyperplasia is observed in lung tissue from smokers ⁶⁴, and according to SCIRA
282 *EHF* regulatory activity was correlated with expression of the goblet cell marker *MUC5AC*
283 (**Fig.5G**), whereas *EHF* expression itself was not, highlighting once again that SCIRA can
284 recapitulate biological differential activity patterns not obtainable via TF-expression alone.
285 Given that there is goblet cell expansion in smokers ⁶⁴, the increased regulatory activity of
286 *EHF* and other TFs like *ELF3* in smokers could reflect this increase. Of note *ELF3* becomes

287 inactivated in lung cancer cells (**Fig.5F**), which is consistent with its role in lung epithelial
288 cell differentiation being impaired in cancer^{70 71}.

289

290 **SCIRA is scalable to millions of cells**

291 Finally, we note that SCIRA can estimate regulatory activity in a manner that scales linearly
292 with the number of profiled cells, thus making it easily scalable to scRNA-Seq studies
293 profiling 100s of thousands to a million cells. In the application to the kidney DropSeq
294 dataset (**SI table S5**) which profiled 9190 cells, runtime was under 4 minutes for 4 processing
295 cores, and under 10 minutes with the regulon-inference step in GTEX included. We
296 performed a subsampling analysis on the kidney set, recording runtimes for manageable
297 numbers of cells, fitted linear functions on a log-log scale, and subsequently estimated
298 runtimes for larger scRNA-Seq studies profiling up to a million cells (**Methods**). In a
299 scRNA-Seq study profiling one million cells, SCIRA would take approximately 100 minutes
300 on 4-cores, or only 4 minutes on a 100-node HPC, whereas other methods would run for
301 months on the same 100-node HPC (**Fig.2D**). Only VIPER-D exhibited a marginally
302 improved computational efficiency compared to SCIRA (**Fig.2D**), owing to the fact that the
303 TF-regulons are derived from a database and are thus precomputed. Thus, SCIRA offers
304 scalability where most competing methods do not.

305

306 **Discussion**

307 Dissecting the cellular heterogeneity of cancer, preinvasive lesions and normal tissue at
308 cancer risk is a critically important task for personalized medicine, and it is clear that
309 mapping such cellular heterogeneity needs to be done at single-cell resolution. In the context
310 of cancer risk prediction, the ability to measure gene expression in single normal cells from
311 individuals exposed to an environmental risk factor, could help identify those at most risk of
312 cancer development. Our rationale was to focus on transcription factors that are important for
313 the specification of a given tissue-type, since there is substantial evidence that
314 inactivation/silencing of these transcription factors is an early event in oncogenesis, present
315 in normal cells at risk of neoplastic transformation and thus preceding cancer development
316 itself^{2-4,7-9,72-74}. It follows that identifying such early “tumor suppressor” inactivation events
317 in normal cells at cancer risk in single-cell data could allow prospective identification of
318 individuals at higher risk of cancer development. As demonstrated here, using scRNA-Seq
319 profiles to identify silencing of tissue-specific TFs lacks sensitivity due to the high dropout
320 rate. Instead, we have presented an alternative strategy called SCIRA, which we have very
321 comprehensively validated on many scRNA-Seq datasets profiling normal cells,
322 demonstrating that it can substantially improve the sensitivity and precision to detect correct
323 dynamic TF-activity changes at single-cell resolution.

324

325 Application of SCIRA to two scRNA-Seq datasets profiling both normal and cancer cells
326 revealed preferential inactivation of tissue-specific TFs in the corresponding cancer cells, an
327 important biological and clinical insight, which we would not have obtained had we used
328 differential expression. These results are not only in line with analogous findings obtained in
329 bulk RNA-Seq cancer studies ⁵, but helps to further establish which key tissue-specific TFs
330 are inactivated in cancer epithelial cells independently of changes in stromal composition,
331 which could otherwise confound results. For instance, in a tissue like lung, at least 40% of
332 cells are stromal cells ⁷⁵, and so differential expression changes seen in bulk cancer tissue
333 may not be observed or may not be due to expression changes in the epithelial compartment.
334 On the other hand, some consistency with observations in bulk data should be expected, and
335 in this regard we stress that, unlike SCIRA, differential expression approaches on single-cell
336 data did not reveal any consistent patterns with those observed at the bulk level. This
337 inconsistency between single cell and bulk differential expression in cancer is therefore
338 another important insight which demonstrates the need and added value of SCIRA to uncover
339 key tumor suppressor events. For instance, many of the lung-specific TFs which SCIRA
340 predicts to be inactivated in lung tumor epithelial cells (e.g. *NKX2-1*, *FOXA2*, *FOXJ1*, *AHR*,
341 *HIF3A*) ³⁰ implicate key cancer-pathways (lung development, alveolarization, ciliogenesis,
342 immune-response, hypoxia-response), and their inactivation likely represent key driver events.
343 Supporting this, epigenetically induced silencing of *NKX2-1* has been proposed to be a key
344 driver event in the development of lung cancer ^{48,76}. In the case of colon, our results in the
345 scRNA-Seq data confirm a tumor suppressor role for TFs like *CDX1/CDX2* ⁷⁷, but also serve
346 to reinforce a novel putative tumor suppressor role for *ATOH1* ⁷⁸, for the autophagy inducer
347 *TRIM31* ⁶² and *KLF5* ⁷⁹. Of note, these last three TFs did not exhibit clear significant
348 differential expression changes, yet they were highly significant via analysis with SCIRA.
349 In the application to normal lung cells from smokers and non-smokers, no preferential
350 inactivation of lung-specific TFs in smokers was observed, consistent with observations
351 derived from buccal (squamous epithelial) cells ⁸. This would suggest that in normal cells
352 exposed to a risk factor, such inactivation events may not yet be under significant selection
353 pressure, yet some of the inactivation events, if present, could be important indicators of
354 future cancer risk. In line with this, out of the 6 lung-specific TFs that were observed to be
355 inactivated in normal lung cells from smokers, all 6 were also inactivated in lung cancer cells.
356 This list included *NKX2-1* and *TBX4*, both of which have tumor suppressor functions ^{67,76}. We
357 also observed 6 lung-specific TFs exhibiting increased regulatory activity in cells from
358 smokers, which included *ELF3*, *XBPI* and *EHF*. Interestingly, *EHF* has been implicated as a
359 driver of goblet hyperplasia ⁶⁹, which is observed in the lung tissue of smokers ⁶⁴. Our data
360 supports the view that *EHF* is a marker of goblet cells and that the increased expression in
361 smokers could be due to an increase in relative goblet cell numbers as observed by Duclos ⁶⁴.
362 Whilst *ELF3* has been reported to be a tumor suppressor in many epithelial cancer types, its

363 function has also been observed to be highly cell-type specific with reported oncogenic roles
364 in lung adenoma carcinoma (LUAD)⁸⁰. Here we observed *ELF3* activation in the lung
365 non-ciliated cells from smokers and overexpression in bulk lung squamous cell carcinoma
366 (LUSC) tissue, but inactivation in single lung cancer cells (predominantly LUAD) and no
367 expression change in bulk-tissue LUAD. Thus, in future it will be important to profile larger
368 numbers of cells in the lung epithelial compartment of healthy smokers and non-smokers,
369 including lung cancer patients from LUAD, LSCC, NSCLC subtypes, to determine if
370 differential activity patterns are specific to individual lung-epithelial cell subtypes.

371 Here, and due to obvious limitations on data availability at single-cell resolution, we could
372 not assess the specific mechanism associated with tissue-specific TF silencing in cancer.
373 However, in the context of bulk-tissue data from the TCGA, we have previously shown that
374 the preferential silencing of tissue-specific TFs in cancer is predominantly associated with
375 promoter DNA hypermethylation⁵. Indeed, inactivation through somatic mutation or
376 copy-number loss/deletion is not a frequent event when considering tissue-specific TFs⁵, in
377 contrast to other gene-families like kinases, epigenetic enzymes or membrane receptors
378 which do exhibit more frequent genetic alterations^{81,82}. Thus, it is very likely that the
379 observed inactivation of tissue-specific TFs in individual cancer cells is also associated with
380 promoter DNA hypermethylation.

381

382 In summary, we have presented and validated a computational strategy called SCIRA that can
383 improve the sensitivity and precision to detect regulatory activity changes of key
384 tissue-specific transcription factors in scRNA-Seq data, and that can reveal tumor suppressor
385 events at single-cell resolution which would otherwise not be possible using differential
386 expression. SCIRA has shown that tissue-specific TFs are preferentially inactivated in
387 corresponding cancer cells, suggesting that these could be tumor suppressor driver events.
388 Importantly, SCIRA also provides a scalable framework in which to infer tissue-specific
389 regulatory activity in scRNA-Seq studies profiling even millions of cells. We envisage that
390 SCIRA will be particularly useful for scRNA-Seq studies aiming to identify altered
391 differentiation programs in normal tissue exposed to cancer risk factors, preinvasive lesions
392 and cancer at single-cell resolution. This is important as this may offer clues and insight into
393 the earliest stages of oncogenesis.

