



**HAL**  
open science

## Single-cell multi-omics identifies chronic inflammation as a driver of TP53-mutant leukemic evolution

Alba Rodriguez-Meira, Ruggiero Norfo, Sean Wen, Agathe L Chédeville, Haseeb Rahman, Jennifer O'sullivan, Guanlin Wang, Eleni Louka, Warren W Kretzschmar, Aimee Paterson, et al.

### ► To cite this version:

Alba Rodriguez-Meira, Ruggiero Norfo, Sean Wen, Agathe L Chédeville, Haseeb Rahman, et al.. Single-cell multi-omics identifies chronic inflammation as a driver of TP53-mutant leukemic evolution. Nature Genetics, 2023, 55 (9), pp.1531-1541. 10.1038/s41588-023-01480-1 . hal-04427133

**HAL Id: hal-04427133**

**<https://hal.science/hal-04427133>**

Submitted on 30 Jan 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Single-Cell Multi-Omics Identifies Chronic Inflammation as a Driver of *TP53* mutant Leukaemic Evolution

Alba Rodriguez-Meira <sup>1,2, #, §, \*</sup>, Ruggiero Norfo <sup>1,2, #, &</sup>, Sean Wen <sup>1,2,3, #</sup>, Agathe L. Chédeville <sup>4,5,6,7, #</sup>, Haseeb Rahman <sup>1,2</sup>, Jennifer O'Sullivan <sup>1,2</sup>, Guanlin Wang <sup>1,2,3</sup>, Eleni Louka <sup>1,2</sup>, Warren W. Kretzschmar <sup>8,9,10</sup>, Aimee Paterson <sup>1,2</sup>, Charlotte Brierley <sup>1,2,11</sup>, Jean-Edouard Martin <sup>4,5,6,7</sup>, Caroline Demeule <sup>12</sup>, Matthew Bashton <sup>13</sup>, Nikolaos Sousos <sup>1,2</sup>, Daniela Moralli <sup>14</sup>, Lamia Subha Meem <sup>14</sup>, Joana Carrelha <sup>1</sup>, Bishan Wu <sup>1</sup>, Angela Hamblin <sup>2</sup>, Helene Guermouche <sup>15</sup>, Florence Pasquier <sup>4,5,6,16</sup>, Christophe Marzac <sup>4,5,6,17</sup>, François Girodon <sup>12,18</sup>, William Vainchenker <sup>4,5,6</sup>, Mark Drummond <sup>19</sup>, Claire Harrison <sup>20</sup>, J. Ross Chapman <sup>21</sup>, Isabelle Plo <sup>4,5,6</sup>, Sten Eirik W. Jacobsen <sup>1,8,9,10</sup>, Bethan Psaila <sup>1,2</sup>, Supat Thongjuea <sup>3</sup>, Iléana Antony-Debré <sup>4,5,6, §, \*</sup>, Adam J Mead <sup>1,2, §, \*</sup>

1. Haematopoietic Stem Cell Biology Laboratory, Medical Research Council Molecular Haematology Unit, Medical Research Council Weatherall Institute of Molecular Medicine, University of Oxford, Oxford OX3 9DS, UK
2. NIHR Biomedical Research Centre, University of Oxford, Oxford, OX3 9DS, United Kingdom
3. Medical Research Council Centre for Computational Biology, Weatherall Institute of Molecular Medicine, University of Oxford, Oxford OX3 9DS, UK
4. INSERM, UMR 1287, Villejuif, France
5. Gustave Roussy, Villejuif, France
6. Université Paris Saclay, Gif-sur-Yvette, France
7. Université Paris Cité, Paris, France
8. Department of Cell and Molecular Biology, Karolinska Institutet, SE-171 77 Stockholm, Sweden
9. Karolinska University Hospital, Stockholm, Sweden
10. Center for Hematology and Regenerative Medicine, Department of Medicine Huddinge, Karolinska Institutet, Karolinska University Hospital, SE-141 86 Stockholm, Sweden.
11. Center for Hematological Malignancies, Memorial Sloan Kettering Cancer Center, New York, NY, USA
12. Laboratoire d'Hématologie, CHU Dijon, Dijon, France
13. The Hub for Biotechnology in the Built Environment, Faculty of Health and Life Sciences, Northumbria University, Newcastle upon Tyne, NE1 8S, UK
14. The Wellcome Centre for Human Genetics, Roosevelt Drive, Oxford, OX3 7BN, UK
15. Sorbonne Université, INSERM, Centre de Recherche Saint-Antoine, AP-HP, Hôpital Saint-Antoine, Service d'hématologie biologique, F-75012, Paris
16. Département d'Hématologie, Gustave Roussy, Villejuif, France
17. Laboratoire d'Immuno-Hématologie, Gustave Roussy, Villejuif, France
18. INSERM, UMR 1231, Centre de Recherche, Dijon, France
19. Beatson Cancer Centre, Glasgow, UK.
20. Guy's and St Thomas' NHS Foundation Trust, London, UK.
21. Genome Integrity Laboratory, Medical Research Council Molecular Haematology Unit, Medical Research Council Weatherall Institute of Molecular Medicine, University of Oxford, Oxford OX3 9DS, UK.

# These authors contributed equally

§ These authors jointly supervised the work

\* Correspondence: [albarmeira@gmail.com](mailto:albarmeira@gmail.com); [ileana.antony-debre@gustaveroussy.fr](mailto:ileana.antony-debre@gustaveroussy.fr); [adam.mead@imm.ox.ac.uk](mailto:adam.mead@imm.ox.ac.uk)

\$ Present address: Department of Cancer Biology, Dana Farber Cancer Institute, Boston, MA, USA and Broad Institute, Cambridge, MA, USA.

## Summary

***TP53* is the most commonly mutated gene in human cancer, typically occurring in association with complex cytogenetics and dismal outcomes. Understanding the genetic and non-genetic determinants of *TP53*-mutation driven clonal evolution and subsequent transformation is a crucial step towards the design of rational therapeutic strategies. Here, we carry out allelic resolution single-cell multi-omic analysis of haematopoietic stem/progenitor cells (HSPCs) from patients with a myeloproliferative neoplasm who transform to *TP53*-mutant secondary acute myeloid leukaemia (AML), a tractable model of *TP53*-mutant cancer evolution. All patients showed dominant *TP53* 'multi-hit' HSPC clones at transformation, with a leukaemia stem cell transcriptional signature strongly predictive of adverse outcome in independent cohorts, across both *TP53*-mutant and wild-type AML. Through analysis of serial samples, antecedent *TP53*-heterozygous clones and *in vivo* perturbations, we demonstrate a hitherto unrecognised effect of chronic inflammation, which suppressed *TP53* wild-type HSPCs whilst enhancing the fitness advantage of *TP53*-mutant cells and promoting genetic evolution. Our findings will facilitate the development of risk-stratification, early detection and treatment strategies for *TP53*-mutant leukaemia, and is of broad relevance to other cancer types.**

## Main Text

Tumour protein 53 (*TP53*) is the most frequently mutated gene in human cancer, typically occurring as a multi-hit process with a point mutation in of one allele and loss of the other wild-type allele<sup>1,2</sup>. *TP53* mutations are also strongly associated with copy number alterations (CNA) and structural variants, reflecting the role of p53 in the maintenance of genomic integrity<sup>2,3</sup>. In myeloid malignancies, presence of a *TP53* mutation defines a distinct clinical entity<sup>1</sup>, associated with complex CNA, lack of response to conventional therapy and dismal outcomes<sup>2,4,5</sup>. Understanding the mechanisms by which *TP53* mutations drive clonal evolution and disease progression is a crucial step towards the development of rational strategies to diagnose, stratify, treat and potentially prevent this condition.

Myeloproliferative neoplasms (MPN) arise in haematopoietic stem cells (HSC) through the acquisition of mutations in JAK/STAT signalling pathway genes (*JAK2*, *CALR* or *MPL*), leading to aberrant proliferation of myeloid lineages<sup>6</sup>. Progression to secondary acute myeloid leukaemia (sAML) occurs in 10-20% of MPN and is characterized by cytopenias, increased myeloid blasts, acquisition of aberrant leukaemia stem cell (LSC) properties by haematopoietic stem/progenitor cells (HSPC) and median survival of less than one year<sup>7,8</sup>. *TP53* mutations are detected in approximately 20-35% of post-MPN sAML<sup>9-11</sup> (collectively termed *TP53*-sAML), often in association with loss of the remaining wild-type allele<sup>12</sup> and multiple CNAs<sup>13</sup>. Furthermore, deletion of *Trp53* combined with *JAK2V617F* mutation leads to a highly penetrant myeloid leukaemia in mice<sup>11,14</sup>.

Notwithstanding the established role of *TP53* mutation in MPN transformation, *TP53*-mutant subclones are also present in 16% of chronic phase MPN (CP-MPN) and in most cases this does not herald the development of *TP53*-sAML<sup>15</sup>. However, little is known about the additional genetic and non-genetic determinants of clonal evolution following the acquisition of a *TP53* mutation. Resolving this question requires unravelling multiple layers of intratumoural heterogeneity, including reliable identification of the *TP53* mutation, loss of the wild-type allele and presence of CNA. Integrating this mutational landscape with cellular phenotype and transcriptional signatures will resolve aberrant haematopoietic differentiation and molecular

35 properties of LSC in *TP53*-sAML. This collectively requires single-cell approaches which combine molecular and phenotypic analysis of HSPCs with allelic-resolution mutation detection, an approach recently enabled by the TARGET-seq technology<sup>16</sup>.

### **Convergent clonal evolution during *TP53*-driven leukaemic transformation**

40 To characterize the genetic landscape of *TP53*-sAML, we analysed 33 *TP53*-sAML patients (Table S1) through bulk-level targeted next generation sequencing and SNP array (Extended Data Fig.1). We detected MPN-driver mutations (*JAK2*, *CALR*) in 28 patients (85%), and co-occurring myeloid driver mutations in 24 patients (73%). Multiple *TP53* mutations were present in one third (n=11) of patients, including 2  
45 patients with 3 *TP53* mutations. 82% (18/22) of patients with a single *TP53* mutation showed a high variant allelic frequency (VAF) of >50%. CNAs were present in all patients analysed, and 87% (20/23) had a complex karyotype ( $\geq 3$  CNA; Extended Data Fig.1a-g). Deletion or copy neutral loss of heterozygosity affecting the *TP53* locus (chr17p13.1) was detectable at the bulk level in 43% of patients (10/23) (Extended  
50 Data Fig.1b-d). Taken together, these findings support that *TP53*-sAML is associated with complex genetic intratumoural heterogeneity.

