

1 **Single-cell Transcriptomics Uncovers Distinct Molecular Signatures of**
2 **Stem Cells in Chronic Myeloid Leukemia**

3
4 Giustacchini, A.,^{1,2*} Thongjuea, S.,^{1,2*} Barkas, N.,^{1,2} Woll, P.,² Povinelli, B.,^{1,2}
5 Booth, C.,^{1,2} Sopp, P.,¹ Norfo, R.,^{1,2} Rodriguez-Meira, A.,^{1,2} Ashley, N.,^{1,2}
6 Jamieson, L.,^{1,2} Vyas, P.,¹ Anderson, K.,³ Segerstolpe, Å.,^{4,5} Qian, H.,⁶
7 Olsson-Strömberg, U.,⁷ Mustjoki, S.,⁸ Sandberg, R.,^{4,9} Jacobsen,
8 S.E.W.^{1,2,4,6,10+} and Mead, A.J.^{1,2,11+}

9
10 1. MRC Molecular Hematology Unit, Weatherall Institute of Molecular
11 Medicine, University of Oxford, Oxford, United Kingdom

12 2. Haemopoietic Stem Cell Laboratory, Weatherall Institute for Molecular
13 Medicine, University of Oxford, Oxford, United Kingdom

14 3. Department of Cellular Therapy, Norwegian Radium Hospital, Oslo
15 University Hospital, Oslo, Norway

16 4. Department of Cell and Molecular Biology, Karolinska Institutet, Stockholm,
17 Sweden

18 5. Integrated Cardio Metabolic Center (ICMC), Karolinska Institutet, Huddinge,
19 Sweden

20 6. Department of Medicine, Center for Hematology and Regenerative
21 Medicine, Karolinska Institutet, Stockholm, Sweden

22 7. Department of Hematology, Uppsala University Hospital, Uppsala, Sweden

23 8. Hematology Research Unit Helsinki, Department of Clinical Chemistry and
24 Hematology, University of Helsinki and Helsinki University Hospital
25 Comprehensive Cancer Center, Helsinki, Finland

26 9. Ludwig Institute for Cancer Research, Stockholm, Sweden

27 10. Karolinska University Hospital, Stockholm, Sweden

28 11. NIHR Biomedical Research Centre, Churchill Hospital, Oxford, United
29 Kingdom

30
31 * or + These authors contributed equally to this work

32 Corresponding authors: Adam J. Mead (adam.mead@imm.ox.ac.uk) and

33 Sten Eirik W. Jacobsen (sten.eirik.jacobsen@ki.se or sten.jacobsen@imm.ox.ac.uk)

34 **Abstract**

35

36 Recent advances in single-cell transcriptomics are ideally placed to unravel
37 intratumoral heterogeneity and selective resistance of cancer stem cell (SC)
38 subpopulations to molecularly targeted cancer therapies. However, current
39 single-cell RNA-sequencing approaches lack the sensitivity required to
40 reliably detect somatic mutations. We developed a method combining high-
41 sensitivity mutation detection with whole-transcriptome analysis of the same
42 single-cell. We applied this technique to analyze over 2000 SCs from chronic
43 myeloid leukemia (CML) patients throughout the disease course, revealing
44 heterogeneity of CML-SCs, including the identification of a subgroup of CML-
45 SCs with a distinct molecular signature that selectively persisted during
46 prolonged therapy. Analysis of non-leukemic SCs from CML patients also
47 provided new insights into cell-extrinsic disruption of hematopoiesis in CML
48 associated with clinical outcome. Furthermore, we used this single-cell
49 approach to identify a blast crisis specific SC population which was also
50 present in a subclone of CML-SCs during chronic phase in a patient who
51 subsequently developed blast crisis. This approach, which might be broadly
52 applied to any malignancy, illustrates how single-cell analysis can identify
53 subpopulations of therapy-resistant SCs that are not apparent through cell-
54 population analysis.

55

56 **Introduction**

57

58 Molecularly targeted therapies for cancer frequently induce impressive
59 remissions, however, complete disease elimination remains rare, and patients
60 remain at risk of disease relapse. At a cellular level this is likely to reflect
61 intratumoral heterogeneity, with differential response to treatment in distinct
62 tumor subpopulations¹. This phenomenon relates to the proposed hierarchical
63 organization of some tumors, with only rare “cancer stem cells” (CSCs) being
64 capable of tumor propagation.²⁻⁴ There is now ample evidence for the
65 existence of such rare CSCs in some tumors, subsets of which are therapy-
66 resistant and persist during remission.²⁻⁴ However, studies characterizing
67 CSCs during remission are lacking, reflecting in part that these residual CSCs
68 are typically rare and outnumbered by their normal tissue counterparts from
69 which they cannot easily be separated.^{5,6}

70

71 Advances in single-cell gene expression techniques offer great potential to
72 study the CSC heterogeneity that might underlie therapy resistance.^{1,7-9}
73 However, to date, the application of single-cell RNA-sequencing in cancer has
74 been relatively limited in patients achieving remission following therapy,^{1,7-12}
75 partly because detection of somatic mutations is grossly underappreciated
76 using current techniques.¹¹ This primarily relates to poor coverage in the RNA
77 sequencing reads from single cells across the specific mutated region of a
78 gene due to both technical dropouts and stochastic gene expression in
79 individual cells.⁸ Thus, it is difficult to simultaneously apply single-cell
80 transcriptome analysis with highly sensitive detection of specific mutations;
81 the latter is essential to reliably distinguish normal cells from somatically
82 mutated cells which form part of the malignant clone. This is of particular
83 importance when analyzing CSC during remission, when malignant cells are
84 rare and may largely share transcriptomic features with normal tissue
85 counterparts.

86

87 Chronic Myeloid Leukemia (CML) is a paradigm for molecularly targeted
88 therapy and an ideal disease to explore the cellular basis of selective
89 resistance to targeted therapy.^{13,14} CML is less genetically complex than most
90 cancers and is defined by presence of the BCR-ABL fusion gene, the product
91 of which is the target of tyrosine kinase inhibitor (TKI) treatments which have
92 dramatically improved outcomes for this disease.¹⁵ However, chronic-phase
93 CML (CP-CML) is propagated by rare CML stem cells (CML-SCs) that are
94 selectively resistant to TKI therapy and are incompletely eradicated in most
95 patients,^{16,17} leading to frequent relapse following treatment discontinuation.¹⁸
96 CML-SCs reside in the same phenotypic compartment as their normal
97 hematopoietic stem cell (HSC) counterparts, and both express a CD34⁺CD38⁻
98 surface phenotype^{5,6}. Techniques to selectively analyze BCR-ABL⁺ SCs
99 throughout the disease course are not currently available. It therefore remains
100 to be established whether therapy-resistant CML-SCs following TKI therapy
101 represent the stochastic persistence of heterogeneous CML-SCs, a selective
102 persistence of a pre-existing distinct therapy-resistant CML-SC subset, or a
103 resistant CML-SC with novel properties that evolved as a result of the

104 therapeutic selection process.

105

106 In addition, there is ample evidence in hematological malignancies that
107 dysregulated hematopoiesis occurs as much through extrinsic disruption of
108 the normal-HSC compartment as through intrinsic expansion of the leukemic
109 clones¹⁹⁻²¹. For example, recent evidence from mouse models supports
110 involvement of non-clonal BCR-ABL⁻ SCs in the CML disease phenotype^{22,23}.
111 However, in the absence of single-cell analysis enabling separation of BCR-
112 ABL⁻ and BCR-ABL⁺ SCs within individual patients, it remains unclear to what
113 degree disruption of BCR-ABL⁻ SCs occurs in CML patients and how
114 disruption of the non-clonal SC compartment might correlate with response to
115 treatment.

116

117 Herein, we developed a new protocol integrating fluorescence activated cell
118 sorting (FACS), high sensitivity single-cell mutation detection and single-cell
119 RNA-sequencing. We apply this method to characterize distinct molecular
120 signatures of SC subpopulations in human CML samples from diagnosis
121 through remission and disease progression.

122

123 **Results**

124

125 **Combined single-cell mutation detection and transcriptomics**

126

127 Presence of the BCR-ABL fusion gene remains the only unequivocal marker
128 of CML-SCs and we therefore first sought to determine the sensitivity of BCR-
129 ABL detection using Smart-seq2, a commonly used single-cell RNA-
130 sequencing approach,^{8,24,25} by analysing the BCR-ABL⁺ K562 cell-line.²⁶
131 BCR-ABL transcripts were not detected in as many as 18/24 cells (75%; Fig.
132 1a), despite generation of satisfactory cDNA libraries as determined through
133 bioanalyser analysis of the size and concentration of amplified cDNA libraries
134 (Supplementary Fig. 1a). We obtained a similar result using a commercial
135 nanofluidic platform²⁷ (Supplementary Fig. 1b) and across a range of other
136 myeloid leukemia mutation hotspots (Supplementary Fig. 1c), validating that
137 current single-cell RNA-sequencing techniques do not enable sensitive

138 mutation detection.¹⁰

139

140 To improve sensitivity of BCR-ABL detection, we developed a BCR-ABL
141 targeted Smart-seq2 protocol (BCR-ABL tSS2) (Supplementary Fig. 2a-d). By
142 multiplexing BCR-ABL specific primers at the reverse transcription and
143 amplification steps, BCR-ABL-detection was improved to 100% of K562 cells
144 in plate-based (Fig. 1b) or microfluidic-based platforms (Supplementary Fig.
145 1d). Importantly, there was no evidence of bias caused by BCR-ABL tSS2 in
146 relation to library quality (Supplementary Fig. 2d), with good correlation
147 between level of expression of 14,240 RefSeq genes (Fig. 1c) generated by
148 Smart-seq2 or BCR-ABL tSS2; these samples also did not show separate
149 clustering (Fig. 1d; Supplementary Fig. 3). BCR-ABL plasmid "spike-in"
150 experiments demonstrated sensitivity to detect single molecules of BCR-ABL
151 with expected Poisson distribution. Importantly, no BCR-ABL amplification
152 was observed from any negative control cells in this (Fig. 1e), or in any
153 subsequent experiment (n=232 cells). This BCR-ABL tSS2 method therefore
154 allows highly specific, sensitive and quantitative BCR-ABL detection with
155 parallel unbiased whole transcriptome analysis from the same single-cell.

156

157 **Single-cell RNA-sequencing and BCR-ABL detection in CML Stem Cells**

158

159 As human HSCs are small and highly quiescent cells compared to K562 cells,
160 we first analyzed 232 Lin⁻CD34⁺CD38⁻ BM cells from five healthy human
161 donors (normal-HSCs) using BCR-ABL tSS2. Satisfactory cDNA libraries
162 were generated (Supplementary Fig. 4a), with a plateau in the numbers of
163 genes detected above 1×10^6 mapped reads/cell (Fig. 2a). With an average
164 sequencing depth of 3.4×10^6 mapped reads, a mean of 3,445 genes were
165 detected in each cell (Supplementary Fig. 4b). 12,018 genes were detected
166 (RPKM ≥ 1) in single-cell ensembles (sequencing data from all 232 cells was
167 pooled *in silico*) correlating well with cell-population data (Supplementary Fig.
168 4c) and with sensitivity to detect low level expressed transcripts
169 (Supplementary Fig. 4d), in line with previous reports.²⁸ Human HSCs
170 clustered separately and were more heterogeneous than K562 cells (Fig. 2b;
171 Supplementary Fig. 4e). Importantly, independently-processed cells from five

172 different donors clustered together, illustrating stability of the data across
173 independent experiments (Fig. 2b).

174

175 We next analyzed 40 Lin⁻CD34⁺CD38⁻ SCs from a CML patient in hematologic
176 remission following 3 months of TKI therapy (OX1407; Supplementary Table
177 1). BCR-ABL was detected in 17/40 cells (43%) by BCR-ABL tSS2 and in
178 7/20 cells (35%; P=0.8) by single-cell fluorescent-in-situ-hybridization. We
179 detected 12,499 genes in data ensembles, which correlated well with bulk
180 analysis data (Fig. 2c). Comparison of BCR-ABL⁺ and BCR-ABL⁻ SCs
181 identified genes showing differential expression (Fig. 2d-e). Level of
182 expression correlated well between single-cell RNA-sequencing and QPCR
183 data (Fig. 2f-h; Supplementary Fig. 5a-c). Together, these data provide proof
184 of principle that BCR-ABL tSS2 can be applied to detect distinct gene
185 expression in BCR-ABL⁺ versus BCR-ABL⁻ SCs from the same patient during
186 TKI treatment.

187

188 **Single-cell RNA-sequencing of CML-SCs at diagnosis**

189

190 We next used BCR-ABL tSS2 to process 2070 Lin⁻CD34⁺CD38⁻ BM SCs from
191 diagnosis samples from 20 patients with CP-CML (Supplementary Table 1).
192 Two of these CP patients developed early progression to blast crisis (BC) and
193 these cells were removed from the current analysis and analyzed in later
194 experiments (Figure 6). As previously reported²⁹, although the progenitor
195 compartment was disrupted in CML patients, the HSC-containing Lin⁻
196 CD34⁺CD38⁻ compartment was relatively intact phenotypically
197 (Supplementary Fig. 6). As expected^{5,6,30}, frequency of BCR-ABL⁺ SCs was
198 variable (median: 69%, 9%-94%; Supplementary Table 1).

199

200 We selected 854 CP-CML-SCs (477 BCR-ABL⁺ and 377 BCR-ABL⁻) for
201 sequencing and detected a mean of 3591 genes/cell (Supplementary Fig. 7a).
202 Read depth and mapped reads/cell were not different between normal-HSCs
203 (n=232), BCR-ABL⁻ SCs and BCR-ABL⁺ SCs (Supplementary Fig. 7b,c).
204 Expression of housekeeping genes e.g. *B2M* was also comparable in the
205 three groups (Supplementary Fig. 7d). In contrast, whilst the mean number of

206 genes detected was comparable between normal-HSCs (n=3,445) and BCR-
207 ABL⁻ SCs (n=3,409), a significantly higher number of genes was detected in
208 BCR-ABL⁺ SCs (n=3,735, P=1.67e-06, Supplementary Fig. 7e). This
209 correlated with BCR-ABL driven proliferation, as markedly increased
210 proliferation gene expression (Fig. 3a) and reduced quiescence-associated
211 gene expression (Fig. 3b) was observed in BCR-ABL⁺ SCs in comparison with
212 normal-HSCs. In contrast, BCR-ABL⁻ SCs showed similar proliferation
213 (Supplementary Fig. 7f) and quiescence-associated gene expression
214 (Supplementary Fig. 7g) as normal-HSCs. Consequently, co-expression of
215 G2M-associated genes was selectively increased in BCR-ABL⁺ SCs
216 (Supplementary Fig. 7h).

217

218 t-Distributed Stochastic Neighbor Embedding (tSNE) analysis using 8,589
219 highly-variable genes revealed distinct clustering of normal-HSCs, BCR-ABL⁺
220 and BCR-ABL⁻ SCs (Fig. 3c). Differentially expressed genes between normal-
221 HSCs, BCR-ABL⁺ and BCR-ABL⁻ SCs included many that were previously
222 implicated in CML pathogenesis (Supplementary Fig. 8a; Supplementary
223 Table 2) but also a number of novel candidate genes of interest such as
224 *RXFP1*, *RAB31*, *SRSF2* and *LGALS1* (Supplementary Fig. 8b;
225 Supplementary Table 2). *In silico* generation of cell-ensemble data
226 demonstrated that very few of these differentially expressed genes would
227 have been revealed without single-cell analysis (Supplementary Fig. 8c).
228 Using the top 245 differentially expressed genes, BCR-ABL⁺ cells clustered
229 separately from BCR-ABL⁻ SCs, importantly without evidence of major
230 patient-specific clustering (Fig. 3d), reflecting consistency of aberrant gene
231 expression in BCR-ABL⁺ SCs across different patients (Supplementary Fig.
232 9a,b).

233

234 Comparison of BCR-ABL⁺ SCs versus normal-HSCs and/or BCR-ABL⁻ SCs
235 showed expected enrichment in BCR-ABL⁺ SCs for the large majority of
236 established CML stem/progenitor gene-sets (Supplementary Tables 3 and 4;
237 Supplementary Fig. 10). Analysis using unbiased gene-sets (Supplementary
238 Table 5; Fig. 3e), uncovered multiple gene-sets selectively enriched in BCR-
239 ABL⁺ SCs (e.g. overexpression of *MTORC*, E2F-targets, G2M-checkpoint,

240 oxidative phosphorylation and glycolysis associated gene expression;
241 Supplementary Table 5; Fig. 3e), none of which showed enrichment through
242 *in silico* bulk analysis of the same dataset.

243

244 Importantly, our single-cell approach also uniquely allowed analysis of BCR-
245 ABL⁻ SCs within the same patients, of relevance for recent evidence that the
246 microenvironment is disrupted in CML mouse models^{22,23}. IL6-associated
247 gene expression and downstream mediators such as STAT5a were indeed
248 significantly enriched in BCR-ABL⁻ SCs in comparison with normal HSCs (Fig
249 3e, Supplementary Fig. 10 and Supplementary Tables 4 and 5). Furthermore,
250 other inflammation associated gene expression, including TGFβ and TNFα
251 pathways were also markedly enriched in BCR-ABL⁻ SCs in comparison with
252 normal HSCs (Fig. 3e). Inflammation is an important suppressor of HSC
253 function^{31,32}, including TGFβ and TNFα which are notably both cell-extrinsic
254 suppressors of HSCs^{33,34}.