394

395

396

397 **Methods**

398

399 **Single cell data and preprocessing**

400 We analyzed scRNA-Seq data from a total of 6 studies:

401 *Lung Differentiation set*: This scRNA-Seq Fluidigm C1 dataset derives from Treutlein et al ³⁵.
402 Normalized (FPKM) data were downloaded from GEO under accession number GSE52583
403 (file: GSE52583.Rda). Data was further transformed using a log₂ transformation adding a
404 pseudocount of 1, so that 0 FPKM values get mapped to 0 in the transformed basis. After
405 quality control, there are a total of 201 single cells assayed at 4 different stages in the
406 developing mouse lung epithelium, including embryonic day E14.5 (n = 45), E16.5 (n = 27),
407 E18.5 (n = 83) and adulthood (n = 46).

408 *Liver Differentiation set*: This scRNA-Seq Fluidigm C1 dataset was derived from Yang et al
409 ³⁶, a study of differentiation of mouse hepatoblasts into hepatocytes and cholangiocytes.
410 Normalized (TPM) data was downloaded from GEO under accession number GSE90047 (file:
411 GSE90047-Singlecell-RNA-seq-TPM.txt). Data was further transformed using a log₂
412 transformation adding a pseudocount of 1, so that 0 TPM values get mapped to 0 in the
413 transformed basis. After quality control, 447 single-cells remained, with 54 single cells at
414 embryonic day 10.5 (E10.5), 70 at E11.5, 41 at E12.5, 65 at E13.5, 70 at 14.5, 77 at 15.5 and
415 70 at E17.5.

416 *Pancreas Differentiation set*: This scRNA-Seq Smart-Seq2 data derives from Yu et al ³⁷,
417 profiling single cells during murine pancreas development, from embryonic stages E9.5 to
418 E17.5. Normalized (TPM) data was downloaded from GEO (GSE115931, file:
419 GSE115931_SmartSeq2.TPM.txt"). Data was further log₂-transformed with a pseudocount of
420 1. After quality control, 2195 cells remained: 113 (E9.5), 211 (E10.5), 263 (E11.5), 252
421 (E12.5), 421 (E13), 338 (E14.5), 242 (E15), 185 (E16.5), 170 (E17.5).

422 *Kidney-organoid Differentiation set*: This scRNA-Seq DropSeq data derives from Wu et al ³⁸,
423 profiling single cells in a kidney organoid differentiation experiment (Takasato protocol)
424 starting out from iPSCs, with 218 cells profiled at day-0, 1741 at day-7, 1169 at day-12, 1097
425 at day-19 and 4965 at day-26. Read count data for all 9190 high quality cells was
426 downloaded from GEO (GSE118184, file: GSE118184_Takasato.iPS.timecourse.txt").
427 Counts were scaled for each cell by the total read count, multiplied by a common scaling
428 factor of 10⁴ and subsequently log₂-transformed with a pseudocount of 1.

429 *Normal and cancer lung tissue dataset*: This scRNA-Seq 10X Chromium dataset was derived
430 from ⁴⁷, a study profiling malignant and non-malignant lung samples from five patients. We
431 downloaded all .Rds files available from ArrayExpress (E-MTAB-6149), which included the
432 processed data and t-SNE coordinates, as well as cluster cell-type assignments. After quality
433 control, a total of 52,698 single-cells remained of which 1709 were annotated as alveolar,
434 5603 as B-cells, 1592 as endothelial cells, 1465 as fibroblasts, 9756 as myeloid cells, 24911
435 as T-cells and 7450 as tumor epithelial cells. A small cluster of 212 cells was annotated as
436 normal epithelial, yet they derived from a malignant sample ⁴⁷, so given this inconsistency we
437 removed these cells from any analysis, as according to us their "normal" nature is far from
438 clear. The alveolar epithelial cell cluster derived mainly from non-malignant samples and was
439 therefore considered most representative of the normal epithelial cells found in lung.

440 *Normal and cancer colon dataset:* This scRNA-Seq Fluidigm C1 dataset is derived from ⁵⁵, a
441 study profiling malignant and non-malignant colon epithelial cells from 11 patients. We
442 downloaded the normal mucosa and tumor epithelial cell FPKM files from GEO under
443 accession number GSE81861. In total there were 160 and 272 normal and tumor epithelial
444 cells, respectively, as determined by the original publication.

445 *Normal lung from smokers and non-smokers.* This scRNA-Seq dataset is derived from ⁶⁴,
446 where FACS sorted lung epithelial cells from 6 never-smokers and 6 smokers were analysed
447 with the CEL-Seq platform. We downloaded the raw UMI counts from GEO under accession
448 number GSE131391. We followed a similar normalization and QC procedure as described in
449 ⁶⁴, although we used a more stringent cell quality criterion, removing any cells with a total
450 UMI count less than 2400. This threshold was chosen because the total UMI count per cell
451 exhibit a natural bimodal distribution, with the value 2400 defining the natural decision
452 boundary between low and high quality cells. This resulted in 564 epithelial cells. For these
453 cells data was further normalized by scaling UMI counts to TPM, adding a pseudocount of 1
454 and finally taking the \log_2 transformation. We note that results reported here were unchanged
455 if not scaling UMI counts, i.e. if using $\log_2(\text{UMI}+1)$.

456

457

458 **Bulk tissue mRNA expression datasets**

459 For applying SCIRA to data from epithelial tissues, we used the bulk RNA-Seq dataset from
460 the GTEX resource ²⁴ to infer regulons. Specifically, the normalized RPKM data was
461 downloaded from the GTEX website and annotated to Entrez gene IDs. Data was then \log_2
462 transformed with a pseudocount of +1. This resulted in a data matrix of 23929 genes and
463 8555 samples, encompassing 30 tissue types (adipose=577, adrenal gland=145, bladder=11,
464 blood=511, blood vessel=689, brain=1259, breast=214, cervix uteri=11, colon=345,
465 esophagus=686, fallopian tube=6, heart=412, kidney=32, liver=119, lung=320, muscle=430,
466 nerve=304, ovary=97, pancreas=171, pituitary=103, prostate=106, salivary gland=57,
467 skin=891, small intestine=88, spleen=104, stomach=193, testis=172, thyroid=323, uterus=83,
468 vagina=96). In addition, we also analyzed the bulk RNA-Seq dataset from the lung TCGA
469 studies ^{53,54}, which was normalized as described in our previous publications ^{5,83}.

470

471 **The SCIRA algorithm**

472 The SCIRA algorithm has two main steps: (i) construction of a tissue-specific regulatory
473 network and (ii) inference of regulatory activity in single cells for the transcription factors
474 (TFs) in the network constructed in step (i).

475 *(i) Construction of tissue-specific regulatory network:* For a given tissue-type, SCIRA infers a
476 corresponding tissue-specific regulatory network using a greedy partial correlation algorithm
477 framework called SEPIRA ⁸. The greedy partial correlation approach is similar in concept to
478 the GENIE3 algorithm ⁸⁴ (which was found to be one of the best performing

479 reverse-engineering methods in the DREAM-5 challenge⁸⁵), in the sense that it infers the
480 candidate regulators for each gene in turn. However, we use partial correlations instead of
481 regression trees. By computing partial correlations over the GTEX dataset, which consists of
482 8555 samples across 30 different tissue-types, it is possible to identify direct regulatory
483 relations that are relevant in the context of differentiation and development. Briefly, having
484 log-transformed the GTEX RNA-Seq set, as described previously⁸, we first select genes with
485 a standard deviation larger than 0.25, so as to remove genes with no significant expression
486 variation across the 8555 samples. A total of 19478 genes with Entrez gene annotation were
487 left after this step. Next, we used a list of 1385 human TFs as defined by the
488 TRANSC_FACT term of the Molecular Signatures Database⁸⁶, of which 1313 had
489 representation in our filtered GTEX set. Genes not annotated as TFs, were considered
490 putative targets, and we first estimated Pearson correlations between the 1313 TFs and the
491 18165 targets. Using a conservative P-value threshold of 1e-6 to define putative interactions
492 between TFs and targets, we next selected TFs with at least 10 putative targets. For each
493 target-gene g and its putative TF regulators f , we then computed partial correlations between
494 g and f , as

$$\tilde{\rho}_{gf} = -\frac{\Omega_{gf}}{\sqrt{\Omega_{gg}\Omega_{ff}}}$$

495 where Ω is the inverse of the expression covariance matrix, which is of dimension
496 $(1+nf)*(1+nf)$ with nf the number of putative TF regulators. Importantly, by estimating the
497 partial correlations in a greedy fashion, i.e. for each target gene separately, the inverse of the
498 covariance matrix is always well defined (no need to estimate a pseudo-inverse) since $nf \ll$
499 8555, i.e much less than the number of samples over which the partial correlations are
500 estimated. In other words, we estimate the partial correlations between each target gene and
501 its candidate regulators from the marginal analysis above, and we do this for each target gene
502 separately, which thus provides a natural regularization. Partial correlation thresholds of +/-
503 0.2, or even +/- 0.1 are statistically significant given the large number of samples (8555) in
504 the GTEX set (as verified by random resampling), so we use either one of these thresholds
505 depending on the number of TFs desired, although the number of resulting TFs is similar for
506 both choices of threshold. This then defines a global regulatory network between TFs and
507 target genes, where indirect dependencies have been removed due to the use of partial
508 correlations⁸⁷.

