

# 1 **A statistical framework for the robust detection of hidden** 2 **variation in single cell transcriptomes**

3  
4 Donghyung Lee<sup>1,\*</sup>, Anthony Cheng<sup>1,2</sup>, Mohan Bolisetty<sup>1</sup>, and Duygu Ucar<sup>1,3,\*</sup>

5  
6 <sup>1</sup> The Jackson Laboratory for Genomic Medicine, Farmington, CT, USA.

7 <sup>2</sup> Department of Genetics and Genome Sciences, University of Connecticut Health Center,  
8 Farmington, CT, USA.

9 <sup>3</sup> Institute of Systems Genomics, University of Connecticut Health Center, Farmington, CT,  
10 USA.

11 \* Correspondence: [donghyung.lee@jax.org](mailto:donghyung.lee@jax.org) and [duygu.ucar@jax.org](mailto:duygu.ucar@jax.org)

## 12 13 **Abstract**

14 Single cell RNA-sequencing (scRNA-seq) precisely characterize gene expression levels and dissect variation in  
15 expression associated with the state (technical or biological) and the type of the cell, which is averaged out in  
16 bulk measurements. Multiple and correlated sources contribute to gene expression variation in single cells,  
17 which makes their estimation difficult with the existing methods developed for bulk measurements (e.g.,  
18 surrogate variable analysis (SVA)) that estimate orthogonal transformations of these sources. We developed  
19 iteratively adjusted surrogate variable analysis (IA-SVA) that can estimate hidden and correlated sources of  
20 variation by identifying a set of genes affected with each hidden factor in an iterative manner. Analysis of  
21 scRNA-seq data from human cells showed that IA-SVA could accurately capture hidden variation arising from  
22 technical (e.g., stacked doublet cells) or biological sources (e.g., cell type or cell-cycle stage). Furthermore, IA-  
23 SVA delivers a set of genes associated with the detected hidden source to be used in downstream data analyses.  
24 As a proof of concept, IA-SVA recapitulated known marker genes for islet cell subsets (e.g., alpha, beta), which  
25 improved the grouping of subsets into distinct clusters. Taken together, IA-SVA is an effective and novel  
26 method to dissect multiple and correlated sources of variation in scRNA-seq data.

## 27 **Introduction**

28 Single-cell RNA-Sequencing (scRNA-seq) enable precise characterization of gene expression  
29 levels, which harbour variation in expression associated with both technical (e.g., biases in  
30 capturing transcripts from single cells, PCR amplifications or cell contamination) and  
31 biological sources (e.g., differences in cell cycle stage or cell types). If these sources are not  
32 accurately identified and properly accounted for, they might confound the downstream  
33 analyses and hence the biological conclusions<sup>1-3</sup>. In bulk measurements, hidden sources of  
34 variation are typically unwanted (e.g., batch effects) and are computationally eliminated from

35 the data. However, in single cell RNA-seq data, variation/heterogeneity stemming from  
36 hidden biological sources can be the primary interest of the study; which necessitate their  
37 accurate detection (i.e., existence of a hidden factor) and estimation (i.e., contribution of this  
38 factor to the gene expression levels) for downstream data analyses and interpretation. For  
39 example, a recent such study uncovered a CD1C+ dendritic cell (DC) subset by profiling  
40 human blood samples<sup>4</sup> and improved the immune monitoring of human DCs in health and  
41 disease. One challenge in detecting hidden sources of variation in scRNA-seq data lies in the  
42 existence of multiple and highly correlated hidden sources, including geometric library size  
43 (i.e., the library size of log-transformed read counts), number of expressed genes in a cell,  
44 experimental batch effects, cell cycle stage and cell type<sup>5-8</sup>. The correlated nature of hidden  
45 sources limits the efficacy of existing algorithms to accurately detect the source and estimate  
46 its contribution to the variation in the data.

47 ‘Surrogate variable analysis’ (SVA)<sup>9-11</sup> is a family of algorithms that are developed to  
48 detect and remove hidden and “unwanted” variation (e.g., batch effect) in gene expression  
49 data by accurately parsing the data into signal and noise. A number of SVA-based methods  
50 have been developed and used for the analyses of microarray, bulk, and single-cell RNA-seq  
51 data including SSVA<sup>11</sup> (supervised surrogate variable analysis), USVA<sup>10</sup> (unsupervised SVA),  
52 ISVA<sup>12</sup> (Independent SVA), RUV (removing unwanted variation)<sup>13,14</sup>, and most recently  
53 scLVM<sup>6</sup> (single-cell latent variable model). These methods primarily aim to remove  
54 ‘unwanted’ variation (e.g., batch or cell-cycle effect) in data while preserving the biological  
55 signal of interest typically to improve downstream differential expression analyses between  
56 cases and controls. For this purpose, they utilize PCA (principal component analysis), SVD  
57 (singular vector decomposition) or ICA (independent component analysis) to infer orthogonal  
58 transformations of hidden factors that can be used as covariates in downstream analysis.  
59 However, this paradigm by definition results in orthogonality between multiple estimated  
60 factors and limits the efficacy of existing SVA-based methods for single-cell data analyses, in  
61 which some of the sources of variation are ‘wanted’ and are highly correlated with each other.

62 To fill this gap, we developed a robust and iterative SVA-based statistical framework:  
63 **Iteratively Adjusted Surrogate Variable Analysis (IA-SVA)** (Figure 1A and Methods for  
64 details), which provides three major advantages. First, it accurately estimates multiple hidden  
65 sources of variation even if the sources are correlated with each other and with known  
66 sources, which is a limitation of existing SVA-based methods. Second, it enables assessing  
67 the significance of each detected factor for explaining the unmodeled variation in the data.  
68 Third, it delivers a set of genes that are significantly associated with the detected hidden

69 source. Application of IA-SVA for scRNA-seq data analyses is diverse including the  
70 detection of “unwanted” variation due to cell contamination or “wanted” variation associated  
71 with rare cell types (Figure 1B). In simulation studies we showed that IA-SVA i) provides  
72 high statistical power in detecting hidden factors; ii) controls Type I error rate at the nominal  
73 level ( $\alpha = 0.05$ ); iii) delivers high accuracy in estimating hidden factors. We evaluated the  
74 efficacy of IA-SVA on scRNA-seq data from human pancreatic islets and brain cells and  
75 showed that IA-SVA is effective in capturing heterogeneity associated with both technical  
76 (e.g., doublet cells) and biological sources (e.g., differences in cell types or cell-cycle stages).  
77 Furthermore, we showed that IA-SVA based gene selection can be further utilized in  
78 downstream analyses such as in data visualization using t-distributed stochastic neighbor  
79 embedding (tSNE)<sup>15</sup> and performs favourably compared to existing methods developed for  
80 gene selection and visualization (e.g., Spectral tSNE<sup>16</sup>).

## 81 **Results**

### 82 **Benchmarking IA-SVA on simulated data.**

83 To assess and compare the detection power, Type I error rate, and the accuracy of hidden  
84 source estimates using IA-SVA and existing state-of-the-art methods (i.e., USVA and SSVA),  
85 we performed simulation studies (see Methods for details) under the null hypothesis (i.e., a  
86 group (case/control) variable affecting 10% of genes and no hidden factor) and under the  
87 alternative hypothesis (i.e., a group variable and three hidden factors affecting 30%, 20%, 10%  
88 of genes, respectively). Under the alternative hypothesis, we considered two correlation  
89 scenarios where the three hidden factors are moderately ( $|r| \sim 0.3-0.6$ ) or weakly ( $|r| < 0.3$ )  
90 correlated with the group variable (i.e., a known factor). Under each simulation scenario, we  
91 generated 1,000 scRNA-seq data sets (10,000 genes and 50 cells) and performed IA-SVA,  
92 USVA and SSVA ( $\alpha = 0.05$ , 50 permutations) on them to detect simulated hidden factors.  
93 Using these simulation results, we assessed the empirical Type I error rate of each method  
94 (i.e., the number of times each method detects a false positive factor under the null hypothesis  
95 at the nominal level of 0.05 divided by the number of simulations ( $n=1,000$ )). Similarly, we  
96 also quantified the empirical detection power rate of each method under different alternative  
97 hypothesis scenarios as the number of times each method detects a simulated factor under the  
98 alternative hypothesis (i.e., a factor actually exists and is detected as significant by the  
99 method) divided by the number of simulations. We used the average of the absolute Pearson

100 correlation coefficients between the simulated and estimated hidden factors to quantify the  
101 accuracy of estimates.