To characterize tumour phylogenies and subclonal structures, we performed TARGET-seq analysis<sup>16</sup>, a technology that allows allelic-resolution genotyping, whole  
55 transcriptome and immunophenotypic analysis from the same single-cell, on 17517 Lin<sup>-</sup>CD34<sup>+</sup> HSPCs from 14 *TP53*-sAML patients (Extended Data Fig.1a), 9 age-matched healthy donors (HD) and 8 previously published myelofibrosis (MF) patients (Fig.1a, gating strategy shown in Extended Data Fig.2a). HSPCs wild-type for all mutations analyzed were present in 10 of 14 patients (Extended Data Fig.2b-o),  
60 providing a valuable population of cells for intra-patient comparison with mutation-positive cells<sup>17</sup>. In all cases, the dominant clone showed loss of wild-type *TP53* through 4 patterns of clonal evolution: (1) biallelic *TP53* mutations by acquisition of a second mutation on the other *TP53* allele, (2) hemizygous *TP53* mutations (deleted *TP53* wild-type allele), (3) parallel evolution with 2 clones harbouring different *TP53* alterations,  
65 (4) a *JAK2* negative dominant clone with biallelic *TP53* mutations in patients with previous *JAK2*-mutant MPN<sup>18</sup> (Fig.1b-e, Extended Data Fig.2b-o). Biallelic mutations were confirmed by single molecule cloning and computational analysis (Extended Data Fig.1h-j). Integration of index-sorting data revealed that dominant *TP53* multi-hit clones

were enriched in progenitor populations as previously described in *de novo* AML<sup>19</sup>,  
70 whereas *TP53*-mutant cells were less frequent in the HSC compartment (Extended  
Data Fig.3a). CNA analysis using single-cell transcriptomes showed that all *TP53*  
multi-hit clones harboured at least one highly clonally-dominant CNA, with very few  
*TP53*-mutant cells without evidence of a CNA ( $3.4 \pm 1.2\%$ ) and an additional 5/14 (36%)  
75 patients also showing cytogenetically-distinct subclones (Fig.1f,g, Extended Data  
Fig.2p,q).

To confirm that dominant HSPC clones were functional LSCs, we established patient-  
derived xenografts (PDX) for 2 *TP53*-sAML patients (Fig.1h). Mice developed  
leukaemia in 27-31 weeks with high numbers of human CD34<sup>+</sup> myeloid blast cells in  
80 the bone marrow (BM) (Extended Data Fig.3b-d), with a progenitor phenotype, *TP53*  
mutations and CNAs similar to the dominant clone from patients' primary cells (Fig. 1i,  
Extended Data Fig.3e-l). In Patient IF0131, a monosomy 7 subclone (Fig.1f)  
preferentially expanded in PDX models (Fig.1i). Monosomy 7 cells showed a distinct  
transcriptional profile with increased WNT, RAS, MAPK signalling and cell cycle  
85 associated transcription (Extended Data Fig.3m,n). Together, these data are  
compatible with a fitness advantage of monosomy 7 cells, a recurrent event in *TP53*-  
sAML (Extended Data Fig.1b,c), driven by activation of signalling pathways which may  
relate to deletion of chromosome 7 genes such as *EZH2*<sup>20</sup>. In summary, the dominant  
leukaemic clones in *TP53*-sAML were invariably characterized by multiple hits affecting  
90 *TP53* (multi-hit state), indicating strong selective pressure for complete loss of wild-  
type *TP53*, together with gain of CNAs and complex cytogenetic evolution, with very  
few *TP53* multi-hit cells with a normal karyotype (Fig.1j).

### **Molecular signatures of *TP53*-mutant mediated transformation**

95 To understand the cellular and molecular framework through which *TP53* mutation  
drives clonal evolution, we next analysed single-cell RNA-seq data from 10459 *TP53*-  
sAML HSPCs alongside 2056 MF and 5002 HD HSPCs passing quality control. Force-  
directed graph analysis revealed separate clustering of *TP53*-mutant HSPC in  
comparison with HD and MF cells, with a high degree of inter-patient heterogeneity  
100 (Extended Data Fig.4a) as observed in other haematopoietic malignancies<sup>21</sup>. This  
could potentially be explained by patient-specific cooperating mutations and  
cytogenetic alterations (Extended Data Fig.1). TARGET-seq analysis uniquely enabled

comparison of *TP53* multi-hit HSPC to *TP53* wild-type preleukaemic stem cells (“preLSCs”) from the same *TP53*-sAML patients as well as HD and MF, to derive a specific *TP53* multi-hit signature including known p53-pathway genes (Extended Data Fig.4b,c).

Integration of single cell transcriptomes and diffusion map analysis of HSPCs from *TP53*-sAML patients showed that *TP53* multi-hit HSPCs clustered separately from *TP53* wild-type preLSCs in two distinct populations with enrichment of LSC and erythroid-associated transcription respectively (Fig.2a, Table S3), and a differentiation trajectory towards the erythroid-biased population (Fig.2b), an unexpected finding given that erythroleukaemia is uncommon in *TP53*-sAML<sup>22,23</sup>. Sorted CD34<sup>+</sup> *TP53*-multi-hit cells exhibited potential for erythroid differentiation *in vivo* and *in vitro*, supporting that this occurs downstream of the LSC population (Extended Data Fig.5a-c). *TP53* multi-hit LSCs showed enrichment of cell cycle, inflammatory, signalling pathways and LSC associated transcription, whereas *TP53* multi-hit erythroid cells were depleted of the latter (Extended Data Fig.4d).

To further explore this erythroid-biased population, we projected *TP53* multi-hit cells onto a previously published healthy donor haematopoietic hierarchy<sup>24</sup>. *TP53*-sAML differed from *de novo* AML with an enrichment into HSC and early erythroid populations, whereas *de novo* AML were enriched in myeloid progenitors (Fig.2c,d)<sup>25</sup>. A similar enrichment was observed for *TP53* multi-hit cells when mapped on a Lin<sup>-</sup> CD34<sup>+</sup> MF cellular hierarchy (Extended Data Fig.5d,e), with erythroid-biased populations being highly enriched in immunophenotypically defined MEPs (Extended Data Fig.5f). Taken together, these findings support an aberrant erythroid-biased differentiation trajectory in *TP53*-sAML.

To determine whether upregulation of erythroid-associated transcription was a more widespread phenomenon in *TP53*-mutant AML, we investigated erythroid-myeloid associated transcription in the BeatAML and TCGA cohorts<sup>26,27</sup>. Erythroid scores were increased in *TP53* mutant compared to *TP53* wild-type AML, whereas there was no significant difference in myeloid scores (Fig.2e-f, Extended Data Fig.5g-j, scores described in Table S3). Concomitantly, patients with high erythroid scores also showed decreased *TP53*-target gene expression (Extended Data Fig.5k). We next investigated

the expression of key transcription factors for erythroid/granulomonocytic commitment and found increased *GATA1* expression in Lin<sup>-</sup>CD34<sup>+</sup> *TP53* multi-hit HSPCs, whereas *CEBPA* was only expressed at low levels (Fig.2g). Analysis of the BeatAML cohort revealed increased *GATA1* and reduced *CEBPA* expression in association with *TP53* mutation (Extended Data Fig.5l), with consequent reduction in the *CEBPA/GATA1* expression ratio (Fig.2h). Similar findings were observed in *TP53* knock-out or mutant isogenic MOLM13 cell lines (Extended Data Fig.5m)<sup>28</sup>. These observations suggest that the *CEBPA/GATA1* expression ratio, an important transcription factor balance which affects erythroid versus myeloid differentiation in leukaemia<sup>29,30</sup> is disrupted by *TP53* mutation.

To determine whether p53 directly influences myeloid-erythroid differentiation, we knocked-down *TP53* in *JAK2V617F* CD34<sup>+</sup> cells from MPN patients (Extended Data Fig.5n). *TP53* knock-down led to increased erythroid (CD71<sup>+</sup>CD235a<sup>+</sup>) and decreased myeloid (CD14<sup>+</sup>/CD15<sup>+</sup>/CD11b<sup>+</sup>) differentiation *in vitro* (Fig.2i) and consequently decreased *CEBPA/GATA1* expression ratio (Fig.2j), suggesting that p53 may directly contribute to the aberrant myelo-erythroid differentiation observed.

As 'stemness scores' have previously been applied to determine prognosis in AML<sup>31</sup>, we next asked whether a single-cell defined *TP53* multi-hit LSC signature might identify AML patients with adverse outcomes. Single cell multi-omics allowed us to derive a 44-gene "p53LSC-signature" (Table S4) by comparing gene expression of HD, *JAK2*-mutant MF HSPC and *TP53* wild-type preLSC to transcriptionally-defined *TP53*-mutant LSCs (Fig.2a,k). High p53LSC-signature score (Extended Data Fig.6a,b) was strongly associated with *TP53* mutation status, although some *TP53* wild-type patients also showed a high p53LSC score. A high p53LSC score predicted for poor survival in the independent BeatAML and TCGA cohorts, irrespective of *TP53* mutational status (Fig.2l, Extended Data Fig.6c-e). The p53LSC signature performed well as a predictor of survival, including in sAML patients, as compared to the previously published LSC17 score<sup>31</sup> and p53-mutant score generated using all *TP53*-mutant HSPC rather than LSCs (Extended Data Fig.6f-g, TableS4), providing a powerful new tool to aid risk stratification in AML.



170 **Preleukaemic *TP53*-wild-type cells display self-renewal and differentiation defects**

TARGET-seq uniquely enabled phenotypic and molecular characterization of rare *TP53* wild-type cells, referred to as preLSCs, which include both residual HSPCs that were wild-type for all mutations analyzed, as well as HSPCs which form part of the antecedent MPN clone. These preLSCs were obtained in sufficient numbers (>20 cells) from 9 of 14 *TP53*-sAML patients, including all patterns of clonal evolution (Fig.3a and Extended Data Fig.7a). PreLSCs representing the antecedent MPN clone (positive for MPN-associated driver mutations) were more frequent (60.5%) than preLSCs that were wild-type for all mutations (39.5%). PreLSCs were enriched in HSC-associated genes, and mapped onto HSC clusters in HD and MF haematopoietic hierarchies (Fig.3a,b). Index sorting revealed that preLSCs were strikingly enriched in the phenotypic HSC compartment, unlike *TP53* multi-hit HSPCs (Fig.3c, Extended Data Fig.3a). Pre-LSCs were rare, as reflected by a reduction in the numbers of phenotypic HSCs present within the Lin<sup>-</sup>CD34<sup>+</sup> HSPC compartment in *TP53*-sAML compared to HD (Extended Data Fig.7b).

We reasoned that the HSC phenotype of preLSCs, with reduced frequency in progenitor compartments, might reflect impaired differentiation. To explore this hypothesis, we carried out scVelo analysis, which showed absence of a transcriptional differentiation trajectory in preLSCs, unlike HD HSCs (Fig.3d). PreLSCs showed increased expression of haematopoietic stem cell and Wnt  $\beta$ -catenin genes and decreased cell cycle genes as compared to HD and MF cells (Fig.3e-g, TableS3). To functionally confirm these findings, we sorted phenotypic HSCs (to purify preLSCs), as well as other progenitor cells, from HD, MF and *TP53*-sAML patients for long term culture initiating cell (LTC-IC) and short-term cultures (Fig.3h; Extended Data Fig.7c). PreLSC LTC-IC activity was similar to HD and increased compared to MF, with preserved terminal differentiation capacity and confirmed *TP53* wild-type genotype (Fig.3i, Extended Data Fig.7d-g). In short-term liquid culture, preLSCs showed reduced clonogenicity, with retained CD34 expression and decreased proliferation (Fig.3j, Extended Data Fig.7h-i). In summary, we identified rare and phenotypically distinct preLSCs from *TP53*-sAML samples which were characterized by differentiation defects and distinct stemness, self-renewal and quiescence signatures. As these cells were *TP53*-wild-type, and showed normal differentiation after prolonged *ex vivo* culture, we

reasoned that these functional and molecular abnormalities are likely to be cell-  
205 extrinsically mediated. Indeed, preLSCs showed enrichment of gene-signatures  
associated with certain cell-extrinsic inflammatory mediators (TNF $\alpha$ , IFN $\gamma$ , TGF $\beta$ , IL2)  
(Fig. 3k).

