

1 Taming cell-to-cell heterogeneity in Acute Myeloid Leukaemia with 2 machine learning

3
4 Yara E. Sánchez-Corrales^{1*}, Ruben V.C. Pohle², Sergi Castellano^{1,3}, Alice Giustacchini^{2,*}

5 ¹Genetics and Genomic Medicine Department, Great Ormond Street Institute of Child Health,
6 University College London, London, UK.

7 ²Molecular and Cellular Immunology Section, Great Ormond Street Institute of Child Health,
8 University College London, London, UK.

9 ³UCL Genomics, Great Ormond Street Institute of Child Health, University College London,
10 London, UK.

11 * Correspondence:

12 Corresponding Authors

13 y.sanchez-corrales@ucl.ac.uk

14 a.giustacchini@ucl.ac.uk

15

16 **Keywords:** AML, Machine learning, classification, clustering, leukaemia

17

18 Abstract

19 Acute Myeloid Leukaemia (AML) is a phenotypically and genetically heterogeneous blood cancer
20 characterised by very poor prognosis, with disease relapse being the primary cause of treatment failure.

21

22 AML heterogeneity arises from different genetic and non-genetic sources, including its proposed
23 hierarchical structure, with leukemic stem cells (LSCs) and progenitors giving origin to a variety of
24 more mature leukemic subsets. Recent advances in single-cell molecular and phenotypic profiling have
25 highlighted the intra and inter-patient heterogeneous nature of AML, which has so far limited the
26 success of cell-based immunotherapy approaches against single targets.

27

28 Machine Learning (ML) can be uniquely used to find non-trivial patterns from high-dimensional
29 datasets and identify rare sub-populations. Here we review some recent ML tools that applied to single-
30 cell data could help disentangle cell heterogeneity in AML by identifying distinct core molecular
31 signatures of leukemic cell subsets. We discuss the advantages and limitations of unsupervised and
32 supervised ML approaches to cluster and classify cell populations in AML, for the identification of
33 biomarkers and the design of personalised therapies.

34

35 1 Introduction

36 AML is an aggressive and fast-progressing leukaemia characterised by the accumulation of myeloid
37 progenitors (Tenen, 2003). Although most patients achieve remission after first line chemotherapy and

Machine learning for taming cell heterogeneity in AML

38 haematopoietic stem cell transplantation, about 40% later relapse (Tsirigotis et al., 2016). Long-term
39 survival following relapse is below 20% with a median survival of 4-6 months, an outcome that has
40 not improved over the last two decades with conventional approaches (Tsirigotis et al., 2016; Medeiros,
41 2018; Lonetti et al., 2019) and novel therapies are therefore urgently needed (Lonetti et al., 2019).

42
43 AML is a molecularly heterogeneous group of diseases with a complex mutational landscape,
44 characterized by intra- and inter-patient variation (Figure 1A). Advances in next-generation sequencing
45 and single-cell technologies have revealed that AML cells display genetic and epigenetic heterogeneity
46 in different patients and even within the same patient multiple sub-clones co-exist, each carrying its
47 own hierarchical structure and possessing distinct immunophenotypes (Miles et al., 2020).

48
49 A non-genetic source of heterogeneity in AML is its proposed hierarchical structure, mimicking the
50 cellular hierarchy in normal hematopoietic development (Figure 1B). In healthy individuals, this
51 involves a stepwise differentiation process, with hematopoietic stem cells (HSCs) giving rise to
52 progressively more mature blood cells (Velten et al., 2017; Karamitros et al., 2018; Liggett & Sankaran,
53 2020). LSCs lie at the top of AML cellular hierarchies, and carry an unlimited ability to self-renew as
54 well as giving origin to a variety of more mature leukemic subsets (Tenen, 2003), each expressing
55 characteristic patterns of cell surface markers. LSCs can persist in a dormant state, making them
56 selectively unresponsive to conventional chemotherapies and allowing them to eventually fuel disease
57 relapse. For these reasons, the effective targeting of LSCs underpins any successful treatment for AML.

58
59 A promising approach is to target LSCs using immunotherapy with autologous T cells genetically
60 redirected to express Chimeric Antigen Receptors (CARs). In fact, CAR-T cells can effectively target
61 tumour cells irrespectively of their quiescent status. However, the lack of surface markers preferentially
62 expressed on LSCs as opposed to healthy HSCs has hindered the development of cell-based
63 immunotherapy strategies for AML, given the high risk of on-target off-tumour toxicity (Perna et al.,
64 2017; Lambie & Tasian, 2019). In addition, some of the targets tested so far (e.g. CD33 or CD123)
65 have heterogenous expression in the LSC compartment, with the risk of relapse due to their incomplete
66 targeting (Mardiana & Gill, 2020). Upon relapse, genetic and immunophenotypic heterogeneity in
67 AML LSCs further increases, complicating the discovery of ‘one fits all’ drug target (Ho et al., 2016).

68
69 As a result of AML’s heterogenous nature, CAR-T cell approaches against a single target are unlikely
70 to be effective, thus the design of combinations of CAR-T cells against multiple targets requires a
71 systematic characterization of the expression levels of surface antigens in AML cell populations at
72 single-cell resolution (Figure 1C) (Perna et al., 2017).

73
74 The unprecedented resolution achieved with single-cell technologies has enabled the dissection of cell
75 populations, including tumour and rare cell types that could not be identified using conventional bulk
76 sequencing (Giustacchini et al., 2017; Aldridge & Teichmann, 2020). In AML, the quantitative
77 phenotyping of leukemic cell profiles has allowed the identification of leukemic subsets without prior
78 knowledge of phenotypic markers for their prospective isolation, opening up new analytical challenges

79 for their clinical interpretation (Van Galen et al., 2019; Petti et al., 2019; Miles et al., 2020; Wu et al.,
80 2020; Velten et al., 2021; Triana et al., 2021).

