**Ensemblex: an accuracy-weighted ensemble genetic demultiplexing framework for population-scale scRNAseq sample pooling**

**Authors:** Michael R. Fiorini[1,2], michael.fiorini@mail.mcgill.ca; Saeid Amiri[2] saeid.amiri@mcgill.ca; Allison A. Dilliott[2,3], allison.dilliot@mcgill.ca; Cristine M. Yde Ohki[4], cristinemarie.ydeohki@uzh.ch; Lukasz Smigielski[4], lukasz.smigielski@kjpd.uzh.ch; Susanne Walitza[4,5,6], susanne.walitza@pukzh.ch; Edward A. Fon[2,3], ted.fon@mcgill.ca; Edna Grünblatt[4,5,6], edna.gruenblatt@kjpd.uzh.ch; Rhalena A. Thomas[2,3*], rhalena.thomas@mcgill.ca; Sali M.K. Farhan[1,2,3*], sali.farhan@mcgill.ca

1. Department of Human Genetics, McGill University, Montreal, Quebec H3A 2B4, Canada

2. The Montreal Neurological Institute-Hospital, McGill University, Montreal, Quebec H3A 2B4, Canada

3. Department of Neurology and Neurosurgery, McGill University, Montreal, Quebec H3A 2B4, Canada

4. Department of Child and Adolescent Psychiatry and Psychotherapy, Psychiatric University Hospital Zurich, University of Zurich, Zurich, Switzerland

5. Neuroscience Center Zurich, University of Zurich and the ETH Zurich, Zurich, Switzerland

6. Zurich Center for Integrative Human Physiology, University of Zurich, Switzerland

* Shared co-senior authorship

**Corresponding authors:** Sali M.K. Farhan, Sali.farhan@mcgill.ca; Rhalena A. Thomas, rhalena.thomas@mcgill.ca

**Abstract**

Multiplexing samples from distinct individuals prior to sequencing is a promising step toward achieving population-scale single-cell RNA sequencing by reducing the restrictive costs of the technology. Individual genetic demultiplexing tools resolve the donor-of-origin identity of pooled cells using natural genetic variation but present diminished accuracy on highly multiplexed experiments, impeding the analytic potential of the dataset. In response, we introduce Ensemblex: an accuracy-weighted, ensemble genetic demultiplexing framework that integrates four distinct algorithms to identify the most probable subject labels. Using computationally and experimentally pooled samples, we demonstrate Ensemblex's superior accuracy and illustrate the implications of robust demultiplexing on biological analyses.

**Keywords:** single-cell RNA sequencing, multiplexing, sample pooling, genetic demultiplexing, induced pluripotent stem cells, differential gene expression, dopaminergic neurons, doublet detection, accuracy-weighted probability, high-throughput sequencing

**Background**

Single-cell RNA sequencing (scRNAseq) continues to revolutionize our molecular understanding of biology by providing unprecedented insight into the transcriptional landscape of individual cells. Unlike bulk RNAseq, where the RNA from all cells within a tissue is sequenced to produce total expressional profiles across all cells, scRNAseq captures transcriptional signatures at a single-cell resolution, elucidating the diverse gene expression across distinct cell types and subtypes. Differential gene expression (DGE) can then be calculated between subgroups of cells to reveal cell type-specific expression changes between patient or treatment groups. However, scRNAseq has come at the expense of increased costs, hindering its application for population-scale analyses, which are critical for deriving clinico-pathological associations and characterizing the genetic heterogeneity of complex diseases in biomedical sciences (1, 2).

In addition to the expense of separately capturing and sequencing cells from individual donors, the costs of scRNAseq are exacerbated for cell cultures, such as those derived from induced pluripotent stem cells (iPSC) (1). In particular, neurological diseases are difficult to study in human tissue because access to post-mortem brains is limited and experimental manipulations are not possible; in contrast, iPSC-derived cultures of neurons and other brain cells grown from reprogrammed skin or blood cells of human donors are an excellent model of the brain (3). However, iPSCs from each donor must be individually plated and differentiated in parallel, presenting prohibitively high consumable and labour costs that render the methodology unfeasible for population-scale analyses. Multiplexing cultures by pooling cells from multiple donors prior to growth and differentiation, droplet capture, and sequencing, is one solution to address this limitation as it reduces costs by a factor of the number of samples multiplexed (4). Similarly,

73  samples such as tumor biopsies can be pooled at acquisition to realize the same benefits. In turn,

74  genetic demultiplexing tools are cost-effective, statistical frameworks that use the natural genetic

75  variation at sites of single-nucleotide polymorphisms (SNP) observed in the transcriptome to

76  cluster cells on the basis of their donor's genotype. Importantly, genetic demultiplexing can be

77  informed by prior genotype information of the donors to improve demultiplexing accuracy and

78  facilitate the assignment of each cell back to its specific donor-of-origin, which is critical for

79  downstream analyses aiming to investigate discrepancies between subjects. At present, six genetic

80  demultiplexing tools have been developed for scRNAseq: Demuxalot (5) and Demuxlet (6) both

81  require prior genotype information as input; Freemuxlet (6) relies entirely on the de novo

82  transcriptome and does not incorporate prior genotype information; and ScSplit (7), Souporcell

83  (8), and Vireo (9) provide versions of the algorithm that can work with and without prior genotype

84  information **(Table 1).**

85

86  A robust genetic demultiplexing tool is tasked with mitigating the addition of technical artifacts

87  into scRNAseq datasets by correctly classifying each pooled cell to its donor-of-origin, correctly

88  identifying heterogenic doublets (erroneous barcodes composed of two or more cells from distinct

89  subjects), and quantifying its confidence in the demultiplexed labels so that low-confidence

90  classifications can be eliminated from downstream analyses. While benchmarking analyses on the

91  available genetic demultiplexing tools have shown effectiveness for demultiplexing small sample

92  sizes, limitations emerge as the number of multiplexed samples approach a population scale (6)

93  (7) (8) (9). For example, using computationally multiplexed samples, Neavin et al. evaluated the

94  performance of genetic demultiplexing tools as the number of samples approached a population

95  scale and observed diminished demultiplexing accuracy with increasing numbers of pooled

4

96    samples, as well as notable classification discrepancies between tools (10). Furthermore, even at

97    small sample sizes, divergent assignments between genetic demultiplexing tools are common (8)

98    (9) (11). Another feature that has been shown to affect genetic demultiplexing performance is the

99    underrepresentation of samples in a pool, which is especially relevant for cell culture-based

100    multiplexed experiments, as variable growth rates *in vitro* across cell lines is common (12) (8) (9).

101    Genetic demultiplexing tools have also shown low concordance for identifying heterogenic

102    doublets, which should be removed prior to downstream analyses to avoid technical noise in the

103    data (10). Importantly, benchmarking analyses have repeatedly highlighted ScSplit's poor

104    performance relative to the remaining tools (9) (10) (8) (11). The sum of these limitations calls to

105    question the robustness of the individual genetic demultiplexing tools for resolving the donor

106    identities of highly multiplexed samples, which represents an important hurdle for feasibly

107    achieving population-scale scRNAseq analysis.

108

109    In response to the divergent assignments commonly observed across tools, a consensus framework,

110    whereby only cells that show matching sample labels across all individual tools are retained for

111    downstream analyses, may appear sufficient to resolve the risk of introducing technical noise into

112    the data from misclassified cells. However, consensus frameworks are restricted to performing

113    only as well as the worst-performing tool, and genetic demultiplexing performance is highly

114    dataset dependent (10); thus, the overall performance of a consensus framework can vary

115    immensely between datasets. To this end, Neavin et al. proposed a majority vote framework for

116    genetic demultiplexing, whereby a cell is assigned to the sample called by the majority of tools

117    (10). However, this approach can be vulnerable to a subset of tools performing poorly on the

118    dataset, does not allocate additional weight to the votes of tools that perform more favourably on

119    the dataset, cannot account for instances when ties occur amongst tools, and cannot capture cells

120    that are correctly classified by only one tool. The sum of these limitations leads to the unnecessary

121    removal of cells from downstream analyses, reducing statistical power, especially for highly

122    multiplexed pools where each donor, on average, will have a lower representation of cells in the

123    pool. Moreover, the ability to capture the transcriptional profiles of rare cell types with scRNAseq

124    provides a notable advancement over bulk RNAseq and can strongly influence biological

125    interpretations (13); thus, investigators are reluctant to discard valuable cells in order to maximize

126    the analytic potential of their dataset.

127

128    To address the need for a robust genetic demultiplexing framework that can maximize the number

129    of confidently classified cells retained for downstream analyses, achieve high demultiplexing

130    accuracy for population-scale scRNAseq sample pooling, and maintain reliability across different

131    datasets, we developed Ensemblex: an accuracy-weighted ensemble genetic demultiplexing

132    framework designed to identify the most probable sample labels from each of its constituent tools

133    — Demuxalot, Demuxlet/Freemuxlet, Souporcell, and Vireo. Our ensemble method capitalizes on

134    combining distinct statistical frameworks for genetic demultiplexing while adapting to the overall

135    performance of its constituent tools on the respective dataset, making it resilient against a poorly

136    performing tool and facilitating a higher yield of cells for downstream analyses. The Ensemblex

137    workflow is assembled into a three-step pipeline — 1) accuracy-weighted probabilistic ensemble;

138    2) graph-based doublet detection; 3) Ensemble-independent doublet detection — and can

139    demultiplex pools with or without prior genotype information.

140

141 Here, we showcase Ensemblex's improved demultiplexing performance across a variety of settings

142 through benchmarking analyses on a total of 141 computationally multiplexed pools with known

143 ground-truth sample labels ranging in size from 4 to 80 samples. We applied the ensemble method

144 to three diverse, experimentally multiplexed datasets: 1) non-small cell lung cancer (NSCLC)

145 dissociated tumor cells from 7 individuals with donor-specific oligonucleotide labels; 2) iPSC-

146 derived dopaminergic neurons (DaN) from 22 healthy individuals; and 3) iPSC-derived neural

147 stem cells (NSC) from 9 individuals with attention deficit hyperactivity disorder (ADHD) and 7

148 healthy controls. We demonstrate Ensemblex's robustness across distinct datasets, its ability to

149 return a high proportion of confidently classified cells for downstream analysis, and the

150 implications that its improved demultiplexing performance has on biological interpretations of

151 multiplexed experiments.

152 **Table 1. Summary of individual genetic demultiplexing tools.**

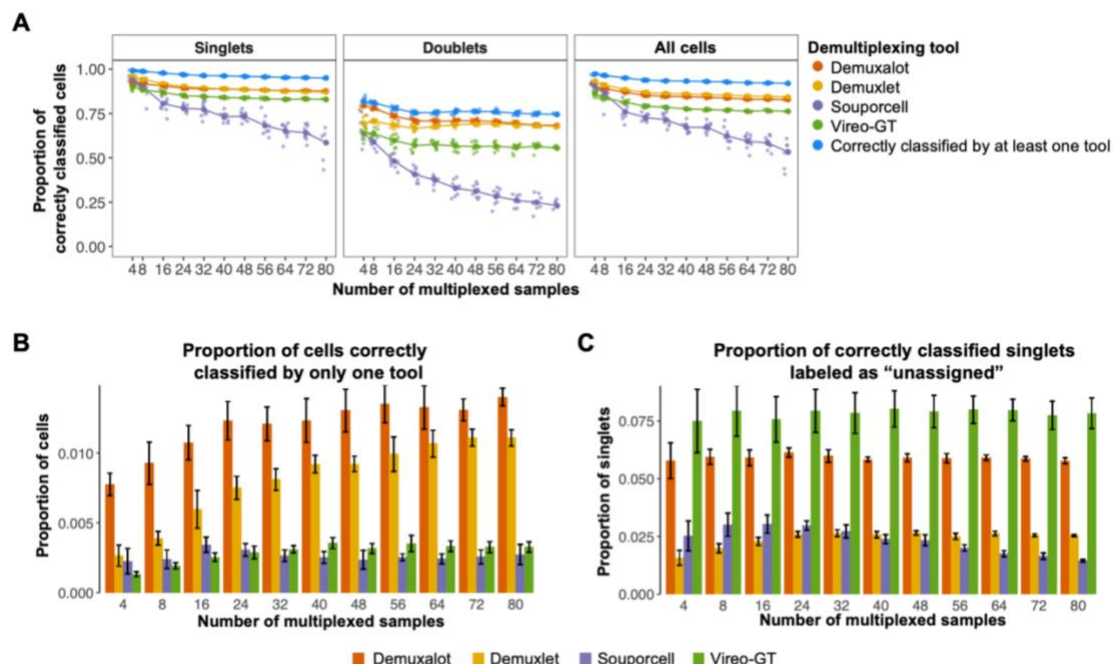| Genetic demultiplexing tool | Prior genotype information for genetic demultiplexing | Included in the Ensemblex framework |
| --- | --- | --- |
| Demuxalot (5) | Required | Yes |
| Demuxlet (6) | Required | Yes |
| Freemuxlet (6) | Not supported | Yes |
| ScSplit (7) | Optional | No |
| Souporcell (8) | Optional | Yes |
| Vireo (9) | Optional | Yes |

153

154

155 **Results and Discussion**

156 *Evaluating the performance of existing individual genetic demultiplexing tools*

7

157   To evaluate the performance of individual genetic demultiplexing tools, we generated

158   computationally multiplexed pools using scRNAseq of 80 different iPSC lines from Parkinson's

159   disease patients and healthy controls, which were differentiated towards a DaN state as part of the

160   Foundational Data Initiative for Parkinson's Disease (FOUNDIN-PD) (14). Processed scRNAseq

161   data from the independent iPSC lines were merged to simulate sample-pooling using a previously

162   described protocol (9), which provided known ground-truth donor and doublet labels. We

163   generated 96 *in silico* pools ranging in size from 4 to 80 multiplexed samples, where each sample

164   corresponded to a unique donor-of-origin. The *in silico* pools averaged 17,396 cells per pool with

165   a constant 15% doublet rate.

166

167   Leveraging whole-genome sequencing (WGS) of the 80 donors from which the iPSC lines were

168   derived and the four genetic demultiplexing tools that can utilize prior genotype information —

169   Demuxalot, Demuxlet, Souporcell, and Vireo-GT — we first investigated the proportion of

170   correctly classified cells by the individual tools (**Figure 1A**). Across the 96 *in silico* pools, all tools

171   showed decreasing demultiplexing performance as the number of samples within the pool

172   increased. Souporcell demonstrated the largest decrease in the proportion of correctly classified

173   cells as the number of multiplexed samples increased from 4 (mean = 90.60%) to 80 (mean =

174   53.27%). In accordance with previous findings (10, 15), the individual genetic demultiplexing

175   tools performed better on singlet classification than doublet detection, highlighting an avenue for

176   improved genetic demultiplexing accuracy by increasing the rate of heterogenic doublet

177   identification (**Figure 1A**).

178

**Figure 1. Evaluation of existing individual genetic demultiplexing tools.** Evaluation of genetic demultiplexing tools with prior genotype information on 96 *in silico* pools with known ground-truth sample labels ranging in size from 4 to 80 multiplexed induced pluripotent stem cell (iPSC) lines from genetically distinct individuals, averaging 17,396 cells per pool and a 15% doublet rate. **A)** Line graphs showing the proportion of correctly classified singlets, doublets, and all cells by each individual genetic demultiplexing tool across varying numbers of multiplexed iPSC lines in a single pool (sample number). The large dots show the mean proportion of correct classifications by an individual tool across replicates at a given sample size (n = 9 per pool size). The blue points show the proportion of cells that were correctly classified by at least one individual genetic demultiplexing tool: Demuxalot, Demuxlet, Souporcell, or Vireo-GT. **B)** Bar chart showing the mean proportion of total cells from an individual pool correctly classified by only one genetic demultiplexing tool. Error bars represent one standard deviation from the mean. (n = 9 per pool size) **C)** Bar chart showing the proportion of correctly classified singlet cells labelled as "unassigned" (ambiguous singlet assignments) due to assignment probabilities below the recommended threshold of the respective genetic demultiplexing tool. Error bars represent one standard deviation from the mean. (n = 9 per pool size).