### Inflammation promotes *TP53*-associated clonal dominance

210

To understand the transcriptional signatures associated with leukaemic progression  
we analysed samples from 5 CP-MPN patients who subsequently developed *TP53*-  
sAML (“pre-*TP53*-sAML”) alongside 6 CP-MPN patients harbouring *TP53* mutated  
clones who remained in chronic phase (“CP *TP53*-MPN”, median 4.43 years [2.62-  
215 5.94] of follow-up, Fig.4a, Extended Data Fig.8). Compared to *TP53*-sAML samples,  
CP *TP53*-MPN had a lower VAF and number of *TP53* mutations (Extended Data  
Fig.8a-d). The type, distribution and pathogenicity score of *TP53* mutations were  
similar between chronic and acute stages (Extended Data Fig.8e,f). All 5 pre-*TP53*-  
sAML samples and 4 of the 6 CP *TP53*-MPN were then analysed by TARGET-seq  
220 (Fig.4a). HSPC immunophenotype was similar for pre-*TP53*-sAML and CP *TP53*-MPN  
patients (Extended Data Fig.9a-c), and clearly distinct from the *TP53*-sAML stage  
(Extended Data Fig.9d). Heterozygous *TP53* clones were identified in 3 pre-*TP53*-  
sAML patients and all 4 CP *TP53*-MPN (Fig.4b, Extended Data Fig.9e-m). A minor  
homozygous *TP53* mutated clone initially present in one CP *TP53*-MPN patient was  
225 undetectable after 4 years (Extended Data Fig.9h). As *TP53*-heterozygous mutant  
HSPCs represent the direct genetic ancestors of *TP53* “multi-hit” LSCs, we compared  
gene expression of heterozygous *TP53* mutant HSPC from pre-*TP53*-sAML (n=296)  
to CP *TP53*-MPN (n=273) (Fig.4b, blue boxes) to identify putative mediators of  
transformation. *TP53*-heterozygous HSPC from pre-*TP53*-sAML patients showed  
230 downregulation of TNF $\alpha$  and TGF $\beta$  associated gene signatures, both of which are  
known to be associated with HSC attrition<sup>32,33</sup>, with upregulated expression of  
oxidative phosphorylation, DNA repair and interferon response genes (Table S5,  
Fig.4c-e), without changes in IFN receptor expression levels or concurrent interferon  
treatment (Extended Data Fig.9n, Table S1). Upregulation of inflammatory signatures  
235 was detected in *TP53*-homozygous cells from the same pre-*TP53*-sAML patients at a  
higher level than in *TP53*-heterozygous cells (Extended Data Fig.9o). Collectively,

these findings raise the possibility that inflammation might contribute to preleukaemic clonal evolution towards *TP53*-sAML.

240 To evaluate the role of inflammation in *TP53*-driven leukaemia progression, we performed competitive mouse transplantation experiments between CD45.1+ Vav-Cre *Trp53*<sup>R172H/+</sup> and CD45.2+ *Trp53*<sup>+/+</sup> BM cells followed by repeated poly(I:C) or LPS intraperitoneal injections, recapitulating chronic inflammation through induction of multiple pro-inflammatory cytokines<sup>34,35</sup> known to be increased in the serum of patients with MPN<sup>36</sup>, including IFN $\gamma$  (Fig.5a, Extended Data Fig.10a). *Trp53* mutant peripheral  
245 blood (PB) myeloid cells, BM HSCs (Lin<sup>-</sup>Sca1<sup>+</sup>c-Kit<sup>+</sup>CD150<sup>+</sup>CD48<sup>-</sup>) and LSKs (Lin<sup>-</sup>Sca1<sup>+</sup>c-Kit<sup>+</sup>) were selectively enriched upon poly(I:C) treatment (Fig.5b,c, Extended Data Fig.10b-e). Crucially, the fitness advantage of *Trp53* mutant HSCs and LSKs was exerted both through an increase in numbers of *Trp53*<sup>R172H/+</sup> HSPCs and reduction in  
250 numbers of wild-type competitors (Fig.5d,e, Extended Data Fig.10f,g). Treatment of chimeric mice with LPS (Fig.5a), which induces an inflammatory response mediated through release of IL1 $\beta$  and IL6<sup>37</sup>, amongst others, lead to a similar increase in the number of *Trp53* mutant PB myeloid cells and LSKs (Fig.5f,g). These results indicate that a variety of inflammatory stimuli can promote expansion of the *Trp53* mutant clone.

255 To determine how inflammation might alter haematopoietic differentiation and exert a selective pressure to drive the expansion of the *Trp53* mutant clone, we established an inducible SCL-CreER<sup>T</sup> *Trp53*<sup>R172H/+</sup> mouse model (Fig.5h). Poly(I:C) treatment led to inflammation-associated changes in blood cell parameters, including anaemia, leucopenia and thrombocytopenia (Extended Data Fig.10h-j). Similarly to the Vav-Cre  
260 model, poly(I:C) treatment promoted the selection of myeloid *Trp53* mutant cells in the PB (Extended Data Fig.10k), with a myeloid bias induced by the inflammatory stimulus in PB leucocytes specifically associated with *Trp53*-mutation (Fig.5i,j, Extended Data Fig.10l). Analysis of HSPCs showed the expected selection for *Trp53* mutant HSCs and LSKs following Poly(I:C) treatment (Extended Data Fig.10m). Numbers of wild-  
265 type competitor erythroid progenitors were reduced upon poly(I:C) treatment as expected<sup>38</sup>, whereas *Trp53*-mutation was associated with an increase in erythroid progenitors that was not impacted by inflammation (Fig.5k, Extended Data Fig.10n) in line with the erythroid bias detected in patient samples. Finally, to determine the

270 mechanisms by which inflammation might promote a fitness advantage for *Trp53*  
mutated cells, we performed cell cycle and apoptosis analysis following chronic  
poly(I:C) treatment. Cell cycle was similarly increased in poly(I:C)-treated WT and  
*Trp53*-mutated LSKs<sup>39</sup>; however, *Trp53*-mutated LSKs were resistant to inflammation-  
induced apoptosis<sup>40</sup> in comparison with their WT counterparts (Fig.5l,m).

275

### **Inflammation promotes genetic evolution of *Trp53* mutant HSPC**

As exit from dormancy promotes DNA-damage-induced HSPCs attrition<sup>41</sup>, we  
reasoned that *Trp53* mutation might rescue HSPCs that acquire DNA damage (and  
would otherwise undergo apoptosis) driven by chronic inflammation-associated  
280 proliferative stress. To explore this possibility, we carried out M-FISH karyotype  
analysis of *Trp53*<sup>+/+</sup> LSKs expanded *in vitro* from mice following poly(I:C) treatment and  
*Trp53*<sup>R172H/+</sup> LSKs from mice with or without poly(I:C) treatment. Wild-type competitor  
LSK-derived cells from poly(I:C) treated mice were karyotypically normal. In contrast,  
we observed a striking increase in the frequency and number of karyotypic  
285 abnormalities in *Trp53* mutated LSK-derived cells upon poly(I:C) treatment (Fig.6a-d).  
Collectively, these results support a model whereby chronic inflammation promotes the  
survival and genetic evolution of *TP53* mutated cells whilst suppressing wild-type  
haematopoiesis, ultimately promoting clonal expansion of *TP53* mutant HSPCs  
(Fig.6e).

290

### **Discussion**

Here, we unravel multi-layered genetic, cellular and molecular intratumoural  
heterogeneity in *TP53* mutation driven disease transformation through single-cell multi-  
omic analysis. Allelic resolution genotyping of leukaemic HSPCs revealed a strong  
295 selective pressure for gain of *TP53* missense mutation, loss of the *TP53* wild-type allele  
and acquisition of complex CNAs, including cases with parallel genetic evolution during  
*TP53*-sAML LSC expansion. Despite the known dominant negative and/or gain of  
function effect of certain *TP53* mutations<sup>28,42</sup>, loss of the *TP53* wild-type allele, a  
genetic event associated with a particularly dismal prognosis<sup>2</sup>, confers additional  
300 fitness advantage to *TP53*-sAML LSCs. As CNA were universally present in *TP53*-  
sAML with a very high clonal burden, it is not possible, even with high-resolution single-  
cell analyses, to disentangle the impact of *TP53*-multi hit mutation versus the effects  
of patient-specific CNA which were inextricably linked in all patients analysed.

305 Three distinct clusters of HSPCs were identified in *TP53*-sAML, including one  
characterized by overexpression of erythroid genes, of particular note as  
erythroleukaemia is a rare entity, associated with adverse outcome and *TP53*  
mutation<sup>43,44</sup>. Analysis of a large AML cohort also revealed overexpression of erythroid  
genes as a more widespread phenomenon in *TP53* mutant AML, with disrupted  
310 balance of *GATA1* and *CEBPA* expression. Notably, *CEBPA* knockout or mutation is  
reported to cause a myeloid to erythroid lineage switch with increased expression of  
*GATA1*<sup>29,30</sup> and, in addition, *GATA1* associates with and inhibits p53<sup>45</sup>. Importantly, our  
data do not distinguish whether this lineage-switch is primarily an instructive versus  
permissive effect of *TP53*-mutation<sup>46</sup>. A second ‘*TP53*-sAML LSC’ cluster allowed us  
315 to establish a novel p53LSC-signature, which we demonstrated to be highly relevant  
to predict outcome in AML, independently of *TP53* status. This powerful approach  
could be more broadly applied in cancer, whereby single multi-omic cell derived gene  
scores can be used to stratify larger patient cohorts using bulk gene expression data.

320 A third *TP53* wild-type ‘preLSC’ HSPC cluster was characterized by quiescence  
signatures and defective differentiation, reflecting the impaired haematopoiesis  
observed in patients with *TP53*-sAML. Through integration of single cell multi-omic  
analysis with *in vitro* and *in vivo* functional assays we show that *TP53*-wild-type  
preLSCs are cell-extrinsically suppressed whilst chronic inflammation promotes the  
325 fitness advantage of *TP53* mutant cells, ultimately leading to clonal selection (Fig.6e).  
Inflammation is a cardinal regulator of HSC function with many effects on HSC fate  
and function<sup>47</sup>, including proliferation-induced DNA-damage and depletion of HSCs<sup>41</sup>.  
There is emerging evidence that clonal HSCs can become inflammation-adapted<sup>47-49</sup>  
and by altering the response to inflammatory challenges, mutations can thus confer a  
330 fitness advantage to HSCs. Here, we demonstrate a hitherto unrecognized effect of  
*TP53* mutations, which conferred a marked fitness advantage to HSPC in the presence  
of chronic inflammation induced with both poly(I:C) as well as LPS. We provide  
evidence that *TP53* mutant HSCs showing dysregulated inflammation-associated gene  
expression are enriched in patients who will develop *TP53*-sAML. We propose that  
335 HSCs that would otherwise undergo inflammation-associated and DNA-damage-  
induced attrition, are rescued by *TP53* mutation, ultimately leading to the accumulation  
of HSCs which have acquired DNA damage, thus promoting genetic evolution that

underlies disease progression. This hypothesis was strongly supported through *in vivo* experiments in which inflammation promoted genetic evolution of *Trp53* mutant mouse HSPCs. Further studies are required to characterize the key inflammatory mediators and molecular mechanisms involved, which we believe are unlikely to be restricted to a single axis, with a myriad of inflammatory mediators overexpressed in MPN<sup>50</sup>. Furthermore, loss of the wild-type *Trp53* allele confers an additional fitness advantage to *Trp53* mutant HSPC following DNA-damage as previously described<sup>28</sup>, providing an explanation for the selection for multi-hit *TP53* mutant clones observed in patients. Consequently, we believe that approaches which target the inflammatory state, rather than a specific cytokine, are likely to be required to restrain disease progression, as reported for bromodomain inhibitors<sup>51</sup>. Collectively, our findings provide a crucial conceptual advance relating to the interplay between genetic and non-genetic determinants of *TP53*-mutation associated disease transformation. This will facilitate the development of early detection and treatment strategies for *TP53*-mutant leukaemia. Since *TP53* is the most commonly mutated gene in human cancer<sup>3,52</sup>, we anticipate that these findings will be of broader relevance to other cancer types.