81
82 Despite Machine Learning (ML) techniques having shown prognostic utility in classifying patients at
83 high risk of relapse and having been applied to risk-adapted treatments (review by Eckardt et al.
84 (2020)), they have only been recently applied to resolve heterogeneity in single-cell datasets from AML
85 patients (Van Galen et al., 2019; Triana et al., 2021). Fortunately, there has been an explosion of new
86 algorithms based on ML for the characterization of cell populations in single-cell datasets (Table 1)
87 that could be applied to identify molecular markers specific to AML subpopulations.

88
89 Here, we review some recent state-of-the-art ML methods with the potential to shed light into cell
90 heterogeneity in AML and identify biomarkers for specific cell populations in single-cell datasets.
91 Benchmarking of some recent methods has been done by Abdelaal et al. (2019) and Zhao et al. (2020).
92 Rather than an extensive discussion of algorithms, we provide a general overview of tools available to
93 identify cell populations in single-cell studies, highlighting ones that have the potential to reveal new
94 and rare cell types in AML and aid the design of personalised treatments.

95 **2 Machine learning for cell type identification in single-cell datasets and biomarker** 96 **discovery for personalized immunotherapy**

97 Single-cell high-throughput techniques, such as scRNA-seq, quantitatively characterise **cell types**
98 within a tissue (Trapnell, 2015). Typical workflows in single-cell transcriptional profiling include
99 dimensionality reduction and clustering of cells based on their gene expression patterns followed by
100 manual annotation of cell clusters from known cell type **markers** (Kolodziejczk et al., 2015). In the
101 context of AML and other cancers, transcriptionally similar malignant cells are expected to group
102 together, and can be unambiguously identified by the expression of certain feature genes that can be
103 used as biomarkers for designing personalised treatments.

104
105 The identification of cell types using typical workflows has several drawbacks: first, rare cell types are
106 easily missed and grouped together with some more prevalent ones; second, cell identity is often not
107 discrete but lies in a continuum (for instance, cells with mixed identities or in transition); and third, the
108 clustering can reflect other sources of variability unrelated to cell types (Kiselev et al., 2019). To
109 address these issues, ML tools have recently been developed allowing quantitative identification and
110 probabilistic assignment of cell types, thus aiding the identification of rare and heterogeneous cell
111 populations.

112
113 In general, ML approaches are either **unsupervised** or **supervised** (Figure 1D). The main difference
114 being the use of prior knowledge. Supervised methods are **trained** on an **annotated reference** with
115 known **classes** of cell types, whereas unsupervised models identify patterns in the data without prior
116 knowledge. A summary of recent methods is shown in Table 1.

117

118 **2.1 Recent ML unsupervised methods**

Machine learning for taming cell heterogeneity in AML

119 A common task for unsupervised methods is to use the intrinsic structure of the data to find clusters of
120 cells. The advantage of these approaches is that cells can be grouped in an automatic and unbiased
121 manner and thus, have the potential to discover unknown cell populations.

122

123 The popular single-cell processing packages Seurat (Butler et al., 2018) and Scanpy (Wolf et al., 2018)
124 use a graph-based clustering approach combined with modularity optimization to group
125 transcriptionally-similar cells together. Markers differentially expressed in each cluster can be found
126 using different methods, including logistic regression. The cell identity of each cluster is assigned
127 manually according to previous knowledge of cell-type specific markers. The main disadvantage of
128 this approach is that the number of clusters depends on a resolution parameter assigned by the user
129 (higher values will lead to a greater number of clusters) and thus, they may not faithfully reflect cell
130 types.

131

132 The recently developed Single-Cell Clustering Assessment Framework (SCCAF) (Miao et al, 2020)
133 generates an optimal number of clusters automatically. After the data has been clustered, SCCAF builds
134 an ML classifier (logistic regression) using part of the data (training). By applying this model to the
135 rest of the dataset (test), it iteratively merges clusters that appear indistinguishable to the ML classifier
136 to produce the final optimum clustering. The output of the model is a weighted list of feature genes
137 characteristic of every cluster that often include known markers for a given cell type and could
138 potentially be used to detect common biomarkers of leukemic cell subsets from AML patients.

139

140 Another unsupervised method, single-cell consensus clustering (SC3) uses the first 4-7% * N (number
141 of cells) **eigenvectors** to build multiple **k-means clustering** solutions (Kiselev et al., 2017). After
142 hierarchical grouping, the final clustering is driven by the combination of multiple clustering solutions.
143 The output is a list of marker genes that define each consensus cluster. While SC3 may not be the most
144 sensitive method to find rare populations (such as LSCs), SC3 was successful in identifying clusters of
145 prevalent genetic subclones with different mutations in myeloproliferative neoplasms (Kiselev et al.,
146 2017). A disadvantage of this method is that it does not scale well for datasets with more than 5,000
147 cells (Andrews et al., 2021).

148

149 A recent unsupervised method, weighted-nearest neighbour (WNN), was used to cluster cells using
150 multiple data modalities (e.g. surface proteins and transcriptomes) measured in the same cell (Hao et
151 al., 2020). This method uses **k-nearest neighbours** (kNN) to learn cell-specific modality “weights”.
152 When applied to a multiomics dataset generated from human bone marrow samples (Stuart et al., 2019),
153 it showed that the combination of surface proteins and gene expression was superior for identifying
154 cell populations than using one data modality alone. Multiomic single-cell technologies quantifying
155 both surface proteins and transcriptomes of individual cells (e.g. CITE-seq), could be ideally applied
156 to the identification of surface targets for the design of cell based immunotherapies (Stoeckius et al.,
157 2017).

