# Iterative transcription factor screening enables rapid generation of microglia-like cells from human iPSC

Songlei Liu[1,2,*], Li Li[1,2,*], Fan Zhang[3,4,5,6,7,*], Björn van Sambeek[1,8], Evan Appleton[1,2], Alex H. M. Ng[1,2,9], Parastoo Khoshakhlagh[1,2,9], Yuting Chen[1], Mariana Garcia-Corral[1,2], Chun-Ting Wu[1,2], Jeremy Y. Huang[1,2], Yuqi Tan[10,11], George Chao[1,2], John Aach[1], Jenny Tam[1,2], Elaine T. Lim[12,13], Soumya Raychaudhuri[3,4,5,6,7,14,#], George M. Church[1,2,#]

[1] Department of Genetics, Blavatnik Institute, Harvard Medical School, Boston, MA, USA.
[2] Wyss Institute for Biologically Inspired Engineering, Harvard University, Boston, MA, USA.
[3] Center for Data Sciences, Brigham and Women's Hospital, Boston, MA, USA.
[4] Division of Genetics, Department of Medicine, Brigham and Women's Hospital, Boston, MA, USA.
[5] Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA.
[6] Broad Institute of MIT and Harvard, Cambridge, MA, USA.
[7] Division of Rheumatology, Inflammation, and Immunity, Brigham and Women's Hospital and Harvard Medical School, Boston, MA, USA.
[8] Radboud University, Nijmegen, the Netherlands.
[9] GC Therapeutics, Inc, Cambridge, MA, USA.
[10] Institute for Cell Engineering, Johns Hopkins University School of Medicine, Baltimore, MD, USA
[11] Department of Molecular Biology and Genetics, Johns Hopkins University School of Medicine, Baltimore, MD, USA
[12] Program in Bioinformatics and Integrative Biology and Departments of Neurology and Molecular, Cell and Cancer Biology, University of Massachusetts Chan Medical School, Worcester, MA, USA.
[13] NeuroNexus Institute, University of Massachusetts Chan Medical School, Worcester, MA, USA.
[14] Centre for Genetics and Genomics Versus Arthritis, Centre for Musculoskeletal Research, The University of Manchester, Manchester, UK.

[*] These authors contributed equally: Songlei Liu, Li Li, Fan Zhang.
[#] Correspondence: soumya@broadinstitute.org; gchurch@genetics.med.harvard.edu.

33 **Abstract**

34       The ability to differentiate stem cells into human cell types is essential to define basic

35   mechanisms and therapeutics, especially for cell types not routinely accessible by biopsies. But

36   while engineered expression of transcription factors (TFs) identified through TF screens has been

37   found to rapidly and efficiently produce some cell types, generation of other cell types that require

38   complex combinations of TFs has been elusive. Here we develop an iterative, pooled single-cell

39   TF screening method that improves the identification of effective TF combinations using the

40   generation of human microglia-like cells as a testbed: Two iterations identified a combination of

41   SPI1, CEBPA, FLI1, MEF2C, CEBPB, and IRF8 as sufficient to differentiate human iPSC into

42   microglia-like cells in 4 days. Characterization of TF-induced microglia demonstrated molecular

43   and functional similarity to primary microglia. We explore the use of single-cell atlas reference

44   datasets to confirm identified TFs and how combining single-cell TF perturbation and gene

45   expression data can enable the construction of causal gene regulatory networks. We describe

46   what will be needed to fashion these methods into a generalized integrated pipeline, further ideas

47   for enhancement, and possible applications.

48 **Introduction**

49       Recent advances and applications of single-cell assays, exemplified by collaborative

50   efforts such as the Human Cell Atlas (HCA)[1], have begun to provide a comprehensive view of cell

51   types and cellular states within the human body. Such maps are crucial for understanding human

52   development and diseases. From a synthetic biology perspective, these maps can be mined for

53   promising targets for cell fate engineering, with significant implications for disease modeling, cell

54   therapy, and regenerative medicine. Previously, we reported the construction of a comprehensive

55   human transcription factor library (TFome)[2]. In that study we used an unbiased approach for

56   screening differentiation and identified 290 transcription factors (TFs) that induced differentiation

57   of human induced pluripotent stem cells (hiPSCs) into various cell types. While the unbiased

58   screening method led to many interesting discoveries, it does not guarantee the generation of any

59   particular cell type. For those wishing to differentiate stem cells into a specific cell type of interest

60   for studying diseases and creating therapeutics, the availability of experimental and computational

61   pipelines for the identification of TFs to produce target cell types would be of great benefit. In this

62   study, we picked a target cell type for which TF-based differentiation method has not yet been

63   found, the microglia, for a proof-of-principal for developing new screening methodologies.

64   Microglia are the resident immune cells of the brain, which originated from erythro-myeloid

65   progenitors (EMPs) in the yolk sac[3,4]. They play important and diverse roles in brain development

66   and maintaining homeostasis[5–9]. Recent studies have demonstrated the link between

67   neuroinflammation and neurodegenerative disease, such as Alzheimer's Disease (AD)[10,11], and

68   along these lines, microglia have been shown to be an important cell type in AD and other

69   neurodegenerative diseases[12–15]. However, functional studies to define therapeutics targeting

70   human microglia have been greatly hindered by the limited availability of human brain biopsies[16,17].

71   The supply issue cannot simply be mitigated by using murine models, because differences

72   between human and mouse microglia limit the transferability of knowledge[18,19]. Producing human

73   microglia-like cells from hiPSCs might fill this gap. Several studies have accomplished this goal

74   through a process of embryoid body formation, growth factor treatment, and, in some cases, co-

75   culturing with neurons[20–28]. These protocols draw inspiration from the natural developmental

76   stages of microglia and have timelines ranging from 30-74 days. As the effects of extrinsic factors

77   on cell fate are frequently mediated by TFs, and building on top of our group's and others' prior

78   research in using TFs to accelerate differentiation[2,29–31], we hypothesized that direct

79   manipulations of TF expression could differentiate hiPSCs to microglia in a shorter timeframe.

80   In this study, we conducted two sequential iterations of pooled TF screening. Each round

81   of screening involves creating a barcoded TF library, pooled transfection into iPSCs for inducing

82    differentiation, and subsequent single-cell transcriptome analysis. From the analysis, TFs are

83    ranked by their ability to induce microglial gene expression, and the top hits then characterized

84    for their ability to induce differentiation into microglia (**Figure 1a, 2a**). We identified a TF

85    combination that produced cells transcriptionally resembled microglial cells in four days without

86    the need for media exchange. We demonstrated that these TF-induced microglia-like cells

87    (TFiMGLs) shared molecular and functional features of human primary microglia. Our barcoding

88    and amplification strategy allow for simultaneous detection of cell and TF barcodes from the high

89    throughput single-cell experiments, thus empowering us to analyze the regulatory relationships

90    between TFs and other genes. We also constructed a human single-cell transcriptome reference

91    by integrating publicly available scRNA-seq datasets of 225 samples representing 59 tissue types,

92    to which TF differentiated single cells could be mapped. We expect the methodology described

93    in this study to be broadly adopted to cell fate engineering and enable researchers to more

94    effectively select TFs to generate novel iPSC-derived cell types.

95    **Results**

96    **First round of pooled screening identified initial TFs for inducing microglia gene**

97    **expression**

98        To identify TFs that differentiate hiPSCs to microglia, our strategy is to first transfect and

99    stochastically integrate a pool of TFs into the cells, followed by differentiation induction and single-

100    cell RNA sequencing (scRNA-seq). From the scRNA-seq data, we use gene-expression profiles

101    to determine which cells are differentiating into microglia, and identify the TFs that were

102    transfected into these cells. To begin with, we surveyed previous literature on microglial

103    development[3,6,8,32,33], epigenetic and transcriptomic patterns[34–38], and gene regulatory networks[39].

104    We shortlisted 40 TFs for the first pooled TF screening (**Supplementary Table S1**). We cloned

105    each TF into the pBAN[2] vector for genomic integration with PiggyBac transposase and

106    doxycycline (Dox)-inducible expression. To distinguish between exogenous and endogenous TF

107 transcripts, we added a 20-nucleotide (nt) barcode between the stop codon and the poly-A

108 sequence of each TF (**Supplementary Figure S1**). We transfected the 40 TF vectors into

109 600,000 hiPSCs from a healthy donor (PGP1) with mass ratio between TF and transposase DNA

110 being 4:1, with the assumption that each cell can uptake and integrate multiple TFs into their

111 genome. After puromycin selection for TF-integrated cells, we induced differentiation by Dox for

112 four days (**Figure 1a**). With flow cytometry we observed that 0.3-0.5% of the cells expressed

113 microglial surface proteins, including CX3CR1, P2RY12, and CD11b (**Figure 1b**).

114  After four days of differentiation, we observed that 30% of the cells lost expression of a

115 stem cell marker, TRA-1-60 (**Figure 1c**). To pinpoint which of the 40 TF(s) were inducing

116 microglial gene expression, we sorted all differentiated (TRA-1-60 negative) cells for scRNA-seq

117 (**Figure 1a**). We performed two independent transfections (**Supplementary Figure S2**), and

118 spiked in 10% non-induced hiPSCs into each replicate as undifferentiated control during scRNA-

119 seq. After observing reproducible differentiation between the two replicates (**Figure 1d**), we

120 pooled the data together for downstream analysis. We observed expression of microglia genes

121 (*ITGAM*, *P2RY12*, *CX3CR1*, *TMEM119, TREM2*), as well as a cluster of cells with high expression

122 of *POU5F1*, marking stem cells (**Figure 1d, e**). Through amplicon sequencing of co-amplified TF

123 and cell barcodes from cDNAs (**Figure 1f**), we quantified expression of exogeneous TF(s) in

124 single cells in parallel (**Supplementary Figure S3**). In the Dox-induced cells, an average of 6.9

125 TFs (median 6, first quantile 4, third quantile 9) were expressed per cell. And 877 (8.5%) out of

126 10285 single cells had no TF expression, consistent with the 10% stem cell spike-in (**Figure 1g**).

127  We compared TF expression levels in cells with or without microglial RNA expression to

128 identify TF(s) correlated with a higher expression of microglial genes. Presumably these TFs were

129 the key TFs that had the potential to drive microglial differentiation. We identified three TFs likely

130 to cause microglial gene expression: *SPI1*, *FLI1*, and *CEBPA* (**Figure 1h**). *SPI1*, which encodes

131 PU.1 protein, is a known TF required for microglia development[32,33]. *CEBPA* is a known critical

132    regulator for myeloid differentiation[40]. *FLI1*, while hasn't yet been reported for microglial

133    differentiation, has been reported to interact with *RUNX1*[41] and *SPI1*[42], where both TFs are

134    indispensable for tissue-resident macrophage development[43].

135       We wanted to understand if these TFs could lead to microglial differentiation individually,

136    or if they needed to be used combinatorically. Individually expression of *CEBPA* and *FLI1* in

137    hiPSCs led to almost complete cell death, indicating the need for additional TFs to stabilize the

138    induced gene expression network. Expression of SPI1 alone was also ineffective, leading to only

139    induction of CD11b in 3% cells (**Supplementary Figure S4**), indicating that multiple TFs are

140    needed for the microglia differentiation.

