1  **Title:** PURE-seq identifies *Egr1* as a Potential Master Regulator in Murine Aging by Sequencing

2  Long-Term Hematopoietic Stem Cells

3

4  **Authors:** Sixuan Pan[†,1], Kai-Chun Chang[†,1], Inés Fernández-Maestre[†,2,3], Stéphane Van Haver[4,5],

5  Matthew G. Wereski[2], Robert L. Bowman[6], Ross L. Levine[#,2,7,8], Adam R. Abate[#,1]

6

7  **Affiliations**

8  [1]Department of Bioengineering, University of California San Francisco, San Francisco, CA 94143,

9  USA.

10  [2]Human Oncology and Pathogenesis Program, Memorial Sloan Kettering Cancer Center, New

11  York, NY, USA.

12  [3]Louis V. Gerstner Jr Graduate School of Biomedical Sciences, Memorial Sloan Kettering Cancer

13  Center, New York, NY, USA.

14  [4]Molecular Pharmacology Program, Memorial Sloan Kettering Cancer Center, New York, NY,

15  USA.

16  [5]Tow Center for Developmental Oncology, Memorial Sloan Kettering Cancer Center, New York,

17  NY, USA.

18  [6]Department of Cancer Biology, Perelman School of Medicine, University of Pennsylvania,

19  Philadelphia, PA, USA.

20  [7]Department of Medicine, Weill Cornell Medical College, New York, NY, USA.

21  [8]Center for Hematologic Malignancies, Memorial Sloan Kettering Cancer Center, New York, NY,

22  USA.

23

24  [†] These authors contributed equally

25  # Corresponding author

26  Correspondence: arabate@gmail.com, leviner@mskcc.org

27

28

## Abstract

Single-cell transcriptomics is valuable for uncovering individual cell properties, particularly in highly heterogeneous systems. However, this technique often results in the analysis of many well-characterized cells, increasing costs and diluting rare cell populations. To address this, we developed PURE-seq (PIP-seq for Rare-cell Enrichment and Sequencing) for scalable sequencing of rare cells. PURE-seq allows direct cell loading from FACS into PIP-seq reactions, minimizing handling and reducing cell loss. PURE-seq reliably captures rare cells, with 60 minutes of sorting capturing tens of cells at a rarity of 1 in 1,000,000. Using PURE-seq, we investigated murine long-term hematopoietic stem cells and their transcriptomes in the context of hematopoietic aging, identifying *Egr1* as a potential master regulator of hematopoiesis in the aging context. PURE-seq offers an accessible and reliable method for isolating and sequencing cells that are currently too rare to capture successfully with existing methods.

## Introduction

Single-cell transcriptomics is powerful for elucidating the properties of individual cells and can discover phenotypes without relying on predetermined genes or markers. This makes it useful in highly heterogeneous systems with unknown cell properties[1–4]. However, its unbiased nature often leads to the analysis of abundant, well-characterized cellular states at the expense of rare cell populations and increased cost[5,6]. An enrichment method that selectively captures rare cell populations while removing unwanted cells can increase the coverage of rare cells, enabling deeper analysis at the same cost.

Several methods exist for enriching target cells before single-cell sequencing, typically using antibody-based capture approaches to label and isolate cells of interest. Techniques such as fluorescence-activated cell sorting (FACS), magnetic-activated cell sorting (MACS), and cell levitation isolate cells based on expression of specific surface markers[7–9]. However, current single-cell methods do not directly integrate with the output of a flow cytometer, necessitating a transfer step that can result in cell loss or degradation, compromising data quality. This is especially problematic for extremely rare cell applications where the number of captured cells may be too low for reliable transfer. Other alternatives, such as direct cell sorting into well plates or using nanowell array chips, involve labor-intensive workflows and have limited throughput capabilities[10,11]. An ideal approach would allow the flow cytometer to directly load cells into the high-throughput single-cell RNA-sequencing (scRNA-seq) apparatus, minimizing handling, ensuring the highest data quality, and capturing rare cell populations; however, this is not possible with existing methods.

In this paper, we introduce PURE-seq (PIP-seq for Rare-cell Enrichment and Sequencing), a method for sequencing rare cells. PURE-seq is based on our development, Pre-templated Instant Partition sequencing (PIP-seq) [12], which allows scalable scRNA-seq without microfluidics using a fully encapsulated Eppendorf tube. The compact nature of the PIP-seq reservoir and its compatibility with standard Eppendorf tubes, commonly used in flow cytometry, enable direct cell loading from the flow cytometer into the PURE-seq reaction. This eliminates additional handling, reducing cell loss and degradation. The tube is vortexed immediately after cell loading,

3

75    encapsulating, and lysing the cells in droplets within one minute for the PIP-seq single-cell

76    barcoding workflow[12]. This simplicity and minimal handling allow reliable capture of extremely

77    rare cells; 60 minutes of sorting can capture tens of cells at a rarity of 1 in 1,000,000. The rarity of

78    cells captured scales with sorting duration, allowing even rarer cells to be sequenced with more

79    sorting time.

80

81    Using PURE-seq, we analyzed murine long-term hematopoietic stem cells (LT-HSCs), a rare and

82    heterogeneous bone marrow (BM) population difficult to recover in sufficient numbers for high-

83    quality scRNA-seq with current methods[13,14]. PURE-seq enabled us to characterize their

84    transcriptomes in low-input samples. Previous studies hinted at the role of *EGR1* in human LT-

85    HSCs[15,16], but its exact function in mice is unclear. These studies demonstrate higher *EGR1*

86    expression in aged human hematopoietic stem and progenitor cells (HSPCs), suggesting EGR1

87    may regulate quiescence, proliferation, and localization. Attenuated expression of EGR1 might

88    decrease senescence and activate aged HSPCs, offering a potential target for rejuvenation

89    strategies[17]. PURE-seq allowed us to recover sufficient cell numbers to identify *Egr1* as a potential

90    master regulator gene in the aging of murine LT-HSCs. Here, we show that PURE-seq provides a

91    simple workflow to sort and sequence rare cell populations, which is arduous with existing

92    methods, and reliably recapitulates data generated by standard 10x Genomics.
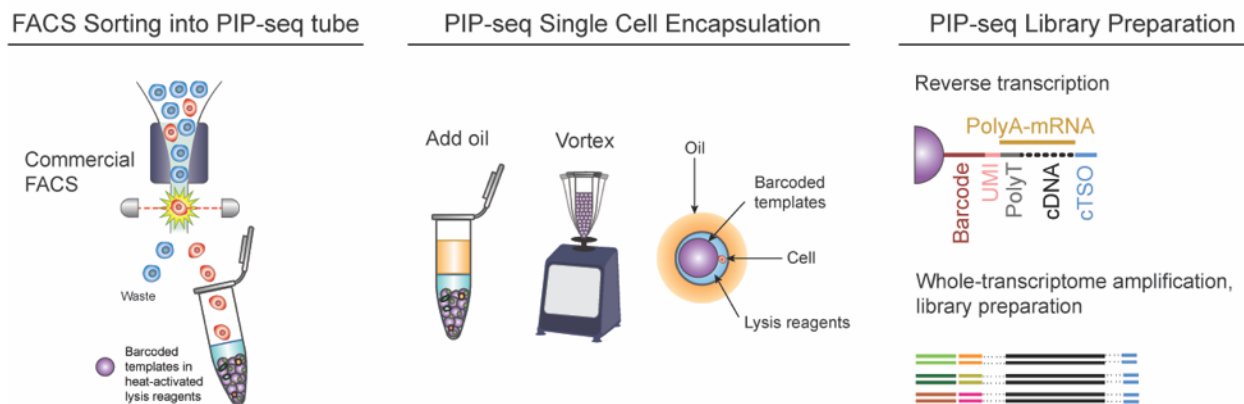
93

94

95 **Results**

96

97 **PURE-seq: Direct FACS sorting of target cells into PIP-seq single-cell RNAseq**

98 **reactions**

99

100 The PURE-seq workflow utilizes readily available commercial platforms, FACS and PIP-seq, to

101 achieve scalable, reliable, and accessible sequencing of extremely rare cells. In PURE-seq, cells

102 are sorted directly into single-cell barcoding reaction tubes. Subsequent cell encapsulation follows

103 the standard PIP-seq protocol[12], which involves adding encapsulation oil, vortexing for one

104 minute, lysing cells, and capturing mRNA (**Figure 1**). To optimize cell viability and capture

105 efficiency, we fine-tuned cell sort stream alignment, sorting speed, and total sorting duration

106 (**Methods**).

107

## Figure 1



**Figure 1. Workflow of PURE-seq with enriched and sorted rare cells from a heterogeneous population.** PURE-seq utilizes a commercial FACS system to sort fluorescently labeled target cells directly into PIP-seq reaction tubes containing barcoded templates in heat-activated lysis reagents. The subsequent single-cell encapsulation in droplets follows the standard PIP-seq protocol[12], which involves adding oil, vortexing, heat-activated lysis, and capturing mRNA on the barcoded templates. After mRNA capture, reverse transcription, and whole-transcriptome amplification are conducted in bulk to prepare barcoded cDNA for Illumina sequencing.

109

110    Fluorescence-activated cell sorters have multiple sorting precision modes. In "single-cell" mode,

111    sorting specificity is prioritized, and ambiguous results due to staining, cell clumping, or

112    coincidences in the detector are discarded. In "yield" mode, ambiguous events are recovered to

113    ensure maximum retrieval of rare cells, even at the cost of capturing some off-target cells. With

114    PURE-seq, we can prioritize capturing rare cells over the purity of the sorted population,

115    leveraging the high single-cell sequencing capacity downstream. For example, PIP-seq reactions

116    can be scaled to accommodate inputs of 2,000, 20,000, and over 100,000 cells[12]. This high capacity

117    is especially useful for sequencing extremely rare cell populations, allowing us to maximize the

118    capture of rare cells during the flow cytometry step. While the final single-cell sequenced

119    population may contain off-target cells, the overall enrichment from pre-sort to post-sort is

120    significant.