509 The final step is the construction of a tissue-specific regulatory network as the subnetwork
510 obtained by identification of tissue-specific TFs, i.e. TFs with significantly higher expression
511 in the given tissue type compared to all other tissues combined. This is done using the
512 empirical Bayes moderated t-test framework (limma)⁸⁸. Importantly, a second limma
513 analysis is performed by comparing the tissue of interest to individual tissue types if these
514 contain cells that are believed to significantly infiltrate and contaminate the tissue of interest.

515 Thus, in the case of liver we perform two limma analyses: comparing liver to all other
516 tissue-types, and separately, liver to only blood and spleen combined, since blood/spleen
517 consists of immune cells which are known to infiltrate liver tissue accounting for
518 approximately 40% of all cells found in liver⁷⁵. We require a liver-specific TF to be one with
519 significantly higher expression in both comparisons: when comparing to all tissues we use an
520 adjusted P-value threshold of 0.05 and a log₂(FC) threshold of log₂(1.5)≈0.58, whereas when
521 comparing to blood/spleen we only use an adjusted P-value threshold of 0.05. This ensures
522 that the identified TFs are not driven by a higher immune cell (IC) infiltration in the tissue of
523 interest compared to an “average” tissue where the IC infiltration may be low. As applied to
524 liver and using a significance threshold on partial correlations of +/- 0.2, SCIRA/SEPIRA
525 inferred a network of 22 liver-specific TFs and their regulons, with the average number of
526 genes per regulon being 41, and with range 10 to 151. This network is available as an Rds file
527 “netLIV.Rds” in **Supplementary File 1**. The same procedure was used for the other
528 tissue-types and the corresponding networks for pancreas (netPANC.Rds), kidney
529 (netKID.Rds) and colon (netCOL.Rds) are also available in **Supplementary File 1**.

530 We note that regulon genes could be selected further based on whether they are direct binding
531 targets of the TF, as for instance determined by a ChIP-Seq assay. However, we did not
532 pursue this strategy here, for a number of good reasons. First, the definition of a regulon, as
533 originally proposed by Andrea Califano’s lab^{31,89}, does not require a member of the regulon
534 to be a direct target of the regulator. Indeed, it could well be that a downstream gene in the
535 pathway is an equally good if not even better marker of upstream regulatory activity. Thus, it
536 makes sense to keep all inferred regulon genes in the regulon, following previous studies. On
537 the other hand, some enrichment for direct targets is to be expected, and we indeed checked
538 enrichment for ChIP-Seq binding targets using data from the ChIP-Seq Atlas³². A second
539 reason is that reducing the number of regulon genes also means a loss of power, specially so
540 if the regulon genes are bona-fide markers of upstream regulatory activity. Third, ChIP-Seq
541 data is still very limited in the number of cell-types profiled, which may not include a
542 representative cell-type of the tissue in question. In other words, the sensitivity of a ChIP-Seq
543 assay is also limited and if a gene is not predicted to be a binding target in cell-type “A” it
544 could still be a direct target in the tissue/cell-type of interest.

545

546 *(ii) Estimation of regulatory activity:* Having inferred the tissue-specific TFs and their
547 regulons, we next estimate regulatory activity of the TFs in each single cell of a scRNA-Seq
548 dataset. This is done by regressing the log-normalized scRNA-Seq expression profile of the
549 cell against the “target-profile” of the given TF, where in the target profile, any regulon
550 member is assigned a +1 for activating interactions, a -1 for inhibitory interactions. All other
551 genes not members of the TF’s regulon are assigned a value of 0. The TF-activity is then
552 defined as the t-statistic of this linear regression. Before applying this procedure the
553 normalized scRNA-Seq dataset is z-score normalized, i.e. each gene is centered and scaled to

554 unit standard deviation.

555 We note that SCIRA relies on the tissue-specific regulatory network inferred in step-1. As
556 such, SCIRA is particularly useful for scRNA-Seq studies that profile cells in the tissue of
557 interest, either as part of a developmental or differentiation timecourse experiment, or in the
558 context of diseases where altered differentiation is a key disease hallmark e.g. cancer and
559 precursor cancer lesions.

560

561 **Pseudocode implementing SCIRA algorithm**

562 The previously described steps implementing SCIRA can be run using the functions provided
563 as part of the SEPIRA Bioconductor package, or preferably from the SCIRA-package:
564 <https://github.com/aet21/scira> . Briefly, assuming the normalized GTEx RNA-Seq dataset
565 matrix is stored in an R-object called “*data.m*”, with rows labeling genes and columns
566 labeling samples, and assuming we choose liver as our tissue of interest, we would run the
567 following set of commands in order to construct the liver-specific regulatory network:

568

```
569     > net.o <- sciraInfReg(data=data.m, sdth=0.25, sigth=1e-6, pcorth=0.2, spTH=0.01,  
570       minNtgts=10, ncores=4)  
571     > livernet.o <- sciraSelReg(net.o, tissue=colnames(data.m), toi="Liver", cft="Blood",  
572       degth=c(0.05,0.05), lfcth=c(log2(1.5),0))
```

573

574 In the above *colnames(data.m)* labels the tissue-type of each sample (column) of the data
575 matrix. Note that the parameter *cft* labels the confounding tissue-type, which in this case is
576 blood, because immune-cells, the main component of blood, is a major contaminant cell-type
577 in liver-tissue⁷⁵. One important parameter in the above function, which directly controls the
578 number of retrieved TFs is *spTH*: this parameter controls the number of significant
579 correlations in the marginal analysis to be included in the subsequent partial correlation
580 analysis. By default this is set at 1% of all possible interactions, but increasing this threshold
581 to 5 or 10% will increase the number of interactions and thus the number of retrieved TFs.
582 The tissue-specific regulatory network can be found in the *livernet.o\$netTOI* entry, which is a
583 matrix with columns labeling the tissue-specific transcription factors and rows labeling gene
584 targets. The entries in this matrix are either 1 for a positive interaction, 0 for no interaction,
585 and -1 for inhibitory associations. This matrix provides the regulons to the function for
586 estimating regulatory activity in a bulk sample or in single-cells. For instance, assuming that
587 we have a log-normalized scRNA-Seq dataset representing liver development in humans,
588 *scRNA.m*, we would obtain regulatory activity estimates for each of the transcription factors
589 present in *livernet.o\$netTOI*, by running:

590

```
591     > actTF.m <- sciraEstRegAct(data=scRNA.m, regnet=livernet.o$netTOI,  
592       norm="z",ncores=4)
```

593

594 where the *norm* argument specifies that genes in the *scRNA.m* data matrix should be z-score
595 normalized, before estimating regulatory activity. We note that the output object *actTF.m*
596 would define a matrix with rows labeling the tissue-specific transcription factors and columns
597 labeling the single-cells, and with matrix entries representing regulatory activities. We further
598 note that the tissue-specific regulatory networks derived from GTEX, as used in this work,
599 are provided in **Supplementary File 1**. Full details of how to run *scira* are provided in the
600 vignette of the *scira* R-package.

601

602

603 **Power-calculation for SCIRA**

604 We derived a formula to estimate the sensitivity (which we shall denote by *SE*) of SCIRA to
605 detect highly expressed cell-type specific TFs in a given tissue, as a function of the
606 corresponding cell-type proportion in the tissue. The main parameters affecting the power
607 estimate include the relative sample sizes of the two groups being compared (n_1 and n_2), the
608 average expression effect size e (in effect the average expression fold-change) of the cell-type
609 specific TFs compared to all other cell-types, which will depend on the proportion of the
610 cell-type (w) within the tissue of interest. Indeed, it is not difficult to prove that under
611 reasonable assumptions⁹⁰, the sensitivity (*SE*) is given by the formula

$$612 \quad SE(t, n_1, n_2, e(FC, w)) \approx 2(1 - \int_{-\infty}^t T_A(t', n_1, n_2, e(FC, w, \sigma)) dt')$$