We also investigated the proportion of cells that were correctly classified by at least one genetic demultiplexing tool to designate the best possible performance of an ensemble method that successfully incorporates every correct classification from its constituent tools (**Figure 1A**). Across the 96 *in silico* pools, an average of 93.64% of cells were correctly classified by at least
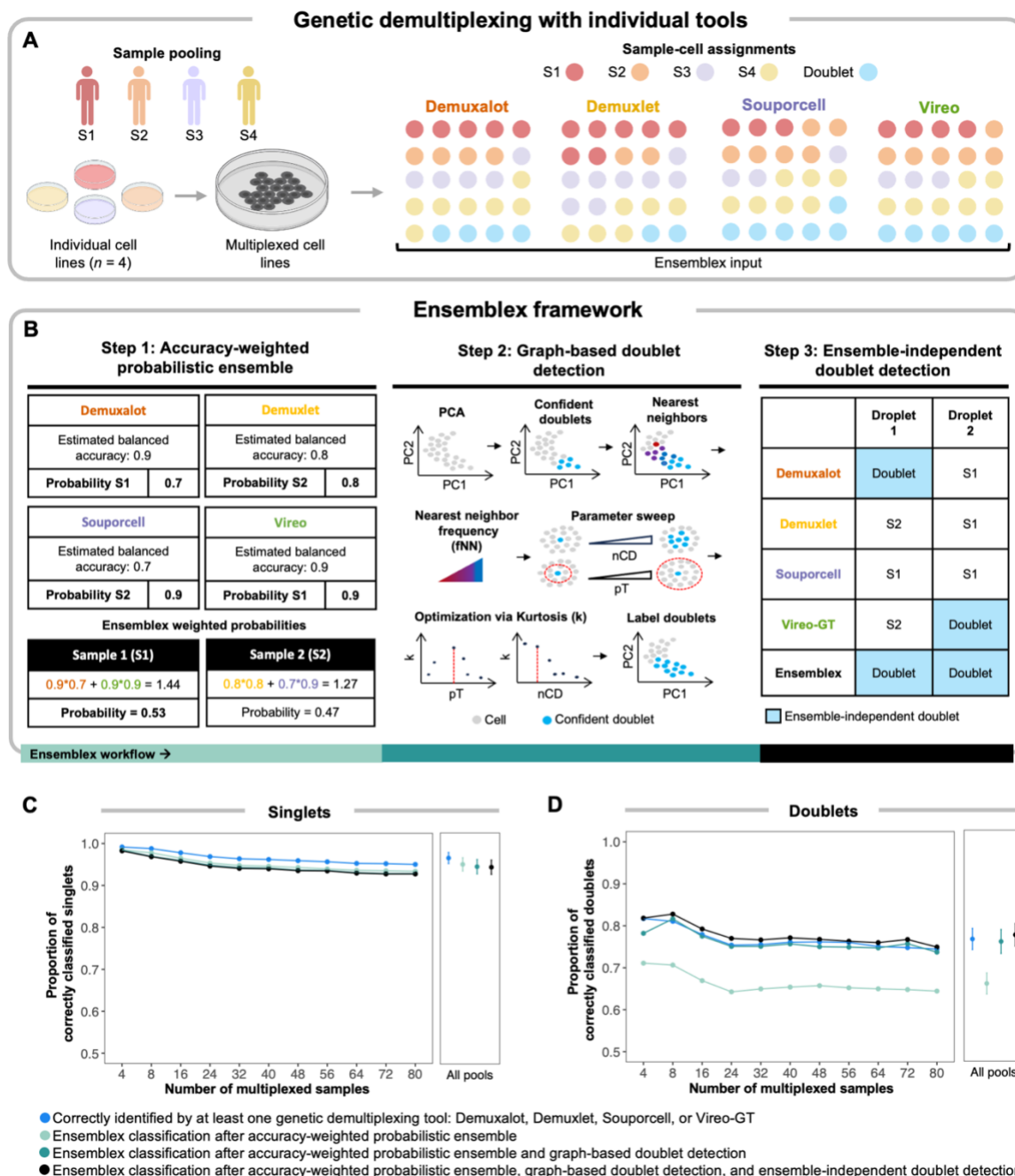
9

200     one tool. In comparison, Demuxlet, which demonstrated the best overall performance amongst

201     individual tools, correctly classified 86.73% of cells, on average. Demuxalot was consistently

202     responsible for the highest proportion of cells correctly classified by only one tool; 1.21% of

203     pooled cells, on average, were correctly classified by Demuxalot only, followed by Demuxlet

204     (mean = 0.83%), Vireo-GT (mean = 0.29%), and Souporcell (mean = 0.26%) (**Figures 1B;**

205     **Additional File 1: Figure S1**). Conversely, a consensus framework, correctly classified only

206     81.06% of cells, on average (data not shown). Based on these results, we reasoned that an ensemble

207     genetic demultiplexing method that can identify the most probable sample label from its

208     constituent tools, independent of a consensus assignment, would increase the yield of correctly

209     classified cells.

210

211     Next, we explored the frequency at which correctly classified singlets were labelled as unassigned

212     because their assignment probability failed to meet the tool's recommended probability threshold.

213     Across the 96 *in silico* pools, Vireo-GT consistently showed the highest proportion of correctly

214     classified singlets with insufficient assignment probabilities (Vireo-GT mean = 7.86%) followed

215     by Demuxalot (mean = 5.91%), Demuxlet (mean = 2.44%) and Souporcell (mean = 2.34%)

216     (**Figure 1C**). While a stringent probability threshold is important to prevent erroneous

217     classifications in downstream analyses, we reasoned that the unnecessary removal of correctly

218     classified cells could be mitigated by a carefully calibrated ensemble method that allocates

219     additional assignment confidence to cells with matching sample labels across constituent tools,

220     despite low internal tool-specific assignment probabilities.

221

222    We repeated the above analyses using the same 96 computationally multiplexed pools and the

223    genetic demultiplexing tools that do not require prior genotype information: Freemuxlet,

224    Souporcell, and Vireo. Here, we observed the same overarching limitations as when

225    demultiplexing with prior genotype information: 1) decreasing demultiplexing performance as the

226    number of multiplexed samples increased; 2) poor doublet detection performance compared to

227    singlet classification; 3) high rates of cells only correctly classified by a single tool; and 4)

228    discarded correctly classified cells due to insufficient assignment probabilities (**Additional File 1:**

229    **Figure S2)**. When we compared demultiplexing with and without prior genotype information, we

230    observed a trend towards a higher proportion of cells being correctly classified when prior

231    genotype information was available, as previously seen in separate benchmarking analyses (9)

232    (**Additional File 1: Figure S3**).

233

234    *Validating the Ensemblex framework on pools with known ground-truth sample labels*

235    To mitigate the limitations of the individual genetic demultiplexing tools and maximize the

236    analytic potential of multiplexed scRNAseq datasets, we developed Ensemblex (**Figure 2A)**. The

237    Ensemblex workflow begins by demultiplexing pooled samples with four distinct demultiplexing

238    algorithms, followed by three steps: 1) accuracy-weighted probabilistic ensemble; 2) graph-based

239    doublet detection; and 3) ensemble-independent doublet detection **(Figure 2B)**. As output,

240    Ensemblex returns its own cell-specific sample labels and corresponding assignment probabilities,

241    as well as the sample labels and corresponding assignment probabilities for each of its constituent

242    tools.

11

**Figure 2. Characterization of the Ensemblex framework.** Ensemblex is a probabilistic-weighted ensemble genetic demultiplexing framework for single-cell RNA sequencing analysis, which was designed to leverage the most probable sample labels from each of its constituent tools: Demuxalot, Demuxlet, Souporcell, and Vireo when using prior genotype information or Demuxalot, Freemuxlet, Souporcell, and Vireo when prior genotype information is not available. **A)** The Ensemblex workflow begins with demultiplexing pooled cells from genetically distinct individuals by each of the constituent tools. The outputs from each individual demultiplexing tool are then used as input into the Ensemblex framework. **B)** The Ensemblex framework comprises

252    three distinct steps that are assembled into a pipeline: 1) accuracy-weighted probabilistic ensemble,
253    2) graph-based doublet detection, and 3) ensemble-independent doublet detection. **C-D)** Line
254    graphs showng the contribution of each step of the Ensemblex framework on 96 *in silico* pools
255    with known ground-truth sample labels ranging in size from 4 to 80 multiplexed induced
256    pluripotent stem cell (iPSC) lines from genetically distinct individuals, averaging 17,396 cells per
257    pool and a 15% doublet rate. The average proportion of correctly classified **C)** singlets and **D)**
258    doublets across replicates at a given pool size is shown after sequentially applying each step of the
259    Ensemblex framework with prior genotype information (n = 9 per pool size). The right panels
260    show the average proportion of correct classifications across all 96 pools; error bars represent one
261    standard deviation from the mean. The blue points show the proportion of cells that were correctly
262    classified by at least one individual genetic demultiplexing tool: Demuxalot, Demuxlet,
263    Souporcell, or Vireo-GT.
264

265    In response to our observation that certain cells are correctly classified by only one tool, we

266    implemented the accuracy-weighted probabilistic ensemble component (Step 1) of the Ensemblex

267    framework. In brief, this unsupervised weighting model identifies the most probable sample label

268    for each cell by assigning weights to each tool's assignment probabilities based on their estimated

269    balanced accuracy for the dataset (see "Methods") (**Figures 2B)** (16). Ensemblex then retains the

270    sample label with the highest cumulative probability across its constituents. However, one

271    challenge for this framework is computing the balanced accuracy of the constituent tools for

272    experimentally multiplexed pools that lack ground-truth labels. Therefore, to estimate the balanced

273    accuracy of a particular constituent tool (e.g., Demuxalot) without ground-truth labels, Ensemblex

274    leverages the cells with a consensus assignment across the three remaining tools (e.g., Demuxlet,

275    Souporcell, and Vireo-GT) as a proxy for ground-truth. To validate this approach, we utilized *in*

276    *silico* pools with known ground truth sample labels to compute the Adjusted Rand Index (ARI)

277    between Ensemblex's sample labels when the balanced accuracy of the constituent tools was

278    computed using consensus labels or ground-truth labels. Here, we consistently observed a mean

279    ARI > 0.99, independent of the number of multiplexed samples in a pool, suggesting high

280    assignment concordance between the two approaches (**Additional File 1: Figure S4**). Applying

13

281    the accuracy-weighted probabilistic ensemble component to the 96 *in silico* pools correctly

282    classified 94.98% of singlets, on average, across all pools, approaching the number of singlets that

283    were correctly classified by at least one constituent tool (mean = 96.48%) (**Figure 2C**). In contrast,

284    only 66.01% of doublets, on average, were correctly identified across all pools after Step 1,

285    compared to 76.59% of doublets that were correctly identified by at least one constituent tool

286    (**Figure 2D**).

287

288    Given that previous analyses have demonstrated strong doublet call discordance across genetic

289    demultiplexing tools (10), it was unsurprising that Step 1 of the Ensemblex framework performed

290    poorly on doublet identification. Therefore, instead of relying on the cell type classifications of the

291    constituent tools (i.e., singlet or doublet), we elected to leverage the doublet-related features (e.g.,

292    doublet probability; see "Methods") returned by the constituent tools to identify the cells with the

293    highest doublet likelihood, independent of the existing classifications. We implemented this

294    approach in the graph-based doublet detection component (Step 2) of the Ensemblex framework,

295    which was specifically designed to increase the rate of true doublet detection. Step 2 begins by

296    identifying the top *n* most confident doublets in the pool (see "Methods"). Then, based on the

297    Euclidean distances in principal component analysis (PCA) space, the cells that appear most

298    frequently amongst the nearest-neighbors of the high confident doublets and exceed the optimized

299    percentile threshold for the nearest-neighbor frequency are labelled as doublets by Ensemblex

300    (**Figure 2B; Additional File 1: Figure S5;** see "Methods"). Upon applying the graph-based

301    doublet detection component to the 96 *in silico* pools following Step 1, Ensemblex correctly

302    identified 76.00% of doublets, on average: a 9.99% increase in doublet detection from Step 1. In

303     turn, the average proportion of correctly classified singlets across all pools (94.43%) decreased by

304     only 0.55% (**Figure 2D**).

305

306     The ensemble-independent doublet detection component (Step 3) of the Ensemblex framework

307     was implemented to further improve doublet detection. Step 3 was motivated by our observation

308     that certain tools, namely Demuxalot and Vireo, showed high doublet detection specificity (mean

309     = 0.99) on *in silico* pools with known ground-truth sample labels, but that Steps 1 and 2 failed to

310     incorporate a subset of these correct doublet calls (**Additional File 1: Figure S6**). Therefore, by

311     default, Ensemblex accepts the doublet calls made by Demuxalot and Vireo-GT (**Figure 2B**).

312     Applying the ensemble-independent doublet detection component to the 96 *in silico* pools

313     following Steps 1 and 2 further increased the average proportion of correctly identified doublets

314     across all pools by 1.58% for a total of 77.63% of doublets detected, while only decreasing the

315     average proportion of correctly classified singlets by 0.13% for a total of 94.30% of singlets

316     correctly classified (**Figures 2C and 2D**). Notably, owing to the graph-based doublet detection

317     component, the average proportion of doublets identified by Ensemblex exceeded the average

318     proportion of doublets that were correctly classified by at least one constituent tool.
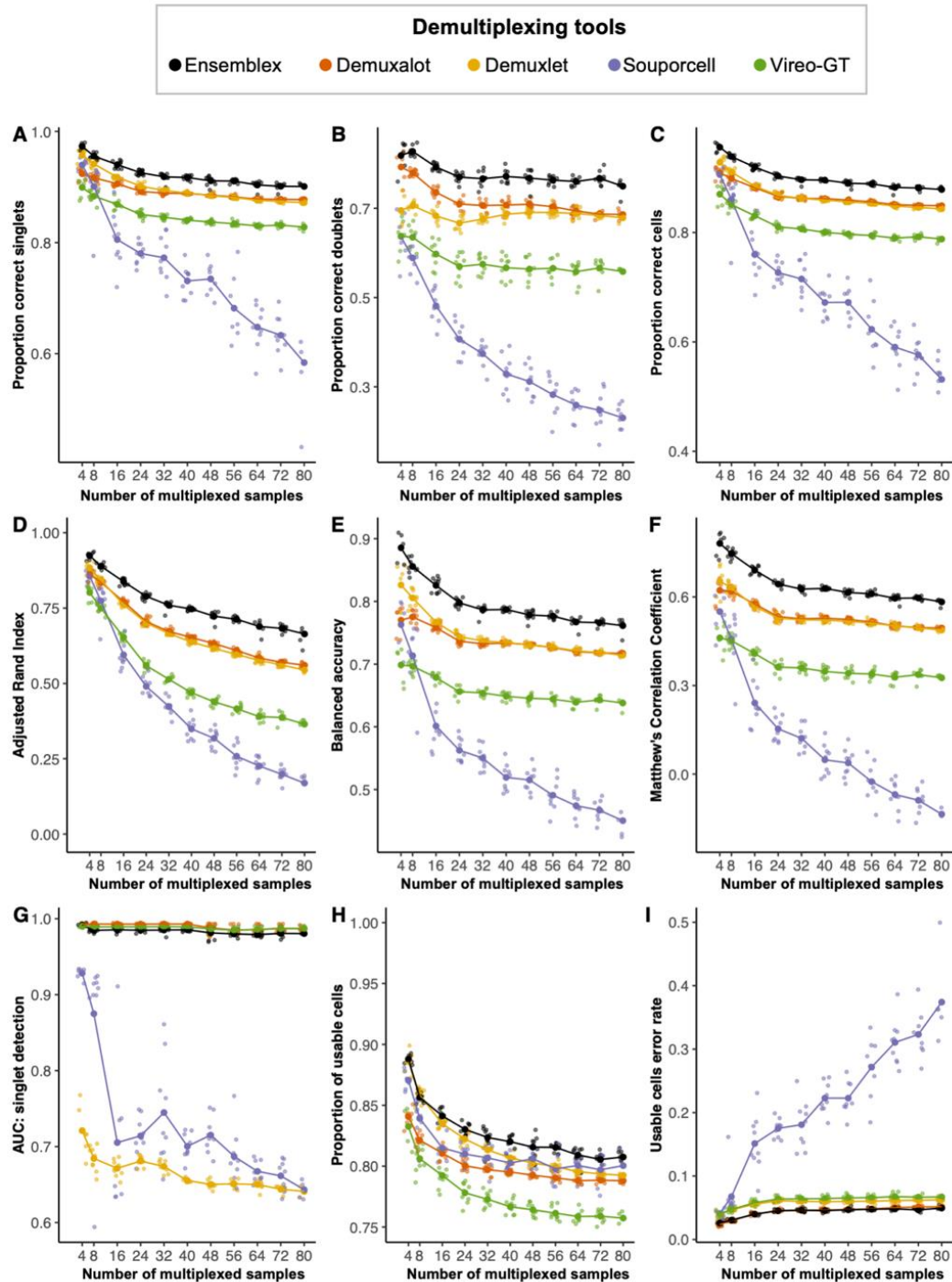
319

320     While the three-step workflow of the Ensemblex pipeline was designed to maximize the balance

321     between singlet classification and doublet identification, we do prioritize the identification of

322     doublets at the expense of a slightly lower singlet yield to minimize technical noise in the data.

323     However, we recognize that different experimental designs will require varying levels of doublet

324     detection stringency; thus, users can modify the percentile thresholds for graph-based doublet

325 detection and nominate different tools for ensemble-independent doublet detection (see

326 "Methods").

327

### *Benchmarking Ensemblex on pools with known ground-truth sample labels*

329 To benchmark Ensemblex against Demuxalot, Demuxlet, Souporcell, and Vireo-GT with prior

330 genotype information, we first utilized the 96 *in silico* pools with known ground-truth sample

331 labels to assess how Ensemblex's demultiplexing performance varied as the number of multiplexed

332 samples approached a cohort scale (4-80 samples). Unlike doublets, singlets were only considered

333 correctly classified if their assignment probability exceeded the recommended threshold of the

334 respective tool. On average across all pools, Ensemblex showed a higher proportion of correctly

335 classified singlets (mean = 92.19%), doublets (mean = 77.63%), and all cells (mean = 90.12%)

336 than the other tools. In comparison, Demuxlet, widely considered the "gold standard" tool,

337 correctly classified 89.72% of singlets, 68.57% of doublets, and 86.73% of all cells, on average

338 (**Figures 3A-3C**). Importantly, the discrepancy in the proportion of correctly classified cells

339 between Ensemblex and the next-best tool was amplified as the number of multiplexed samples

340 increased from 4 (2.78%) to 80 (3.52%), demonstrating that our ensemble method was able to

341 partially mitigate decreased demultiplexing accuracy as the pools approach a population scale.