121

122    To assess the efficacy of PURE-seq, we conducted a human-mouse species-mixing experiment,

123    introducing human HEK 293T cells into mouse NIH 3T3 cells at a dilution of 1 in 1,000. The

124    human (HEK 293T) cells served as the representative target cells within a background population.

125    We labeled the human and mouse cells with different Calcein dyes (**Methods**) and processed the

126    sample using the BD FACS Aria system. We instructed the instrument to sort the first 2,500 human

127    cells into the PIP-seq reaction. In parallel, we used PIP-seq to sequence the unsorted population.

128    For the unsorted population, we recovered no human HEK 293T cells since the rarity was 1 in

129    1,000, and sequencing just 2,500 cells resulted in no random capture of human cells. By contrast,

130    in the sorted reaction, we recovered 584 human (HEK 293T) cells and 112 off-target mouse (NIH

131    3T3) cells, illustrating significant enrichment for the target population (**Figure 2A**).

132

133    To confirm successful scRNA-seq, we generated barnyard plots, plotting the number of mouse

134    reads for each cell versus the number of human reads it contains. The two populations aligned

135    along the axes, illustrating that most captured cells had either pure mouse or human transcriptomes.

136    We observed some mixed transcriptomes along the diagonal, consistent with co-encapsulation of

137    mouse and human cells during the PIP-seq barcoding step, as is typical in single-cell reactions

138    relying on limiting dilution. These results demonstrate that PURE-seq allows reliable single-cell

139    sequencing of the target cell population for the spiked population at the 1:1,000 rarity level.

140

141     A major strength of flow cytometry is its capacity for high-throughput cell sorting, allowing the

142     screening of vast populations to identify rare cellular states. In this experiment, we sought to

143     determine the maximum rarity compatible with PURE-seq. Therefore, we tuned sorting parameters

144     to maximize the total number of cells that could be sorted while minimizing the impact on the

145     cells. We set a maximum sorting duration of 60 minutes and speed of 8 kHz to prevent perturbation

146     of gene expression due to long waiting times and high shear forces in the sorter, respectively,

147     allowing 28.8 million cells to be sorted per run. At peak efficiency, this setup can, therefore,

148     recover cells with a rarity of approximately 1 in 1 million, delivering tens of target cells to the PIP-

149     seq reaction. Thus, the sequencing reaction must be exceptionally efficient to reliably barcode such

150     a tiny number of inputs; typical cell inputs for commercial single-cell instruments exceed 1,000

151     cells per reaction. Since the maximum input volume for the PIP-seq T2 kit is 5 µL, we also

152     restricted the maximum number of sorted cells to 2,500 based on the droplet volume of the BD

153     FACS Aria system (1.81 nL/drop). If more sorted cells are desired, multiple PIP-seq T2 tubes can

154     be used, or larger PIP-seq kits, such as the T20 (20,000 cells) and T100 (100,000 cells)[12], can be

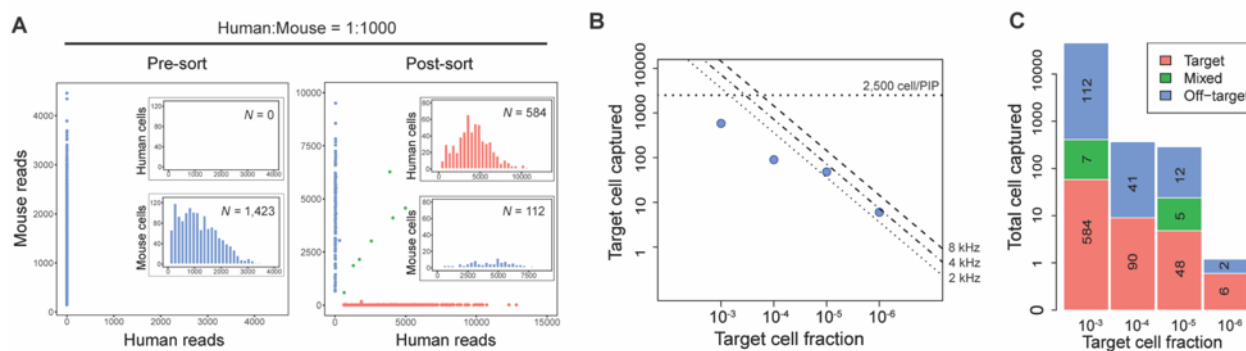155     utilized instead.

156

157     With the abovementioned parameters, we assessed the limits of enrichment possible with PURE-

158     seq by conducting sorting experiments at different target cell rarity (**Figure 2B, Supplementary**

159     **Figure 1**). We confirmed that for target cell rarity ranging from $10^{-3}$ to $10^{-6}$, between 564 and 6

160     target cells, respectively, could be captured and sequenced with 75% or greater purity (**Figure**

161     **2C**). This purity level can be increased to 98% by switching from "yield" sorting precision mode

162     to "single-cell" mode, although this reduces the number of recovered cells by ~33%

163     (**Supplementary Figure 2**).

164

165

166

**Figure 2. PURE-seq efficiently captures and sequences rare cells isolated by FACS. A)** Barnyard plots of mixed human-mouse population (Human:Mouse = 1:1000) sequenced before (left) and after sorting (right). Inserts are histograms of read distribution for sequenced human or mouse cells. Cells are colored by cell type (blue, mouse reads; red, human reads; green, mixed reads). **B)** Number of target cells captured as a function of target cell fraction. The dashed lines mark the theoretical limit of the captured cells. A maximum of 2,500 cells are sorted into each T2 PIP-seq reaction. Contour lines are the theoretical numbers of target cells that can be sorted within 60 minutes with different sorting rates (8 kHz, 4 kHz, and 2 kHz). Blue dots are the actual number of cells sequenced for the mixed human-mouse population with target cell fractions of $10^{-3}$, $10^{-4}$, $10^{-5}$, and $10^{-6}$. **C)** Number of target and off-target (mis-sort) cells sequenced for each rarity group.

## PURE-seq significantly increases the capture of LT-HSCs compared to the pre-sort control

LT-HSCs are a rare population in the mouse BM and lie at the top of the hematopoietic hierarchy[18]. Profiling LT-HSCs in scRNA-seq studies has been especially challenging due to their rarity and heterogeneity, which makes it difficult to capture enough true LT-HSCs for detailed analysis[13,14]. To demonstrate the utility of PURE-seq for the analysis of primary samples, we used it to investigate murine LT-HSCs sorted from Lineage⁻Sca-1⁺c-Kit⁺ (LSK) cells based on the expression of SLAM markers, which enrich for HSCs (CD150⁺CD48⁻ LSK cells)[19]. Specifically, to demonstrate how PURE-seq can increase the capture of LT-HSCs compared to a pre-sort control and provide a high-quality dataset to gain biological insights, we studied LT-HSCs throughout murine aging. We harvested whole BM cells from young (2-3 months old), middle-aged (12-14

8

181     months old), and old (18-20 months old) C57BL/6 mice. We removed lineage-positive cells to

182     enrich for hematopoietic stem/progenitor cells (HSPCs) before starting the PURE-seq workflow,

183     which encompassed LT-HSC sorting from BM pools (n=2-3 mice/pool) followed by the PIP-seq

184     pipeline and Illumina sequencing (**Figure 3A**). After processing and SCT-transforming the

185     samples with Seurat v4, our analysis revealed that 19.37% expressed both *Sca-1* and *c-Kit* and that

186     7.27% could be considered LT-HSCs by including the expression of *Slamf1*, which encodes for

187     the phenotypic cell surface marker CD150[20] (**Figure 3B**).

188

189     We observed that LT-HSCs did not express CD48, consistent with our FACS gating strategy,

190     which excluded CD48[+] cells (**Supplementary Figure 3A**). Our analysis also revealed that the

191     percentage of LT-HSCs increased with age (**Supplementary Figures 3A, B**), which aligns with

192     previous studies demonstrating an increase in their percentage within the aged BM[22,23]. This was

193     further confirmed by the generation of Uniform Manifold Approximation and Projection (UMAP)

194     plots that showed a higher number of hematopoietic cells expressing *Kit*, *Sca-1*, and *Slamf1* genes

195     in the middle-aged and old samples compared to their young counterparts (**Figure 3C**). *Kit*[+], *Sca-*

196     *1*[+], *Slamf1*[+] cells clustered in the head region of the UMAP plot, co-localizing with the expression

197     of key LT-HSC genes such as myeloproliferative leukemia virus oncogene (*Mpl*), endoglin (*Eng*),

198     MDS1 (*Mecom*), Meis homeobox 1 (*Meis1*), and homeobox genes (*Hoxb4* and *Hoxb5*) (**Figure**

199     **3D**).