613 where t is the statistic value (we assume a t-statistic) dictating the significance threshold, and
614 where T_A denotes the non-central Student's t-distribution with non-centrality parameter μ
615 equal to

$$\mu = \sqrt{\frac{n_1 n_2}{n_1 + n_2}} e(FC, w, \sigma)$$

616 We note that the effect size e is of the form $|\bar{x}_1 - \bar{x}_2|/\sigma$, i.e. the ratio of the difference in
617 average expression between the two groups divided by a common pooled standard deviation
618 that reflects the intrinsic variance in each group. We note that we are assuming that the bulk
619 RNA-Seq data has been log-normalized so that e is derived from the log-transformed data.
620 For instance, if a gene (say a TF) shows the same expression distribution for all cell-types in
621 the tissue of interest compared to all other tissues, then $\bar{x}_1 - \bar{x}_2 = \log_2(\bar{I}_1/\bar{I}_2)$ where \bar{I}_i
622 denotes the average intensity (i.e. FPKM/TPM) value in group- i . Assuming that the given TF
623 is only more highly expressed in a cell-type that makes up only a proportion w of the cells in
624 the tissue of interest, then $e = \log_2[FC * w + 1 * (1 - w)]/\sigma$ where FC is the average
625 fold-change. To estimate the sample sizes for the power calculation, we note that the median
626 number of samples per tissue-type in GTEX is approximately 170. We took a more
627 conservative value of $n_i=150$ to represent the number of samples in our tissue of interest,

628 with the rest of samples in GTEX, i.e. $n_2=8555-150=8405$, defining the number of samples
629 from other tissue-types. To estimate the average expression fold-change FC for top DEGs
630 between single-cell types in a tissue, we analysed expression data from purified FACS sorted
631 luminal and basal cells from the mammary epithelium⁹¹. Because FACS sorted cell
632 populations are still heterogeneous, we thus expect the resulting fold change estimates to be
633 conservative. Using limma⁸⁸, we estimated FC to be 8 for the highest ranked DEGs, and
634 approximately 6 for the top 200-300 DEGs. We note that these estimates are for a scaled basis
635 where $\sigma=1$. Thus, we approximate the effect size $e \approx \log_2[FC * w + 1 * (1 - w)]$ with
636 $FC=8$ or 6 , so as to consider two different effect size scenarios. For the proportion w we
637 assumed two values: $w=0.05$ and $w=0.2$ representing 5 and 20% of the cells in the tissue of
638 interest. Note that if $w=1$, all cells within the issue of interest exhibit differential expression
639 at magnitude FC and if $w=0$, no cell is differentially expressed. Finally, to compute the
640 sensitivity as a function of the significance level threshold t , we used the parameters above as
641 input to the TOC function of the OCplus R-package⁹⁰.

642

643

644 **Implementation of scImpute, MAGIC and SCRABBLE**

645 scImpute (version 0.0.9)³⁹ was run with default parameters (labeled =FALSE, drop_thre=0.5)
646 in all analysis, with the exception of the Kcluster parameter, which was chosen to reflect the
647 number of underlying cell-types in each tissue analysed: Liver=3, Lung=4, Pancreas=15,
648 Kidney= 14, i.e this parameter was set for each tissue following the numbers of cell-types as
649 specified in the original papers. For MAGIC (version 1.4.0)⁴⁰ in liver, lung and pancreas, we
650 used the following parameters: $k = 15$, $\alpha= 5$, $t = \text{"auto"}$, $\text{knn_dist.method} = \text{"euclidean"}$.
651 For kidney, because of the much larger number of cells, we chose larger values for $k=30$ and
652 $\alpha=10$. The number of PCs (npca) was determined in all tissues as the number of PCs
653 explaining 70% of variation in the data, as recommended⁴⁰. For SCRABBLE (version 0.0.1)
654⁴¹, the average bulk RNA-seq expression vector was computed using the corresponding
655 tissue-type samples from the GTEX dataset. The alpha parameter in the function was chosen
656 for each tissue-type, following the recommendations given in the paper: Liver=1, Lung=1,
657 Pancreas=0.1, Kidney=0.1. The other parameter values were $\beta = 1e-5$, $\gamma = 0.01$. For
658 all other parameters, we used the default choices: $nIter = 20$, $\text{error_out_threshold} = 1e-04$,
659 $nIter_inner = 20$, $\text{error_inner_threshold} = 1e-04$.

660

661

662 **Implementation of GENIE3 and SCENIC**

663 SCENIC is a pipeline of 3 distinct methods (GENIE3, RcisTarget, AUCCell), each with its own
664 Bioconductor package. We used the following versions: GENIE3_1.4.0, RcisTarget_1.2.0 and
665 AUCCell_1.4.1. Because the lung, liver and pancreas scRNA-Seq sets are from mice, we used
666 as regulators a list of 1686 mouse TF from the RIKEN lab (<http://genome.gsc.riken.jp/TFdb/>)

667 together with the homologs of the human TFs in our lung, liver and pancreas specific
668 networks if these were not in the RIKEN lab list. GENIE3 was run with default parameter
669 choices (treeMethod="RF", K="sqrt", nTrees=1000) but on a reduced data matrix where
670 genes with a standard deviation less than 0.5 were removed. Regulons of TFs were obtained
671 from GENIE3 using a threshold on the inferred weights (representing the regulatory strength
672 and termed "importance measure" in GENIE3) of 0.01, and only positively correlated targets
673 were selected using a Spearman correlation coefficient threshold > 0 . In SCENIC, the targets
674 are then scanned for enriched binding motifs using RcisTarget. We used the 7species.mc9nr
675 feather files for both 500bp upstream of the TSS and also for a 20kb window centered on the
676 TSS. Any enriched motifs in both analyses were combined to arrive at a single list of
677 enriched motifs and associated TFs. We then found the overlap with the annotated TFs from
678 GENIE3, and only those that overlapped were considered valid TF regulons. For these we
679 then estimated a regulatory activity score using an approach similar to the one implemented
680 in AUCell, but one that is threshold independent, and therefore an improvement over the
681 method used in AUCell. Specifically, the activity score was defined as the AUC of a
682 Wilcoxon rank sum test, whereby in each single cell, genes are first ranked in decreasing
683 order of expression, and the AUC-statistic is then derived by comparing the ranks of the
684 regulon (all positively correlated) genes to the ranks of all other genes.

685

686 **Implementation of VIPER-D**

687 In order to assess the importance of the tissue-specific regulons used in SCIRA, we compared
688 SCIRA to a method that uses non tissue-specific TF-regulons. We note that there are tools
689 like PAGODA⁹² that can infer activity scores from gene sets, yet a regulon also entails
690 directionality (i.e. positive or inhibitory interaction) information, which also needs to be
691 assessed. Hence, motivated by the recent work by Holland et al⁴⁶, we decided to test SCIRA
692 against the combined use of VIPER⁴³ and the dorothea TF-regulon database⁴⁵. Of note,
693 VIPER infers regulatory activity in any given sample/cell given a TF-regulon, and that the
694 dorothea TF-regulon database is not tissue-specific, although one of the sources in building
695 dorothea is the same GTEX dataset used by SCIRA to build its tissue-specific regulons. We
696 ran viper with the following argument choices: dnull = NULL, pleiotropy = FALSE, nes =
697 TRUE, method = c("none"), bootstraps = 0, minsize = 5, adaptive.size = FALSE, eset.filter =
698 TRUE, mvws = 1, cores = 4. Dorothea also provides likelihood information that a given
699 regulatory interaction in the database is true, and VIPER allows such likelihood information
700 to be used when inferring regulatory activity. We ran VIPER-D in two ways: (i) assigning the
701 same likelihood to all listed regulatory interactions (ie equal weights), and (ii) by using the
702 likelihood information. In Dorothea, the likelihood is encoded as an ordinal categorical
703 variable: A, B, C, D, E, with A indicating highest confidence. In order to run this with VIPER,
704 we transformed these categories into confidence weights using the mapping: A=1, B=0.8,
705 C=0.6, D=0.4, E=0.2 . Results in this manuscript are reported for the case of equal weights.

706 We note that these likelihoods vary mostly between TFs, and not between the targets of a
707 given TF, which is why results are largely unchanged had we used the likelihood information.

708

709

710 **Differential Expression (DE) analysis**

711 In this work we compare SCIRA to ordinary DE analysis, as implemented using a Wilcoxon
712 rank sum test for binary phenotypes, or using non-parametric Spearman rank correlations for
713 ordinal phenotypes (e.g. multiple timepoints or stages). The use of a non-parametric test,
714 which is distribution assumption free, works well for scRNA-Seq with high dropout rates.
715 When comparing statistics of differential activity from SCIRA to those from DE analysis, we
716 transform Wilcoxon rank sum or Spearman test P-values into z-statistics using a quantile
717 normal distribution, taking into account the magnitude of the AUC value from the Wilcoxon
718 test (i.e. AUC values > 0.5 correspond to higher expression in one group compared to other,
719 whereas AUC < 0.5 represents the opposite case), or the sign of the Spearman correlation
720 coefficient in the case of ordinal phenotypes.