158

159

160 Other unsupervised methods rely on Non-negative matrix factorization (NMF) methods (Kotliar et al.,
161 2019; Stein-O'Brien et al., 2019). These methods allow for the identification of cell types and,
162 simultaneously, **cell states**. Given the great transcriptional heterogeneity seen in AML even within
163 clonal populations carrying the same mutational patterns (Petti et al., 2019), it may be helpful to
164 consider cell identities and activities separately when clustering leukemic populations. Moreover, NMF
165 is potentially useful to identify LSC populations in AML, where the classical surface proteins defining
166 primitive cell types are present in highly similar patterns to healthy HSCs, but a 'malignant stem-like'
167 profile can still be identified (Levine et al., 2015).

168

169 **2.2 Recent ML supervised methods**

170 Supervised methods to classify cell types exploit previously identified cell types and use either known
171 marker genes or annotated reference datasets as an input to probabilistically assign new cells to a given
172 category.

173

174 Some methods take a list of markers for each cell type as input (Lee & Hemberg, 2019). For example,
175 CellAssign (Zhang et al., 2019) uses predefined cell types input as a marker gene list to build a
176 hierarchical model that produces a statistical classification of cells. This approach was used to delineate
177 the composition of the tumour microenvironment in serial samples (treatment and relapse) from
178 follicular lymphoma. Garnett (Pliner et al., 2019) also takes as input a list of markers. The format of
179 the input list permits accounting for cellular hierarchy (i.e, cell subtypes) and can include positive and
180 negative markers to define cell types (Pliner et al., 2019).

181

182 Other supervised methods use an annotated reference dataset to classify cell types but differ in the
183 features and the ML methods used to train models (see Table 1). For instance, SingleCellNet (Tan et
184 al., 2019) uses the most discriminative **gene pairs** (top pair transformation) to build a **random forest**
185 classifier while methods such as scPred (Alquicira-Hernandez et al., 2019) and Moana (Wagner &
186 Yanai, 2018) use principal components as features to fit a **support vector machine** (SVM). Some
187 methods rely on one or several similarity metrics (such as SingleR Aran et al. (2019)) and **k-nearest**
188 **neighbours** (kNN) to map query datasets into a known reference (e.g. scmap (Kiselev et al., 2018) and
189 scClassify (Lin et al., 2020)). Other methods use the training dataset to build an **Artificial Neural**
190 **Network** (ANN) model such as SuperCT (Xie et al., 2019) and ACTINN (Ma & Pellegrini, 2019) with
191 an input layer containing as many nodes as the number of genes in the training set and an output layer
192 with nodes equal to the number of cell types. Interestingly, both ANN methods provide pre-trained
193 models that could be used to classify new AML datasets.

194

195 An advantage of supervised ML approaches is that cell types are assigned probabilistically and some
196 approaches allow for the possibility of an "unassigned" category (Kiselev et al., 2018, Zhang et al.,
197 2019, Tan et al., 2019, Pliner et al., 2019, Ma & Pellegrini, 2019). The unassigned label for cells that
198 are absent or are very different in the reference dataset is key to limit misclassification and to allow the
199 discovery of new cell types.

200

201 Algorithms such as CHETAH (de Kanter et al., 2019) and scClassify (Lin et al., 2020) allow for
202 intermediate categories that can highlight populations with a mixture of identities as previously
203 reported in AML (Smith et al., 1983). These methods are based on hierarchical correlation trees to
204 classify test datasets (de Kanter et al., 2019, Lin et al., 2020).

205

206 As more annotated single-cell datasets become available, the primary advantage of supervised methods
207 is leveraging previous knowledge. Reference datasets of human bone marrow cells from healthy
208 individuals are available from resources such as the Human Cell Atlas (Regev et al., 2017). Distinct
209 cell populations or patient-specific tumour clones could be identified as unknown (because they are
210 very different or absent in the reference data sets). As AML single-cell datasets become more abundant,
211 they can be integrated with healthy single or multimodal references using ML methods (Hao et al.,
212 2020).

213

214 A disadvantage of supervised methods is that they rely on known markers or accurate cell type
215 annotations to build classification models. Often, markers for rare cell populations, such as LSCs, are
216 unknown, not robust (Pollyea & Jordan, 2017) or can be expressed by more than one cell type (Van
217 Galen et al., 2019). Further, in many cases, annotation of single-cell datasets requires additional
218 standardisation (de Kanter et al., 2019).

219

220 **3 Discussion**

221

222 ML techniques are able to find non-trivial patterns in high-dimensional data (Geron, 2019). In fact,
223 ML has already proven useful in identifying markers in bulk studies in prospectively isolated leukemic
224 sub-populations (Ng et al., 2016; Li et al., 2020). However, ML has not reached its full potential for
225 the characterisation of AML cell populations at single-cell resolution, partly due to the recent
226 development of large datasets (Van Galen et al., 2019; Petti et al., 2019; Miles et al., 2020; Triana et
227 al., 2021; Velten et al., 2021).

228

229 Here we have reviewed tools to aid biomarker discovery using ML at single-cell level resolution. Many
230 ML models explicitly quantify the contribution of individual features (genes) for a given classification.
231 Importantly, genes identified in microarray data as important for classifying samples into “AML” or
232 “no-AML” were not always differentially expressed (Warnat-Herresthal et al., 2020). This means that
233 traditional differential expression analysis could fail to identify biomarkers that are good predictors for
234 assigning a given group of cells (Alquicira-Hernandez et al., 2019). Thus, ML algorithms can find
235 biomarkers that otherwise will be missed, expediting the design of suitable target combinations for
236 immunotherapy.

