1 **gCAnno: a graph-based single cell type annotation method**

2 Xiaofei Yang[1,2¶], Shenghan Gao[2,3¶], Tingjie Wang[2,3¶], Boyu Yang[2,3], Ningxin Dang[4], Kai

3 Ye[2,3,4,5*]

4 [1] School of Computer Science and Technology, Faculty of Electronic and Information

5 Engineering, Xi'an Jiaotong University, Xi'an, Shaanxi, China

6 [2] MOE Key Lab for Intelligent Networks & Networks Security, Faculty of Electronic

7 and Information Engineering, Xi'an Jiaotong University, Xi'an, Shaanxi, China

8 [3] School of Automation Science and Engineering, Faculty of Electronic and Information

9 Engineering, Xi'an Jiaotong University, Xi'an, Shaanxi, China

10 [4] Genome Institute, the First Affiliated Hospital of Xi'an Jiaotong University, Xi'an,

11 Shaanxi, China

12 [5] The School of Life Science and Technology, Xi'an Jiaotong University, Xi'an, Shaanxi,

13 China

14 * Corresponding author

15 E-mail: kaiye@xjtu.edu.cn (Kai Ye)

16 ¶These authors contributed equally to this work.

17

18  **Abstract**

19  **Background**

20  Current single cell analysis methods annotate cell types at cluster-level rather than

21  ideally at single cell level. Multiple exchangeable clustering methods and many tunable

22  parameters have a substantial impact on the clustering outcome, often leading to

23  incorrect cluster-level annotation or multiple runs of subsequent clustering steps. To

24  address these limitations, methods based on well-annotated reference atlas has been

25  proposed. However, these methods are currently not robust enough to handle datasets

26  with different noise levels or from different platforms.

27  **Results**

28  Here, we present gCAnno, a graph-based Cell type Annotation method. First, gCAnno

29  constructs cell type-gene bipartite graph and adopts graph embedding to obtain cell type

30  specific genes. Then, naïve Bayes (gCAnno-Bayes) and SVM (gCAnno-SVM)

31  classifiers are built for annotation. We compared the performance of gCAnno to other

32  state-of-art methods on multiple single cell datasets, either with various noise levels or

33  from different platforms. The results showed that gCAnno outperforms other state-of-

34  art methods with higher accuracy and robustness.

35  **Conclusions**

36    gCAnno is a robust and accurate cell type annotation tool for single cell RNA analysis.

37    The source code of gCAnno is publicly available at https://github.com/xjtu-

38    omics/gCAnno.

39    **Keywords:** graph embedding, cell type annotation, single cell RNA analysis


40    **Background**

41    Bulk RNA sequencing measures average gene expression level in a large population of

42    cells, hindering dissection of heterogeneous cell types [1]. In 2009, single cell RNA

43    sequencing (scRNA-seq) technology was developed to provide valuable insights into

44    cell heterogeneity [2].

45       In general, accurate cell type annotation for single cell data is a prerequisite for

46    any further investigation of cell heterogeneous [3-6]. The commonly used cell type

47    annotation methods, including Seurat [7], SCANPY [8] and SINCERA [9], adopts a

48    similar procedure of data quality control, reads mapping, UMI quantification,

49    expression normalization, clustering, differentially expressed genes (DEGs) of each

50    cluster identification and cell type assignment based on biomarker genes [10]. However,

51    those methods report cluster-level rather than truly single cell-level annotation results,

52    masking subtle differences within each cluster. In addition, different clustering methods

53    and many tunable parameters led to uncertain clustering outcome. These above two

54    factors cause incorrect cluster-level annotations or multiple runs of subsequent

55    clustering steps [10].

56        To overcome the above issues, two distinct strategies, namely biomarker-based

57    and reference-based approaches, have been proposed. The biomarker-based methods,

58    such as Garnett [11] and CellAssign [12], aim to establish mappings between the query

59    dataset and the well-studied biomarkers. In particular, Garnett trains a classifier based

60    on the user defined markup language. CellAssign builds a probabilistic model that

61    leverages prior knowledge of cell-type marker genes for annotation. However,

62    collecting a comprehensive biomarker set of different cell types is cumbersome, time-

63    consuming and subjective [13]. Thus recently reference-based approaches, such as

64    Scmap [14], Chetah [15] and scPred [16] have been developed and are gaining

65    popularity after a number of well-annotated single cell data were published, especially

66    the datasets released by human cell atlas (HCA) [17]. The reference-based methods

67    follow data-driven strategy and construct mappings between query dataset and the well-

68    annotated reference datasets. For example, Scmap uses drop-based method to select

69    feature genes as variables and constructs mapping by distance and correlation

70    coefficient. Another method, scPred selects differential principle components (PCs)

71    calculated by gene expression value between cell types and trains an SVM model with

72 these PCs. However, these methods are sensitive to experiment batches, sequencing

73 platforms and noises, all of which are intrinsic properties of the single cell datasets.