355

360

365

370

## Methods

### Banking and processing of human samples

Primary human samples (peripheral blood or bone marrow, described in Table S1) were analysed with approvals from the Inserm Institutional Review Board Ethical Committee (project C19-73, agreement 21-794, CODECOH n°DC-2020-4324); and  
375 from the INForMeD Study (REC: 199833, 26 July 2016, University of Oxford). Patients and normal donors provided written informed consent in accordance with the Declaration of Helsinki for sample collection and use in research. For secondary AML patients, we specifically selected samples from patients with known *TP53*-mutation.

380 Cells were subjected to Ficoll gradient centrifugation and for some samples, CD34 enrichment was performed using immunomagnetic beads (Miltenyi). Total mononuclear cells (MNCs) or CD34<sup>+</sup> cells were frozen in FBS supplemented with 10% DMSO for further analysis.

### Targeted bulk sequencing

385 Bulk genomic DNA from patient samples' mononuclear or CD34<sup>+</sup> cells was isolated using DNeasy Blood & Tissue Kit (Qiagen) or QIAamp DNA Mini Kit (Qiagen) as per manufacturer's instructions. Targeted sequencing was performed using a TruSeq Custom Amplicon panel (Illumina) or a Haloplex Target Enrichment System (Agilent technologies) with amplicons designed around 32, 44 or 77 genes<sup>53</sup>. Targets were  
390 chosen based on the genes/exons most frequently mutated and/or likely to alter clinical practice (diagnostic, prognostic, predictive or monitoring capacity) across a range of myeloid malignancies (e.g. MDS/AML/MPN). Targets covered in all panels include *ASXL1*, *CALR*, *CBL*, *CEBPA*, *CSF3R*, *DNMT3A*, *EZH2*, *FLT3*, *HRAS*, *IDH1*, *IDH2*, *JAK2*, *KIT*, *KRAS*, *MPL*, *NPM1*, *NRAS*, *PHF6*, *RUNX1*, *SETBP1*, *SF3B1*, *SRSF2*,  
395 *TET2*, *TP53*, *U2AF1*, *WT1*, *ZRSR2*. Sequencing was performed with a MiSeq sequencer (Illumina), according to the manufacturer's protocols. Results were analysed after alignment of the reads using two dedicated pipelines, SOPHiA DDM<sup>®</sup> (Sophia Genetics) and an in-house software GRIIO-Dx<sup>®</sup>. For all samples, an average depth exceeding 200X for > 90% of the target regions was required, or as previously  
400 described<sup>16</sup>. All pathogenic variants were manually checked using Integrative Genomics Viewer software. Analysis is presented in Extended Data Fig.1a and Extended Data Fig.8a.

Pathogenic scores for each *TP53* variant (Extended Data Fig.8e) were derived from COSMIC (Catalogue Of Somatic Mutations In Cancer) using the FATHMM-MKL  
405 algorithm. The FATHMM-MKL algorithm integrates functional annotations from ENCODE with nucleotide-based hidden Markov models to predict whether a somatic mutation is likely to have functional, molecular and phenotypic consequences. Scores greater than 0.7 indicate that a somatic mutation is likely pathogenic, whilst scores less than 0.5 indicate a neutral classification.

410 The type and location of *TP53* mutations from this study, *de novo* AML patients and CHIP individuals represented in Extended Data Fig.8f were generated using Pecan Portal<sup>54</sup>. *De novo* AML *TP53* mutations were downloaded from Papaemmanuil, *et al.*<sup>55</sup> and Ley, *et al.*<sup>27</sup>; CHIP associated *TP53* mutations were obtained from Coombs, *et al.*, Desai, *et al.*, Young, *et al.*<sup>56-58</sup>

#### 415 **Sanger sequencing of patient-associated mutations in PDX models**

Genomic DNA from PDX sorted populations (LMPP: hCD45<sup>+</sup>Lin<sup>-</sup>CD34<sup>+</sup>CD38<sup>-</sup>CD45RA<sup>+</sup>CD90<sup>-</sup> and GMP: hCD45<sup>+</sup>Lin<sup>-</sup>CD34<sup>+</sup>CD38<sup>+</sup>CD45RA<sup>+</sup>CD123<sup>+</sup>) was extracted using QIAamp DNA Mini Kit (Qiagen). Sanger sequencing was performed with forward or reverse primers (TableS6a) targeting mutations identified by targeted bulk  
420 sequencing in the corresponding primary samples using Mix2seq kit (Eurofins Genomics) and sequences were analysed with the ApE editor.

#### **Single Nucleotide Polymorphism Array sample preparation, Copy Number Variant and Loss of Heterozygosity Analysis**

Bulk genomic DNA from patients' mononuclear cells was isolated using DNeasy Blood  
425 & Tissue Kit (Qiagen) as per manufacturer's instructions. 250 ng of gDNA were used for hybridization on an Illumina Infinium OmniExpress v1.3 BeadChips platform.

To call mosaic copy number events in primary patient samples, genotyping intensity data generated was analysed using the Illumina Infinium OmniExpress v1.3 BeadChips platform. Haplotype phasing, calculation of log R ratio (LRR) and B-allele  
430 frequency (BAF) and calling of mosaic events was performed using Mocha (Mocha: A BCFtools extension to call mosaic chromosomal alterations starting from phased VCF files with either B Allele Frequency (BAF) and Log R Ratio (LRR) or allelic depth (AD)), as previously described<sup>59,60</sup>. In brief, Mocha comprises the following steps: (1) filtering



of constitutional duplications; (2) use of a parameterized hidden Markov model to  
435 evaluate the phased BAF for variants on a per-chromosome basis; (3) deploying a  
likelihood ratio test to call events; (4) defining event boundaries; (5) calling copy  
number; (6) estimating the cell fraction of mosaic events. A series of stringent filtering  
steps was applied to reduce the rate of false positive calls. To eliminate possible  
440 constitutional and germline duplications, excluding calls with  $\text{lod\_baf\_phase} < 10$ , those  
with length  $< 500\text{kbp}$  and  $\text{rel\_cov} > 2.5$ , and any gains with estimated cell fraction  $> 80\%$ ,  
 $\text{logR} > 0.5$  or length  $< 24\text{Mb}$ . Given that interstitial LOH are rare and likely artefactual, all  
LOH events  $< 8\text{Mb}$  were filtered<sup>59</sup>. Events on genomic regions reported to be prone to  
recurrent artefact<sup>59</sup> (chr6  $< 58\text{Mb}$ , chr7  $> 61\text{Mb}$ , and chr2  $> 50\text{Mb}$ ) were also filtered, and  
those where manual inspection demonstrated noise or sparsity in the array.

445

To find common genomic lesions on a focal and arm level, Infinium OmniExpress  
arrays were initially processed with Illumina Genome Studio v2.0.4. Following this, Log  
R Ratio (LRR) data was extracted for all probes and array annotation obtained from  
Illumina (InfiniumOmniExpress-24v1-3\_A1). LRR data was then smoothed and  
450 segmentation called using the CBS algorithm from the DNACopy<sup>61,62</sup> v1.60.0 package  
in R. A minimum number of 5 probes was required to call a segment, and segments  
where analysed using GenomicRanges<sup>63</sup> v1.38.0. Definitions of amplification, gain,  
loss and deletion events where as outlined in Bashton, *et al.*<sup>64</sup>. Segmentation data was  
then analysed in GISTIC<sup>65</sup> v2.023.

455 For PDX models, genomic DNA from sorted populations (LMPP: hCD45<sup>+</sup>Lin<sup>-</sup>  
CD34<sup>+</sup>CD38<sup>-</sup>CD45RA<sup>+</sup>CD90<sup>-</sup> and GMP: hCD45<sup>+</sup>Lin<sup>-</sup>CD34<sup>+</sup>CD38<sup>+</sup>CD45RA<sup>+</sup>CD123<sup>+</sup>)  
was extracted using QIAamp DNA Mini Kit (Qiagen). SNP-CGH array hybridization  
was performed using the Affymetrix Cytoscan® HD (Thermo Fisher Scientific)  
according to the manufacturer's recommendations. DNA amplification was checked  
460 using BioSpec-nano<sup>TM</sup> spectrophotometer (Shimadzu) with expected concentrations  
between 2,500 and 3,400ng/ $\mu\text{L}$ . DNA length distribution post-fragmentation was  
checked using D1000 ScreenTapes on TapeStation 4200 instrument (Agilent  
Technologies). Cytoscan HD array includes 2.6 million markers combining SNP and  
non-polymorphic probes for copy number evaluation. Raw data CEL files were  
465 analysed using the Chromosome Analysis Suite software package (v4.1, Affymetrix)  
with genome version GRCh37 (hg19) only if achieving the manufacturer's quality cut-

offs. Only CNAs > 10kb were reported in the analysis presented in Extended Data Fig.3k,l.

### **Single-molecule cloning and sequencing of patient-derived cDNA**

470 To experimentally verify the biallelic nature of *TP53* mutations in *TP53*-sAML patients, cDNA from a selected patient with putative *TP53* biallelic status (Patient ID GR004) was PCR-amplified using cDNA-specific primers spanning both *TP53* mutations (Fwd: 5'-GACCCTTTTTGGACTTCAGGTG-3', Rev: 5'-CCATGAGCGCTGCTCAGATAG-3'). PCR amplification was performed with KAPA 2X Ready Mix (Roche), a Taq-derived  
475 enzyme with A-tailing activity, for direct cloning into a TA vector (pCR2.1 TOPO vector, TOPO® TA Cloning® Kit, Invitrogen) as per manufacturer's instructions. Sanger sequencing for 10 different colonies was performed using M13 forward and reverse primers; a representative example is shown in Extended Data Fig.1h.

### **Fluorescent activated cell sorting (FACS) and single-cell isolation**

480 Single cell FACS-sorting was performed as previously described<sup>16</sup>, using BD Fusion I and BD Fusion II instruments (Becton Dickinson) for 96-well plate experiments or bulk sorting experiments, and SH800S or MA900 (SONY) for 384-well plate experiments. Experiments involving isolation of human haematopoietic stem and progenitor cells (HSPCs) included single colour stained controls (CompBeads, BD Biosciences) and  
485 Fluorescence Minus One controls (FMOs). Antibodies used for HSPC staining are detailed in TableS7a (Panel A or B).

Briefly, single cells directly sorted into 384-well plates containing 2.07 µL of TARGET-seq lysis buffer<sup>66</sup>. Lineage<sup>-</sup>CD34<sup>+</sup> cells were indexed for CD38, CD90, CD45RA, CD123 and CD117 markers, which allowed us to record the fluorescence levels of  
490 each marker for each single cell. 7- aminoactinomycin D (7-AAD) was used for dead cell exclusion. Flow cytometry profiles of the human HSPC compartment (Extended Data Fig.2, Fig.9) were analysed using FlowJo software (version 10.1, BD Biosciences).

### **Single-cell TARGET-seq cDNA synthesis.**

495 RT and PCR steps were performed as previously described<sup>66</sup>, using 24 cycles of PCR amplification. Target-specific primers spanning patient-specific mutations were added to RT and PCR steps (TableS6a). After cDNA synthesis, cDNA from up to 384 single-

cell libraries was pooled, purified using Ampure XP Beads (0.6:1 beads to cDNA ratio; Beckman Coulter) and resuspended in a final volume of 50  $\mu$ L of EB buffer (Qiagen).  
500 The quality of cDNA traces was checked using a High Sensitivity DNA Kit in a Bioanalyzer instrument (Agilent Technologies).