237

238 Recently, it was shown that single-cell transcriptomics is capable of dissecting genetic subclones in
239 AML, such as $GATA2^{R361C}$, which cluster separately from normal hematopoietic cell types (Petti et
240 al., 2019). This observation suggests that subclonal diversity in AML could be associated with distinct
241 gene expression profiles which ML techniques can leverage to identify mutated populations. Some
242 AML mutations create subtle differences in expression profiles (Van Galen et al., 2019; Petti et al.,

243 2019; Velten et al., 2021) and isolating these populations represents an analytical challenge
244 contemporary ML methods could address.

245

246 Moreover, recent experimental innovations allowing for the simultaneous quantitative assessment of
247 cellular and molecular information at single-cell resolution promise to better dissect cell heterogeneity
248 in AML. Particularly important is the ability to detect mutations in single cells combined with their
249 transcriptional profiling, offering an unprecedented opportunity to identify specific leukemic cell
250 populations (Giustacchini et al., 2017; Rodriguez-Meira et al., 2019; Van Galen et al., 2019; Petti et
251 al., 2019; Ludwig et al., 2019; Velten et al., 2021). For instance, the combination of single-cell
252 transcriptomics and mutational profiles allowed the distinction of pre-leukemic clones, LSC and
253 healthy HSC (Velten et al., 2021). ML such as SVM could be used next to identify molecules that
254 maximise this classification as done before for bulk RNA-seq and microarray data (Li et al., 2020).

255

256 In addition, the identification of mutant and non-mutant cells allows for applying ML methods to both
257 all and only mutated cells to further characterise subpopulations (Petti et al., 2019), and can be used to
258 fine-tune ML classification algorithms. For instance, a two-step ML classification strategy was applied
259 to bone marrow samples of AML patients (Van Galen et al., 2019). First, a fraction of mutant cells was
260 identified by genotyping and these were classified into one of six normal haematopoietic cell types
261 (monocyte-like, progenitor-like, etc.). Subsequently, these malignant cell types were incorporated as
262 additional classes in a second classifier that successfully identified mutant and normal cells from their
263 transcriptome profiles.

264

265 The simultaneous characterization of surface proteins at single-cell resolution (Stoeckius et al., 2017)
266 is especially important for isolation of heterogeneous cell populations. There are some analytical
267 challenges with the integration of multiple data modalities (Efremova & Teichmann, 2020), but
268 combining different data types from the same cell has already shown to improve cell population
269 identification in AML datasets (Petti et al., 2019; Triana et al., 2021) and healthy bone marrow samples
270 (Hao et al., 2020), thus we anticipate that multimodal datasets will improve the performance of ML
271 models in isolating specific cell populations and may facilitate the identification of relevant surface
272 targets for precision immunotherapy.

273

274 All the methods reviewed here will incur a certain degree of **underfitting** and **overfitting**. Thus, it is
275 wise to compare algorithms in the initial cell composition assessment. Some, such as hierarchical
276 methods, are potentially more suitable for AML samples, where there is an intrinsic hierarchy shared
277 with normal hematopoietic development (Figure 1B). Also, methods that enable the recognition of
278 intermediate cell types, mixed identities or different cell states would be more suitable for the
279 identification of abnormally differentiated leukemic cells, known to be characteristic of AML (Smith
280 et al., 1983).

281

282 Finally, we anticipate that single-cell resolution phenotyping will be important for the design of cell-
283 based immunotherapy combinatorial strategies accounting for clonality and differentiation states of

284 AML populations, with ML likely playing a pivotal role in the selection of optimal therapeutic targets
285 for the design of personalised workflows tailored to each patient.
286

287 **4 Conflict of Interest**

288 The authors declare that the research was conducted in the absence of any commercial or financial
289 relationships that could be construed as a potential conflict of interest.

290 **5 Author Contributions**

291 All authors contributed to the conception and editing of this review and approved the final manuscript.
292 YSC conducted literature review and wrote the manuscript in consultation with RP. SG and AG
293 critically revised the work.

294 **6 Funding**

295 This work was supported by the NIHR GOSH BRC.

296 **7 Acknowledgments**

297 This work was supported by the NIHR GOSH BRC, the views expressed are those of the authors and
298 not necessarily those of the NHS, the NIHR or the Department of Health. We thank George Hall for
299 helpful feedback to the manuscript.
300

301 **8 References**

- 302 Abdelaal, T., Michielsen, L., Cats, D., Hoogduin, D., Mei, H., Reinders, M. J. T., et al. (2019).
303 A comparison of automatic cell identification methods for single-cell RNA sequencing
304 data. *Genome Biology*, 20(1). doi:10.1186/s13059-019-1795-z
- 305 Aldridge, S., and Teichmann, S. A. (2020). Single cell transcriptomics comes of age. *Nature*
306 *Communications*, 11(1). doi:10.1038/s41467-020-18158-5
- 307 Alquicira-Hernandez, J., Sathe, A., Ji, H. P., Nguyen, Q., and Powell, J. E. (2019). scPred:
308 accurate supervised method for cell-type classification from single-cell RNA-seq data.
309 *Genome Biology*, 20(1). doi:10.1186/s13059-019-1862-5
- 310 Andrews, T. S., Kiselev, V. Y., McCarthy, D., and Hemberg, M. (2021). Tutorial: guidelines for
311 the computational analysis of single-cell RNA sequencing data. *Nature Protocols*,
312 16(1), 1-9. doi:10.1038/s41596-020-00409-w
- 313 Aran, D., Looney, A. P., Liu, L., Wu, E., Fong, V., Hsu, A., et al. (2019). Reference-based
314 analysis of lung single-cell sequencing reveals a transitional profibrotic macrophage.
315 *Nature Immunology*, 20(2), 163-172. doi:10.1038/s41590-018-0276-y
- 316 Butler, A., Hoffman, P., Smibert, P., Papalexi, E., and Satija, R. (2018). Integrating single-cell
317 transcriptomic data across different conditions, technologies, and species. *Nature*
318 *Biotechnology*, 36(5), 411-420. doi:10.1038/nbt.4096
- 319 de Kanter, J., Lijnzaad, P., Candelli, T., Margaritis, T., and Holstege, F. C. P. (2019). CHETAH:
320 a selective, hierarchical cell type identification method for single-cell RNA sequencing.
321 *Nucleic Acids Research*, 47(16), e95-e95. doi:10.1093/nar/gkz543