721

722 **Comparative sensitivity and precision analysis**

723 We compared SCIRA to seven other methods in their sensitivity and precision to identify
724 gold-standard sets of tissue-specific TFs. These gold-standard sets were constructed from
725 GTEX and validated in orthogonal bulk tissue gene expression datasets from NormalAtlas³³
726 and Roth et al³⁴. The number tissue-specific TFs for liver, lung, pancreas and kidney were 22,
727 38, 30 and 38, respectively. The seven other methods were ordinary differential DE analysis,
728 scImpute+DE, MAGIC+DE, Scrabble+DE, GENIE3, SCENIC and VIPER-D. We note that
729 SCENIC runs GENIE3 as a first step and then selects TF-regulons for which corresponding
730 TF-binding motifs are enriched. So, for the method denoted “GENIE3” we drop the
731 requirement of TF-binding motif enrichment. For SCIRA, GENIE3, SCENIC and VIPER-D
732 we obtain TF-activity estimates, whereas the other methods rely on direct gene expression,
733 measured or imputed. Sensitivity (SE) was estimated as the fraction of gold-standard TFs
734 which exhibited significant increased activation/expression with differentiation timepoint, as
735 determined using a Bonferroni adjusted $P < 0.05$ threshold. Precision equals $1 - \text{FDR}$ (false
736 discovery rate), with the FDR defined by the ratio of significantly inactivated TFs to the total
737 number of significantly differentially active TFs, since inactivation of these TFs is
738 inconsistent with known biology and therefore represent false positives. Correspondingly, for
739 methods relying on differential expression, the FDR is defined by the ratio of significantly
740 downregulated TFs to the total number of significantly differentially expressed TFs.

741

742 **Comparative runtime and scalability analysis**

743 Objective comparison of runtimes of the different algorithms is hard because each method
744 has different requirements for input, and because runtimes depend critically on the choice of

745 method-specific parameters. Nevertheless, we compared runtimes for 5 important algorithms
746 (SCIRA, MAGIC, Scrabble, GENIE3/SCENIC and VIPER-D), both in terms of their actual
747 implementations on the liver, lung, and pancreas and kidney sets, but also in a scaling
748 analysis with largely default parameters, where we applied all 5 methods to varying subsets
749 of the kidney scRNA-Seq set (total 9190 cells). Briefly, we processed the scRNA-Seq kidney
750 DropSeq data as described earlier, and filtered genes with sufficient variance resulting in
751 12596 genes. We then constructed subsets with variable cell numbers by randomly
752 subsampling 200, 400, 600, 800, 1000 and 1500 cells, and ran each of these methods on each
753 of these subsampled datasets. In the case of SCIRA, MAGIC, GENIE3/SCENIC and
754 VIPER-D we ran the algorithms with 4 processing cores on a Dell PowerEdge server with
755 Intel Xeon CPU E5-4660 v4 and clock speed of 2.20GHz. Unfortunately, Scrabble does not
756 offer a parallelizable option and is excruciatingly slow for larger e.g. a 10,000 cell dataset.
757 Thus, for each method, we obtained runtimes as a function of cell-number, and fitted a linear
758 regression to the data on a log-log scale. On a log-log scale where both runtime and
759 cell-number are logged, the relation is generally linear. Next, we imputed runtimes for much
760 larger datasets up to a million cells.

761

762 **Data Availability:** Data analyzed in this manuscript is already publicly available from the
763 following GEO (www.ncbi.nlm.nih.gov/geo/) accession numbers: GSE52583, GSE90047,
764 GSE115931, GSE118184, GSE81861, GSE131391, and from ArrayExpress
765 (www.ebi.ac.uk/arrayexpress) under accession number E-MTAB-6149.

766

767 **Code Availability:** The scira R-package is freely available from
768 <https://github.com/aet21/scira>

769

770 **Additional Files:** Supplementary File 1 contains R (.Rds) object files, containing the inferred
771 regulatory networks for colon (netCOL.Rds), kidney (netKID.Rds), liver (netLIV.Rds), lung
772 (netLUNG.Rds) and pancreas (netPANC.Rds). Supplementary Information File contains all
773 Supplementary Figures and Supplementary Tables.

774

775 **Ethics:** All data analyzed in this manuscript is freely available in the public domain, and so
776 no Ethics statement is required as all primary data was already presented elsewhere.

777

778 **Competing Interests:** The authors declare that there are no competing interests.

779

780 **Author Contribution:** Study was conceived and designed by AET. Statistical analyses were
781 performed by AET and replicated by NW. Software package was prepared by NW.
782 Manuscript was written by AET.

783

784 **Acknowledgements**

785 This work was supported by NSFC (National Science Foundation of China) grants, grant
786 numbers 31571359 and 31771464 and by a Royal Society Newton Advanced Fellowship
787 (NAF award number: 164914). We would also like to thank Peter Kharchenko for useful
788 discussions.

789

790 **References**

- 791 1. Heinaniemi, M. *et al.* Gene-pair expression signatures reveal lineage control. *Nat Methods* **10**, 577-83
792 (2013).
- 793 2. Ohm, J.E. *et al.* A stem cell-like chromatin pattern may predispose tumor suppressor genes to DNA
794 hypermethylation and heritable silencing. *Nat Genet* **39**, 237-42 (2007).
- 795 3. Baylin, S.B. & Ohm, J.E. Epigenetic gene silencing in cancer - a mechanism for early oncogenic pathway
796 addiction? *Nat Rev Cancer* **6**, 107-16 (2006).
- 797 4. Schlesinger, Y. *et al.* Polycomb-mediated methylation on Lys27 of histone H3 pre-marks genes for de
798 novo methylation in cancer. *Nat Genet* **39**, 232-6 (2007).
- 799 5. Teschendorff, A.E. *et al.* The multi-omic landscape of transcription factor inactivation in cancer.
800 *Genome Med* **8**, 89 (2016).
- 801 6. Hanahan, D. & Weinberg, R.A. Hallmarks of cancer: the next generation. *Cell* **144**, 646-74 (2011).
- 802 7. Feinberg, A.P., Ohlsson, R. & Henikoff, S. The epigenetic progenitor origin of human cancer. *Nat Rev*
803 *Genet* **7**, 21-33 (2006).
- 804 8. Chen, Y., Widschwendter, M. & Teschendorff, A.E. Systems-epigenomics inference of transcription
805 factor activity implicates aryl-hydrocarbon-receptor inactivation as a key event in lung cancer
806 development. *Genome Biol* **18**, 236 (2017).
- 807 9. Zheng, S.C., Widschwendter, M. & Teschendorff, A.E. Epigenetic drift, epigenetic clocks and cancer risk.
808 *Epigenomics* **8**, 705-19 (2016).
- 809 10. Spira, A. *et al.* Precancer Atlas to Drive Precision Prevention Trials. *Cancer Res* **77**, 1510-1541 (2017).
- 810 11. Regev, A. *et al.* The Human Cell Atlas. *Elife* **6**(2017).
- 811 12. Rozenblatt-Rosen, O., Stubbington, M.J.T., Regev, A. & Teichmann, S.A. The Human Cell Atlas: from
812 vision to reality. *Nature* **550**, 451-453 (2017).
- 813 13. Stubbington, M.J.T., Rozenblatt-Rosen, O., Regev, A. & Teichmann, S.A. Single-cell transcriptomics to
814 explore the immune system in health and disease. *Science* **358**, 58-63 (2017).
- 815 14. Moris, N., Pina, C. & Arias, A.M. Transition states and cell fate decisions in epigenetic landscapes. *Nat*
816 *Rev Genet* **17**, 693-703 (2016).
- 817 15. Puram, S.V. *et al.* Single-Cell Transcriptomic Analysis of Primary and Metastatic Tumor Ecosystems in
818 Head and Neck Cancer. *Cell* **171**, 1611-1624 e24 (2017).
- 819 16. Tirosh, I. *et al.* Single-cell RNA-seq supports a developmental hierarchy in human oligodendrogloma.
820 *Nature* **539**, 309-313 (2016).
- 821 17. Patel, A.P. *et al.* Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma.
822 *Science* **344**, 1396-401 (2014).
- 823 18. Tirosh, I. *et al.* Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-seq.
824 *Science* **352**, 189-96 (2016).
- 825 19. Todorov, H., Cannoodt, R., Saelens, W. & Saeys, Y. Network Inference from Single-Cell Transcriptomic