255

256 **Single-cell RNA-sequencing of CML-SCs predicts TKI response**

257

258 Next, to establish the potential clinical utility of CML-SC single-cell gene
259 expression signatures, in line with current guidelines,³⁵ we stratified patients
260 with sufficient response data available as good (n=11) or poor (n=5)
261 responders on the basis of subsequent achievement of a major molecular
262 response (MMR) to TKI, defined as a BCR-ABL transcript level <0.1%
263 (Supplementary Table 1). There was no significant difference in the frequency
264 of BCR-ABL⁺ SCs between good (61%) and poor (58%) responders (P=0.7).
265 While BCR-ABL⁺ SCs at diagnosis did not clearly cluster according to
266 response category (Fig.4a), BCR-ABL⁻ SCs from poor-responder patients
267 showed highly distinct clustering using 5,611 highly-variable genes (Fig. 4b).
268 Notably, in all 5 CML patients failing to achieve MMR, the frequency of BCR-
269 ABL⁻ SCs contained within the poor-responder cluster were higher than for all
270 11 patients that achieved MMR, in 4 cases with virtually all BCR-ABL⁻ SCs
271 falling within the poor-responder cluster (Fig. 4c). The five patients with >10%
272 of BCR-ABL⁻ SCs falling within the poor-responder cluster had a markedly
273 inferior likelihood of achieving MMR (Fig. 4d, P<0.01).

274
275 GSEA also showed enrichment at diagnosis of signaling pathways,
276 inflammation, TGF β and TNF α associated gene expression in BCR-ABL⁻ SCs
277 from poor as compared to good responders (Fig. 4e and Supplementary
278 Table 6). In contrast, both BCR-ABL⁻ and BCR-ABL⁺ SCs from good
279 responders showed enrichment of MYC, E2F and G2M checkpoint gene
280 expression, associated with increased proliferation (Fig. 4e and
281 Supplementary Table 6). These data demonstrate that BCR-ABL⁺ as well as
282 BCR-ABL⁻ SCs in poor responding patients are already at diagnosis
283 expressing more quiescence associated genes than in patients who will later
284 achieve MMR; as this was observed for both BCR-ABL⁺ as well as BCR-
285 ABL⁻ SCs, this may reflect differences in cell-extrinsic, microenvironmental
286 factors in good versus poor responders.

287
288 As poor responders showed upregulation of TGF β and TNF α associated gene
289 expression, combined with a highly quiescent CML-SC signature, we
290 reasoned that TGF β and TNF α might promote quiescence in the CML-SC
291 compartment and thereby confer TKI resistance. We therefore cultured single
292 normal-HSCs and CML-SCs *in vitro* with or without TGF β or TNF α and
293 tracked the time taken for the SCs to divide. TNF α promoted quiescence of
294 both CML-SCs and normal-HSCs (Supplementary Fig.11). Notably, TGF β
295 more strongly impacted on the rate of cell division of CML-SCs compared to
296 normal-HSCs (Supplementary Fig.11). Together, these data highlight the
297 power of single-cell (unlike bulk) RNA-sequencing of CML-SCs at diagnosis to
298 reveal gene expression patterns in leukemic as well as non-leukemic SCs
299 within the same patient.

300

301 **Characterisation of quiescent CML-SCs persisting during TKI therapy**

302

303 We next analyzed 19 patients who had already commenced TKI therapy and
304 had achieved at least a hematological remission (normalization of blood
305 counts) and with additional cytogenetic response in most patients
306 (Supplementary Table 1). In 11 of these patients, paired diagnosis and follow-
307 up BM samples were available following either 3 or 6 months of TKI (Table 1).

308 In follow-up samples, the percent of BCR-ABL⁺ SCs (median: 9%, 0%-82%)
309 was lower than in diagnosis samples from the same patient (P=0.0001;
310 Supplementary Table 1). From a total of 3,306 cells processed, we selected
311 245 BCR-ABL⁺ SCs and 420 BCR-ABL⁻ SCs for single-cell sequencing.
312 Notably, unlike in diagnosis samples, the average number of genes detected
313 in each cell was similar between BCR-ABL⁺ (n=3,284) and BCR-ABL⁻ SCs
314 (n=3,196) in follow-up samples.

315

316 Using the top 500 genes informative for distinguishing normal-HSCs from
317 BCR-ABL⁺ SCs at diagnosis and during remission (Supplementary table 7),
318 tSNE analysis revealed two distinct clusters of remission BCR-ABL⁺ SCs
319 (group-A and group-B; Fig. 5a). Group-A remission BCR-ABL⁺ SCs were
320 enriched for quiescence and HSC-associated gene expression whereas
321 group-B showed enrichment of MYC, E2F and proliferation-associated gene
322 sets (Fig. 5b and Supplementary Table 8). Group-A cells were progressively
323 enriched with more prolonged TKI treatment, accounting for 43% of BCR-
324 ABL⁺ SCs at 3 months and 84% at ≥1year (P<0.01; Fig. 5c). This enrichment
325 for group-A cells was even more striking when only including patients
326 subsequently achieving MMR with 65% and 91% of BCR-ABL⁺ SCs falling
327 within group-A at 3 months and 1 year following initiation of TKI treatment
328 respectively (P<0.01; Fig. 5d; Supplementary Fig. 12a). The only exceptions
329 were one patient who temporarily interrupted TKI therapy and a patient failing
330 to achieve therapeutic imatinib levels, in both cases showing predominantly
331 group-B SCs at 3 months (Supplementary Fig. 12b). This supports the
332 concept that an excess of group-B cells during TKI therapy identifies patients
333 with inadequate BCR-ABL inhibition. We also noted that in 15 of 18 (83%)
334 diagnosis samples a minority of BCR-ABL⁺ SCs clustered within group-A
335 (26% of all diagnosis BCR-ABL⁺ SCs; Fig. 5a, c and d and Supplementary
336 Fig. 12), although the frequency of group A cells at diagnosis did not correlate
337 with response to TKI in this small cohort of patients. Together, these data
338 suggest that prolonged TKI treatment results in the selective persistence of a
339 distinct and highly quiescent BCR-ABL⁺ CML-SC subset (group-A) already
340 present at diagnosis, rather than a stochastic persistence of heterogeneous

341 CML-SCs or a resistant CML-SC with novel properties. To better understand
342 the selective persistence of quiescent CML-SCs during long-term TKI
343 treatment we therefore focused subsequent analysis on remission group-A
344 cells.

345

346 Most group-A BCR-ABL⁺ SCs clustered separately from normal-HSCs (Fig.
347 5a); we detected 1,086 differentially expressed genes in group-A remission
348 BCR-ABL⁺ SCs in comparison with normal-HSCs (Supplementary Table 9).
349 We also detected 1681 and 1348 differentially expressed genes in group-A
350 remission BCR-ABL⁺ SCs in comparison with group-B remission BCR-
351 ABL⁺ SCs and BCR-ABL⁺ SCs at diagnosis respectively (Supplementary
352 Table 9). In comparison with normal HSCs, Group-A BCR-ABL⁺ SCs showed
353 enrichment of TGFβ, TNFα via NFκB and IL6-JAK-STAT associated gene
354 expression whereas E2F, G2M checkpoint and MYC associated gene
355 expression was enriched in normal-HSCs (Fig. 5e and Supplementary Fig. 13
356 and Supplementary Table 10). Similar findings were obtained by comparing
357 group-A remission cells with BCR-ABL⁻ SCs during TKI treatment
358 (Supplementary Fig. 13). These findings support that group A remission CML-
359 SCs, characterized by marked quiescence-associated gene expression,
360 selectively evade eradication by TKI. These cells show more quiescence-
361 associated gene expression than normal HSCs or BCR-ABL⁻ SCs during
362 remission; likely because the latter are intrinsically much less sensitive to TKIs
363 due to absence of BCR-ABL expression. TGFβ and TNFα via NFκB
364 associated gene expression was progressively more enriched within
365 remission group-A BCR-ABL⁺ SCs during the course of TKI treatment (Fig.
366 5f), supporting that these pathways may be important to sustain this resistant
367 and quiescent CML-SC population during TKI treatment. Remission group-A
368 BCR-ABL⁺ SCs, also showed overexpression of Wnt/β-Catenin pathway
369 genes (*GAS2* and *CTNNB1*), the TGF-β pathway gene *SKIL*, regulators of
370 NF-κB (*NFKB1A* and *p62/SQSTM1*), the hypoxia factors *HIF1A* and the WT1
371 partner *WTAP* as well as downregulation of the chemokine receptor *CXCR4*
372 and the transcription factor *FOS* in comparison with normal HSCs (Fig. 5g and
373 Supplementary Table 9). This single-cell analysis provides insight into

374 pathways that may be involved in promoting the selective persistence of
375 distinct BCR-ABL⁺ SCs following TKI treatment.

376

377 **Analysis of CML-SC heterogeneity during blast crisis**

378

379 We next analyzed three patients with lymphoid (n=2) or myeloid (n=1) BC
380 transformation of CML (Supplementary Table 1), to explore the possibility that
381 single-cell sequencing of BCR-ABL⁺ CML-SCs could already in CP predict a
382 subsequent BC transformation. At the time of BC, tSNE analysis of CML-SCs
383 revealed a separate cluster of BCR-ABL⁺ SCs, clearly distinct from both
384 normal-HSCs, BCR-ABL⁺ SCs from 18 CP-CML patients at diagnosis and
385 K562 cells (Fig. 6a). Notably, myeloid and lymphoid blast crisis BCR-
386 ABL⁺ SCs clustered together. Comparison of gene expression of the BC and
387 CP BCR-ABL⁺ SC clusters revealed 1,166 differentially expressed genes (Fig.
388 6b and Supplementary Table 11), including overexpression of HGF³⁶ and
389 reduced expression of the Wnt pathway negative-regulator EAF2³⁷.

390

391 Two of the patients who developed BC following TKI initiation also had
392 samples available from diagnosis, when the patients presented in CP (pre-
393 BC) 12 months and 3 months before transformation to myeloid and lymphoid
394 BC, respectively (Supplementary Table 1). All pre-BC cells from the patient
395 transforming to myeloid-BC 12 months later clustered with other CP-CML SCs
396 (CP-CML cluster, Supplementary Fig. 14a). However, the pre-BC SCs from
397 the patient who 3 months later developed lymphoid BC fell into 2 distinct
398 groups, one clustering close to the BC-SCs (BC cluster, n= 124), but notably
399 with a minority (n=8) clustering separately from the BC cluster within the CP-
400 CML cluster (Fig. 6c), providing direct evidence of evolution from CP to BC
401 within the SC compartment of this patient, before any clinical or morphological
402 evidence of development of BC. In further support of this, the pre-BC single
403 BCR-ABL⁺ SCs cells falling within the BC cluster showed aberrant co-
404 expression of myeloid and lymphoid genes in comparison with normal HSCs
405 or CP-CML-SCs, as did cells within the BC cluster from all the 3 investigated
406 BC patients, whereas none of the pre-BC BCR-ABL⁺ SCs cells clustering with
407 the CP CML-SCs showed this aberrant co-expression pattern (Fig. 6d) with

408 validation of a number of aberrant expressed genes by single-cell QPCR
409 (Supplementary Fig. 14b). Moreover, index-sorting analysis (allowing specific
410 FACS data of individual cells to be linked with gene expression data from the
411 same cell) of the rare pre-BC cells in the CP-SC cluster showed that they all
412 resided within the normal $\text{Lin}^- \text{CD34}^+ \text{CD38}^- \text{CD90}^+ \text{CD45RA}^-$ HSC compartment,
413 whereas in contrast 62 of 68 of the pre-BC SCs falling in the BC-SC cluster
414 had a distinct $\text{Lin}^- \text{CD34}^+ \text{CD38}^- \text{CD90}^- \text{CD45RA}^+$ phenotype (Fig. 6e). Indeed, in
415 contrast to CP-CML patients (Supplementary Fig. 6), all BC patients analyzed
416 showed a marked expansion of $\text{Lin}^- \text{CD34}^+ \text{CD38}^- \text{CD90}^- \text{CD45RA}^+$ lymphoid-
417 primed multipotent progenitor (LMPP) like cells (Supplementary Fig. 15), a
418 population previously implicated to propagate acute leukemia³⁸.

419

420 Finally, to explore a possible genetic basis for clonal evolution within the pre-
421 BC SCs, we carried out exome sequencing of the patient with early lymphoid
422 BC, which revealed a somatic *RUNX1* mutation (c.G521A, Supplementary
423 Fig. 16a). In order to track acquisition of the *RUNX1* mutation within the BCR-
424 *ABL*⁺ SC compartment, we carried out parallel targeted amplification of both
425 BCR-ABL and the *RUNX1* mutation. All 4 pre-BC cells falling in the CP-SC
426 cluster, were *RUNX1* wild-type. In contrast, all *RUNX1* mutated pre-BC SCs
427 (n=43) were found within the BC-SC cluster (P<0.01). This distinct distribution
428 of the *RUNX1* mutation was confirmed by single-cell QPCR (Supplementary
429 Fig. 16b,c). Furthermore, differentially expressed genes between the pre-BC
430 CP-SC and BC-SC clusters were typically *RUNX1* target genes (Fig. 6f).
431 These findings are consistent with acquisition of a *RUNX1* mutation as a key
432 genomic event occurring during pre-BC, driving subsequent BC-
433 transformation at least in this one patient, with expansion of lympho-myeloid
434 transcriptionally primed LMPP-like SCs preceding the clinical BC. This is also
435 consistent with our observed expansion of this distinct SC population in both
436 myeloid and lymphoid BC patients (Supplementary Fig. 15). These data
437 illustrate how integrated single-cell gene expression, mutational profiling and
438 index sorting can be used to unravel CSC heterogeneity and reveal insights
439 that may help predict and understand subsequent disease progression.

440

441 **Discussion**

442

443 Single-cell gene expression approaches offer great promise to explore the
444 cellular heterogeneity that might underlie therapy resistance and disease
445 progression in cancer,^{1-3,8,10,11,16,17}, not-the-least in rare CSC populations, of
446 crucial importance as therapeutic elimination of all CSCs is not only required
447 but might also be sufficient to cure cancers³. However, lack of coverage in the
448 RNA-sequencing data has precluded parallel mutation analysis^{10,11},
449 representing a major limitation with current techniques.

450

451 We used CML as the disease model for single-cell CSC analysis as the
452 identity of the CSC-compartment is well-established³⁹, and rare CML-SCs
453 persisting during therapy remain a key challenge¹⁶. Although certain cell
454 surface markers have been proposed to allow for selective enrichment of
455 CML-SCs⁴⁰⁻⁴², they are not reproducible across all patients nor do they allow
456 effective purification of BCR-ABL⁺ CML-SCs during remission. In reality, the
457 presence of BCR-ABL remains the only unequivocal marker of CML-SC.
458 Therefore, we herein established a method for single-cell RNA-sequencing
459 with markedly improved sensitivity for BCR-ABL detection compared to
460 standard techniques. This new technique uniquely allowed us to selectively
461 analyse aberrant gene expression in BCR-ABL⁻ SCs at diagnosis, of
462 relevance in view of recent findings that cell-extrinsic factors disrupt normal
463 stem/progenitor cells in CML mouse models, and other hematological
464 malignancies²¹⁻²³. Our analysis revealed marked dysregulation of TGFβ and
465 TNFα pathways in BCR-ABL⁻ (as well as BCR-ABL⁺) SCs, associated with
466 increased SC quiescence. Moreover, we uncovered heterogeneity of BCR-
467 ABL⁻ SCs in CML patients, with a distinct cluster of BCR-ABL⁻ SCs
468 dominating already at diagnosis in patients who later failed to achieve MMR
469 on TKI treatment. Indeed, elevated serum levels of TNFα and TGFα also
470 correlates with poor treatment-response in CML.⁴³ Further validation studies
471 in larger patient cohorts will be required to determine whether gene
472 expression signatures of BCR-ABL⁻ SCs might have utility as a clinically
473 predictive biomarker. Furthermore, targeting inflammatory pathways such as
474 TGFβ and TNFα might also be of therapeutic value, by reducing
475 microenvironment-induced quiescence of CML-SCs, although further

476 preclinical evidence of the feasibility of such an approach is needed before
477 this could be taken forward into a clinical trial setting.

478

479 Our single-cell method also provided a unique opportunity to assess rare
480 BCR-ABL⁺ SCs persisting during TKI-induced remission^{16,17}. It was not
481 possible to analyse resistant CML-SCs in patients who had already achieved
482 deep molecular remissions due to very low frequency of BCR-ABL⁺ SCs in
483 these patients. However, analysis of samples from patients established on
484 TKI, including serial samples and patients on long term TKI (>1 year),
485 identified a distinct subpopulation of highly quiescent BCR-ABL⁺ SCs, already
486 present at diagnosis, that is markedly selected for during otherwise clinically
487 effective TKI treatment. Quiescence is a hallmark of many normal SCs,
488 including HSCs, conferring selective resistance to therapeutic targeting^{44,45}.
489 Crucially, our data using a whole-transcriptome approach, support that TKI-
490 resistant CML-SCs are transcriptionally distinct from quiescent normal-HSCs,
491 with dysregulation of specific genes and pathways (TGFβ, TNFα, JAK/STAT,
492 *CTNNB1*, *NFKB1A*) that might be selectively targeted in CML-SCs. Another
493 recent study applied a single-cell targeted gene expression analysis of BCR-
494 ABL⁺ CML-SCs⁴⁶, rather than unbiased single-cell global RNA-sequencing.
495 While the much more restricted gene expression analysis was focused at
496 looking at the heterogeneity of lineage programs in BCR-ABL⁺ stem cells and
497 improved strategies for prospective purification of CML-SCs, also the findings
498 in those studies supported a TKI-induced enrichment of quiescent BCR-ABL⁺
499 stem cells, although only investigated following short-term TKI treatment.