- 322 Eckardt, J.-N., Bornhäuser, M., Wendt, K., and Middeke, J. M. (2020). Application of machine
323 learning in the management of acute myeloid leukemia: current practice and future
324 prospects. *Blood Advances*, 4(23), 6077-6085.
325 doi:10.1182/bloodadvances.2020002997
- 326 Efremova, M., and Teichmann, S. A. (2020). Computational methods for single-cell omics
327 across modalities. *Nature Methods*, 17(1), 14-17. doi:10.1038/s41592-019-0692-4
- 328 Geron, A. (2019). *Hands-on Machine Learning with Scikit-Learn, Keras, and TensorFlow:
329 Concepts, Tools, and Techniques to Build Intelligent Systems* (Second edition ed.):
330 O'Reilly.
- 331 Giustacchini, A., Thongjuea, S., Barkas, N., Woll, P. S., Povinelli, B. J., Booth, C. A. G., et al.
332 (2017). Single-cell transcriptomics uncovers distinct molecular signatures of stem cells
333 in chronic myeloid leukemia. *Nature Medicine*, 23(6), 692-702. doi:10.1038/nm.4336
- 334 Hao, Y., Hao, S., Andersen-Nissen, E., Mauck, W. M., Zheng, S., Butler, A., et al. (2020).
335 *Integrated analysis of multimodal single-cell data*. Cold Spring Harbor Laboratory.
336 Retrieved from <https://dx.doi.org/10.1101/2020.10.12.335331>
- 337 Ho, T.-C., Lamere, M., Stevens, B. M., Ashton, J. M., Myers, J. R., O'Dwyer, K. M., et al.
338 (2016). Evolution of acute myelogenous leukemia stem cell properties after treatment
339 and progression. *Blood*, 128(13), 1671-1678. doi:10.1182/blood-2016-02-695312
- 340 Karamitros, D., Stoilova, B., Aboukhalil, Z., Hamey, F., Reinisch, A., Samitsch, M., et al.
341 (2018). Single-cell analysis reveals the continuum of human lympho-myeloid
342 progenitor cells. *Nature Immunology*, 19(1), 85-97. doi:10.1038/s41590-017-0001-2
- 343 Kiselev, V. Y., Andrews, T. S., and Hemberg, M. (2019). Challenges in unsupervised clustering
344 of single-cell RNA-seq data. *Nature Reviews Genetics*, 20(5), 273-282.
345 doi:10.1038/s41576-018-0088-9
- 346 Kiselev, V. Y., Kirschner, K., Schaub, M. T., Andrews, T., Yiu, A., Chandra, T., et al. (2017).
347 SC3: consensus clustering of single-cell RNA-seq data. *Nature Methods*, 14(5), 483-
348 486. doi:10.1038/nmeth.4236
- 349 Kiselev, V. Y., Yiu, A., and Hemberg, M. (2018). scmap: projection of single-cell RNA-seq data
350 across data sets. *Nature Methods*, 15(5), 359-362. doi:10.1038/nmeth.4644
- 351 Kolodziejczk, A. A., Kim, J. K., Svensson, V., Marioni, J. C., and Teichmann, S. A. (2015). The
352 Technology and Biology of Single-Cell RNA Sequencing. *Molecular Cell*, 58(4), 610-
353 620. doi:10.1016/j.molcel.2015.04.005
- 354 Kotliar, D., Veres, A., Nagy, M. A., Tabrizi, S., Hodis, E., Melton, D. A., et al. (2019). Identifying
355 gene expression programs of cell-type identity and cellular activity with single-cell
356 RNA-Seq. *eLife*, 8. doi:10.7554/elife.43803
- 357 Lamble, A. J., and Tasian, S. K. (2019). Opportunities for immunotherapy in childhood acute
358 myeloid leukemia. *Blood Advances*, 3(22), 3750-3758.
359 doi:10.1182/bloodadvances.2019000357
- 360 Lee, J. T. H., and Hemberg, M. (2019). Supervised clustering for single-cell analysis. *Nature
361 Methods*, 16(10), 965-966. doi:10.1038/s41592-019-0534-4
- 362 Levine, J. H., Simonds, E. F., Bendall, S. C., Davis, K. L., Amir, E.-a. D., Tadmor, M. D., et
363 al. (2015). Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells
364 that Correlate with Prognosis. *Cell*, 162(1), 184-197. doi:10.1016/j.cell.2015.05.047
- 365 Li, J., Lu, L., Zhang, Y.-H., Xu, Y., Liu, M., Feng, K., et al. (2020). Identification of leukemia
366 stem cell expression signatures through Monte Carlo feature selection strategy and
367 support vector machine. *Cancer Gene Therapy*, 27(1-2), 56-69. doi:10.1038/s41417-
368 019-0105-y