141       We then tested combinations of TFs. Pooled transfection of *CEBPA+FLI1* ("C+F pool") or

142    *CEBPA+SPI1* ("C+S pool") led to improved microglial marker expression, while

143    *CEBPA+FLI1+SPI1* ("C+F+S pool") produced the most positive cells, reaching 14% CD11b+, 54%

144    P2RY12+ after four days (**Figure 1i**). However, we observed no expression of CX3CR1, a

145    chemokine receptor important for microglia activation and migration[44,45]. Because pooled

146    transfection and PiggyBac integration of three plasmids does not guarantee that every cell

147    expressed all three TFs, we built polycistronic expression cassettes by linking the TFs with 2A

148    peptides. A previous study reported that the gene position in the cassette affects their relative

149    expression level, with the first gene being the highest-expressed[46]. Therefore, we arranged the

150    TFs in different orders (**Supplementary Figure S5**). We named the construct by ordering letters

151    corresponding to each TF corresponding to their order on the plasmid. For example, **S***PI1*-T2A-

152    **F***LI1*-P2A-**C***EBPA* was named "MG3.1-SFC". Transfection and induction of the MG3.1-CFS and

153    FCS cassettes both led to dramatic cell death by day 4, consistent with the previous observation

154    that sole CEBPA and FLI1 expression caused cell death. MG3.1-SFC, which positioned *SPI1* at

155    the front, produced cells expressing two microglial genes, CD11b+ and P2RY12+ cells (37% and

156    6% of cells, respectively), but no *CX3CR1* expression (**Figure 1i**). The differences between cells

157    derived using the MG3.1-SFC cassette and the microglia-like cells from the C+F+S pool are

158    potentially due to different dosages of the TFs. While individual cells within the C+F+S pool may

159    have expressed variable dosage combination of the three TFs, MG3.1-SFC likely induced a fixed

160    dosage ratio for all cells. Critically, the lack of CX3CR1 expression and the low percentage of

161    CD11b- and P2RY12-positive cells from all 3-TF conditions indicated that additional TFs were

162    needed for efficient microglia differentiation from hiPSCs.

163

**164    Second iteration of pooled TF screen using MG3.1-SFC as baseline identified additional**

**165    TFs for improved microglia differentiation**

166            Recognizing that the three TFs identified in the first iteration were in themselves

167    inadequate to differentiate microglia, we pursued a second iteration of our screen. To build upon

168    the hits from the first pooled screen and identify additional TFs essential for microglia

169    differentiation, we performed a second pooled screen. For this iteration, we used the top three

170    TFs from the first iteration as baseline and tested the addition of other TFs (3+X) (**Figure 2a**). To

171    determine what TFs should be included in the second pool, we performed a bulk RNA-seq

172    analysis for MG3.1-SFC and compared it with published data on human primary microglia

173    (GSE89189, GSE99074)[21,35]. Based on differential gene expression analysis using DESeq2[47], we

174    first picked 25 TFs included in the first pool that still had lower expression levels in MG3.1-SFC

175    than primary microglia. We then included six new TFs that are significantly higher in primary

176    microglia. By looking at which TFs regulate the genes that have lower expression in MG3.1-SFC

177    using Molecular Signatures Database (MSigDB)[48] regulatory target gene sets, we included *IRF2*

178    and *ELF1*. Additionally, using CellNet[49], a computational tool that can classify bulk transcriptomic

179    data and predict missing gene regulators, we added six more TFs to the pool. Lastly, referring to

180    a recent single-cell study on fetal microglia development[50], we added *SPIB*, *ETS1* and *ELK3*.

181    Thus, the second TF pool contained a total of 42 TFs (**Supplementary Table S2**).

182    To ensure each cell expressed both the *SPI1*-T2A-*FLI1*-P2A-*CEBPA* cassette and

183    additional TFs from the second pool, we cloned the SFC cassette into a bleomycin-resistant

184    vector and the new TF pool into a puromycin-resistant vector (**Figure 2b**). We transfected 600,000

185    PGP1 hiPSCs in duplicates and performed dual-drug selection for genomic integration. After four

186    days of TF expression, we performed the same process of scRNA-seq and TF barcode amplicon

187    sequencing as in the first iteration (**Supplementary Figure S6**). As controls, we spiked in 5%

188    undifferentiated hiPSCs and 10% MG3.1-SFC during single-cell encapsulation to mark the

189    differentiation starting point of two iterations. We applied this approach to two pools of cells.

190    When we analyzed TF barcodes in this experiment, we observed two clusters of cells on

191    UMAP that corresponded to hiPSCs and MG3.1-SFC, while also showed new clusters of cells

192    that express additional TFs (**Figure 2c, d**). When we counted TF barcodes, we observed that out

193    of the total 8051 single cells from two independent transfections, 284 (3.5%) cells had no TF

194    barcode and 613 (7.6%) cells had only the barcode for MG3.1-SFC. On average each cell

195    expressed five TFs (median 4, first quantile 3, third quantile 7) (**Figure 2e, f, Supplementary**

196    **Figure S7**), with most cells (88.9%) expressing the SFC cassette plus at least one other TF.

197    To determine which of the new TFs lead to improved microglia differentiation, we analyzed

198    their effects on microglial gene expression. We were especially interested in increasing

199    expression of *CX3CR1*, which was not expressed in MG3.1-SFC. We observed significantly

200    higher (p < 0.01) number of *MEF2C* and *KLF6* barcode in cells expressing *CX3CR1* (**Figure 2g**),

201    suggesting their ability to induce *CX3CR1* expression. It's worth noting that *MEF2C* was also

202    present in the first screening but failed to reach significance for upregulating *CX3CR1*, indicating

203    the use of the SFC cassette as baseline enabled other influential TFs to be found. *MEF2C* also

204    reached high ranking for *TMEM119* (**Figure 2g**). The SFC cassette ranked at the top for inducing

205    *ITGAM* and *P2RY12* expression, which is expected from the results of the first iteration. In this

206    round of screening, *CEBPB* and *IRF8* emerged as high-potential TFs for promoting *ITGAM* or

207    *P2RY12* expression (**Figure 2g**). From this second pooled TF screening, additional TFs of interest

208    found were *MEF2C*, *CEBPB*, *IRF8*, *KLF6*, and *BHLHE41*.

209        To validate that these additional TFs can promote microglial gene expression, we

210    individually expressed them in addition to SFC (SFC+1). When compared with MG3.1-SFC,

211    SFC+CEBPB increased the percentage of CD11b+ cells from 37% to 98% (**Figure 2h**) but led to

212    more cell death at day 4. SFC+MEF2C and SFC+IRF8 increased P2RY12 expression from 6%

213    to 45% (**Figure 2h**). Most importantly, SFC+MEF2C and SFC+KLF6 increased CX3CR1+ cells

214    from 0% to 20% and 2% respectively (**Figure 2h**). These results agreed well with the predictions

215    from single-cell TF barcode analysis, indicating the validity of using pooled TF screening for

216    inferring causality between TF and target gene expression.

217        To see if microglia differentiation can be further promoted by delivering more TFs to each

218    cell, we chose the three TFs from the SFC+1 experiment that led to highest increase in

219    percentage of microglial gene-expressing cells, CEBPB, IRF8, and MEF2C, to add to the SFC

220    set. We combined MEF2C, CEBPB and IRF8 into polycistronic cassettes. Because MEF2C

221    demonstrated ability to induce both CX3CR1 and P2RY12, we put it in the first place and varied

222    the position of CEBPB and IRF8, producing two cassettes: MIC and MCI (**Figure 2i**). We also

223    varied the position of FLI1 and CEBPA in the first construct to produce SFC and SCF, keeping

224    SPI1 in the front to avoid excessive cell death during differentiation. We tested all four

225    combinations of the two 3-TF cassettes (SFC-MIC, SFC-MCI, SCF-MIC, SCF-MCI) for their ability

226    to induce microglia differentiation (**Figure 2i**). Encouragingly, all 6-TF cocktails produced cell

227    pools with increased expression of microglial proteins when compared with MG3.1-SFC (**Figure**

228    **2j**). We observed that the most effective combination was MG6.4-SCF-MCI, resulting in 66%

229    CD11b+, 93% P2RY12+ and 16% CX3CR1+ cells at day 4, compared with 37%, 6%, and 0%

230    respectively for MG3.1-SFC, the baseline of the second iteration. These results highlighted the

231    value of the second iteration and demonstrated the utility of iterative TF screening for cell fate

232    engineering. We term cells differentiated using MG6.4-SCF-MCI Transcription Factor-induced

233    MicroGlial-Like cells, or TFiMGLs.

234

235    **TFiMGLs are phagocytic, responsive to disease-relevant stimulation, and share molecular**

236    **signatures with primary microglia**

237         To determine the differentiation dynamics of TFiMGLs, we performed bulk RNA-seq

238    analysis of the cells on 0, 1, 2, 3, 4, 6 days post induction of the six TFs. We observed a rapid

239    induction of all six TFs on day 1, and they reached plateau on day 2 (**Figure 3a**). This

240    accompanied a quick downregulation of *POU5F1* on day 1 and followed by upregulation of

241    microglial genes from day 2 and onwards (**Figure 3b**). Principal component analysis (PCA)

242    reflected a similar trend, where a rapid differentiation occurred on day 1 and 2, followed by a

243    gradual deceleration from day 3 to day 6 (**Figure 3c**). In a plot of PC1 and PC2, the day 4 and

244    day 6 transcriptomes are close together, suggesting a stable window for functional investigation.

245    For downstream characterizations of TFiMGLs, we chose to differentiate cells for 4 days.

246         Brightfield microscopy analysis of TFiMGLs revealed rapid morphological change from

247    day 1 to day 6 (**Supplementary Figure S8**). Immunofluorescence analysis confirmed the loss of

248    pluripotency marker OCT4 and expression of key microglial proteins: CD11b, P2RY12, and

249    CX3CR1 (**Figure 3d**). TFiMGLs demonstrated reproducible differentiation between replicates,

250    with 53.9 ± 0.57% (SD, n=3) CD11b+, 93.1 ± 0.50% (SD, n=3) P2RY12+ and 14.8 ± 0.68% (SD,

251    n=3) CX3CR1+ cells (**Figure 3e**).

252         As brain resident macrophages, microglia play important roles in brain development and

253    homeostasis. Microglia's abilities to respond to signals related to degenerating neurons and

254    phagocytose are integral parts of their function[21]. To investigate if TFiMGLs could mimic the

255    phagocytosis function of microglia, we incubated TFiMGLs with pHrodo green labeled S. aureus

256    particles for 0.5, 2 and 4 hours, and performed flow cytometry and microscopy. While 0.5-hour

257    incubation showed minimal phagocytosis activity, nearly all cells were positive for pHrodo green

258    at 2 hours and the intensity grew even stronger at 4 hours (**Figure 3f, Supplementary Video S1-**

259    **2**). Microscopy analysis at 4 hours with co-staining of microglia surface proteins confirmed the

260    intracellular position of these particles (**Figure 3g**). ADP is one of the substances released from

261    injured neurons and works as an signal to stimulate microglial responses[51,52]. To study if TFiMGLs

262    are responsive to ADP stimulation, we incubated TFiMGLs with calcium indicator Fluo-4 and then

263    stimulated with ADP containing media. We imaged the cells at a three-second interval. We

264    observed a rapid increase in calcium signal when ADP was added (**Figure 3h, i, j;**

265    **Supplementary Video S3**), suggesting TFiMGLs are responsive to ADP stimulation.

266        To assess how accurately TFiMGLs recapitulate the transcriptome of human microglia,

267    we compared TFiMGLs bulk-RNA-seq data to previously published bulk RNA-seq data for human

268    primary microglia and iPSC-derived microglia (GSE89189, GSE99074)[21,34,35]. To minimize

269    potential batch effects that might hinder meaningful comparison between datasets, we aligned all

270    raw FASTQ files to the same reference genome and applied a negative binomial regression-

271    based batch effect correction method, ComBat-seq[53], before downstream analysis. From the PCA

272    analysis, we observed that the transcriptomes of TFiMGLs from days 2-6 more closely resembled

273    primary microglia of different sources than to iPSCs or hematopoietic progenitors (HPCs),

274    suggesting a successful microglial fate induction (**Figure 4a**). We also observed that TFiMGLs

275    were distinct from monocytes or dendritic cells, two related cell types from the myeloid lineage.