74 Here, we propose a reference-based method, gCAnno, using graph representation

75 feature selection strategy to comprehensively represent the global view of associations

76 between cell types and genes for robust and high accuracy single cell-level annotation.

77 Our gCAnno method starts with construction of a weighted cell type-gene bipartite

78 graph. Then, graph embedding is applied to capture the cell type specific genes and

79 naïve Bayes (gCAnno-Bayes) and SVM (gCAnno-SVM) classifiers are built for further

80 annotation (Fig. 1). We compared gCAnno with the state-of-the-art methods on four

81 published datasets as the basic test [3-6]. We also reported the performance comparison

82 on large dataset with deep annotation level [18], different single cell platforms,

83 simulated datasets with either various cell type imbalance situations and different

84 dropout noise levels as the advanced test. Finally, runtime is summarized to

85 demonstrate the efficiency of gCAnno.

86 **Methods**

87 Here we summarized the framework of gCAnno. gCAnno adopts graph structure for

88 cell type specific gene set detection and accurate cell type annotation. Firstly, gCAnno

89 builds cell type-gene bipartite graph based on gene expression abundances and

4

90    intensities, in which gene expression abundance is the proportion of cells expressing

91    the gene in a given cell type while intensity is the average expression in cells expressing

92    the gene. Then, graph embedding is adopted to obtain the embedding vectors of gene

93    nodes and cell type nodes. Next, gCAnno selects a set of genes for each cell type with

94    similar profiles in the embedding space. Finally, based on the detected cell type specific

95    genes, gCAnno trains naïve Bayes and SVM classifiers. The workflow of gCAnno is

96    depicted in Fig. 1.

97    **Cell type-gene bipartite graph construction**

98    Starting from the well-annotated reference scRNA-seq data, we constructed a weighted

99    cell type-gene bipartite graph (wCGBG) containing both cell type nodes (CTN) and

100   gene nodes (GN). Edges between CTN and GN indicate the correlation of a gene and a

101   cell type while weight $W$ measures significance of correlation. The weight is

102   calculated by:

103
$$w_{k,j} = \begin{cases} \dfrac{m_{k,j}}{n_k} \times mean\left(\overrightarrow{g_{k,j}}\right), & if \ \ n_k \neq 0 \\ 0 & , \ \ others \end{cases}$$
(1)

104   where $n_k$ is the cell count of cell type $k$, $m_{j,k}$ is the number of cells expressed gene

105   $j$ in cell type $k$. $\overrightarrow{g_{j,k}}$ is the expression vector of gene $j$ in cell type $k$. $W$ is

106   the product of the gene expression abundance and intensity. We use gene expression

5

107 abundance and intensity to establish a relationship between cell types and genes in the

108 form of proportion to reduce the impact of individual gene loss (dropout) or cell number

109 imbalance.

110 **Graph embedding and cell type-gene specific relation detection**

111 After wCGBG construction, we used node2vec to obtain the low dimensional vectors

112 (the embedding vectors) of gene nodes and cell type nodes. The first step is construction

113 of a neighborhood set $N(u)$ of each node $u$ (either gene or cell type node) by a

114 probability walk [19]. Then, we optimized the following objective function $f(u)$ by

115 maximizing the log-probability of observing a neighborhood set.

116
$$\max_f \sum_{u \in V} \log P\big(N(u) \mid f(u)\big) \tag{2}$$

117 This optimization step enables the embedding vectors to capture the specificity

118 and strength of interactions between cell node and gene node, e.g. if one gene is specific

119 and highly expressed in one cell type, the corresponding two embedding vectors are

120 similar. Then, we calculated Euclidean distance between the vector of genes and cell

121 types. We selected top $n$ (a user defined parameter, default $n = 65$, Additional file

122 1:Figure S1) closest genes for each cell type as the cell type specific gene set based on

123 the overall performance on the five datasets we used [3-6][18].

124 **Classifier construction**

6

125      After obtaining the cell type specific gene set, we build naïve Bayes (gCAnno-Bayes)

126      and SVM (gCAnno-SVM) classifiers for annotation. For gCAnno-SVM, we directly

127      use the expression of cell type specific genes as features to train an SVM classifier. For

128      gCAnno-Bayes, we build a binary matrix to presents cell type and its corresponding

129      specific genes, e.g. the element $b_{ij} = 1$ indicates gene $j$ is one of the specific genes in

130      cell type $i$. We train a Bernoulli Naïve Bayes to get genes' conditional probability in

131      each cell type and the prior probability of cell types. The query dataset is binarized and

132      the annotation is based on maximum posterior probability of single cell's cell type

133      specific genes expression.