500

501 CML is an ideal tractable disease model to apply this single-cell technique
502 due to its relative genomic simplicity¹⁵, however, a number of our findings may
503 also be more generally applicable to other malignant disease. For example,
504 although limited by relatively small numbers of blast crisis patient samples
505 available, our analysis of BC-CML patients support that a single-cell approach
506 may prove powerful towards predicting imminent disease progression in CSC
507 populations. Specifically, our ability to detect RUNX1 mutations in distinct BC
508 CML-SCs subclones shows how a single cell approach can help to unravel
509 the mechanisms underlying clonal progression associated with certain

510 mutations at the CSC level. However, further work is required to determine
511 feasibility of applying this new method to detect a range of other mutations
512 and a number of possible limitations in the approach need to be considered:
513 Some tumors are characterized by exceedingly complex clonal heterogeneity.
514 It is likely that there will be a limitation in relation to the number of mutations
515 that could be simultaneously detected by targeted amplification in individual
516 cells before this impacts on the complexity of the RNA-seq library generated,
517 although this remains to be determined. Our technique also relies on
518 expression of the mutation of interest, and with increasing interest in
519 mutations in the non-coding space⁴⁷, further modifications to this approach
520 will be required, for example to allow parallel gDNA analysis. Furthermore, in
521 order to obtain a high level of sensitivity for BCR-ABL detection, the amplicon
522 size used in this study was short and did not encompass the kinase domain of
523 ABL. Longer BCR-ABL amplicons were less efficient at BCR-ABL detection
524 (data not shown). We were therefore unable to detect presence of kinase
525 domain mutations in individual cells, of relevance for TKI resistance⁴⁸. Further
526 modification to our technique will be required to detect multiple, distantly
527 located mutations occurring in *C/S* within the same allele.

528

529 In summary, we present a novel method allowing simultaneous single-cell
530 RNA-sequencing and high-sensitivity targeted mutation detection. We
531 demonstrate how this technique can be applied to unravel heterogeneity in
532 clonal CSCs as well as in co-existing and frequently suppressed normal SCs,
533 to provide novel insights into cellular and molecular mechanisms of therapy
534 resistance and clonal evolution. In principle, this approach can be applied
535 across a broad range of clonal disorders. Although considerable technical
536 challenges remain in relation to standardization of single cell genomics
537 techniques, we anticipate that the next few years will see major inroads
538 towards clinical application of this powerful new technology.

539

540 **Accession codes**

541 Gene Expression Omnibus (GEO) accession code GSE76312.

542

543 **Data Availability Statement**

544 Single-cell RNA sequencing data are available at NCBI's Gene Expression
545 Omnibus (GEO) data repository with the accession code GSE76312.

546

547 **Acknowledgments**

548

549 This work was funded by a Medical Research Council Senior Clinical
550 Fellowship (MR/L006340/1), MRC Confidence in Concept award
551 (MC_PC_13073) and Rosetrees Trust award (M435) to A.J.M., the MRC
552 Molecular Haematology Unit core award (A.J.M. and S.E.J.;
553 MC_UU_12009/5), a MRC programme grant to SEJ (G0801073), an
554 international recruitment award from the Swedish Research Council (SEJ),
555 and grants from the Tobias Foundation (SEJ) and the Center for Innovative
556 Medicine (CIMED) at the Karolinska Institute (SEJ). This work was also
557 supported by The MRC funded Oxford Consortium for Single-cell Biology
558 (MR/M00919X/1) and the Oxford NIHR Biomedical Centre based at Oxford
559 University Hospitals NHS Trust and University of Oxford (131/030). The views
560 expressed are those of the author(s) and not necessarily those of the NHS,
561 the NIHR, the Department of Health or the NIH. The work was also supported
562 by an educational grant from Novartis. The authors acknowledge the
563 contributions of the WIMM Flow Cytometry Facility, supported by the MRC
564 HIU; MRC MHU (MC_UU_12009); NIHR Oxford BRC and John Fell Fund
565 (131/030 and 101/517), the EPA fund (CF182 and CF170) and by the WIMM
566 Strategic Alliance awards G0902418 and MC_UU_12025. N.A. was supported
567 by the Oxford-Wellcome Trust Institutional Strategic Support Fund. S.M. is
568 supported by the Finnish Cancer Institute and the Finnish Cancer
569 Organizations.

570

571 **Author Contribution**

572

573 A.G. designed, performed and analyzed experiments and contributed to
574 writing the manuscript. S.T. designed and performed bioinformatic analyses
575 and contributed to writing the manuscript. N.B. and B.P. performed analyses
576 of RNA-sequencing and QPCR results. P.W. and P.S. were involved in FACS

577 analysis/sorting. R.N., A.R.M., C.B., L.J. performed experiments. N.A.
578 maintained single-cell facility infrastructure. P.V., S.M. and H.Q. provided
579 infrastructure for sample banking and provided input on experimental design
580 and analysis. K.A. performed FISH experiments. A.S. was involved in RNA-
581 sequencing experiments. S.U. collected clinical information. R.S. provided
582 input on RNA-sequencing experiments. A.J.M. and S.E.W.J. conceived and
583 supervised the project, designed and analyzed experiments and wrote the
584 manuscript.

585

586 **Competing Financial Interests**

587

588 A.J.M. has received honoraria and research funding from Novartis.

589

590 **References**

591

- 592 1. McGranahan, N. & Swanton, C. Biological and therapeutic impact of
593 intratumor heterogeneity in cancer evolution. *Cancer cell* **27**, 15-26
594 (2015).
- 595 2. Tehranchi, R., *et al.* Persistent malignant stem cells in del(5q)
596 myelodysplasia in remission. *The New England journal of medicine* **363**,
597 1025-1037 (2010).
- 598 3. Magee, J.A., Piskounova, E. & Morrison, S.J. Cancer stem cells: impact,
599 heterogeneity, and uncertainty. *Cancer cell* **21**, 283-296 (2012).
- 600 4. Woll, P.S., *et al.* Myelodysplastic syndromes are propagated by rare and
601 distinct human cancer stem cells in vivo. *Cancer cell* **25**, 794-808 (2014).
- 602 5. Sloma, I., *et al.* Genotypic and functional diversity of phenotypically
603 defined primitive hematopoietic cells in patients with chronic myeloid
604 leukemia. *Exp Hematol* **41**, 837-847 (2013).
- 605 6. Mustjoki, S., *et al.* Impact of malignant stem cell burden on therapy
606 outcome in newly diagnosed chronic myeloid leukemia patients.
607 *Leukemia* **27**, 1520-1526 (2013).
- 608 7. Alizadeh, A.A., *et al.* Toward understanding and exploiting tumor
609 heterogeneity. *Nature medicine* **21**, 846-853 (2015).
- 610 8. Wills, Q.F. & Mead, A.J. Application of single-cell genomics in cancer:
611 promise and challenges. *Human molecular genetics* (2015).
- 612 9. Wang, Y. & Navin, N.E. Advances and applications of single-cell
613 sequencing technologies. *Molecular cell* **58**, 598-609 (2015).
- 614 10. Patel, A.P., *et al.* Single-cell RNA-seq highlights intratumoral heterogeneity
615 in primary glioblastoma. *Science* **344**, 1396-1401 (2014).
- 616 11. Miyamoto, D.T., *et al.* RNA-Seq of single prostate CTCs implicates
617 noncanonical Wnt signaling in antiandrogen resistance. *Science* **349**,
618 1351-1356 (2015).

- 619 12. Tirosh, I., *et al.* Dissecting the multicellular ecosystem of metastatic
620 melanoma by single-cell RNA-seq. *Science* **352**, 189-196 (2016).
- 621 13. Druker, B.J., *et al.* Effects of a selective inhibitor of the Abl tyrosine kinase
622 on the growth of Bcr-Abl positive cells. *Nature medicine* **2**, 561-566
623 (1996).
- 624 14. Goldman, J.M. & Melo, J.V. Targeting the BCR-ABL tyrosine kinase in
625 chronic myeloid leukemia. *The New England journal of medicine* **344**,
626 1084-1086 (2001).
- 627 15. Longo, D.L. Imatinib Changed Everything. *The New England journal of*
628 *medicine* **376**, 982-983 (2017).
- 629 16. Gallipoli, P., Abraham, S.A. & Holyoake, T.L. Hurdles toward a cure for
630 CML: the CML stem cell. *Hematology/oncology clinics of North America* **25**,
631 951-966, v (2011).
- 632 17. Chu, S., *et al.* Persistence of leukemia stem cells in chronic myelogenous
633 leukemia patients in prolonged remission with imatinib treatment. *Blood*
634 **118**, 5565-5572 (2011).
- 635 18. Mahon, F.X., *et al.* Discontinuation of imatinib in patients with chronic
636 myeloid leukaemia who have maintained complete molecular remission
637 for at least 2 years: the prospective, multicentre Stop Imatinib (STIM)
638 trial. *Lancet Oncol* **11**, 1029-1035 (2010).
- 639 19. Schepers, K., *et al.* Myeloproliferative neoplasia remodels the endosteal
640 bone marrow niche into a self-reinforcing leukemic niche. *Cell Stem Cell*
641 **13**, 285-299 (2013).
- 642 20. Colmone, A., *et al.* Leukemic cells create bone marrow niches that disrupt
643 the behavior of normal hematopoietic progenitor cells. *Science* **322**, 1861-
644 1865 (2008).
- 645 21. Schepers, K., Campbell, T.B. & Passegue, E. Normal and leukemic stem cell
646 niches: insights and therapeutic opportunities. *Cell Stem Cell* **16**, 254-267
647 (2015).
- 648 22. Welner, R.S., *et al.* Treatment of chronic myelogenous leukemia by
649 blocking cytokine alterations found in normal stem and progenitor cells.
650 *Cancer cell* **27**, 671-681 (2015).
- 651 23. Reynaud, D., *et al.* IL-6 controls leukemic multipotent progenitor cell fate
652 and contributes to chronic myelogenous leukemia development. *Cancer*
653 *cell* **20**, 661-673 (2011).
- 654 24. Picelli, S., *et al.* Full-length RNA-seq from single cells using Smart-seq2.
655 *Nature protocols* **9**, 171-181 (2014).
- 656 25. Wilson, N.K., *et al.* Combined Single-Cell Functional and Gene Expression
657 Analysis Resolves Heterogeneity within Stem Cell Populations. *Cell Stem*
658 *Cell* **16**, 712-724 (2015).
- 659 26. Wu, S.Q., *et al.* Extensive amplification of bcr/abl fusion genes clustered
660 on three marker chromosomes in human leukemic cell line K-562.
661 *Leukemia* **9**, 858-862 (1995).
- 662 27. Wu, A.R., *et al.* Quantitative assessment of single-cell RNA-sequencing
663 methods. *Nat Methods* **11**, 41-46 (2014).
- 664 28. Islam, S., *et al.* Highly multiplexed and strand-specific single-cell RNA 5'
665 end sequencing. *Nature protocols* **7**, 813-828 (2012).
- 666 29. Bruns, I., *et al.* The hematopoietic stem cell in chronic phase CML is
667 characterized by a transcriptional profile resembling normal myeloid

- 668 progenitor cells and reflecting loss of quiescence. *Leukemia* **23**, 892-899
669 (2009).
- 670 30. Miyawaki, K., *et al.* The Expansion Of CML Clones Initiates At The CMP
671 Stage, and Is Associated With The Down-Regulation Of IRF8 and GFI1.
672 *Blood* **122**(2013).
- 673 31. Schuettpelez, L.G. & Link, D.C. Regulation of hematopoietic stem cell
674 activity by inflammation. *Front Immunol* **4**, 204 (2013).
- 675 32. King, K.Y. & Goodell, M.A. Inflammatory modulation of HSCs: viewing the
676 HSC as a foundation for the immune response. *Nat Rev Immunol* **11**, 685-
677 692 (2011).
- 678 33. Pronk, C.J., Veiby, O.P., Bryder, D. & Jacobsen, S.E. Tumor necrosis factor
679 restricts hematopoietic stem cell activity in mice: involvement of two
680 distinct receptors. *The Journal of experimental medicine* **208**, 1563-1570
681 (2011).
- 682 34. Cashman, J.D., Eaves, A.C., Raines, E.W., Ross, R. & Eaves, C.J. Mechanisms
683 that regulate the cell cycle status of very primitive hematopoietic cells in
684 long-term human marrow cultures. I. Stimulatory role of a variety of
685 mesenchymal cell activators and inhibitory role of TGF-beta. *Blood* **75**,
686 96-101 (1990).
- 687 35. Baccarani, M., *et al.* European LeukemiaNet recommendations for the
688 management of chronic myeloid leukemia: 2013. *Blood* **122**, 872-884
689 (2013).
- 690 36. Kentsis, A., *et al.* Autocrine activation of the MET receptor tyrosine kinase
691 in acute myeloid leukemia. *Nature medicine* **18**, 1118-1122 (2012).
- 692 37. Liu, J.X., *et al.* Eaf1 and Eaf2 negatively regulate canonical Wnt/beta-
693 catenin signaling. *Development* **140**, 1067-1078 (2013).
- 694 38. Goardon, N., *et al.* Coexistence of LMPP-like and GMP-like leukemia stem
695 cells in acute myeloid leukemia. *Cancer cell* **19**, 138-152 (2011).
- 696 39. Jamieson, C.H. Chronic myeloid leukemia stem cells. *Hematology / the*
697 *Education Program of the American Society of Hematology. American*
698 *Society of Hematology. Education Program*, 436-442 (2008).
- 699 40. Zhao, K., *et al.* IL1RAP as a surface marker for leukemia stem cells is
700 related to clinical phase of chronic myeloid leukemia patients. *Int J Clin*
701 *Exp Med* **7**, 4787-4798 (2014).
- 702 41. Herrmann, H., *et al.* Dipeptidylpeptidase IV (CD26) defines leukemic stem
703 cells (LSC) in chronic myeloid leukemia. *Blood* **123**, 3951-3962 (2014).
- 704 42. Gerber, J.M., *et al.* Genome-wide comparison of the transcriptomes of
705 highly enriched normal and chronic myeloid leukemia stem and
706 progenitor cell populations. *Oncotarget* **4**, 715-728 (2013).
- 707 43. Nievergall, E., *et al.* TGF-alpha and IL-6 plasma levels selectively identify
708 CML patients who fail to achieve an early molecular response or progress
709 in the first year of therapy. *Leukemia* **30**, 1263-1272 (2016).
- 710 44. Trumpp, A., Essers, M. & Wilson, A. Awakening dormant haematopoietic
711 stem cells. *Nat Rev Immunol* **10**, 201-209 (2010).
- 712 45. Clevers, H. The cancer stem cell: premises, promises and challenges.
713 *Nature medicine* **17**, 313-319 (2011).
- 714 46. Warfvinge, R., *et al.* Single-cell molecular analysis defines therapy
715 response and immunophenotype of stem cell subpopulations in CML.
716 *Blood* (2017).

- 717 47. Mansour, M.R., *et al.* Oncogene regulation. An oncogenic super-enhancer
718 formed through somatic mutation of a noncoding intergenic element.
719 *Science* **346**, 1373-1377 (2014).
720 48. Soverini, S., De Benedittis, C., Mancini, M. & Martinelli, G. Present and
721 future of molecular monitoring in chronic myeloid leukaemia. *Br J*
722 *Haematol* **173**, 337-349 (2016).
723

724

725 **Figure Legends**

726

727 **Figure 1.** High sensitivity single-cell detection of BCR-ABL with parallel
728 unbiased whole transcriptome analysis. **(a, b)** Detection of BCR-ABL and
729 GAPDH by QPCR in libraries from single K562 cells processed by Smart-
730 seq2 **(a)**, or BCR-ABL tSS2 **(b)**. Values shown are the gene expression levels
731 relative to the limit of detection (LOD), indicated by the dashed horizontal line.
732 The box plot shows median and quartile values and whiskers show outlier
733 values within 1.5 interquartile range of the quartiles. Numbers below each plot
734 show the frequency of cells showing expression above the LOD. **(c)**
735 Correlation of expression data for 14,240 RefSeq genes generated from K562
736 single cells using Smart-seq2 (n=38) or BCR-ABL tSS2 (n=38). **(d)** RNA-
737 sequencing results from single K562 cells processed with Smart-seq2 (blue
738 n=38) or by BCR-ABL tSS2 (red n=38) shown by t-Distributed Stochastic
739 Neighbor Embedding (tSNE) using 3,368 highly variable genes (see
740 methods). **(e)** Dot plot illustrating the sensitivity to detect specific copy
741 numbers of BCR-ABL spiked-in before BCR-ABL tSS2 amplification of single
742 BM cells from a healthy donor. The Y-axis indicates the gene expression level
743 of BCR-ABL relative to the LOD. The X-axis indicates the absolute number of
744 copies of BCR-ABL expected to be present in each reaction, calculated using
745 a commercial standard. The table above shows the numbers of wells that
746 would be expected to contain at least one copy of BCR-ABL by Poisson
747 distribution and the actual frequency of amplification following BCR-ABL tSS2.