- 369 Liggett, L. A., and Sankaran, V. G. (2020). Unraveling Hematopoiesis through the Lens of
370 Genomics. *Cell*, 182(6), 1384-1400. doi:10.1016/j.cell.2020.08.030
- 371 Lin, Y., Cao, Y., Kim, H. J., Salim, A., Speed, T. P., Lin, D. M., et al. (2020). scClassify: sample
372 size estimation and multiscale classification of cells using single and multiple
373 reference. *Molecular Systems Biology*, 16(6). doi:10.15252/msb.20199389
- 374 Lonetti, A., Pession, A., and Masetti, R. (2019). Targeted Therapies for Pediatric AML: Gaps
375 and Perspective. *Frontiers in Pediatrics*, 7. doi:10.3389/fped.2019.00463
- 376 Ludwig, L. S., Lareau, C. A., Ulirsch, J. C., Christian, E., Muus, C., Li, L. H., et al. (2019).
377 Lineage Tracing in Humans Enabled by Mitochondrial Mutations and Single-Cell
378 Genomics. *Cell*, 176(6), 1325-1339.e1322. doi:10.1016/j.cell.2019.01.022
- 379 Ma, F., and Pellegrini, M. (2019). ACTINN: automated identification of cell types in single cell
380 RNA sequencing. *Bioinformatics*. doi:10.1093/bioinformatics/btz592
- 381 Mardiana, S., and Gill, S. (2020). CAR T Cells for Acute Myeloid Leukemia: State of the Art
382 and Future Directions. *Frontiers in Oncology*, 10. doi:10.3389/fonc.2020.00697
- 383 Medeiros, B. C. (2018). Is there a standard of care for relapsed AML? *Best Practice &*
384 *Research Clinical Haematology*, 31(4), 384-386. doi:10.1016/j.beha.2018.09.006
- 385 Miao, Z., Moreno, P., Huang, N., Papatheodorou, I., Brazma, A., and Teichmann, S. A. (2020).
386 Putative cell type discovery from single-cell gene expression data. *Nature Methods*,
387 17(6), 621-628. doi:10.1038/s41592-020-0825-9
- 388 Miles, L. A., Bowman, R. L., Merlinsky, T. R., Csete, I. S., Ooi, A. T., Durruthy-Durruthy, R.,
389 et al. (2020). Single-cell mutation analysis of clonal evolution in myeloid malignancies.
390 *Nature*, 587(7834), 477-482. doi:10.1038/s41586-020-2864-x
- 391 Ng, S. W. K., Mitchell, A., Kennedy, J. A., Chen, W. C., McLeod, J., Ibrahimova, N., et al.
392 (2016). A 17-gene stemness score for rapid determination of risk in acute leukaemia.
393 *Nature*, 540(7633), 433-437. doi:10.1038/nature20598
- 394 Perna, F., Berman, S. H., Soni, R. K., Mansilla-Soto, J., Eyquem, J., Hamieh, M., et al. (2017).
395 Integrating Proteomics and Transcriptomics for Systematic Combinatorial Chimeric
396 Antigen Receptor Therapy of AML. *Cancer Cell*, 32(4), 506-519.e505.
397 doi:10.1016/j.ccell.2017.09.004
- 398 Petti, A. A., Williams, S. R., Miller, C. A., Fiddes, I. T., Srivatsan, S. N., Chen, D. Y., et al.
399 (2019). A general approach for detecting expressed mutations in AML cells using
400 single cell RNA-sequencing. *Nature Communications*, 10(1). doi:10.1038/s41467-019-
401 11591-1
- 402 Pliner, H. A., Shendure, J., and Trapnell, C. (2019). Supervised classification enables rapid
403 annotation of cell atlases. *Nature Methods*, 16(10), 983-986. doi:10.1038/s41592-019-
404 0535-3
- 405 Pollyea, D. A., and Jordan, C. T. (2017). Therapeutic targeting of acute myeloid leukemia stem
406 cells. *Blood*, 129(12), 1627-1635. doi:10.1182/blood-2016-10-696039
- 407 Regev, A., Teichmann, S. A., Lander, E. S., Amit, I., Benoist, C., Birney, E., et al. (2017). The
408 Human Cell Atlas. *eLife*, 6. doi:10.7554/elife.27041
- 409 Rodriguez-Meira, A., Buck, G., Clark, S.-A., Povinelli, B. J., Alcolea, V., Louka, E., et al.
410 (2019). Unravelling Intratumoral Heterogeneity through High-Sensitivity Single-Cell
411 Mutational Analysis and Parallel RNA Sequencing. *Molecular Cell*, 73(6), 1292-
412 1305.e1298. doi:10.1016/j.molcel.2019.01.009
- 413 Smith, L., Curtis, J., Messner, H., Senn, J., Furthmayr, H., and McCulloch, E. (1983). Lineage
414 infidelity in acute leukemia. *Blood*, 61 (66): 1138-1145.
415 doi:<https://doi.org/10.1182/blood.V61.6.1138.1138>
- 416 Stein-O'Brien, G. L., Clark, B. S., Sherman, T., Zibetti, C., Hu, Q., Sealfon, R., et al. (2019).
417 Decomposing Cell Identity for Transfer Learning across Cellular Measurements,