- 826 Data. *Methods Mol Biol* **1883**, 235-249 (2019).
- 827 20. Stegle, O., Teichmann, S.A. & Marioni, J.C. Computational and analytical challenges in single-cell
828 transcriptomics. *Nat Rev Genet* **16**, 133-45 (2015).
- 829 21. Eling, N., Morgan, M.D. & Marioni, J.C. Challenges in measuring and understanding biological noise.
830 *Nat Rev Genet* **20**, 536-548 (2019).
- 831 22. Chen, S. & Mar, J.C. Evaluating methods of inferring gene regulatory networks highlights their lack of
832 performance for single cell gene expression data. *BMC Bioinformatics* **19**, 232 (2018).
- 833 23. Grun, D. Revealing dynamics of gene expression variability in cell state space. *Nat Methods* **17**, 45-49
834 (2020).
- 835 24. Consortium, G.T. The Genotype-Tissue Expression (GTEx) project. *Nat Genet* **45**, 580-5 (2013).
- 836 25. Teschendorff, A.E. & Relton, C.L. Statistical and integrative system-level analysis of DNA methylation
837 data. *Nat Rev Genet* **19**, 129-147 (2018).
- 838 26. Gao, F., Foat, B.C. & Bussemaker, H.J. Defining transcriptional networks through integrative modeling
839 of mRNA expression and transcription factor binding data. *BMC Bioinformatics* **5**, 31 (2004).
- 840 27. Ludtke, T.H. *et al.* Tbx2 controls lung growth by direct repression of the cell cycle inhibitor genes
841 Cdkn1a and Cdkn1b. *PLoS Genet* **9**, e1003189 (2013).
- 842 28. Wan, H. *et al.* Foxa2 is required for transition to air breathing at birth. *Proc Natl Acad Sci U S A* **101**,
843 14449-54 (2004).
- 844 29. Wan, H. *et al.* Foxa2 regulates alveolarization and goblet cell hyperplasia. *Development* **131**, 953-64
845 (2004).
- 846 30. Herriges, M. & Morrisey, E.E. Lung development: orchestrating the generation and regeneration of a
847 complex organ. *Development* **141**, 502-13 (2014).
- 848 31. Margolin, A.A. *et al.* ARACNE: an algorithm for the reconstruction of gene regulatory networks in a
849 mammalian cellular context. *BMC Bioinformatics* **7 Suppl 1**, S7 (2006).
- 850 32. Oki, S. *et al.* ChIP-Atlas: a data-mining suite powered by full integration of public ChIP-seq data. *EMBO*
851 *Rep* **19**(2018).
- 852 33. Uhlen, M. *et al.* Proteomics. Tissue-based map of the human proteome. *Science* **347**, 1260419 (2015).
- 853 34. Roth, R.B. *et al.* Gene expression analyses reveal molecular relationships among 20 regions of the
854 human CNS. *Neurogenetics* **7**, 67-80 (2006).
- 855 35. Treutlein, B. *et al.* Reconstructing lineage hierarchies of the distal lung epithelium using single-cell
856 RNA-seq. *Nature* **509**, 371-5 (2014).
- 857 36. Yang, L. *et al.* A single-cell transcriptomic analysis reveals precise pathways and regulatory
858 mechanisms underlying hepatoblast differentiation. *Hepatology* **66**, 1387-1401 (2017).
- 859 37. Yu, X.X. *et al.* Defining multistep cell fate decision pathways during pancreatic development at
860 single-cell resolution. *EMBO J* **38**(2019).
- 861 38. Wu, H. *et al.* Comparative Analysis and Refinement of Human PSC-Derived Kidney Organoid
862 Differentiation with Single-Cell Transcriptomics. *Cell Stem Cell* **23**, 869-881 e8 (2018).
- 863 39. Li, W.V. & Li, J.J. An accurate and robust imputation method scImpute for single-cell RNA-seq data. *Nat*
864 *Commun* **9**, 997 (2018).
- 865 40. van Dijk, D. *et al.* Recovering Gene Interactions from Single-Cell Data Using Data Diffusion. *Cell* **174**,
866 716-729 e27 (2018).
- 867 41. Peng, T., Zhu, Q., Yin, P. & Tan, K. SCRABBLE: single-cell RNA-seq imputation constrained by bulk
868 RNA-seq data. *Genome Biol* **20**, 88 (2019).
- 869 42. Aibar, S. *et al.* SCENIC: single-cell regulatory network inference and clustering. *Nat Methods* **14**,

- 870 1083-1086 (2017).
- 871 43. Alvarez, M.J. *et al.* Functional characterization of somatic mutations in cancer using network-based
872 inference of protein activity. *Nat Genet* **48**, 838-47 (2016).
- 873 44. Ding, H. *et al.* Quantitative assessment of protein activity in orphan tissues and single cells using the
874 metaVIPER algorithm. *Nat Commun* **9**, 1471 (2018).
- 875 45. Garcia-Alonso, L., Holland, C.H., Ibrahim, M.M., Turei, D. & Saez-Rodriguez, J. Benchmark and
876 integration of resources for the estimation of human transcription factor activities. *Genome Res* **29**,
877 1363-1375 (2019).
- 878 46. Holland, C.H. *et al.* Robustness and applicability of transcription factor and pathway analysis tools on
879 single-cell RNA-seq data. *Genome Biol* **21**, 36 (2020).
- 880 47. Lambrechts, D. *et al.* Phenotype molding of stromal cells in the lung tumor microenvironment. *Nat*
881 *Med* **24**, 1277-1289 (2018).
- 882 48. Teixeira, V.H. *et al.* Deciphering the genomic, epigenomic, and transcriptomic landscapes of
883 pre-invasive lung cancer lesions. *Nat Med* (2019).
- 884 49. Zhu, Y., Li, Y., Jun Wei, J.W. & Liu, X. The role of Sox genes in lung morphogenesis and cancer. *Int J Mol*
885 *Sci* **13**, 15767-83 (2012).
- 886 50. Li, Q.F., Wang, X.R., Yang, Y.W. & Lin, H. Hypoxia upregulates hypoxia inducible factor (HIF)-3alpha
887 expression in lung epithelial cells: characterization and comparison with HIF-1alpha. *Cell Res* **16**,
888 548-58 (2006).
- 889 51. Boule, L.A. *et al.* Activation of the aryl hydrocarbon receptor during development enhances the
890 pulmonary CD4+ T-cell response to viral infection. *Am J Physiol Lung Cell Mol Physiol* **309**, L305-13
891 (2015).
- 892 52. Beamer, C.A. & Shepherd, D.M. Role of the aryl hydrocarbon receptor (AhR) in lung inflammation.
893 *Semin Immunopathol* **35**, 693-704 (2013).
- 894 53. Cancer Genome Atlas Research, N. Comprehensive genomic characterization of squamous cell lung
895 cancers. *Nature* **489**, 519-25 (2012).
- 896 54. Cancer Genome Atlas Research, N. Comprehensive molecular profiling of lung adenocarcinoma.
897 *Nature* **511**, 543-50 (2014).
- 898 55. Li, H. *et al.* Reference component analysis of single-cell transcriptomes elucidates cellular
899 heterogeneity in human colorectal tumors. *Nat Genet* **49**, 708-718 (2017).
- 900 56. Lynch, J., Keller, M., Guo, R.J., Yang, D. & Traber, P. Cdx1 inhibits the proliferation of human colon
901 cancer cells by reducing cyclin D1 gene expression. *Oncogene* **22**, 6395-407 (2003).
- 902 57. McConnell, B.B., Ghaleb, A.M., Nandan, M.O. & Yang, V.W. The diverse functions of Kruppel-like
903 factors 4 and 5 in epithelial biology and pathobiology. *Bioessays* **29**, 549-57 (2007).
- 904 58. Yang, Q., Bermingham, N.A., Finegold, M.J. & Zoghbi, H.Y. Requirement of Math1 for secretory cell
905 lineage commitment in the mouse intestine. *Science* **294**, 2155-8 (2001).
- 906 59. Ishibashi, F. *et al.* Contribution of ATOH1(+) Cells to the Homeostasis, Repair, and Tumorigenesis of the
907 Colonic Epithelium. *Stem Cell Reports* **10**, 27-42 (2018).
- 908 60. Kazanjian, A. & Shroyer, N.F. NOTCH Signaling and ATOH1 in Colorectal Cancers. *Curr Colorectal Cancer*
909 *Rep* **7**, 121-127 (2011).
- 910 61. Nakaya, T. *et al.* KLF5 regulates the integrity and oncogenicity of intestinal stem cells. *Cancer Res* **74**,
911 2882-91 (2014).
- 912 62. Ra, E.A. *et al.* TRIM31 promotes Atg5/Atg7-independent autophagy in intestinal cells. *Nat Commun* **7**,
913 11726 (2016).