102 Simulation studies showed that IA-SVA performs as well or better than USVA and  
103 SSVA in terms of detection power and accuracy of the estimate while controlling the Type I  
104 error rate (0.04 for IA-SVA versus 0.09 for USVA and SSVA) (Table 1). In particular, IA-  
105 SVA was more effective when a hidden factor affected a small percentage of genes and when  
106 the factors were correlated ( $|r|=0.3-0.6$ ) with the known factor (i.e., group variable). For  
107 example IA-SVA detected Factor3, which affected only 10% genes, 87% of the time,  
108 whereas USVA and SSVA detected this factor 78% of the simulations (first three columns in  
109 Table 1). More importantly, IA-SVA correctly inferred the correlations among multiple  
110 hidden factors while USVA and SSVA delivered biased estimates due to their orthogonality  
111 assumption (Supplementary Figure S1).

### 112 **IA-SVA captures variation stemming from a small number of alpha cells.**

113 To test whether IA-SVA is effective in capturing variation within a homogenous cell  
114 population, we analysed scRNA-seq data generated from human alpha cells ( $n=101$ , marked  
115 with glucagon (*GCG*) expression) obtained from three diabetic patients<sup>17</sup> using the Fluidigm  
116 C1 platform<sup>18</sup>, for which the original study did not report any separation of these alpha cells .  
117 Using geometric library size and patient ID as known factors, significant surrogate variables  
118 (SVs) were inferred using IA-SVA ( $\alpha=0.05$ , 50 permutations) on the data (14,416 genes and  
119 101 cells). For comparison, we applied PCA, USVA, and tSNE on this data. In USVA  
120 analysis, similarly geometric library size and patient ID were used as known factors and  
121 significant SVs were obtained ( $\alpha=0.05$ , 50 permutations). In the PCA analysis, PC1 was  
122 discarded since it is highly correlated ( $r=0.99$ ) with the geometric library size.

123 Top two significant SVs inferred by IA-SVA clearly separated alpha cells into two  
124 groups (six outlier cells marked in red vs. the rest marked in grey at  $SV2 > 0.1$ ) (Figure 2A).  
125 27 genes significantly associated with second SV (SV2) (Benjamini-Hochberg q-value (FDR)  
126  $< 0.05$ , coefficient of determination ( $R^2$ )  $> 0.6$ ), which included genes expressed in fibroblasts  
127 such as *COL4A1* and *COL4A2*. These genes were exclusively expressed in six outlier cells  
128 and clearly separated alpha cells into two clusters (Figure 2B). A larger set of SV2-associated  
129 genes ( $n = 108$ ,  $FDR < 0.05$ ,  $R^2 > 0.3$ ) was used for pathway and GO enrichment analyses  
130 and uncovered that these genes are associated with extracellular matrix receptors  
131 (Supplementary Table S1). Hence, these outlier cells likely arise from cell contamination  
132 (e.g., fibroblasts contaminating islet cells) or cell doublets (e.g., two cells captured together)

133 — a known problem in early Fluidigm C1 experiments<sup>20,21</sup>. Alternative methods (i.e., PCA,  
134 USVA, tSNE) failed to clearly detect these outlier cells (Figure 2C-E).

135 We next studied whether this source of heterogeneity can be recapitulated in an  
136 independent and bigger human islet scRNA-seq dataset<sup>18</sup>, using gene expression profiles  
137 (17,168 genes) of 569 alpha cells from six diabetic patients. Using geometric library size and  
138 patient ID as known factors we identified top 2 significant SVs using IA-SVA and USVA.  
139 For comparison, we also conducted PCA and tSNE analyses on this data. In PCA, PC1 was  
140 discarded since it matched the geometric library size, which is adjusted for in IA-SVA and  
141 USVA analyses. IA-SVA's SV2 separated alpha cells into two groups (Supplementary  
142 Figure S2A) and as in the previous case it was associated with fibrotic response genes  
143 including *SPARC*, *COL4A1*, *COL4A2* (n=81, FDR < 0.05 and  $R^2 > 0.3$ ) (Supplementary  
144 Figure S2B, GO/pathway results in Supplementary Table S2). These results highlight IA-  
145 SVA's ability to detect variation among alpha cells potentially due to cell contamination or  
146 cell doublets. PCA, USVA, and tSNE failed to clearly separate these compromised alpha  
147 cells (Supplementary Figure S2C-E) from the rest of the cells.

#### 148 **IA-SVA accurately detects variation arising from cell-cycle stage differences.**

149 Differences in cell-cycle stages lead to variation in single cell gene expression data<sup>3</sup>.  
150 Supervised methods based on SVA have been developed to detect and correct for cell cycle  
151 stage differences, most notably the scLVM algorithm. scLVM implements a Bayesian latent  
152 variable model to infer hidden cell-cycle factors by using known cell cycle genes<sup>6</sup>. IA-SVA  
153 can provide an unsupervised alternative by accurately capturing cell-cycle related variation in  
154 single cell data. To show this, we analyzed scRNA-seq data (21,907 genes and 74 cells)  
155 obtained from human glioblastomas that has an established cell-cycle signature<sup>22</sup>. We  
156 conducted IA-SVA analyses by using geometric library size as a known factor and extracted  
157 top 2 significant SVs ( $\alpha=0.05$ , 50 permutations). For comparison, we applied PCA, USVA  
158 and tSNE analyses on this data, where for USVA geometric library size is used as a known  
159 factor.

160 IA-SVA's SV1 clearly separated 13 cells from the rest (cells marked in red in Figure  
161 3A), which was associated with 119 genes (FDR < 0.05 and  $R^2 > 0.3$ ). Hierarchical clustering  
162 (ward.D2, cutree\_cols=2) using these genes confirmed the separation of cells into two groups  
163 (Figure 3B), whereas alternative methods failed to clearly separate these two groups of cells  
164 (Figure 3C-E). Pathway and GO enrichment analyses of these genes<sup>23,24</sup> revealed significant  
165 enrichment for cell-cycle process related GO terms and pathways (Supplementary Table S3),

166 suggesting that this hidden variation is stemming from cell-cycle stage differences. Indeed,  
167 cell-cycle-stage predictions of cells using the SCRAN R package<sup>25</sup> showed that cells in  
168 different cell-cycle stages have different SV1 values (Figure 3F). We noted that SV1 is  
169 highly correlated ( $|r|= 0.44$ ) with the geometric library size (typically the top contributor to  
170 the variation in single cell data), which might explain why alternative methods failed to  
171 clearly detect this variation in the data. These results demonstrate that IA-SVA can  
172 effectively detect variation stemming from cell-cycle differences in an unsupervised manner  
173 from single cell transcriptomes, even if this factor is highly correlated with known factors.

#### 174 **IA-SVA based gene selection improves single cell data visualization.**

175 tSNE and other dimension reduction algorithms (e.g., Spectral tSNE implemented in Seurat<sup>16</sup>)  
176 are frequently used to visualize single cell data since they group together cells with similar  
177 gene expression patterns. However, variation introduced by technical or biological factors  
178 can confound the signal of interest and generate spurious clustering of data. IA-SVA can be  
179 particularly effective in handling this problem by estimating hidden factors of interest  
180 accurately while adjusting for all known factors of no interest. Moreover, IA-SVA identifies  
181 genes associated with each detected hidden factor, which could be biologically relevant such  
182 as marker genes for different cell types. The genes inferred by IA-SVA can significantly  
183 improve the performance of data visualization methods (e.g., tSNE<sup>15</sup>). To illustrate this, we  
184 studied single cell gene expression profiles (16,005 genes) of alpha (n=101, marked with  
185 glucagon (*GCG*) expression), beta (n=96, marked with insulin (*INS*) expression), and ductal  
186 (n=16, marked with *KRT19* expression) cells obtained from three diabetic patients<sup>17</sup>. First, we  
187 applied tSNE on all genes (n=16,005) and color-coded genes based on the reported cell type  
188 assignments<sup>17</sup>, which failed to separate cells from different origins (Figure 4A). Next, we  
189 applied IA-SVA on this data using patient ID, batch ID and the number of expressed genes as  
190 known factors and obtained significant SVs. SV1 and SV2 separated cells into distinct  
191 clusters (Supplementary Figure S3), suggesting that these SVs might be associated with cell  
192 type differences. Indeed, genes associated with SV1 and SV2 (n=92, FDR < 0.05 and  $R^2 >$   
193 0.5) included known marker genes used in the original study (*INS*, *GCG*, *KRT19*) and  
194 uncovered alternative marker genes associated with alpha, beta and ductal cells (Figure 4B).  
195 These genes were annotated with diabetes and insulin processing related GO terms and  
196 pathways (Supplementary Table S4). As expected, tSNE analyses based on these 92 genes  
197 improved data visualization and clearly grouped together cells with respect to their cell type  
198 assignments (Figure 4C). Such improved analyses can be instrumental in discovering cells

199 that might be incorrectly labelled based on a single marker gene. For example, our analyses  
200 revealed a beta cell that is labelled as a ductal cell in the original study (one green cell  
201 clustered with blue cells in Figure 4C). For comparison, we applied recently developed  
202 visualization methods, CellView<sup>26</sup> and Spectral tSNE<sup>16</sup>, on the same data with their  
203 recommended settings. CellView identified the 1000 most over-dispersed genes and  
204 conducted tSNE on these genes. Spectral tSNE detected 2,933 most over-dispersed genes and  
205 performed tSNE on significant principal components of these genes. On this small dataset,  
206 both methods managed to group cells of different types into distinct groups (Supplementary  
207 Figure S4), suggesting that existing methods for gene selection and visualization are effective  
208 when datasets are small in size and are not confounded with multiple factors.