134

135      **Performance measurement and dataset**

136      **Performance assessment and comparison.** Cell type annotation is a typical multi-

137      classification problem. We applied kappa coefficient as the performance measurement

138      of classification, defined as equation (3).

139
$$\kappa = \frac{p_o - p_e}{1 - p_e}, \quad p_o = \frac{N_{corr}}{N_t}, \quad p_e = \frac{\sum_{i=1}^{K} a_i \times b_i}{N_t \times N_t} \tag{3}$$

140      where $N_{corr}$ is the ratio of total number of cells with corrected cell type annotation,

141      $N_t$ is the total number of cells in the dataset, $K$ is the number of truly cell types, $a_i$

142      is the number of corrected annotated cells in the $i$-th cell type, and $b_i$ is the number

7

143　　of cells in the $i$-th cell type，$p_o$ is the accuracy, $a_i \times b_i$ is the product of the actual

144　　and predicted quantity, $p_e$ punishes bias for unbalance evaluation.

145　　　　To evaluate the performance of gCAnno, we performed both cross-validation test

146　　and independent heterogeneous test (cross-platform test). First, we adopted the five-

147　　fold cross-validation strategy following recent single cell analysis comparison

148　　published earlier [13, 16] on four published datasets and simulated noise datasets to

149　　evaluate the overall and robustness performance (Additional file 2: File S1). Then, we

150　　performed independent test on datasets from different sequencing platforms (the cross-

151　　platform testing) to evaluate the generalization capability of gCAnno.

152　　**Tools in comparison.** The calculation results of Scmap, Chetah and scPred were

153　　obtained from the corresponding publications [14-16]. For SVM, we followed the

154　　previous report [13] which is using drop-based method [20] for feature selection.

155　　**Datasets used in basic overall performance test.** To illustrate the stable performance

156　　of gCAnno across various species and tissue types, we compared gCAnno with other

157　　methods using four published datasets, including liver, pancreas, Arabidopsis thaliana

158　　root (AT root), hepatocellular carcinoma and intrahepatic cholangiocarcinoma (HCC

159　　and ICCA) datasets (Table 1; Additional file 2: File S1; Additional file 3: Figure S2;

8

160 Additional file 4: Table S1). The true labels of the cells in each dataset are obtained

161 from the corresponding publications.

162 **Table 1 The list of scRNA-seq datasets in overall performance test.**

| Dataset | #Cells | #Genes | # Cell types |
|---|---|---|---|
| Liver [4] | 8,444 | 20,007 | 14 |
| Pancreas [3] | 8,562 | 20,126 | 13 |
| AT root [6] | 7,053 | 32,833 | 19 |
| HCC, ICCA [5] | 4,729 | 19,379 | 8 |

163 Note: # means the number of

164 **Large dataset with deep annotation level.** To demonstrate the performance of

165 gCAnno in large dataset (cell number more than 50,000) with deep annotation level

166 (more than 20 cell types). We compared gCAnno with other methods in 20 mouse

167 organs dataset with 54,246 cells, 29 cell types and 23,433 genes. The true labels of the

168 cells in each dataset are also obtained from the original publications [18]. (Additional

169 file 2: File S1; Additional file 5: Figure S3; Additional file 6: Table S2).

170 **Simulated dropout and imbalance datasets.** To evaluate the robustness of gCAnno

171 in the presence of dropout noise, we simulated different dropout rates in four above

172 datasets (Table 1), by modifying the expression level of a random gene subset (10%,

173 20%, 30%, 40% and 50% of all genes) to zero (Additional file 2: File S1). Similarly,

174 we used five-fold cross validation to evaluate its performance. In each validation, we

175 simulated the dropout noise in either training group (reference dropout) or test group

176 (query dropout), and calculated the kappa coefficient for each method.

177     To simulate the cell number imbalance noise, we randomly sampled different

178 proportions (0.1:1, 0.3:1, 0.5:1, 0.7:1, 0.9:1, 1:0.9, 1:0.7, 1:0.5, 1:0.3 and 1:0.1) of cell

179 count in two cell types (Hepatocyte and GamaDetaT) in liver dataset as the reference

180 data for classifier constructing. To get more accuracy testing, this simulation was

181 repeated five times (Additional file 2: File S1).