200

201     As a control, we sequenced pre-sort samples using the PIP-seq pipeline and found that only 0.78%

202     of the cells co-expressed *Kit*, *Sca-1,* and *Slamf1*, indicating that with PURE-seq, we were able to

203     increase the percentage of LT-HSCs by 9.3-fold. Regarding the pre-sort control, we also detected

204     that even though the samples were enriched for HSPCs, there were still differentiated immune

205     cells and non-hematopoietic BM cell types, such as endothelial cells and fibroblasts

206     (**Supplementary Figure 3C**), which highlights the inefficiency of cell enrichment methods, such

207     as MACS for lineage-positive hematopoietic cell depletion (as we used in our experiment). In

208     terms of the post-sort samples, 6,725 cells that passed the Seurat quality control were captured,

209     with an average of 841 cells per sample after sorting 2,500 cells with the single-cell mode

210     (**Supplementary Figure 3D**). This demonstrates that 33.64% of the sorted cells were of high
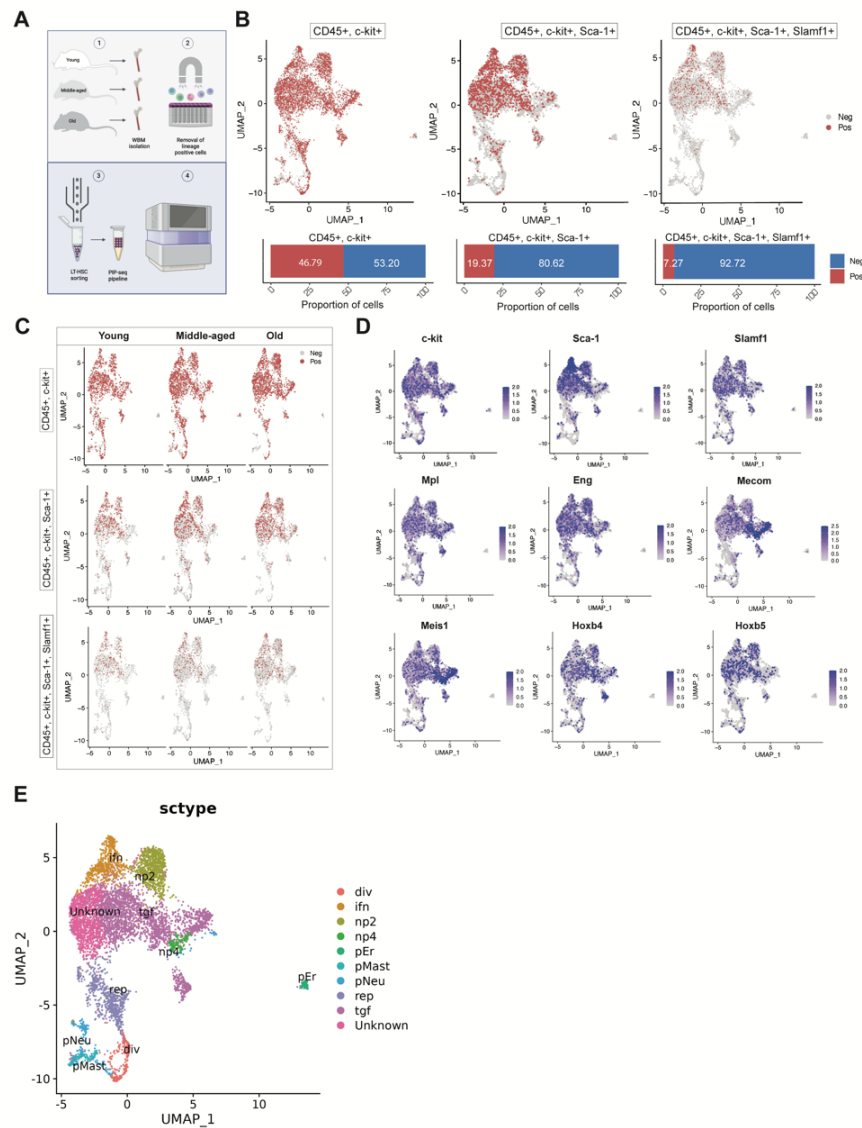
211    quality, a percentage that can be increased using the yield mode, as shown in our sorting precision

212    modes experiment (**Figure 2, Supplementary Figure 2**).

213

214    After integrating all the samples, we identified 12 clusters based on transcriptomic differences

215    (**Supplementary Figure 3D**). Next, we used a publicly available dataset from Héuralt *et al.*[21] to

216    compare their signatures with ours (**Supplementary Table 1**). Similarly, they analyzed LT-HSCs

217    from pooled FACS-sorted LT-HSC samples of old and young mice after the removal of lineage-

218    positive cells, using 10x Genomics instead. They characterized their cell clusters based on

219    differential gene expression analysis in combination with gene set enrichment analysis and gene

220    signatures related to hematopoiesis. Based on their gene markers, we were able to identify 9 out

221    of their 15 cell types, mostly coinciding with non-primed clusters, thus classified because of their

222    lack of expression of lineage-restricted genes (i.e., interferon response (ifn), non-primed (np)2,

223    growth factor signaling (tgf), np4, replicative (rep), and dividing (div)). These non-primed clusters

224    were in the head of the UMAP plot, except for an unknown cluster that did not match any of their

225    signatures, possibly due to the lack of the middle-aged group or other experimental variations in

226    their dataset. We also detected three lineage-primed clusters that were enriched for cells with

227    neutrophil (pNeu) and mastocyte (pMast) or erythroid (pEr) commitment gene markers, but these

228    were located either at the very end of the tail (pMast and pNeu) or clustered completely separately

229    from the bulk of cells (pEr) (**Figure 3E**).

230

231    Our dataset was largely comparable to datasets generated with 10X Genomics Chromium, with a

232    predominance of non-primed hematopoietic cell clusters[21]. Furthermore, the good quality metrics

233    across our 12 identified clusters (**Supplementary Figure 3F**), the clear split by biological

234    condition (i.e., age group) with concomitant detection of differences in cell numbers across clusters

235    in our integrated dataset (**Supplementary Figure 3G**), indicated the suitability of PURE-seq as a

236    reliable alternative pipeline to isolate a rare cell population and analyze their single-cell

237    transcriptomes to study their heterogeneity in complex biological phenomena such as

238    hematopoietic aging.

239

240

**Figure 3. PURE-Seq isolates murine long-term repopulating hematopoietic stem cells and enables single-cell sequencing via PIP-seq and analysis throughout aging. A)** Schematic of the PURE-seq pipeline for sorting murine LT-HSCs from young, middle-aged, and old mice after depleting lineage-positive cells for scRNA-seq library preparation using PIP-seq and Illumina sequencing. **B)** Comparison of hematopoietic cells (CD45+) expressing *c-Kit* only; *c-Kit* and *Sca-1*; or *c-Kit*, *Sca-1*, and *Slamf1*, simultaneously in the integrated UMAP plot from the dataset (top) and breakdown bar graphs of the total percentages of positive and negative cells (bottom) **C)** UMAP plots showing differences in the numbers of *c-kit* only; *c-Kit* and *Sca-1*-double positive; or *c-Kit*, *Sca-1*, and *Slamf1*-triple positive cells across murine aging. **D)** UMAP plots from the integrated dataset showing cells expressing key LT-HSC signature genes. **E)** UMAP displaying identified cell populations in the integrated dataset annotated according to Hérault *et al.* [21]

241

242

11

**Subsetting LT-HSCs from the bulk sample allows for analysis of age-related cell cycle and transcriptomic differences**

Next, we evaluated the purity of LT-HSCs in our data using the scGate package[24] (**Supplementary Table 2**). We confirmed that LT-HSCs were indeed dispersed throughout the UMAP plot, with the highest concentration in the head and middle regions of the tail (**Figure 4A**). This aligned with previous findings using Seurat (**Figures 3B, C**). Notably, the distinct cluster that stood apart did not contain any LT-HSCs. Additionally, the end of the tail of the central projection had minimal LT-HSC numbers, which was consistent with the Héuralt *et al*. integration that revealed erythroid, neutrophil, and mastocyte commitment gene expression in these clusters[21], suggesting that they likely consisted of committed progenitors or were possibly contaminated with differentiated cells.