276    To investigate if TFiMGLs express genes that are enriched in primary microglia, we performed

277    Gene Set Enrichment Analysis (GSEA)[54] on TFiMGLs versus iPSCs using two microglial gene

278    collections from the MSigDB derived from human brain scRNA-seq (M40168; M39077). We

279    observed significant positive microglial gene enrichment scores using both collections (M40168:

280    score = 0.72, p-value = 9.01e-10, gene set size = 313; M39077: score = 0.66, p-value = 9.01e-

281    10, gene set size = 391), indicating TFiMGLs upregulated those microglia-enriched genes (**Figure**

282    **4b**). To further investigate if TFiMGLs achieve transcriptomic similarity to primary microglia, we

283    also used a previously published collection of 881 microglia-enriched genes[34] to cluster the

284    samples from **Figure 4a**. While the transcriptome of day-1 TFiMGLs clustered closer to iPSCs,

285    day-2 and later TFiMGLs clustered closer to primary microglia, with key microglial genes

286    increasingly upregulated by the day (**Figure 4c**). Similar to what we saw in **Figure 4a**, TFiMGLs

287    were distinct from monocytes or dendritic cells by measuring the 881 genes. These analyses

288    demonstrate that the TFiMGLs have lost iPSC-like identity and now closely resemble microglia.

289         Microglia are able to respond to signals indicating brain infection and inflammation, such

290    as IFNγ, beta amyloid, and TDP-43. IFNγ is a known activator of microglia secreted by T

291    lymphocyte[55]. Beta amyloid (Aβ) is a key molecule in AD pathology which has been shown to

292    elicit microglia response[56]. TDP-43, whose aggregation is considered a hallmark of ALS and is

293    present in the vast majority of amyotrophic lateral sclerosis (ALS) patient, had also been also

294    shown to activate microglia[57]. To investigate how TFiMGLs respond to IFNγ, fibrillar Aβ (fAβ) and

295    TDP-43, we treated TFiMGLs in triplicates with each of the three molecules for 24 hours and

296    harvested cells for RNA-seq. PCA analysis revealed transcriptomic changes in the IFNγ and TDP-

297    43 treated group, while the fAβ-treated group showed minimal differences (**Figure 4d,**

298    **Supplementary Figure S9**). We confirmed fAβ formation by conducting an in vitro amyloid

299    fibrillation experiment that showed the Aβ peptide could form fibrils after 1 hour of incubation

300    (**Supplementary Figure S10**). Pathway analysis of differentially expressed genes from the IFNγ

301    treated group included "response to virus" and "response to bacterium" (**Figure 4e,**

302    **Supplementary Figure S11**), corresponding to the role of IFNγ production as a response to

303    infection. Top upregulated genes by IFNγ included CXCL10, CXCL11, IRF1 and IL18BP

304    (**Supplementary Figure S11**), which aligns with the IFNγ response genes revealed by an

305    independent single-cell level human-derived macrophage stimulation study[58]. For the TDP-43

306    treated cells, top differentially regulated pathway included "myeloid leukocyte mediated immunity"

307    and "myeloid cell activation involved in immune response" (**Figure 4f, Supplementary Figure**

308    **S12**), demonstrating that TFiMGLs were activated by the TDP-43 treatment. Collectively, these

309    results suggest that TFiMGLs exhibited microglia-like responses to infection- and ALS-related

310    stimulations.

311

312    **Single-cell atlas reference mapping confirmed microglia-like fate induction**

313            There is the opportunity to map cell states precisely leveraging advances in single-cell

314    analysis technologies. Up until this point, we have been using a group of cellular markers to

315    determine cell type. While this is a common practice in both primary human tissue and stem cell

316    differentiation studies, it is now possible to define cell types based on more comprehensive

317    molecular profiles, including the whole transcriptome. We aspired to develop a strategy that might

318    be generally applicable to leveraging single-cell atlas data to guide cell fate engineering efforts.

319            To achieve this goal, there are several important prerequisites: 1) reference single-cell

320    data sets from all developing human tissue types, capturing different developmental stages; 2)

321    data integration methods to combine these different reference data sets to create a

322    comprehensive cell atlas; 3) reference mapping methods to project iPSC-derived cell data onto

323    the combined reference and quantitative assessment of similarity to differentiating cell classes; 4)

324    existence of perturbation libraries that are scRNA-seq compatible, which include but not limited

325    to open reading frame (ORF) and CRISPR libraries. While advances have been made and

326    continue for #2-4, incomplete reference data continues to be an acute problem.

327            To explore this idea of atlas-guided cell fate engineering, despite limitations in available

328    reference data, we re-examined the previously described two pooled TF screen for microglia

329    differentiation. We built a refence data set by compiling scRNA-seq data from published datasets

330     generated through 10x Chromium platform. In total, the final single-cell atlas contains 225

331     samples from 59 organ or tissue types, with a total of 1,004,650 single cells (**Supplementary**

332     **Table S3**). The majority of the data was obtained from PanglaoDB[59], where all raw reads from

333     different studies were aligned and processed together. We added other data sets representing

334     brain[60] and endometrium[61] which were under-represented in PanglaoDB. We carefully filtered all

335     raw data downloaded from PanglaoDB for cell, gene, UMI number and mitochondria gene ratio

336     (**Supplementary Table S4**). In its current form, the cells in the atlas are annotated according to

337     their organ or tissue of origin; cellular level annotation for all 59 studies have yet to be defined.

338         To reduce batch variability across different studies, the data were integrated with

339     Harmony[62] (**Figure 5a**; **Supplementary Figure S13**). Qualitative assessment of the UMAP plots

340     post-integration indicated co-clustering of cells in related tissues from different studies,

341     exemplified by the overlap between "Primary brain" with "Embryo forebrain" datasets, as well as

342     between the "Pancreatic islets" with "Pseudoislets" datasets (**Supplementary Figure S13**). To

343     visualize where microglia are located on this map, we acquired cell annotation information for the

344     "Primary brain" dataset[60]. We observed a cluster containing microglia on the right of the UMAP

345     plot (**Figure 5b**). To see if any of the TF differentiated cells can be mapped closely to microglia,

346     we projected the scRNA-seq data from the two pooled TF screens onto the integrated atlas using

347     Symphony[63] (**Figure 5c, d**). In the projections, we observed 4.3% and 26.5% cells from the first

348     and second screen being projected to the microglia-containing cluster. This increase in

349     percentage is like because the top hits from the first screen, SPI1, FLI1, and CEBPA were the

350     baseline for the second screen. Because we have identified SPI1, CEBPA, FLI1, MEF2C, CEBPB,

351     and IRF8 as the inducers for differentiating iPSCs to microglia-like cells, we wanted to see if their

352     barcode expression lead to co-localization with microglia on the reference atlas. We saw CEBPA

353     barcode had high expression in the microglia cluster, as well as a few others (**Figure 5e**),

354     indicating its ability to activate a broad range of genes. SPI1 barcode, on the other hand, had a

355    distinct enrichment in the microglia cluster (**Figure 5f**), suggesting a microglia specific gene

356    induction. We observed the co-expression of SPI1, FLI1, and CEBPA by the SFC cassette led to

357    a strong mapping to the microglia cluster (**Figure 5g**). We also visualized the expression of FLI1,

358    MEF2C, CEBPB, and IRF8 individually (**Supplementary Figure S14**). While they did not show

359    obvious enrichment in microglia cluster themselves, combined reads from the all six TFs still

360    showed a strong induction towards microglia (**Figure 5h**). This single-cell atlas reference mapping

361    analysis confirmed microglial fate induction by the six TFs with a comprehensive comparison with

362    1,004,650 single cells from 59 organ or tissue types.

363

364    **Regression analysis revealed causal TF-gene regulatory relationships**

365        Accumulating scRNA-seq data of human cell types provide ever-expanding information

366    about what genes define a cell fate. As a result, a complete knowledge map of what TFs turn on

367    what genes is instrumental for cell fate engineering. A tremendous amount of insight on this

368    subject  have been produced by computational methods for inferring TF-gene regulatory network

369    (GRN)[64,65] and databases based on TF-binding sites[66,67], TF-gene co-expression[68,69], and protein-

370    protein interaction[70,71]. However, the only way to acquire a definitive causal TF-gene regulatory

371    relationship map is by introducing the TF perturbations and observing their effects in a highly

372    multiplexed way. Our two pooled TF screens combined with scRNA-seq readout enabled us to

373    explore this exact idea. Taking advantage of our experimental design that captured both TF

374    barcode and cellular RNA expression, and by implementing a stepwise regression model (Online

375    Methods), we were able to construct GRNs from the two pooled screens (**Figure 6**). In our dataset,

376    each TF transgene was represented by two distinct values: counts from barcode amplicon

377    sequencing, and counts from their RNA molecules in scRNA-seq. Although the counts from

378    scRNA-seq might contain reads from endogenous TF expression, the two measurements

379    correlated well for most TFs (**Supplementary Figure S15**), the exception being a few TFs with

380 low expression. We reasoned that TFs that had higher correlation between their barcode and

381 RNA measurements demonstrated higher consistency between experiments, which were more

382 reliable to produce accurate regression results using the two matrices. Thus, by selecting TFs

383 had a correlation coefficient greater than 0.3 between the two measurements, we selected 18

384 TFs from the first iteration and 21 TFs from the second iteration for regression analysis. We

385 observed extensive gene expression changes caused by CEBPA and the triple TF cassette

386 MG3.1_SFC expression (**Figure 6a, f**), as well as slightly smaller networks from CIITA, SPI1,

387 ERG2, JUN, CEBPB, ZFP36, and BHLHE41 (**Figure 6b-e, g, h, Supplementary Figure S16**).

388 Among 672 edges, we observed 76% to be positive regulations. Some TFs (CIITA, SPI1, JUN)

389 only showed positive edges in current thresholding conditions (Abs(coefficient) > 0.1 & -log10(p-

390 value) > 20), indicating they were mostly activating other genes. Other TFs (CEBPB, ZFP36,

391 BHLHE41) showed negative edges, indicating their repressive roles. We also observed several

392 genes simultaneously connected with more than one TFs (**Supplementary Figure S17**). For

393 example, RAB13, a membrane trafficking regulator, was upregulated by both CEBPA and CEBPB;

394 HMGA1, a master regulator of chromatin structure, was downregulated by both BHLHE41 and

395 CEBPA; FLNC, an actin crosslinking protein, was upregulated by JUN while downregulated by

396 CEBPA. There are many more regulatory relationships we listed in detail from these two pooled

397 screens (**Supplementary Table S5-6**). With larger perturbation libraries, higher-throughput

398 scRNA-seq, and more scalable regression analysis methods, we believe a more complete

399 knowledge map of causal TF-gene regulatory relationships could be built and greatly facilitate cell

400 fate engineering efforts.

401 **Discussion**

402 Differentiating human cell types from stem cells provides is essential for basic research

403 and therapeutics development, especially when the desired cell types are not easily obtainable

404 from accessible human tissues. Advances in the understanding of developmental biology has

405    fueled the discovery and application of protocols to differentiate specific cell types from iPSCs.