- 914 63. Burada, F. *et al.* Autophagy in colorectal cancer: An important switch from physiology to pathology.
915 *World J Gastrointest Oncol* **7**, 271-84 (2015).
- 916 64. Duclos, G.E. *et al.* Characterizing smoking-induced transcriptional heterogeneity in the human
917 bronchial epithelium at single-cell resolution. *Sci Adv* **5**, eaaw3413 (2019).
- 918 65. van der Maaten, L. Visualizing Data using t-SNE. *J Mach Learn Res* **9**, 2579-2605 (2008).
- 919 66. Ester, M., Kriegel, H.P., Sander, J. & Xu, X. A Density-Based Algorithm for Discovering Clusters in Large
920 Spatial Databases with Noise. in *2nd International Conference on Knowledge Discovery and Data
921 Mining (KDD-96)* (Institute for Computer Science, University of Munich, 1996).
- 922 67. Nehme, E. *et al.* Epigenetic Suppression of the T-box Subfamily 2 (TBX2) in Human Non-Small Cell Lung
923 Cancer. *Int J Mol Sci* **20**(2019).
- 924 68. Lai, I.L. *et al.* Male-Specific Long Noncoding RNA TTTY15 Inhibits Non-Small Cell Lung Cancer
925 Proliferation and Metastasis via TBX4. *Int J Mol Sci* **20**(2019).
- 926 69. Fossum, S.L. *et al.* Ets homologous factor (EHF) has critical roles in epithelial dysfunction in airway
927 disease. *J Biol Chem* **292**, 10938-10949 (2017).
- 928 70. Oliver, J.R. *et al.* Elf3 plays a role in regulating bronchiolar epithelial repair kinetics following Clara
929 cell-specific injury. *Lab Invest* **91**, 1514-29 (2011).
- 930 71. Luk, I.Y., Reehorst, C.M. & Mariadason, J.M. ELF3, ELF5, EHF and SPDEF Transcription Factors in Tissue
931 Homeostasis and Cancer. *Molecules* **23**(2018).
- 932 72. Widschwendter, M. *et al.* Epigenetic stem cell signature in cancer. *Nat Genet* **39**, 157-8 (2007).
- 933 73. Teschendorff, A.E. *et al.* Age-dependent DNA methylation of genes that are suppressed in stem cells is
934 a hallmark of cancer. *Genome Res* **20**, 440-6 (2010).
- 935 74. Teschendorff, A.E. *et al.* Epigenetic variability in cells of normal cytology is associated with the risk of
936 future morphological transformation. *Genome Med* **4**, 24 (2012).
- 937 75. Zheng, S.C. *et al.* A novel cell-type deconvolution algorithm reveals substantial contamination by
938 immune cells in saliva, buccal and cervix. *Epigenomics* **10**, 925-940 (2018).
- 939 76. Winslow, M.M. *et al.* Suppression of lung adenocarcinoma progression by Nkx2-1. *Nature* **473**, 101-4
940 (2011).
- 941 77. Hryniuk, A., Grainger, S., Savory, J.G. & Lohnes, D. Cdx1 and Cdx2 function as tumor suppressors. *J Biol
942 Chem* **289**, 33343-54 (2014).
- 943 78. Bossuyt, W. *et al.* Atonal homolog 1 is a tumor suppressor gene. *PLoS Biol* **7**, e39 (2009).
- 944 79. Diakiv, S.M., D'Andrea, R.J. & Brown, A.L. The double life of KLF5: Opposing roles in regulation of
945 gene-expression, cellular function, and transformation. *IUBMB Life* **65**, 999-1011 (2013).
- 946 80. Enfield, K.S.S. *et al.* Epithelial tumor suppressor ELF3 is a lineage-specific amplified oncogene in lung
947 adenocarcinoma. *Nat Commun* **10**, 5438 (2019).
- 948 81. Plass, C. *et al.* Mutations in regulators of the epigenome and their connections to global chromatin
949 patterns in cancer. *Nat Rev Genet* **14**, 765-80 (2013).
- 950 82. Vogelstein, B. *et al.* Cancer genome landscapes. *Science* **339**, 1546-58 (2013).
- 951 83. Yang, Z., Jones, A., Widschwendter, M. & Teschendorff, A.E. An integrative pan-cancer-wide analysis of
952 epigenetic enzymes reveals universal patterns of epigenomic deregulation in cancer. *Genome Biol* **16**,
953 140 (2015).
- 954 84. Huynh-Thu, V.A., Irrthum, A., Wehenkel, L. & Geurts, P. Inferring regulatory networks from expression
955 data using tree-based methods. *PLoS One* **5**(2010).
- 956 85. Marbach, D. *et al.* Wisdom of crowds for robust gene network inference. *Nat Methods* **9**, 796-804
957 (2012).

- 958 86. Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for interpreting
959 genome-wide expression profiles. *Proc Natl Acad Sci U S A* **102**, 15545-50 (2005).
- 960 87. Opgen-Rhein, R. & Strimmer, K. From correlation to causation networks: a simple approximate
961 learning algorithm and its application to high-dimensional plant gene expression data. *BMC Syst Biol* **1**,
962 37 (2007).
- 963 88. Smyth, G.K. Linear models and empirical bayes methods for assessing differential expression in
964 microarray experiments. *Stat Appl Genet Mol Biol* **3**, Article3 (2004).
- 965 89. Basso, K. *et al.* Reverse engineering of regulatory networks in human B cells. *Nat Genet* **37**, 382-90
966 (2005).
- 967 90. Pawitan, Y., Michiels, S., Koscielny, S., Gusnanto, A. & Ploner, A. False discovery rate, sensitivity and
968 sample size for microarray studies. *Bioinformatics* **21**, 3017-24 (2005).
- 969 91. Shehata, M. *et al.* Phenotypic and functional characterization of the luminal cell hierarchy of the
970 mammary gland. *Breast Cancer Res* **14**, R134 (2012).
- 971 92. Fan, J. *et al.* Characterizing transcriptional heterogeneity through pathway and gene set
972 overdispersion analysis. *Nat Methods* **13**, 241-4 (2016).
- 973
- 974
- 975

976 **Figure Legends**

977 **Figure-1: SCIRA rationale and workflow.** **A)** Since bulk RNA-Seq data does not suffer
978 from technical dropouts and is much more reliable than scRNA-Seq data, for a given choice
979 of tissue, we use the high-powered GTEX bulk RNA-Seq expression set (>20,000 genes,
980 8555 samples, 30 tissue-types) to derive a corresponding tissue-specific regulatory network,
981 consisting of a gold-standard list of tissue-specific transcription factors (TFs) and their targets
982 (regulons). The inference of the network uses a greedy partial correlation framework, whilst
983 also adjusting for stromal (immune cell) contamination within the tissue. **B)**
984 Power/Sensitivity (SE) estimates to detect tissue-specific TFs in the GTEX bulk RNA-Seq
985 dataset as a function of the minor cell-type fraction (MCF) (left), number of samples in the
986 tissue of interest (middle) and average fold change of differential expression between the
987 tissue of interest and the rest of tissues in GTEX (right). In left panel, we depict SE curves for
988 4 tissue types in GTEX (number of samples in each tissue is given) and for an average FC=8.
989 In the middle panel we depict SE curves for two MCF values, as indicated. In the right panel,
990 we assume a sample size of 150. A MCF value of 0.05 means we assume that the
991 tissue-specific TFs is only highly expressed in 5% of the tissue resident cells. **C)** Given the
992 high technical dropout rate and overall noisy nature of scRNA-Seq data, it may not be
993 possible to reliably infer regulatory activity from the TF expression profile alone. However,
994 using the TF regulons derived in A), and using the genes within the regulon that are not
995 strongly affected by dropouts, we can estimate regulatory activity across single-cells.
996 Depicted is an example with 3 lung-specific TFs (*Sox18*, *Tbx4*, *Foxa2*), as well as the
997 expression pattern of the regulon genes for *Tbx4*, in the context of a lung development study

998 from embryonic day-10 to adult stage (Treutlein dataset). We use linear regressions between
999 the expression values of all the genes in a given cell and the corresponding TF-regulon
1000 profile, to derive the activity of the TF as the t-statistic of the estimated regression coefficient,
1001 resulting in a regulatory activity map over the tissue-specific TFs and single cells. The same
1002 tissue-specific TFs and their regulons can be applied to normal-cancer scRNA-Seq datasets to
1003 infer regulatory activity maps across normal and cancer cells.

1004

1005 **Figure-2: SCIRA displays improved sensitivity, precision and scalability.** **A)** Barplots
1006 with 95% confidence intervals included displaying the sensitivity (SE) to detect increased
1007 activity or expression for a gold-standard set of tissue-specific TFs in a corresponding
1008 timecourse differentiation scRNA-Seq study. Methods represented are SCIRA, ordinary
1009 differential expression (DE), imputation with scImpute, MAGIC or Scrabble following by
1010 DE, SCENIC, running SCENIC without the TF-binding motif enrichment step (denoted
1011 “GENIE3”) and VIPER using the dorothea regulon database (denoted “VIPER-D”). **B)**
1012 Barplots and 95% confidence intervals displaying the false discovery rate (FDR) of each
1013 method in the same scRNA-Seq datasets. Precision is defined as 1-FDR and is the fraction of
1014 true positives among all positives. In this case, tissue-specific TFs predicted to be
1015 significantly downregulated/inactivated during the timecourse were identified as false
1016 positives with FDR defined as the fraction of false positives among all significantly
1017 differentially expressed (or activated) TFs. **C)** Heatmap of P-values assessing the
1018 improvement of SCIRA over the other 7 methods, in terms of both sensitivity (left) and FDR
1019 (right). P-values for each tissue were derived from a one-tailed Binomial test. The P-values
1020 for the meta-analysis (“Meta”) were derived using Fisher’s method. The FDR for SCENIC in
1021 liver could not be defined as the number of positives was zero. **D)** A plot of run times (y-axis,
1022 log-scale) for 5 methods (SCIRA, MAGIC, Scrabble, GENIE3/SCENIC and VIPER-D)
1023 against the number of single-cells profiled (x-axis, log-scale). Filled symbols represent times
1024 estimated from actual runs, unfilled symbols are imputed estimates obtained by extrapolation
1025 of fitted linear functions (on a log-scale). Run times were estimated using 4 processing cores
1026 (SCIRA, MAGIC, GENIE3/SCENIC, VIPER-D) and 1 core for Scrabble (as Scrabble offers
1027 no option for parallelization).