- 418 Platforms, Tissues, and Species. *Cell Systems*, 8(5), 395-411.e398.
419 doi:10.1016/j.cels.2019.04.004
- 420 Stoeckius, M., Hafemeister, C., Stephenson, W., Houck-Loomis, B., Chattopadhyay, P. K.,
421 Swerdlow, H., et al. (2017). Simultaneous epitope and transcriptome measurement in
422 single cells. *Nature Methods*, 14(9), 865-868. doi:10.1038/nmeth.4380
- 423 Stuart, T., Butler, A., Hoffman, P., Hafemeister, C., Papalexi, E., Mauck, W. M., et al. (2019).
424 Comprehensive Integration of Single-Cell Data. *Cell*, 177(7), 1888-1902.e1821.
425 doi:10.1016/j.cell.2019.05.031
- 426 Tan, Y., and Cahan, P. (2019). SingleCellNet: A Computational Tool to Classify Single Cell
427 RNA-Seq Data Across Platforms and Across Species. *Cell Systems*, 9(2), 207-
428 213.e202. doi:10.1016/j.cels.2019.06.004
- 429 Tenen, D. G. (2003). Disruption of differentiation in human cancer: AML shows the way.
430 *Nature Reviews Cancer*, 3(2), 89-101. doi:10.1038/nrc989
- 431 Trapnell, C. (2015). Defining cell types and states with single-cell genomics. *Genome*
432 *Research*, 25(10), 1491-1498. doi:10.1101/gr.190595.115
- 433 Triana, S. H., Vonficht, D., Jopp-Saile, L., Raffel, S., Lutz, R., Leonce, D., et al. (2021). *Single-*
434 *cell proteo-genomic reference maps of the hematopoietic system enable the*
435 *purification and massive profiling of precisely defined cell states*. Cold Spring Harbor
436 Laboratory. Retrieved from <https://dx.doi.org/10.1101/2021.03.18.435922>
- 437 Tsigotis, P., Byrne, M., Schmid, C., Baron, F., Ciceri, F., Esteve, J., et al. (2016). Relapse
438 of AML after hematopoietic stem cell transplantation: methods of monitoring and
439 preventive strategies. A review from the ALWP of the EBMT. *Bone Marrow*
440 *Transplantation*, 51(11), 1431-1438. doi:10.1038/bmt.2016.167
- 441 Van Galen, P., Hovestadt, V., Wadsworth li, M. H., Hughes, T. K., Griffin, G. K., Battaglia, S.,
442 et al. (2019). Single-Cell RNA-Seq Reveals AML Hierarchies Relevant to Disease
443 Progression and Immunity. *Cell*, 176(6), 1265-1281.e1224.
444 doi:10.1016/j.cell.2019.01.031
- 445 Velten, L., Haas, S. F., Raffel, S., Blaszkiewicz, S., Islam, S., Hennig, B. P., et al. (2017).
446 Human haematopoietic stem cell lineage commitment is a continuous process. *Nature*
447 *Cell Biology*, 19(4), 271-281. doi:10.1038/ncb3493
- 448 Velten, L., Story, B. A., Hernández-Malmierca, P., Raffel, S., Leonce, D. R., Milbank, J., et al.
449 (2021). Identification of leukemic and pre-leukemic stem cells by clonal tracking from
450 single-cell transcriptomics. *Nature Communications*, 12(1). doi:10.1038/s41467-021-
451 21650-1
- 452 Wagner, F., and Yanai, I. (2018). *Moana: A robust and scalable cell type classification*
453 *framework for single-cell RNA-Seq data*. Cold Spring Harbor Laboratory. Retrieved
454 from <https://dx.doi.org/10.1101/456129>
- 455 Warnat-Herresthal, S., Perrakis, K., Taschler, B., Becker, M., Baßler, K., Beyer, M., et al.
456 (2020). Scalable Prediction of Acute Myeloid Leukemia Using High-Dimensional
457 Machine Learning and Blood Transcriptomics. *iScience*, 23(1), 100780.
458 doi:10.1016/j.isci.2019.100780
- 459 Wolf, F. A., Angerer, P., and Theis, F. J. (2018). SCANPY: large-scale single-cell gene
460 expression data analysis. *Genome Biology*, 19(1). doi:10.1186/s13059-017-1382-0
- 461 Wu, J., Xiao, Y., Sun, J., Sun, H., Chen, H., Zhu, Y., et al. (2020). A single-cell survey of
462 cellular hierarchy in acute myeloid leukemia. *Journal of Hematology & Oncology*,
463 13(1). doi:10.1186/s13045-020-00941-y
- 464 Xie, P., Gao, M., Wang, C., Zhang, J., Noel, P., Yang, C., et al. (2019). SuperCT: a
465 supervised-learning framework for enhanced characterization of single-cell

466 transcriptomic profiles. *Nucleic Acids Research*, 47(8), e48-e48.
 467 doi:10.1093/nar/gkz116
 468 Zhang, A. W., O’Flanagan, C., Chavez, E. A., Lim, J. L. P., Ceglia, N., McPherson, A., et al.
 469 (2019). Probabilistic cell-type assignment of single-cell RNA-seq for tumor
 470 microenvironment profiling. *Nature Methods*, 16(10), 1007-1015. doi:10.1038/s41592-
 471 019-0529-1
 472 Zhao, X., Wu, S., Fang, N., Sun, X., and Fan, J. (2020). Evaluation of single-cell classifiers for
 473 single-cell RNA sequencing data sets. *Briefings in Bioinformatics*, 21(5), 1581-1595.
 474 doi:10.1093/bib/bbz096
 475
 476