### **Whole transcriptome library preparation and sequencing**

Pooled and bead-purified cDNA libraries were diluted to 0.2 ng/ $\mu$ L and used for  
tagmentation-based library preparation using a custom P5 primer and 14 cycles of  
505 PCR amplification<sup>66</sup>. Each indexed library was purified twice with Ampure XP beads  
(0.7:1 beads to cDNA ratio), quantified using Qubit dsDNA HS Assay Kit (Invitrogen,  
Cat# Q32854) and diluted to 4 nM. Libraries were sequenced on a HiSeq4000, HiSeqX  
or NextSeq instrument using a custom sequencing primer for read1 (P5\_seq:  
GCCTGTCCGCGGAAGCAGT GGTATCAACGCAGAGTTGC\*T, PAGE purified) with  
510 the following sequencing configuration: 15 bp R1; 8 bp index read; 69 bp R2 (NextSeq)  
or 150 bp R1; 8 bp index read; 150 bp R2 (HiSeq).

### **TARGET-seq single-cell genotyping**

After RT-PCR, cDNA+amplicon mix was diluted 1:2 by adding 6.25  $\mu$ L of  
DNase/RNase free water to each well of each 384-well plate. Subsequently, a 1.5  $\mu$ L  
515 aliquot from each single cell derived library was used as input to generate a targeted  
and Illumina-compatible library for single cell genotyping<sup>66</sup>. In the first PCR step, target-  
specific primers containing a plate-specific barcode (TableS6b) were used to amplify  
the target regions of interest. In a subsequent PCR step, Illumina compatible adaptors  
(PE1/PE2) containing single-direction indexes (Access Array™ Barcode Library for  
520 Illumina® Sequencers-384, Single Direction, Fluidigm) were attached to pre-amplified  
amplicons, generating single-cell barcoded libraries. Amplicons from up to 3,072  
libraries were pooled and purified with Ampure XP beads (0.8:1 ratio beads to product;  
Beckman Coulter). These steps were performed using Biomek FxP (Beckman Coulter),  
Mosquito (TTP Labtech) and VIAFLO 96/384 (INTEGRA Biosciences) liquid handling  
525 platforms. Purified pools were quantified using Qubit dsDNA HS Assay Kit (Invitrogen,  
Cat# Q32854) and diluted to a final concentration of 4 nM. Libraries were sequenced  
on a MiSeq or NextSeq instrument using custom sequencing primers as previously  
described<sup>66</sup> with the following sequencing configuration: 150 bp R1; 10 bp index read;  
150 bp R2.

## 530 Targeted single-cell genotyping analysis

### *Data pre-processing*

For each cell, the FASTQ file containing both targeted gDNA and cDNA-derived sequencing reads was aligned to the human reference genome (GRCh37/hg19) using Burrow-Wheeler Aligner (BWA v0.7.17)<sup>31</sup> and STAR (v2.6.1d)<sup>67</sup>. Custom perl scripts  
535 were used to demultiplex the gDNA and mRNA reads in the BAM file into separate SAM files based on targeted-sequencing primer coordinates (<https://github.com/albarmeira/TARGET-seq>). Next, Samtools (v1.9)<sup>68</sup> was used to concatenate the BAM header to the resulting SAM files before re-converting the SAM file to BAM format, which was subsequently sorted by genomic coordinates and  
540 indexed. Both gDNA and mRNA reads were tagged with the cell's unique identifier using Picard (v2.3.0) “*AddOrReplaceReadGroups*” and duplicate reads were subsequently marked using Picard “*MarkDuplicates*”. The sequencing reads overhanging into intronic regions in the mRNA reads were additionally hard-clipped using GATK (v4.1.2.0) *SplitNCigarReads*<sup>69,70</sup>.

### 545 *Variant calling*

Variants were called from the processed BAM files using GATK *Mutect2* with the options [–*tumor-lod-to-emit* 2.0 –*disable-read-filter* *NotDuplicateReadFilter* –*max-reads-per-alignment-start*] to increase the sensitivity of detecting low-frequency variants. The frequency of each nucleotide (A, C, G, T) and indels at each pre-defined  
550 variant site were also called using a Samtools *mpileup* as previously described<sup>16</sup>. Lastly, the coverage at each pre-defined variant site were computed using Bedtools (v2.27.1)<sup>71</sup>.

To determine the coverage threshold of detection for each variant site, the coverage for “blank” controls (empty wells) were first tabulated. A cut-off coverage outlier value  
555 was computed as having a coverage exceeding 1.5 times the length of the interquartile range from the 75th percentile. Next, a value of 30 was added to this outlier value to yield the final coverage threshold to be used for genotype assignment.

### *Genotype assignment*

For each pre-defined variant site, the number of reads representing the reference and  
560 alternative (variant) alleles for indels (insertion and deletions) and SNVs (single

nucleotide variants) were tabulated from the outputs of GATK *Mutect2* and Samtools *mpileup*, respectively.

Here, a genotype scoring system was introduced to assign each variant site into one of three possible genotypes: wildtype, heterozygous, or homozygous mutant. Chi-square ( $\chi^2$ ) test was first used to compare the observed frequency of reference and alternative alleles against the expected fraction of reference and alternative alleles corresponding to the three genotypes. The expected fraction of the reference alleles was 0.999, 0.5, and 0.001, and the expected fraction of the alternative alleles was 0.001, 0.5, and 0.999 for wildtype, heterozygous, and homozygous mutant genotype, respectively. The  $\chi^2$  statistics were then tabulated for each fitted model and converted to genotype scores using the following formula:

$$Score_{genotype} = \frac{1}{\log_{10}(\chi^2 + 1)}$$

The genotype assigned to the variant site was based on the genotype model with the highest score.

Next, the variant (alternative) allele frequency (VAF) was computed and variant sites with  $2 < \text{VAF} < 4$  and  $96 < \text{VAF} < 98$  were reassigned as “ambiguous”. For cells with no variants detected at the specific variant sites by the mutation callers (either due to the absence of the variants, i.e. wild-type genotype, or that such variants were present below the detection limit), a “wild-type” genotype was assigned to those cells with a coverage above the specific threshold and “low coverage” to those cells with coverage below such threshold.

Taken together, each variant site was assigned one of the five following genotypes: wildtype, heterozygous, homozygous mutant, ambiguous, or low coverage. Variants with ambiguous or low coverage assignments for a particular cell were excluded from the analysis.

### **Computational reconstruction of clonal hierarchies**

Genotypes for each single cell were recoded for input to SCITE in a manner inspired by Morita *et al*<sup>72</sup>: each mutation in each gene was coded as two loci, representing two different alleles. In the first recorded loci, all homozygous calls from each mutation

where coded as heterozygous genotype calls. In the second recorded loci, all heterozygous and homozygous genotype calls in the original mutation matrix were coded as homozygous reference (i.e. WT) and heterozygous, respectively. For example, if for a certain mutation 0 represents WT status, 1 encodes heterozygous and 2 refers to homozygous status, these would be encoded as (0,0), (1,0) and (1,1) respectively, where the first term in the parenthesis corresponds to the first loci and the subsequent, to the second loci.

Then, SCITE was used (git revision 2016b31, downloaded from <https://github.com/cbg-ethz/SCITE.git><sup>73</sup>) to sample 1000 mutation trees from the posterior for every single-cell genotype matrix corresponding to a particular patient, where all possible mutation trees are equally likely *a priori*. For patients in which several disease timepoints were available, all timepoints were merged for SCITE analysis. As parameters for every SCITE run “-fd 0.01” (corresponding to the allelic dropout rate of reference allele in our adapted SCITE model), “-ad 0.01” (corresponding to the allelic dropout of the alternate allele), a chain length (-l) of 1e6, and a thinning interval of 1 while marginalizing out cell attachments (-p 1 -s) were used.

To summarize the posterior tree sample distribution, the number of times a particular sample matched each tree was computed. For each patient, the most common tree topology in the posterior tree samples is reported (Extended Data Fig.2b-o, Fig.9e-m), where “pp” is the proportion of samples that match this tree. For each clade in the most common posterior tree, clade probabilities were estimated as the proportion of trees in the posterior that contained the clade. These are indicated in each square for each mutation in (Extended Data Fig.2b-o, Fig.9e-m).

### *Clone assignment*

For every patient’s most common posterior tree, we assigned every cell to the tree node that matches the genotype of that particular cell. If an exact match was not found, then for every tree node the loss of assigning a cell to that node was calculated using the following loss function:

$$l(m) = \log(\text{ADO})(m[1, 2] + m[3, 2]) \\ + \log(\text{FD})(m[2, 1] + m[2, 3]) \\ + \log(\text{ADO}^2\text{FD})(m[1, 3] + m[3, 1])$$

620 where  $m$  is a confusion matrix generated across all loci of a cell in which the first index represents the genotype that was measured for that particular cell (1 = homozygous reference, 2 = heterozygous, 3 = homozygous alternate), and the second index represents the genotype implied by the tree node. ADO = 0.01 and FD = 0.001 were used. Every cell was assigned to the node with the lowest loss  $l$ . For the trees  
625 presented in Extended Data Fig.2b-o and Extended Data Fig.9e-m only the numbers of cells with exact genotype matches were reported.

#### *Testing for evidence of homozygous genotypes*

Due to the nature of our loci-specific mutation encoding (each gene is encoded as two loci), homozygous mutations are placed in the clonal hierarchy independently of their  
630 accuracy. Therefore, for every patient and at every locus with observed homozygous alternate genotype calls, the tested null hypothesis was that all homozygous alternate genotype calls are due to allelic dropout at a level not exceeding 0.05 using a one-tailed binomial test. The total number of draws for the test is the number of heterozygous and homozygous alternate genotype calls at the locus, the number of  
635 successful draws is the number of homozygous alternate calls, and the success rate is 0.05. Only homozygous alternate genotype calls below this 0.05 cut-off were reported in Extended Data Fig.2b-o and Extended Data Fig.9e-m; the results of the binomial test are reported for each patient and mutation in TableS8.

#### **Computational validation of *TP53* biallelic status from single-cell targeted genotyping datasets**

640

To further validate the biallelic status of *TP53* mutations in our dataset, the patterns of allelic dropout in TARGET-seq single-cell genotyping data from patient carrying at least 2 different *TP53* mutations were investigated (n=6; Extended Data Fig.1j).

To test the hypothesis that the observed *TP53*-WT/*TP53*-homozygous (*TP53*-WT/HOM; or (0,2)) cells are the result of a chromosomal loss (and therefore, in different alleles), the following null hypothesis ( $H_0$ ) was formulated: observed *TP53*-WT/HOM cells are double allelic dropout events. Under  $H_0$ , every *TP53*-WT/HOM cell (0,2), *TP53*-HOM/WT cell (2,0), *TP53*-HOM/HOM (2,2) as well as an unknown number of *TP53*-WT/WT (0,0) are the result of a *TP53*-HET/HET (1,1) cell undergoing allelic  
645 dropout (ADO) at both sites. The following assumptions were made: (a) ADO is unbiased towards HOM or WT and (b) ADO events at each *TP53* site are independent.

650

The null hypothesis was then tested with a binomial test, where the number of (2,2) events should be half the sum of (0,2) + (2,0) events (Extended Data Fig.1j). (0,0) events were disregarded.

655 If *TP53* mutations are biallelic, the expected number of WT/HOM and HOM/WT would be higher than HOM/HOM cells considering TARGET-seq expected allelic dropout rates (1-5%).