**Figure 3. Ensemblex ground-truth benchmarking on computationally multiplexed pools.** The genetic demultiplexing tools with prior genotype information were evaluated on 96 *in silico* pools with known ground-truth sample labels ranging in size from 4 to 80 multiplexed induced pluripotent stem cell (iPSC) lines from genetically distinct individuals, averaging 17,396 cells per pool and a 15% doublet rate. A singlet was considered correctly classified if the assigned sample label matched the ground-truth sample label and the assignment probability exceeded the

17

349 recommended threshold for the respective tool; a doublet was considered correctly classified if the
350 assigned sample label matched the ground-truth sample label, regardless of the assignment
351 probability. **A-I)** Line graphs showing the performance of Ensemblex and the individual genetic
352 demultiplexing tools across evaluation metrics. The large dots show the mean value for each tool
353 across replicates at a given sample size (n = 9 per pool size). **A)** Proportion of correctly classified
354 singlets. **B)** Proportion of correctly classified doublets. **C)** Proportion of correctly classified cells.
355 **D)** Adjusted Rand Index between each tool's sample labels and the ground-truth sample labels. **E)**
356 Balanced accuracy of each tool. **F)** Matthew's Correlation Coefficient of each tool. **G)** Area under
357 the receiver operating characteristic curve (AUC) of the singlet assignment probability for each
358 tool. **H)** Proportion of usable cells returned by each tool. Usable cells were defined as cells
359 classified by singlets with an assignment probability exceeding the recommended threshold of the
360 respective tool. **I)** Error rate amongst the usable cells returned by each tool; erroneous
361 classifications comprised of true doublets labeled as singlets or true singlets assigned to the wrong
362 sample.
363

364 Next, we applied evaluation metrics for classification models to gauge the overall performance of

365 the genetic demultiplexing tools. We first computed the ARI to evaluate the similarity between the

366 demultiplexed sample labels and the ground-truth sample labels. Here, Ensemblex showed the

367 highest ARI with the ground truth sample labels across all pools (mean = 0.76), followed by

368 Demuxalot (mean = 0.67) and Demuxlet (mean = 0.66) (**Figure 3D**). We then computed the

369 balanced accuracy to evaluate the binary classification performance — singlet or doublet — of

370 each genetic demultiplexing tool as well as the Matthew's Correlation Coefficient (MCC), which

371 previous work has suggested is more reliable and informative for classification cases where

372 positive (singlet) and negative (doublet) cases have the same analytic importance (17). Across all

373 pools, Ensemblex showed the highest balanced accuracy (mean = 0.80) and MCC (mean = 0.64),

374 whereas Demuxalot and Demuxlet showed average balanced accuracies of 0.74 and 0.75,

375 respectively, and both tools showed an average MCC of 0.54 (**Figures 3E and 3F**). To evaluate

376 how well Ensemblex's confidence score (see "Methods") and each constituent tool's assignment

377 probability corresponded to the accuracy of their singlet classification, we plotted the area under

18

378     the receiver operating characteristic curve (AUC). Although Demuxalot (mean = 0.99) and Vireo-

379     GT (mean = 0.99) showed the highest AUC across all pools on average, Ensemblex's AUC was

380     comparable (mean = 0.98) (**Figure 3G**).

381

382     Finally, we investigated the proportion of usable cells returned by each demultiplexing tool and

383     the error rate amongst usable cells. We define usable cells as singlet classifications exceeding the

384     recommended probability threshold of the respective tool, while the error rate amongst usable cells

385     constituted incorrectly classified singlets to the wrong donor-of-origin or true doublets incorrectly

386     classified as singlets. We observed that, on average, Ensemblex returned the highest proportion of

387     usable cells across all pools (82.66%), followed by Demuxlet (81.66%), Souporcell (81.01%),

388     Demuxalot (79.99%), and Vireo-GT (77.53%) (**Figure 3H**). Importantly, Ensemblex showed the

389     lowest error rate amongst usable cells (4.34%), followed by Demuxalot (4.43%), Demuxlet

390     (5.77%), Vireo-GT (6.16%), and Souporcell (21.82%) (**Figure 3I**).

391

392     Using computationally multiplexed pools comprised of 24 iPSC lines, we further assessed how the

393     performance of Ensemblex varied as a function of the number cells in a pool when prior genotype

394     information was available. Here, we observed that our ensemble method consistently outperformed

395     the individual demultiplexing tools (**Additional File 1: Figure S7**). When cells are pooled

396     experimentally, it is reasonable to expect some iPSC lines to be underrepresented in the pool.

397     Therefore, to assess Ensemblex's demultiplexing performance in the presence of an

398     underrepresented iPSC line, we produced computationally multiplexed pools comprising of 24

399     samples, with one sample showing varying degrees of under representation. Again, we observed

400     that Ensemblex consistently outperformed the individual tools (**Additional File 1: Figure S8**).
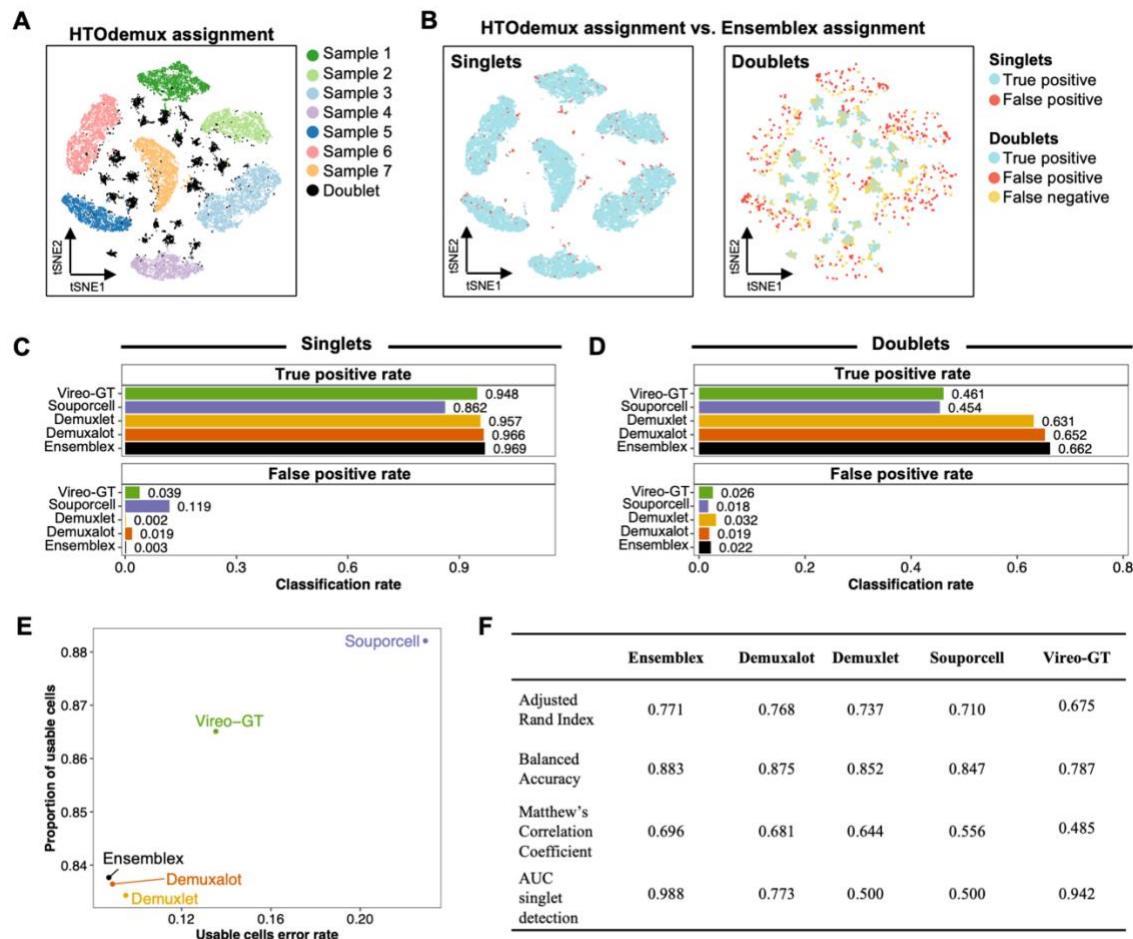
19

401 Finally, we repeated the above analyses to assess whether the benefits of using Ensemblex to

402 demultiplex with prior genotype information extended to cases where prior genotype information

403 is not available. In doing so, we observed a trend towards better overall performance by

404 Ensemblex; however, the discrepancy between Ensemblex and the top-performing individual

405 tools, namely Freemuxlet and Souporcell, was less pronounced than when demultiplexing with

406 prior genotype information (**Additional File 1: Figures S9-S11).**

407

408 Taken together, these results indicate that the Ensemblex framework mitigates the limitations of

409 the individual tools, leading to greater overall demultiplexing performance across computationally

410 multiplexed pools with known ground-truth labels. Ultimately, Ensemblex's improved

411 demultiplexing performance translates to a higher recovery of usable cells for downstream

412 analyses as well as a higher accuracy amongst usable cells, limiting the unnecessary removal of

413 cells from the dataset and mitigating the introduction of technical artifacts into biological analyses.

414

415 ***Evaluating Ensemblex on experimentally pooled samples with donor-specific oligonucleotide***

416 ***labels***

417 To determine whether Ensemblex's improved performance across the *in silico* pools is reflected in

418 real-world multiplexed experiments, we applied Ensemblex to an experimentally multiplexed pool

419 composed of NSCLC dissociated tumor cells from 7 donors, hereafter referred to as the NSCLC

420 dataset (18). Importantly, these NSCLC cells were labelled with donor-specific Cell Multiplexing

421 Oligonucleotides (CMOs), providing a proxy for ground-truth sample labels to evaluate the

422 performance of the genetic demultiplexing tools. For this experiment, we used HTOdemux (19) to

423 assign the cells back to their donor-of-origin based on the CMO expression profiles. HTOdemux

424    confidently assigned 19,695 cells, of which 15,534 (78.87%) were assigned to individual donors

425    and 4,161 (21.13%) were assigned as doublets; 769 cells (3.76%) were unassignable at a positive

426    quantile of 0.99 and were excluded from downstream analyses (**Figures 4A)**. Application of the

427    Ensemblex framework with prior genotype information to the NSCLC dataset achieved a singlet

428    true positive (TP) rate of 96.92% and doublet TP rate of 66.21% (**Figure 4B**). To evaluate the

429    benefits of utilizing the entire Ensemblex workflow (Steps 1-3), we investigated the contribution

430    of each step of the Ensemblex framework to the overall demultiplexing accuracy. Applying graph-

431    based doublet detection (Step 2) and ensemble-independent doublet detection (Step 3) to the

432    accuracy weighted assignments obtained from Step 1 increased the proportion of correctly

433    identified doublets by 14%, while slightly decreasing the proportion of correctly classified singlets

434    by 0.05% (**Additional File 1: Table S1**). Although users can elect to utilize different step-

435    combinations of the Ensemblex pipeline, these results reaffirm that leveraging the entire workflow

436    maximizes the overall demultiplexing accuracy by achieving a meticulous balance between singlet

437    classification and doublet identification.

**Figure 4. Evaluating Ensemblex on experimentally multiplexed cells using donor-specific oligonucleotide labels as a proxy for ground-truth.** Non-small cell lung cancer (NSCLC) dissociated tumor cells from 7 individuals were pooled and labelled with donor-specific oligonucleotide-labels. Cells were demultiplexed according to their expression of donor-specific oligonucleotide labels by HTOdemux; HTOdemux's sample labels were used as a proxy for ground truth. True positives (TP) singlets were defined as cells classified as singlets by both HTOdemux and Ensemblex with matching sample labels; false positives (FP) singlets were defined as cells classified as singlets by both HTOdemux and Ensemblex but assigned to different donors. TP doublets were defined as cells classified as doublets by both HTOdemux and Ensemblex; FP doublets were defined as cells classified as singlets by HTOdemux and doublets by Ensemblex; false negatives (FN) doublets were defined as cells classified as doublets by HTOdemux and singlets by Ensemblex. **A)** T-distributed Stochastic Neighbor Embedding (t-SNE) visualization of HTOdemux's sample labels. **B)** T-SNE visualization of Ensemblex's demultiplexing performance using HTOdemux's sample labels as ground truth for singlets (left) and doublets (right). **C)** Bar plots showing the singlet TP and FP rates for each genetic demultiplexing tool using HTOdemux's sample labels as ground truth. **D)** Bar plots showing the doublet TP and FP rates for each genetic demultiplexing tool using HTOdemux's sample labels as
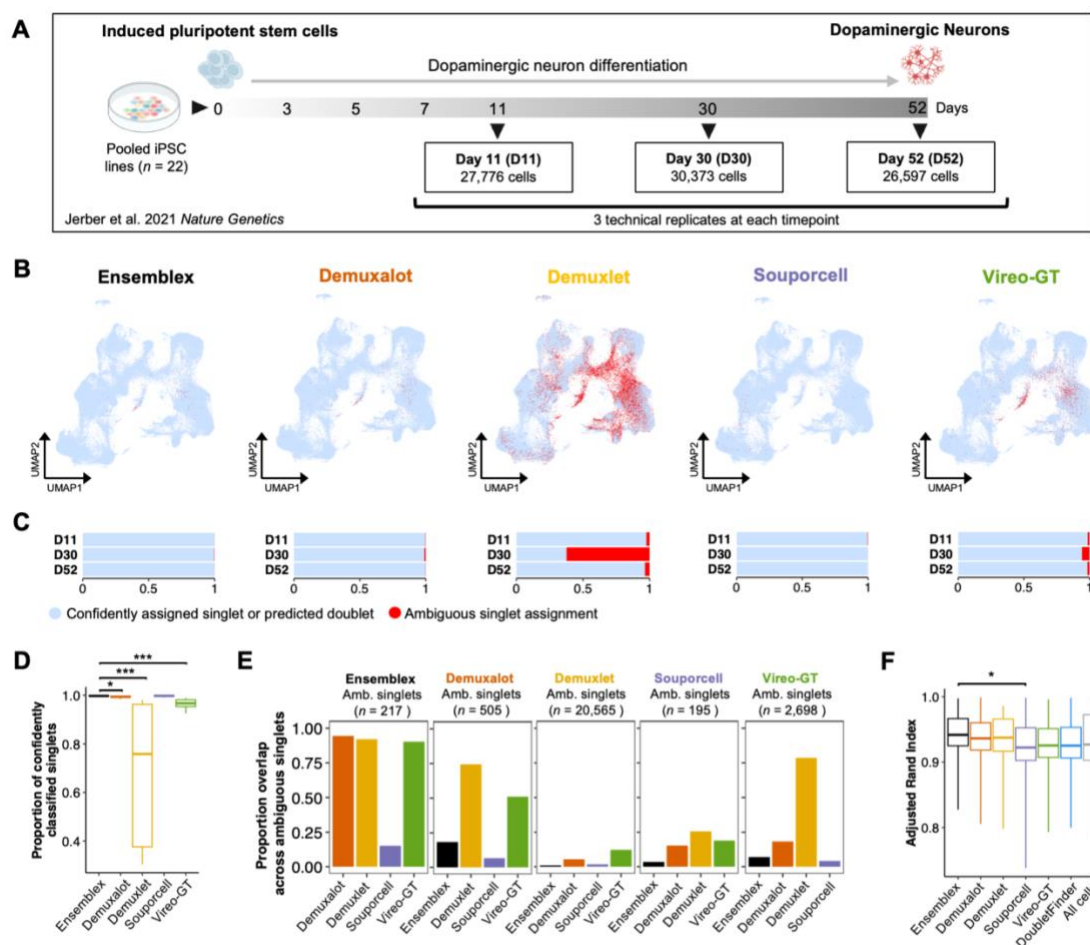
22

456   ground truth. **E)** Scatter plot showing the proportion of usable cells (confidently classified singlets)

457   and the corresponding usable cell error rate for each genetic demultiplexing tool. **F)** Adjusted Rand

458   Index, balanced accuracy, Matthew's Correlation Coefficient, and area under the receiver operating

459   characteristic curve (AUC) of the singlet assignment probability for each genetic demultiplexing

460   tool.

461

462   Upon comparing Ensemblex's demultiplexing performance with prior genotype information on

463   the NSCLC dataset to the individual genetic demultiplexing tools, it emerged that our ensemble

464   method obtained the highest singlet and doublet TP rates (**Figures 4C and 4D**). Ensemblex and

465   Demuxlet also showed the lowest singlet false positive (FP) rates (0.25% and 0.21%, respectively),

466   indicating that singlets were least frequently assigned to the wrong donor-of-origin by these two

467   methods compared to Demuxalot (1.87%), Vireo-GT (3.91%), and Souporcell (11.94%).

468   Souporcell and Vireo-GT returned the highest proportion of usable cells (confidently classified

469   singlets; 88.21% and 86.51%, respectively); albeit, at the expense of high usable cell error rates

470   (22.91% and 13.53%, respectively) (**Figure 4E**). In turn, Ensemblex, Demuxalot, and Demuxlet

471   showed lower error rates across the usable cells (8.75%, 8.91%, and 9.51%, respectively), amongst

472   which Ensemblex returned the highest proportion of usable cells (83.77%) compared to Demuxalot

473   (83.64%) and Demuxlet (83.43%). Here, the relatively high error rate amongst usable cells

474   returned by each demultiplexing tool is attributed to true doublets classified as singlets. Finally,

475   we computed the ARI, balanced accuracy, MCC, and AUC for singlet detection for each tool and

476   observed that Ensemblex again outperformed the remaining tools (**Figure 4F**). We repeated the

477   above analyses without prior genotype information and observed a similar trend towards better

478   overall performance by Ensemblex (**Additional File 1: Table S2 and Figure S12**). Together, these

479   results corroborate that Ensemblex's improved performance on the *in silico* pools extends to

480   experimentally multiplexed samples.

481

23

482    *Application of Ensemblex to experimentally pooled, highly multiplexed subjects*

483    To evaluate Ensemblex's demultiplexing performance on experimentally pooled, highly

484    multiplexed scRNAseq datasets with prior genotype information, we used pools containing iPSC

485    lines from 22 donors that were differentiated towards DaN by Jerber et al., hereafter referred to as

486    the DaN dataset (12) (**Figure 5A**). To capture the transcriptional changes throughout neurogenesis,

487    Jerber et al. performed scRNAseq of the iPSC lines grown in pooled cultures at days 11, 30, and

488    52 of differentiation (**Figure 5A**). Using three technical replicates from each timepoint, we

489    obtained 84,746 cells after performing quality control as previously described (12) (**Additional**

490    **File 1: Table S3**). Each technical replicate was demultiplexed independently by Ensemblex and

491    its constituent tools.