748

749 **Figure 2.** Single-cell whole transcriptome analysis and BCR-ABL detection in
750 single CML stem cells. **(a)** Box plot illustrating the number of genes detected
751 (RPKM ≥ 1) in relation to depth of sequencing in normal-HSC samples (shown

752 as million reads/cell). (b) RNA-sequencing results from single K562 cells
753 processed by BCR-ABL tSS2 (purple n=38) and normal-HSCs (n=232) shown
754 by t-Distributed Stochastic Neighbor Embedding (tSNE) using 7,428 highly-
755 variable genes. (c) Correlation between the merged data from 40 single-cells
756 from CML patient OX1407 (“ensemble”) and the bulk (100 cells sorted
757 together) RNA-sequencing measurement of gene expression from the same
758 patient. The ensemble was created by computationally pooling all the reads
759 obtained from the 40 single Lin⁻CD34⁺CD38⁻ cells from patient OX1407. Some
760 of the genes shown in panel f are highlighted. (d) Correlation between the
761 levels of gene expression of BCR-ABL⁺ ensemble and BCR-ABL⁻ ensemble
762 data. Some of the genes shown in panel f are highlighted. (e) Heat map
763 illustrating the hierarchical clustering of BCR-ABL⁺ SCs (red, n=17) or BCR-
764 ABL⁻ SCs (Blue, n=23) showing the top 75 differentially expressed genes. (f)
765 Correlation of log₂(FC) by RNA-sequencing (y-axis) and by QPCR (x-axis)
766 between BCR-ABL⁺ and BCR-ABL⁻ SCs for selected genes. Differentially
767 expressed genes (red dots) were selected by setting a fold change cutoff >8
768 and selected 12 genes of potential biologic interest. Non-differentially
769 expressed genes (grey dots, n=12) were selected as housekeeping genes or
770 relevant genes for the cell type analyzed. (g, h) Beeswarm plots for 6 of the
771 12 selected differentially expressed genes between BCR-ABL⁺ and BCR-
772 ABL⁻ SCs showing RNA-sequencing (g) and QPCR (h) data. Numbers of cells
773 analyzed and numbers showing amplification for the selected gene are shown
774 below the plot. Nonparametric Wilcoxon test p-values are shown on top of
775 each bar graph. Fisher’s exact test p-values are shown below the graph. The
776 average gene expression levels are indicated by red squares, the median and
777 quartiles of gene expression levels are represented by the boxes. The dashed
778 lines represent the LOD.

779

780 **Figure 3.** Single-cell RNA-sequencing reveals distinct molecular signatures of
781 BCR-ABL⁺ CML-SCs at diagnosis. (a, b) Gene-set enrichment analysis on
782 477 BCR-ABL⁺ single-cells from 18 chronic phase CML patients at diagnosis
783 versus 232 normal-HSCs from 5 normal donors. Gene-sets shown are (a) cell
784 proliferation and (b) quiescence associated genes. (c) tSNE visualization of
785 single normal-HSCs (Grey circles; n=232), BCR-ABL⁻ SCs (blue diamonds;

786 n=377) and BCR-ABL⁺ SCs (red triangles; n=477) using 8,589 highly-variable
787 genes. **(d)** Hierarchical clustering analysis of the same 1086 cells. The
788 heatmap is built using Pearson correlation generated using the top 245
789 differentially expressed genes. The horizontal color bar on top of the heatmap
790 indicates the sample from which each single SC was purified (upper bar,
791 individual color for each patient) and the cell ID (lower bar): normal-HSCs
792 (black), BCR-ABL⁻ SCs (blue) and BCR-ABL⁺ SCs (red). **(e)** GSEA of
793 unbiased HALLMARK gene-sets for 1) Normal-HSCs (n=6) vs BCR-ABL⁺ SCs
794 (n=18) as an *in silico* bulk analysis; 2) Single-cell analysis of normal-HSCs
795 (n=232) vs BCR-ABL⁻ SCs (n=377); 3) Single-cell analysis of normal-HSCs
796 (n=232) vs BCR-ABL⁺ SCs (n=477); 4) Single-cell analysis of BCR-ABL⁻ SCs
797 (n=377) vs BCR-ABL⁺ SCs (n=477). A false discovery rate (FDR) cut-off of
798 0.25 was used.

799

800 **Figure 4.** Single-cell RNA-sequencing of SCs in CML patients at diagnosis
801 predicts molecular response to TKI. **(a)** tSNE visualization of single BCR-
802 ABL⁺ SCs (from 16 CP-CML patients with molecular follow-up data available,
803 n=436) using 5,011 highly-variable genes. Color indicates if cells were
804 isolated from good responders (n=11 patients achieving MMR, blue) or poor
805 responders (n=5 patients not achieving MMR, red). **(b)** tSNE visualization of
806 single BCR-ABL⁻ SCs from 16 patients with molecular follow-up data available
807 (n=356) using 5,611 highly-variable genes. Color indicates if cells were
808 isolated from good responders (11 patients achieving MMR, blue) or poor
809 responders (5 patients not achieving MMR, red). **(c)** The dot plot shows the
810 proportion (%) of BCR-ABL⁻ SCs falling in the “poor responders” cluster for
811 individual patients (n=16, in red patients with >10% of cells in poor responder
812 cluster, in blue patients with <10% of cells in poor responder cluster; squares
813 represent patients failing to achieve MMR and circles patients who achieved
814 MMR). **(d)** Kaplan Meier curves showing time for MMR achievement for
815 patients with >10% (red, n=5) or <10% (blue, n=11;) of BCR-ABL⁻ SCs falling
816 in the poor responder cluster. P-value represents the logrank test. **(e)** GSEA
817 of unbiased HALLMARK gene sets comparing BCR-ABL⁻ SCs (n=356) and
818 BCR-ABL⁺ SCs (n=436) from good (n=11) and poor (n=5) TKI responders.

819

820 **Figure 5.** Single-cell analysis reveals distinct molecular signatures of
821 quiescent CML-SCs persisting during TKI therapy. **(a)** tSNE visualization of
822 single normal-HSCs (black circles; n=232, 5 donors), BCR-ABL⁺ SCs from
823 patients at diagnosis (grey circles; n=477, 18 donors) and BCR-ABL⁺ SCs
824 from patients at remission (light-blue diamonds and dark-blue triangles,
825 n=245, 16 donors). Remission BCR-ABL⁺ SCs clustering closer to normal-
826 HSCs (light-blue diamonds, n=122) are defined as group-A BCR-ABL⁺ SCs,
827 while those cells clustering with most diagnostic BCR-ABL⁺ SCs are defined
828 as group-B (dark-blue triangles, n=123). **(b)** Gene-set enrichment analysis of
829 group-A vs group-B BCR-ABL⁺ SCs at remission (n=122 and n=123,
830 respectively). Gene-sets shown are cell proliferation, quiescence and HSC-
831 associated genes. **(c)** Bar graph showing the proportion (%) of group-A BCR-
832 ABL⁺ SCs and group-B BCR-ABL⁺ SCs for all patients analyzed at diagnosis
833 (n=18), at 3 months (n=11) and more than 1 year (n=4) after TKI initiation.
834 Chi-square and Fisher's exact test p<0.01 for comparison of diagnosis versus
835 1-year samples. **(d)** The bar graph shows same results as in panel **c** but only
836 for patients eventually achieving MMR, samples taken at diagnosis (n=11), at
837 3 months (n=6) and more than 1 year (n=2) after TKI initiation. Chi-square
838 and Fisher's exact test p<0.01 for comparison of diagnosis versus 1-year
839 samples. **(e)** Gene-set enrichment analysis of TNF α , TGF β and IL6-JAK-
840 STAT pathways comparing group-A BCR-ABL⁺ SCs at remission (n=122) vs
841 normal-HSCs (n=232). **(f)** Gene-set enrichment analysis of TNF- α and TGF
842 β pathways performed on normal-HSCs (n=232) vs group-A BCR-ABL⁺ SCs
843 at 3 months (n=72), 6 months (n=24) and over 1 year after TKI initiation
844 (n=27). **(g)** Beeswarm plots for 10 selected differentially expressed genes
845 between normal-HSCs (black; n=232, 5 donors), BCR-ABL⁺ SCs from patients
846 at diagnosis (red; n=477, 18 donors) and BCR-ABL⁺ SCs from patients at
847 remission group-A (light blue, n=122, 16 donors). Numbers of cells analyzed
848 and numbers showing amplification for the selected gene are shown below
849 the plot. The average gene expression levels are indicated by red squares,
850 the boxes represent the median and quartiles of gene expression levels.
851 Nonparametric Wilcoxon test p-values are shown on top of each bar graph.
852 Fisher's exact test p-values are shown below the graph.

853

854 **Figure 6.** Single-cell RNA-sequencing reveals heterogeneity of CML stem
855 cells associated with disease progression in CML. **(a)** tSNE visualization of
856 single normal-HSCs from 5 donors (grey circles; n=232), BCR-ABL⁺ SCs from
857 18 CP-CML patients (red triangles; n=477) BCR-ABL⁺ SCs from 3 patients at
858 the time of BC (CML1931, light blue squares, n=85; CML1266, purple
859 squares, n=63; CML1203, pink squares, n=7 and K562 cells (brown circles,
860 n=53). The tSNE has been generated using 207 differentially expressed
861 genes as described in the methods. **(b)** The heatmap shows the top 40 genes
862 differentially expressed between BCR-ABL⁺ SCs falling in the CP-CML cluster
863 (n=477) and BCR-ABL⁺ SCs falling in the BC cluster (n=155). The bar above
864 the heatmap indicates CP-CML cluster in red, BC-CML cluster in purple. **(c)**
865 tSNE visualization as shown in panel **a** but with BC (light-blue squares, n=85)
866 and pre-BC (orange diamonds, n=132) cells from patient 1931 highlighted.
867 Arrow indicates 8 pre-BC cells clustering separately from remaining pre-BC
868 cells and together with BCR-ABL⁺ CP-SCs. **(d)** Heatmap of log₂(RPKM) of
869 selected lymphoid and myeloid genes in BCR-ABL⁺ CML-SCs at BC (n=155,
870 3 donors), pre-BC (n=185, 2 donors with SCs from patient OX1931 annotated
871 according to those falling within the CP-CML cluster in yellow or BC-CML
872 cluster in orange) and CP-CML SCs at diagnosis (n=477, 18 donors) and
873 normal-HSCs (n=232, 5 donors), showing aberrant co-expression of lymphoid
874 and myeloid genes in SCs falling within the BC cluster. **(e)** Dot plot showing
875 index sort results corresponding to individual BCR-ABL⁺ SCs from pre-BCs
876 OX1931. The color and shape of the dots indicate if the SC clustered with CP-
877 CML BCR-ABL⁺ SCs (blue triangles) or with BC-CML BCR-ABL⁺ SCs (red
878 circles) according to the RNA-seq results presented in the tSNE analysis in
879 panel **a**. The value is expressed as fluorescent intensity for CD90 and
880 CD45RA antigens (y and x axis, respectively). **(f)** Histogram (left panel)
881 shows the frequency of differentially expressed genes between pre-BC
882 OX1931 BCR-ABL⁺ SCs clustering with CP-CML or BC-CML BCR-ABL⁺ SCs
883 (Y axis) with respect to the distance (Kb) of RUNX1 binding sites from the
884 respective transcription start site (TSS; X axis). The box plot (right panel)
885 shows the fraction of RUNX1 binding sites/window found in the genes that are
886 differentially expressed between the CP-CML and the BC-CML SC clusters

887 (red box) versus those found in background genes (grey box). P value
888 =0.0018 by Wilcoxon test.

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907 **Online Methods**

908 Cell Lines

909 Authenticated K562 and mycoplasma negative (Chronic Myeloid Leukemia,
910 human cell line) were obtained from ATCC and grown in IMDM, 10% Fetal
911 Bovine Serum (FBS).

912

913 Samples and bone marrow mononuclear cells processing

914 CML patients included in the study and their clinical details are listed in
 915 Supplementary Table 1. Patients provided written informed consent in
 916 accordance with the Declaration of Helsinki for sample collection and use in
 917 research under Oxford University ethics committee approval (MREC
 918 06/Q1606/110). Bone marrow (BM) mononuclear cells (MNCs) were isolated
 919 using Ficoll density gradient. Cryopreserved BM MNCs were thawed and
 920 processed for flow cytometry analysis as previously described⁴.

921

922 FACS staining and single-cell sorting

923 All Fluorescence Activated Cell Sorting (FACS) experiments included single-
 924 color stained CompBeads (BD Biosciences) and fluorescent-minus-one
 925 (FMO) controls. Live cells were selected based on their non-permeability and
 926 subsequent lack of fluorescence associated with 7AAD or DAPI. The
 927 combination of monoclonal antibodies used to identify hematopoietic stem
 928 and progenitor cell populations was previously described⁴ and is listed below.
 929 The cocktail of lineage markers used (Lin) was: CD2, CD3, CD4, CD7, CD8a,
 930 CD10, CD11b, CD14, CD19, CD20, CD56, CD235ab. Single cells were
 931 isolated from BM samples of healthy controls or CML patients. Single cells
 932 were FACS sorted as Lin-CD34⁺CD38⁻. For some experiments, index sort
 933 data of the mean fluorescence intensities (MFI) of CD90, CD45RA and
 934 CD123 were also recorded for each individual cell isolated.

935

936 Anti-human Antibodies

Antigen	Clone	Conjugate	Company
CD34	8G12	APC	Biologend
CD38	HIT2	PETXR	Life Technologies
CD90	5E10	PE	Biologend
CD45RA	MEM56	FITC	Life Technologies
CD123	6H6	PECy7	Biologend
CD2	RPA-2.10	PECy5	Biologend
CD3	HIT3a	PECy5	Biologend
CD4	RPA-T4	PECy5	Biologend
CD7	CD7-6B7	PECy5	Biologend
CD8a	RPA-T8	PECy5	Biologend
CD10	HI10a	PECy5	Biologend

CD11b	ICRF44	PECy5	Biolegend
CD14	RMO52	PECy5	Biolegend
CD19	HIB19	PECy5	Biolegend
CD20	2H7	PECy5	Biolegend
CD56	B159	PECy5	BD
CD235ab	HIR2	PECy5	Biolegend

937

938 Single-cell sorting was performed on FACS ARIA II, FACS ARIA III or FACS
939 ARIA Fusion (Becton Dickinson) directly into 96-well plates (PCR micro plate
940 Thermo-Fast 96 well, semi-skirted). To check the correct alignment of the
941 sorter, BD FACS™ Accudrop Beads (BD Biosciences) were deposited initially
942 onto the lid or film cover of a setup plate. After this sort, 50 beads were sorted
943 into several wells of a clean PCR plate where it was checked that the beads
944 formed a discrete drop in the center of the bottom of the well. If any splashing
945 was noticed on the sides of the well, the alignment was adjusted. Single 488-
946 Flow-Check Fluorospheres (Beckman Coulter) were then deposited into each
947 well of a flat bottomed 96 well tissue culture plate and single-cell mode sorting
948 was verified by checking the presence of 1 fluorosphere/well using a
949 conventional fluorescence microscope. The investigators were not blinded
950 when performing this and following steps of the experiments. Experiments
951 were not randomized.

952

953

954 *Short-term culture from single cells for first division measurement.*

955 Single Lin⁻CD34⁺CD38⁻ cells from normal BM donors or from CML patients
956 were sorted using FACS Aria III into 60-well Terasaki plates containing 25µl of
957 Stemspan SFEM (Stemcell Technologies) medium supplemented with 10%
958 BIT 9500 serum substitute (Stemcell Technologies), 2 mM L-Glutamine (P A
959 A Laboratories), 10⁻⁴ M 2-mercaptoethanol (Sigma), 100 U/mL
960 penicillin/streptomycin (PAA), 100 ng/mL rhSCF (Amgen), 100 ng/mL rhFLT3-
961 ligand (FL; Immunex), 50 ng/mL rhTPO (Peprotech), 10 ng/mL rhIL-3
962 (Peprotech), 10 ng/mL rhG-CSF (Amgen), 10 ng/mL rhIL-6 (Peprotech).
963 rhTNFα (Miltenyi Biotec) or rhTGFβ (Miltenyi Biotec) were added to the

964 culture at 20 ng/ml as indicated. Single cells were scored microscopically for
965 number of cells that had reached time of first division after 96 hours of culture.

966

967 Fluorescence in situ hybridization (FISH)

968 For interphase FISH, Lin⁻CD34⁺CD38⁻ cells were cytocentrifuged onto slides
969 and hybridized with the LSI BCR/ABL Dual Color, Dual Fusion Translocation
970 Probe (Abbot Molecular) spanning the ABL1 and BCR respective breakpoints
971 involved in the t(9;22) translocation (ABL1: 9q34, BCR: 22q11.2).
972 Fluorescence images were obtained with the use of fluorescence microscopy.
973 In nuclei from normal cells lacking the t(9;22) translocation the probe
974 hybridizing to ABL1 region appears as 2 orange signals while the probe
975 hybridizing to BCR region appears as 2 green signals. Nuclei containing a
976 balanced t(9;22) will display one orange and one green signal from the normal
977 9 and 22 chromosomes and two orange/green (yellow) fusion signals, one
978 each from the derivative 9 and 22 chromosomes.