406    Some of this work has been translated into treatment strategies that are now being investigated

407    with clinical trials for devastating diseases like age-related macular degeneration[72] and type 1

408    diabetes[73]. Differentiated iPSCs have now also become routinely used in laboratories for studying

409    disease mechanisms and testing drugs[74]. Recent global efforts on building single cell atlases of

410    cellular development have expanded the knowledge of human development and diseases, but

411    also present a key resource for cell fate engineering. Combined with technological advancements

412    in genetic library construction, high-throughput screening and sequencing technologies, to the

413    field is primed for investigating how to engineer cell fate in a systematic and multiplexed fashion,

414    as performed in this work.

415        Our study demonstrated the feasibility of combining an iterative genetic library screen with

416    high-throughput scRNA-seq for cell fate engineering. Using microglia, a cell type which previously

417    did not have a TF-driven differentiation protocol, as a model target, we performed two iterations

418    of our design-screen-validate workflow and identified SPI1, CEBPA, FLI1, MEF2C, CEBPB, and

419    IRF8 as a potent recipe for driving microglia differentiation from hiPSCs within 4 days, a dramatic

420    reduction from the standard 35 days through growth factor-based protocols[26]. Characterizations

421    of TFiMGLs indicated that they possessed transcriptomic and functional resemblance to primary

422    human microglia.  We also explored the possibility of using single-cell atlas for guiding cell fate

423    engineering by building a single-cell reference and mapping our pooled screen scRNA-seq data

424    to it. We observed an increased percentage of microglia mapping cells from the second iteration

425    when compared to the first. The high expression of key TFs like CEBPA and SPI1 in the microglia

426    containing cluster also confirmed their ability to drive microglia differentiation from iPSCs. By

427    doing genome wide regression analysis between TF barcode counts and gene expression levels,

428    we revealed TF-gene regulatory relationships present in these pooled screens.

429        During this study, we noted several technological challenges that could be addressed in

430    future studies. Most current TF-based differentiation protocols rely on a one-time induction of TF

431    expression, while lacking the capability for sequential induction. Despite this, current strategies

432    have successfully generated certain cell types, including multiple types of neurons[75], endothelial

433    cells[2], and the induction of iPSC itself[76]. However, during development *in vivo* coordinated gene

434    programs are sequentially activated, as observed in time-resolved transcriptomic analysis of

435    developing tissues[77]. This feature could be re-created by identifying orthogonal induction system

436    with a comparable strength to the doxycycline-inducible system or developing tunable gene

437    circuits. With these tools, it will be possible to test whether sequential TF expression can lead to

438    improved differentiation accuracy. In addition to temporal control of TF expression, the ability to

439    regulate expression levels of individual TFs could also lead to improvements in differentiation.

440    There were two manifestations of this pattern in the current study. In the first case, although

441    CEBPA and FLI1 expressed by themselves led to cell death, the presence in the SFC cassette

442    enabled cell survival and differentiation. The reduced expression levels of CEBPA and FLI1 likely

443    and potential interactions with other TFs could also explain why we were able to observe an

444    extensive GRN for CEBPA from the first pooled screen, which would not be possible due to toxicity

445    in CEBPA-expressing cells. The second case can be observed in the effect that different

446    sequential arrangements of TFs in the polycistronic cassettes led to different levels of downstream

447    microglial protein expression. The effect of TF stoichiometry on differentiation efficiency has also

448    been observed for cardiac myocyte programming[78]. While the positional effects in a polycistronic

449    cassette offer one way to explore the stoichiometry space, development of new titratable

450    promoters allowing turning of individual genes could be an important tool for cell fate engineering.

451        We demonstrated the potential for using primary single-cell atlas to guide iPSC

452    differentiation efforts. We note that these approaches will improve with better quality atlases. A

453    cell atlas with wide representation of human tissues, deep sequencing coverage, and reliable

454  cellular-level annotation is ideal for guiding cell fate engineering. A number of methods for scRNA-

455  seq datasets integration have been developed[79], with the goal of enabling comparison between

456  batches of data. However, because construction of a comprehensive cell atlas would need to

457  integrate data from tens to hundreds of separately acquired datasets, the accuracy and

458  computational efficiency of current integration strategies require improvement. We also note, that

459  current reference mapping methods were designed to project new data from the same tissue to

460  old datasets, or projecting data from the same tissue but acquired by different modalities. As a

461  result, these methods are not optimized for iPSC-derived cells with incomplete conversions,

462  leading to partially resemble multiple primary cells types. A reference mapping method that

463  provides probability-based quantitative measurement and rejection options is needed to address

464  this issue. Furthermore, we believe experimental strategies might be devised to improve mapping

465  of differentiated iPSCs to single cell reference data sets. For example, spiking in a standard set

466  of differentiated cell types could be used as landmarks to validate single cell profile mapping.

467  As demonstrated in this study, we used microglia as a model target to develop iterative

468  screening methods for identifying TFs that drive iPSC differentiation towards a specific cell type.

469  We found that the combination of SPI1, CEBPA, FLI1, MEF2C, CEBPB, and IRF8 could produce

470  microglia-like cells from iPSCs in as quickly as four days. We built computationally a

471  comprehensive single-cell reference atlas and used it to validate the results from our iterative

472  screening. We also used a stepwise regression model to discover TF-gene regulatory relationship

473  from the scRNA-seq data. We believe the TFiMGLs can be used as a model for human microglia

474  and facilitate basic and translational research that would benefit from a rapid turnaround time. We

475  also envisage that the methods of iterative pooled TF screen, single-cell reference atlas mapping,

476  and TF-gene regulation analysis will find their utility in other cell type targets for iPSC
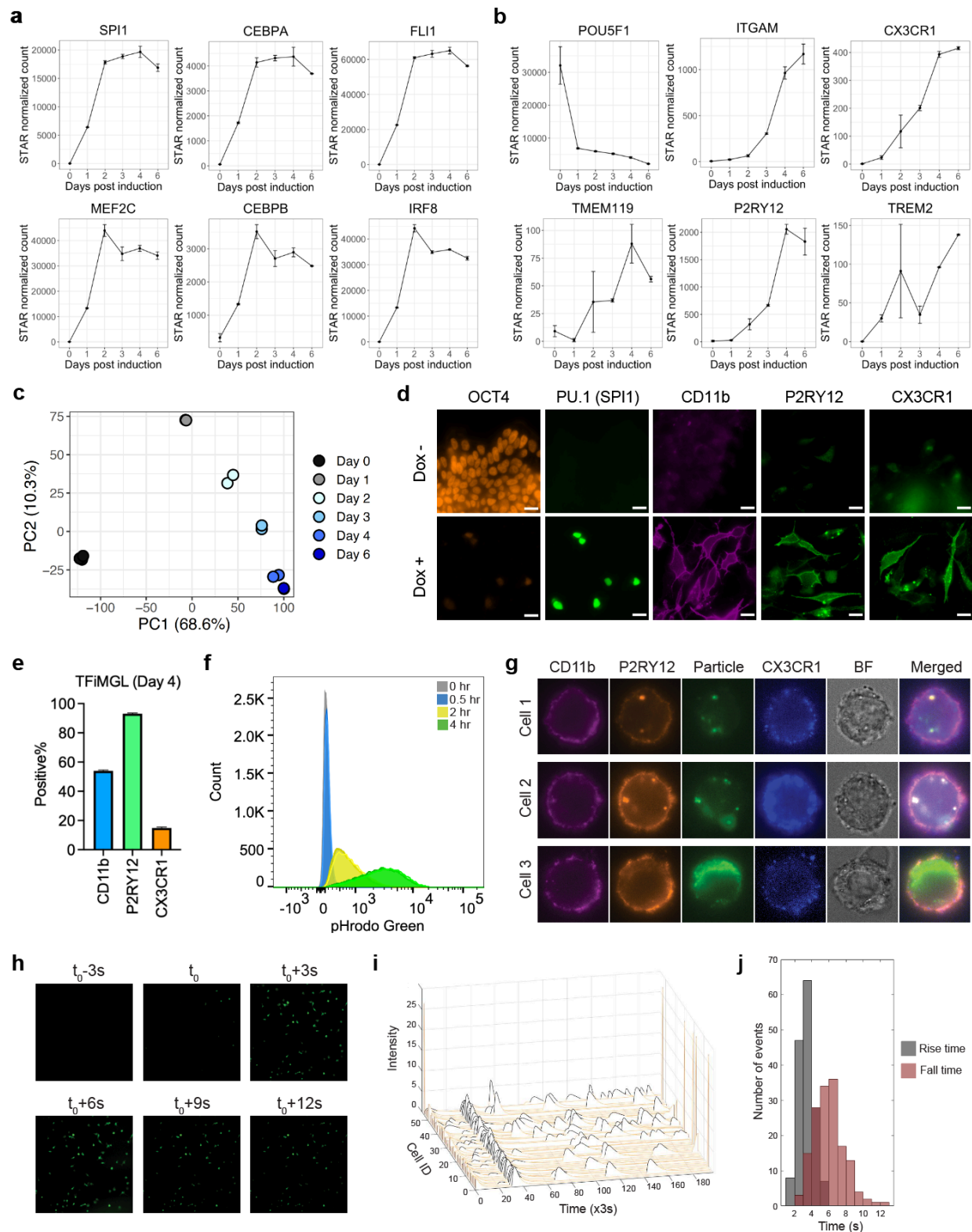
477  differentiation.

478

479

480

481 **Figure 1.** First round of pooled screening identified initial TFs for inducing microglia gene
482 expression. (**a**) Workflow of the first pooled TF screen. (**b**) Flow cytometry analysis of stem cell
483 (TRA-1-60) and microglia (P2RY12, CD11b, CX3CR1) proteins in the PGP1 + 40 TF pool before
484 and after Dox induction. (**c**) Cells with low TRA-1-60 expression in the Dox+ group were sorted
485 for scRNA-seq. (**d**) Clustering of two independently transfected and differentiated PGP1 iPSC
486 pools. Colors represent clusters identified by Seurat at 0.3 resolution. (**e**) Expression of microglia
487 (*ITGAM*, *CX3CR1*, *TMEM119*, *P2RY12*, *TREM2*) and spiked-in stem cell (*POU5F1*) gene in
488 scRNA-seq. (**f**) Primer designs for co-amplification of TF and cell barcodes in 10x Genomics 3'
489 workflow. (**g**) Number of TFs per cell counted from normalized and binarized TF expression matrix.
490 (**h**) Ranking of the 40 TFs after Wilcoxon rank sum test with the two tested groups being with or
491 without microglia gene expression. Blue highlights top-ranking TFs. (**i**) Flow cytometry validation
492 of top-ranking TFs for inducing microglia protein expression. C = CEBPA, F = FLI1, S = SPI1.
493 "Pool" means pooled transfection, no polycistronic cassette used.