1028

1029 **Figure-3: SCIRA predicts inactivation of lung-specific TFs in lung tumor epithelial cells.**
1030 **A)** t-SNE scatterplot of approximately 52,000 single cells from 5 lung cancer patients, with a
1031 common non-malignant alveolar and (tumor) epithelial clusters highlighted in blue and red,
1032 respectively. **B)** Corresponding t-SNE scatterplot with cells colored-labeled by expression of
1033 an alveolar marker *CLDN18*. **C)** As B), but with cells colored according to expression of the
1034 epithelial marker *EPCAM*. Right panel depicts boxplots of the $\log_2(\text{counts per million} + 1)$ of
1035 *EPCAM* for cells in the non-malignant alveolar cluster, the tumor epithelial clusters and all
1036 other cell clusters combined (T-cells, B-cells, endothelial, myeloid and fibroblast cells). In

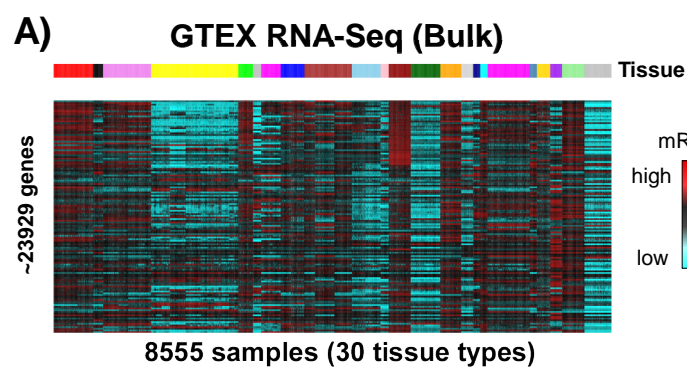
1037 boxplot, horizontal lines describe median, interquartile range and whiskers extend to
1038 $1.5 \times$ inter-quartile range. **D)** Barplot displaying the number of TFs (y-axis) passing a
1039 Bonferroni adjusted < 0.05 threshold and exhibiting decreased (DN) or increased activity (UP)
1040 in tumor epithelial cells (SCIRA & VIPER-D) indicated in darkgreen and darkred,
1041 respectively, and correspondingly the same numbers for differential expression (DE).
1042 P-values are from a Binomial-test, to test if there is a skew towards
1043 inactivation/downregulation in cancer. Right panel depicts the Monte-Carlo (n=1000 runs)
1044 significance analysis with grey curve denoting the null distribution for the fraction of TFs
1045 exhibiting significant inactivation in tumor epithelial cells, and darkgreen line labeling the
1046 observed fraction (0.92=35/38). Empirical P-value derived from the 1000 Monte-Carlo runs
1047 is given. **E)** Scatterplot as in A), but now with cells color-labeled according to activation of
1048 *FOXA2* as estimated using SCIRA. Beanplots of the predicted SCIRA activity level of
1049 *FOXA2* between normal alveolar, tumor epithelial and all other cells. P-value is from a t-test
1050 between normal alveolar and tumor epithelial cell clusters. **F)** Pattern of differential activity
1051 (SCIRA & VIPER-D) and differential expression for the 38 lung-specific TFs. Darkgreen
1052 denotes significant inactivation or underexpression in tumor epithelial cells compared to
1053 normal alveolar, brown denotes significant activation or expression. Grey=no-change (NC)
1054 and white indicates missing regulon information (VIPER-D). **G)** Pattern of differential
1055 expression for the same 38 lung-specific TFs in the bulk RNA-Seq lung cancer datasets
1056 (LUSC=lung squamous cell carcinoma, LUAD=lung adenoma carcinoma). **H)** Barplot
1057 displaying the number of lung-specific TFs displaying significant and directionally consistent
1058 changes in both single-cell and bulk RNA-Seq datasets. In the single-cell data we use
1059 differential activity for SCIRA and VIPER-D, whereas for DE we use differential expression.
1060

1061 **Figure-4: Inactivation of colon-specific TFs in colorectal cancers at single-cell resolution**
1062 **A)** Heatmaps of TF-activity (left panel) and TF expression (right panel), with cells ordered by
1063 hierarchical clustering over the 56 colon-specific TFs. TFs undergoing significant
1064 inactivation/underexpression in cancer cells are labeled in blue, whilst those undergoing
1065 activation/overexpression are labeled in darkred. **B)** Heatmap of differential TF-activity
1066 (SCIRA & VIPER-D) and TF-expression (DE) between cancer and normal cells, with colors
1067 indicating statistical significance (Bonferroni $P < 0.05$) and directionality of change:
1068 blue=significant inactivation/underexpression in cancer, brown=significant
1069 activation/overexpression in cancer, grey=no-change. Barplots compare the number of
1070 inactivated/underexpressed (blue) TFs to the number that are activated/overexpressed
1071 (brown). P-values derive from a one-tailed Binomial test to assess significance of skew. **C)**
1072 Boxplots displaying TF-activity and TF-expression between normal epithelial and cancer
1073 cells for two representative TFs where there is substantial discordance between differential
1074 activity and differential expression. P-values for differential TF-activity and TF-expression
1075 derive from a t-test and a Wilcoxon rank sum test, respectively. **D)** Heatmaps of TF-activity

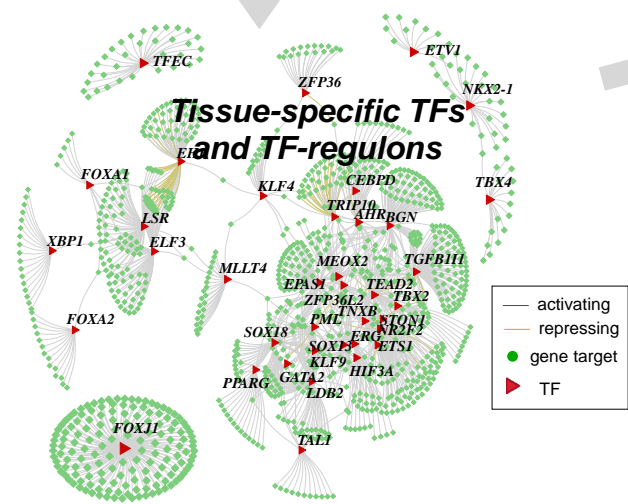
1076 for the normal and cancer cells from each of 3 patients, and displaying only the subset of the
1077 56 colon-TFs which exhibit significant inactivity in the cancer cells (Bonferroni $P < 0.05$).

1078

1079 **Figure-5: SCIRA reveals smoking-associated tumor suppressor events. A)** tSNE diagrams
1080 of normal lung epithelial cells obtained by application to the SCIRA-derived regulatory
1081 activity estimates for the 38 lung-specific TFs. Left panel depicts the two main clusters
1082 inferred using DBSCAN, whilst right panels depict the TF-activity levels for 4 of the
1083 lung-specific TFs. **B)** As A), but now displaying the mRNA expression levels of 4 markers,
1084 one for each of ciliated cells (*FOXJ1*), goblet cells (*MUC5AC*), club cells (*SCBG1A1*) and
1085 basal cells (*KRT5*). **C)** As B), but now for 5 lung-specific TFs. **D)** Left panel: As A), but now
1086 with cells color-labeled according to whether they derived from a smoker or non-smoker.
1087 Right panel: PCA scatterplot (PC1 vs PC2) obtained from a PCA on all non-ciliated cells,
1088 plus associated density plots along PC1 for cells stratified according to smoking status.
1089 P-value is from a two-tailed Wilcoxon rank sum test. **E)** Hierarchical clustering heatmap over
1090 12 lung-specific TFs exhibiting significant (Bonferroni adjusted $P < 0.05$) activity changes
1091 according to smoking-status. Color bar to the right indicates whether TF is more or less active
1092 in cells exposed to smoking. **F)** Color bar indicating the pattern of differential regulatory
1093 activity for the same 12 TFs in lung cancer cells. **G)** Density distribution of *EHF* activity (left)
1094 and *EHF* expression (right) for cells expressing *MUC5AC* (*MUC5AC+*), a goblet cell marker,
1095 and cells not expressing *MUC5AC* (*MUC5AC-*). P-values derive from a two tailed Wilcoxon
1096 rank sum test.



Derive tissue-specific (e.g lung) regulatory network, adjusting for stromal heterogeneity



C) scRNA-Seq dataset (e.g. lung development)