979

980 Generation of single-cell cDNA libraries using Smart-seq2 protocol

981 Single K562 or Lin⁻CD34⁺CD38⁻ cells were FACS sorted into 96-well plates
982 (Thermo) containing 4 µL of a lysis mix including oligo dT (Biomers), RNase
983 inhibitor (Takara), dNTPs mix (Fermentas) at concentrations described in the
984 original Smart-seq2 protocol and listed in the lysis mix box below²⁴. ERCC
985 spikes (AMBion) were pre-diluted to 1:400,000 from stock concentration and
986 added to the lysis mix at a final dilution of 1:40,000,000. ERCC spikes were
987 not included in the analysis of patient samples as a number of samples were
988 analyzed before ERCC spikes were routinely included in the reaction.
989 Retrotranscription and PCR amplification steps were performed following the
990 Smart-seq2 protocol using reagents concentrations optimized for small cells
991 (see box RT mix and PCR mix for individual reagents concentration). The
992 thermal conditions for RT and PCR reactions were according to the original
993 Smart-seq2 protocol. The number of cycles used for PCR amplification was
994 22. After PCR amplification, cDNA libraries from single-cells were purified

995 using Ampure XP magnetic beads according to the manufacturer instructions
 996 in a ratio of 0.8 to 1 with cDNA. After purification the libraries were
 997 resuspended in 17.5 μ L of buffer EB (Qiagen) and stored at -20 °C. Quality
 998 and concentration of the cDNA libraries generated was assessed using High
 999 Sensitivity Bioanalyzer (Agilent).

1000

1001 Generation of single-cell cDNA libraries using BCR-ABL tSS2 protocol

1002 BCR-ABL targeted amplification Smart-seq2 protocol (tSS2) was implemented
 1003 during both RT and PCR steps of Smart-seq2 as described in Supplementary
 1004 Figure 2. During RT and PCR-amplification a pair of primers recognizing the
 1005 BCR and the ABL portion of the fusion transcripts (sequences indicated in
 1006 Supplementary Fig. 2b) were added to the RT and PCR mixes respectively at
 1007 the concentrations indicated in Supplementary Figure 2c in the condition 6.
 1008 The primers pair was designed to give rise to a PCR amplicon of 505 bp for
 1009 BCR-ABL e14a2 transcript and 430 for BCR-ABL e13a2 transcript.

1010

1011 Lysis Mix reagents

Reagent	Volume for 1 cell (μ L)
0.4% Triton X + RNase Inhibitor (1:20)	2
dNTPS (10 mM)	1
Oligo dT (10 μ M)	1
ERCC (pre-diluted 1:400,000)	0.1
TOTAL	4 (4.1 with ERCC)

1012

1013 RT mix reagents

Reagent	Volume for 1 cell (μ L)
Superscript II first strand buffer (5x)	2
DTT (100 mM)	0.5
Betaine (5 M)	2
MgCl ₂ (1 M)	0.1
RNase Inhibitor (40 U/ μ L)	0.25
TSO (100 μ M)	0.1
BCR-ABL Primer set #1 F+R (200 μ M)	0.07
Superscript II (200 U/ μ L)	0.25
Water	0.33
TOTAL	5.6

1014

1015 PCR mix reagents

Reagent	Volume for 1 cell (µL)
KAPA Hifi HS Ready Mix (2x)	12.5
ISPCR oligo (10 µM)	0.125
BCR-ABL Primer set#2 F+R (20 µM)	0.07
Water	2.305
TOTAL	15

1016

1017 Generation of single-cell cDNA libraries using BCR-ABL targeted amplification
1018 protocol with C1 microfluidic platform

1019 K562 cells were captured on a large-sized (17-25 µm cell diameter) C₁TM
1020 Single-Cell Auto Prep IFC for mRNA Sequencing (Fluidigm) using the
1021 Fluidigm C1 system. Cells were loaded onto the chip at a concentration of
1022 ~250K cells/mL and imaged by phase-contrast microscopy to check single-
1023 cell per capture site. Cells were lysed and cDNA prepared on the C1 Fluidigm
1024 chip according to manufacturer's protocol, using SMARTer Ultra Low RNA kit
1025 for Illumina (Clontech). BCR-ABL targeted amplification in the C1 setting was
1026 performed using the modified C1 PCR MIX protocol described in the table
1027 below. BCR-ABL Taqman assay (20X) was included in the C1 PCR MIX at a
1028 final dilution of 1:495. Any other step in the C1 protocol was performed
1029 following manufacturer's indications.

1030

1031 Modified C1 PCR MIX for BCR-ABL targeted amplification

Reagent	Volume (µL)
PCR Water (Advantage 2 Kit)	59.5
10X Advantage 2 PCR Buffer (Advantage 2 Kit)	10
50X dNTP Mix (Advantage 2 Kit)	4
IS PCR primer (Clontech SMARTer)	4
50X Advantage 2 Polymerase Mix (Advantage 2 Kit)	4
C1 Loading Reagent (Fluidigm)	4.5
BCR-ABL Taqman assay pre-diluted 1:22 (Hs03024541_ft)	4

1032

1033 Illumina library preparation and sequencing

1034 1.25 µL of cDNA was used for tagmentation reaction carried out with Nextera
1035 XT DNA Sample Preparation kit (Illumina) according to manufacturer's
1036 instruction but using one-fourth of the volumes. Purification of the product
1037 was done with a 1:1 ratio of AMPure XP beads with a final elution in 17.5 µl in
1038 resuspension buffer provided from the Nextera kit. Samples were loaded on a
1039 High-Sensitivity DNA chip to check the size and quality of the indexed library
1040 while the concentration was measured with Qubit High-Sensitivity DNA kit
1041 (Invitrogen). BCR-ABL⁺ or BCR-ABL⁻ SCs eligible for sequencing were
1042 selected based on the quality of their indexed cDNA libraries (size 400-900bp;
1043 concentration > 4 ng/ml). The number of BCR-ABL⁺ or BCR-ABL⁻ SCs to be
1044 sequenced per patient was determined by availability of SCs from each
1045 sample and space available per flow cell to ensure sufficient depth of
1046 sequencing. Libraries were pooled to a final concentration ranging between 3
1047 and 10 nM and were sequenced with Illumina HiSeq 2000 and Illumina HiSeq
1048 4000 (51 bp single-end read) at The Wellcome Trust Centre for Human
1049 Genetics in Oxford.

1050

1051 *BCR-ABL genotyping of single-cell cDNA libraries*

1052 *BCR-ABL* genotyping of cDNA libraries from single-cells was performed using
1053 QPCR reaction in a 384-well plate (Roche, Lightcycler). QPCR was performed
1054 in duplicate using 1.5 µL of the cDNA library for each reaction. The
1055 expression of *BCR-ABL* and *GAPDH* were measured using the following
1056 Taqman FAM-MGB assays *BCR-ABL*: Hs03024541_ft and *GAPDH*:
1057 Hs02758991_g1 (Life Technologies). The reactions were performed using a
1058 minimum of 60 cycles of amplification. We used QPCR and not raw
1059 sequencing reads to genotype cells for presence of BCR-ABL due to the low
1060 coverage of BCR-ABL in the sequencing data.

1061 *Exome Sequencing.*

1062 Genomic DNA was extracted from unfractionated BM MNCs from patient
1063 OX1931 at both pre-BC and BC stages using QIAamp DNA Blood Mini Kit
1064 (Qiagen) according to manufacturer's instructions. Exome capture was
1065 performed from GATC Biotech, using INVIEW Human Exome □Library

1066 preparation □ Enrichment with SureSelectXT Human All Exon Kit for Illumina
 1067 Paired-End Sequencing (Read length: 2 x 125 bp). The number of PCR
 1068 cycles performed for the amplification of the adaptor-ligated library was 5. The
 1069 number of cycles used for the post- hybridization captured library amplification
 1070 step was 12. The enriched exome fragments were pooled and paired-end
 1071 sequenced on a HiSeq 2000 platform (Illumina). From this we obtained >60x
 1072 on target coverage for the majority of positions for each of the samples.

1073

1074 Assessment of BCR-ABL tSS2 sensitivity using plasmid spike-in

1075 BCR-ABL breakpoint region (e14a2) was PCR-amplified from cDNA of K562
 1076 cells using specific BCR-ABL primer set #1 described in Supplementary
 1077 Figure 2b. The resulting PCR-amplicon was Sanger sequenced before being
 1078 cloned into the pcr™- Blunt II-TOPO® vector using Zero Blunt® TOPO PCR
 1079 Cloning Kit (Thermo Fisher). Correct size of the BCR-ABL insert was verified
 1080 by PCR. Concentration of the resulting BCR-ABL plasmid was measured
 1081 using Qubit (Invitrogen) and the absolute number of plasmid copies/μL was
 1082 calculated. Several plasmid pre-dilutions were produced in order to be able to
 1083 spike in the retro-transcription reaction of single BCR-ABL⁻ SCs (HSCs from a
 1084 normal donor), the desired amount of plasmid copies (1, 2, 5, 10, 20, 50, 100,
 1085 1000) always at a volume of 1 μL. The PCR step was performed according to
 1086 the standard BCR-ABL tSS2 protocol. Quantification of absolute number of
 1087 BCR-ABL amplified copies after BCR-ABL tSS2 reaction was carried out by
 1088 QPCR (Roche, Lightcycler) using a commercial BCR-ABL standard curve as
 1089 a reference (Ipsogen BCR-ABL1 Mbcr, Qiagen).

1090

1091 RT mix reagents for BCR-ABL plasmid standard curve

Reagent	Volume for 1 cell (μL)
Superscript II first strand buffer (5x)	2
DTT (100 mM)	0.5
Betaine (5 M)	2
MgCl ₂ (1 M)	0.1
RNAse Inhibitor (40 U/μL)	0.25
TSO (100 μM)	0.1
BCR-ABL Primer set #1 (200 μM)	0.07
Superscript II (200 U/μL)	0.25

Plasmid (serial pre-dilutions)	1
Water	0.33
TOTAL	6.6

1092

1093 *RUNX1 c.G521A detection with single-cell QPCR*

1094 RUNX1 c.G521A mutation detected in patient OX1931 by exome sequencing
1095 was PCR amplified and validated by Sanger sequencing.

1096 (Fw:GGCTGGCAATGATGAAAAC and

1097 Rev:CAATGGATCCCAGGTATTGG). A SNP genotyping Taqman assay

1098 specific for RUNX1 c.G521A was designed using the Custom Assay Design

1099 tool (ThermoFisher) and validated on positive (OX1931) and negative

1100 controls.

1101

1102 *Single-cell gene expression analysis*

1103 For single-cell gene expression analysis, single-cells isolated by FACS were

1104 collected in each well of a 96-well plate containing 5 µl Cells Direct One-Step

1105 qRT-PCR (Invitrogen) mix and pre-amplified as previously described⁴⁹. Pre-

1106 amplified samples were diluted 1:5 with TE before analysis of gene

1107 expression analysis on either a Fluidigm 96.96 or 192.24 Dynamic array using

1108 gene-specific Taqman assays (Life technologies). No template and no reverse

1109 transcriptase were included as negative controls.

1110

1111 *Taqman assays used for Dynamic Array*

1112

Gene Symbol	Taqman Assay ID
ABL1	Hs00245443_m1
ATG3	Hs00223937_m1
B2M	Hs00984230_m1
BCL2	Hs00608023_m1
BCR	Hs00244731_m1
BCR-ABL	Hs03024541_ft
BLNK	Hs00179459_m1
CD33	Hs01076281_m1
CD34	Hs00990732_m1
CD79A	Hs00998119_m1
CD79B	Hs01058826_g1

CD164	Hs00174789_m1
CDK6	Hs01026371_m1
CKLF	Hs03047057_s1
CLU	Hs00156548_m1
CSF1R	Hs00911250_m1
CTNNB1	Hs00355049_m1
CXCR4	Hs00607978_s1
DNTT	Hs00172743_m1
FCER1A	Hs00758600_m1
GAPDH	Hs02758991_g1
GAS2	Hs01086684_m1
GOLGA8A	Hs01104342_m1
HPRT	Hs02800695_m1
HSP90A1	Hs03043878_g1
IFITM1	Hs00705137_s1
IGF1R	Hs00609566_m1
IGJ	Hs00376160_m1
ITGA6	Hs01041011_m1
MEIS1	Hs01017441_m1
MLLT3	Hs00180312_m1
MMRN1	Hs00201182_m1
MPL	Hs00180489_m1
MZB1	Hs00414907_m1
PTRF	Hs00396859_m1
RGS2	Hs01009070_g1
RXFP1	Hs01073141_m1
SAT1	Hs00161511_m1
SELL	Hs00174151_m1
SELP	Hs00927900_m1
SOD2	Hs00167309_m1
TESPA1	Hs00207702_m1
VWF	Hs01109446_m1

1113

1114 *Analysis of quantitative PCR single-cell gene expression data*

1115 We calculated ΔC_t values, which are relative to the mean expression level of
1116 two housekeeping genes (B2M and GAPDH). As previously described^{25,50}, C_t
1117 values were subtracted from the limit of detection (CT=30) followed by
1118 subtraction of the mean C_t value of housekeeping genes for each cell. $C_t=40$
1119 was used for the comparative analysis of the detection of BCR-ABL and
1120 GAPDH in K562 cells between Smart-seq2 and BCR-ABL tSS2 protocols.
1121 Cells not expressing at the 15th percentile of all genes, or two housekeeping
1122 genes were removed from the analysis. Analysis of differential gene
1123 expression between BCR-ABL⁺ and BCR-ABL⁻ SCs was performed using the

1124 Wilcoxon test and Fisher's exact test to compare expression level and
1125 expression frequency respectively.

1126

1127 Analysis of single-cell RNA sequencing

1128 Short reads (51-bp) were aligned to the human genome (GRCh37 assembly
1129 (hg19)) using Tophat⁵¹ with a supplied set of known RefSeq transcripts as the
1130 input. The mapping parameters '-g 1' was used to allow one alignment to the
1131 reference for a given read. Expression values were quantified as read per
1132 kilobase of transcript length per million mapped reads (RPKM) based on the
1133 RefSeq gene model using the rpkmforgenes⁵². As previously demonstrated
1134 the reliable classification of cell types at a sequencing depth of 50,000 reads
1135 per cell,^{53,54} we used cells with higher than 50,000 mapping reads and 1,000
1136 detected genes (RPKM ≥ 1) for the downstream analysis. We used the genes
1137 that were highly expressed in more than 50% of each population of cells to
1138 identify the candidate outliers based on gene expression level, similar to the
1139 method previously described in Singular™ from Fluidigm, using the standard
1140 method for the outlier detection⁵⁵. The modified Z-scores were calculated
1141 using the formula $0.6745(x_i - \tilde{x})/MAD$; MAD denoting the median absolute
1142 deviation and \tilde{x} denoting the median. Cells with the absolute modified Z-score
1143 greater than 3 were considered as candidate outliers (28 out of 2,287 cells),
1144 and these cells were monitored during the analysis. We found that excluding
1145 or including them in our analysis did not have any significant impact on the
1146 results.

1147

1148 Analysis of the effective sequencing depth

1149 To examine the effective sequencing depth, we selected 12 normal-HSCs
1150 with a sequencing depth larger than 6 million mapped reads. We randomly
1151 sampled reads in the range of 0.1 to 6 million mapped reads and calculated a
1152 number of detected genes with RPKM ≥ 1 in each category. We observed that
1153 detected number of genes plateaued at a sequencing depth of beyond 1
1154 million mapped reads per cell (Fig. 2a).

1155

1156 *T-distributed stochastic neighbor embedding (tSNE) analysis*

1157 As previously described, the advantage of using t-distributed stochastic
1158 neighbor embedding (tSNE) over a traditional principal component analysis is
1159 to visualize the projection of high dimensional single-cell gene expression
1160 data into a low dimensional space.^{25,56,57} We selected genes expressed in ≥ 10
1161 cells with a coefficient of variation score (CV), “standard deviation/mean”, ≥ 1
1162 and the summed up of genes expression values in log2 scale ≥ 1 for the tSNE
1163 analysis. We normalized the RPKM values into log2(RPKM) scale and set up
1164 the limit of detection at 1 RPKM. Log2 scale of genes expressed < 1 RPKM
1165 was set up to 0. Possible batches from processing samples in different dates
1166 were removed from expression values using the function “removeBatchEffect”
1167 in Limma package⁵⁸. We then downloaded the tSNE software from
1168 <https://lvdmaaten.github.io/tsne/> to perform the analysis using the Matlab
1169 implementation with “initial dims=20” and “perplexity=20” parameters.