209 To test the efficacy of these methods on a bigger and more complex dataset, we  
210 conducted similar analyses on scRNA-seq data (19,226 genes) of 1,600 islet cells including  
211 alpha (n=946), beta (n=503), delta (n=58), and PP (n=93) cells from 6 diabetic and 12 non-  
212 diabetic individuals, where the study includes multiple confounding factors (e.g., ethnicity,  
213 disease state)<sup>18</sup>. We noted that original cell type assignments significantly correlate with  
214 patient identifications ( $C=0.48$ ,  $C$ =Pearson's contingency coefficient) and with ethnicity  
215 ( $C=0.25$ ), which would reduce the ability of existing methods to detect variation associated  
216 with cell types. In such complex datasets, failing to properly adjust for potential confounding  
217 factors prior to data analyses can lead to spurious grouping of cells, which might mislead the  
218 biological conclusions. Indeed, when these cells were visualized using tSNE using all genes  
219 (n=19,226) and were color-coded with respect to the original cell-type assignments<sup>18</sup>, cell  
220 types did not separate from each other and spurious clusters were observed within each cell  
221 type (Figure 4D). As suspected, potential confounding factors (i.e., patient ID and ethnicity)  
222 explained this grouping of cells (Supplementary Figure S5), which might be misleading as  
223 researchers are looking for alpha and beta cell subtypes that can be related to Type 2 Diabetes  
224 pathogenesis<sup>27</sup>. To eliminate spurious clusters stemming from known factors, existing  
225 methods (e.g., Seurat) simply regress out all known factors prior to visualization. However,  
226 this might affect the signal of interest (i.e., cell type assignment), due to high correlation  
227 between known factors (i.e., patient ID) and the hidden factor (i.e., cell types).

228 We applied IA-SVA on this complex data, while accounting for known factors (i.e.,  
229 the number of expressed genes and patient ID) and extracted top four significant SVs  
230 (Supplementary Figure S6A and B). We identified 57 genes associated with the most  
231 significant SV (SV1) ( $FDR < 0.05$  and  $R^2 > 0.5$ ), which included known marker genes (i.e.,  
232 *INS* and *GCG*) (Supplementary Figure S7, Supplementary Table S5) and revealed novel

233 marker genes for these cell types. tSNE analyses using these 57 IA-SVA detected genes  
234 clearly separated different cell types into discrete groups and reinforced the importance of  
235 properly adjusting for known factors prior to data analyses (Figure 4E). For comparison, we  
236 applied CellView and Spectral tSNE on this data with recommended settings; however they  
237 failed to accurately group cells into distinct cell types (Figure 4F and G). Similar analyses  
238 were conducted using PCA and USVA on the same data, where top surrogate factors  
239 obtained with both methods failed to separate different cell types into distinct groups  
240 (Supplementary Figure S6C and D). Combined together these analyses suggest that IA-SVA  
241 is particularly effective in the analyses of complex datasets, which include the measurements  
242 of many cells that are affected by diverse confounding factors.

## 243 **Discussion**

244 Surrogate variable analyses based methods are effective in detecting and eliminating hidden  
245 and unwanted variation in bulk gene expression data (such as batch effects). By using  
246 dimensionality reduction algorithms (e.g., PCA or SVD), these methods infer linear  
247 transformations of hidden factors and utilize these factor estimates as additional covariates in  
248 downstream analyses to eliminate unwanted variation<sup>14</sup>. However, measurement of gene  
249 expression levels at single cell resolution pose novel challenges in the detection and  
250 adjustment of hidden sources of data variation. First, single cell transcriptomes harbour  
251 hidden variation that can be biologically interesting (hence ‘wanted’) and can be the major  
252 goal of the study, for example detection of rare cells within a tissue<sup>28</sup> or detection of a cell’s  
253 subtypes that can be linked to health or disease<sup>27</sup>. Second, since single cell data do not  
254 average out variation as in the case of bulk profiling, the data reflect variation arising from  
255 diverse biological and technical sources some of which could be highly correlated. Existing  
256 SVA-based methods do not readily apply to the unique needs of single cell data analyses. To  
257 fill this gap, we developed IA-SVA, where the objective is the accurate estimation of hidden  
258 factors even if these factors are correlated with each other or with the known factors. Unlike  
259 other SVA-based methods, IA-SVA focuses more on the accurate detection and estimation of  
260 hidden factors rather than their elimination since these factors can be biologically interesting,  
261 e.g., identification of a new cell type and its marker genes. Indeed, analyses on simulated  
262 scRNA-seq data showed that IA-SVA outperforms existing supervised (i.e., SSVA) and  
263 unsupervised (i.e., USVA) state-of-the-art methods in the estimation of hidden factors (not  
264 necessarily in their elimination). Furthermore, we noted that IA-SVA is particularly effective



265 (i.e., high detection power and accuracy, and Type I error rate controlled under the nominal  
266 level of 0.05) in detecting correlated factors that affect a small fraction of genes. Therefore  
267 IA-SVA is an effective unsupervised alternative to existing SVA-based algorithms when the  
268 goal is to accurately estimate hidden factors (and their marker genes) rather than to eliminate  
269 these factors.

270 Through analyses of diverse human datasets from multiple studies, we established that  
271 IA-SVA can effectively detect hidden heterogeneity in scRNA-seq data arising from a small  
272 number of cells either due to technical (i.e., contamination or doublets) or biological (i.e., a  
273 rare cell type) sources. In two independently generated islet scRNA-seq datasets, we showed  
274 that IA-SVA detects heterogeneity stemming from compromised alpha cells (contaminated or  
275 stacked), which should be excluded from the downstream analyses (Figure 2 and  
276 Supplementary Figure S2). Therefore, IA-SVA provides an easy-to-apply statistical  
277 framework to uncover variation in scRNA-seq data even if it is stemming from only a  
278 handful of cells. This ability of IA-SVA can be effective in identifying rare cells within a  
279 population of cells, where genes associated with the detected factor can uncover relevant  
280 marker genes for the rare population of cells. In addition, IA-SVA can be effective in  
281 detecting heterogeneity associated with cell-cycle stages without prior knowledge, therefore  
282 providing an unsupervised solution to this common problem in single cell data analyses  
283 (Figure 3).

284 An important feature of IA-SVA is its ability to uncover genes associated with  
285 detected hidden factors. This feature can be used to detect marker genes associated with  
286 different cell types. As a proof-of-concept we demonstrated this in pancreatic islet cells,  
287 where we captured known marker genes (e.g., *INS*, *GCG*) in an unsupervised manner.  
288 Moreover, genes captured by IA-SVA can be used to improve the visualization of single cells  
289 into their respective clusters, as demonstrated with the analyses of islet cells from two  
290 separate studies (Figure 4). Spectral tSNE<sup>16</sup> is a commonly used method for scRNA-seq data  
291 visualization especially in the existence of confounding factors. This method regresses out  
292 variation associated with known factors before data visualization. However, when a hidden  
293 factor is ‘wanted’ (e.g., cell types) and is highly correlated with known factors, removing the  
294 known factors will also diminish the ability to detect the wanted hidden factor and the genes  
295 associated with this factor (e.g., marker genes for different cell types). Indeed, our analyses  
296 using islet cells emphasized the importance of properly adjusting the data for known factors  
297 prior to further analyses, such as data visualization (e.g., tSNE) to prevent spurious clustering

298 of cells due to the confounding factors (Figure 4E). IA-SVA is an alternative method that can  
299 effectively handle data with multiple confounding factors.