494

495

496

497 **Figure 2.** Second iteration of pooled TF screen using MG3.1-SFC as baseline identified additional
498 TFs for improved microglia differentiation. (**a**) Workflow of the second pooled TF screen. (**b**)
499 Polycistronic cassette design for performing dual-drug selection to achieve 3+X TF screen. (**c**)
500 Normalized mRNA expression from the polycistronic cassette (*SPI1*, *FLI1*, *CEBPA*) and stem
501 cells (*POU5F1*). (**d**) TF barcode counting enabled identification of stem cells ("No TF BC"),
502 MG3.1-SFC and cells with additional TFs ("SFC+X"). (**e**) Example histograms of TF barcode raw
503 counts in single cells. (**f**) Number of TFs per cell counted from normalized and binarized TF
504 expression matrix. (**g**) Ranking of the 42 TFs after Wilcoxon rank sum test with the two tested
505 groups being with or without microglia gene expression. Blue highlights top-ranking TFs. Grey
506 highlights the SFC polycistronic cassette. (**h**) Flow cytometry validation of top-ranking TFs for
507 improving microglia protein expression. (**i**) Polycistronic cassettes design for varying TF orders.
508 (**j**) Flow cytometry analysis of different arrangements of the six-TF recipe in comparison with
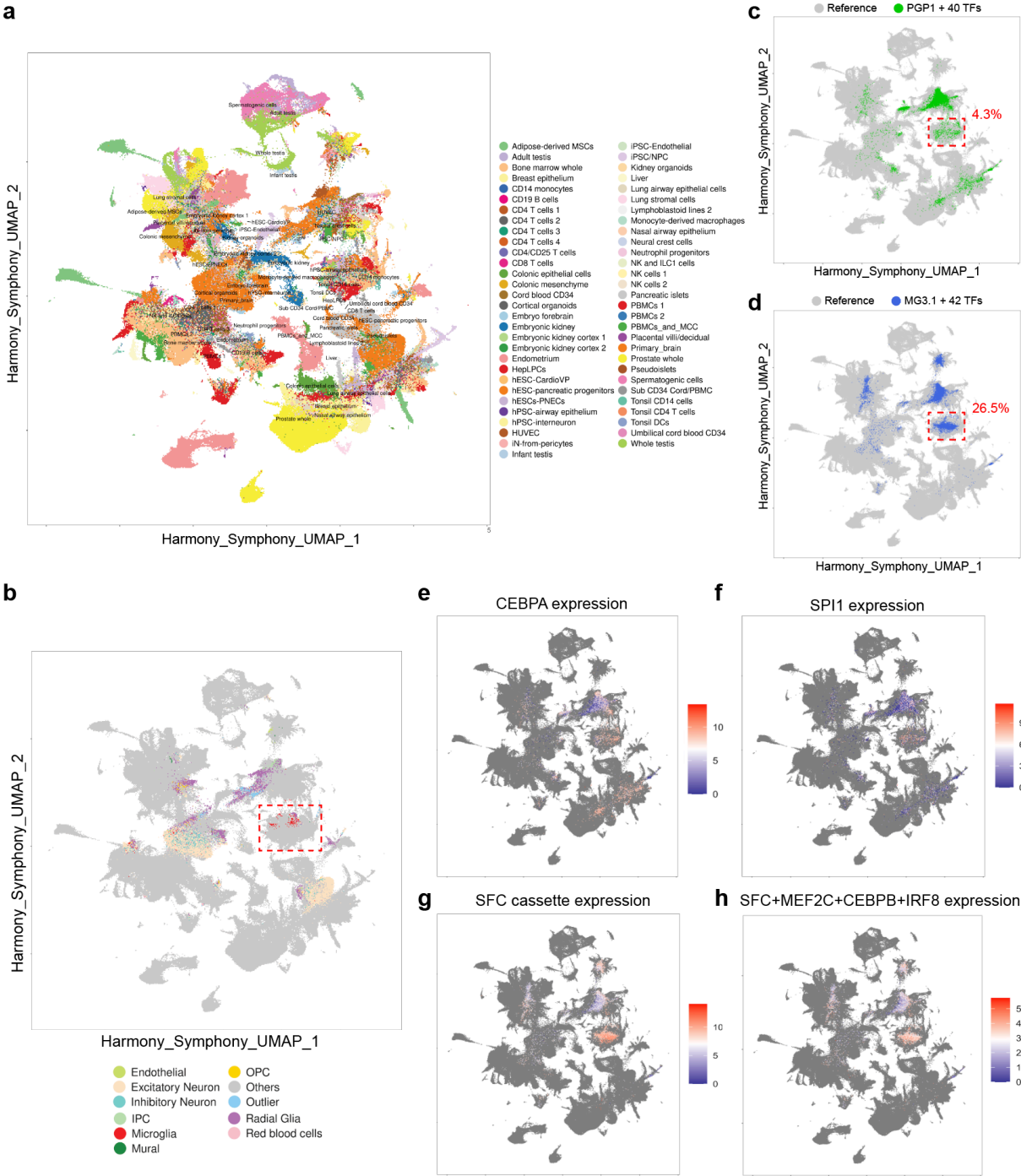509 MG3.1-SFC.

510

511

512

513 **Figure 3.** TFiMGLs differentiate quickly, are phagocytic and responsive to ADP stimulation. (**a**)
514 Expression of the six induced TFs over time measured by bulk RNA-seq. n = 2 for each day. Error
515 bar represents standard deviation. (**b**) Expression of stem cell (*POU5F1*) and microglia (*ITGAM*,
516 *CX3CR1*, *TMEM119*, *P2RY12*, *TREM2*) genes over time measured by bulk RNA-seq. (**c**) PCA
517 plot for the transcriptome of TFiMGLs (MG6.4) over time. (**d**) Immunofluorescence of stem cell
518 (OCT4), Dox-induced (PU.1), and microglia (CD11b, P2RY12, CX3CR1) proteins on day 4. Scale
519 bar: 20 μm. (**e**) Flow cytometry quantification of microglia protein expression on day 4 (n=3). (**f**)
520 Flow cytometry analysis of the uptake of pHrodo-labeled S. aureus Bioparticles over time (n=3).
521 (**g**) Microscopy analysis of particle uptake combined with microglia surface protein staining. (**h**)
522 Calcium imaging with Fluo-4 after stimulation with 150 μM ADP and peak quantification. Images
523 taken once every three seconds. ADP was added at $t_0$. (**i**) Quantification of fluorescent signals
524 from all cells in the field of view in panel h over a period of 10 minutes. (**h**) Peak dynamics analysis
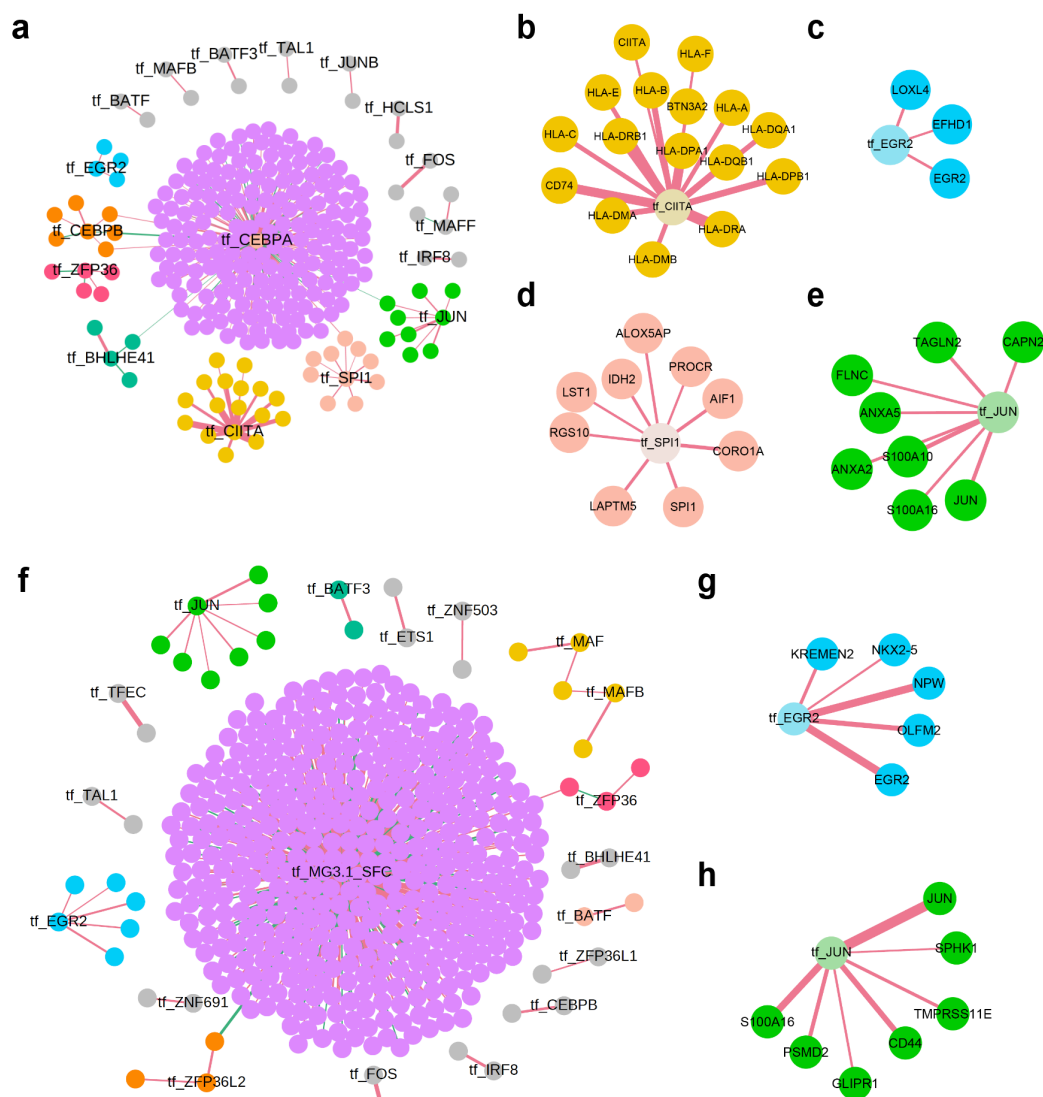525 shows a fast rise and slow decay pattern of the intracellular calcium concentration.

526

527

528

529   **Figure 4.** Transcriptome analysis of TFiMGLs on different days and under disease relevant
530   stimulations. (**a**) PCA of bulk RNA-seq data from multiple sources containing primary microglia.
531   MG: Microglia; DC: dendritic cell; HPC: hematopoietic progenitor; iMGL: growth factor-induced
532   microglia-like cell; Mono: monocyte; iPS: induced pluripotent stem cell. (**b**) GSEA of TFiMGLs
533   versus iPS using two microglia marker gene sets from MSigDB: M40168 and M39077. (**c**)
534   Heatmap and clustering with 881 microglia specific genes previously reported (Ref. Gosselin). (**d**)
535   PCA of TFiMGLs' transcriptome after 24 hours treatment with IFNγ, fAβ, or TDP43. (**e-f**) Pathway
536   analysis of significantly differentially expressed genes after treatment with IFNγ or TDP43.

537

538

539

540 **Figure 5.** Single-cell atlas reference mapping confirmed microglia-like fate induction. (**a**)
541 Harmony-integrated single-cell reference atlas with 225 samples from 59 organ or tissue types,
542 with a total of 1,004,650 single cells. UMAP plot is colored by sample. Sample annotation is
543 acquired by manual curation of each study. (**b**) UMAP is colored by the cellular level annotation
544 from the "Primary_brain" dataset downloaded from "Organoid Report Card". Red cells are
545 microglia. The red dashed box highlighted the cluster containing microglia. (**c**) Symphony
546 projection of cells from the first pooled screen (PGP1+40TFs) to the reference atlas. The red
547 dashed box highlighted the 4.3% cells mapped to the microglia-containing cluster. (**d**) Symphony
548 projection of cells from the second pooled screen (MG3.1+42TFs) to the reference atlas. The red
549 dashed box highlighted the 26.5% cells mapped to the microglia-containing cluster. (**e**-**h**)
550 Expression of key TF barcodes in the projected cells. CEBPA, SPI1, and the SFC cassette
551 showed high expression in the microglia-containing cluster.