477 **9 Figure legend**

478 **Figure 1. The high cell-to-cell heterogeneity in AML tumours can be dissected using machine**
 479 **learning methods.** **A)** The schematic representing clonal diversity in two putative AML patients
 480 highlights the complex intra and inter-patient variation of cell diversity (schematics adapted from Petti
 481 et al., 2019). Importantly, each clone carries its own hierarchical structure (here shown for one clone
 482 as an example). **B)** Leukemic populations share the hierarchical organization of normal hematopoietic
 483 development, where hematopoietic stem cells (HSCs) differentiate into multiple cell lineages, giving
 484 rise to all mature blood cells (blue lineages). Genetic mutations induce malignant transformation and
 485 give rise to leukemic stem cells (LSCs) that share some characteristics of their normal counterparts
 486 such as unlimited ability to self-renew and the potential to give origin to a variety of more mature
 487 leukemic subsets (red lineages). **C)** Ideal targets for immunotherapy with engineered T cells are those
 488 present in both leukemic blast and LSC cells and absent in healthy cell types. Targets that are
 489 ubiquitously expressed will fail to target specific leukemic populations and will be toxic for normal
 490 cells (on target off, tumour toxicity). Targets that are absent from LSC will render the treatment prone
 491 to relapse. Due to the high cell heterogeneity in AML more than one molecule is likely to fulfil these
 492 requirements. **D)** Machine learning methods to identify cell populations can be unsupervised and
 493 supervised. The former uses the intrinsic structure of the data to cluster cells in an automatic fashion.
 494 The second uses a predefined set of groups to classify unknown cells, leveraging previous knowledge.
 495

496 **10 Tables**

497 **Table 1. Summary of recent ML-based methods to identify cell types.**

Algorithm name	Classification type	Method	Input data	Important contribution	Reference
SC3	Unsupervised	Consensus clustering and hierarchical clustering	Normalised expression matrix	Transcriptome-based identification of genetic subclones in myeloproliferative neoplasms	Kiselev et al. (2017)

Machine learning for taming cell heterogeneity in AML

cNMF	Unsupervised	Non-negative matrix factorization	Expression matrix and several parameters	Identification of previously misclassified immature skeletal muscle cells in a published dataset from brain organoids	Kotliar et al. (2019)
scCOGAPS	Unsupervised	Non-negative matrix factorization	Normalised and log-scaled expression matrix	Identification of gene expression signatures characteristic of discrete cell types in the developing retina	Stein-O'Brien et al. (2019)
SCCAF	Unsupervised	Logistic Regression and self-projection	Expression matrix and several parameters	Identification of cell states associated with different stages of erythroid maturation in mouse	Miao et al. (2020)
WWN	Unsupervised	K-nearest neighbours and Jaccard distance	Expression matrix and protein matrix (or any other single-cell measurement)	Single-cell multimodal analysis improves resolution of cell states in the immune system and identify previously unreported subpopulations	Hao et al. (2020)
CellAssign	Supervised	Expectation-Maximization hierarchical model	List of cell markers, subset of expression matrix containing the marker	Resolution of malignant and non-malignant cells and their molecular dynamics during disease progression	Zhang et al. (2019)

Machine learning for taming cell heterogeneity in AML

			genes and some parameters	in follicular lymphoma	
Garnett	Supervised	Multinomial elastic-net regression	Hierarchical list of cell markers (positive and negative) and expression matrix	The model trained on a mouse lung dataset is successfully applied to detect both healthy cell types and tumor cells in a human lung cancer dataset	Pliner et al. (2019)
scmap	Supervised	k-means (scmap-cluster) and k-nearest-neighbour (scmap-cell)	Annotated reference dataset and query expression matrix	Cell types in a test datasets are annotated with high accuracy irrespectively of batch effect	Kiselev et al. (2018)
CHETAH	Supervised	Hierarchical Spearman correlation	Annotated reference dataset and query expression matrix (both normalised and log – scaled)	The cell type identification algorithm correctly identifies cancer cells absent in the reference dataset as “unassigned” or “intermediate”	de Kanter et al. (2019)
scClassify	Supervised	Hierarchical ordered partitioning, ensemble learning and weighted k-nearest-neighbour	Annotated reference dataset and query expression matrix (both log – transformed)	Identification of cell types from the Tabula Muris single cell dataset that were unidentified in the original publication, including very rare populations	Lin et al. (2020)

Machine learning for taming cell heterogeneity in AML

SingleR	Supervised	Correlation to training set	Annotated reference dataset and query expression matrix (both normalised and log-transformed)	Identification of a subgroup of macrophages whose molecular markers are upregulated in samples from patients with idiopathic pulmonary fibrosis.	Aran et al. (2019)
SingleCellNet	Supervised	Random Forest	Annotated reference dataset and expression matrix (both raw)	Cells from pancreatic tissue that were “unclassified” in the original study are identified as Schwann cells and gamma cells	Tan, and Cahan (2019)
SuperCT	Supervised	Artificial Neural Network	Pre-trained ANN model and a query expression matrix	The model predicts cell types with high accuracy in multiple single cell test datasets including cord blood mononuclear cells and mouse pancreatic cancer.	Xie et al. (2019)
ACTINN	Supervised	Artificial Neural Network	Annotated reference dataset and query expression matrix	Model trained on a T cell subtype reference accurately predicts T cell subtypes from an independent peripheral blood mononuclear cells dataset	Ma, and Pellegrini (2019)

Machine learning for taming cell heterogeneity in AML

Moana	Supervised	Support Vector Machine	Pre-trained model and raw query expression matrix	Identification of common and cell type-specific gene expression responses to IFN- β treatment in peripheral blood cells	Wagner, and Yanai (2018)
scPred	Supervised	Support Vector Machine	Annotated reference dataset and query expression matrix (both normalised)	Prediction of pathological cell states in gastric and colorectal cancer	Alquicira-Hernandez et al. (2019)

498

499

500 **11 Supplementary Material**

501 A glossary is included as a Supplementary Material.

502