300 In summary, IA-SVA is an SVA-based unsupervised method designed to accurately  
301 estimate hidden factors (sources of variation) in single cell gene expression data while  
302 adjusting for known factors. The iterative and flexible framework of IA-SVA allows the  
303 accurate estimation of multiple and potentially correlated factors along with their statistical  
304 significance, which is the main advantage of IA-SVA over existing methods. This flexibility  
305 is more realistic given the confounded nature of known and unknown factors in single cell  
306 gene expression measurements. Therefore, IA-SVA has an improved performance over  
307 existing SVA-based methods in terms of estimating hidden sources of variation when they  
308 are correlated with each other and with known variables. IA-SVA is an effective alternative  
309 to methods developed for single cell data analyses (e.g., CellView and Seurat), especially for  
310 the analyses of complex data (i.e., data with multiple confounding and correlated factors).  
311 With the increasing amount of single cell studies and the increasing complexity of human  
312 cohorts, IA-SVA will serve as an effective statistical framework specifically designed to  
313 handle unique challenges of scRNA-seq data analyses.

## 314 **Methods**

### 315 **IA-SVA framework.**

316 We model the log-transformed sequencing read counts for  $m$  cells and  $n$  genes (i.e.,  $Y_{m \times n}$ ) as  
317 a combination of known and unknown variables as follows:

318

$$319 \quad Y_{m \times n} = X_{m \times p} \beta_{p \times n} + Z_{m \times k} \delta_{k \times n} + \varepsilon_{m \times n},$$

320

321 where  $X_{m \times p}$  is a matrix for  $p$  known variable(s) (e.g., group assignment for cases and  
322 controls, sex or ethnicity),  $Z_{m \times k}$  is a matrix for  $k$  unknown variables and  $\varepsilon_{m \times n}$  is the error  
323 term. With this model, we can account for any clinical/experimental information about  
324 samples (e.g., sex, ethnicity, age, BMI or batch) as known factors ( $X_{m \times p}$ ) and dissect  
325 unaccounted variation in the read count data that is attributable to hidden factors ( $Z_{m \times k}$ ).

326 Existing unsupervised SVA-based methods (e.g., USVA<sup>10</sup>, RUV<sup>13</sup>, ISVA<sup>12</sup>) obtain the  
327 residual matrix ( $Y'_{m \times n}$ ) by regressing read counts ( $Y_{m \times n}$ ) on all known factors ( $X_{m \times p}$ ). Then,  
328 they infer the hidden factors from this residual matrix ( $Y'_{m \times n}$ ) using dimensionality reduction  
329 algorithms (e.g., PCA, SVD or ICA). Thus, by definition, multiple hidden factors captured by

330 these methods are orthogonal to each other and to known variables. Therefore, if hidden  
331 factors are correlated with each other and with known factors, the direct inference from the  
332 residual matrix leads to biased estimates of hidden factors due to the orthogonality  
333 assumption.

334 In contrast, IA-SVA does not impose orthogonality between factors (hidden or known)  
335 and allows an unbiased estimation of correlated factors via a novel iterative framework  
336 (Figure 1). At each iteration, IA-SVA first obtains residual matrix ( $Y'_{m \times n}$ ), i.e., read counts  
337 adjusted for all known factors ( $X_{m \times p}$ ) including surrogate variables of unknown factors  
338 estimated from previous iterations and extracts the first principal component (PC1) from the  
339 residuals ( $Y'_{m \times n}$ ) using SVD. Next it tests the significance of PC1 in terms of its contribution  
340 to the unmodeled variation (i.e., residual variance). Using this PC1 as a surrogate variable (as  
341 in the case of existing methods) implicitly imposes orthogonality between known and hidden  
342 factors. Instead, IA-SVA uses PC1 to infer gene weights, which are also used to infer genes  
343 associated with the hidden factor. IA-SVA relies on the fact that the first principal component  
344 (PC1) of the residual matrix is highly correlated with the hidden factor that contributes the  
345 most to the unmodeled variation in data, and thus, PC1 can be used to sort genes in terms of  
346 their relative association strength with the hidden factor. To infer these genes, IA-SVA  
347 regresses  $Y$  on PC1 and calculates the coefficient of determination ( $R^2$ ) for each gene. Genes  
348 with high  $R^2$  scores can be treated as marker genes for the factor. These  $R^2$  scores are further  
349 utilized for an unbiased inference of the hidden factor while retaining the correlation structure  
350 between known and hidden factors. For this, IA-SVA first obtains a weighted read count  
351 matrix ( $Y''_{m \times n}$ ) by weighing all genes with respect to their  $R^2$  scores (i.e.,  $Y''_{m \times n} = Y_{m \times n} W_{n \times n}$ ,  
352 where  $W$  is a diagonal matrix of  $R^2$  values). Then it conducts a SVD on  $Y''_{m \times n}$  and obtains the  
353 PC1 to be used as a surrogate variable (SV) for the hidden factor. In the next iteration, IA-  
354 SVA uses this SV as an additional known factor to identify further significant hidden factors.  
355 The iterative procedure of IA-SVA is composed of six major steps as summarized in Figure  
356 1A and below:

357 **[Step 1]** Regress  $Y_{m \times n}$  on all known factors ( $X_{m \times p}$ ), including SVs obtained from previous  
358 iterations, to obtain residuals ( $Y'_{m \times n}$ ).

359 **[Step 2]** Conduct a SVD on the obtained residuals ( $Y'_{m \times n}$ ) to extract the first PC (PC1).

360 **[Step 3]** Test the significance of the contribution of PC1 to the variation in residuals ( $Y'_{m \times n}$ )  
361 using a non-parametric permutation-based assessment<sup>9,10,29</sup> as explained further in the next  
362 section.

363 **[Step 4]** If PC1 is significant, regress  $Y_{m \times n}$  (in this case do not use the residual matrix to be  
364 able to capture factors correlated with known factors) on PC1 to compute the coefficient of  
365 determination ( $R^2$ ) for every gene. If PC1 is not significant, stop the iteration and conduct  
366 subsequent down stream analysis using previously obtained significant SVs.

367 **[Step 5]** Weigh each gene in  $Y_{m \times n}$  with respect to its  $R^2$  value by multiplying a gene's read  
368 counts ( $Y_{m \times n}$ ) with its  $R^2$  values ( $Y''_{m \times n} = Y_{m \times n} W_{n \times n}$ ). The highly weighted genes in this  
369 framework serve as the genes affected by the hidden factor.

370 **[Step 6]** Conduct a second SVD on this weighted read counts matrix ( $Y''_{m \times n} = Y_{m \times n} W_{n \times n}$ ) to  
371 obtain PC1, which will be used as the surrogate variable (SV) for the hidden factor.

372 At the end of this six-step procedure, IA-SVA uses the detected SV (if significant) as  
373 an additional known factor in the next iteration. The algorithm stops, when no more  
374 significant PC1s are detected in Step 3. Significant SVs obtained via IA-SVA can be used in  
375 subsequent analyses. If an SV arises from an unwanted factor (e.g., cell contamination), these  
376 SVs can be included as covariates in the model to remove the unwanted variation or to filter  
377 out contaminated cells. In single cell data significant SVs could also explain 'wanted'  
378 biological factors (e.g., different cell types) and genes associated with such SVs can be  
379 further evaluated to discover novel biology from these complex datasets.

### 380 **Assessing the significance of hidden factors.**

381 To test the significance of the contribution of a hidden factor estimate (i.e., PC1 obtained in  
382 Step 2) to the residual variation, we used the permutation based significance test as  
383 previously applied in the surrogate variable analysis<sup>10,29</sup>. Unlike SVA<sup>10</sup>, which tests all  
384 putative hidden factors at once, IA-SVA assesses the significance of hidden factors one at a  
385 time during the corresponding iteration (always for the PC1 detected in that iteration). Briefly,  
386 IA-SVA i) conducts a SVD on the residual matrix obtained from Step 1, ii) computes the  
387 proportion of variation in this matrix explained by the first singular vector (i.e., PC1) and iii)  
388 compares this proportion against the values obtained from permuted residual matrices, as  
389 further explained below:

390 **[Step 1]** Conduct a SVD on the residual matrix ( $Y'_{m \times n}$ ).

391 **[Step 2]** Calculate the proportion of residual variance explained by the first singular vector  
392 (PC1) using the test statistic:  $T_{obs} = \frac{\lambda_1^2}{\sum_k \lambda_k^2}$ , where  $\lambda_k$  is the  $k$ -th singular value.

393 **[Step 3]** Generate a permuted residual matrix by i) permuting each row of the log-  
394 transformed read count matrix  $Y_{m \times n}$  and regressing the permuted read count matrix on all  
395 known factors ( $X_{m \times p}$ ) to obtain fitted residuals.