492

493 **Figure 5. Application of Ensemblex to highly multiplexed, experimentally pooled cultures of**
494 **differentiated dopaminergic neurons.** **A)** Time line of iPSC pooling, dopaminergic neuron
495 (DaN) differentiation, and sample collection from the DaN dataset by Jerber et al. (12). Three
496 technical replicates at each time point (days 11, 30 and, 52 of differentiation) from pools containing
497 22 individual iPSC lines were used in the analysis. Across all timepoints and technical replicates,
498 84,746 cells were obtained for analysis. **B)** Uniform manifold approximation and projection
499 (UMAP) plots showing confidently assigned singlets or predicted doublets (blue) and ambiguous
500 singlets (singlet assignments with insufficient assignment probabilities; red) returned by each
501 demultiplexing tool. **C)** Stacked bar chart showing the proportion of confidently assigned singlets
502 or predicted doublets (blue) and ambiguous singlets (red) across technical replicates at each time
503 point returned by each demultiplexing tool. **D)** Boxplot showing the proportion of confidently
504 classified singlets across technical replicates and time points by each demultiplexing tool.
505 Wilcoxon rank-sum tests were used to compare the proportion of confidently classified singlets by
506 Ensemblex to that of its constituents (n = 9 pools). **E)** Bar chart showing the proportion of
507 overlapping ambiguous singlet assignments amongst demultiplexing tools across technical
508 replicates and time points (n = 9 pools). **F)** Boxplot showing the Adjusted Rand Index (ARI)
509 assessing cluster stability across a range of 11 clustering resolutions (*n* clustering iterations = 25)
510 after removing doublets identified by each demultiplexing tool. Wilcoxon rank-sum tests were
511 used to compare the clustering ARI after removing Ensemblex doublets to the clustering ARI after
512 removing doublets identified by each constituent tool. * Adjusted P-value < 0.05; ** adjusted P-
513 value < 0.01; *** adjusted P-value < 0.001

514

515 To characterize the relationship between Ensemblex and its constituent demultiplexing tools, we

516 computed the ARI between Ensemblex's sample labels and those of its constituent as well as the

517 percent contribution of each tool to Ensemblex's final sample labels (**Table 2**). Notably, we

518 observed that across day 30 technical replicates Demuxlet showed an ARI of 0.063 with

519 Ensemblex and only contributed 29.74% to Ensemblex's final sample labels. In contrast, across

520 day 11 and 52 technical replicates Demuxlet showed an ARI of 0.928 and 0.884, respectively, and

521 contributed 95.91% and 90.55%, respectively, to Ensemblex's final sample labels. Importantly,

522 Demuxlet's variable contribution to Ensemblex's sample labels across sequencing time points

523 demonstrates our ensemble method's ability to adapt to the relative performance of its constituent

524 tools and override the classifications of a poorly performing tool on the respective dataset.

525 **Table 2. Application of Ensemblex to pooled cultures of dopaminergic neurons from 22**
526 **healthy controls.**

| | ARI between Ensemblex and constituent tool assignments | | | Percent contribution to Ensemblex assignments | | | *n* usable cells | *n* doublets |
|---|---|---|---|---|---|---|---|---|
| | Day 11 | Day 30 | Day 52 | Day 11 | Day 30 | Day 52 | | |
| **Demuxalot** | 0.987 | 0.955 | 0.982 | 97.29% | 94.75% | 97.57% | 75,962 | 8,279 |
| **Demuxlet** | 0.928 | 0.062 | 0.884 | 95.91% | 29.74% | 90.55% | 57,567 | 6,614 |
| **Souporcell** | 0.883 | 0.876 | 0.912 | 91.62% | 91.82% | 93.84% | 76,811 | 7,740 |
| **Vireo-GT** | 0.961 | 0.879 | 0.958 | 95.95% | 88.80% | 95.16% | 75,933 | 6,115 |
| **Ensemblex** | NA | NA | NA | NA | NA | NA | 76,222 | 8,307 |
| **DoubletFinder** | NA | NA | NA | NA | NA | NA | NA | 4,597 |

527 Pooled cultures of induced pluripotent stem cell (iPSC) lines from 22 healthy donors were
528 differentiated towards a dopaminergic neuron (DaN) fate and sequenced on days 11, 30, and 52 of
529 differentiation by Jerber et al. (12). For the analysis we used three technical replicates for each
530 sequencing timepoint. Each pool was demultiplexed independently by Ensemblex and its
531 constituent tools with prior genotype information. The Adjusted Rand Index (ARI) between
532 Ensemblex's assignments and those of the constituent tools was computed across technical
533 replicates corresponding to each differentiation timepoint. The percent contribution represents the
534 proportion of assignments from each constituent tool that matched Ensemblex's assignments.
535 Usable cells were defined as singlet classifications whose assignment probability exceeded the
536 recommended threshold of the respective tool. Abbreviations: NA = Not applicable.

537

538 To elucidate the discrepancy in Demuxlet's contribution to Ensemblex's sample labels across

539 sequencing time points, we investigated the proportion of ambiguous singlet assignments from

540 Ensemblex and its constituents. Ambiguous singlets are defined as singlet classifications whose

541 assignment probabilities failed to meet the recommended threshold of the respective tool, leaving

542 the identity of the pooled cell unresolved. Across 84,746 cells, Souporcell (195 singlets; 0.23% of

543 cells) and Ensemblex (217 singlets; 0.26% of cells) showed the lowest proportion of ambiguous

544 singlet assignments, followed by Demuxalot (505 singlets; 0.60% of cells) and Vireo-GT (2,698

26

545      singlets; 3.18% of cells). Strikingly, Demuxlet showed 20,565 ambiguous singlet assignments

546      (24.27% of cells), with 92.04% derived from day 30 technical replicates, reflecting Demuxlet's

547      remarkably low contribution to Ensemblex's sample labels for cells sequenced at this timepoint

548      (**Figures 5B and 5C**). In accordance with previous analyses (9, 10), Demuxlet was consistently

549      amongst the top performing constituent tools throughout our benchmarking analyses. Yet, its poor

550      performance across day 30 technical replicates illustrates how the accuracy of individual tools can

551      vary greatly between datasets, highlighting the importance of utilizing multiple distinct algorithms

552      for genetic demultiplexing. We compared the mean proportion of confidently classified singlets

553      across technical replicates from each time point (*n = 9)* between Ensemblex (99.72%) and each

554      constituent demultiplexing tool using a Wilcoxon rank-sum test. After correction for multiple

555      hypothesis testing, we observed that the mean proportion of confidently classified singlets by

556      Ensemblex was significantly higher than Demuxalot (mean = 99.36%, P-value = 3.55e-3),

557      Demuxlet (mean = 75.82%, P-value = 1.55e-5), and Vireo-GT (mean = 96.71%, P-value = 1.55e-

558      5) (**Figure 5D**). Thus, despite Demuxlet's unusually poor performance across day 30 technical

559      replicates, Ensemblex still confidently classified 27,520 singlets (99.61% of singlet assignments)

560      from these pools. Indeed, our ensemble method mitigates the consequences of a poorly performing

561      constituent tool by outweighing the erroneous classifications. In contrast, using a consensus

562      framework returned only 7,446 confidently classified singlets from day 30 technical replicates

563      (20,074 fewer cells than Ensemblex), limiting the availability of data for downstream analyses.

564

565      To further evaluate the ambiguity amongst singlet classification, we investigated the intersection

566      of ambiguous singlets across demultiplexing tools, reasoning that cells that are most challenging

567      to demultiplex would be labelled as ambiguous across all tools (**Figure 5E**). The singlets that were

568 assigned as ambiguous by Ensemblex showed the highest ambiguous singlet rate across the

569 remaining tools (mean across all constituent tools = 73.04%; mean across Demuxalot, Demuxlet,

570 and Vireo-GT = 92.32%). In contrast, while Souporcell showed the lowest ambiguous singlet rate

571 overall, only 15.90% of its unassigned singlets, on average, were ambiguous across the remaining

572 tools. These results indicate that the cells labelled as ambiguous by Ensemblex represent the cells

573 that are most challenging to classify across the distinct demultiplexing algorithms. Indeed, limiting

574 Ensemblex's ambiguous singlet assignments to those that are most difficult to classify is critical

575 for maintaining a balance between maximizing the number of usable cells and minimizing the

576 introduction of technical artifacts into downstream analyses from misclassified cells.
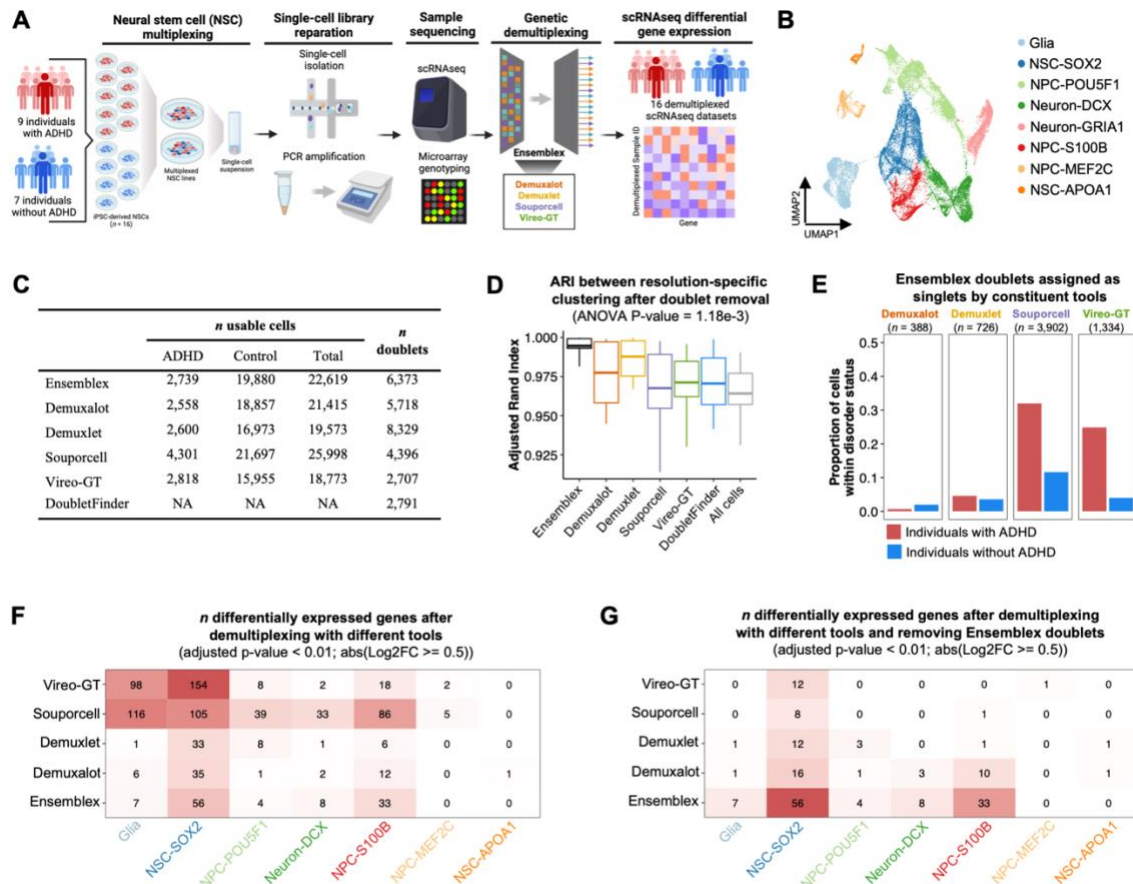
577

578 Next, we compared the doublet predictions made by each genetic demultiplexing tool and

579 DoubletFinder, a doublet detection tool that predicts doublets by estimating the similarity of the

580 transcriptional profile of a pooled cell to artificial doublets generated by combining the

581 transcriptional profiles of randomly selected cell pairs (20). Although the average number of

582 unique molecular identifiers (UMI) per cell across doublets identified by each tool was

583 significantly higher than the consensus singlets (**Additional File 1: Figure S13**), we observed a

584 notable discrepancy in the number of doublets identified by each tool; DoubletFinder identified

585 the fewest doublets ($n = 4,597$), while Ensemblex identified the most doublets ($n = 8,307$) (**Table**

586 **1**). Accordingly, all tools identified doublets that every other tool assigned as singlets (**Additional**

587 **File 1: Figure S13**). While Ensemblex identified the highest number of doublets, it still returned

588 a higher number of confidently classified singlets ($n = 76,222$) than Demuxalot ($n = 75,962$),

589 Demuxlet ($n = 57,567$), and Vireo-GT ($n = 75,933$). Thus, even though the Ensemblex framework

28

590    prioritizes the identification of doublets at the expense of a slightly lower singlet classification

591    rate, our ensemble method still returns a high proportion of usable cells for downstream analyses.

592

593    To evaluate the impact of doublet removal on the stability of clusters in the DaN dataset, we

594    performed 25 different random start iterations of the Louvain network detection at various

595    clustering resolutions after removing the doublets identified by each tool (21). Removing the

596    doublets identified by Ensemblex resulted in the highest ARI (mean ARI = 0.942), on average,

597    across clustering resolutions (**Figure 5F**), suggesting the greatest cluster stability. However,

598    Wilcoxon rank-sum tests only revealed a statistically significant difference in the cluster

599    assignment ARI between Ensemblex and Souporcell (mean ARI = 0.922, P-value = 1.08e-2) after

600    correction for multiple hypothesis testing. Nonetheless, the highest cluster stability after removal

601    of Ensemblex's putative doublets illustrates how improved doublet detection can translate to

602    improved biological analyses and is reflective of its superior doublet identification performance

603    on the benchmarking analyses.

604

605    ***Evaluating the impact of demultiplexing tools on differential gene expression analysis***

606    To evaluate the impact of genetic demultiplexing tools on scRNAseq DGE analysis, we

607    multiplexed iPSC-derived NSCs from individuals with ADHD and controls (**Figure 6A**). NSCs

608    were pooled and cultured until 100% confluence was reached. Two multiplexing experiments were

609    performed: Experiment 1 (*n* ADHD = 7; *n* control = 6) and Experiment 2 (*n* ADHD = 9; *n* control

610    = 7). After filtering cells for > 500 total and unique RNA transcripts, we obtained 30,433 cells

611    across both pools. Louvain clustering on the integrated scRNAseq dataset identified 12 clusters,

612    which were annotated as eight putative cell types (**Figure 6B**).

**Figure 6. Evaluating the impact of discordant assignments between genetic demultiplexing tools on differential gene expression analysis. A)** Schematic illustrating the workflow for the neural stem cell (NSC) dataset. Pooled induced pluripotent stem cell (iPSC)-derived neural stem cell cultures from individuals with attention deficit hyperactivity disorder (ADHD) and controls were collected in two separate experiments. NSCs were dissociated for single-cell RNA sequencing and prior genotype information of the pooled subjects was obtained through microarray genotyping. The pools were demultiplexed by Ensemblex and its constituents with prior genotype information and differential gene expression (DEG) was computed between ADHD and controls. **B)** Uniform manifold approximation and projection (UMAP) plot showing the putative cell types. **C)** Summary of the number of usable cells — singlets above the recommended probability threshold of the respective demultiplexing tool — assigned to ADHD donors and controls and the number of identified doublets by each demultiplexing tool. **D)** Boxplot showing the Adjusted Rand Index (ARI) assessing cluster stability across a range of 11 clustering resolutions (*n* clustering iterations = 25) after removing doublets identified by each demultiplexing tool. A one-way Analysis of Variance (ANOVA) test comparing the ARI after removing doublets identified by each tool revealed a significant difference between tools (n = 11 clustering resolutions; P-value = 1.18e-3). **E)** Proportion of ADHD and control cells identified as putative doublets by Ensemblex that were assigned as singlets by the constituent demultiplexing tools. **F)**

632   Heatmap showing the number of cell-type specific DEGs between ADHD and controls using the

633   subject labels of each demultiplexing tool. **G)** Heatmap showing the number of cell-type specific

634   DEGs between ADHD and controls using the subject labels of each demultiplexing tool and

635   removing putative doublets identified by Ensemblex. Cell-types not shown in the heatmaps had no

636   DEGs passing the adjusted P-value < 0.01 and |Log2FC >= 0.5| threshold across all tools.

637

638   We independently demultiplexed both pools using Ensemblex and its constituents to assign the

639   cells back to their donor-of-origin with prior genotype information (**Figure 6C).** The number of

640   cells assigned to ADHD and control donors by each genetic demultiplexing tool is shown in

641   **Additional File 1: Table S6**. Importantly, the NSC dataset provides a valuable illustration of the

642   consequences of unnecessarily discarding cells from downstream analyses. For example,

643   Ensemblex and Vireo-GT returned 2,387 and 882 confidently assigned $GRIA1^{high}$ neurons,

644   respectively, whereas a consensus approach would have confidently assigned only 563 $GRIA1^{high}$

645   neurons (**Additional File 1: Table S6**).

646

647   Each genetic demultiplexing tool predicted the ADHD cells to be vastly underrepresented

648   compared to the control cells; Ensemblex assigned 2,739 cells to individuals with ADHD and

649   19,880 cells to controls, suggesting that the ADHD iPSC lines were lost throughout the culturing

650   and sequencing process (**Figure 6C).** Additionally, we observed a notable difference in the number

651   of identified doublets across the tools; Vireo-GT identified the fewest doublets ($n = 2,707$), while

652   Demuxlet identified the most doublets ($n = 8,329$) (**Figure 6C**). We aimed to characterize the

653   change in cluster stability after removing the doublets identified by each tool and observed that

654   removing the doublets identified by Ensemblex resulted in the highest ARI (mean ARI = 0.995),

655   on average, across clustering resolutions (**Figure 6D**). A one-way ANOVA test comparing the

656   clustering ARI after removal of doublets identified by each tool revealed a significant difference

657   between tools (P-value = 1.18e-3). Demuxlet ($n = 8,329$) identified more doublets than Ensemblex

658    (n = 6,373), but exhibited lower cluster stability (ARI), suggesting that increased cluster stability

659    is not merely representative of the number of doublets removed but rather the quality of doublet

660    removal.

661

662    Given the underrepresentation of ADHD cells across the dataset, we elected to investigate the cells

663    that were identified as doublets by Ensemblex but assigned as singlets by the constituent tools and

664    how these putative doublets were distributed across samples according to disorder status.

665    Demuxalot ($n = 388$) and Demuxlet ($n = 726$) assigned a relatively low number of Ensemblex's

666    doublets as singlets, which represented 0.66% and 4.58% of ADHD sample assignments,

667    respectively, and 1.97% and 3.58% of control sample assignments, respectively (**Figure 6E**). In

668    contrast, Souporcell ($n = 3,902$) and Vireo-GT ($n = 1,334$) assigned a relatively high number of

669    Ensemblex's doublets as singlets, which represented 31.97% and 24.88% of ADHD sample

670    assignments, respectively, and 11.65% and 3.97% of control sample assignments, respectively,

671    illustrating how variable doublet detection can impact the assembly of cells assigned to donor

672    categories and which cells are retained for downstream analyses.