To further validate our dataset, we set out to determine whether age-related cell cycle differences could be detected across the UMAP plot, as such changes are expected with hematopoietic aging. Using the Seurat pipeline, we found that most of the LT-HSC signature overlapped with the G1 phase classification and that the number of cells at the G1 phase appeared to increase with aging (**Figure 4B**). To further examine these differences in LT-HSCs, we extracted the pure LT-HSC population for re-embedding and re-clustering. We identified three distinct clusters (**Supplementary Figure 4A**) where nearly 100% of the cells were labeled LT-HSCs (**Figure 4C**, **Supplementary Figure 4B**). After successfully running a second post-clustering quality control check (**Supplementary Figure 4C**), we observed that G1 phase cells dominated the top clusters (clusters #0 and #1 in **Supplementary Figure 4A**), while cells at G2/M and S phases appeared to preferentially locate within the bottom cluster (cluster #2 in **Supplementary Figure 4A**) (**Figure 4D**). As we observed in the overall integrated sample before extracting the LT-HSCs subset, the proportion of LT-HSCs at the G1 phase increased at the expense of the G2/M and S phases, showing a more significant trend throughout aging compared to that of the larger dataset (**Figure 4E**). We then analyzed the gene expression signatures provided by Héuralt *et al*.[21], focusing on the LT-HSC subset. We observed that these corresponded to non-primed gene expression, specifically tgf, np1, and rep (**Figure 4F**). The rep signature, characterized by DNA repair and replication genes, had the highest number of cells at the G2/M and S phases, coinciding with cluster #2 (**Supplementary Figure 4A**). These findings support the notion that, despite an increase in their
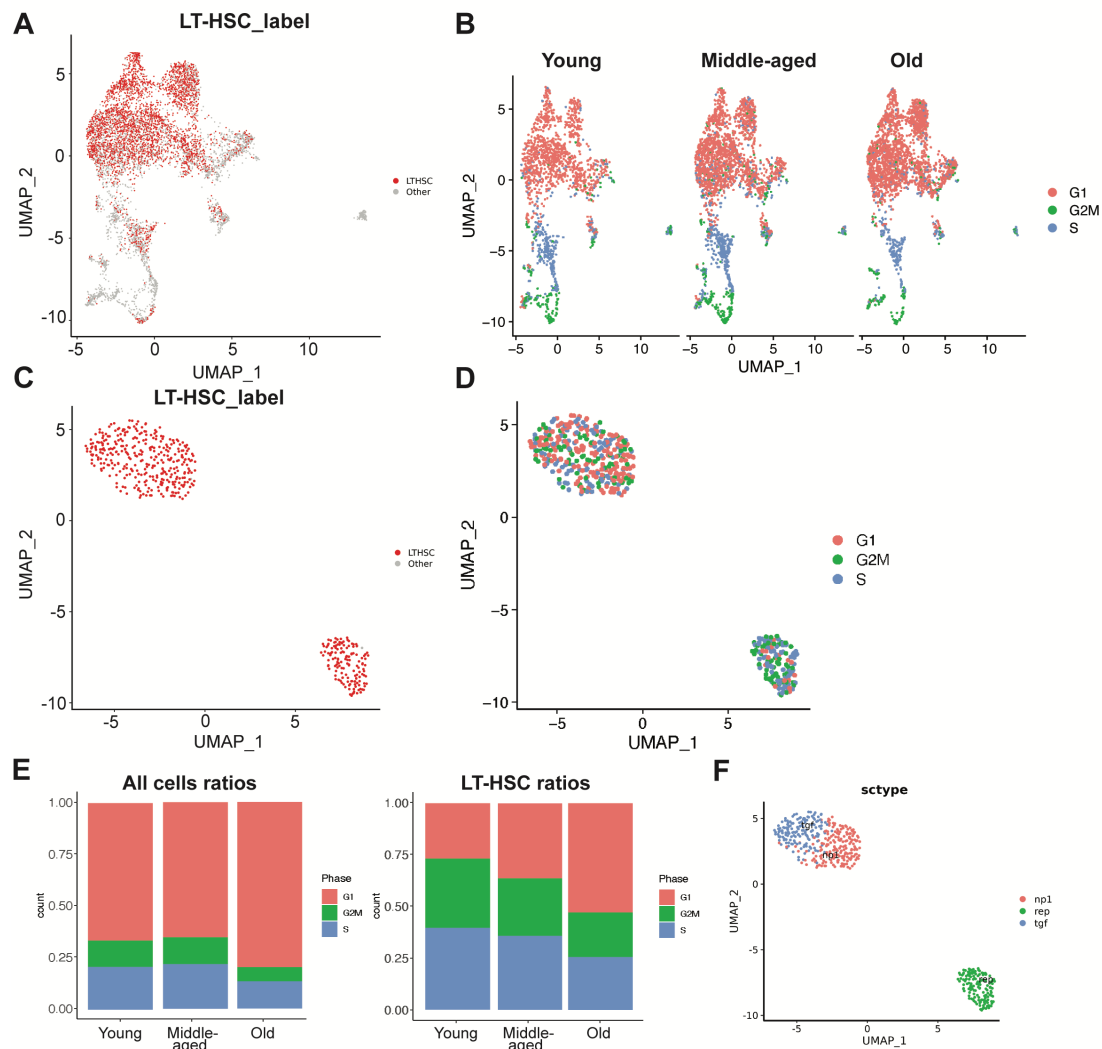
274    numbers, LT-HSCs have a gradual loss of self-renewal with aging, which has been extensively

275    reported[25].

276

277    Although refining the dataset was possible by extracting and re-clustering the LT-HSC

278    transcriptomes, the use of the overall integrated sample showed enough LT-HSC purity to conduct

279    a representative analysis, as shown by the scGate LT-HSC label (**Figure 4A**), and the expression

280    of relevant LT-HSC genes (**Figure 3D**), as well as markers of undifferentiated HSPCs (e.g., *Procr*)

281    and regeneration/myelosuppression following injury (e.g., *Notch2*), in combination with the

282    nonexistent or low expression of lineage-specific genes, such as the lymphoid-associated

283    interleukin 7 receptor (*Il7r*) and CD79A antigen (*Cd79a*), which drive differentiation towards T/B

284    lymphoid cell lineages (**Supplementary Figure 4**). Additionally, both the overall integrated

285    dataset and the LT-HSC subset allowed for the detection of age-related cell number differences

286    across all the Seurat clusters (**Supplementary Figures 4E, F**). The cross-comparison with the

287    Héuralt *et al.* dataset [21] demonstrates that the PURE-seq pipeline can obtain similar results while

288    analyzing over half the number of cells (6,725 versus 15,000 cells) while allowing for the inclusion

289    of an extra condition (the middle-aged group); this ability is especially valuable in sample scarcity

290    scenarios where cell numbers are limiting.

291

**Figure 4**



292

293

**Figure 4. scGATE marker-based purification, cell cycle analysis, and re-clustering of LT-HSCs. A)** UMAP plot indicating the purity of LT-HSCs using scGate. **B)** Analysis of cell cycle phases in the integrated UMAP plot. **C)** UMAP plot of re-clustered LT-HSCs as per the scGate label. **D)** Analysis of cell cycle phases in the re-clustered (purified) LT-HSC population. **E)** Stacked bar graphs showing the ratios of all cells (left) or LT-HSCs (right) in different phases of the cell cycle. **F)** UMAP plot of LT-HSCs labeled by cell types as annotated by Hérault *et al.* [21]

294

295

14

**Identification of EGR1 as a transcription factor determining LT-HSC gene upregulation during aging**

Aging causes genetic and epigenetic changes that lead to a decline in HSPC function and self-renewal[26]. Recent studies have identified genes that may regulate hematopoietic aging, revealing differences in gene expression and aging biomarkers, as well as an inclination towards myeloid-biased hematopoiesis as early as middle-age in mice[27,28]. In this context, single-cell transcriptomics has been useful in identifying crucial genes that could be targeted in potential hematopoietic rejuvenation strategies. To explore whether we could identify a relevant gene determining LT-HSC gene upregulation in aging from our dataset, we performed differential gene expression analysis and generated a bubble plot with top-downregulated or upregulated genes during LT-HSC aging (**Figure 5A**). Although most differences laid in the expression of genes involved in fundamental cellular processes, including DNA synthesis (e.g., *Rrm2b*), autophagy (e.g., *Vmp1*), and transcription (e.g., *Cnot6*), we observed that there was an overall elevated expression of genes regulating the immune system and inflammatory responses with aging, as previously shown[27,29]. These genes included jun B proto-oncogene (*Junb*), suppressor of cytokine signaling 3 (*Socs3*), metallothionein (*Mt1*), immediate early response 2 (*Ier2*), Krüppel-like transcription factor 4 (*Klf4*), death-associated protein kinase 1 (*Dapk1*) and genes encoding for members of the S100 protein family (e.g., *S100a6*, *S100a9*). We also found that metabolic genes showed noteworthy differences, including the upregulation of genes implicated in lipid metabolism (e.g., *Slc22a27*), glycogenesis (e.g., *Phkg1*), and growth factor signaling, such as the early growth receptor 2 (*Egr2*) and 3 (*Egr2*), and the expression of *Egr1*, Insulin growth factor 1 receptor (*Igf1r*) and transforming growth factor, beta receptor I (*Tgfbr1*); interestingly, with the latter three peaking in middle age (**Figure 5A**).

Next, we performed ShinyGO Pathway Analysis[30] to identify significantly enriched cellular pathways in aged LT-HSCs in an unbiased manner. We utilized the complete list of upregulated genes in old LT-HSCs compared to their middle-aged and young counterparts, respectively. The gene ontology category "ribosome" was the most significantly enriched gene set, which was an expected finding given the known altered upregulation in ribosomal gene transcription with hematopoietic aging, from which others have inferred that old HSPCs may be aberrantly activated

15

327 through ribosomal biogenesis despite cycling less than younger cells[32]. The rest of the enriched

328 pathways were mainly metabolism-related or linked to the pathogenesis of age-related diseases,

329 such as cardiovascular or degenerative disorders (**Figure 5B**). Using a web-based transcription

330 factor (TF) enrichment analysis tool, ChEA3[31], we identified EGR1 as the core member of the

331 most probable TF network responsible for the shift in the gene transcription profile of old LT-

332 HSCs (**Figure 5C**).

333

334 UMAP analysis revealed that although the expression of *Egr1* was not restricted to middle-aged

335 and old LT-HSCs, its expression level notably increased in middle age, as observed in the bubble

336 plot (**Figure 5A**). Furthermore, *Egr1* was widely expressed within the single LT-HSC cluster seen

337 in older mice. These UMAP plots also showed that while the young and middle-aged groups had

338 the same three clusters, the old LT-HSCs (**Figure 5D**) were absent in the bottom cluster (cluster

339 rep in **Figure 4F**, which had an enriched expression of genes involved in DNA repair and

340 replication). This might be a consequence of the age-related DNA repair defects and subsequent

341 downregulation of genes involved in such pathways or merely an observation derived from a loss

342 of heterogeneity in old LT-HSCs driven by age-related clonal hematopoiesis. Indeed, the

343 expression level of *Egr1* was found to be statistically significant when comparing young versus

344 middle-aged or young versus old LT-HSCs (**Figure 5E**). These results suggest that the

345 upregulation of *Egr1* in middle age might be responsible for a subsequent gene program

346 upregulation promoting murine LT-HSC aging, with widespread *Egr1* constitutive expression in

347 old age to maintain it.

348

349 Overall, these data demonstrate that the PURE-seq pipeline can enrich and sequence rare cell

350 populations, such as murine LT-HSCs, to generate high-quality single-cell transcriptomes and, in

351 so doing, give valuable insights into complex biological processes, as it is hematopoietic aging.