182 **Cross platform datasets.** To compare cross platform performance (various studies

183 using different sequencing platforms), we searched and identified four datasets suitable

184 for this purpose, including two liver datasets from 10x and mCel-seq2 platforms and

185 two pancreas datasets from drop-seq and smart-seq2 platforms (Table 2). We noticed

186 that the cell type annotation labels of the same tissue from different platforms are not

187 identical. Thus, we unified the labels by removing cell types absent in either of the

188 datasets (Additional file 7: Figure S4; Additional file 8: Table S8; Additional file 2: File

189 S1).

190 **Table 2 The list of scRNA-seq datasets in cross platform test**

| Dataset | #Cells | #Genes | # Cell types | Platform |
|---|---|---|---|---|
| Liver [4] | 8,103 | 20,007 | 7 | 10x |
| Pancreas [3] | 8,037 | 20,126 | 9 | Drop-seq |

10

| Liver [21] | 7,130 | 33,941 | 7 | mCel-seq2 |
| Pancreas [22] | 2,068 | 25,526 | 9 | Smart-seq2 |

191 Note: # means the number of

## Results

193 To evaluate the performance of gCAnno, we first evaluated the cell type-gene specific

194 relation, and then compared gCAnno with five state-of-art methods, including Scmap-

195 cell, Scmap-cluster, Chetah, scPred and SVM, in the following four aspects: 1) cell type

196 specificity of gCAnno detected genes, 2) overall performance on different scRNA-seq

197 datasets, 3) robustness test on simulated drop-out and imbalance noise data, 4) cross

198 platform annotation.

**Cell type specificity of gene sets detected by gCAnno**

200 After graph embedding step, gCAnno selects cell type specific gene sets, which largely

201 determines the performance of our approach. Thus, we first evaluated the cell type

202 specificity of gene sets detected in the four datasets. We noticed that clear cell type

203 specific expression patterns are observed for these selected genes (Additional file 9:

204 Figure S5; Additional file 10: Figure S6; Additional file 11: Figure S7). Among the

205 reported marker genes from the corresponding publications, gCAnno is able to capture

206 an average of 57% of them, indicating gCAnno's effectiveness of cell type specific

207 gene identification (Additional file 12: Figure S8; Additional file 13: Table S4).

11

**Overall and large dataset performance evaluation**

We next evaluated and compared overall performance of gCAnno, Scmap, scPred, Chetah and SVM with four published scRNA-seq datasets (Table 1). We found that the comprehensive kappa coefficient of both gCAnno was consistently much higher than those of Scmap-cluster, Scmap-cell and scPred, respectively ($p < 0.05$, Wilcoxon rank sum test) (Fig. 2a-2d) (Additional file 14: Table S5), hinting gCAnno's better performance than other methods on cell type annotation across different species (e.g. human or plant), organs (e.g. liver or pancreases), or disease states (e.g. health or cancer). In 20 mouse organs dataset, the comprehensive kappa coefficient of both gCAnno were 0.74 (gCAnno-Bayes) and 0.94 (gCAnno-SVM), and other methods achieve 0.16 (Scmap-cluster), 0.18 (Scmap-cell), 0.80 (Chetah), 0.63 (scPred) and 0.92 (SVM), respectively. We found that gCAnno-SVM achieved highest performance than other methods in large dataset with deep annotation level (Additional file 6: Table S2; Additional file 15: Figure S9).

**Robustness on dropout and imbalance noisy data**

Besides basic accuracy, we examined its robustness in the presence of different types of noises. Dropout and cell count imbalance noises are two major types and the most challenging in scRNA-seq data. Dropout is a technical noise in the form of missing

12

226   value in gene expression [10], while cell number imbalance among cell types is coming

227   from biology itself. We found gCAnno achieved the highest and rather stable kappa

228   coefficients for both reference dropout and query dropout tests in four datasets (Fig. 3;

229   Additional file 16: Figure S10; Additional file 17: Table S6; Additional file 18: Figure

230   S11). Remarkably, gCAnno achieved average kappa coefficients of 0.88 (gCAnno-

231   SVM) and 0.79 (gCAnno-Bayes) even when dropout rate was as high as 50%, while

232   other methods achieve 0 (Scmap-cluster), 0.44 (Scmap-cell), 0.37 (Chetah), 0.25

233   (scPred) and 0.79 (SVM), respectively. Moreover, we found gCAnno, SVM and

234   Scmap-cell achieved the highest and stable kappa coefficients (average values are about

235   0.99) for different cell count imbalance ratios (Additional file 16: Figure S10;

236   Additional file 19: Table S7). All of these results show gCAnno is better than other

237   methods for dropout and cell count imbalance noises and achieved the best performance

238   on highly noisy data (e.g. 50% dropout rate and 1:0.1 imbalance rate), suggesting the

239   effectiveness of the wCGBG in selecting accurate features in the presence of high noise.