673

674    Finally, we used the model-based analysis of single-cell transcriptomics (MAST) statistical

675    framework to compute cell-type specific DGE between individuals with ADHD and controls using

676    the demultiplexed sample labels from each tool (22). We observed a significant discrepancy in the

677    number of cell type-specific differentially expressed genes (DEGs; adjusted P-value < 0.01;

678    absolute log2 fold change > 0.5) depending on the demultiplexing tool used (**Figure 6F**). Most

679    notably, for glia cells Souporcell identified 116 DEGs; Vireo-GT identified 98 DEGs; Ensemblex

680    identified 7 DEGs; Demuxalot identified 6 DEGs; and Demuxlet identified 1 DEG. Similar

681    patterns were observed across $SOX2^{high}$ NSCs, $POU5F1^{high}$ neural progenitor cells (NPC),

682    $S100B^{high}$ NPCs, and $DCX^{high}$ neurons, whereby Souporcell or Vireo-GT's sample labels resulted

683    in a remarkably high number of DEGs compared to Ensemblex, Demuxalot, and Demuxlet. Given

684    that Souporcell and Vireo-GT made relatively few doublet calls and that 31.97% and 24.88% of

685    ADHD sample assignments made by Souporcell and Vireo-GT, respectively, were putative

686    doublets identified by Ensemblex, we elected to repeat the DGE analysis using the demultiplexed

687    sample labels from each tool but this time we removed all putative doublets identified by

688    Ensemblex. In doing so, we observed a decrease in the number of DEGs identified by Souporcell

689    and Vireo-GT across cell types, suggesting that the putative doublets identified by Ensemblex,

690    which were classified as singlets by Souporcell and Vireo-GT, were driving the initial signals

691    (**Figure 6G**). For example, the number of glia-specific DEGs decreased from 116 to 0 with

692    Souporcell's sample labels, and 98 to 0 with Vireo-GT's sample labels. Given that the NSC dataset

693    lacked ground-truth sample labels, we could not definitively determine which cells were true

694    doublets; however, the increase in clustering ARI after removal of Ensemblex's putative doublets

695    (**Figure 6D**), coupled with Ensemblex's improved doublet identification performance on pools

696    with known ground-truth sample labels (**Figure 2B**), afforded confidence to assume that our

697    ensemble method performed favorably. Nonetheless, this analysis reveals that the choice of

698    demultiplexing tool can greatly impact biological analyses.

699

700    **Conclusion**

701    Multiplexing protocols, coupled with the introduction of genetic demultiplexing tools constituted

702    a significant advancement for scRNAseq by providing a feasible means to dramatically increase

703    the throughput of biological replicates. As the demand for population-scale scRNAseq analysis

704    continues to grow with the maturation of singe-cell technologies, the prospect of multiplexing

705    entire cohorts has emerged. However, the realization of this goal is impeded by the limitations of

706    the current genetic demultiplexing tools. These include decreasing demultiplexing performance as

707    the number of multiplexed samples increases (9, 10), relatively poor doublet detection

708    performance (10), relatively high rates of cells that can only be correctly classified by single

709    algorithms, the unnecessary removal of correctly classified cells due to insufficient assignment

710    probabilities, and highly variable demultiplexing performance between datasets (10). In this work

711    we presented Ensemblex, which offers a unique solution to these limitations by meticulously

712    implementing distinct demultiplexing algorithms into a robust, accuracy-weighted ensemble

713    framework that is exceptionally equipped to classify highly multiplexed pools.

714

715    We applied Ensemblex to a diverse array of computationally and experimentally multiplexed

716    scRNAseq datasets. Benchmarking analyses on pools with known ground-truth sample labels

717    revealed Ensemblex's superior demultiplexing performance across pools reaching 80 multiplexed

718    samples, which translated to a higher proportion of cells retained for downstream analyses and

719    lower error rates amongst classified cells. Ensemblex also demonstrated a notable advancement

720    for identifying heterogenic doublets, which is a well-documented limitation of the genetic

721    demultiplexing tools currently available (9, 10, 15). While previous analyses indicated that the

722    number of multiplexed samples in a pool directly impacted doublet detection efficiency (15), we

723    showed that Ensemblex's ability to identify doublets remained relatively constant when >24

724    samples were multiplexed. Our findings suggest that super loading cells prior to sequencing —

725    which will result in a higher number of usable cells but a higher a doublet rate (6) — followed by

726    heterogenic doublet detection by Ensemblex, may be a viable approach for implementing

34

727     population-scale multiplexing in practice. We also demonstrated that the performance of individual

728     genetic demultiplexing tools can be highly dataset-dependent, reflecting the findings of previous

729     work (10). However, due to its unsupervised weighting model, we showed that Ensemblex is

730     resistant to poorly performing constituent tools, maximizing the consistency of its demultiplexing

731     performance. Nonetheless, if each constituent tool performs poorly on a given dataset, the poor

732     performance will be reflected in Ensemblex's demultiplexing accuracy. Finally, we illustrated that

733     discordant sample assignments amongst genetic demultiplexing tools can greatly impact DGE

734     analyses, necessitating that investigators carefully consider their choice of genetic demultiplexing

735     tool. Although untested, we anticipate that the impacts of discordant sample assignments amongst

736     genetic demultiplexing tools on biological interpretations would be exacerbated for computational

737     analyses that consider the specific donor identity of the pooled cells, such as expression

738     quantitative trait loci (eQTL) analyses, as opposed to donor groups (i.e., case and control). Due to

739     Ensemblex's ability to seamlessly integrate multiple algorithms into an adaptable framework, we

740     argue that our ensemble method achieves unmatched reliability for experimentally multiplexed

741     pools that lack ground truth sample labels.

742

743     Undoubtedly, a limitation of utilizing an ensemble method for genetic demultiplexing is the

744     necessity to run each individual demultiplexing algorithm, which can be computationally

745     expensive. Yet, in the absence of comparing demultiplexed sample labels across tools, poor

746     performance by a given individual algorithm on experimentally multiplexed pools is undetectable,

747     and the risk of introducing technical artifacts and losing usable cells for downstream analyses is

748     prominent. As such, we believe that the relatively high computational cost of Ensemblex is a

749     worthwhile investment to maximize the biological insight obtained from multiplexed scRNAseq

35

750    datasets. To mitigate the burden of genetic demultiplexing by multiple individual tools, we provide

751    a coherent pipeline that runs each constituent demultiplexing tool in parallel and seamlessly

752    processes the respective output files with the Ensemblex algorithm.

753

754    Compared to when demultiplexing was informed by prior genetic data of the pooled samples, the

755    improvement of Ensemblex over its constituent tools was far less pronounced for genotype-free

756    demultiplexing cases. All demultiplexing tools, including Ensemblex, showed drops in

757    demultiplexing performance when >16 samples were multiplexed in a pool without prior genotype

758    information. Nonetheless, Ensemblex still constitutes an advancement over the individual tools for

759    genotype-free demultiplexing cases due to the robustness achieved by incorporating distinct

760    demultiplexing algorithms, which protects against the prospect of poorly performing individual

761    tools on the respective dataset. Furthermore, an intrinsic limitation of demultiplexing without prior

762    genotype information is that samples cannot be directly linked to metadata, leaving the sample

763    identity of the inferred clusters unresolved (9). Although challenging, this limitation can be

764    mitigated by identifying a small subset of discriminatory variants from the reconstructed genotypes

765    of the constituent demultiplexing tools, which could be used to manually assign the computed

766    clusters to samples if such discriminatory variants are known by the investigator. While the

767    Ensemblex pipeline provides users the option to demultiplex pools with or without prior genotype

768    information, we assert that users take caution when electing to perform population-scale

769    multiplexing experiments without using prior genetic data.

770

771    Genetic demultiplexing tools have been used extensively for scRNAseq analysis across many

772    disciplines in the biological sciences, including microbiology (8), model organisms (15), cancer

36

773     biology (23), and neurodegenerative disease (12). Recent work has also evaluated the utility of

774     genetic demultiplexing tools for different single-cell, read-based modalities such as single-nuclei

775     RNA sequencing (snRNAseq) and single-nuclei assay for transposase-accessible chromatin

776     sequencing (scATACseq) (24). Although untested, we expect Ensemblex to prove beneficial in

777     demultiplexing for these assays, but comprehensive benchmarking with the appropriate datasets is

778     required and was not explored here.

779

780     We expect numerous biological fields to exploit the benefits of Ensemblex through its application

781     to highly multiplexed pools comprising cells from many genetically distinct individuals.

782     Specifically for biomedical sciences, the preparation and labour costs of scRNAseq remains

783     prohibitively expensive for analyzing entire cohorts of patients, which is critical for characterizing

784     the genetic heterogeneity and etiological diversity of disease, and for maintaining sufficient

785     statistical power for detecting associations between transcriptional changes and clinical or

786     pathological observations (1). By increasing the throughput of biological replicates, multiplexing

787     has rendered the prospect of analyzing entire patient cohorts with single-cell transcriptomics

788     feasible. Highly-multiplexed scRNAseq experiments have already been presented in the literature

789     and, to the best of our knowledge, have pooled up to 24 samples in a single dish (12). However,

790     we demonstrated that Ensemblex's demultiplexing accuracy remains relatively constant when >24

791     samples are multiplexed at concentrations that abide by the current limitations of experimental

792     protocols, suggesting that Ensemblex equips the research community with the necessary

793     computational framework to expand the upper limits of the number of genetically distinct

794     individuals in a single pool.

795

37

796    While multiplexing mitigates the labour and consumable costs of scRNAseq analysis, the cost of

797    sequencing remains expensive and the increasing number of genetically distinct individuals in a

798    single pool necessitates that a greater number of cells must be sequenced to ensure adequate

799    representation. Accordingly, Ensemblex is equipped to demultiplex pools comprising cells from

800    more genetically distinct individuals than is feasible with the current laboratory technologies.

801    However, we expect that the cost of sequencing will continue to decrease with the maturation of

802    the technology, and our tool will be in place for when the anticipated wet lab advancements are

803    realized. Overall, we conclude that Ensemblex constitutes a notable advancement towards the

804    pressing demand for population-scale single-cell transcriptomics.

805

806    **Methods**

807    **Ensemblex framework overview**

808    Ensemblex is an ensemble genetic demultiplexing framework for scRNAseq sample pooling that

809    was designed to identify the most probable sample labels from each of its constituent tools:

810    Demuxalot (5), Demuxlet (6), Souporcell (8), and Vireo (9) when demultiplexing with prior

811    genotype information or Demuxalot, Freemuxlet (6), Souporcell, and Vireo when demultiplexing

812    without prior genotype information. After running each constituent demultiplexing tool in parallel,

813    Ensemblex merges the output files containing the sample-cell assignments from each tool and

814    performs three distinct steps of the Ensemblex pipeline:

815        1.  Accuracy-weighted probabilistic ensemble;

816        2.  Graph-based doublet detection;

817        3.  Ensemble-independent doublet detection.

818     Upon obtaining the final Ensemblex sample labels (donor-of-origin identity of the pooled cells),

819     the singlet assignment confidence score is computed.

820

821     ***Step 1: Accuracy-weighted probabilistic ensemble***

822     Ensemblex utilizes an unsupervised weighting model to identify the most probable sample

823     label for each cell. Ensemblex weighs each constituent tool's assignment probability

824     distribution by its estimated balanced accuracy for the dataset in a framework adapted from

825     the work of Large et al. (16). To estimate the balanced accuracy of a particular constituent tool

826     (e.g., Demuxalot) for experimentally multiplexed datasets lacking ground-truth labels,

827     Ensemblex uses the cells with a consensus assignment across the three remaining tools (e.g.,

828     Demuxlet, Souporcell, and Vireo-GT) as a proxy for ground-truth. The balanced accuracy for

829     each tool is calculated using equation 1:

830

831         $(1)\ Balanced\ accuracy = \frac{1}{2}\left(\left(\frac{TP}{TP+FN}\right)+\left(\frac{TN}{TN+FP}\right)\right)$

832

833     Where TP is the number of correctly classified singlets; true-negative (TN) is the number of

834     correctly classified doublets; FP is the number of incorrectly classified singlets; false- negative

835     (FN) is the number of incorrectly classified doublets. The probability distribution of each

836     constituent tool $(\hat{p}_j)$ is then weighted by its estimated balanced accuracy $(w_j)$ to produce an

837     accuracy-weighted ensemble probability for each cell:

838

839         $(2)\ \hat{p}(y = i|E) \propto \sum_{j=1}^{k} w_j \hat{p}_j(y = i|M_j)$

840

39

841      Where $\hat{p}$ is the probability that a barcode belongs to class $i$; $y$ is the class variable with $c$

842      possible values, $y \in (1, \ldots, c)$; $c$ is the number of pooled samples plus 1 to account for

843      doublets; $E$ is a vector of the results of $M$ classifiers, $E = (M_1, \ldots, M_k)$; $M$ is the individual

844      constituent demultiplexing output from each tool. Given $\hat{p}$, Ensemblex assigns each barcode's

845      sample identity ($\hat{y}$) as the class (sample label) with the maximum probability:

846

847          (3) $\hat{y} = \arg max_{i \in (1, \ldots, c)}\, \hat{p}(y = i | E)$

848

849      ***Step 2: Graph-based doublet detection***

850      Ensemblex employs a graph-based approach to identify doublets that are incorrectly labeled as

851      singlets by the accuracy-weighted probabilistic ensemble component (Step 1). For graph-based

852      doublet detection, Ensemblex leverages pre-defined features returned from each constituent

853      tool:

854          1. Demuxalot: doublet probability;

855          2. Demuxlet/Freemuxlet: singlet log likelihood – doublet log likelihood;

856          3. Demuxlet/Freemuxlet: number of single nucleotide polymorphisms (SNP) per cell;

857          4. Demuxlet/Freemuxlet: number of reads per cell;

858          5. Souporcell: doublet log probability;

859          6. Vireo: doublet probability;

860          7. Vireo: doublet log likelihood ratio.

861      For each feature independently, the pooled cells are ordered from the most to the least probable

862      doublet and are assigned a percentile rank. Beginning with a percentile threshold of 99.99,

863      Ensemblex screens each cell to identify those that exceed the percentile threshold across all

864 features; cells that exceed the percentile threshold across all features are labeled as "confident

865 doublets". For each iteration, Ensemblex decreases the percentile threshold by 0.01 and repeats

866 the screening process until it has identified $n$ confident doublets (nCD). Ensemblex performs

867 a parameter sweep to determine the optimal nCD to use for graph-based doublet detection (see

868 below).

869

870 Next, the above features are input into a PCA using the *stats* (v3.6.2) R package (25) and a

871 Euclidean distance matrix is generated from the first two principal components (PC). For each

872 confident doublet independently, the remaining cells in the pool are assigned a percentile rank

873 based on their proximity in Euclidean space to the confident doublet and the cells that exceed

874 the designated nearest neighbour percentile threshold (pT) are identified. For all cells that

875 exceeded the designated pT for any confident doublet (putative doublets), Ensemblex

876 computes the number of times the putative doublet was amongst the nearest neighbours of any

877 confident doublet (fNN); an fNN equal to nCD indicates that a putative doublet was amongst

878 the top nearest neighbours for each confident doublet.

879

880 To optimize the nCD and pT parameters for experimentally pooled samples lacking ground-

881 truth labels, Ensemblex performs an automated parameter sweep at each pairwise combination

882 of nCD and pT values; nCD values range from 50 to 300, in increments of 50, while pT values

883 depend on the expected doublet rate (exDR) and range from $1 - \frac{exDR}{6}$ to $1 - exDR$, in

884 intervals of $\frac{1-exDR}{6}$. The distribution of fNN values for each combination of nCD and pT

885 parameters are plotted and Pearson's measure of kurtosis (K), is used to predict which

886 combination of pT and nCD values optimize the identification of true doublets while

41

887 minimizing the rate of incorrectly labelled true singlets as doublets. Ensemblex screens for

888 combinations of nCD and pT values that result in negatively skewed fNN distributions with

889 high K, signifying high peakedness and heavy tails. High peakedness indicates that cells

890 exceeding the designated pT concentrated around nCD, reflecting their proximity in Euclidean

891 space to all high confident doublets, while heavy tails indicate that even cells with lower fNN

892 values were identified as nearest neighbour to many confident doublets. Ensemblex first

893 identifies the pT that returns the highest K, on average, across nCD values tested in the

894 parameter sweep using equation 4:

895

896 $$(4)\ \widehat{pT} = \arg max_{pT \in \left\{1 - \frac{exDR}{6},....,1-exDR\right\}} \left(\frac{\sum_{nCD \in \{50,100,150,200,250,300\}} K(y=pT)}{2}\right)$$

897

898 Where K of the distribution of fNN values of the putative doublets is defined as:

899

900 $$(5)\ K(\text{fNN}) = E\left[\left(\frac{X-\mu}{\sigma}\right)^4\right]$$

901

902 Where $\mu$ is the mean of the distribution and $\sigma$ is the standard deviation. Upon identifying the

903 optimal pT value ($\widehat{pT}$), Ensemblex plots the K corresponding to $\widehat{pT}$ across all nCD values

904 tested in the parameter sweep. If an inflection point is identifiable, Ensemblex identifies $\widehat{nCD}$

905 as the nCD value corresponding to the point of inflection on the curve. Otherwise, Ensemblex

906 identifies $\widehat{nCD}$ as the nCD value corresponding to the highest K. Cells flagged as putative

907 doublets identified using $\widehat{pT}$ and $\widehat{nCD}$ are labelled as doublets by Ensemblex.