396 **[Step 4]** Repeat Step 3  $M$  times and generate an empirical null distribution of the test statistic  
397 by calculating  $(T_i^0, i = 1, \dots, M)$  for the  $M$  permuted residual matrices.

398 **[Step 5]** Compute the empirical p-value for the first singular vector (PC1) by counting the  
399 number of times the null statistics  $(T_i^0)$  exceeds the observed one ( $T_{obs}$ ) divided by the  
400 number of permutations ( $M$ ).

#### 401 **scRNA-seq data simulation.**

402 To eliminate the potential bias in data simulations and make simulation studies more  
403 objective<sup>30</sup>, we used a third-party simulation software (Polyester R package<sup>31</sup>) and study  
404 design (<http://jtleek.com/svaseq/simulateData.html>) and simulated scRNA-seq data to test IA-  
405 SVA's performance. The original simulation design is slightly modified to reflect  
406 characteristics of scRNA-seq data for high dropout rate (i.e., excessive number of zeros in the  
407 data) and multiple hidden factors highly correlated with known factors. First, to simulate high  
408 dropout rates (proportion of zero counts = ~70%), we estimated Polyester's zero-inflated  
409 negative binomial model parameters (i.e.,  $p_0$ : probabilities that the count will be zero,  $\mu$ :  
410 mean of the negative binomial,  $size$ : size of the negative binomial) from real-world scRNA-  
411 seq data from human pancreatic islets using the Fluidigm's C1 platform<sup>17</sup>. Using these  
412 estimated model parameters, we simulated expression data for  $m$  cells and  $n$  genes under two  
413 hypotheses: 1) the null hypothesis: no hidden sources of variation, and 2) the alternative  
414 hypothesis: three hidden factors with two values (-1 vs. 1). Under both scenarios, we  
415 simulated a primary variable of interest (i.e., case vs. control) and simulated 10% of genes to  
416 be differentially expressed between the two groups. Under the alternative hypothesis, we  
417 simulated three hidden factors that affect 30%, 20% and 10% of randomly chosen genes  
418 respectively and simulated two different scenarios where these factors are moderately  
419 correlated ( $|r| \sim 0.3-0.6$ ) or weakly correlated ( $|r| < 0.3$ ) with the group variable.

#### 420 **Data processing and normalization.**

421 In all analyses, we filtered out low-expressed genes with read counts  $\leq 5$  in less than three  
422 cells and normalized the retained gene expression counts using SCnorm<sup>19</sup> with default  
423 settings for further analyses. For single cell data visualization examples, we normalized gene

424 read counts by dividing each cell column by its total counts then multiplying median of  
425 library size, which is similar to the default normalization method “LogNormalize”  
426 implemented in Seurat <sup>16</sup>.

#### 427 **Availability of data and methods.**

428 An R package for IA-SVA with example case scenarios is freely available from  
429 <https://github.com/UcarLab/IA-SVA>. The published data sets analyzed in this study including  
430 single-cell RNA sequencing read counts and annotations describing samples and experiment  
431 settings are included in an R data package (iasvaExamples) deposited at  
432 <https://github.com/dleelab/iasvaExamples>.

#### 433 **References**

- 434 1 Tung, P.-Y. *et al.* Batch effects and the effective design of single-cell gene expression studies.  
435 *bioRxiv*, 062919 (2016).
- 436 2 Kowalczyk, M. S. *et al.* Single-cell RNA-seq reveals changes in cell cycle and differentiation  
437 programs upon aging of hematopoietic stem cells. *Genome research* **25**, 1860-1872 (2015).
- 438 3 Stegle, O., Teichmann, S. A. & Marioni, J. C. Computational and analytical challenges in  
439 single-cell transcriptomics. *Nat Rev Genet* **16**, 133-145, doi:10.1038/nrg3833 (2015).
- 440 4 Villani, A. C. *et al.* Single-cell RNA-seq reveals new types of human blood dendritic cells,  
441 monocytes, and progenitors. *Science* **356**, doi:10.1126/science.aah4573 (2017).
- 442 5 Ilicic, T. *et al.* Classification of low quality cells from single-cell RNA-seq data. *Genome Biol* **17**,  
443 29, doi:10.1186/s13059-016-0888-1 (2016).
- 444 6 Buettner, F. *et al.* Computational analysis of cell-to-cell heterogeneity in single-cell RNA-  
445 sequencing data reveals hidden subpopulations of cells. *Nat Biotechnol* **33**, 155-160,  
446 doi:10.1038/nbt.3102 (2015).
- 447 7 McDavid, A., Finak, G. & Gottardo, R. The contribution of cell cycle to heterogeneity in  
448 single-cell RNA-seq data. *Nat Biotechnol* **34**, 591-593, doi:10.1038/nbt.3498 (2016).
- 449 8 Hicks, S. C., Teng, M. & Irizarry, R. A. On the widespread and critical impact of systematic  
450 bias and batch effects in single-cell RNA-Seq data. *bioRxiv*, doi:10.1101/025528 (2015).
- 451 9 Leek, J. T. & Storey, J. D. Capturing heterogeneity in gene expression studies by surrogate  
452 variable analysis. *PLoS Genet* **3**, 1724-1735, doi:10.1371/journal.pgen.0030161 (2007).
- 453 10 Leek, J. T. & Storey, J. D. A general framework for multiple testing dependence. *Proc Natl*  
454 *Acad Sci U S A* **105**, 18718-18723, doi:10.1073/pnas.0808709105 (2008).
- 455 11 Leek, J. T. svaseq: removing batch effects and other unwanted noise from sequencing data.  
456 *Nucleic Acids Res* **42**, doi:10.1093/nar/gku864 (2014).
- 457 12 Teschendorff, A. E., Zhuang, J. & Widschwendter, M. Independent surrogate variable  
458 analysis to deconvolve confounding factors in large-scale microarray profiling studies.  
459 *Bioinformatics* **27**, 1496-1505, doi:10.1093/bioinformatics/btr171 (2011).
- 460 13 Risso, D., Ngai, J., Speed, T. P. & Dudoit, S. Normalization of RNA-seq data using factor  
461 analysis of control genes or samples. *Nat Biotechnol* **32**, 896-902, doi:10.1038/nbt.2931  
462 (2014).
- 463 14 Gagnon-Bartsch, J. A. & Speed, T. P. Using control genes to correct for unwanted variation in  
464 microarray data. *Biostatistics* **13**, 539-552, doi:10.1093/biostatistics/kxr034 (2012).
- 465 15 Maaten, L. V. D. Accelerating t-SNE using tree-based algorithms. *J. Mach. Learn. Res.* **15**,  
466 3221-3245 (2014).