### Single cell 3'-biased RNA-sequencing data pre-processing

660 FASTQ files for each single cell were generated using bcl2fastq (version 2.20) with default parameters and the following read configuration: Y8N\*, I8, Y63N\*. Read 1 corresponds to a cell-specific barcode, index read correspond to an i7 index sequence from each cDNA pool, and read 2 corresponds to the cDNA molecule. PolyA tails were trimmed from demultiplexed FASTQ files with TrimGalore (version 0.4.1). Reads were  
665 then aligned to the human genome (hg19) using STAR (version 2.4.2a) and counts for each gene were obtained with FeatureCounts (version 1.4.5-p1; options --primary). Counts were then normalized by dividing each gene count by the total library size of each cell and multiplying this value by the median library size of all cells processed, as implemented in the “*normalize\_UMIs*” function from the SingCellaR package<sup>74</sup>  
670 (<https://github.com/supatt-lab/SingCellaR>). A summary of the pre-processing pipeline can be found in <https://github.com/albarmeira/TARGET-seq-WTA>.

Quality control was performed using the following parameters: number of genes detected >500, percentage of ERCC derived reads <35%, percentage of mitochondrial reads <0.25%, percentage of unmapped reads <75%. Cells with less than 2000 reads  
675 in batch1, 5000 reads in batch2 and 20000 reads in batch3 were further excluded. This QC step was performed independently for each sequencing batch owing to differences in sequencing depth (mean library size: 42949 batch 1, 93580 batch 2 and 171393 batch3). After these QC steps, 7123 cells passed QC for batch1, 5779 for batch2 and 6319 for batch 3 (79.3%, 68.9% and 80.3% of cells processed, respectively). Then,  
680 2734 cells from a previously published study<sup>16</sup> corresponding to 8 myelofibrosis patients and 2 normal donor controls were further integrated, encompassing a final dataset of 21955 cells in total.

### Identification of highly variable genes



Highly variable genes above technical noise were identified by fitting a gamma  
685 generalized linear model (GLM) model of the  $\log_2$ (mean expression level) and  
coefficient of variation for each gene, using the  
“*get\_variable\_genes\_by\_fitting\_GLM\_model*” from SingCellaR package and the  
following options: *mean\_expr\_cutoff* = 1, *disp\_zscore\_cutoff* = 0.1,  
690 *quantile\_genes\_expr\_for\_fitting* = 0.6, *quantile\_genes\_cv2\_for\_fitting* = 0.2. Those  
genes with a coefficient of variation above the fitted model and expression cut-off were  
selected for further analysis, excluding those annotated as ribosomal or mitochondrial  
genes.

### **CNA inference from single cell transcriptomes**

InferCNV was used to identify CNAs in single-cell transcriptomes<sup>75</sup>  
695 (<https://github.com/broadinstitute/inferCNV/wiki>). Briefly, inferCNV creates genomic  
bins from gene expression matrices and computes the average level of expression for  
each of these bins. The expression across each bin is then compared to a set of normal  
control cells, and CNAs are predicted using a hidden markov model. For each patient,  
inferCNV was performed with the following parameters: “*cutoff=0.1, denoise=T,*  
700 *HMM=T*”, compared to the same set of normal donor control cells (n=992). To identify  
CNA subclones, inferCNV in *analysis\_mode='subclusters'* was used. CNAs identified  
by inferCNV were manually curated by removing those with size<10kb, merging  
adjacent CNA calls with identical CNA status into larger CNA intervals and comparing  
them to SNP-Array bulk CNA calls. Finally, to generate combined TARGET-seq single-  
705 cell genotyping and CNA-based clonal hierarchies, the CNA status from each inferCNV  
cluster was assigned to its predominant genotype.

### **Dimensionality reduction, data integration and clustering**

PCA was performed using “*runPCA*” function from the *SingCellaR* R package, and  
Force-directed graph analysis was subsequently performed using the  
710 “*runFA2\_ForceDirectedGraph*” with the top 30 PCA dimensions to generate the plots  
in Extended Data Fig.4a.

For the analysis of patient IF0131 presented in Extended Data Fig.3m, PCA was  
performed using “*runPCA*” function from the *SingCellaR* R package and then UMAP  
was performed using the “*runUMAP*” function with the top 10 PCA dimensions and the

715 following options: *n.neighbors*=20, *uwot.metric* = "correlation", *uwot.min.dist*=0.30, *n.seed* = 1.

Integration of TARGET-seq single-cell transcriptomes from 10459 cells corresponding to 14 *TP53*-sAML samples was performed using “*runHarmony*” function implemented in the SingCellaR package, using the patient ID as covariate and the following options:  
720 *n.dims.use*=20, *harmony.theta* = 1, *n.seed* = 1. Diffusion map analysis was performed using “*runDiffusionMap*” with the integrative Harmony embeddings and the following parameters: *n.dims.use*=20, *n.neighbors*=5, *distance*="euclidean". Signature scores were calculated using “*plot\_diffusionmap\_label\_by\_gene\_set*” to generate the plots in Fig.2a and Fig.3a.

### 725 **Pseudotime trajectory analysis**

Monocle3<sup>76</sup> (<https://cole-trapnell-lab.github.io/monocle3/>) was used to infer differentiation trajectories from single cell transcriptomes. Raw UMI count matrix and clustering annotations were extracted from the SingCellaR object to build a Monocle3 ‘*cds*’ object. Next, we retrieved the first two components of the diffusion map (DC1 and DC2), and the ‘*learn\_graph*’ function was then used calculate the trajectory on the two-  
730 dimensional (2D) diffusion map, using *TP53*-WT preleukemic cell cluster as the root node. Pseudotime was calculated using ‘*order\_cells*’ function and overlaid on the diffusion map embeddings to generate the plot in Fig.2b.

### 735 **Differential expression analysis**

Differentially expressed genes from TARGET-seq datasets were identified using a combination of non-parametric Wilcoxon test, to compare the expression values for each group, and Fisher’s exact test, to compare the frequency of expression for each group, as previously described<sup>17</sup>. Logged normalized counts were used as input for  
740 this comparison, including genes expressed in at least 2 cells. Combined p-values were calculated using Fisher’s method and adjusted p-values were derived using Benjamini & Hochberg procedure. Significance level was set at p-adjusted<0.05. For the analysis presented in Extended Data Fig.4b and TableS2, the top 100 differentially expressed genes with  $\log_2(\text{fold-change}) > 0.3$  and at least 20% expressing cells are  
745 shown. Differentially expressed genes identified between *TP53*-multi-hit versus *TP53*-WT cells were further assessed for the enrichment of known p53 target genes (337

750 curated p53 target genes from Fisher *et al*<sup>77</sup>) for the analysis presented in Extended Data Fig.4c. We assessed the extent of overlap of these gene lists using the R package GeneOverlap. The overlapping genes were further assessed for enrichment of p53-related pathways using the R package clusterProfiler.

755 For the analysis presented in Fig.2k,l, only genes overexpressed in *TP53* multi-hit cells and  $\log_2(\text{fold-change}) > 0.75$  were included; for Fig.4d, only those with  $\log_2(\text{fold-change}) > 1$  were considered. Violin plots (Fig.4e and Extended Data Fig.9n) from selected differentially expressed genes were generated using “ggplot2” package in R.

### Gene-Set Enrichment analysis

760 For analysis involving <600 cells (Fig.4c, TableS5) GSEA was performed using GSEA software (<https://www.gsea-msigdb.org/gsea/index.jsp>) with default parameters and 1000 permutations on the phenotype, using  $\log_2(\text{normalized counts})$ .

765 For analysis involving >600 cells per group (Fig.3k, Extended Data Fig.4d and Fig.9o), GSEA was performed with “*identifyGSEAPrerankedGene*” function from *SingCellaR* R package with default options. Briefly, differential expression analysis was performed between two cell populations using Wilcoxon rank sum test and the resulting p-values were adjusted for multiple testing using the Benjamini-Hochberg approach. Prior to the differential expression analysis, down-sampling was performed so that both cell populations had the same number of cells. Next,  $-\log_{10}(\text{p-value})$  transformation was performed and the resulting p-values were multiplied by +1 or -1 if the corresponding  $\log_2\text{FC}$  was  $> 0.1$  or  $< -0.1$ , respectively. The genelist was ranked using this statistic in ascending order and used as input for GSEA analysis using “*fgsea*” function from the *fgsea* R package with default options.

775 MSigDB HALLMARK v7.4 50-gene sets or previously published signatures ([https://www.gsea-msigdb.org/gsea/msigdb/cards/GENTLES\\_LEUKEMIC\\_STEM\\_CELL\\_UP](https://www.gsea-msigdb.org/gsea/msigdb/cards/GENTLES_LEUKEMIC_STEM_CELL_UP)) were used for all analysis. Normalised enrichment scores (NES) were displayed in a heatmap using *pheatmap* R package. Gene sets with False Discovery Rate (FDR) q-value lower than 0.25 were considered significant.

### Projection of single cell transcriptomes

A previously published human haematopoietic atlas was downloaded from <https://github.com/GreenleafLab/MPAL-Single-Cell-2019> and used as a normal haematopoietic reference to project *TP53*-sAML and *de novo* AML transcriptions using Latent Semantic Index Projection (LSI)<sup>78</sup>. Common genes to all datasets were selected and then, *TP53*-sAML or previously published *de novo* AML cells<sup>25</sup> were projected using “*projectLSI*” function for the analysis presented in Fig.2c,d. A previously published human myelofibrosis atlas<sup>79</sup> was used as a reference to project *TP53*-sAML multi-hit cells in the analysis presented in Extended Data Fig.5d,e, using previously defined force-directed graph embeddings.

### Velocyto analysis

Loom files were generated for each single cell using *velocyto* (v0.17.13) with options *-c* and *-U*, to indicate that each BAM represents an independent cell and reads are counted instead of molecules (UMIs), respectively<sup>80</sup>. The individual loom files were subsequently merged using the *combine* function from the *loompy* python module.

Healthy donors with at least 300 cells with RNA-sequencing data and patients with at least 300 cells consisting of >50 preleukemic (*TP53* wildtype) cells and > 50 *TP53* multi-hit cells were included for analysis. For each individual, Seurat object was created from the merged loom file and processed for downstream RNA-velocity analysis<sup>81</sup>. Specifically, for each patient, the spliced RNA counts were normalised using regularised negative binomial regression with the *SCTransform* function<sup>82</sup>. Next, linear dimension reduction was performed using *RunPCA* function and the first 30 principal components were further used to perform non-linear dimension reduction using the *RunUMAP* function. Ninety-six multiple rate kinetics (MURK) genes previously shown to possess coordinated step-change in transcription and hence violate the assumptions behind scVelo were removed<sup>83</sup>. The processed and MURK gene-filtered Seurat object was then saved as h5Seurat format using the *SaveH5Seurat* function and finally converted to h5ad format using the *Convert* function.

AnnData object was created from the h5ad file using the *scvelo* python module for RNA velocity analysis<sup>84</sup>. Highly variable genes were identified and the corresponding spliced and unspliced RNA counts were normalized and log2-transformed using the

815 *scvelo.pp.filter\_and\_normalize* function. Next, the 1<sup>st</sup> and 2<sup>nd</sup> order moments were computed for velocity estimation using the *scvelo.pp.moments* function. The velocities (directionalities) were computed based on the stochastic model as defined in the *scvelo.t1.velocity* function, and the velocities was subsequently projected on the UMAP embeddings generated from Seurat above. Finally, the UMAP embeddings were annotated using the HSPC and erythroid lineage signature scores <sup>74</sup>, and *TP53* mutation status. For each cell, the cell lineage signature score was computed using the average *SCTransform* expression values of the individual cell lineage genes.