352 Compared with existing pipelines, PURE-seq offers a user-friendly solution requiring significantly

353 fewer cells while delivering comparable quality data, which is suitable for biological analyses of

354 rare cell populations.

355

**Figure 5. Identifying *Egr1* as a potential master regulator gene in the gene expression signature of aged murine long-term repopulating hematopoietic stem cells. A)** Bubble plot of the top downregulated/upregulated gene signature of old LT-HSCs compared to their young and middle-aged counterparts. The color of the spheres indicates the average gene expression, and their size represents the percentage of cells expressing each gene. **B)** The ShinyGO Pathway Analysis[30] illustrates the top enriched pathways in aged LT-HSCs compared to their young and middle-aged counterparts. The circle size represents the number of differentially expressed genes classified into one specific pathway category. **C)** Transcription factor network derived from the top upregulated genes in aged LT-HSCs based on the ChEA3 analysis [31]. **D)** UMAP plots showing *Egr1*-expressing cells in young, middle-aged, and old LT-HSC samples. **E)** Violin plots showing *Egr1* expression in young, middle-aged, and old LT-HSC samples; p-values from two-tailed unpaired Student's t-test, indicating a p-value less than 0.0001 (****) or no significance (ns).

356

357

17

## Discussion

PURE-seq enables the recovery and sequencing of rare cells from complex cellular populations by integrating two commercially available platforms: FACS and PIP-seq. PIP-seq allows cell barcoding within standard Eppendorf tubes—commonly used vessels for cell recovery in FACS protocols. This direct integration eliminates additional cell transfer steps, significantly reducing cell loss and enabling the reliable capture and sequencing of extremely rare cells.

Our study demonstrates that PURE-seq can enrich and analyze murine LT-HSCs comparably to current methods, such as 10X Genomics, even when using only half of the input cells. This approach is cost-effective, compatible with readily available commercial systems, and opens doors for proteomic analysis, including technologies like CyTOF[33] and Abseq[34], as well as multiomics through CITE-seq[35]. PURE-seq has the potential to significantly contribute to genomic and proteomic investigations, particularly those focusing on extremely rare cell populations that can be enriched using flow cytometry. Furthermore, PIP-seq can be combined with antibody-based cell hashing[12]. Although we did not perform hashing in this study, it can be used to further increase the number of cells and conditions processed in the PIP-seq pipeline. In this study, we applied PURE-seq to study hematopoietic aging in murine LT-HSCs. Our results show that LT-HSC heterogeneity is similar in young and middle age but decreases in old mice. We also found that old LT-HSCs exhibit reduced cycling and remain primarily in the G1 phase at the expense of the G2/M and S phases, as previously shown by Hérault *et al*.[21] Furthermore, our results suggest that EGR1 may be a key TF regulating LT-HSC gene expression during aging, thereby controlling the upregulation of an age-related gene program. Interestingly, *Egr1* expression increases in middle age, potentially indicating its role as an early master regulator of LT-HSC aging, further reinforcing the notion that hematopoietic aging starts in middle age[27].

While prior studies have shed some light on LT-HSCs[36,37], the role of *Egr1* in murine LT-HSC aging has not yet been fully elucidated. Recent studies involving scRNA-seq and bulk RNA sequencing have indicated increased *EGR1* expression in aged human HSPCs[15,16]. EGR1 may regulate HSPC quiescence, proliferation, and localization, making it crucial in determining their function and fate. It has been suggested that reducing EGR1 expression may decrease senescence

389    and re-activate aged HSPCs, potentially improving their function and offering a target for

390    hematopoietic rejuvenation strategies[17]. Using PURE-seq, we have identified that *Egr1* may

391    indeed be a master regulator gene of LT-HSC aging in mice, aligning with emerging research in

392    the field and providing a basis for subsequent genomic, epigenomic, and mechanistic studies.

393

394    PURE-seq offers significant potential for studying circulating tumor cells (CTCs), which are

395    valuable for research and diagnostics but challenging to capture due to their rarity[38–41]. While

396    positive enrichment using markers like EpCAM, HER2, and MUC1 is common[40,41], PURE-seq's

397    throughput enables negative enrichment, allowing it to capture CTCs that may not express these

398    markers. This capability can help discover novel or unexpected CTC types that current methods

399    might miss. With PURE-seq, sufficient CTCs can be captured for meaningful analysis. Using the

400    yield sorting precision mode, we can leverage high-throughput single-cell sequencing downstream

401    of FACS isolation to recover single CTC transcriptomes, even when mixed with non-CTCs.

402    Although this approach may increase false positives, scalable single-cell sequencing can still

403    identify the relevant CTCs, offering a less biased and useful method for diagnostics and monitoring

404    measurable residual disease at low levels.

405

406

## Methods:

### PURE-seq workflow

PURE-seq combines Fluorescence-activated cell sorting (FACS) and Particle-templated instant partition sequencing (PIP-seq) in an integrated workflow. For the mouse-human mixing experiments described herein, the BD FACS Aria system was used for sorting, and "Sweetspot" was turned on to ensure a stable stream during the sorting. The cooling unit was set to 4°C to keep the collection unit with PIP-seq reaction tube cold throughout the sort. A 0.5 mL tube adapter (Cole-Parmer, EW-17414-73) was inserted into the Aria 1.5 mL collection tube holder to hold the PIP-seq T2 tube. Then, we fine-tuned cell sort stream alignment by using an empty 0.5 mL Eppendorf tube to make sure the test sort droplet was located at the center of the lid when the lid was closed and at the center of the tube bottom when the lid was open. For quality control of each sorting session, we quantified the sorting recovery rate by sorting 100 Calcein labeled cells into a 0.5 mL Eppendorf tube pre-loaded with 10 µL media and counted the number of cells collected under the microscope. The recovery rate is calculated as # Target cells counted under the microscope / # Target cells reported to have been sorted by the instrument. To optimize cell viability and capture efficiency, we capped the total sorting duration to 60 minutes and the total sorted volume to 5 µL (2,500 drops with 85 µM nozzle). Based on BD FACS Aria's instrument specifications, we limited the flow rate to no more than 8 kHz to minimize shear stress during sheath flow focusing (i.e., 8,000 events per second with 85 µm nozzle). Once the sorting was complete, the PIP-seq T2 tube was unloaded to proceed to the standard PIP-seq protocol from Cell Capture and Lysis after the cell loading step to the preparation of the scRNA-seq library.

### Mouse-human mixing experiment

Human HEK 293T and mouse NIH 3T3 cells (ATCC) were cultured in Dulbecco's modified Eagle's medium (DMEM, Thermo Fisher, 11995073) supplemented with 10% fetal bovine serum (FBS; Gibco, 10082147) and 1× Antibiotic-Antimycotic (Gibco, 15240062) at 37°C and 5% $CO_2$. Cells were treated with 0.05% Trypsin-EDTA with Phenol red (Gibco, 25200114) for 3 min, quenched with growth medium, and centrifuged for 3 min at 300$g$. The supernatant was removed,

20

438  and the cells were resuspended in 1X DPBS without calcium or magnesium. Fresh-frozen human

439  peripheral blood mononuclear cells (PBMCs) were obtained from STEMCELL Technologies.

440  DMEM with 10% FBS was warmed up to 37°C, and the frozen PBMCs were thawed by adding

441  1 mL of warm media on top of the frozen cells and immediately transferring the media to a 15-mL

442  conical. This process was repeated until all PBMCs were thawed and transferred. Cells were

443  centrifuged for 3 min at 300$g$ and resuspended in 1X DPBS. For the $10^{-3}$, $10^{-4}$, and $10^{-5}$ target cell

444  fraction samples, human HEK 293T cells were the target population mixed with mouse NIH 3T3

445  cells background population. For the $10^{-6}$ target cell fraction sample, mouse NIH 3T3 cells were

446  the target population mixed with the human PBMCs background population. The target population

447  was treated with 1 μg/mL Calcein Red-Orange (Invitrogen, C34851), and the background

448  population was treated with 1 μg/mL Calcein Green (Invitrogen, C34852) for 15 min at 37°C,

449  followed by washing and dilution to the final concentration in 1× DPBS with 0.1% BSA. The

450  viability and cell concentration were evaluated by an automated cell counter (Bio-Rad, TC20) after

451  adding Trypan Blue (Gibco, 15250061). The mixed cell suspension was filtered through a 40 μm

452  cell strainer (Flowmi, BAH136800040) and processed through the PURE-seq workflow described

453  above to enrich for Calcein Red-Orange labeled cells. For this experiment, we selected the "yield"

454  sorting mode to ensure as many rare cells were sorted, set the flow rate to 8 kHz, and restricted the

455  sorting duration to 60 minutes or if the total sorted volume of 5 μL (2,500 drops with 85 μm nozzle)

456  was reached. In the sequenced libraries, cell transcriptomes were aligned to human or mouse

457  genome to quantify for PURE-seq sensitivity and specificity.

458

459  **Sorting precision modes experiment**

460

461  Calcein Red-Orange labeled human HEK 293T cells and Calcein Green labeled mouse NIH 3T3

462  cells were mixed at a ratio of 1:1000. The mixed sample volume was controlled at 1mL. Each

463  sample was processed through the PURE-seq workflow described above using "yield" or "single-

464  cell" sorting precision mode until depletion of sample.