- 467 16 Macosko, E. Z. *et al.* Highly Parallel Genome-wide Expression Profiling of Individual Cells  
468 Using Nanoliter Droplets. *Cell* **161**, 1202-1214, doi:10.1016/j.cell.2015.05.002 (2015).
- 469 17 Lawlor, N. *et al.* Single cell transcriptomes identify human islet cell signatures and reveal cell-  
470 type-specific expression changes in type 2 diabetes. *Genome Res*,  
471 doi:10.1101/gr.212720.116 (2016).
- 472 18 Xin, Y. *et al.* RNA Sequencing of Single Human Islet Cells Reveals Type 2 Diabetes Genes. *Cell*  
473 *Metab* **24**, 608-615, doi:10.1016/j.cmet.2016.08.018 (2016).
- 474 19 Bacher, R. *et al.* SCnorm: robust normalization of single-cell RNA-seq data. *Nat Methods* **14**,  
475 584-586, doi:10.1038/nmeth.4263 (2017).
- 476 20 Xin, Y. *et al.* Use of the Fluidigm C1 platform for RNA sequencing of single mouse pancreatic  
477 islet cells. *Proc Natl Acad Sci U S A* **113**, 3293-3298, doi:10.1073/pnas.1602306113 (2016).
- 478 21 Wang, Y. J. *et al.* Single-Cell Transcriptomics of the Human Endocrine Pancreas. *Diabetes* **65**,  
479 3028-3038, doi:10.2337/db16-0405 (2016).
- 480 22 Patel, A. P. *et al.* Single-cell RNA-seq highlights intratumoral heterogeneity in primary  
481 glioblastoma. *Science* **344**, 1396-1401, doi:10.1126/science.1254257 (2014).
- 482 23 Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference  
483 resource for gene and protein annotation. *Nucleic Acids Res* **44**, D457-462,  
484 doi:10.1093/nar/gkv1070 (2016).
- 485 24 Gene Ontology, C. Gene Ontology Consortium: going forward. *Nucleic Acids Res* **43**, D1049-  
486 1056, doi:10.1093/nar/gku1179 (2015).
- 487 25 Lun, A. T., McCarthy, D. J. & Marioni, J. C. A step-by-step workflow for low-level analysis of  
488 single-cell RNA-seq data with Bioconductor. *F1000Res* **5**, 2122,  
489 doi:10.12688/f1000research.9501.2 (2016).
- 490 26 Bolisetty, M. T., Stitzel, M. L. & Robson, P. CellView: Interactive Exploration Of High  
491 Dimensional Single Cell RNA-Seq Data. *bioRxiv*, doi:10.1101/123810 (2017).
- 492 27 Lawlor, N., Khetan, S., Ucar, D. & Stitzel, M. L. Genomics of Islet (Dys)function and Type 2  
493 Diabetes. *Trends Genet* **33**, 244-255, doi:10.1016/j.tig.2017.01.010 (2017).
- 494 28 Proserpio, V. & Lonnberg, T. Single-cell technologies are revolutionizing the approach to rare  
495 cells. *Immunol Cell Biol* **94**, 225-229, doi:10.1038/icb.2015.106 (2016).
- 496 29 Buja, A. & Eyuboglu, N. Remarks on Parallel Analysis. *Multivariate Behav Res* **27**, 509-540,  
497 doi:10.1207/s15327906mbr2704\_2 (1992).
- 498 30 Gelman, A. & Hennig, C. Beyond subjective and objective in statistics. *Journal of the Royal*  
499 *Statistical Society: Series A (Statistics in Society)*, n/a-n/a, doi:10.1111/rssa.12276 (2017).
- 500 31 Frazee, A. C., Jaffe, A. E., Langmead, B. & Leek, J. T. Polyester: simulating RNA-seq datasets  
501 with differential transcript expression. *Bioinformatics* **31**, 2778-2784,  
502 doi:10.1093/bioinformatics/btv272 (2015).

503

## 504 **Acknowledgements**

505 We thank the Jackson Laboratory Computational Science group, Ucar and Stitzel lab  
506 members for constructive feedback throughout this project. We thank Jane Cha, JAX  
507 scientific illustrator, for her help with Figure 1. This study was made possible by generous  
508 financial support of the National Institute of General Medical Sciences (NIGMS) under  
509 award number GM124922 (to DU) and by the Jackson Laboratory Scientific Services  
510 Innovation Fund (to D.L. and D.U.) Opinions, interpretations, conclusions, and

511 recommendations are solely the responsibility of the authors and do not necessarily represent  
512 the official views of the National Institutes of Health (NIH).

### 513 **Author Contributions**

514 D.L. and D.U. designed the project, generated the figures and wrote the manuscript. D.L.  
515 developed the statistical framework and run the data analyses. M.B. provided advice on data  
516 analyses and interpretation of results. A.C. contributed to the data pre-processing and the  
517 generation of the R package. All authors read and approved this manuscript prior to  
518 submission.

### 519 **Additional Information**

520 Supplementary information accompanies this paper at <http://www.nature.com/srep>.  
521 Competing financial interests: The authors declare no competing financial interests.

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538



## 539 **Figure Legends**

540 **Figure 1. IA-SVA is a robust statistical framework to detect and estimate multiple and**  
541 **correlated hidden sources of variation. (A)** Six-step IA-SVA procedure. IA-SVA computes  
542 the first principal component (PC1) from read counts adjusted for all known factors and tests  
543 its significance [Steps 1-3]. If significant, IA-SVA uses this PC1 to infer a set of genes  
544 associated with the hidden factor [Steps 4-5] and obtain a surrogate variable (SV) to represent  
545 the hidden factor using these genes [Step 6]. **(B)** IA-SVA uses single-cell gene expression  
546 data matrix and known factors to detect hidden sources of variation (e.g., cell contamination,  
547 cell-cycle status, and cell type). If these factors match to a biological variable of interest (e.g.,  
548 cell type assignment), genes highly correlated with the factor can be detected and used in  
549 downstream analyses (e.g., data visualization).

550 **Figure 2. IA-SVA can detect heterogeneity originating from potentially contaminated**  
551 **alpha cells. (A)** Outlier alpha cells captured using IA-SVA and same cells marked in  
552 respective **(C)** PCA, **(D)** USVA, and **(E)** tSNE analyses. Cells are clustered into two groups  
553 (red vs. gray dots) based on IA-SVA's surrogate variable 2 ( $SV_2 > 0.1$ ). In PCA, PC1 was  
554 discarded since it explains the geometric library size. **(B)** Hierarchical clustering of alpha  
555 cells using 27 genes significantly associated with  $SV_2$  ( $FDR < 0.05$  and  $R^2 > 0.6$ ) (ward.D2  
556 and  $cutree\_cols = 2$ ). 6 cells clearly separated from the rest of the alpha cells based on the  
557 expression of these 27 genes.

558 **Figure 3. IA-SVA can detect heterogeneity stemming from differences in cell-cycle stage.**  
559 **(A)** Visualization of glioblastoma cells based on IA-SVA-detected factors ( $SV_1$  and  $SV_2$ ).  
560 Same cells are marked in respective analyses with **(C)** PCA, **(D)** USVA, and **(E)** tSNE  
561 analyses. IA-SVA's  $SV_1$  clearly separates cells into two groups (red vs. blue dots,  $SV_1 > 0.1$ ).  
562 Other methods failed to clearly detect this cell-cycle related heterogeneity. **(B)** Hierarchical  
563 clustering on 119 genes significantly associated ( $FDR < 0.05$  and  $R^2 > 0.3$ ) with IA-SVA's  
564  $SV_2$  confirms the separation of cells based on these genes (ward.D2 and  $cutree\_cols = 2$ ). **(F)**  
565 IA-SVA's  $SV_1$  can segregate cells based on their cell-cycle-stage as predicted by SCRAN.

566 **Figure 4. IA-SVA based gene selection enhances single cell data visualization. (A)** tSNE  
567 analyses using all expressed genes in human islet data (tSNE). Cells are color-coded based on  
568 the original cell-type assignments. Note that cells are not effectively clustered with respect to  
569 their assigned cell types. **(B)** Hierarchical clustering using genes ( $n=92$ ) selected by IA-SVA  
570 clearly separate cell types (ward.D2 and  $cutree\_cols=3$ ). Known marker genes (e.g., *INS*) are

571 highlighted in red color. **(C)** tSNE analyses using the 92 IA-SVA genes (IA-SVA+tSNE).  
572 Note the improved grouping of cell types into discrete clusters. **(D)** tSNE analyses using top  
573 variable genes in a second and bigger islet scRNA-seq data. Note that cells are not effectively  
574 clustered with respect to their assigned cell types just using tSNE. **(E)** tSNE analyses  
575 repeated using genes (n=57) obtained via IA-SVA (IA-SVA+tSNE). Note the improved  
576 clustering of different cell types into discrete clusters. **(F)** tSNE analyses using 1000 most  
577 over-dispersed genes (CellView). **(G)** tSNE analyses on significant PCs obtained from highly  
578 over-dispersed genes (Spectral tSNE).

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600 **Tables**

601

	USVA	SSVA	IA-SVA	USVA	SSVA	IA-SVA
	$ r  = 0.3 \sim 0.6$			$ r  < 0.3$		
Power*(F1**)	1	1	1	1	1	1
Power (F2)	1	1	1	1	1	1
Power (F3)	0.78	0.78	0.87	1	1	1
Cor***(F1)	0.93	0.95	0.95	0.98	0.98	1
Cor (F2)	0.72	0.75	0.94	0.94	0.94	0.99
Cor (F3)	0.75	0.78	0.95	0.93	0.93	0.98
	USVA	SSVA	IA-SVA			
Type I error*	0.09	0.09	0.04			

\*Nominal Type I error rate: 0.05

\*\*F1, F2, F3 refers to Factor1, Factor2, and Factor3

\*\*\*Average of the absolute Pearson correlation coefficients between the true factor and the estimated factor is used as the accuracy measure.