1170

1171 To identify variable genes, similarly to a previously described approach⁵⁶, we
1172 fitted a simple noise model using the lowess model of mean expression level
1173 and the coefficient of variation (CV) to estimate the high variable genes from
1174 each type of cells. The lowess model predicted 3,368, 5,611, and 5,011 and
1175 5,522 genes from K562, BCR-ABL⁻ and BCR-ABL⁺ SCs (from diagnosis), and
1176 normal-HSCs that show high variation compared to the whole genes set with
1177 mean of expression log2(RPKM) higher than 0. We next used these genes for
1178 the tSNE analysis, and compared the tSNE results to the previous tSNE
1179 results of different gene sets. We found the same pattern of clustering,
1180 suggesting reproducibility of our results. We then selected 3,368, 7,428
1181 (combined variable genes from K562 and normal-HSCs) and 8,589 (combined
1182 variable genes from BCR-ABL⁻ SCs, BCR-ABL⁺ SCs and normal-HSCs)
1183 genes to generate Figure 1d, Figures 2b and 3c respectively. 5,011 and 5,611
1184 variable genes in BCR-ABL⁺ and BCR-ABL⁻ SCs were used to generate
1185 Figure 4a and 4b for the good and poor responder classification.

1186

1187 For the tSNE analysis of samples following TKI therapy, we performed the
1188 random forests analysis of normal-HSCs, BCR-ABL⁺ SCs (diagnosis), and
1189 BCR-ABL⁺ SCs (remission) cells using the “randomForest” package in R
1190 (ntree parameter = 2,000). We obtained top 500 important genes, measured
1191 by the Gini index (Supplementary Table 7). These genes were used for
1192 distinguishing normal-HSCs from BCR-ABL⁺ SCs at diagnosis and during
1193 remission. We next used this gene set for the tSNE analysis of the remission
1194 cells (Fig. 5a). We applied K-means clustering (k=3) based on tSNE analysis
1195 results (from dimensions 1 and 2) to assign remission cells to group A and
1196 group B (Fig. 5a).

1197

1198 For the tSNE analysis of normal-HSCs, K562, and BCR-ABL⁺ SCs from
1199 diagnosis, pre-blast crisis and blast crisis samples, we obtained combined
1200 differentially expressed genes from the multiple ways comparison. 207 genes
1201 shown to be differentially expressed between BCR-ABL⁺ SCs from 18 chronic-
1202 phase CML patients (n=477), BCR-ABL⁺ SCs at BC (n=148), BCR-ABL⁺ SCs
1203 at pre-BC (n=185), and normal-HSCs (n=232). We next used this gene set for
1204 the tSNE to generate Figures 6a, 6c, and Supplementary Figure 14a. We note
1205 that pre-BC cells were involved in the tSNE analysis but were not shown in
1206 the Figure 6a.

1207

1208 Cell to cell variation analysis

1209 To analyze the variation within K562 and normal-HSCs, the Pearson
1210 correlation was calculated based on log₂(RPKM) expression values among
1211 the cells of each group using the same set of genes from the tSNE analysis.
1212 Kolmogorov–Smirnov test was used to test the difference of correlation score
1213 distribution (Supplementary Fig. 4e).

1214

1215 Differentially expressed gene analysis

1216 Differentially expressed gene analysis was performed using the
1217 nonparametric Wilcoxon test on log₂(RPKM) expression values for the

1218 comparison of expression level and Fisher's exact test for the comparison of
1219 expressing cell frequency. P-values generated from both tests were then
1220 combined using Fisher's method and were adjusted using Benjamini-
1221 Hochberg (BH). Differentially expressed genes were selected based on the
1222 absolute log2 fold change ≥ 1 and the adjusted p-value < 0.05 . Selected genes
1223 were subjected to the hierarchical clustering analysis using Pearson
1224 correlation as a distance with the complete clustering method performed in R
1225 with the "pheatmap" function. Beeswarm plots from selected genes were
1226 generated using the "beeswarm" package in R. We determined the top 100
1227 differentially expressed genes ranked by adjusted p-values from normal-HSCs
1228 against BCR-ABL⁻ and BCR-ABL⁺ SCs, and BCR-ABL⁻ against BCR-ABL⁺
1229 SCs at diagnosis. We then used the 245 unique genes from this analysis to
1230 make the heatmap of Figure 3d.

1231

1232 Gene Set Enrichment Analysis (GSEA)

1233 GSEA⁵⁹ was performed using GSEA software
1234 (<http://www.broadinstitute.org/gsea>) with permutation on the phenotype, 1000
1235 permutations, and default values for other parameters. Gene sets used in this
1236 study were selected from the CML, proliferation, quiescence, and HSC related
1237 pathways shown in Supplementary Table 3 and the MSigDB hallmark gene
1238 sets (Supplementary Table 5 and
1239 <http://www.broadinstitute.org/gsea/msigdb/collections.jsp>).

1240

1241 Comparison of bulk and single-cells analysis

1242 To compare differentially expressed genes identified between analysis
1243 performed at the bulk level and at the single-cell level, in *silico* bulk data were
1244 generated by generating a data "ensemble", by combining mapped reads per
1245 gene for all cells, for 5 normal donors (5 replicates from the ensemble of
1246 single cells from 5 donors plus 1 set of normal-HSCs from a sixth donor, that
1247 were isolated as a bulk population of 100 cells rather than as single-cells). 18
1248 replicates were generated from the ensemble of single-cells from each CML
1249 patient. DESeq2⁶⁰ was then performed from the raw read count of the

1250 ensemble to get differentially expressed genes. We then performed
1251 differentially expressed genes analysis for the single-cell analysis as
1252 described above using 232 single normal-HSCs against 477 BCR-ABL⁺ SCs
1253 and 377 BCR-ABL⁻ SCs. To make a comparison of differentially expressed
1254 genes, we applied the same cutoff (adjusted p-value < 0.05 and the absolute
1255 log₂ fold change ≥0.5) to get the number of differentially expressed genes
1256 from both bulk and single-cells analyses.

1257

1258 Co-expression analysis of G2M, lymphoid and myeloid genes

1259 We selected a gene set from gene ontology “G2M transition of mitotic cell
1260 cycle (GO:0000086)” from “amigo2.berkeleybop.org/amigo/term/GO:0000086”
1261 to analyze co-expression. Gene would be called “present” when the quantile
1262 normalized RPKM value ≥1 and “absent” when RPKM values <1. We counted
1263 the frequency of genes that expressed in the same cells (Supplementary Fig.
1264 7h). We calculated the co-expression frequency of random genes set
1265 excluding cell cycle related genes as the background. Kolmogorov–Smirnov
1266 test was used to test the difference of correlation score distribution. For
1267 analysis of co-expression of lymphoid and myeloid genes, we selected known
1268 lymphoid and myeloid signature genes to show the co-expression analysis in
1269 the CP-CML, BC-CML, and normal-HSC clusters. The heatmap was
1270 generated using the log₂(RPKM) with the pheatmap function in R (Fig. 6d).

1271

1272 ChIP-Seq analysis

1273 RUNX1 and IgG control ChIP-Seq data (CD34⁺ HSPCs)⁶¹ were downloaded
1274 from the GEO database (GSE45144). Raw reads were mapped to the human
1275 genome (GRCh37 assembly (hg19)) using bowtie2⁶² with default parameters.
1276 The peak calling was performed by MACS2⁶³ using default parameters with
1277 IgG ChIP-Seq data as a control. 8,706 RUNX1 binding sites were identified
1278 with 5% FDR. We next calculated the distribution of distances between
1279 RUNX1 binding sites and transcription start sites (TSS) of differentially
1280 expressed genes that are between the CP-CML and the BC-CML clusters.
1281 We further analyzed the fraction of RUNX1 binding sites/window (ranging

1282 from ± 0.5 to ± 10 kb windows around a TSS of a gene) found in the
1283 differentially expressed genes in comparison to randomly selected 2,000
1284 background genes (Fig. 6f).

1285

1286 **Code availability**

1287 R and MATLAB scripts used for data analyses are available on request.

1288

1289 **Statistical analysis**

1290 All statistical analyses were performed in R and GraphPad Prism 6
1291 (GraphPad Software, San Diego, CA). For single-cell expression levels,
1292 nonparametric Wilcoxon test was used, and Fisher's exact test was used to
1293 compare expression frequencies at the single cell level between defined
1294 populations. No statistical method was used to predetermine sample size, and
1295 experiments were not randomized. The Investigators were not blinded to
1296 allocation during experiments or outcome assessment.

1297

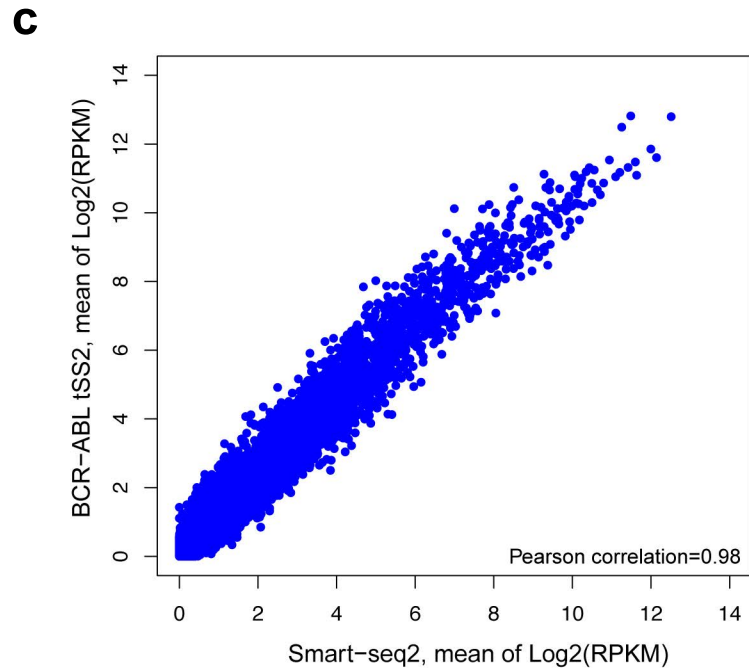
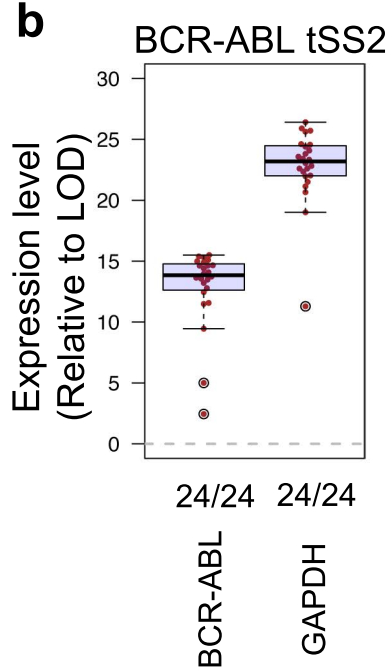
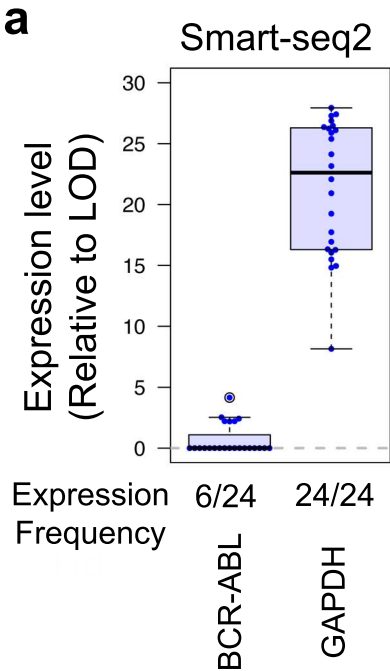
1298 **Methods-only References**

1299

1300

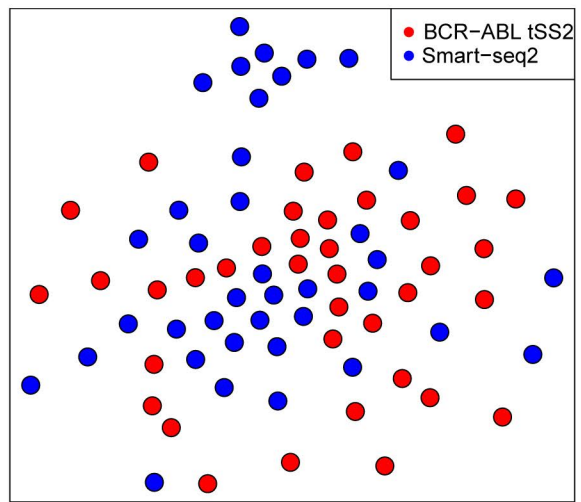
- 1301 49. Sanjuan-Pla, A., *et al.* Platelet-biased stem cells reside at the apex of the
1302 haematopoietic stem-cell hierarchy. *Nature* **502**, 232-236 (2013).
- 1303 50. Guo, G., *et al.* Resolution of cell fate decisions revealed by single-cell gene
1304 expression analysis from zygote to blastocyst. *Dev Cell* **18**, 675-685
1305 (2010).
- 1306 51. Kim, D., *et al.* TopHat2: accurate alignment of transcriptomes in the
1307 presence of insertions, deletions and gene fusions. *Genome Biol* **14**, R36
1308 (2013).
- 1309 52. Ramskold, D., Wang, E.T., Burge, C.B. & Sandberg, R. An abundance of
1310 ubiquitously expressed genes revealed by tissue transcriptome sequence
1311 data. *PLoS Comput Biol* **5**, e1000598 (2009).
- 1312 53. Pollen, A.A., *et al.* Low-coverage single-cell mRNA sequencing reveals
1313 cellular heterogeneity and activated signaling pathways in developing
1314 cerebral cortex. *Nat Biotechnol* **32**, 1053-+ (2014).
- 1315 54. Streets, A.M. & Huang, Y.Y. How deep is enough in single-cell RNA-seq?
1316 *Nat Biotechnol* **32**, 1005-1006 (2014).
- 1317 55. Iglewicz, B. & Hoaglin, D.C. *How to detect and handle outliers*, (ASQC
1318 Quality Press, Milwaukee, Wis., 1993).
- 1319 56. Zeisel, A., *et al.* Brain structure. Cell types in the mouse cortex and
1320 hippocampus revealed by single-cell RNA-seq. *Science* **347**, 1138-1142
1321 (2015).

- 1322 57. Saadatpour, A., Guo, G., Orkin, S.H. & Yuan, G.C. Characterizing
1323 heterogeneity in leukemic cells using single-cell gene expression analysis.
1324 *Genome Biol* **15**, 525 (2014).
- 1325 58. Ritchie, M.E., *et al.* limma powers differential expression analyses for
1326 RNA-sequencing and microarray studies. *Nucleic Acids Res* **43**, e47
1327 (2015).
- 1328 59. Subramanian, A., *et al.* Gene set enrichment analysis: a knowledge-based
1329 approach for interpreting genome-wide expression profiles. *Proc Natl*
1330 *Acad Sci U S A* **102**, 15545-15550 (2005).
- 1331 60. Love, M.I., Huber, W. & Anders, S. Moderated estimation of fold change
1332 and dispersion for RNA-seq data with DESeq2. *Genome biology* **15**, 550
1333 (2014).
- 1334 61. Beck, D., *et al.* Genome-wide analysis of transcriptional regulators in
1335 human HSPCs reveals a densely interconnected network of coding and
1336 noncoding genes. *Blood* **122**, e12-22 (2013).
- 1337 62. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2.
1338 *Nature methods* **9**, 357-359 (2012).
- 1339 63. Zhang, Y., *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome biology*
1340 **9**, R137 (2008).
1341

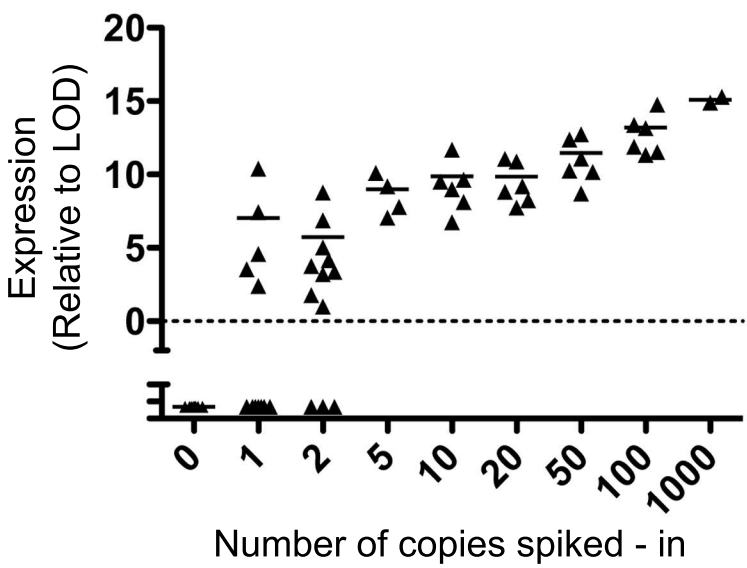


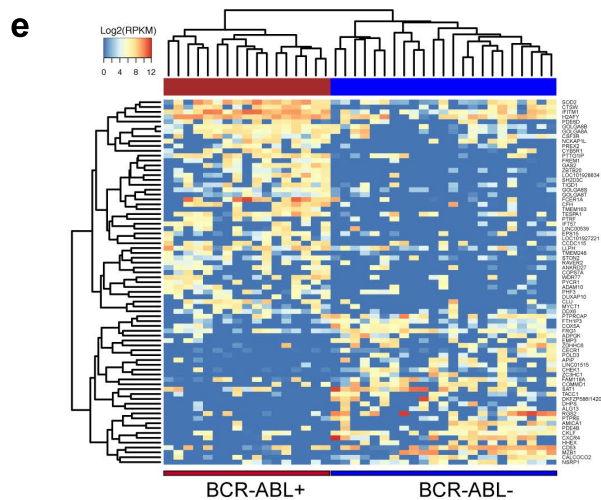
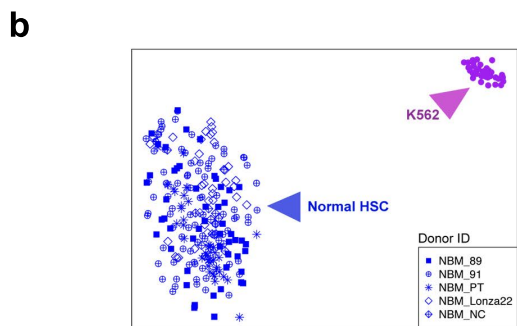
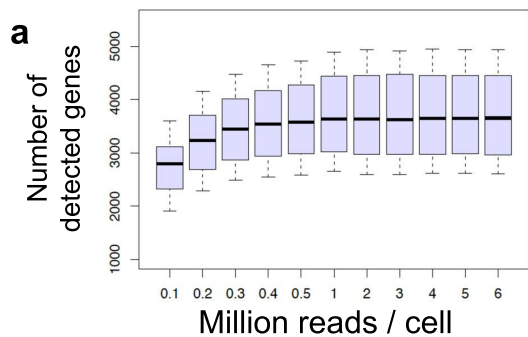
d