465

466  **Experimental animals**

467

468 The study with primary mice was performed in accordance with institutional guidelines established
469 by Memorial Sloan Kettering Cancer Center under the Institutional Animal Care and Use
470 Committee-approved animal protocol (#07-10-016) and the Guide for the Care and Use of
471 Laboratory Animals (National Academy of Sciences 1996). Mice were maintained under specific
472 pathogen-free conditions in a controlled environment that maintained a 12-hour light-dark cycle,
473 and food and water were provided *ad libitum*. The following mice were used: young (2-3 months
474 old), middle-aged (12-14 months old), and old (18-20 months old) female C57BL/6 mice. Young
475 mice were purchased from the Jackson Laboratories and either used when young or aged in-house
476 until middle age. Old mice were obtained from the National Institute of Aging (NIA) and
477 acclimatized for at least 2 weeks at our facility before use. Mice were healthy, had intact immune
478 systems, and had not undergone any prior procedures before euthanasia. For each cohort, 4-6 mice
479 were used to make 2-3 pooled age-matched bone marrow (BM) samples per group prior to sorting.
480

**Mouse bone marrow harvesting and sample processing for sorting**

482

483 Mice were humanely euthanized using $CO_2$. BM cells from their limb bones were isolated and
484 resuspended in FACS buffer (PBS + 2% FBS) by centrifugation at 8,000 × g for 1 minute. After
485 removing red blood cells (RBC) with a commercial lysis buffer (BioLegend, 420302), diluted to
486 1X with distilled water, single-cell suspensions were depleted of hematopoietic cells committed to
487 a specific lineage using a Lineage Cell Depletion Kit (EasySep, StemCell Technologies, Inc.,
488 19856A), according to the manufacturer's instructions. To label LT-HSC cells, the following
489 fluorophore-conjugated antibodies were used at the indicated dilutions: CD117 (c-Kit) BV785
490 (clone 2B8, BioLegend; 1:200 dilution), Ly-6A/E (Sca-1) PE/Cy7 (clone D7, BioLegend; 1:1000
491 dilution), CD48 PerCP/Cy5.5 (clone HM48-1, BioLegend; 1:100 dilution) and CD150 (SLAM)
492 APC (clone TC15-12F12.2, BioLegend; 1:50 dilution). After adding the rat serum and isolation
493 cocktail of the Lineage Cell Depletion Kit, the LT-HSC-labeling antibodies were also added for a
494 30-minute-long incubation in the dark at 4°C. Following the removal of lineage-positive cells,
495 samples were spun down in FACS buffer and subsequently resuspended in 200-300 μL of FACS
496 buffer containing DAPI at a final concentration of 1 μg/mL. Cells from 2/3 age-matched mice were
497 combined to generate each pool sample, with a total of 2 replicates for the young condition and 3
498 replicates for the middle-aged and old conditions, respectively (total n=10 mice). Before sorting,

499    we also performed the Rmax method to calculate the maximum recovery of the sample sort and a

500    sorting test with horseradish peroxidase (HRP) using a 0.5 mL collection tube containing a drop

501    of a 3,3',5,5'-tetramethylbenzidine (TMB), which turned blue if the HRP fell directly into the tube

502    center. Leveraging this HRP-TMB reaction, we ensured that the instrument alignment was correct

503    so that the sample was sorted straight into the PIP-seq T2 reaction. All the mouse primary samples

504    were sorted using a Spectrally Enabled (SE) five-laser BD FACSymphony™ S6, following the

505    protocol described in the "Pure-seq workflow" section and using the "single-cell" sorting precision

506    mode to maximize the purity level.

507

## scRNA-seq library preparation and sequencing

509    Single cells were processed for scRNA-seq using the PIP-seq T2 3' Single Cell RNA kit (v3.0)

510    according to the manufacturer's protocol (Fluent Biosciences, FB0001026). cDNA and final

511    library DNA quality were confirmed using a 2100 Bioanalyzer Instrument (Agilent Technologies).

512    Libraries were pooled at equimolar ratios and sequenced on an Illumina NovaSeq 6000 S4

513    platform at PE100 (200 cycles), targeting >50,000 reads per cell. Library demultiplexing, read

514    alignment, identification of empty droplets, and UMI quantification were performed with

515    PIPseeker 1.0.0 (Fluent BioSciences) with default parameters.

516

## scRNA-seq data analysis in mice

518    Filtered feature matrices were imported into Seurat, and all downstream analyses were performed

519    using Seurat v4.3.0[42]. For quality control, data were filtered to remove outliers in gene count, UMI

520    count, mitochondrial genes, and ribosomal genes. The 8 samples (young 1-2, middle-aged 1-3, and

521    old 1-3) were normalized by SCTranform and then integrated by Seurat integration using default

522    parameters (SelectIntegrationFeatures and FindIntegrationAnchors), succeeded by normalization

523    and scaling steps[42]. The combined post-sort dataset contained 6,725 cells (Figure), while the pre-

524    sort sample had 40,137 cells. On the complete data, a PCA was estimated, and clustering was

525    performed on 20 principal component dimensions (selected by visual analysis of an Elbowplot)

526    with a resolution of 0.9. A uniform manifold approximation and projection (UMAP) embedding

527    was calculated using the selected 20 principal components as input. Cell cycle was not regressed.

23

528 As LT-HSCs were of interest in this study, hematopoietic cells co-expressing the developmental
529 markers *c-Kit*, *Ly6a*, and *Slamf1* were extracted, re-embedded, and re-clustered, followed by a
530 second post-clustering quality control step for further in-depth analysis. From the identified
531 clusters, differential gene expression analysis was conducted using the Seurat function
532 FindAllMarkers to identify genes that were significantly up/downregulated in specific cell clusters
533 compared to others.

534 After Seurat integration and clustering, different cell types were annotated using the ScType
535 automated cell type classification[43] with custom markers from the previously published dataset
536 generated by Héuralt *et al.*[21], where they identified a total of 15 subtypes of LT-HSCs, including
537 6 primed types (pMast, pNeu, pEr, pL2, pL1, pMk) and 9 non-primed types (div, rep, diff, np4,
538 np3, ifn, np2, np1, tgf). We input the gene markers of these 15 subtypes as a custom marker set to
539 score the cluster markers in our dataset using the ScType R package. Low ScType score clusters
540 (i.e., less than a quarter of the number of cells in a cluster) were considered low-confident and thus
541 designated as "unknown" cell types.

542 The purity of LT-HSCs in the data was evaluated using the scGate R package[24]. We manually
543 defined a gating model based on the LT-HSC features (*Ptprc* (*CD45*)$^+$, c-*Kit*$^+$, *Ly6a*$^+$, *Slamf1*$^+$).
544 The model annotated cells as either "pure" or "impure" based on each cell gene expression. No
545 mouse sample was excluded from these scRNA-seq analyses.

## Data availability

546
547
548 Sequencing data were deposited into the NCBI Gene Expression Omnibus under GSE273803.
549

## Code availability

550
551
552 The open-source software, tools, and packages used for data analysis in this study, as well as the
553 version of each program, were R (v3.6.1), PIPseeker (v1.0.0), Seurat R package (v4.3.0), scGate
554 R package (v1.6), ScType R package (v1.0), SingleR R package (v1.0). No custom software,
555 tools, or packages were used.

## Acknowledgments

## Contributions

S.P, K.C, A.R.A designed the study; S.P and K.C optimized the PURE-seq workflow; I.F.-M. designed and performed the experiments for all mouse studies; I.F.-M, K.C., S.P analyzed scRNA-seq data; S.V.H. provided bioinformatic and data curation support; M.G.W. assisted with mouse dissections and sample processing; R.L.B. provided input on data visualization; S.P., I.F.-M., A.R.A wrote the manuscript; R.L.L. revised the manuscript; all authors read, reviewed, and approved the manuscript.

## References:

1. Aldridge, S. & Teichmann, S. A. Single cell transcriptomics comes of age. *Nat Commun* **11**, 4307 (2020).

2. Jovic, D. *et al.* Single-cell RNA sequencing technologies and applications: A brief overview. *Clin Transl Med* **12**, (2022).

3. Mathys, H. *et al.* Single-cell atlas reveals correlates of high cognitive function, dementia, and resilience to Alzheimer's disease pathology. *Cell* **186**, 4365-4385.e27 (2023).

4. Allen, W. E., Blosser, T. R., Sullivan, Z. A., Dulac, C. & Zhuang, X. Molecular and spatial signatures of mouse brain aging at single-cell resolution. *Cell* **186**, 194-208.e18 (2023).

5. Conte, M. I., Fuentes-Trillo, A. & Domínguez Conde, C. Opportunities and tradeoffs in single-cell transcriptomic technologies. *Trends in Genetics* **40**, 83–93 (2024).

6. Ding, J. *et al.* Systematic comparison of single-cell and single-nucleus RNA-sequencing methods. *Nat Biotechnol* **38**, 737–746 (2020).

7. Cohn, L. B. *et al.* Clonal CD4+ T cells in the HIV-1 latent reservoir display a distinct gene profile upon reactivation. *Nat Med* **24**, 604–609 (2018).

8. Pauken, K. E. *et al.* Single-cell analyses identify circulating anti-tumor CD8 T cells and markers for their enrichment. *Journal of Experimental Medicine* **218**, (2021).

9. Lindner, S. *et al.* Altered microbial bile acid metabolism exacerbates T cell-driven inflammation during graft-versus-host disease. *Nat Microbiol* **9**, 614–630 (2024).

10. Hagemann-Jensen, M. *et al.* Single-cell RNA counting at allele and isoform resolution using Smart-seq3. *Nat Biotechnol* **38**, 708–714 (2020).

11. Cheng, Y.-H. *et al.* Hydro-Seq enables contamination-free high-throughput single-cell RNA-sequencing for circulating tumor cells. *Nat Commun* **10**, 2163 (2019).

12. Clark, I. C. *et al.* Microfluidics-free single-cell genomics with templated emulsification. *Nat Biotechnol* **41**, 1557–1566 (2023).