240   **Cross platform annotation**

241   Different single cell sequencing platforms have platform specific features or bias [23],

242   limiting cross platform cell type annotation. We evaluated the platform compatibility

243   of gCAnno on two liver datasets and two pancreas datasets from four platforms (10x,

244   mCel-seq2, Drop-seq, and Smart-seq2) (Table 2). We used one platform dataset as the

245   training data and the other as the testing data. For the performance comparison, gCAnno

246   achieved consistently high kappa coefficient values for liver dataset tests (Fig. 4a and

247   Fig. 4b) and for pancreas dataset tests (Fig. 4c and Fig. 4d) (Additional file 20: Table

248   S8). These results show gCAnno is able to maintain high annotation accuracy for real

249   heterogeneous and cross platform data in the presence of systematic platform specific

250   bias.

251   **Runtime evaluation**

252   Finally, we evaluated the runtime of gCAnno based on datasets in above tests

253   (Additional file 21: Table S9; Additional file 22: Figure S12). We found that the time

254   takes in model building (including graph construction and embedding) step is positive

255   correlated with the number of graph nodes (Pearson's correlation is 0.94). Once the

256   model has been built, the annotation step only takes less than 1 minute (e.g. for mCel-

257   seq2 platform liver dataset with 8103 cells only takes 48 seconds).

258   **Discussion**

259   In this study, we present gCAnno, a novel graph-based cell type identification method

260   for scRNA-seq data. The most significant feature of gCAnno is the construction of

261   wCGBG, enabling gCAnno to capture the global characteristics of association between

14

262    cell types and genes. This feature allows gCAnno to detect accurate feature genes for

263    each cell type, leading to accurate annotation results and robustness for different noise

264    types and rates. In addition, gCAnno is able to annotate not only human scRNA-seq,

265    but also plant scRNA-seq (e.g. Arabidopsis data). Its stable and high performance

266    across platforms, indicates wide application as a "pan-platform" method.

267        gCAnno contains SVM version (gCAnno-SVM) and naïve Bayes version

268    (gCAnno-Bayes). The SVM version takes into account the effect of expression value

269    while naïve Bayes version only considers the existence of cell type specific genes. From

270    the evaluation result, the SVM version seems suitable for the dataset with deep

271    annotation level and contains largely similar cell types between training and test sets.

272    However, in cross platform datasets from different studies and different sequencing

273    platforms, gene expression value might fluctuate significantly, rendering better

274    performance of naïve Bayes version than SVM version.

275        Since gCAnno is a reference-based cell type annotation method, it lacks the ability

276    to identify novel cell types. For novel type cells, gCAnno assigns the closest cell types

277    with the most similar expression profiles to them, which might be reasonable in most

278    of applications but probably require further improvement. Integrating the biomarker-

279 based method for novel cell type annotation and reference-based method for accurate

280 pre-defined cell type annotation, we think, will be one direction to explore.


281 **Conclusion**

282 We have implemented a stable and high-performance automated cell type annotation

283 tool, gCAnno, for scRNA-seq datasets. With an easy use Python running script as an

284 example, we hope gCAnno will be useful for the scRNA-seq data analysis.


285 **Abbreviations**

286 **kappa:** kappa coefficient

287 **DEGs:** differentially expressed genes

288 **UMI:** Unique Molecular Identifier

289 **scRNA-seq:** Single-cell RNA-seq

290 **PCs:** principle components

291 **HCA:** human cell atlas

292 **wCGBG:** weighted cell type-gene bipartite graph

293 **10x:** 10x Genomics platform

294 **SVM:** Support Vector Machine


295 **Declarations**

**Ethics approval and consent to participate**

This study used previously published data and did not obtain any new data directly

involved humans, plants or animals.

**Consent for publication**

Not applicable.

**Availability of data and materials**

Datasets used for the analyses in this study are summarized in Additional file 2: File S1.

The source code of gCAnno is publicly available at https://github.com/xjtu-

omics/gCAnno.

**Competing interests**

The authors declare that they have no competing interests.

314 **Authors' contributions**

315 KY and XY conceived the study. SG designed and performed the experiments. SG and

316 BY analysed the data. SG developed the program. XY and SG wrote the manuscript.

317 SG and ND completed figures of manuscript. All authors read and approved the final

318 manuscript.