% Positive	0	41	75	100	100	100	100	100	100
% Positive expected by Poisson distribution	0	63	86	99	100	100	100	100	100



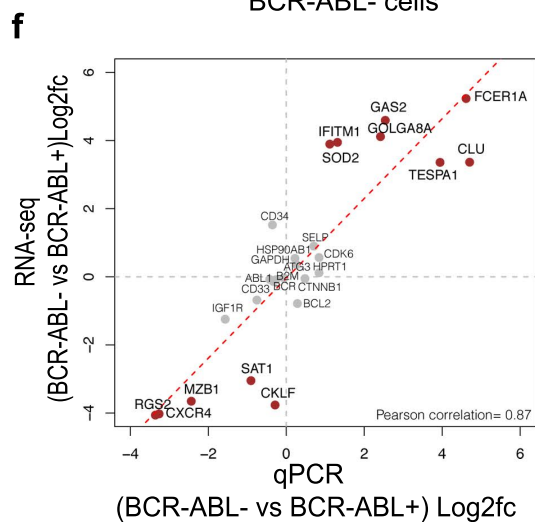
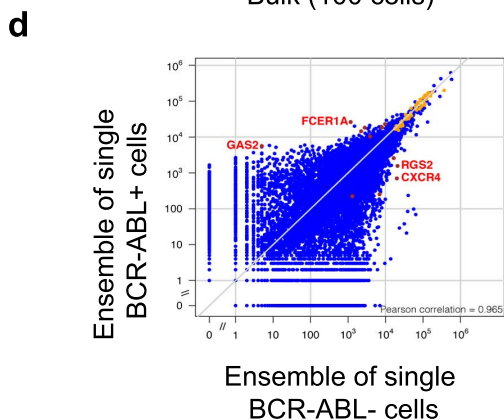
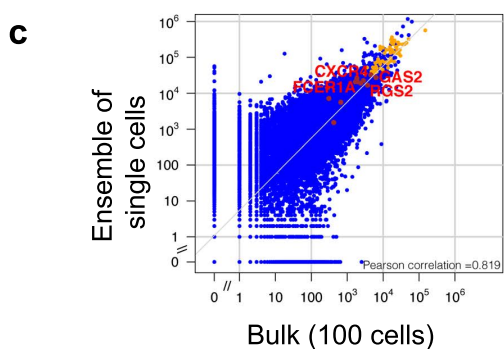
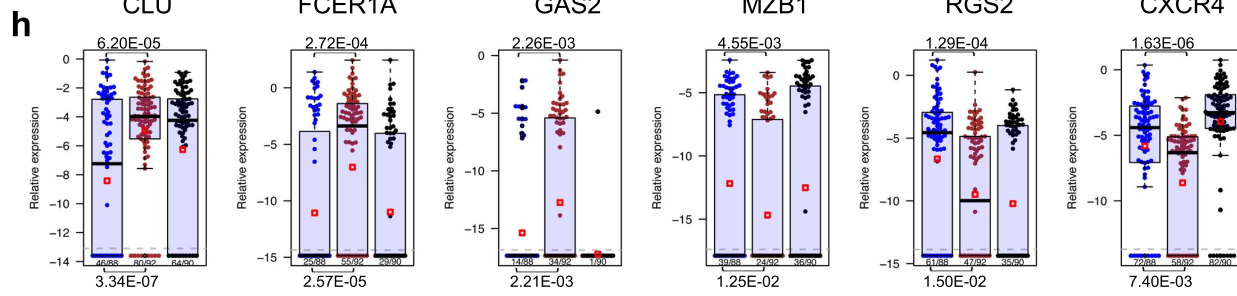
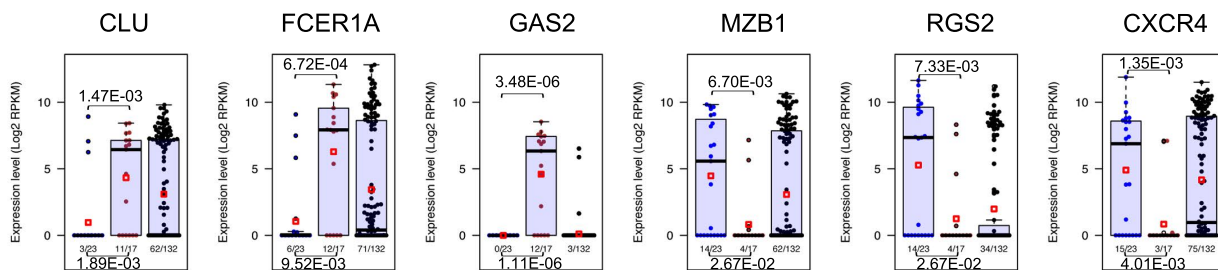
e

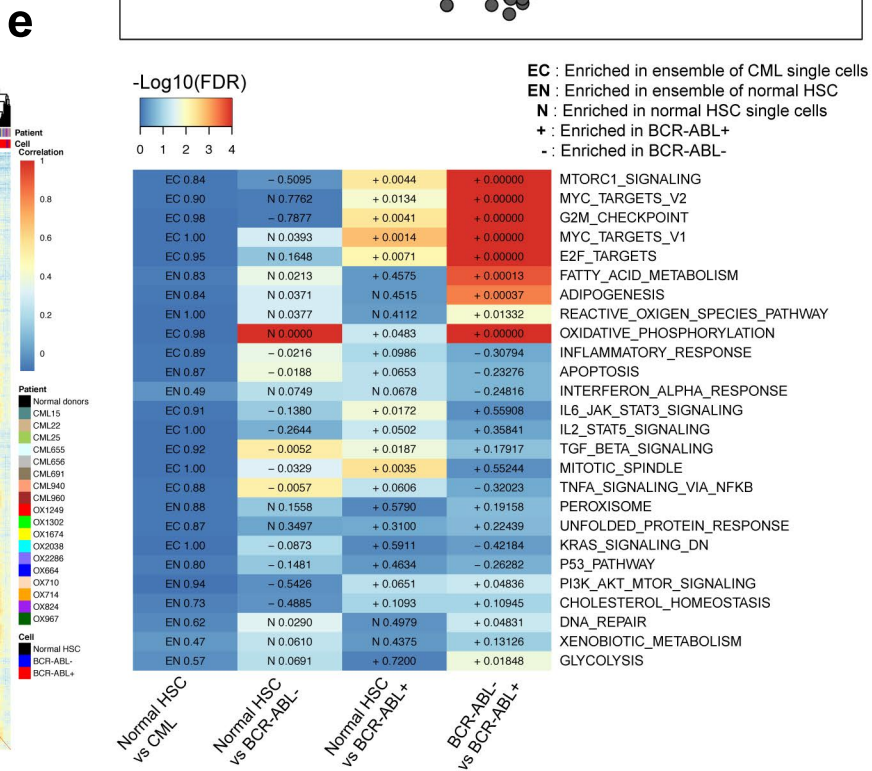
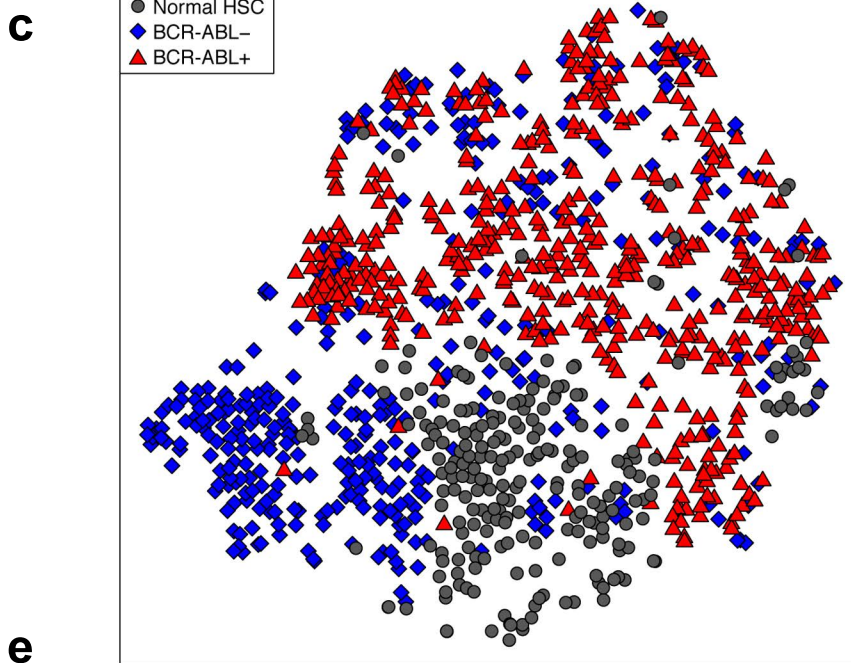
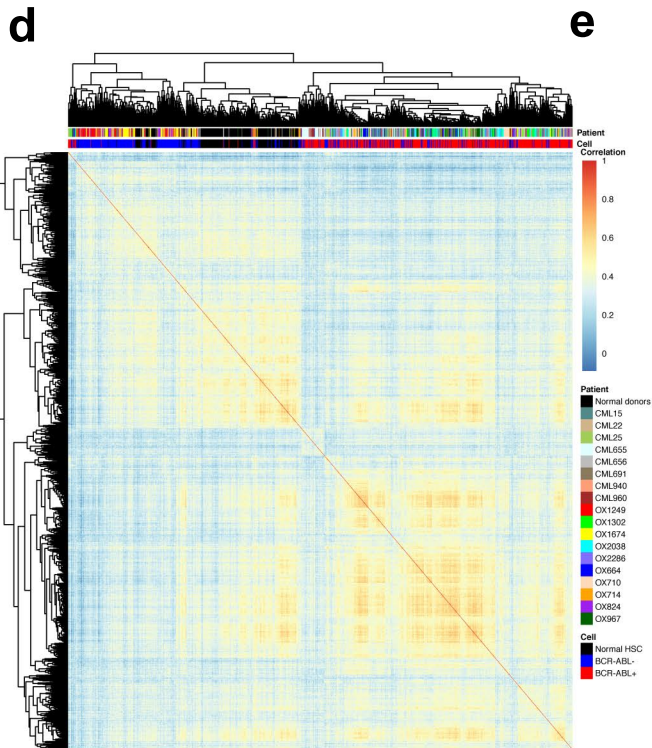
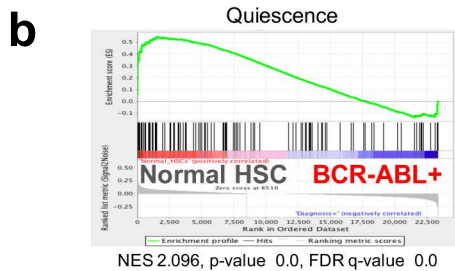
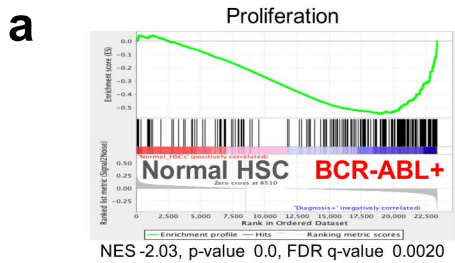


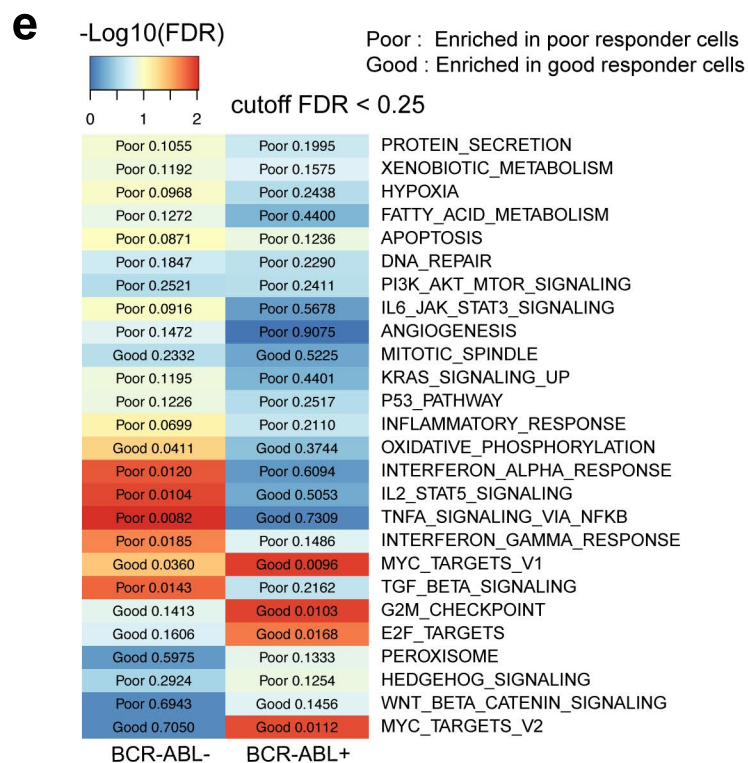
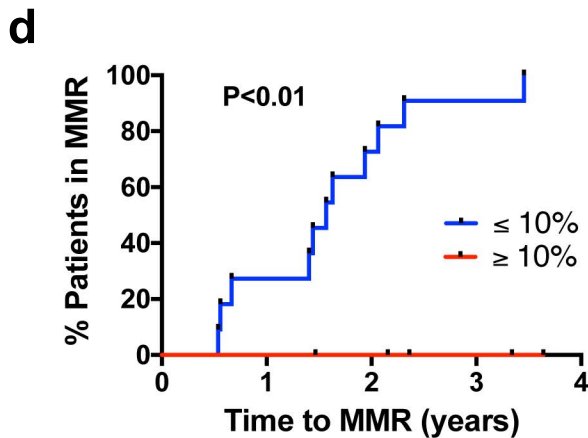
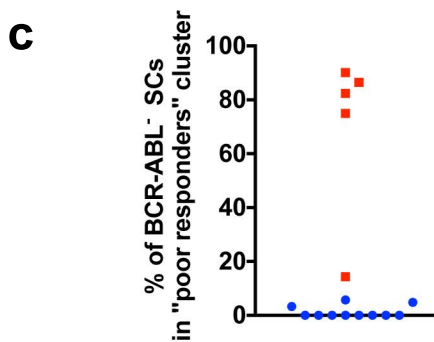
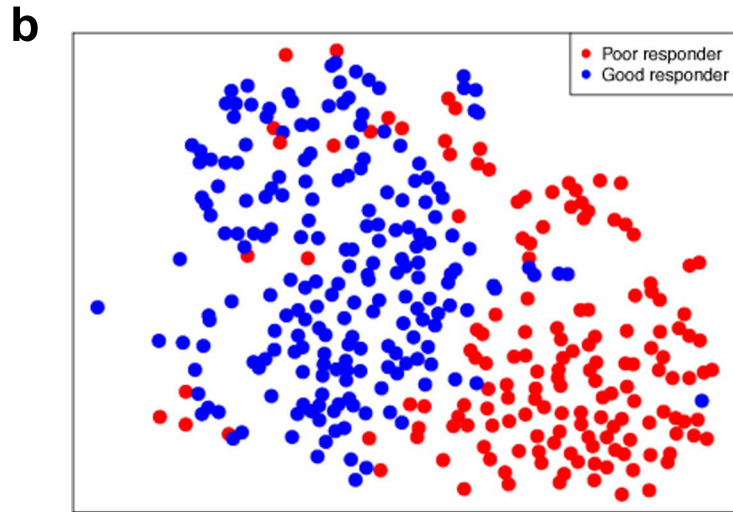
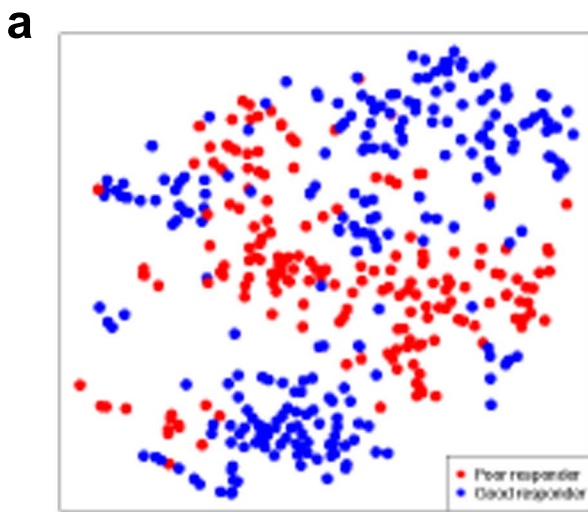