908

42

909    *Step 3: Ensemble-independent doublet detection*

910    Benchmarking on computationally multiplexed pools with known ground-truth sample labels

911    revealed that certain genetic demultiplexing tools, namely Demuxalot and Vireo, showed high

912    doublet detection specificity, but that Steps 1 and 2 of the Ensemblex workflow failed to

913    correctly label a subset of doublet calls by these tools. To mitigate this issue and maximize the

914    rate of doublet identification, Ensemblex labels the cells that are identified as doublets by Vireo

915    or Demuxalot as doublets by default; however, users can nominate different tools for the

916    ensemble-independent doublet detection component depending on the desired doublet

917    detection stringency. Doublet specificity was computed using equation 6:

918

919    $$(6)\ Doublet\ specificity = \left(\frac{TN}{TN+FP}\right)$$

920

921    Where TN is the number of correctly classified doublets; FP is the number of true singlets

922    incorrectly classified as doublets.

923

924    *Ensemblex singlet assignment confidence score*

925    Ensemblex computes a singlet confidence score to inform which cells should be discarded to

926    avoid misclassification in downstream analyses. First, Ensemblex evaluates how well an

927    individual constituent tool's assignment probability (e.g., Demuxalot) corresponded to the

928    accuracy of their assignment, using consensus cells across the three remaining tools (e.g.,

929    Demuxlet, Souporcell, Vireo) as a proxy for ground-truth, by fitting a binary logistic regression

930    model to compute the odds that a singlet was correctly classified given its corresponding

931    probability. Using the binary logistic regression models, Ensemblex computes the AUC using

43

932    the empirical method implemented in the *ROCit* (v2.1.1) R package for each tool (26). Then,

933    for each cell, if Ensemblex's sample label matches that of a constituent tool, and if the

934    assignment probability of the constituent tool supersedes its probability threshold, the tool's

935    computed AUC is added to the accuracy-weighted probabilistic ensemble probability produced

936    in Step 1 to yield the confidence score. By default, singlet assignments with a confidence score

937    less than 1.00 are labelled as unassigned by Ensemblex. Ensemblex's confidence score and the

938    designated threshold is a successful predictor of accurately classified singlets because singlets

939    will only achieve a confidence score $\geq 1$ if:

940    1.  All constituent tools show the same sample label (accuracy-weighted probabilistic

941        ensemble probability = 1.00);

942    2.  At least one constituent tool confidently assigns the cell to an individual donor and the

943        constituent tool's probability assignment adequately corresponds to the overall

944        accuracy of their singlet assignment.

945

946    ***Application of Ensemblex with and without prior genotype information***

947    Given the dependencies of certain tools on prior genotype information, there are notable

948    differences between the Ensemblex workflows for demultiplexing with and without prior

949    genotype information. When demultiplexing with prior genotype information, Ensemblex

950    leverages the sample labels from Demuxalot, Demuxlet, and Vireo-GT with prior genotype

951    information, and Souporcell without prior genotype information. When demultiplexing

952    without prior genotype information, Ensemblex leverages the sample labels from Demuxalot,

953    Freemuxlet, Souporcell, and Vireo. However, given that Demuxalot requires prior genotype

44

954   information, Ensemblex uses the estimated donor .vcf file generated by Freemuxlet for input

955   into the Demuxalot algorithm as prior genetic data.

956

957   *Running the Ensemblex pipeline*

958   A complete user guide for running the Ensemblex pipeline can be found at the Ensemblex

959   GitHub site: https://neurobioinfo.github.io/ensemblex/site/. We provide two distinct yet highly

960   comparable pipelines depending on the availability of prior genotype information. Both

961   pipelines can be downloaded as a singularity image and are comprised of four steps:

962       1. Establish the pipeline and working directory;

963       2. Prepare input files for constituent genetic demultiplexing tools;

964       3. Parallel demultiplexing by constituent genetic demultiplexing tools;

965       4. Application of the Ensemblex algorithm for ensemble classification.

966

967   As input into the Ensemblex pipeline, users must provide a .tsv file describing the barcodes of

968   the pooled cells, a. bam sequencing file for the pool, a reference genotype .vcf file (e.g., 1000

969   Genome Project) (27), a reference genome sequence .fasta file (e.g., 10X Genomics), and, if

970   demultiplexing with prior genotype information, a .vcf file describing the genetic data of the

971   pooled samples.

972

973   **Genetic demultiplexing by constituent tools**

974   Genetic demultiplexing by the constituent demultiplexing tools was performed following best

975   practices as defined by the authors of the respective tools using Python (v3.8.10).

976   *Demuxalot*

45

977    CellRanger-generated .bam file, filtered barcode .tsv file, and the corresponding donor .vcf file

978    were used as input into the Demuxalot workflow. Candidate variants for scRNAseq genotyping

979    were retained if the minimum coverage was > 200 and minimum alternative coverage was >

980    10. The top 100 SNPs per donor were retained to cluster the cells by genotype. Doublet calls

981    were made with a prior strength of 0.25.

982

983    ***Demuxlet***

984    We used the popscle suite (https://github.com/statgen/popscle) for Demuxlet. CellRanger-

985    generated .bam file, filtered barcode .tsv file, and the corresponding donor .vcf file were used

986    as input into the Demuxlet workflow. The *dsc-pileup* function was first used to pileup candidate

987    variants around known variant sites with the following parameters: --cp-BQ 40 --min-BQ 13 -

988    -min-MQ 20 --minTD 0 --min-total 0 --min-uniq 0 --min-snp 0. The Demuxlet algorithm was

989    then applied to cluster the cells by genotype with the following parameters: --geno-error-offset

990    0.10 --geno-error-coeff 0.00 --min-callrate 0.50 --doublet-prior 0.50 --cap-BQ 40 --min-BQ 13

991    --min-MQ 20 --min-TD 0 --min-total 0 --min-uniq 0 --min-snp 0.

992

993    ***Freemuxlet***

994    We used the popscle suite (https://github.com/statgen/popscle) for Freemuxlet. CellRanger-

995    generated .bam file, filtered barcode .tsv file, and reference genotype .vcf file from the 1000

996    Genomes Project, phase 3 (27), were used as input into the Freemuxlet workflow. The *dsc-*

997    *pileup* function was first used to pileup candidate variants around known variant sites with the

998    following parameters: --cp-BQ 40 --min-BQ 13 --min-MQ 20 --minTD 0 --min-total 0 --min-

999    uniq 0 --min-snp 0. The Freemuxlet algorithm was then applied to cluster the cells by genotype

46

1000    with the following parameters: --doublet-prior 0.50 --bf-thres 5.41 --frac-init-clust 0.50 --inter-

1001    init 10 --cap-BQ 40 --min-BQ 13 --min-total 0 --min-uniq 0 --min-snp 0.

1002

1003    ***Souporcell***

1004    CellRanger-generated .bam file, filtered barcode .tsv file, 10X Genomics reference .fasta file,

1005    and the corresponding donor .vcf file when demultiplexing with prior genotype information

1006    were used as input into the Souporcell workflow. A FASTQ file was first generated from the

1007    .bam file using the *renamer.py* script. These reads were mapped to the reference genome using

1008    minimap2 with the following parameters: --ax splice –t 8 –G50k –k 21 –w 11 –sr --A2 –B8 –

1009    O12,32 –E2,1 –r200 –p.5 –N20 –f1000,5000 –n2 –m20 –s40 –g200 –2k50m –secondary=no.

1010    The barcodes and UMI were added back to the .sam file using the *retag.py* script and the

1011    resulting .bam file was sorted and indexed with Samtools. Variants were called using Freebayes

1012    with the following parameters: --iXu –C 2 –q 20 –n 3 –E 1 –m 30 –min-coverage 6. Vartix was

1013    used to compute the number of alleles for each cell using the following parameters: --umi –

1014    mapq 30 –scoring-method coverage. The Souporcell algorithm was then applied to cluster the

1015    cells by genotype; when demultiplexing with prior genotype information the --

1016    known_genotypes and --known_genotypes_sample_names parameters were included.

1017    Troublet was used to identify doublets and the *consensus.py* script was used for genotype and

1018    ambient RNA co-inference.

1019

1020    ***Vireo***

1021    CellRanger-generated .bam file, filtered barcode .tsv file, reference genotypes from the 1000

1022    Genomes Project, phase 3 (27), and the corresponding donor .vcf file when demultiplexing

47

1023      with prior genotype information were used as input to the Vireo workflow. CellSNP was used

1024      to identify candidate variants for scRNAseq genotyping with the following parameters: --

1025      minMAF 0.1 and --minCOUNT 100. The Vireo algorithm was then applied to cluster the cells

1026      by genotype with the --forceLearnGT parameter; when demultiplexing with prior genotype

1027      information (Vireo-GT) the --d and --t GT parameters were used.

1028

1029      *Consensus demultiplexing framework*

1030      For the consensus demultiplexing framework, singlets were considered confidently classified

1031      if Demuxalot, Demuxlet, Vireo, and Souporcell assigned a cell to the same donor-of-origin.

1032      Cells classified as "ambiguous" or doublet by at least one tool were discarded.

1033

1034      **Generation of computationally pooled samples for ground-truth benchmarking**

1035 To benchmark Ensemblex on computationally pooled samples with known ground-truth sample

1036 labels, we leveraged 80 independently sequenced iPSC lines from Parkinson's disease patients and

1037 healthy controls, which were differentiated towards a dopaminergic neuronal state and sequenced

1038 after 65 days of differentiation as part of the FOUNDIN-PD (14). Controlled access FASTQ files

1039 from the independently sequenced iPSC lines were obtained from https://www.ppmi-info.org/

1040 (accessed 09-17-2023) and processed by the CellRanger *counts* pipeline (v3.1.0) with default

1041 parameters and aligned to GRCh38 reference genome. The CellRanger-generated .bam and filtered

1042 barcode files were used as input into the *synth_pool.py* script produced by the authors of Vireo to

1043 simulate sample pooling (9). In brief, reads from a subset of cells from the iPSC line-specific .bam

1044 files were merged and doublets were generated by combining the reads from random cell pairs.

48

1045    Sample identities were added to each cell's barcode, revealing the ground-truth sample labels for

1046    benchmarking procedures.

1047

1048    To evaluate how genetic demultiplexing performance varied as a function of the number of

1049    multiplexed samples, we generated 96 computationally multiplexed pools using the 80

1050    FOUNDIN-PD lines with sample sizes of 4, 8, 16, 24, 32, 40, 48, 56, 64, 72, and 80. An equal

1051    number of cells from each line were used in the *in silico* pool. For the sample size of four we

1052    generated six replicates; for the sample sizes of 8-80 we generated nine replicates each. Replicates

1053    were produced with different sample and cell combinations. The 96 *in silico* pools averaged 17,396

1054    cells (minimum = 8,696; maximum = 26,087). For this experiment, we maintained a 15% doublet

1055    rate as previously described (9).

1056

1057    To evaluate how genetic demultiplexing performance varied as a function of the number of cells

1058    in a pool, we generated 18 computationally multiplexed pools using the 80 FOUNDIN-PD lines

1059    with 8,000, 16,000, 24,000, 32,000, 40,000, and 48,0000 pooled cells; we generated three

1060    replicates per pool size. Twenty-four samples were multiplexed for each pool and an equal number

1061    of cells from each sample were used. Replicates were produced with different sample and cell

1062    combinations. For this experiment, we simulated a doublet rate of 6% per 8,000 pooled cells.

1063

1064    To evaluate if the overall demultiplexing performance varied due to the underrepresentation of a

1065    cell line, we generated 15 computationally multiplexed pools using the 80 FOUNDIN-PD lines

1066    comprising 23 multiplexed samples with 1,000 cells and one randomly selected sample that

1067    showed various degrees of underrepresentation, including 100 cells (10%), 300 cells (30%), 500

1068  cells (50%), 700 cells (70%), or 900 cells (90%). Three replicates were generated for each degree

1069  of underrepresentation. Replicates were produced with different sample and cell combinations. For

1070  this experiment, we maintained a 18% doublet rate.

1071

1072  WGS for the 80 donors from which the FOUNDIN-PD lines were derived was performed on whole

1073  blood-extracted DNA as previously described by the Parkinson's Progression Markers Initiative

1074  (PPMI) (28). The controlled-access WGS .vcf files were obtained from https://www.ppmi-

1075  info.org/ (accessed 09-17-2023). Genotypes of common variants (minor allele frequency > 5%)

1076  were used as prior genotype information for the genetic demultiplexing tools in the benchmarking

1077  analyses.

1078

1079  **Preparation, processing, and analysis of experimentally pooled samples**

1080  Unless specified otherwise, experimentally pooled samples were processed with the CellRanger

1081  *counts* pipeline (v5.0.1) and analyzed with the *Seurat* (v5.0.0) R package (29), using the

1082  scRNAbox analytical pipeline (30).

1083

1084  *Non-small cell lung cancer dataset*

1085  NSCLC dissociated tumor cells from seven donors were labelled with TotalSeq-B Human

1086  TBNK Cocktail (18). Multiplexed cells were then sequenced on an Illumina NovaSeq 6000 to

1087  an average read depth of approximately 70,000 reads per cell for gene expression and 25,000

1088  reads per cell for CellPlex. Publicly available gene expression .bam and barcode .tsv files

1089  returned from the CellRanger *multi* pipeline (v6.1.2) were obtained from the 10X Genomics

1090  Datasets portal (10X Genomics Datasets) and used as input into the Ensemblex pipeline. We

1091   used the sample-specific gene expression .bam files and the BCFtools (v1.16) *mpielup*

1092   function to generate genotype likelihoods for prior genotype information (31).

1093

1094   We used HTOdemux to assign the cells back to their donor-of-origin based on the CMO

1095   expression profiles as a proxy for ground-truth sample labels (19). Publicly available feature-

1096   barcode expression matrices were filtered to only include CMO labels used for multiplexing

1097   — CMO301, CMO302, CMO303, CMO304, CMO306, CMO307, and CMO308 — and

1098   barcodes with a CMO count > 0. The CMO expression profiles were normalized with Seurat's

1099   *NormalizeData* function using the CLR normalization method and HTOdemux was applied to

1100   the CMO assay using a positive quantile of 0.99.

1101

1102   ***Dopaminergic neuron dataset***

1103   Jerber et al. sequenced multiplexed experiments comprising 22 healthy donor iPSC lines from

1104   the HipSci project (32) (http://www.hipsci.org) on days 11, 30, and 52 of DaN differentiation

1105   using Illumina HiSeq 4000 to an average depth of 40,000-60,000 reads per cell (12). We used

1106   three technical replicates for each timepoint, which are comprehensively described in

1107   **Additional File 1: Table S3**. Publicly available gene expression .fastq files were obtained from

1108   the European Nucleotide Archive (ENA) with accession number ERP121676 and processed

1109   with the CellRanger *counts* pipeline (v5.0.1) with default parameters using the GRCh37

1110   reference genome. The CellRanger-generated. bam files, filtered barcode .tsv files, and .vcf

1111   files describing the pooled samples (see below) were used as input into the Ensemblex pipeline

1112   for each technical replicate independently. Filtering of the scRNAseq data was performed as

1113   described by Jerber et al. (12). Genes with non-zero counts in at least 0.05% of cells were

51

1114    retained. DoubletFinder (v2.0.4) was applied independently to each technical replicate. Time-

1115    point specific replicates were integrated with Seurat's integration algorithm (33) and clustered

1116    by the Louvain network detection using the top 50 PCs and 10 nearest neighbours.

1117

1118    Whole-exome sequencing (WES) .vcf files corresponding to the 22 pooled HipSci lines were

1119    obtained from the ENA with accession number PRJEB7243 (34). Genotypes of common

1120    variants (minor allele frequency > 1%) were used as prior genotype information for the genetic

1121    demultiplexing tools (12).

1122

1123    *Neural stem cell dataset*

1124    We performed two multiplexed experiments comprising iPSCs from individuals with ADHD

1125    and heathy controls differentiated into NSCs: Experiment 1 (*n* ADHD = 7; *n* control = 6) and

1126    Experiment 2 (*n* ADHD = 9; *n* control = 7).

1127

1128    **Subject recruitment**

1129    Patients diagnosed with ADHD and matching healthy controls between 6−18 years old

1130    were recruited by the Department of Child and Adolescent Psychiatry and Psychotherapy

1131    of the University of Zurich, as described previously (35). Inclusion and exclusion criteria

1132    for recruitment of these individuals described previously (35). **Additional File 1: Table**

1133    **S4** provides a list of the individual subjects and their derived cell lines included in this

1134    study. Salivary DNA from ADHD patients and controls was genotyped using the Infinium

1135    Global Screening Array (Illumina), as previously described, and used as prior genotype

1136    information for genetic demultiplexing (35).

1137

1138 **Neural stem cell culture**

1139 The generation and characterization of iPSC used in this study and the NSCs differentiation

1140 protocols were previously described in (35) (36). NSCs cultures were seeded in two

1141 independent experiments (designated as "1" and "2"), each of them consisting of NSCs

1142 pooled together into two culture dishes and maintained as NSCs until 100% confluence,

1143 when all iPSC lines were combined into one sample for sequencing. For most cell lines

1144 different clones for each iPSC line were used in the two experiments **Additional File 1:**

1145 **Table S5**. When applicable, the second clones of the same NSCs lines were cultured

1146 separately (designated as ".1" and ".2") in a second experiment. In the first experiment,

1147 56,250 cells per cell line were seeded in the pooled dishes. In the second experiment the

1148 proportions of cells seeded we adjusted to their proliferation profile assessed in (36). Upon

1149 reaching 100% confluence, cells were dissociated for scRNAseq experiments and

1150 combined to a single sample for sequencing as described below.

1151

1152 **Dissociation of pooled neural stem cell cultures for single-cell RNA sequencing**

1153 Cells were washed in PBS and then incubated with 1 mL of StemPro Accutase (Gibco) for

1154 3 minutes at 37°C. After incubation, 2 mL of PBS, stopping the Accutase reaction, and cells

1155 were gently pipetted up and down between 5 to 10 times to break up clumps of cells before

1156 transfer to a 15 mL conical tube. The cells were centrifuged at 300 x g for 5 minutes and

1157 the supernatant was removed. Following, 334 μL of Neural Expansion Media (NEM) was

1158 added to each cell pellet using a 1000 μL pipette tip until cells were completely

1159 resuspended. An additional 666 μL of NEM was added to each well and gently pipette

53

1160    mixed 5 times. A 100-µm cell strainer was used to filter the cell suspension before

1161    centrifugation at 300 x g for 4 minutes. The supernatant was carefully removed, and the

1162    pellet was resuspended in 3 mL of PBS 1x containing 0.04% Bovine Serum Albumin

1163    (BSA) by pipetting up and down 5 times using a 5 mL serological pipette. The cells were

1164    centrifuged at 300 x g for 10 minutes and further submitted to live cell sorting with the

1165    Magnetic Dead Cell Removal Kit (Miltenyi Biotec, 130-090-101), according to the

1166    manufacturer. The resulting flow-through containing live cells was centrifuged for 300 x g

1167    for 5 minutes and the supernatant was removed carefully to not disturb the cell pellet. Cells

1168    were resuspended in 1 mL of PBS 1x containing 0.04% BSA for automated cell counting.