319 **Acknowledgements**

322 **Authors' information**

323 **School of Computer Science and Technology, Faculty of Electronic and**
324 **Information Engineering, Xi'an Jiaotong University, Xi'an, Shaanxi, China**
325 Xiaofei Yang
326 **MOE Key Lab for Intelligent Networks & Networks Security, Faculty of**
327 **Electronic and Information Engineering, Xi'an Jiaotong University, Xi'an,**
328 **Shaanxi, China**
329 Kai Ye, Xiaofei Yang, Shenghan Gao, Tingjie Wang, Boyu Yang
330 **School of Automation Science and Engineering, Faculty of Electronic and**
331 **Information Engineering, Xi'an Jiaotong University, Xi'an, Shaanxi, China**
332 Kai Ye, Shenghan Gao, Tingjie Wang, Boyu Yang
333 **Genome Institute, the First Affiliated Hospital of Xi'an Jiaotong University, Xi'an,**
334 **Shaanxi, China**
335 Kai Ye, Ningxin Dang
336 **The School of Life Science and Technology, Xi'an Jiaotong University, Xi'an,**
337 **Shaanxi, China**

338 Kai Ye

**References**

1. Kolodziejczyk AA, Kim JK, Svensson V, Marioni JC, Teichmann SA: **The technology and biology of single-cell RNA sequencing**. *Mol Cell* 2015, **58**(4):610-620.

2. Tang F, Barbacioru C, Wang Y, Nordman E, Lee C, Xu N, Wang X, Bodeau J, Tuch BB, Siddiqui A *et al*: **mRNA-Seq whole-transcriptome analysis of a single cell**. *Nat Methods* 2009, **6**(5):377-382.

3. Baron M, Veres A, Wolock SL, Faust AL, Gaujoux R, Vetere A, Ryu JH, Wagner BK, Shen-Orr SS, Klein AM *et al*: **A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas Reveals Inter- and Intra-cell Population Structure**. *Cell Syst* 2016, **3**(4):346-360 e344.

4. MacParland SA, Liu JC, Ma XZ, Innes BT, Bartczak AM, Gage BK, Manuel J, Khuu N, Echeverri J, Linares I *et al*: **Single cell RNA sequencing of human liver reveals distinct intrahepatic macrophage populations**. *Nat Commun* 2018, **9**(1):4383.

5. Ma L, Hernandez MO, Zhao Y, Mehta M, Tran B, Kelly M, Rae Z, Hernandez JM, Davis JL, Martin SP *et al*: **Tumor Cell Biodiversity Drives Microenvironmental Reprogramming in Liver Cancer**. *Cancer Cell* 2019, **36**(4):418-430 e416.

6. Zhang TQ, Xu ZG, Shang GD, Wang JW: **A Single-Cell RNA Sequencing Profiles the Developmental Landscape of Arabidopsis Root**. *Mol Plant* 2019, **12**(5):648-660.

7. Butler A, Hoffman P, Smibert P, Papalexi E, Satija R: **Integrating single-cell transcriptomic data across different conditions, technologies, and species**. *Nat Biotechnol* 2018, **36**(5):411-420.

8. Wolf FA, Angerer P, Theis FJ: **SCANPY: large-scale single-cell gene expression data analysis**. *Genome Biol* 2018, **19**(1):15.

9. Guo M, Wang H, Potter SS, Whitsett JA, Xu Y: **SINCERA: A Pipeline for Single-Cell RNA-Seq Profiling Analysis**. *PLoS Comput Biol* 2015, **11**(11):e1004575.

10. Kiselev VY, Andrews TS, Hemberg M: **Challenges in unsupervised clustering of single-cell RNA-seq data**. *Nat Rev Genet* 2019, **20**(5):273-282.