13. Zhang, P. *et al.* Single-cell RNA sequencing to track novel perspectives in HSC heterogeneity. *Stem Cell Res Ther* **13**, 39 (2022).

14. Hérault, L., Poplineau, M., Remy, E. & Duprez, E. Single Cell Transcriptomics to Understand HSC Heterogeneity and Its Evolution upon Aging. *Cells* **11**, 3125 (2022).

15. Desterke, C., Bennaceur-Griscelli, A. & Turhan, A. G. EGR1 dysregulation defines an inflammatory and leukemic program in cell trajectory of human-aged hematopoietic stem cells (HSC). *Stem Cell Res Ther* **12**, 419 (2021).

16. Adelman, E. R. *et al.* Aging Human Hematopoietic Stem Cells Manifest Profound Epigenetic Reprogramming of Enhancers That May Predispose to Leukemia. *Cancer Discov* **9**, 1080–1101 (2019).

17. Kulkarni, R. Early Growth Response Factor 1 in Aging Hematopoietic Stem Cells and Leukemia. *Front Cell Dev Biol* **10**, (2022).

18. Rossi, L. *et al.* Less Is More: Unveiling the Functional Core of Hematopoietic Stem Cells through Knockout Mice. *Cell Stem Cell* **11**, 302–317 (2012).

19. Challen, G. A., Pietras, E. M., Wallscheid, N. C. & Signer, R. A. J. Simplified murine multipotent progenitor isolation scheme: Establishing a consensus approach for multipotent progenitor identification. *Exp Hematol* **104**, 55–63 (2021).

20. Gordiienko, I., Shlapatska, L., Kovalevska, L. & Sidorenko, S. P. SLAMF1/CD150 in hematologic malignancies: Silent marker or active player? *Clinical Immunology* **204**, 14–22 (2019).

21. Hérault, L. *et al.* Single-cell RNA-seq reveals a concomitant delay in differentiation and cell cycle of aged hematopoietic stem cells. *BMC Biol* **19**, 19 (2021).

22. Kim, K. M. *et al.* Taz protects hematopoietic stem cells from an aging-dependent decrease in PU.1 activity. *Nat Commun* **13**, 5187 (2022).

23. de Haan, G., Nijhof, W. & Van Zant, G. Mouse strain-dependent changes in frequency and proliferation of hematopoietic stem cells during aging: correlation between lifespan and cycling activity. *Blood* **89**, 1543–50 (1997).

24. Andreatta, M., Berenstein, A. J. & Carmona, S. J. scGate: marker-based purification of cell types from heterogeneous single-cell RNA-seq datasets. *Bioinformatics* **38**, 2642–2644 (2022).

25. de Haan, G. & Lazare, S. S. Aging of hematopoietic stem cells. *Blood* **131**, 479–487 (2018).

26. Zhang, L., Mack, R., Breslin, P. & Zhang, J. Molecular and cellular mechanisms of aging in hematopoietic stem cells and their niches. *J Hematol Oncol* **13**, 157 (2020).
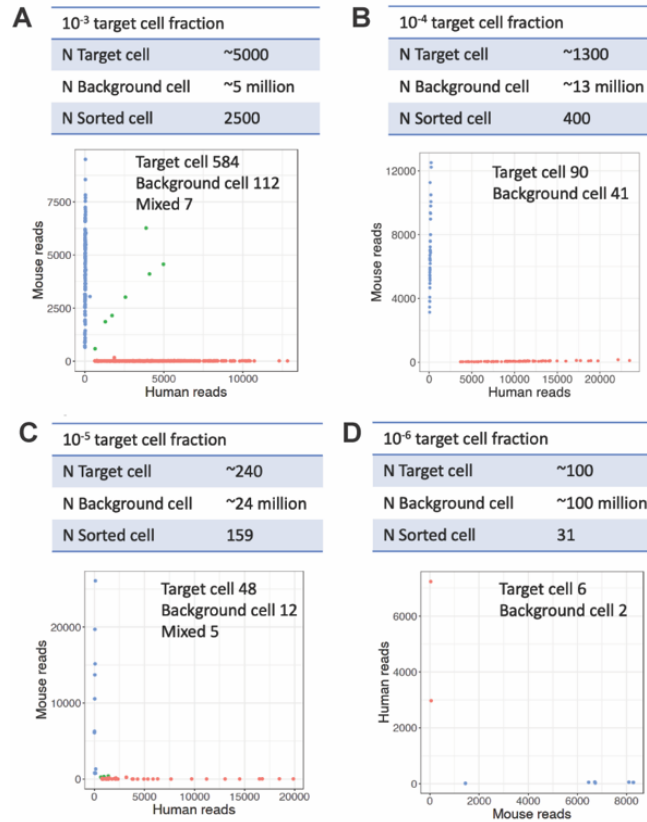
644   27.   Young, K. *et al.* Decline in IGF1 in the bone marrow microenvironment initiates

645         hematopoietic stem cell aging. *Cell Stem Cell* **28**, 1473-1482.e7 (2021).

646   28.   Bogeska, R. *et al.* Inflammatory exposure drives long-lived impairment of hematopoietic

647         stem cell self-renewal activity and accelerated aging. *Cell Stem Cell* **29**, 1273-1284.e8

648         (2022).

649   29.   Mann, M. *et al.* Heterogeneous Responses of Hematopoietic Stem Cells to Inflammatory

650         Stimuli Are Altered with Age. *Cell Rep* **25**, 2992-3005.e5 (2018).

651   30.   Ge, S. X., Jung, D. & Yao, R. ShinyGO: a graphical gene-set enrichment tool for animals

652         and plants. *Bioinformatics* **36**, 2628–2629 (2020).

653   31.   Keenan, A. B. *et al.* ChEA3: transcription factor enrichment analysis by orthogonal omics

654         integration. *Nucleic Acids Res* **47**, W212–W224 (2019).

655   32.   Sun, D. *et al.* Epigenomic Profiling of Young and Aged HSCs Reveals Concerted

656         Changes during Aging that Reinforce Self-Renewal. *Cell Stem Cell* **14**, 673–688 (2014).

657   33.   Bandura, D. R. *et al.* Mass Cytometry: Technique for Real Time Single Cell Multitarget

658         Immunoassay Based on Inductively Coupled Plasma Time-of-Flight Mass Spectrometry.

659         *Anal Chem* **81**, 6813–6822 (2009).

660   34.   Shahi, P., Kim, S. C., Haliburton, J. R., Gartner, Z. J. & Abate, A. R. Abseq: Ultrahigh-

661         throughput single cell protein profiling with droplet microfluidic barcoding. *Sci Rep* **7**,

662         44447 (2017).

663   35.   Stoeckius, M. *et al.* Simultaneous epitope and transcriptome measurement in single cells.

664         *Nat Methods* **14**, 865–868 (2017).

665   36.   Min, I. M. *et al.* The transcription factor EGR1 controls both the proliferation and

666         localization of hematopoietic stem cells. *Cell Stem Cell* **2**, 380–91 (2008).

667   37.   Stoddart, A., Fernald, A. A., Davis, E. M., McNerney, M. E. & Le Beau, M. M. EGR1

668         Haploinsufficiency Confers a Fitness Advantage to Hematopoietic Stem Cells Following

669         Chemotherapy. *Exp Hematol* **115**, 54–67 (2022).

670   38.   Habli, Z., AlChamaa, W., Saab, R., Kadara, H. & Khraiche, M. L. Circulating Tumor Cell

671         Detection Technologies and Clinical Utility: Challenges and Opportunities. *Cancers*

672         *(Basel)* **12**, 1930 (2020).

673   39.   Keller, L. & Pantel, K. Unravelling tumour heterogeneity by single-cell profiling of

674         circulating tumour cells. *Nat Rev Cancer* **19**, 553–567 (2019).

675   40.   Lin, D. *et al.* Circulating tumor cells: biology and clinical significance. *Signal Transduct*
676         *Target Ther* **6**, 404 (2021).
677   41.   Ju, S. *et al.* Detection of circulating tumor cells: opportunities and challenges. *Biomark*
678         *Res* **10**, 58 (2022).
679   42.   Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573-3587.e29
680         (2021).
681   43.   Ianevski, A., Giri, A. K. & Aittokallio, T. Fully-automated and ultra-fast cell-type
682         identification using specific marker combinations from single-cell transcriptomic data. *Nat*
683         *Commun* **13**, 1246 (2022).
684   44.   Aran, D. *et al.* Reference-based analysis of lung single-cell sequencing reveals a
685         transitional profibrotic macrophage. *Nat Immunol* **20**, 163–172 (2019).
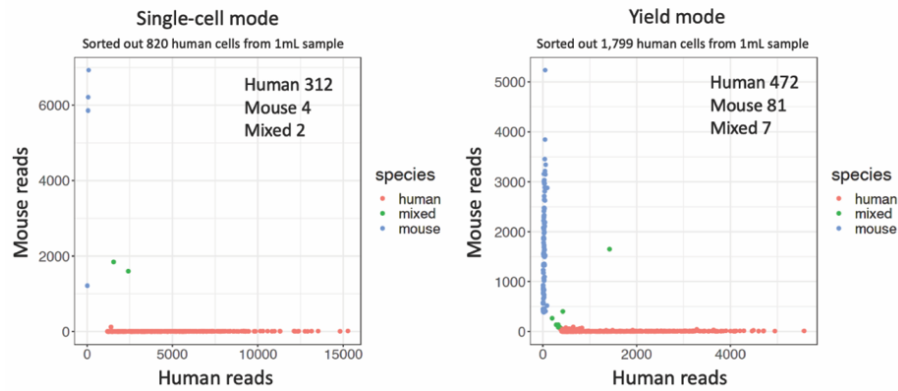686
687