1169    For each experiment, the cells from the two culture dishes were processed in parallel. Equal

1170    counts of cells were combined for the final cell suspension for scRNAseq preparation at

1171    the Functional Genomics Center Zurich at the University of Zurich.

1172

1173    **Library processing and sequencing**

1174    All samples were processed using the 10x Genomics Chromium 3' Single Cell Protocol

1175    and sequenced using NovaSeq 6000 S1 (Illumina). For the first sample containing NSC

1176    pools 1.1 and 1.2, 18,000 NSCs were loaded into one single 10x Genomics Lane to target

1177    13,000 cells. For the second sample containing NSC pools 2.1 and 2.2, 29,000 NSCs were

1178    loaded to target 18,000 cells.

1179

1180    **Demultiplexing and scRNAseq analysis**

1181    FASTQ files were processed with the CellRanger *counts* pipeline (v5.0.1) with default

1182    parameters and aligned to the GRCh37 reference genome. The CellRanger-generated. bam

54

1183    files, filtered barcode .tsv files, and .vcf files describing the pooled samples were used as

1184    input into the Ensemblex pipeline. Genotypes of common variants (minor allele frequency

1185    > 1%) were used as prior genotype information for the genetic demultiplexing tools. The

1186    filtered feature-barcode expression matrices were used to analyze the pooled cells

1187    following a standard scRNAseq analysis workflow using Seurat (30). Cells were filtered

1188    for > 500 total and unique RNA transcripts. Doublets were removed using DoubletFinder

1189    (v2.0.4). The two NSC samples were integrated using Seurat's integration algorithm (33).

1190    The top 25 PCs were selected for Louvain network detection to identify clusters using 65

1191    nearest neighbours. Twelve clusters were identified at a clustering resolution of 0.25, which

1192    were assigned as eight putative cell types using a combination of known markers and gene

1193    enrichment analysis. The top marker genes from each cluster were identified using Seurat's

1194    *FindAllMarkers* with the Wilcoxon rank-sum test. Significant DEGs (log2 fold change >

1195    0.25 and P-value < 0.05 ) were input into EnrichR (37) and cell types were predicted with

1196    the *Cell Marker Augmented 2021* (38) and *Azimuth Cell Types 2021* (39) libraries. Multiple

1197    clusters showed expression profiles for similar broad cell types — Neurons, NPCs, and

1198    NSCs. We used Seurat's *FindMarkers* function to identify differentially expressed marker

1199    genes between the clusters of the same broad cell type and top marker genes were selected

1200    to identify the cell subtypes.

1201

1202    For each putative cell type, DGE was calculated between ADHD and controls using the

1203    MAST statistical framework (22, 40). Pooled cells were assigned as ADHD or control

1204    based on the demultiplexed sample labels from each of the individual genetic

1205    demultiplexing tools. Cells labeled as "ambiguous singlets" or doublets by the individual

1206      tools were excluded from their respective DGE analysis. P-values were corrected for

1207      multiple hypothesis testing using the Bonferroni method. A gene was considered

1208      differentially expressed if the adjusted P-value was $\leq 0.01$ and the absolute value of the

1209      Log2 fold-change was $\geq 0.5$. To compute DGE using the sample labels from the individual

1210      tools after the removal of Ensemblex's putative doublet calls, we repeated the above

1211      procedures but this time all cells labeled as doublets by the respective tool or Ensemblex

1212      were excluded from the DGE analysis.

1213

1214 **Performance metrics and statistical analyses**

1215 We performed all statistical analyses using the R statistical software (v4.2.2) (41). We used the

1216 *ggplot2* R package (v3.4.2) for data visualization (42).

1217

1218      ***Singlet classification***

1219 A singlet was considered correctly classified if the demultiplexed sample label matched the

1220 ground-truth sample label (i.e., specific sample ID) and the assignment probability exceeded

1221 the recommended threshold for the respective tool. For computationally multiplexed pools, the

1222 proportion of correctly classified singlets was computed as:

1223

1224      $(7)\ Proportion\ correct\ singlets = \dfrac{TP}{n\ true\ singlets}$

1225

1226 For the NSCLC dataset, HTOdemux's sample labels were considered ground-truth, and the

1227 singlet TP and FP rate were computed as:

1228

56

1229 $$(8)\ Singlet\ TP\ rate = \frac{TP}{n\ HTOdemux\ singlets}$$

1230 $$(9)\ Singlet\ FP\ rate = \frac{FP}{n\ HTOdemux\ singlets}$$

1231

1232 ***Doublet identification***

1233 A doublet was considered correctly classified if the demultiplexed sample label matched the

1234 ground-truth sample label, independent of the assignment probability. For computationally

1235 multiplexed pools, the proportion of correctly classified doublets was computed as:

1236

1237 $$(10)\ Proportion\ correct\ doublets = \frac{TN}{n\ true\ doublets}$$

1238

1239 For the NSCLC dataset, TP doublets were defined as cells classified as doublets by both

1240 HTOdemux and Ensemblex; FP doublets were defined as cells classified as singlets by

1241 HTOdemux and doublets by Ensemblex; FN doublets were defined as cells classified as

1242 doublets by HTOdemux and singlets by Ensemblex. The doublet TP, FP, and FN rates were

1243 computed as:

1244

1245 $$(11)\ Doublet\ TP\ rate = \frac{TP}{n\ HTOdemux\ doublets}$$

1246 $$(12)\ Doublet\ FP\ rate = \frac{FP}{n\ pooled\ droplets}$$

1247 $$(13)\ Doublet\ FN\ rate = 1 -\ Doublet\ TP\ rate$$

1248

1249 ***Adjusted Rand Index***

1250    To evaluate the similarity between two distinct sample clusterings we computed the ARI using

1251    the *pdfCluster* (v1.0.4) R package (43). For the benchmarking analyses, we computed the ARI

1252    between the demultiplexed sample labels by each genetic demultiplexing tool and the ground-

1253    truth sample labels (computationally pooled samples) or HTOdemux's sample labels (NSCLC

1254    dataset). We followed the same procedure when computing the ARI between Ensemblex's

1255    sample labels and those of its constituent tools (DaN and NSC datasets); however, the ground-

1256    truth sample labels were replaced by Ensemblex's sample labels for these analyses. For

1257    experiments evaluating the impact of doublets on the stability of clusters in gene expression

1258    space, we computed the ARI between clusters at a given clustering resolution after removing

1259    doublets identified by each genetic demultiplexing tool. Clustering stability was computed at

1260    resolutions of 0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, and 1.0. For each clustering

1261    resolution, 25 iterations of Louvain clustering were performed while shuffling the order of the

1262    nodes in the graph. The ARI between clustering pairs at each clustering resolution was then

1263    computed.

1264

1265    ***Balanced accuracy***

1266    Balanced accuracies were computed to evaluate the binary classification performance of each

1267    genetic demultiplexing tool on imbalanced datasets, where doublets represented a minority

1268    class compared to singlets. The balanced accuracy of each genetic demultiplexing tool was

1269    computed against the ground-truth sample labels (computationally pooled samples) or

1270    HTOdemux's sample labels (NSCLC dataset) using equation 1.

1271

1272    ***Matthew's correlation coefficient (MCC)***

58

1273    The MCC was used as a second metric for evaluating the binary classification performance of

1274    the genetic demultiplexing tool. The MCC of each genetic demultiplexing tool was computed

1275    against the ground-truth sample labels (computationally pooled samples) or HTOdemux's

1276    sample labels (NSCLC dataset) using equation 14:

1277

1278    $$(14)\ MCC = \frac{TN \times TP - FN \times FP}{\sqrt{(TP+FP)(TP-FN)(TN+FP)(TN+FN)}}$$

1279

1280    ***Area under the receiver operating characteristic curve for singlet detection***

1281    To evaluate how well each genetic demultiplexing tool's assignment probability corresponded

1282    to the accuracy of their singlet assignments when ground-truth sample labels were known, we

1283    fit a binary logistic regression model to compute the odds that a singlet was correctly classified

1284    by a tool given the corresponding confidence score or probability. Correctly and incorrectly

1285    classified singlets were set as the positive and negative references, respectively. We then used

1286    the binary logistic regression model to compute the receiver operating characteristic curve for

1287    each tool, which plots the singlet TP and FP rates across classification thresholds, and

1288    calculated the AUC using the empirical method implemented in the *ROCit* (v2.1.1) R package

1289    (26).

1290

1291    **Abbreviations**

1292    ADHD, attention deficit hyperactivity disorder; ANOVA, Analysis of variance; ARI, Adjusted

1293    Rand Index; AUC, area under the receiver operating characteristic curve; BSA, Bovine Serum

1294    Albumin; CMO, Cell Multiplexing Oligonucleotides; DaN, dopaminergic neurons; DGE,

1295    differential gene expression; DEG, differentially expressed genes; ENA, European Nucleotide

59

1296 Archive; eQTL, expression quantitative trait loci; FN, false-negative; fNN, nearest neighbour

1297 frequency; FOUNDIN-PD; Foundational Data Initiative for Parkinson's Disease; FP, false

1298 positive; iPSC, induced pluripotent stem cell; K, kurtosis; MAST, model-based analysis of single-

1299 cell transcriptomics; MCC, Matthew's Correlation Coefficient; nCD, number of confident

1300 doublets; NEM, neural expansion media; NPC, neural progenitor cell; NSC, neural stem cell;

1301 NSCLC, non-small cell lung cancer; PC, principal component; PCA principal component analysis;

1302 PPMI, Parkinson's Progression Markers Initiative; pT, nearest neighbour percentile threshold;

1303 scATACseq, single-cell assay for transposase-accessible chromatin sequencing; scRNAseq,

1304 single-cell RNA sequencing; SNP, single nucleotide polymorphism; snRNAseq, single-nuclei

1305 RNA sequencing; TN, true-negative; TP, true-positive; UMI, unique molecular identified; WES,

1306 whole-exome sequencing; WGS, whole-genome sequencing.

1307

1308 **Declarations**

1309 ***Ethics approval and consent to participate***

1310 The iPSC lines (ADHD & controls) used in this project were approved by the Cantonal Ethics

1311 Committee Zurich (BASEC-Nr.-2016-00101 & BASEC-Nr.-201700825) and followed the latest

1312 version of the Declaration of Helsinki, as previously reported (35). The subjects and/or parents

1313 have voluntarily consented to participate in this study.

1314

1315 ***Consent for publication***

1316 Not applicable.

1317

1318 ***Availability of data and materials***

60

1319 Transcriptional data for the 80 independently sequenced iPSC lines and the corresponding WGS

1320 data are available from the PPMI database (www.ppmi-info.org/access-dataspecimens/download-

1321 data), RRID:SCR 006431. For up-to-date information on the study, visit www.ppmi-info.org.

1322 Processed transcriptional data for the NSCLC dataset are available from the 10X Genomics

1323 Datasets Portal (https://www.10xgenomics.com/datasets/20k-mixture-of-nsclc-dtcs-from-7-

1324 donors-3-v3-1-with-intronic-reads-3-1-standard). Transcriptional data for the DaN datasets are

1325 available from the ENA with accession number ERP121676. WES data for the 22 HipSci lines

1326 pooled in the DaN datasets are available from the ENA with accession number PRJEB7243.

1327 Processed scRNAseq data for the NSC dataset are available from the corresponding author upon

1328 reasonable request. The code used for the analyses presented in the work are available at

1329 https://github.com/neurobioinfo/ensemblex. Ensemblex is freely available under an MIT open-

1330 source license at https://zenodo.org/records/11639103.

1331

1332 ***Competing interests***

1333 The authors declare that they have no competing interests.

1334

1335 ***Funding***

1353

1354     ***Authors' contributions***

1355     MRF, EAF, RAT, and SMKF conceived the study. MRF developed the Ensemblex framework and

1356     wrote the corresponding R code. MRF performed the analyses and produced the figures. MRF and

1357     SA developed the Ensemblex pipeline and created the GitHub site. MRF and SA tested the

1358     Ensemblex pipeline. MRF wrote the Ensemblex documentation. MRF, AAD, RAT and SMKF

1359     interpreted the datasets. CMYO performed all cell cultures and sequencing preparation for the

1360     NSC dataset. MRF, CMYO, and RAT performed the cell type annotations for the NSC dataset. LS

1361     and EG provided the NSC genetic data. SW recruited the subjects for the NSC dataset. MRF wrote

1362     the manuscript with input from all authors. EG supervised the NSC data collection. RAT and

1363     SMKF supervised the project.

1364

1380

**Figure legends**

1381    **Figure 1. Evaluation of existing individual genetic demultiplexing tools.** Evaluation of genetic

1383    demultiplexing tools with prior genotype information on 96 *in silico* pools with known ground-

1384    truth sample labels ranging in size from 4 to 80 multiplexed induced pluripotent stem cell (iPSC)

1385    lines from genetically distinct individuals, averaging 17,396 cells per pool and a 15% doublet rate.

1386    **A)** Line graphs showing the proportion of correctly classified singlets, doublets, and all cells by

1387    each individual genetic demultiplexing tool across varying numbers of multiplexed iPSC lines in

1388 a single pool (sample number). The large dots show the mean proportion of correct classifications

1389 by an individual tool across replicates at a given sample size (n = 9 per pool size). The blue points

1390 show the proportion of cells that were correctly classified by at least one individual genetic

1391 demultiplexing tool: Demuxalot, Demuxlet, Souporcell, or Vireo-GT. **B)** Bar chart showing the

1392 mean proportion of total cells from an individual pool correctly classified by only one genetic

1393 demultiplexing tool. Error bars represent one standard deviation from the mean. (n = 9 per pool

1394 size) **C)** Bar chart showing the proportion of correctly classified singlet cells labelled as

1395 "unassigned" (ambiguous singlet assignments) due to assignment probabilities below the

1396 recommended threshold of the respective genetic demultiplexing tool. Error bars represent one

1397 standard deviation from the mean. (n = 9 per pool size).

1398

1399 **Figure 2. Characterization of the Ensemblex framework.** Ensemblex is a probabilistic-

1400 weighted ensemble genetic demultiplexing framework for single-cell RNA sequencing analysis,

1401 which was designed to leverage the most probable sample labels from each of its constituent tools:

1402 Demuxalot, Demuxlet, Souporcell, and Vireo when using prior genotype information or

1403 Demuxalot, Freemuxlet, Souporcell, and Vireo when prior genotype information is not available.

1404 **A)** The Ensemblex workflow begins with demultiplexing pooled cells from genetically distinct

1405 individuals by each of the constituent tools. The outputs from each individual demultiplexing tool

1406 are then used as input into the Ensemblex framework. **B)** The Ensemblex framework comprises

1407 three distinct steps that are assembled into a pipeline: 1) accuracy-weighted probabilistic ensemble,

1408 2) graph-based doublet detection, and 3) ensemble-independent doublet detection. **C-D)** Line

1409 graphs showng the contribution of each step of the Ensemblex framework on 96 *in silico* pools

1410 with known ground-truth sample labels ranging in size from 4 to 80 multiplexed induced

1411    pluripotent stem cell (iPSC) lines from genetically distinct individuals, averaging 17,396 cells per

1412    pool and a 15% doublet rate. The average proportion of correctly classified **C)** singlets and **D)**

1413    doublets across replicates at a given pool size is shown after sequentially applying each step of the

1414    Ensemblex framework with prior genotype information (n = 9 per pool size). The right panels

1415    show the average proportion of correct classifications across all 96 pools; error bars represent one

1416    standard deviation from the mean. The blue points show the proportion of cells that were correctly

1417    classified by at least one individual genetic demultiplexing tool: Demuxalot, Demuxlet,

1418    Souporcell, or Vireo-GT.

1419

1420    **Figure 3. Ensemblex ground-truth benchmarking on computationally multiplexed pools.** The

1421    genetic demultiplexing tools with prior genotype information were evaluated on 96 *in silico* pools

1422    with known ground-truth sample labels ranging in size from 4 to 80 multiplexed induced

1423    pluripotent stem cell (iPSC) lines from genetically distinct individuals, averaging 17,396 cells per

1424    pool and a 15% doublet rate. A singlet was considered correctly classified if the assigned sample

1425    label matched the ground-truth sample label and the assignment probability exceeded the

1426    recommended threshold for the respective tool; a doublet was considered correctly classified if the

1427    assigned sample label matched the ground-truth sample label, regardless of the assignment

1428    probability. **A-I)** Line graphs showing the performance of Ensemblex and the individual genetic

1429    demultiplexing tools across evaluation metrics. The large dots show the mean value for each tool

1430    across replicates at a given sample size (n = 9 per pool size). **A)** Proportion of correctly classified

1431    singlets. **B)** Proportion of correctly classified doublets. **C)** Proportion of correctly classified cells.

1432    **D)** Adjusted Rand Index between each tool's sample labels and the ground-truth sample labels. **E)**

1433    Balanced accuracy of each tool. **F)** Matthew's Correlation Coefficient of each tool. **G)** Area under

1434    the receiver operating characteristic curve (AUC) of the singlet assignment probability for each

1435    tool. **H)** Proportion of usable cells returned by each tool. Usable cells were defined as cells

1436    classified by singlets with an assignment probability exceeding the recommended threshold of the

1437    respective tool. **I)** Error rate amongst the usable cells returned by each tool; erroneous

1438    classifications comprised of true doublets labeled as singlets or true singlets assigned to the wrong

1439    sample.