552

553

554

555
556

557     **Figure 6.** TF-gene regression analysis for studying causal gene-regulatory networks. Nodes' label
558     starts with "tf_" represent the TF captured by single-cell barcode sequencing. All other nodes
559     represent genes captured by scRNA-seq. The width of edges is correlated with coefficient values.
560     The larger the value the wider the edge. A red edge means upregulation while a green edge
561     means downregulation. Edges were selected with these criteria: Abs(coefficient) > 0.1 & $-\log_{10}$(p-
562     value) > 20. (**a**) Global network for the first pooled screen. Sub-network for (**b**) tf_CIITA, (**c**)
563     tf_EGR2, (**d**) tf_SPI1, (**e**) tf_JUN. (**f**) Global network for the second pooled screen. Sub-network
564     for (**g**) tf_EGR2 and (**h**) tf_JUN.

565

573

574 **Author contributions:** G.M.C. and S.L. conceived the project. S.L., with guidance from A.H.M.N.
575 and P.K., performed early exploratory experiments. S.L., G.M.C., L.L., F.Z. designed overall
576 experimental and analytical strategies. S.L. performed experiments with help from B.S., Y.C., E.A.,
577 M.G-C., C-T.W., J.H., Y.T., and G.C. L.L. helped with single-cell TF barcode quantification and
578 regression analysis. F.Z. and S.R. helped with cell atlas integration and mapping. J.A., J.T., E.L.
579 provided significant discussion and input over project design. S.L. wrote the manuscript with help
580 from L.L., and with input and feedback from all authors.

581

582 **Competing financial interests:** G.M.C, P.K., and A.H.M.N. are co-founders of and have equity
583 in GC Therapeutics, Inc. Full disclosure for GMC is available at
584 arep.med.harvard.edu/gmc/tech.html.

585

586 **Code availability:** Code used in this study for bulk RNA-seq, scRNA-seq, reference mapping,
587 and regression analysis will be made available upon publication.

588

589 **Data availability:** Bulk RNA-seq data, scRNA-seq data, and the single-cell reference atlas will
590 be made available upon publication.

591
592

## Online Methods

**Barcoded TF expression vector construction.** All TFs used in this study were obtained from the TFome collection[2] in pDONR format. For expression in hiPSCs, a PiggyBac integrating Dox-inducible vector pBAN[2] was used. To create barcoded pBAN expression vector (pBAN-BC), the original pBAN was digested with AgeI and KpnI, followed by ligation of a gBlock (IDT DNA) containing the same excised piece with an additional 20-bp random barcode. After bacteria transformation, individual colonies were expanded and extracted for plasmid DNA. Gateway cloning was used to transfer each TF from pDONR to pBAN-BC vector. Barcode sequence for each TF was confirmed by Sanger sequencing.

**Cell culture.** hiPSCs were culture in mTeSR Plus media (Stemcell Technologies, 100-0276) on multi-wells plates coated with Matrigel (Corning, 354277) or Cultrex (Bio-Techne Corporation, 3434-005-02). For passaging, cells were dissociated with TrypLE Express (Life Technologies, 12604013) and seeded into fresh plate and media containing 10 µM Y-27632 ROCK inhibitor (Millipore, 688001) for 24 hours. Daily media change was performed until cells were ready for another passaging or downstream experiments.

**Nucleofection, TF integration and differentiation.** TF (pBAN-TF-BC) and Super PiggyBac (SPB) Transposase (System Biosciences, PB210PA-1) expression vectors were mixed at a mass ration of 4:1 and transfected into hiPSCs using P3 Primary Cell 4D-Nucleofector X Kit L (Lonza, V4XP-3024) on a 4D-Nucleofector X Unit (Lonza, AAF-1002X) following manufacturer's instructions. For the two pooled TF screenings, 600,000 cells were transfected with 5 µg of DNA and seeded into one well of a 6-well plate. For individual TF combinations, 120,000 cells were transfected with 2.5 µg of DNA and seeded into one well of a 12-well plate. Program CB150 was used for the nucleofections. 48 hours after nucleofection, 1 µg/mL of puromycin (Gibco, A1113803) or 50 µg/mL of zeocin (Gibco, R25001) was added to the culture for the selection of TF-integrated cells. Cells were passaged again when reaching 80% confluency. For induction of TF expression, cells were seeded into mTeSR Plus media containing 0.5 µg/mL doxycycline (Sigma-Aldrich, D3072) and 10 µM Y-27632 ROCK inhibitor and were changed into media only containing doxycycline after 24 hours.

**Flow cytometry and sorting.** For cytometry analysis, cells were dissociated with TrypLE Express for 5 minutes at 37 degree, diluted with twice the volume of Cell Staining Buffer (Biolegend, 420201) and centrifuged at 200 g for 3 minutes to remove the digesting enzyme. Cells were then incubated with 25 µg/mL of Human Fc Block (BD Biosciences, 564219) diluted in Cell Staining Buffer for 15 minutes on ice, followed immediately by staining with fluorescently conjugated antibodies or isotype controls at proper dilution for 30 minutes on ice. Antibodies were diluted in Cell Staining Buffer and Human Fc Block was not removed from the mixture. After antibody staining, cells were washed twice with Cell Staining Buffer before being put through 35 µm nylon mesh into a 5 mL round bottom polystyrene tube (Falcon, 352235). Flow cytometry data was acquired on a BD LSRFortessa Cell Analyzer. For cell sorting, the staining protocol was the same except for that Cell Staining Buffer was replaced with mTeSR Plus media in order to maintain the

632  best viability of cells. Cell sorting was performed on a BD FACSAria Cell Sorter. Flow cytometry
633  antibodies used in this study were: FITC-TRA-1-60 (BD Biosciences, 560380), BV421-CX3CR1
634  (Biolegend, 341620), PE-P2RY12 (Biolegend, 392104), APC-CD11b (Biolegend, 101212).
635  Isotype controls used were: BV421- Rat IgG2b (Biolegend, 400640), PE-Mouse IgG2a (Biolegend,
636  400214), APC- Rat IgG2b (Biolegend, 400612).

637  **scRNA-seq library preparation.** scRNA-seq experiments were performed using 10x Genomics
638  Chromium Single Cell 3' Reagent Kits v3 or v3.1 following the manufacturer's instruction. 5000
639  single cells were calculated as targeted input for each sample. For the first iteration, 10% of stem
640  cells were spiked in as undifferentiated control. For the second iteration, 5% stem cells and 5%
641  MG3.1-SFC were spiked in as undifferentiated and initial differentiation control. The only
642  modification made to the protocol was at the Sample Index PCR step, where 5 µL of the PCR mix
643  was taken out and mixed with 0.5 µL 1000X SYBR Gold (Invitrogen, S11494) for a qPCR reaction.
644  The optimal amplification cycle was determined as the cycle just before half maximum of the total
645  signal. Final libraries were sequenced on NextSeq 500 or NovaSeq with a goal of at least 30,000
646  reads per cell.

647  **TF barcode amplicon library preparation.** Because after the cDNA amplification step in the 10x
648  scRNA-seq protocol the amplicons contained cell barcodes, UMIs, and TF barcodes, these
649  cDNAs could be used as the template for further amplification of TF-cell barcodes. Two sequential
650  PCR reactions were performed, each was accompanied by a SYBR Gold spike-in qPCR to
651  determine the optimal cycle number as described in "scRNA-seq library preparation". For PCR1,
652  NGS10x-F-i7-BC-PCR1F and i5000 were used as primers. A 50 µL PCR1 reaction contains 25
653  µL Q5 Hot Start High-Fidelity 2X Master Mix (New England Biolabs, M0494L), 5 µL amplified
654  cDNA, 2.5 µL of both primers at 10 µM stock concentration, and 15 µL nuclease-free water. PCR1
655  program was initial denaturation, 98 degrees, 30 seconds; 11-13 cycles (qPCR determined) of 98
656  degrees, 10 seconds, 67 degrees, 30 seconds, 72 degrees, 30 seconds; final extension, 72
657  degrees, 2 minutes. PCR1 reaction was purified with 1.2X SPRIselect beads (Beckman Coulter,
658  B23318) following standard protocol. The sample was eluted in 20 µL water. For PCR2, i7000,
659  P5, and P7 were used as primers. A 50 µL PCR2 reaction contains 25 µL Q5 Hot Start High-
660  Fidelity 2X Master Mix, 10 µL PCR1 product, 2.5 µL of all three primers at 10 µM stock
661  concentration, and 7.5 µL nuclease-free water. PCR2 program was initial denaturation, 98
662  degrees, 30 seconds; 4-5 cycles (qPCR determined) of 98 degrees, 10 seconds, 67 degrees, 30
663  seconds, 72 degrees, 30 seconds; final extension, 72 degrees, 2 minutes. PCR2 product was
664  purified the same as PCR1. Final libraries were submitted for MiSeq v3 with paired-end reads of
665  80 cycles from either direction.

666  Primer sequences:

667  NGS10x-F-i7-BC-PCR1F:
668  GGAGTTCAGACGTGTGCTCTTCCGATCTCTTTTCCAAGCACCTGCTACATAG

669     i5000:

670     AATGATACGGCGACCACCGAGATCTACACaactcgctACACTCTTTCCCTACACGACGCTCTTC

671     CGATCT (lower case region represents a sample-specific barcode)

672     i7000:

673     CAAGCAGAAGACGGCATACGAGATtcgccttaGTGACTGGAGTTCAGACGTGTGCTCTTCCGA

674     TCT (lower case region represents a sample-specific barcode)

675     P5: AATGATACGGCGACCACCGA

676     P7: CAAGCAGAAGACGGCATACGA

677     **Analysis of scRNA-seq and TF barcode-seq data.** For scRNA-seq, raw FASTQ files were
678     aligned to GRCh38 and quantified using Cell Ranger. Detailed information about cell number,
679     read depth and gene detected is visualized in **Supplementary Figure S2** and **S6**. Seurat was
680     used to performed cell filtering, data normalization and clustering. The generated Seurat object
681     also contained the single-cell raw expression matrix for all genes. For TF barcode-seq, in the
682     paired-end MiSeq data, one of the read pair contains the 20 bp TF barcode while the other one
683     contains the 16 bp cell barcode and the 12 bp UMI. By matching the names of the reads within
684     the pair, three sequences were compiled into one table with three columns: TF-BC, cell-BC, UMI.
685     To remove duplicated reads from the same molecule, duplicated rows that has the same value
686     for all three columns were removed. Then the table was counted and reshaped into a frequency
687     table where the row names represent cell and column names represent TF. This table contains
688     the raw counts of each TF barcode in all single cells. Because the TF barcodes were amplified
689     from the cDNA during library preparation, we normalized the TF barcode count with the number
690     of total RNA UMIs detected in each cell, reasoning that cells with more total UMIs were likely to
691     have more reads for TF barcode. The raw gene expression matrix and normalized TF count matrix
692     were used to identify which TF barcodes were likely to induce microglial gene expression.
693     Specifically, the expression of microglial genes was binarized, with any cell had a non-zero
694     expression being 1. Then between the two groups of cells 0 or 1 microglial gene expression, a
695     Wilcoxon rank sum test was performed for all barcoded TFs to determine which TF(s) had a higher
696     expression in cells expressing microglial genes. The TFs were ranked by -log10(p-value).

697     **Bulk RNA-seq library preparation.** Cultured cells were dissolved directly with TRIzol (Thermo
698     Fisher Scientific, 15596018) for total RNA purification with Direct-zol RNA MiniPrep Kit (Zymo
699     Research, R2050). RNA concentration was quantified with Qubit RNA HS Assay Kit (Thermo
700     Fisher Scientific, Q32852). RNA integrity was confirmed by presence of 18S and 28S bands on a
701     2% E-Gel EX Agarose Gel (Thermo Fisher Scientific, G402002). Between 100 ng to 1000 ng total
702     RNA was used as input for mRNA enrichment using NEBNext Poly(A) mRNA Magnetic Isolation
703     Module (New England Biolabs, E7490), followed by library construction with NEBNext Ultra II
704     Directional RNA Library Prep Kit (New England Biolabs, E7760S) following the manufacturer's
705     instructions. Biopolymers Facility at Harvard Medical School performed library QC and
706     sequencing.