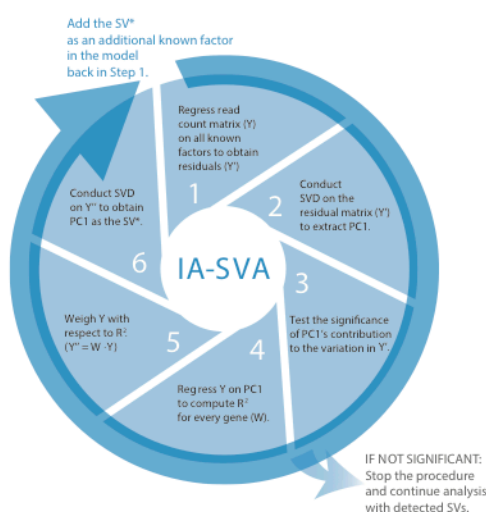
602

603 **Table 1. IA-SVA accurately captures unknown sources of variation while controlling**  
 604 **Type I error rate at a nominal level.** Empirical power, Type I error rate, and the accuracy  
 605 of estimates for IA-SVA, SSVA, and USVA assessed using simulated single-cell gene  
 606 expression data. Alternative scenarios are simulated in which hidden factors are moderately  
 607 ( $|r| \sim 0.3-0.6$ , first three columns) or weakly ( $|r| < 0.3$ , last three columns) correlated with the  
 608 group variable. IA-SVA outperforms alternative methods especially while detecting variation  
 609 stemming from a smaller fraction of genes (10%) and especially when factors are correlated.

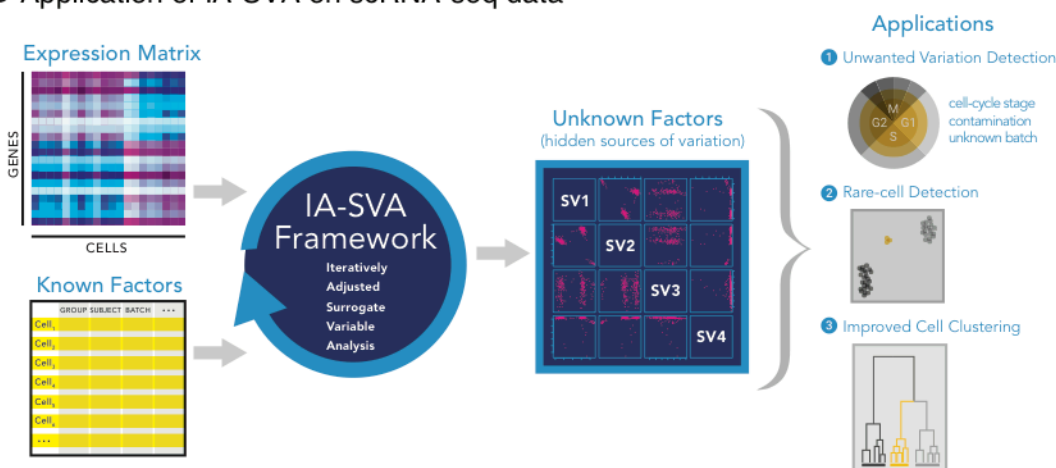
610

611

## A Iterative IA-SVA framework



## B Application of IA-SVA on scRNA-seq data



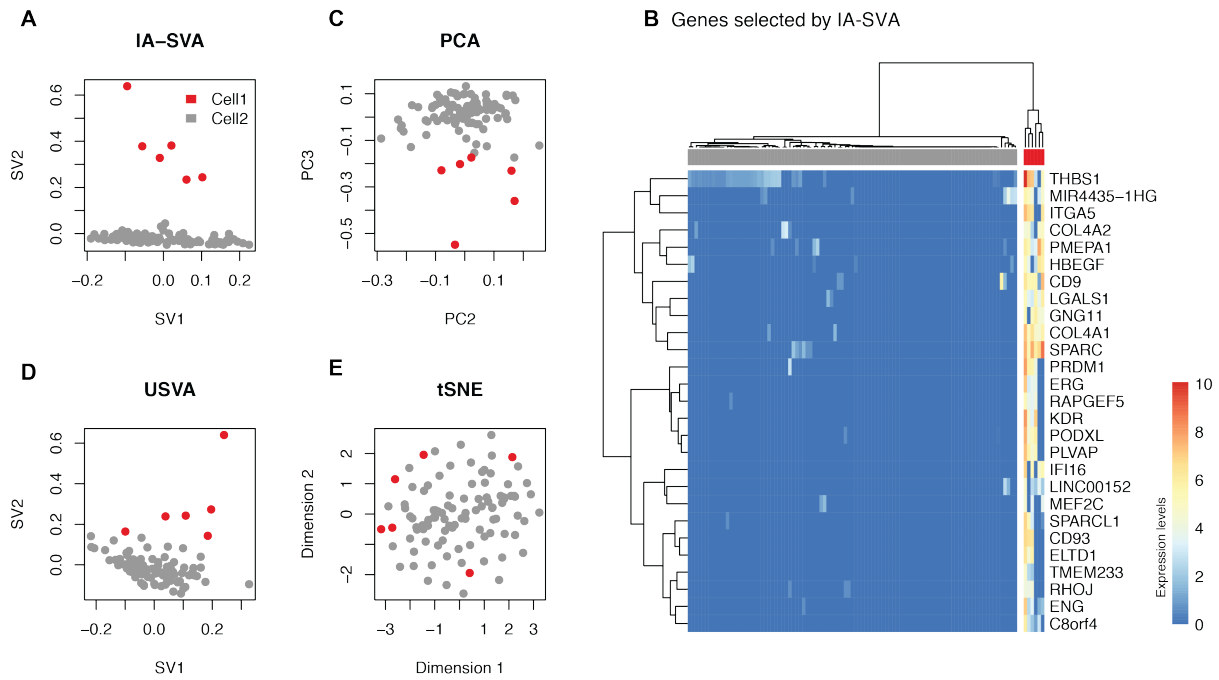
612  
 613 **Figure 1. IA-SVA is a robust statistical framework to detect and estimate multiple and**  
 614 **correlated hidden sources of variation.** (A) Six-step IA-SVA procedure. IA-SVA  
 615 computes the first principal component (PC1) from read counts adjusted for all known factors  
 616 and tests its significance [Steps 1-3]. If significant, IA-SVA uses this PC1 to infer a set of  
 617 genes associated with the hidden factor [Steps 4-5] and obtain a surrogate variable (SV) to  
 618 represent the hidden factor using these genes [Step 6]. (B) IA-SVA uses single-cell gene  
 619 expression data matrix and known factors to detect hidden sources of variation (e.g., cell  
 620 contamination, cell-cycle status, and cell type). If these factors match to a biological variable  
 621 of interest (e.g., cell type assignment), genes highly correlated with the factor can be detected  
 622 and used in downstream analyses (e.g., data visualization).

623  
 624

625

626

627



628

629 **Figure 2. IA-SVA can detect heterogeneity originating from potentially contaminated**  
630 **alpha cells. (A)** Outlier alpha cells captured using IA-SVA and same cells marked in  
631 respective (C) PCA, (D) USVA, and (E) tSNE analyses. Cells are clustered into two groups  
632 (red vs. gray dots) based on IA-SVA's surrogate variable 2 ( $SV2 > 0.1$ ). In PCA, PC1 was  
633 discarded since it explains the geometric library size. (B) Hierarchical clustering of alpha  
634 cells using 27 genes significantly associated with SV2 ( $FDR < 0.05$  and  $R^2 > 0.6$ ) (ward.D2  
635 and  $cutree\_cols = 2$ ). 6 cells clearly separated from the rest of the alpha cells based on the  
636 expression of these 27 genes.