820

## **Analysis of bulk BeatAML and TCGA gene expression datasets**

### *Data retrieval and pre-processing*

825 Two publicly available AML cohorts with genetic mutation and RNA-sequencing data available were used to validate findings from our single-cell analysis, namely BeatAML<sup>26</sup> and The Cancer Genome Atlas (TCGA)<sup>27</sup>. Gene expression values in FPKM (fragments per kilobase of transcript per million mapped reads) were retrieved from the National Cancer Institute (NIH) Genomic Data Commons (GDC)<sup>85</sup>. Gene expression values were then offset by 1 and log2-transformed. *TP53* point mutation status was retrieved from the cBio Cancer Genomics Portal (cBioPortal)<sup>86</sup>. Clinical data 830 including survival data for BeatAML and TCGA was retrieved from the BeatAML data viewer (Vizome) and NIH GDC, respectively.

835 We selected samples from the BeatAML cohort with an AML diagnosis (540 *de novo* AML and 96 secondary AML) collected within 1 month of the patient's enrolment in the study, and with both *TP53* mutation status and RNA-sequencing data available. For patients in which multiple samples were available, samples were collapsed to obtain patient-level data. Specifically, the mean gene expression value for each gene from multiple samples was used to represent patient-level gene expression value. Furthermore, patients with at least one sample with a *TP53* mutation were considered *TP53*-mutant. Analysis of *TP53* variant allele frequency and reported karyotypic 840 abnormalities indicated that the vast majority of patients could be classified as "multi-hit", and therefore patients were classified as *TP53*-mutant or WT without taking into account *TP53* allelic status. In total, 360 patients with *TP53* mutation status (329 *TP53*-WT and 31 *TP53*-mutant) and RNA-sequencing data available were included for

analysis. Of these, 322 patients had concomitant survival data available (294 *TP53*-  
845 WT and 28 *TP53*-mutant).

The TCGA cohort consisted for 200 *de novo* AML patients represented by one sample  
each, out of which 151 patients had *TP53* mutation status (140 *TP53*-WT and 11 *TP53*-  
mutant) and RNA-sequencing data available, and were included for analysis. Of these,  
850 132 patients had concomitant survival data available (124 *TP53*-WT and 8 *TP53*-  
mutant).

#### *Cell lineage gene signature scores*

For each sample, a given cell lineage gene signature score was computed as the mean  
855 expression values of the individual genes belonging to the cell lineage gene signature.  
Here, the gene signature scores for two cell lineages were computed, namely myeloid  
and erythroid populations. Two gene sets for each cell lineage were compiled. The first  
gene set was based on cell lineage markers previously reported in the literature  
whereas the second gene set was based on cell lineage markers derived from  
860 analysing a published single-cell dataset<sup>78</sup>. Genes from each score are described in  
TableS3.

For the former approach, six erythroid genes (*KLF1*, *GATA1*, *ZFPM1*, *GATA2*, *GYPA*,  
*TFRC*; Fig.2e, Extended Data Fig.5k, 5m) and seven myeloid genes (*FLI1*, *SFPI1*,  
*CEBPA*, *CEBPB*, *CD33*, *MPO*, *IRF8*; Fig.2f) were identified. For the latter approach,  
865 the expression values of erythroid and myeloid cell clusters were first compared  
separately against all other cell clusters using Wilcoxon ranked sum test. The erythroid  
cluster consisted of the early and late erythroid populations while the myeloid cluster  
consisted of granulocyte, monocyte, and dendritic cell populations. Erythroid and  
myeloid-specific gene signatures were defined as genes having FDR values < 0.05  
870 and log<sub>2</sub> fold change > 0.5 in  $\geq 20$  and 17 comparisons, respectively. In total, 100  
erythroid genes and 135 myeloid genes were identified from this single-cell dataset  
(TableS3), and were used to compute the scores presented in Extended Data Fig.5g-  
j.

#### *TP53 target gene score*

875 Genes downregulated in *TP53*-multi-hit compared to *TP53*-WT cells (defined as per  
“Differential expression analysis” section above; related to Figure S4b) and p53 targets

positively regulated from Fisher *et al*<sup>77</sup> were used to compute a *TP53*-target gene-score presented in Extended Data Fig.5k.

## 880 **Prognostic signatures and Cox-regression survival models**

### *Leukaemic stem cell (LSC) signature score*

The 17-gene leukaemic stem cell (LSC17) gene set was retrieved from Ng *et al*<sup>31</sup>. For each sample, the LSC17 score was defined as the linear combination of gene expression values weighted by their respective regression coefficients.

885 To identify *TP53*-sAML leukaemic stem cell signatures from our TARGET single-cell dataset, two different approaches were used. First, differentially expressed genes were identified as overexpressed in all Lin<sup>-</sup>CD34<sup>+</sup> *TP53* multi-hit cells regardless of their transcriptional classification (“p53-all-cells”) versus myelofibrosis, healthy donor and *TP53*-WT preleukaemic cells; this gene-set consists of 29 genes (TableS4a). For the  
890 second approach, the same analysis was performed, but *TP53* multi-hit cells transcriptionally defined as leukaemic stem cells (falling in the leukaemic stem cell-like cluster, Fig.2a, middle) were specifically selected; this gene-set is comprised of 51 genes (“p53LSC”; TableS4a).

Next, lasso cox regression with 10-fold cross-validation implemented in the *glmnet* R  
895 package was used to identify p53-all-cells and p53-LSC genes that were associated with survival and to estimate their respective regression coefficients<sup>87</sup>. Specifically, Harrel’s concordance measure (C-index) was used to assess the performance of each fitted model during cross-validation. The best model was defined as the fitted model with the highest C-index. Subsequently, the coefficient for each gene estimated using  
900 the best model was used to compute the gene signature scores. Only genes with non-zero coefficient values were included in the final gene set. In total, 9 and 44 genes were retained from the p53-all-cells and p53-LSC gene sets, respectively. For each sample, the gene signature score for each gene set was defined as the linear combination of gene expression values weighted by their respective regression  
905 coefficient<sup>31,87</sup>. The list of p53-LSC and p53-all-cells gene signatures is provided in TableS4b.

### *Survival analysis*

For each gene expression signature, patients were first split using the median gene expression signature score. This resulted in two groups of patients, namely patients with high expression scores (greater than or equal to the median) and patients with low expression scores (lower than the median), exemplified in Extended Data Fig.6a,b.

The Cox proportional hazards regression model implemented by the *survival* R package was fitted to estimate the hazard ratio associated with each feature. Log-rank test was used to test the differences between survival curves. The features analysed here were LSC17, p53-all-cells and p53-LSC signatures. Patients with low gene expression signature scores (below median) and patients with *TP53* wildtype status were specified as the reference groups in the model. Kaplan-Meier curves were plotted using the *survminer* R package to visualize the probability of survival and sample size at a respective time interval.

920

### ***In vitro* assays**

#### *Short-term liquid culture experiments*

For short-term liquid culture differentiation experiments (Fig.3j, Extended Data Fig.7h,i), single cells from different Lineage<sup>-</sup>CD34<sup>+</sup> HSPC populations (HSC: CD34<sup>+</sup>CD38<sup>-</sup>CD45RA<sup>-</sup>CD90<sup>+</sup>, MPP: CD34<sup>+</sup>CD38<sup>-</sup>CD45RA<sup>-</sup>CD90<sup>-</sup>, LMPP: CD34<sup>+</sup>CD38<sup>-</sup>CD45RA<sup>+</sup> and more committed progenitors CD34<sup>+</sup>CD38<sup>+</sup>) were directly sorted into a 96-well tissue culture plate containing 100 µL of differentiation media: StemSpan (Catalog #09650, StemCell Technologies), 1% Penicillin+Streptomycin, 20 % BIT9500 (Cat# 9500, StemCell Technologies), 10 ng/mL SCF (Cat #300-07, Peprotech), 10 ng/mL FLT3L (Cat# 300-19, Peprotech), 10 ng/mL TPO (Cat# 300-18-10, Peprotech), 5 ng/mL IL3 (Cat # 200-03, Peprotech), 10 ng/mL G-CSF (Cat# 300-23, Peprotech), 10 ng/mL GM-CSF (Cat# 300-03, Peprotech), 1 IU/mL EPO (Janssen, erythropoietin alpha, clinical grade) and 10 ng/mL IL6 (Cat# 200-06, Peprotech).

For all liquid culture experiments, 50 µL of fresh 1X differentiation media was added at day 4. Readout was performed by flow cytometry after 12 days of culture using the antibodies detailed in TableS7.c (Panel D).

#### *Long-term culture initiating-cell (LTC-IC) assay*



50 cells from each Lin<sup>-</sup>CD34<sup>+</sup> population (HSC; MPP; LMPP; CD38<sup>+</sup>) and donor type (HD, MF, *TP53*-sAML) were sorted in triplicate. Cells were resuspended in 100  $\mu$ L of myelocult (Stem Cell Technologies, #H5150) supplemented with Hydrocortisone (10<sup>-6</sup>M; Stem Cell Technologies, Cat#74142) and plated into an irradiated supportive stromal cell layer (5000 SI/SI cells and 5000 M2-10B4 cells per well) in a 96-well tissue-culture plate coated with Collagen type I (CORNING; Cat#354236).

Medium was changed weekly and after 6 weeks of culture, cells were washed in IMDM+20%FCS and plated into 1.4 mL of cytokine-rich methylcellulose (Methocult H4435, Stem Cell Technologies). Colonies were scored 14 days later under an inverted microscope, and each colony was classified according to its morphology as BFU-E (Burst-forming unit erythroid), CFU-G (granulocyte), CFU-GM (granulocyte-macrophage), CFU-M (macrophage) or CFU-GEMM (granulocyte, erythrocyte, macrophage, megakaryocyte). Selected colonies were used for cytopsin and genotyping as outlined below.

#### *LTC-IC colony genotyping*

LTC-IC colonies were picked from methylcellulose media, washed, resuspended in 10  $\mu$ L of PBS and transferred to individual wells in a 96-well PCR plate. 15  $\mu$ L of lysis buffer (Triton X-100 0.4%, Qiagen Protease 0.1 AU/mL) were added to each well and samples were incubated at 56 °C for 10 minutes and 72 °C for 20 minutes. A 3  $\mu$ L aliquot from each lysate was used as input to generate a targeted and Illumina-compatible library for colony genotyping. The preparation of single cell genotyping libraries involves 3 PCR steps. In the first PCR step, target-specific primers spanning each mutation of interest are used for amplification (TableS6a); in the second PCR step, nested target-specific primers (TableS6b) attached to universal CS1 / CS2 adaptors (Forward adaptor, CS1: ACACTGACGACATGGTTCTACA; Reverse adaptor, CS2: TACGGTAGCAGAGACTTGGTCT) further enrich for target regions and in the third PCR step, Illumina-compatible adaptors containing sample-specific barcodes are used to generate sequencing libraries.

#### *TP53 knockdown and differentiation of human CD34<sup>+</sup> cells*

shRNA sequence for p53 knockdown has been previously cloned into the lentiviral vector pRRLsin-PGK-eGFP-WPRE and validated<sup>88</sup>. Primary human CD34<sup>+</sup> cells from patients with MPN (Table S1) were infected twice with scramble (shCTL) or shTP53

970 with a MOI (Multiplicity of Infection) of 15 and sorted 48h later on CD34 and GFP  
expression. Cells were cultured in serum-free medium with a cocktail of human  
recombinant cytokines containing EPO (1 U/mL, Amgen), FLT3-L (10 ng/mL, Celldex  
Therapeutics, Inc.), G-CSF (20 ng/mL, Pfizer), IL-6 (10 ng/mL, Miltenyi), GM-CSF (5  
975 ng/mL, Peprotech), IL-3 (10 ng/mL, Miltenyi), TPO (10 ng/mL, Kirin Brewery) and SCF  
(25 ng/mL, Biovitrum AB).