**Analysis of bulk RNA-seq data.** For both in-house generated sample and datasets downloaded from GEO, raw FASTQ files were aligned to GRCh38 and quantified using STAR 2.5.2b. Regularized-logarithm (rlog) transformation was applied to the raw counts before visualization using PCA. For analysis where data from multiple sources were involved, ComBat-seq was used for batch correction before PCA. Differential gene expression analysis was conducted with DESeq2[47]. Pathway enrichment and GSEA analysis were performed with clusterProfiler[80].

**Immunofluorescence (IF).** IF experiments were performed in µ-Plate 96 Well Black plate (ibidi, 89626). After media removal, cells were fixed with 4% paraformaldehyde (Electron Microscopy Sciences, 15710) in 1x phosphate buffered saline (PBS) (Thermo Fisher Scientific, 10010072) for 15 minutes at room temperature (RT). Cells were rinsed three times with PBS before proceeding to permeabilization or blocking. For staining of Oct-3/4 and PU.1, cells were permeabilized, while not for cell surface proteins' staining. Permeabilization was conducted with 0.25% Triton-X-100 (Thermo Fisher Scientific, 85111) in 1x PBS for 15 minutes at RT followed by three rinses with PBS. Cells were then blocked with 1% bovine serum albumin (BSA) in PBS for one hour at RT. For primary and secondary antibody staining, antibodies were diluted in PBS with 1% BSA and incubated with cells for one hour at RT. Three 5-minute washes with PBS were used to remove excessive antibodies after staining. Cells were directly imaged in plate on a Nikon Ti2 Eclipse inverted microscope with a Plan Apo Lambda DM 60× (1.4 NA, Ph3) oil objective and an Andor Zyla sCMOS camera. Images were acquired by NIS-Element AR software. All antibodies were used at 1:200 dilution. Primary IF antibodies used in this study were: Oct-3/4 (Santa Cruz Biotechnology, sc-5279), PU.1 (Thermo Fisher Scientific, PA5-17505), CD11b (BioLegend, 101202), P2RY12 (Thermo Fisher Scientific, 702516), CX3CR1 (Abcam, ab8021).

**Phagocytosis assay.** Differentiated cells were incubated with 20 µg/mL of pHrodo Green S. aureus BioParticles (Thermo Fisher Scientific, P35382) for 0-4 hours in mTeSR Plus media in the presence of 100 µg/ml Penicillin-Streptomycin (Corning, 30-002-CI). After removal of excessive particles with PBS washes, cells were harvested for antibody (CX3CR1, P2RY12, CD11b) staining and flow cytometry analysis as described in previous section. Remaining stained cells after flow cytometry was transferred into µ-Plate 96 Well Black plate for fluorescence microscopy to confirm the intracellular localization of the particles. This step needs to be conducted swiftly after flow cytometry in order to avoid changing of cellular morphology due to cell death.

**Calcium imaging.** Calcium imaging experiment was conducted in standard 12-well cell culture plates. Differentiated cells were incubated with 1 µg/mL Fluo-4 AM calcium indicator (Thermo Fisher Scientific, F23917) in 1 mL of mTeSR Plus media for 30 minutes in a cell culture incubator. Excessive dye was washed away with two 1 mL media washes. After adding 1 mL of fresh mTeSR Plus, the cells were put on stage in a microscope inside the incubator. Images acquisition started without stimulation for 90 seconds to determine baseline signal. One image was acquired every three seconds, the fastest possible on the instrument. After 90 seconds 1 mL of media containing 150 µM ADP was added to the cells while imaging was continuing. The total length of imaging was 10 minutes. Fluorescent signal was quantified and plotted in MATLAB.

746 **Amyloid fibrillation.** Aβ fibrillation experiments were performed using SensoLyte Thioflavin T β-
747 Amyloid (1-42) Aggregation Kit (AnaSpec, AS-72214) according to manufacturer's instruction.
748 The reaction was set up in μ-Plate 96 Well Black plate. Data was acquired on a plate reader with
749 excitation/emission = 440 nm/484 nm at 37 degree once every 5 minutes for 3 hours.

750 **Preparation of datasets for building single-cell reference atlas.** Files containing raw counts
751 of 10x Genomics Chromium scRNA-seq data for different human tissues were download from
752 PanglaoDB (https://panglaodb.se/index.html). Human primary brain single-cell data from
753 gestational weeks 6-22 were downloaded from Organoid Report Card
754 (https://cells.ucsc.edu/?ds=organoidreportcard). Human endometrium single-cell data were
755 download from GEO GSE111976. All sample went through manual cell filtering using Seurat with
756 different filters on number of gene, UMI, and percentage of mitochondria genes (**Supplementary**
757 **Table S4**). Tissue annotation was compiled through manual curation of each study by checking
758 what tissue/cell types were used (**Supplementary Table S3**). All raw counts were merged into
759 one sparse matrix, which was then used as input for data integration.

760 **Single-cell atlas integration and mapping.** Data integration and projection using Harmony and
761 Symphony was carried out following instructions from the authors on GitHub
762 (https://github.com/immunogenomics/harmony; https://github.com/immunogenomics/symphony).
763 For Harmony integration, RunHarmony function was used. Parameters different from default
764 settings were epsilon.cluster=-Inf, epsilon.harmony=-Inf. Batch correction were performed based
765 on tissue types labeld in Figure 5a. For reference mapping with Symphony,
766 buildReferenceFromSeurat function was used to create the reference object, and mapQuery
767 function was used to map iPSC-derived cells to the reference atlas. Due to the size of the data,
768 these steps were performed on the O2 cluster of Harvard Medical School with at least 180 Gb
769 memory and 8 cores. Most R objects along the pipeline could be saved as standard R files for
770 repeated analysis, except for the UMAP model file, which required saving and loading through
771 the "uwot" package[81]. Code used for integration and projection, together with key reference and
772 annotations files that could be of use for future explorations are shared along with this manuscript.

773 **TF-gene stepwise regression model construction.** The stepwise regression (or stepwise
774 selection) is a regression model that iteratively adds and removes predictors in the predictive
775 model to find the subset of variables in the data set resulting in the best performance, and
776 consequently lowering the perdition error in the model. During the process, the value of the
777 statistical test is used to screen the variables. If the value is less than or equal to 0.05, then the
778 variable enters the regression model, and the selected variable is the independent variable of the
779 regression model. For the construction of the model:

780 Step 1: Establish $P$ regression models between the independent variables $X_1, X_2, \ldots, X_p$
781 ($number = P$) and the dependent variable $Y$ respectively,

782 $$Y = \beta_0 + \beta_i X_i + \epsilon, i = 1, \ldots p$$

783 Calculate the statistical value of the F-test with the regression coefficient $F_1^{(1)}, \ldots, F_p^{(1)}$, and take
784 the maximum value $F_{i1}^{(1)}$,

785
$$F_{i1}^{(1)} = \max\left\{F_1^{(1)}, \dots, F_p^{(1)}\right\}$$

786  For a given significance level $\alpha$, the threshold value is $F^1$. If $F_{i1}^{(1)} > F^1$, then $X_{i1}$ will be included
787  in the regression model and recorded as the set of selected variable indicators as $I_1$.
788  Step 2: Establish a binary regression model of the dependent variable $Y$ and the independent
789  variable subset $\{X_{i1}, \dots, X_1\}$, $\{X_{i1}, \dots, X_{i1-1}\}$, $\{X_{i1}, \dots, X_{i1+1}\}$, calculate the statistical value of the F-
790  test with the regression coefficient $F_k^{(2)}$ and take the maximum value $F_{i2}^{(2)}$,

791
$$F_{i2}^{(2)} = \max\left\{F_1^{(2)}, \dots F_{i1-1}^{(2)}, F_{i1+1}^{(2)}, \dots, F_p^{(2)}\right\}$$

792  For a given significance level $\alpha$, record the corresponding critical value as $F^{(2)}$. If $F_{i2}^{(2)} > F^{(2)}$,
793  then the variable is introduced into the regression model. Otherwise, the variable introduction
794  process is terminated.
795  Step 3: Repeat Step 2 with the subset of variables $\{X_{i1}, X_{i2}, X_k\}$. This step is repeated by selecting
796  an independent variable that is not introduced into the regression model until the test does not
797  introduce any variables.
798  **TF-gene network visualization.** Both p-values and coefficients in the regression analysis work
799  together to represent relationships in the model about the significant factors. The coefficients
800  describe the mathematical relationship between each independent variable and the dependent
801  variable. The p-values for the coefficients indicate whether these relationships are statistically
802  significant. We selected the 250 TF-gene combinations from the first pooled screen and 422 in
803  the second with the criteria Abs(coefficient) > 0.1 & -$\log_{10}$(p-value) > 20, then visualized them in
804  GEPHI. Nodes are the genes and TFs, edges are the regression coefficients (an unstandardized
805  effect size because they indicate the strength of the relationship between variables), and colors
806  are based on the modularity from the network module algorithm[82].
807

## References

1. Regev, A., Teichmann, S. A., Lander, E. S., Amit, I. & Benoist, C. Science forum: the human cell atlas. *elife* (2017).

2. Ng, A. H. M. *et al.* A comprehensive library of human transcription factors for cell fate engineering. *Nat. Biotechnol.* **39,** 510–519 (2021).

3. Ginhoux, F. *et al.* Fate mapping analysis reveals that adult microglia derive from primitive macrophages. *Science* **330,** 841–845 (2010).

4. Hoeffel, G. *et al.* C-Myb(+) erythro-myeloid progenitor-derived fetal monocytes give rise to adult tissue-resident macrophages. *Immunity* **42,** 665–678 (2015).

5. Crotti, A. & Ransohoff, R. M. Microglial Physiology and Pathophysiology: Insights from Genome-wide Transcriptional Profiling. *Immunity* **44,** 505–515 (2016).

6. Nayak, D., Roth, T. L. & McGavern, D. B. Microglia development and function. *Annu. Rev. Immunol.* **32,** 367–402 (2014).

7. Nimmerjahn, A., Kirchhoff, F. & Helmchen, F. Resting microglial cells are highly dynamic surveillants of brain parenchyma in vivo. *Science* **308,** 1314–1318 (2005).

8. Matcovitch, O. Microglia development follows a stepwise program to regulate brain homeostasis. *Natan*

9. Salter, M. W. & Beggs, S. Sublime microglia: expanding roles for the guardians of the CNS. *Cell* **158,** 15–24 (2014).

10. Colonna, M. & Butovsky, O. Microglia function in the central nervous system during health and neurodegeneration. *Annu. Rev. Immunol.* **35,** 441–468 (2017).

11. Salter, M. W. & Stevens, B. Microglia emerge as central players in brain disease. *Nat. Med.* **23,** 1018–1027 (2017).

12. Mathys, H. *et al.* Single-cell transcriptomic analysis of Alzheimer's disease. *Nature* **570,** 332–337 (2019).

13. Keren-Shaul, H. *et al.* A Unique Microglia Type Associated with Restricting Development of Alzheimer's Disease. *Cell* **169,** 1276–1290.e17 (2017).

14. Yeh, F. L., Hansen, D. V. & Sheng, M. TREM2, microglia, and neurodegenerative diseases. *Trends Mol. Med.* **23,** 512–533 (2017).

15. Ulland, T. K. *et al.* TREM2 maintains microglial metabolic fitness in alzheimer's disease. *Cell* **170,** 649–663.e13 (2017).