688     **Supplementary Figures:**

689



690

691

**Figure S1. Barnyard plots of $10^{-3}$, $10^{-4}$, $10^{-5}$ and $10^{-6}$ target cell fractions after sorting.** In each table, cell numbers for the corresponding dilution experiment sample are shown (N Target cell and N Background cell) and the number of sorted cells reported by FACS software is noted (N Sorted cell). In each barnyard plot, cells are colored by cell type (blue, mouse reads; red, human reads; green, mixed reads). **A-C)** Human HEK 293T cells and mouse NIH 3T3 cells were stained with Calcein Red-Orange and Calcein Green, respectively. Calcein Red-Orange-positive HEK 293T cells were sorted into PIPseq tubes. **D)** Mouse NIH 3T3 cells and human PBMCs were stained with Calcein Red-Orange and Calcein Green, respectively. Calcein Red-Orange-positive NIH 3T3 cells were sorted out as target cells.
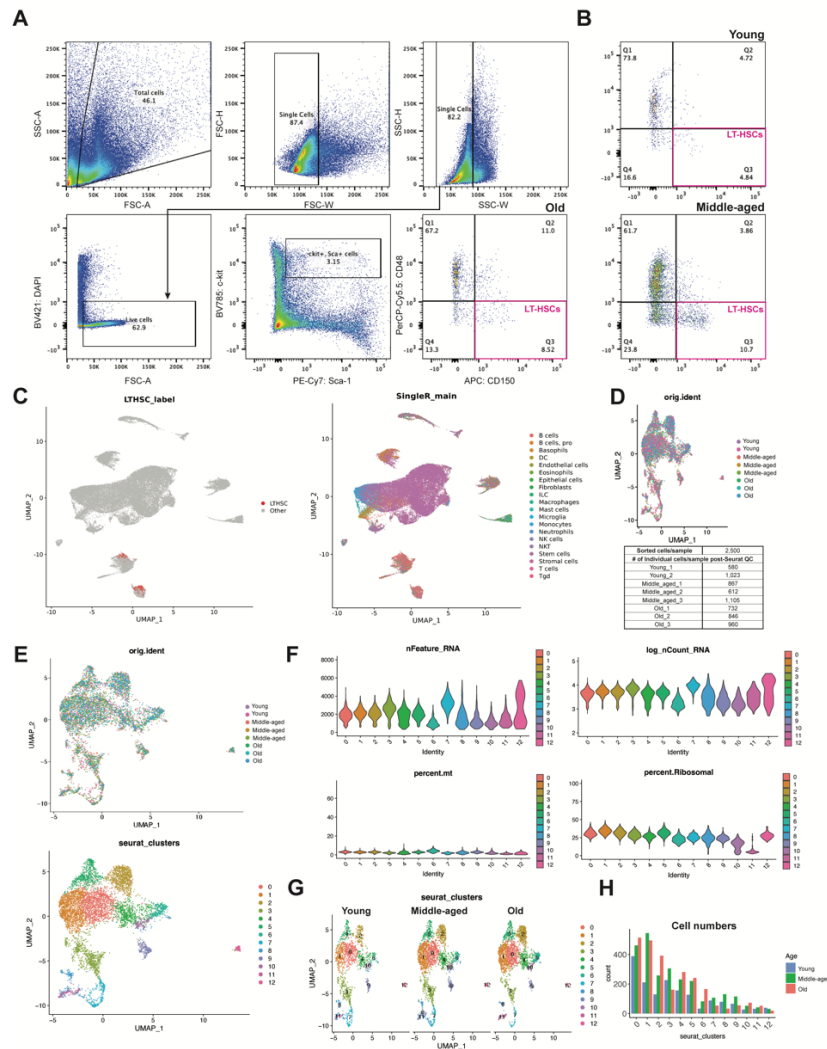
692

**Figure Supplementary 2**

693

694

**Figure S2. Recovery rate comparison of single-cell and yield sorting precision modes of FACS.** In each barnyard plot, cells are colored by cell type (blue, mouse reads; red, human reads; green, mixed reads). Target cell fraction was $10^{-3}$ and the sample volume was controlled at 1mL. Compared with single-cell mode, yield mode sorted out 2-fold the number of total cells, and sequenced 1.5-fold the number of target rare cells from identical spike-in samples. The purities of single-cell and yield modes were 98% and 84%, respectively.
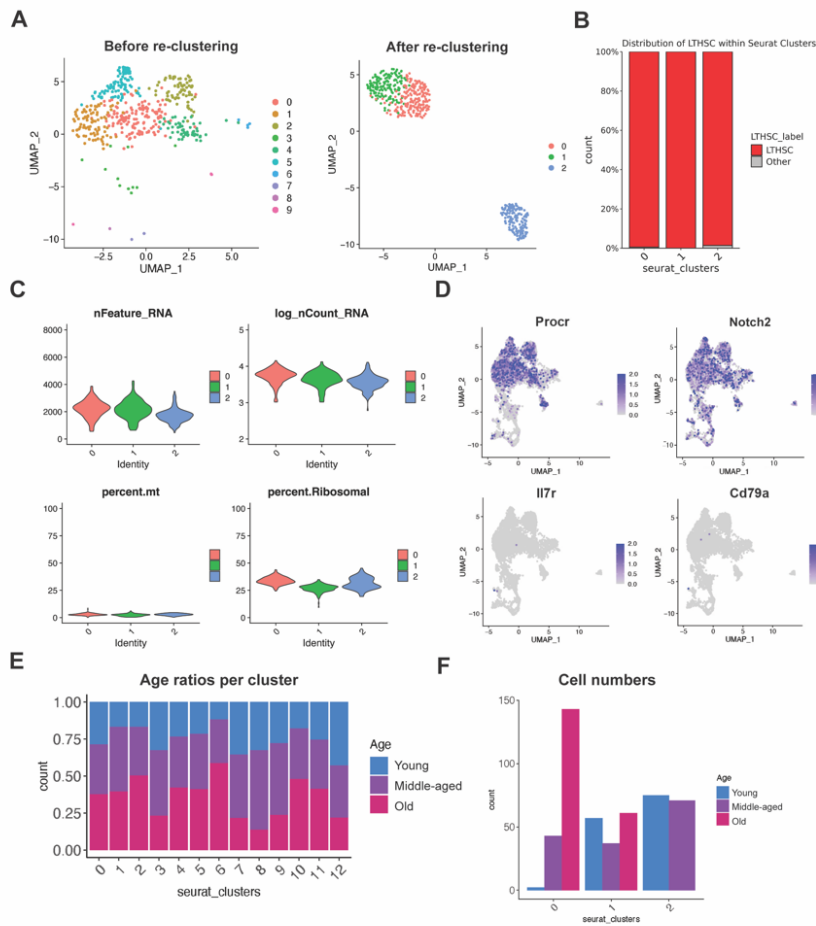
**Figure Supplementary 3**



695

**Figure S3. Sorting of murine long-term repopulating hematopoietic stem cell and quality control analysis A)** Representative FACS plots using the gating strategy to sort LT-HSCs using old cells as an example. **B)** Representative FACS plots for young (top) and middle-aged (bottom) LT-HSCs. **C)** UMAP plots of pre-sort samples, indicating LT-HSCs as labeled by scGate (left) and unbiased clustering by cell type using the SingleR package[44] (right). **D)** Integrated UMAP plot of samples from young (n=2), middle-aged (n=3), and old (n=3) mice (top) and the number of sorted cells per sample (n=2,500) and the number of cells recovered after passing quality control standards using the Seurat v4 pipeline, totaling 6,725 cells. **E)** Larger view of the integrated UMAP plot of samples from young (n=2), middle-aged (n=3), and old (n=3) samples, with each age group combining 4-6 mice. Colors indicate the age of the source mice (top) and the clustering of the 6,725 cells using the Seurat v4 pipeline (bottom). **F)** The number of unique genes (nFeature RNA), transcripts (nCount RNA as a logarithmic value), percent mitochondrial reads (percent.mt), and percent ribosomal reads (percent. Ribosomal) as a function of the cluster. **G)** Seurat clustering of young, middle-aged, and old samples. **H)** Bar graph illustrating the cell count for each age group within each Seurat cluster.

**Figure S4. Re-clustering of murine long-term repopulating hematopoietic stem cells, their distribution within Seurat clusters, and quality control post-re-clustering. A)** Before (left) and after (right) Seurat re-clustering of purified LT-HSCs according to scGate. **B)** Percentages of LT-HSCs defined by scGate within the Seurat clusters following re-clustering. **C)** The number of unique genes (nFeature RNA), transcripts (nCount RNA as a logarithmic value), percent mitochondrial reads (percent.mt), and percent ribosomal reads (percent. Ribosomal) as a function of the cluster after LT-HSC re-clustering. **D)** UMAP plots colored by expression of selected markers, including undifferentiated HSPC markers (*Procr*, *Notch2*) and markers of lineage bias/commitment (*Il7r*, *Cd79a*). **E-F)** Bar graphs illustrating the cell ratios (left) or counts (right) for each age group within each Seurat cluster subsequent to the re-clustering of LT-HSC.

696

697

33

698 **Supplementary Tables:**

699

700 **Supplementary Table 1. A**) Cluster marker gene list for integrated dataset after PURE-seq

701 enrichment of LT-HSCs from young (n=2), middle-aged (n=3), and old (n=3) mice samples. **B**)

702 Marker genes for cluster cell-type identification from the Hérault *et al.* dataset.

703 **Supplementary Table 2.** LT-HSCs identification using scGate analysis for **A**) pre-sort HSC

704 control samples, **B**) PURE-seq enriched LT-HSCs, and **C**) PURE-seq enriched LT-HSCs after

705 reclustering.