374 11. Pliner HA, Shendure J, Trapnell C: **Supervised classification**
375 **enables rapid annotation of cell atlases**. *Nat Methods* 2019,
376 **16**(10):983-986.
377 12. Zhang AW, O'Flanagan C, Chavez EA, Lim JLP, Ceglia N,
378 McPherson A, Wiens M, Walters P, Chan T, Hewitson B *et al*:
379 **Probabilistic cell-type assignment of single-cell RNA-seq for**
380 **tumor microenvironment profiling**. *Nat Methods* 2019,
381 **16**(10):1007-1015.
382 13. Abdelaal T, Michielsen L, Cats D, Hoogduin D, Mei H, Reinders
383 MJT, Mahfouz A: **A comparison of automatic cell identification**
384 **methods for single-cell RNA sequencing data**. *Genome Biol*
385 2019, **20**(1):194.
386 14. Kiselev VY, Yiu A, Hemberg M: **scmap: projection of single-cell**
387 **RNA-seq data across data sets**. *Nat Methods* 2018, **15**(5):359-
388 362.
389 15. de Kanter JK, Lijnzaad P, Candelli T, Margaritis T, Holstege FCP:
390 **CHETAH: a selective, hierarchical cell type identification**
391 **method for single-cell RNA sequencing**. *Nucleic Acids Res* 2019,
392 **47**(16):e95.
393 16. Alquicira-Hernandez J, Sathe A, Ji HP, Nguyen Q, Powell JE:
394 **scPred: accurate supervised method for cell-type classification**
395 **from single-cell RNA-seq data**. *Genome Biol* 2019, **20**(1):264.
396 17. Regev A, Teichmann SA, Lander ES, Amit I, Benoist C, Birney E,
397 Bodenmiller B, Campbell P, Carninci P, Clatworthy M *et al*: **The**
398 **Human Cell Atlas**. *Elife* 2017, **6**.
399 18. Tabula Muris C, Overall c, Logistical c, Organ c, processing,
400 Library p, sequencing, Computational data a, Cell type a, Writing g
401 *et al*: **Single-cell transcriptomics of 20 mouse organs creates a**
402 **Tabula Muris**. *Nature* 2018, **562**(7727):367-372.
403 19. Grover A, Leskovec J: **node2vec: Scalable Feature Learning for**
404 **Networks**. *KDD* 2016, **2016**:855-864.
405 20. Andrews TS, Hemberg M: **M3Drop: dropout-based feature**
406 **selection for scRNASeq**. *Bioinformatics* 2019, **35**(16):2865-2867.
407 21. Aizarani N, Saviano A, Sagar, Mailly L, Durand S, Herman JS,
408 Pessaux P, Baumert TF, Grun D: **A human liver cell atlas reveals**
409 **heterogeneity and epithelial progenitors**. *Nature* 2019,

410     **572**(7768):199-204.

411 22.   Segerstolpe A, Palasantza A, Eliasson P, Andersson EM,

412      Andreasson AC, Sun X, Picelli S, Sabirsh A, Clausen M, Bjursell

413      MK *et al*: **Single-Cell Transcriptome Profiling of Human**

414      **Pancreatic Islets in Health and Type 2 Diabetes**. *Cell Metab*

415      2016, **24**(4):593-607.

416 23.   Ziegenhain C, Vieth B, Parekh S, Reinius B, Guillaumet-Adkins A,

417      Smets M, Leonhardt H, Heyn H, Hellmann I, Enard W:

418      **Comparative Analysis of Single-Cell RNA Sequencing**

419      **Methods**. *Mol Cell* 2017, **65**(4):631-643 e634.

420

421 **Figure legends**

422 **Fig 1. Overview of gCAnno. (a)** Cell type-gene graph building. The graph contains

423 gene nodes (gray circles) and cell type nodes (other color circles). **(b)** Graph embedding

424 converts graphs into low dimensional vectors. Genes are selected based on the distance

425 between the two types of vectors. **(c)** Training Naïve Bayes and SVM classifiers for

426 annotation. **(d)** Cell type annotation for new query dataset.

427 **Fig. 2. Overall performance evaluation.** Comparisons of gCAnno with Scmap-

428 Cluster, Scmap-Cell, scPred, Chetah and SVM based on kappa coefficient on **(a)** liver

429 dataset, **(b)** pancreas dataset, **(c)** HCC & ICCA dataset, and **(d)** AT root dataset. *: *p*-

430 values < 0.1; **: *p*-values < 0.05; ***: *p*-values < 0.01, Wilcoxon rank sum test. The

431 number is the mean of five cross validation. The error bar is the standard deviation. The

432 *y*-axis is the kappa coefficient.

433     **Fig. 3. Robustness performance evaluation.** Robustness of dropout noise

434     comparisons of gCAnno with Scmap-Cluster, Scmap-Cell, scPred, Chetah and SVM

435     on **(a)** liver reference dropout dataset, **(b)** liver query dropout dataset, **(c)** pancreas

436     reference dropout dataset, **(d)** pancreas query dropout dataset. The middle point is the

437     mean kappa coefficients of five-fold cross validation. The error bar is the standard

438     deviation. The *y*-axis is the kappa coefficient and the *x*-axis is the dropout rate.

439     **Fig 4. Platform compatibility evaluation.** Performance comparisons of gCAnno with

440     Scmap-Cluster, Scmap-Cell, scPred, Chetah and SVM on cross platform datasets. **(a)**

441     liver datasets, where reference is mCel-seq2 and query is 10x; **(b)** liver datasets, where

442     reference is 10x and query is mCel-seq2; **(c)** pancreas dataset, where reference is drop-

443     seq and query is smart-seq2 **(d)** pancreas datasets, where reference is smart-seq2 and

444     query is drop-seq. The reference is the training data and the query is the testing data.