16. Dello Russo, C. *et al.* The human microglial HMC3 cell line: where do we stand? A systematic literature review. *J. Neuroinflammation* **15,** 259 (2018).

17. Timmerman, R., Burm, S. M. & Bajramovic, J. J. An Overview of in vitro Methods to Study Microglia. *Front. Cell Neurosci.* **12,** 242 (2018).

18. Smith, A. M. & Dragunow, M. The human side of microglia. *Trends Neurosci.* **37,** 125–135 (2014).

19. Watkins, L. R. & Hutchinson, M. R. A concern on comparing'apples' and'oranges' when differences between microglia used in human and rodent studies go far, far beyond simply species …. *Trends in neurosciences* (2014).

20. Muffat, J. *et al.* Efficient derivation of microglia-like cells from human pluripotent stem cells. *Nat. Med.* **22,** 1358–1367 (2016).

21. Abud, E. M. *et al.* iPSC-Derived Human Microglia-like Cells to Study Neurological Diseases. *Neuron* **94,** 278–293.e9 (2017).

22. Pandya, H. *et al.* Differentiation of human and murine induced pluripotent stem cells to microglia-like cells. *Nat. Neurosci.* **20,** 753–759 (2017).

23. Haenseler, W. *et al.* A Highly Efficient Human Pluripotent Stem Cell Microglia Model Displays a Neuronal-Co-culture-Specific Expression Profile and Inflammatory Response. *Stem Cell Rep.* **8,** 1727–1742 (2017).

24. Douvaras, P. *et al.* Directed differentiation of human pluripotent stem cells to microglia. *Stem Cell Rep.* **8,** 1516–1524 (2017).

25. Takata, K. *et al.* Induced-Pluripotent-Stem-Cell-Derived Primitive Macrophages Provide a Platform for Modeling Tissue-Resident Macrophage Differentiation and Function. *Immunity* **47,** 183–198.e6 (2017).

26. McQuade, A. *et al.* Development and validation of a simplified method to generate human microglia from pluripotent stem cells. *Mol. Neurodegener.* **13,** 67 (2018).

27. Speicher, A. M., Wiendl, H., Meuth, S. G. & Pawlowski, M. Generating microglia from human pluripotent stem cells: novel in vitro models for the study of neurodegeneration. *Mol. Neurodegener.* **14,** 46 (2019).

28. Xu, R. *et al.* Human iPSC-derived mature microglia retain their identity and functionally integrate in the chimeric mouse brain. *Nat. Commun.* **11,** 1577 (2020).

29. Vierbuchen, T. *et al.* Direct conversion of fibroblasts to functional neurons by defined factors. *Nature* **463,** 1035–1041 (2010).

30. Tsunemoto, R. *et al.* Diverse reprogramming codes for neuronal identity. *Nature* **557,** 375–380 (2018).

31. Busskamp, V. *et al.* Rapid neurogenesis through transcriptional activation in human stem cells. *Mol. Syst. Biol.* **10,** 760 (2014).

32. Kierdorf, K. *et al.* Microglia emerge from erythromyeloid precursors via Pu.1- and Irf8-dependent pathways. *Nat. Neurosci.* **16,** 273–280 (2013).

33. Smith, A. M. *et al.* The transcription factor PU.1 is critical for viability and function of human brain microglia. *Glia* **61,** 929–942 (2013).

34. Gosselin, D. *et al.* An environment-dependent transcriptional network specifies human microglia identity. *Science* **356,** (2017).

35. Galatro, T. F. *et al.* Transcriptomic analysis of purified human cortical microglia reveals age-associated changes. *Nat. Neurosci.* **20,** 1162–1171 (2017).

36. Zhong, S. *et al.* A single-cell RNA-seq survey of the developmental landscape of the human prefrontal cortex. *Nature* **555,** 524–528 (2018).

37. Olah, M. *et al.* A transcriptomic atlas of aged human microglia. *Nat. Commun.* **9,** 539 (2018).

38. Butovsky, O. *et al.* Identification of a unique TGF-β-dependent molecular and functional signature in microglia. *Nat. Neurosci.* **17,** 131–143 (2014).

39. Wehrspaun, C. C., Haerty, W. & Ponting, C. P. Microglia recapitulate a hematopoietic master regulator network in the aging human frontal cortex. *Neurobiol. Aging* **36,** 2443.e9-2443.e20 (2015).

40. Avellino, R. & Delwel, R. Expression and regulation of C/EBPα in normal myelopoiesis and in malignant transformation. *Blood* **129,** 2083–2091 (2017).

41. Lichtinger, M. *et al.* RUNX1 reshapes the epigenetic landscape at the onset of haematopoiesis. *EMBO J.* **31,** 4318–4333 (2012).

42. Starck, J. *et al.* Spi-1/PU.1 is a positive regulator of the Fli-1 gene involved in inhibition of erythroid differentiation in friend erythroleukemic cell lines. *Mol. Cell. Biol.* **19,** 121–135 (1999).

43. Hoeffel, G. & Ginhoux, F. Ontogeny of Tissue-Resident Macrophages. *Front. Immunol.* **6,**

486 (2015).

44. Bazan, J. F. *et al.* A new class of membrane-bound chemokine with a CX3C motif. *Nature* **385,** 640–644 (1997).

45. Hughes, P. M., Botham, M. S., Frentzel, S., Mir, A. & Perry, V. H. Expression of fractalkine (CX3CL1) and its receptor, CX3CR1, during acute and chronic inflammation in the rodent CNS. *Glia* **37,** 314–327 (2002).

46. Liu, Z. *et al.* Systematic comparison of 2A peptides for cloning multi-genes in a polycistronic vector. *Sci. Rep.* **7,** 2193 (2017).

47. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15,** 550 (2014).

48. Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set collection. *Cell Syst.* **1,** 417–425 (2015).

49. Cahan, P. *et al.* CellNet: network biology applied to stem cell engineering. *Cell* **158,** 903–915 (2014).

50. Kracht, L. *et al.* Human fetal microglia acquire homeostatic immune-sensing properties early in development. *Science* **369,** 530–537 (2020).

51. Inoue, K. Purinergic systems in microglia. *Cell Mol. Life Sci.* **65,** 3074–3080 (2008).

52. Di Virgilio, F., Ceruti, S., Bramanti, P. & Abbracchio, M. P. Purinergic signalling in inflammation of the central nervous system. *Trends Neurosci.* **32,** 79–87 (2009).

53. Zhang, Y., Parmigiani, G. & Johnson, W. E. ComBat-seq: batch effect adjustment for RNA-seq count data. *NAR Genom. Bioinform.* **2,** lqaa078 (2020).

54. Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. USA* **102,** 15545–15550 (2005).

55. Ivashkiv, L. B. IFNγ: signalling, epigenetics and roles in immunity, metabolism, disease and cancer immunotherapy. *Nat. Rev. Immunol.* **18,** 545–558 (2018).

56. Zhong, L. *et al.* Amyloid-beta modulates microglial responses by binding to the triggering receptor expressed on myeloid cells 2 (TREM2). *Mol. Neurodegener.* **13,** 15 (2018).

57. Zhao, W. *et al.* TDP-43 activates microglia through NF-κB and NLRP3 inflammasome. *Exp. Neurol.* **273,** 24–35 (2015).

58. Zhang, F. *et al.* IFN-γ and TNF-α drive a CXCL10+ CCL2+ macrophage phenotype expanded in severe COVID-19 lungs and inflammatory diseases with tissue inflammation. *Genome Med.* **13,** 64 (2021).

59. Franzén, O., Gan, L.-M. & Björkegren, J. L. M. PanglaoDB: a web server for exploration of mouse and human single-cell RNA sequencing data. *Database (Oxford)* **2019,** (2019).

60. Bhaduri, A. *et al.* Cell stress in cortical organoids impairs molecular subtype specification. *Nature* **578,** 142–148 (2020).

61. Wang, W. *et al.* Single-cell transcriptomic atlas of the human endometrium during the menstrual cycle. *Nat. Med.* **26,** 1644–1653 (2020).

62. Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat. Methods* **16,** 1289–1296 (2019).

63. Kang, J. B. *et al.* Efficient and precise single-cell reference atlas mapping with Symphony. *Nat. Commun.* **12,** 5890 (2021).

64. Van de Sande, B. *et al.* A scalable SCENIC workflow for single-cell gene regulatory network analysis. *Nat. Protoc.* **15,** 2247–2276 (2020).

65. Hecker, M., Lambeck, S., Toepfer, S., van Someren, E. & Guthke, R. Gene regulatory

network inference: data integration in dynamic models-a review. *Biosystems* **96,** 86–103 (2009).

66. Matys, V. *et al.* TRANSFAC and its module TRANSCompel: transcriptional gene regulation in eukaryotes. *Nucleic Acids Res.* **34,** D108-10 (2006).

67. Castro-Mondragon, J. A. *et al.* JASPAR 2022: the 9th release of the open-access database of transcription factor binding profiles. *Nucleic Acids Res.* **50,** D165–D173 (2022).

68. Kanehisa, M., Furumichi, M., Sato, Y., Ishiguro-Watanabe, M. & Tanabe, M. KEGG: integrating viruses and cellular organisms. *Nucleic Acids Res.* **49,** D545–D551 (2021).

69. Liu, Z.-P., Wu, C., Miao, H. & Wu, H. RegNetwork: an integrated database of transcriptional and post-transcriptional regulatory networks in human and mouse. *Database (Oxford)* **2015,** (2015).

70. Oughtred, R. *et al.* The BioGRID database: A comprehensive biomedical resource of curated protein, genetic, and chemical interactions. *Protein Sci.* **30,** 187–200 (2021).

71. Szklarczyk, D. *et al.* The STRING database in 2021: customizable protein-protein networks, and functional characterization of user-uploaded gene/measurement sets. *Nucleic Acids Res.* **49,** D605–D612 (2021).

72. Maeda, T., Sugita, S., Kurimoto, Y. & Takahashi, M. Trends of Stem Cell Therapies in Age-Related Macular Degeneration. *J Clin Med* **10,** (2021).

73. de Klerk, E. & Hebrok, M. Stem Cell-Based Clinical Trials for Diabetes Mellitus. *Front. Endocrinol. (Lausanne)* **12,** 631463 (2021).

74. Shi, Y., Inoue, H., Wu, J. C. & Yamanaka, S. Induced pluripotent stem cell technology: a decade of progress. *Nat. Rev. Drug Discov.* **16,** 115–130 (2017).

75. Flitsch, L. J., Laupman, K. E. & Brüstle, O. Transcription Factor-Based Fate Specification and Forward Programming for Neural Regeneration. *Front. Cell Neurosci.* **14,** 121 (2020).

76. Takahashi, K. & Yamanaka, S. Induction of pluripotent stem cells from mouse embryonic and adult fibroblast cultures by defined factors. *Cell* **126,** 663–676 (2006).

77. Haniffa, M. *et al.* A roadmap for the Human Developmental Cell Atlas. *Nature* **597,** 196–205 (2021).

78. Wang, L. *et al.* Stoichiometry of Gata4, Mef2c, and Tbx5 influences the efficiency and quality of induced cardiac myocyte reprogramming. *Circ. Res.* **116,** 237–244 (2015).

79. Tran, H. T. N. *et al.* A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome Biol.* **21,** 12 (2020).

80. Wu, T. *et al.* clusterProfiler 4.0: A universal enrichment tool for interpreting omics data. *Innovation (N Y)* **2,** 100141 (2021).

81. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv* (2018). doi:10.48550/arxiv.1802.03426

82. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in large networks. *J. Stat. Mech.* **2008,** P10008 (2008).