637

638

639

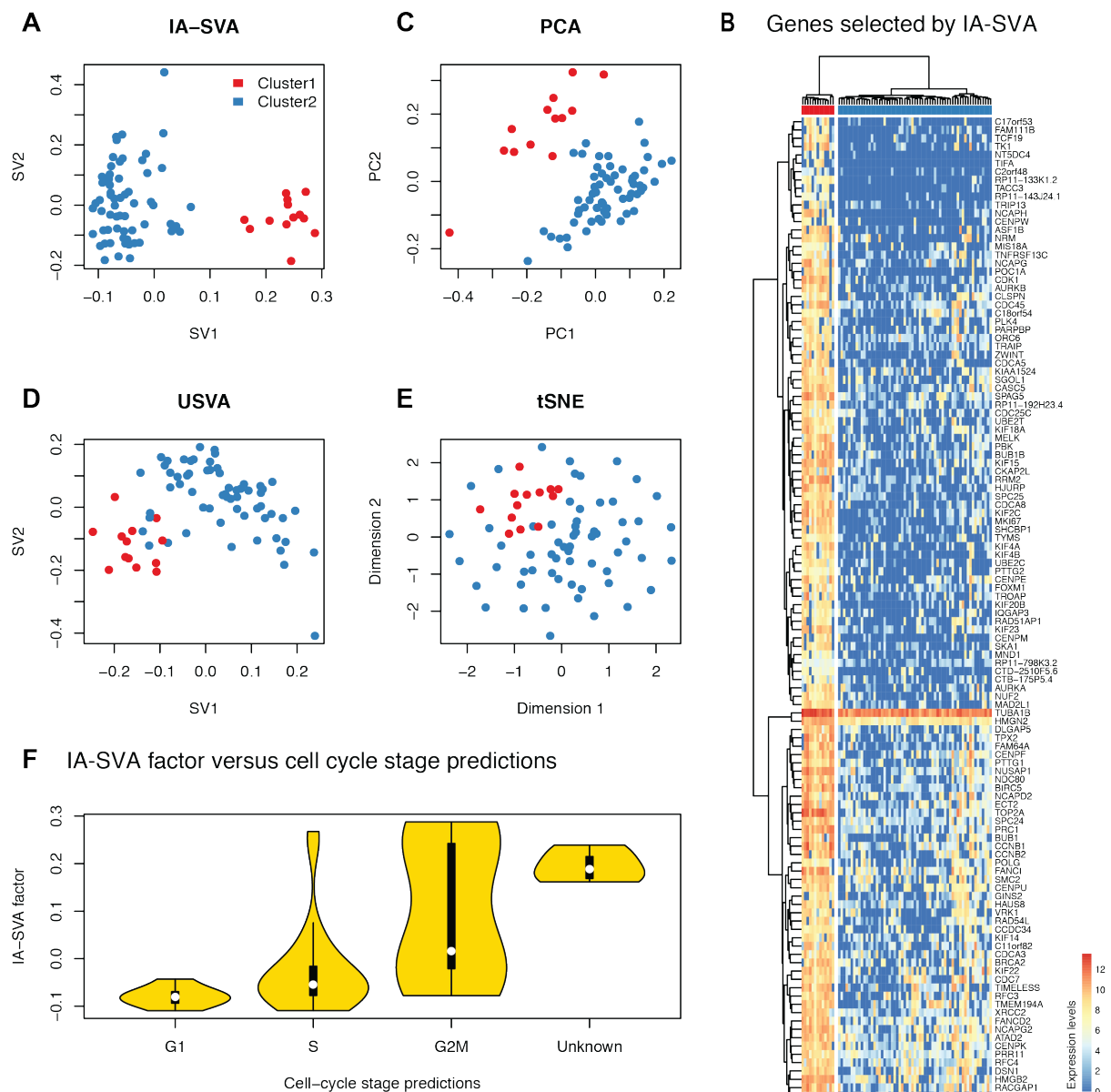
640

641

642

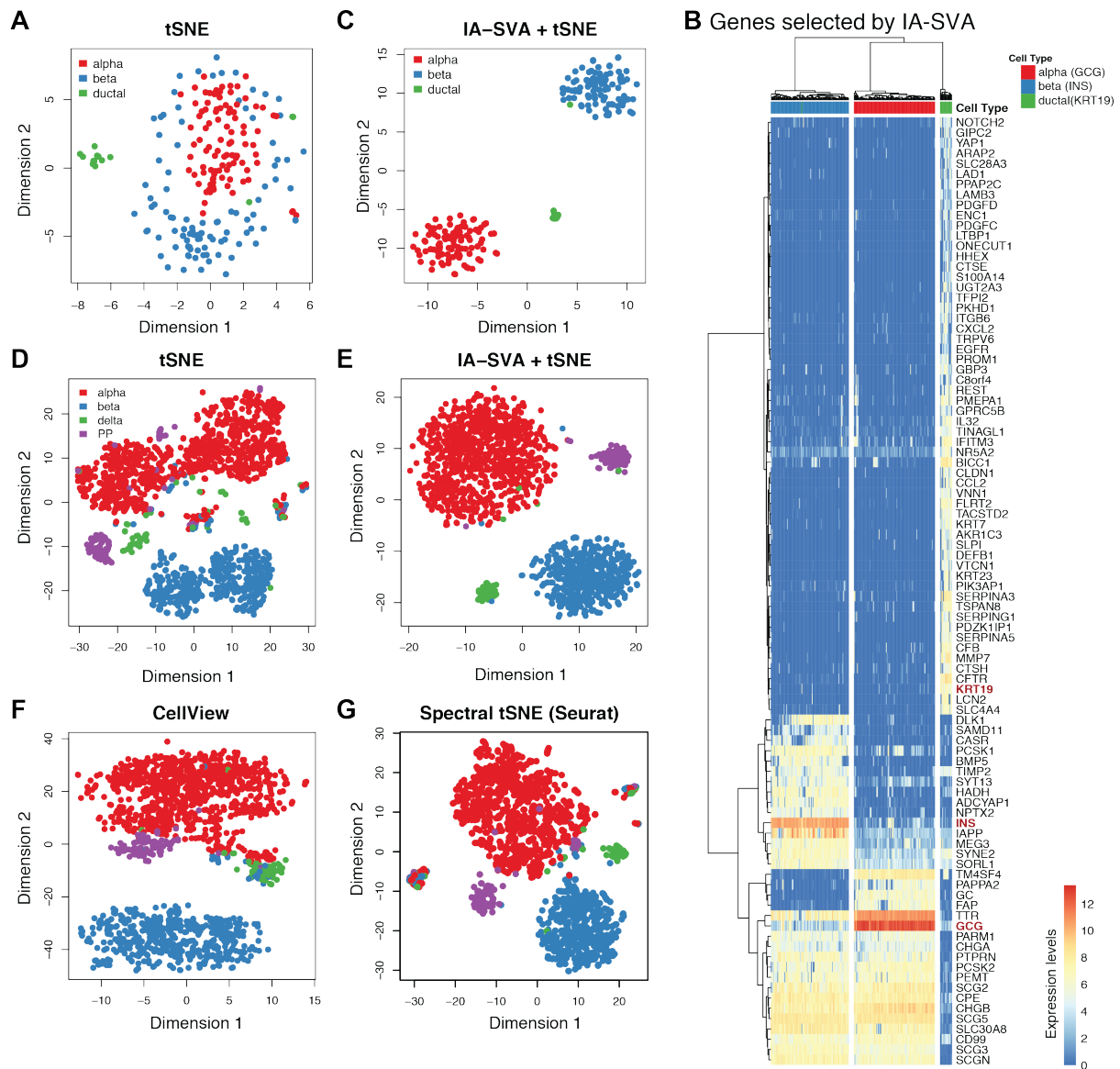
643

644



645  
646  
647  
648  
649  
650  
651  
652  
653  
654

**Figure 3. IA-SVA can detect heterogeneity stemming from differences in cell-cycle stage.** (A) Visualization of glioblastoma cells based on IA-SVA-detected factors (SV1 and SV2). Same cells are marked in respective analyses with (C) PCA, (D) USVA, and (E) tSNE analyses. IA-SVA's SV1 clearly separates cells into two groups (red vs. blue dots, SV1 > 0.1). Other methods failed to clearly detect this cell-cycle related heterogeneity. (B) Hierarchical clustering on 119 genes significantly associated (FDR < 0.05 and  $R^2 > 0.3$ ) with IA-SVA's SV2 confirms the separation of cells based on these genes (ward.D2 and cutree\_cols = 2). (F) IA-SVA's SV1 can segregate cells based on their cell-cycle-stage as predicted by SCRAN.



655  
 656 **Figure 4. IA-SVA based gene selection enhances single cell data visualization.** (A) tSNE  
 657 analyses using all expressed genes in human islet data (tSNE). Cells are color-coded based on  
 658 the original cell-type assignments. Note that cells are not effectively clustered with respect to  
 659 their assigned cell types. (B) Hierarchical clustering using genes (n=92) selected by IA-SVA  
 660 clearly separate cell types (ward.D2 and cutree\_cols=3). Known marker genes (e.g., *INS*) are  
 661 highlighted in red color. (C) tSNE analyses using the 92 IA-SVA genes (IA-SVA+tSNE).  
 662 Note the improved grouping of cell types into discrete clusters. (D) tSNE analyses using top  
 663 variable genes in a second and bigger islet scRNA-seq data. Note that cells are not effectively  
 664 clustered with respect to their assigned cell types just using tSNE. (E) tSNE analyses  
 665 repeated using genes (n=57) obtained via IA-SVA (IA-SVA+tSNE). Note the improved  
 666 clustering of different cell types into discrete clusters. (F) tSNE analyses using 1000 most  
 667 over-dispersed genes (CellView). (G) tSNE analyses on significant PCs obtained from highly  
 668 over-dispersed genes (Spectral tSNE).  
 669