445     **Additional files**

446     **Additional file 1.pdf: Figure S1.** The test of gCAnno parameter top closest genes in

447     five evaluation datasets. The parameter is stable in 25 to 85. When top gene select less

448     than 5 (in all datasets) and more than 125 (in Arabidopsis and liver datasets), the

449     performance are not well. In our evaluation, the default top closest genes in each cell

450     type is 65 and user can adjustment by themselves.

22

451    **Additional file 2.docx: File S1.** Supplementary Materials, including data preparation,

452    cell type information of each datasets, and supplementary methods.

453    **Additional file 3.pdf: Figure S2.** The tSNE plot of (a) liver, (b) pancreas, (c) HCC &

454    ICCA and (d) AT root datasets.

455    **Additional file 4.xlsx: Table S1.** The tSNE result, cell barcodes and cell type labels of

456    (a) liver, (b) pancreas, (c) HCC & ICCA and (d) AT root datasets.

457    **Additional file 5.pdf: Figure S3.** The tSNE plot of a large dataset with deep annotation

458    level (20 mouse organs).

459    **Additional file 6.pdf: Table S2.** The large dataset kappa coefficient result (Fig 2e) and

460    tSNE result.

461    **Additional file 7.pdf: Figure S4.** The tSNE plot of (a) mCel-seq2 liver, (b) 10x liver,

462    (c) Drop-seq pancreas and (d) Smart-seq2 pancreas.

463    **Additional file 8.xlsx: Table S3.** The tSNE result, cell barcodes and cell type labels of

464    (a) mCel-seq2 liver, (b) 10x liver, (c) Drop-seq pancreas and (d) Smart-seq2 pancreas.

465    **Additional file 9.pdf: Figure S5.** The heatmap of each cell type specific genes

466    expression in four datasets (top closest gene number 65). It shows an obvious pattern

467    in diagonal.

468  **Additional file 10.pdf: Figure S6.** The tSNE of embedding vectors of cell type nodes

469  and gene nodes in four datasets. The selecting gene nodes are in red color and non-

470  selecting gene nodes are in grey. The cell type nodes are blue triangles.

471  **Additional file 11.pdf: Figure S7.** An example of top 2 specific genes in each cell type

472  of liver dataset. In tSNE plot, each gene specific expressed in corresponding cell type

473  in red color. The shade of color means the expression value.

474  **Additional file 12.pdf: Figure S8.** The overlap of reported marker genes from the

475  corresponding publications in four datasets with selected genes. The circle is selected

476  genes and the square is not selected genes. The marker genes have different color and

477  non-marker genes are gray.

478  **Additional file 13.xlsx: Table S4.** The statistic of select state of reported marker genes

479  from the corresponding publications in four datasets with selected genes.

480  **Additional file 14.xlsx: Table S5.** The statistic of kappa coefficient in overall

481  performance test (Fig 2a-d).

482  **Additional file 15.pdf: Figure S9.** The heatmap of each cell type specific genes

483  expression in large dataset (top closest gene number 65). It shows an obvious pattern in

484  diagonal.

485    **Additional file 16.pdf: Figure S10.** Comparisons of gCAnno with Scmap-Cluster,

486    Scmap-Cell, scPred, Chetah and SVM on (a) HCC and ICCA reference dropout dataset,

487    (b) HCC and ICCA query dropout dataset, (c) AT root reference dropout dataset, (d)

488    AT root query dropout dataset and (e) imbalance dataset.

489    **Additional file 17.xlsx: Table S6.** The statistic of kappa coefficient in dropout test.

490    **Additional file 18.pdf: Figure S11.** An example of the existence of selected cell type

491    specific genes in liver ref dropout test dataset. The red color in more than one type

492    means these types shared this gene. With the increasing of dropout rate, the degree of

493    shared specific genes increased a little, but the specific pattern is still strong even in

494    dropout rate 50%.

495    **Additional file 19.xlsx: Table S7.** The statistic of kappa coefficient in imbalance test.

496    **Additional file 20.xlsx: Table S8.** The statistic of kappa coefficient in cross platform

497    test.

498    **Additional file 21.xlsx: Table S9.** Runtime statistic of gCAnno.

499    **Additional file 22.pdf: Figure S12.** The plot of building model time and graph scale.

500    The building model time is correlated with graph node number (correlation coefficient

501    is 0.94).