1440

1441    **Figure 4. Evaluating Ensemblex on experimentally multiplexed cells using donor-specific**

1442    **oligonucleotide labels as a proxy for ground-truth.** Non-small cell lung cancer (NSCLC)

1443    dissociated tumor cells from 7 individuals were pooled and labelled with donor-specific

1444    oligonucleotide-labels. Cells were demultiplexed according to their expression of donor-specific

1445    oligonucleotide labels by HTOdemux; HTOdemux's sample labels were used as a proxy for

1446    ground truth. True positives (TP) singlets were defined as cells classified as singlets by both

1447    HTOdemux and Ensemblex with matching sample labels; false positives (FP) singlets were

1448    defined as cells classified as singlets by both HTOdemux and Ensemblex but assigned to different

1449    donors. TP doublets were defined as cells classified as doublets by both HTOdemux and

1450    Ensemblex; FP doublets were defined as cells classified as singlets by HTOdemux and doublets

1451    by Ensemblex; false negatives (FN) doublets were defined as cells classified as doublets by

1452    HTOdemux and singlets by Ensemblex. **A)** T-distributed Stochastic Neighbor Embedding (t-SNE)

1453    visualization of HTOdemux's sample labels. **B)** T-SNE visualization of Ensemblex's

1454    demultiplexing performance using HTOdemux's sample labels as ground truth for singlets (left)

1455    and doublets (right). **C)** Bar plots showing the singlet TP and FP rates for each genetic

1456    demultiplexing tool using HTOdemux's sample labels as ground truth. **D)** Bar plots showing the

66

1457     doublet TP and FP rates for each genetic demultiplexing tool using HTOdemux's sample labels as

1458     ground truth. **E)** Scatter plot showing the proportion of usable cells (confidently classified singlets)

1459     and the corresponding usable cell error rate for each genetic demultiplexing tool. **F)** Adjusted Rand

1460     Index, balanced accuracy, Matthew's Correlation Coefficient, and area under the receiver operating

1461     characteristic curve (AUC) of the singlet assignment probability for each genetic demultiplexing

1462     tool.

1463

1464     **Figure 5. Application of Ensemblex to highly multiplexed, experimentally pooled cultures of**

1465     **differentiated dopaminergic neurons.** **A)** Time line of iPSC pooling, dopaminergic neuron

1466     (DaN) differentiation, and sample collection from the DaN dataset by Jerber et al. (12). Three

1467     technical replicates at each time point (days 11, 30 and, 52 of differentiation) from pools containing

1468     22 individual iPSC lines were used in the analysis. Across all timepoints and technical replicates,

1469     84,746 cells were obtained for analysis. **B)** Uniform manifold approximation and projection

1470     (UMAP) plots showing confidently assigned singlets or predicted doublets (blue) and ambiguous

1471     singlets (singlet assignments with insufficient assignment probabilities; red) returned by each

1472     demultiplexing tool. **C)** Stacked bar chart showing the proportion of confidently assigned singlets

1473     or predicted doublets (blue) and ambiguous singlets (red) across technical replicates at each time

1474     point returned by each demultiplexing tool. **D)** Boxplot showing the proportion of confidently

1475     classified singlets across technical replicates and time points by each demultiplexing tool.

1476     Wilcoxon rank-sum tests were used to compare the proportion of confidently classified singlets by

1477     Ensemblex to that of its constituents (n = 9 pools). **E)** Bar chart showing the proportion of

1478     overlapping ambiguous singlet assignments amongst demultiplexing tools across technical

1479     replicates and time points (n = 9 pools). **F)** Boxplot showing the Adjusted Rand Index (ARI)

67

1480    assessing cluster stability across a range of 11 clustering resolutions (*n* clustering iterations = 25)

1481    after removing doublets identified by each demultiplexing tool. Wilcoxon rank-sum tests were

1482    used to compare the clustering ARI after removing Ensemblex doublets to the clustering ARI after

1483    removing doublets identified by each constituent tool. * Adjusted P-value < 0.05; ** adjusted P-

1484    value < 0.01; *** adjusted P-value < 0.001

1485

1486    **Figure 6. Evaluating the impact of discordant assignments between genetic demultiplexing**

1487    **tools on differential gene expression analysis. A)** Schematic illustrating the workflow for the

1488    neural stem cell (NSC) dataset. Pooled induced pluripotent stem cell (iPSC)-derived neural stem

1489    cell cultures from individuals with attention deficit hyperactivity disorder (ADHD) and controls

1490    were collected in two separate experiments. NSCs were dissociated for single-cell RNA

1491    sequencing and prior genotype information of the pooled subjects was obtained through

1492    microarray genotyping. The pools were demultiplexed by Ensemblex and its constituents with

1493    prior genotype information and differential gene expression (DEG) was computed between ADHD

1494    and controls. **B)** Uniform manifold approximation and projection (UMAP) plot showing the

1495    putative cell types. **C)** Summary of the number of usable cells — singlets above the recommended

1496    probability threshold of the respective demultiplexing tool — assigned to ADHD donors and

1497    controls and the number of identified doublets by each demultiplexing tool. **D)** Boxplot showing

1498    the Adjusted Rand Index (ARI) assessing cluster stability across a range of 11 clustering

1499    resolutions (*n* clustering iterations = 25) after removing doublets identified by each demultiplexing

1500    tool. A one-way Analysis of Variance (ANOVA) test comparing the ARI after removing doublets

1501    identified by each tool revealed a significant difference between tools (n = 11 clustering

1502    resolutions; P-value = 1.18e-3). **E)** Proportion of ADHD and control cells identified as putative

1503     doublets by Ensemblex that were assigned as singlets by the constituent demultiplexing tools. **F)**

1504     Heatmap showing the number of cell-type specific DEGs between ADHD and controls using the

1505     subject labels of each demultiplexing tool. **G)** Heatmap showing the number of cell-type specific

1506     DEGs between ADHD and controls using the subject labels of each demultiplexing tool and

1507     removing putative doublets identified by Ensemblex. Cell-types not shown in the heatmaps had no

1508     DEGs passing the adjusted P-value < 0.01 and |Log2FC >= 0.5| threshold across all tools.

1509

1510     **Tables**

1511     **Table 1. Summary of individual genetic demultiplexing tools.**

| Genetic demultiplexing tool | Prior genotype information for genetic demultiplexing | Included in the Ensemblex framework |
|---|---|---|
| Demuxalot (5) | Required | Yes |
| Demuxlet (6) | Required | Yes |
| Freemuxlet (6) | Not supported | Yes |
| ScSplit (7) | Optional | No |
| Souporcell (8) | Optional | Yes |
| Vireo (9) | Optional | Yes |

1512

1513     **Table 2. Application of Ensemblex to pooled cultures of dopaminergic neurons from 22**

1514     **healthy controls.**

| | ARI between Ensemblex and constituent tool assignments | | | Percent contribution to Ensemblex assignments | | | $n$ usable cells | $n$ doublets |
|---|---|---|---|---|---|---|---|---|
| | **Day 11** | **Day 30** | **Day 52** | **Day 11** | **Day 30** | **Day 52** | | |
| **Demuxalot** | 0.987 | 0.955 | 0.982 | 97.29% | 94.75% | 97.57% | 75,962 | 8,279 |
| **Demuxlet** | 0.928 | 0.062 | 0.884 | 95.91% | 29.74% | 90.55% | 57,567 | 6,614 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Souporcell** | 0.883 | 0.876 | 0.912 | 91.62% | 91.82% | 93.84% | 76,811 | 7,740 |
| **Vireo-GT** | 0.961 | 0.879 | 0.958 | 95.95% | 88.80% | 95.16% | 75,933 | 6,115 |
| **Ensemblex** | NA | NA | NA | NA | NA | NA | 76,222 | 8,307 |
| **DoubletFinder** | NA | NA | NA | NA | NA | NA | NA | 4,597 |

Pooled cultures of induced pluripotent stem cell (iPSC) lines from 22 healthy donors were differentiated towards a dopaminergic neuron (DaN) fate and sequenced on days 11, 30, and 52 of differentiation by Jerber et al. (12). For the analysis we used three technical replicates for each sequencing timepoint. Each pool was demultiplexed independently by Ensemblex and its constituent tools with prior genotype information. The Adjusted Rand Index (ARI) between Ensemblex's assignments and those of the constituent tools was computed across technical replicates corresponding to each differentiation timepoint. The percent contribution represents the proportion of assignments from each constituent tool that matched Ensemblex's assignments. Usable cells were defined as singlet classifications whose assignment probability exceeded the recommended threshold of the respective tool. Abbreviations: NA = Not applicable.

**References**

1.      Fiorini MR, Dilliott AA, Thomas RA, Farhan SMK. Transcriptomics of Human Brain Tissue in Parkinson's Disease: a Comparison of Bulk and Single-cell RNA Sequencing. Mol Neurobiol. 2024.

2.      Ringman JM, Goate A, Masters CL, Cairns NJ, Danek A, Graff-Radford N, et al. Genetic heterogeneity in Alzheimer disease and implications for treatment strategies. Curr Neurol Neurosci Rep. 2014;14(11):499.

3.      McKinney CE. Using induced pluripotent stem cells derived neurons to model brain diseases. Neural Regen Res. 2017;12(7):1062-7.

4.      Howitt G, Feng Y, Tobar L, Vassiliadis D, Hickey P, Dawson MA, et al. Benchmarking single-cell hashtag oligo demultiplexing methods. NAR Genomics and Bioinformatics. 2023;5(4):lqad086.

5.      Rogozhnikov A, Ramkumar P, Shah K, Bedi R, Kato S, Escola GS. Demuxalot: scaled up genetic demultiplexing for single-cell sequencing. bioRxiv. 2021:2021.05. 22.443646.

6.      Kang HM, Subramaniam M, Targ S, Nguyen M, Maliskova L, McCarthy E, et al. Multiplexed droplet single-cell RNA-sequencing using natural genetic variation. Nat Biotechnol. 2018;36(1):89-94.

1543    7.      Xu J, Falconer C, Nguyen Q, Crawford J, McKinnon BD, Mortlock S, et al. Genotype-free
1544    demultiplexing of pooled single-cell RNA-seq. Genome Biol. 2019;20(1):290.
1545    8.      Heaton H, Talman AM, Knights A, Imaz M, Gaffney DJ, Durbin R, et al. Souporcell: robust
1546    clustering of single-cell RNA-seq data by genotype without reference genotypes. Nat Methods.
1547    2020;17(6):615-20.
1548    9.      Huang Y, McCarthy DJ, Stegle O. Vireo: Bayesian demultiplexing of pooled single-cell
1549    RNA-seq data without genotype reference. Genome Biol. 2019;20(1):273.
1550    10.     Neavin D, Senabouth A, Hang Lee JT, Ripoll A, Consortium s-e, Franke L, et al. Demuxafy:
1551    Improvement in droplet assignment by integrating multiple single-cell demultiplexing and
1552    doublet detection methods. BioRxiv. 2022:2022.03. 07.483367.
1553    11.     Cardiello JF, Joven Araus A, Giatrellis S, Helsens C, Simon A, Leigh ND. Evaluation of
1554    genetic demultiplexing of single-cell sequencing data from model species. Life Sci Alliance.
1555    2023;6(8).
1556    12.     Jerber J, Seaton DD, Cuomo AS, Kumasaka N, Haldane J, Steer J, et al. Population-scale
1557    single-cell RNA-seq profiling across dopaminergic neuron differentiation. Nature genetics.
1558    2021;53(3):304-12.
1559    13.     Jindal A, Gupta P, Jayadeva, Sengupta D. Discovery of rare cells from voluminous single
1560    cell expression data. Nature communications. 2018;9(1):4719.
1561    14.     Bressan E, Reed X, Bansal V, Hutchins E, Cobb MM, Webb MG, et al. The Foundational
1562    Data Initiative for Parkinson Disease: Enabling efficient translation from genetic maps to
1563    mechanism. Cell Genom. 2023;3(3):100261.
1564    15.     Cardiello JF, Araus AJ, Giatrellis S, Helsens C, Simon A, Leigh ND. Evaluation of genetic
1565    demultiplexing of single-cell sequencing data from model species. Life Science Alliance.
1566    2023;6(8).
1567    16.     Large J, Lines J, Bagnall A. A probabilistic classifier ensemble weighting scheme based on
1568    cross-validated accuracy estimates. Data Min Knowl Discov. 2019;33(6):1674-709.
1569    17.     Chicco D, Tötsch N, Jurman G. The Matthews correlation coefficient (MCC) is more
1570    reliable than balanced accuracy, bookmaker informedness, and markedness in two-class
1571    confusion matrix evaluation. BioData mining. 2021;14:1-22.
1572    18.     20k Mixture of NSCLC DTCs from 7 donors, 3' v3.1 (with intronic reads) [Internet]. 10X
1573    Genomics. 2022 [cited January 8th, 2024]. Available from:
1574    https://www.10xgenomics.com/datasets/20k-mixture-of-nsclc-dtcs-from-7-donors-3-v3-1-with-
1575    intronic-reads-3-1-standard.
1576    19.     Stoeckius M, Zheng S, Houck-Loomis B, Hao S, Yeung BZ, Mauck WM, 3rd, et al. Cell
1577    Hashing with barcoded antibodies enables multiplexing and doublet detection for single cell
1578    genomics. Genome Biol. 2018;19(1):224.
1579    20.     McGinnis CS, Murrow LM, Gartner ZJ. DoubletFinder: Doublet Detection in Single-Cell
1580    RNA Sequencing Data Using Artificial Nearest Neighbors. Cell Syst. 2019;8(4):329-37 e4.
1581    21.     Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, et al. Highly Parallel
1582    Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. Cell.
1583    2015;161(5):1202-14.
1584    22.     Finak G, McDavid A, Yajima M, Deng J, Gersuk V, Shalek AK, et al. MAST: a flexible
1585    statistical framework for assessing transcriptional changes and characterizing heterogeneity in
1586    single-cell RNA sequencing data. Genome Biol. 2015;16:278.

1587    23.    Weber LM, Hippen AA, Hickey PF, Berrett KC, Gertz J, Doherty JA, et al. Genetic
1588    demultiplexing of pooled single-cell RNA-sequencing samples in cancer facilitates effective
1589    experimental design. Gigascience. 2021;10(9).
1590    24.    Bose A, Beal MF. Mitochondrial dysfunction in Parkinson's disease. Journal of
1591    neurochemistry. 2016;139:216-31.
1592    25.    Wickham H. ggplot2. Wiley interdisciplinary reviews: computational statistics.
1593    2011;3(2):180-5.
1594    26.    Khan MRAA. Rocit-an r package for performance assessment of binary classifier with
1595    visualization. 2019.
1596    27.    Genomes Project C, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, et al. A
1597    global reference for human genetic variation. Nature. 2015;526(7571):68-74.
1598    28.    Parkinson Progression Marker I. The Parkinson Progression Marker Initiative (PPMI). Prog
1599    Neurobiol. 2011;95(4):629-35.
1600    29.    Hao Y, Stuart T, Kowalski MH, Choudhary S, Hoffman P, Hartman A, et al. Dictionary
1601    learning for integrative, multimodal and scalable single-cell analysis. Nature biotechnology.
1602    2023:1-12.
1603    30.    Thomas RA, Fiorini MR, Amiri S, Fon EA, Farhan SM. ScRNAbox: Empowering Single-Cell
1604    RNA Sequencing on High Performance Computing Systems. bioRxiv. 2023:2023.11. 13.566851.
1605    31.    Li H. A statistical framework for SNP calling, mutation discovery, association mapping
1606    and population genetical parameter estimation from sequencing data. Bioinformatics.
1607    2011;27(21):2987-93.
1608    32.    Kilpinen H, Goncalves A, Leha A, Afzal V, Alasoo K, Ashford S, et al. Common genetic
1609    variation drives molecular heterogeneity in human iPSCs. Nature. 2017;546(7658):370-5.
1610    33.    Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, 3rd, et al.
1611    Comprehensive Integration of Single-Cell Data. Cell. 2019;177(7):1888-902 e21.
1612    34.    Streeter I, Harrison PW, Faulconbridge A, The HipSci C, Flicek P, Parkinson H, et al. The
1613    human-induced pluripotent stem cell initiative-data resources for cellular genetics. Nucleic
1614    Acids Res. 2017;45(D1):D691-D7.
1615    35.    Yde Ohki CM, Grossmann L, Doring C, Hoffmann P, Herms S, Werling AM, et al.
1616    Generation of integration-free induced pluripotent stem cells from healthy individuals. Stem
1617    Cell Res. 2021;53:102269.
1618    36.    Yde Ohki CM, Walter NM, Bender A, Rickli M, Ruhstaller S, Walitza S, et al. Growth rates
1619    of human induced pluripotent stem cells and neural stem cells from attention-deficit
1620    hyperactivity disorder patients: a preliminary study. J Neural Transm (Vienna). 2023;130(3):243-
1621    52.
1622    37.    Chen EY, Tan CM, Kou Y, Duan Q, Wang Z, Meirelles GV, et al. Enrichr: interactive and
1623    collaborative HTML5 gene list enrichment analysis tool. BMC Bioinformatics. 2013;14:128.
1624    38.    Zhang X, Lan Y, Xu J, Quan F, Zhao E, Deng C, et al. CellMarker: a manually curated
1625    resource of cell markers in human and mouse. Nucleic Acids Res. 2019;47(D1):D721-D8.
1626    39.    Hao Y, Hao S, Andersen-Nissen E, Mauck WM, 3rd, Zheng S, Butler A, et al. Integrated
1627    analysis of multimodal single-cell data. Cell. 2021;184(13):3573-87 e29.
1628    40.    Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, et al. Highly parallel
1629    genome-wide expression profiling of individual cells using nanoliter droplets. Cell.
1630    2015;161(5):1202-14.

1631    41.     Ihaka R, Gentleman R. R: a language for data analysis and graphics. Journal of
1632    computational and graphical statistics. 1996;5(3):299-314.
1633    42.     Wickham H, Wickham H. Data analysis: Springer; 2016.
1634    43.     Azzalini A, Menardi G. Clustering via nonparametric density estimation: The R package
1635    pdfCluster. arXiv preprint arXiv:13016559. 2013.
1636