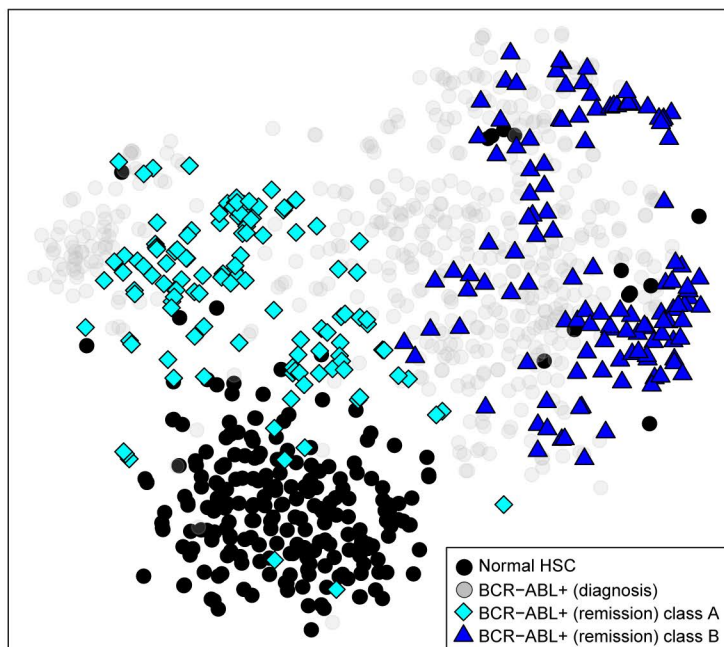
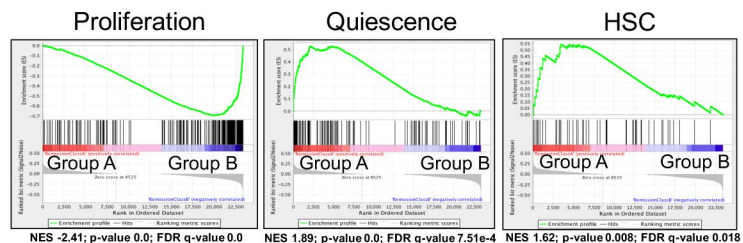
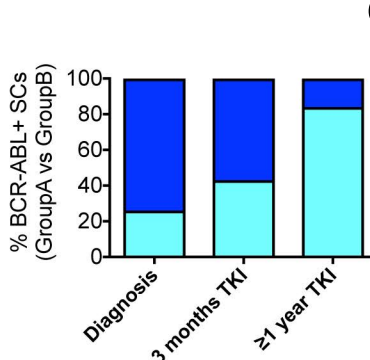
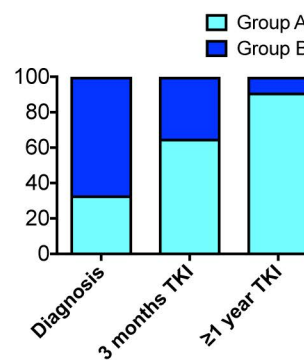
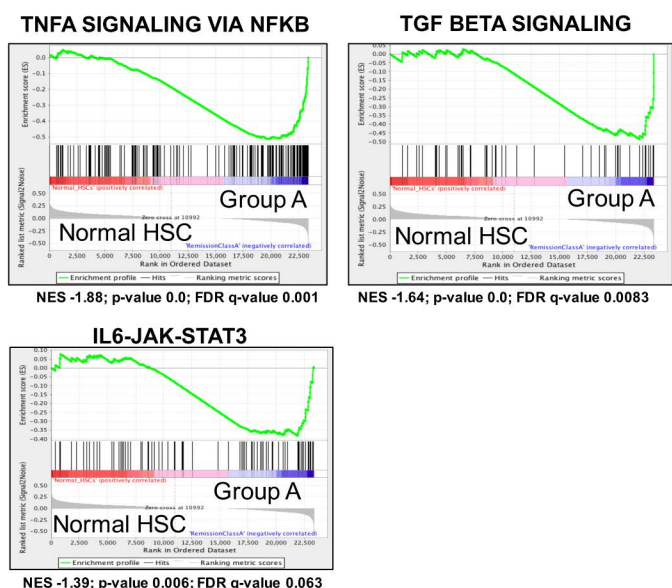
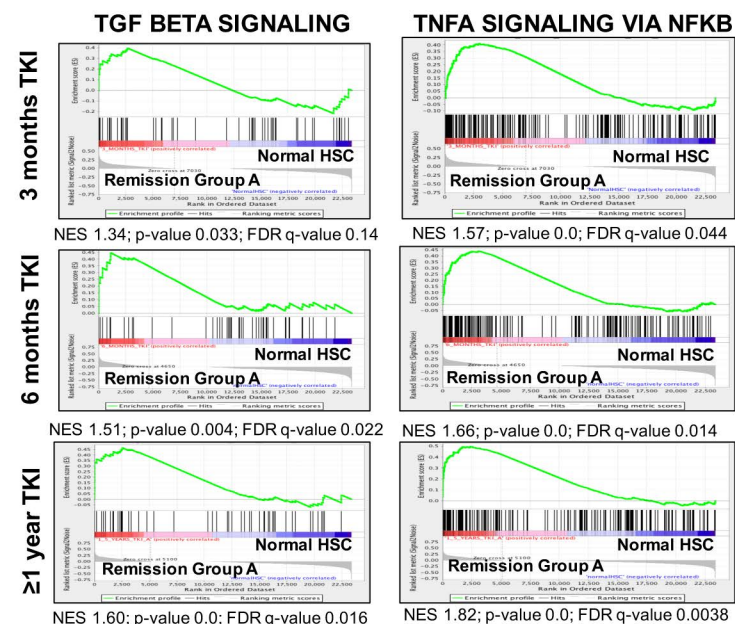
g

- BCR-ABL-
- BCR-ABL+
- Normal HSC

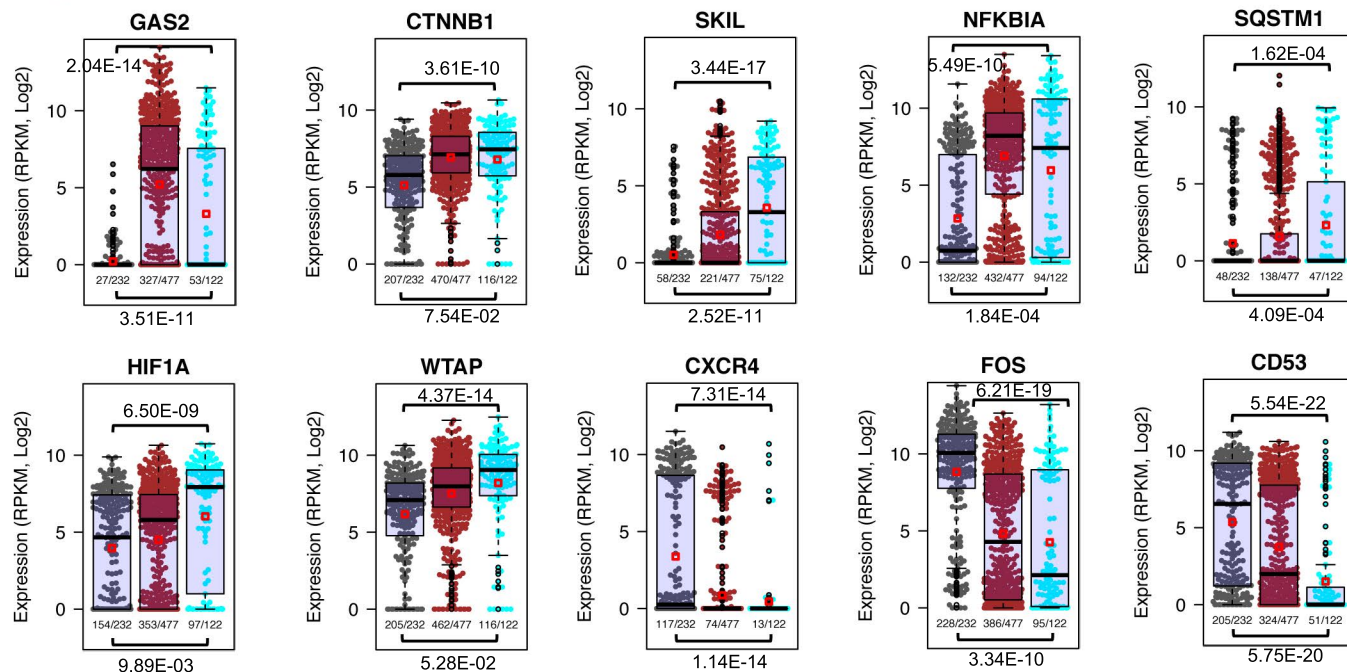


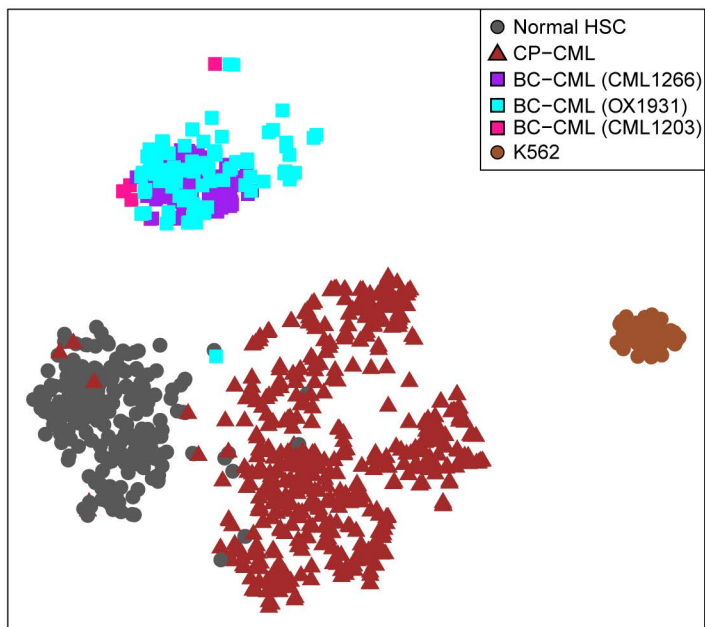
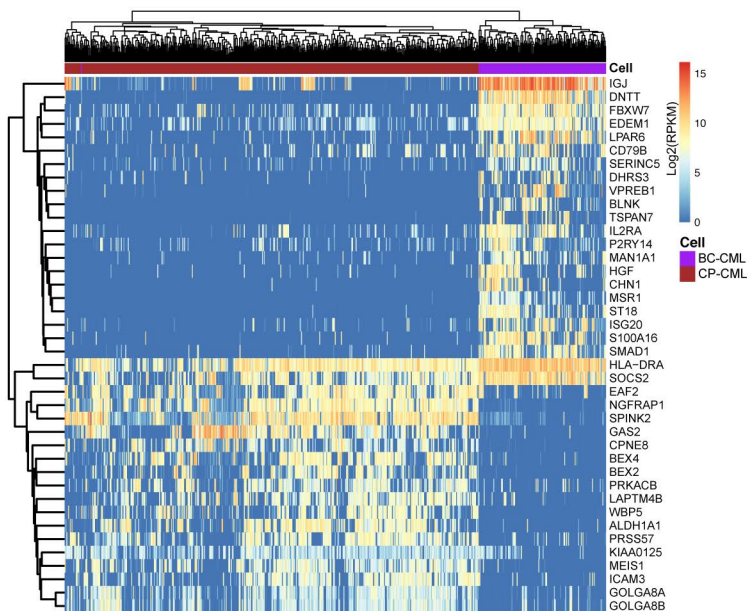
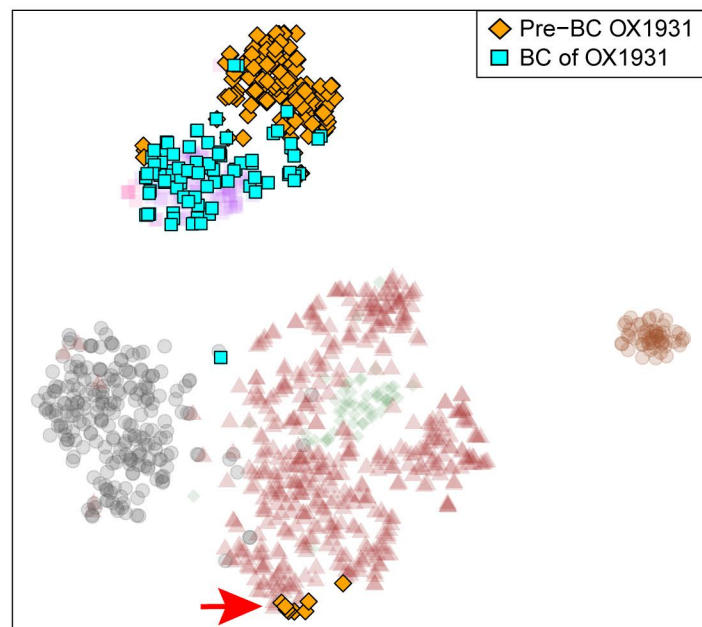
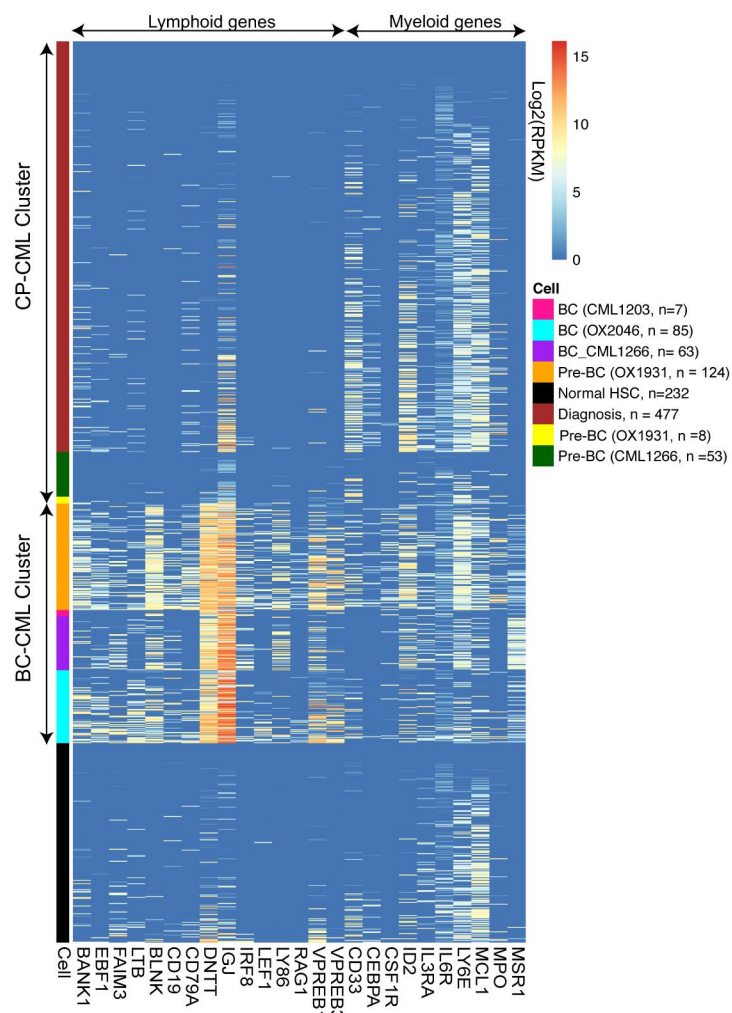
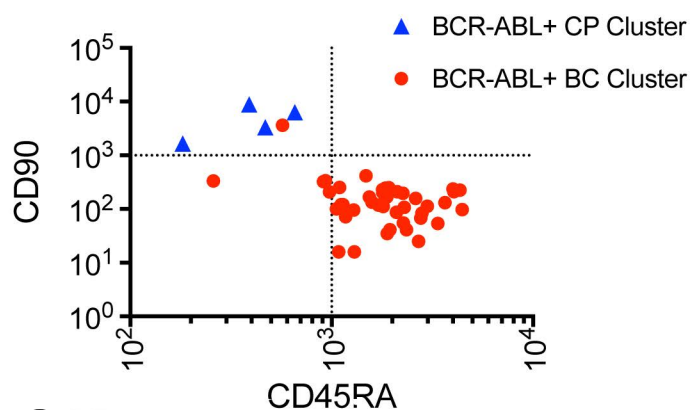




a**b****c****d****e****f****g**

● Normal HSC
● BCR-ABL+ CP-CML
● BCR-ABL+ Group A



a**b****c****d****e****f**