At day 12 of culture, cells were stained with the antibodies detailed in TableS7.c (Panel  
C). DAPI was used for dead cell exclusion before acquisition on a FACSCanto II (BD  
Biosciences) instrument. Analysis of FACS data was performed using Kaluza  
(Beckman Coulter) software.

### 980 **Quantitative real time PCR in shRNA experiments**

In *TP53* knockdown experiments, RNA from either CD34<sup>+</sup> cells sorted after  
transduction or bulk cells at day 12 of culture was extracted using Direct-Zol RNA  
MicroPrep Kit (Zymo Research) and reverse transcription was performed with  
SuperScript Vilo cDNA Synthesis Kit (Invitrogen). Quantitative RT-PCR was performed  
985 on a 7500 Real-Time PCR Machine using SYBR-Green PCR Master Mix (Applied  
Biosystems). Expression levels were normalized to *PPIA* (housekeeping gene).  
Primers used are listed in TableS6c.

### **Xenotransplantation**

Purified CD34<sup>+</sup> cells from AML patients were transplanted via retroorbital vein injection  
990 in sublethally irradiated (1.5Gy) NOD.CB17-*Prkdcscid IL2rgtm1/Bcgen* mice (B-NDG,  
Envigo). All experiments were approved by the French National Ethical Committee on  
Animal Care (n° 2020-007-23589). Blood cell counts were performed monthly by  
submandibular sampling of mice with blood chimerism assessed by flow cytometry  
using hCD34, hCD45 and mCD45 antibodies (TableS7.b; PDX PB panel). At sacrifice,  
995 human BM was stained with the antibodies listed in TableS7.b (PDX BM panel) and  
HSPC fractions were sorted on an Influx Cell sorter (BD Biosciences).

### **Evaluation of cell morphology**

Cell morphology from PDX models (Extended Data Fig.3d) and *in vitro* LTC-IC cultures  
(Extended Data Fig.7f) was assessed after cytopspin of 50-100,000 cells onto a glass

1000 slide (5 min at 500 rpm) and May-Grünwald Giemsa staining, according to standard protocols. Images were obtained using an AxioPhot microscope (Zeiss).

### Mouse Bone Marrow Chimaeras

*Trp53*<sup>tm2Tyj</sup> *Commd10*<sup>Tg(Vav1-icre)A2Kio</sup> or *Trp53*<sup>tm2Tyj</sup> *Tg*<sup>(Tal1-cre/ERT)42-056Jrg</sup> mice (hereafter referred to as Vav-Cre *Trp53*<sup>R172H/+</sup> or SCL-CreER<sup>T</sup> *Trp53*<sup>R172H/+</sup> respectively) and wild-  
1005 type mice used for BM chimera experiments were bred and maintained in accordance to UK and France Home Office regulations. All experiments carried out were performed under Project License P2FF90EE8 approved by the University of Oxford Animal Welfare and Ethical Review Body or under the Project License n° 2020-007-23589, approved by the French National Ethical Committee on Animal Care. *Trp53*<sup>tm2Tyj</sup> <sup>89</sup>,  
1010 *Commd10*<sup>Tg(Vav1-icre)A2Kio</sup> <sup>90</sup> (Jackson laboratory stock number #008610) and *Tg*<sup>(Tal1-cre/ERT)42-056Jrg</sup> <sup>91</sup> have been previously described.

For *in vivo* experiments, two different chimera settings were used. For the first setting (Fig.5a), 1 million bone marrow (BM) cells from Vav-Cre *Trp53*<sup>R172H/+</sup> CD45.1 mice and 1 million BM CD45.2 wild-type cells from competitor mice were transplanted intra-  
1015 venously into lethally irradiated (10 Gy total body irradiation, split dose) congenic CD45.2 mice. For the second setting (Fig.5h), 0.9 million bone marrow (BM) cells from *Trp53*<sup>LSL-R172H/+</sup> CD45.2 mice and 2.1 million BM CD45.1 wild-type competitor mice were transplanted intra-venously into lethally irradiated (9.5 Gy total body irradiation) congenic CD45.2 mice and *Trp53* mutation was induced 4 weeks after transplantation  
1020 by tamoxifen (gavage 200 mg/kg, Sigma) during 4 days, followed by tamoxifen feeding during 2 weeks (Ssniff Diet). In each cohort, a selection of mice were injected intra-peritoneally with 3 rounds of 6 injections each of 200µg poly(I:C) (first setting) or 100µg poly(I:C) (second setting) (GE Healthcare, #27-4732-01) or placebo (PBS1X). Alternatively, Vav-Cre *Trp53*<sup>R172H/+</sup> mice were injected with 3 rounds of 8 injections  
1025 each of 35µg Lipopolysaccharide from Escherichia Coli O111:B4 (LPS; Cat. #L4391-1MG and #L5293-2ML; Sigma-Aldrich).

Poly(I:C) and LPS were administered during weeks 6-7-8, 10-11-12, 14-15-16 (setting 1), or during weeks 7-8, 11-12, 15-16 (setting 2) post-transplantation. Within each round, injections were spaced one or two days apart. Blood cell counts and analysis of  
1030 peripheral blood chimerism along with mature lymphoid and myeloid populations (PB) were performed every 2-4 weeks by submandibular sampling of mice; while BM

chimerism and HSPC populations were analysed 18-20 weeks after transplantation. The antibodies used are detailed in TableS7.d. 7AAD (Sigma) or DAPI (BD Biosciences) were used for dead cell exclusion. FACS analyses were carried out on  
1035 BD Fortessa or BD Fortessa X20 (BD Biosciences) and profiles were later analysed using FlowJo (version 10.1, BD Biosciences) or Kaluza (Beckman Coulter) softwares.

### **LSK apoptosis and cell cycle**

BM LSK cells (setting 2) were stained with Annexin-V and DAPI in Annexin V binding buffer 1X (BD Biosciences) for apoptosis analysis. BM LSK cell cycle was assessed  
1040 by flow cytometry using Ki-67 and DAPI staining, after fixation and permeabilization (BD Cytofix/Cytoperm and Permeabilization Buffer Plus, BD Biosciences).

### **Multiplex in-situ hybridization (M-FISH)**

50 CD45.1 (*Trp53<sup>R172H/+</sup>*) or CD45.2 (wild-type) LSK (Lin-Sca1<sup>+</sup>c-Kit<sup>+</sup>) cells from poly(I:C)-treated and control recipient mice were sorted and cultured for one week into  
1045 Complete X-vivo15 media (Cat. #BE-04-418Q, Lonza) supplemented with 10% Fetal Calf Serum (FCS, #F9665, Sigma-Aldrich), 0.1 mM 2-mercaptoethanol (#21985023, Gibco), 1% penicillin-streptomycin (PAA laboratories), 2ng/ml mouse stem cell factor (mSCF, #250-03, PeproTech), 10ng/ml mouse granulocyte-monocyte colony-stimulating factor (mGM-CSF, Immunex), 5ng/ml human thrombopoietin (hTPO, Cat#  
1050 300-18-10 PeproTech), 10ng/ml human granulocyte colony-stimulating factor (hG-CSF, Neopogen) 5ng/ml human FLT3 ligand (hFL, Cat# 300-19, Immunex), 5ng/ml mouse interleukin 3 (mIL-3, #213-13, PeproTech). Cells were cultured at 37°C 5% CO<sub>2</sub>. On day seven of culture, metaphase spreads were harvested following synchronisation with Colcemid (KaryoMAX™; Cat # 11519876, ThermoFisher  
1055 Scientific) 50 ng/ml, for 3 hours at 37°C. The cells were then incubated with KCl 75mM for 15 minutes at 37°C and spun down. Following this, the cells were fixed in a methanol-acetic acid and then dropped onto glass slides.

M-FISH was performed with the 21XMouse- Multicolor FISH probe kit (Cat #D-0425-  
1060 060-DI, Metasystem Probes), following the manufacturer's instructions. For microscopy analysis, slides were mounted in Vectashield Antifade Mounting Medium with DAPI (Cat. H-1200 2BScientific). Images were acquired and analysed using Leica Cytovision software, on an Olympus BX-51 epifluorescence microscope equipped with

1065 a JAI CVM4+ progressive-scan 24 fps B&W fluorescence CCD camera. All cells were karyotyped, excluding metaphases severely damaged for technical reasons.

1070 The analysis of the M-FISH hybridised cells was blinded. The cells on each slide were scored for the presence of structural aberrations (translocations, and/ or derivative chromosomes and fragments) and/or numerical abnormalities. The presence of more than 40 chromosomes per cell was considered a numerical abnormality, except for cases where it could clearly be attributed to the presence of adjacent metaphases. Chromosome counts lower than 40 were not scored as numerical abnormalities for the impossibility to rule out technical issues (i.e. metaphases bursting at the hypotonic step). We scored as follows: translocations and presence of one chromosome plus one or more extra chromosomal fragment(s)/derivative(s) as “structural abnormalities” (except for sex chromosomes); presence of two chromosomes (or one in case of sex chromosomes) plus one or more extra chromosomal fragment(s)/derivative(s) as “partial chromosome gains”; two chromosomes (or one in case of sex chromosomes) plus one or more extra chromosomes as “whole chromosome gains”; two chromosomes plus two chromosomes with at least 5 different chromosomes present in number=4n as “tetraploidy or sub-tetraploidy”. Counts of numbers of karyotypic aberrations per cell were performed scoring every type of event occurring on one chromosome as single event (i.e., presence of four chromosomes is counted as one aberration).

#### 1085 **IFN $\gamma$ ELISA assay**

1090 Wild-type mice were injected intra-peritoneally with a single dose of 200  $\mu$ g poly(I:C) and spleens were collected from injected mice and non-treated controls 4 hours after injection. Spleens were processed into a single-cell suspension in 200  $\mu$ l PBS, spun down at 500g for 5 minutes and supernatant was collected and used as spleen serum. IFN $\gamma$  levels were assessed using mouse IFN $\gamma$  Quantikine ELISA assay (R&D Systems, cat MIF00) following the manufacturer’s instructions. 450nm and 540nm optical densities were determined using Clariostar microplate reader (BMG Labtech).

#### **Statistical analysis**

1095 Statistical analyses are detailed in Figure Legends and performed using GraphPad Prism software (7 or later version) or R (version 3.6.1) software. Number of independent experiments, donors and replicates for each experiment are detailed in Figure Legends.

### **Data and code availability**

1100 Scripts to reproduce all figures will be available in GitHub upon publication (<https://github.com/albarmeira/>).

The dataset generated in this paper is also available as an interactive vignette [https://wenweixiong.shinyapps.io/TP53\\_MPN\\_AML\\_Single\\_Cell\\_Atlas/](https://wenweixiong.shinyapps.io/TP53_MPN_AML_Single_Cell_Atlas/).

Raw sequencing data is available through GEO (GSE226340) and targeted single-cell genotyping data is publicly available through SRA (PRJNA930152).

### **Acknowledgments**

1105 We are grateful to patients and donors; without their generosity, this study would not have been possible. We also thank Steve Knapper, clinical study teams and other investigators involved in supporting sample collection, and King's Health Partners Biobank for providing access to samples. We thank William Vainchenker for his  
1110 scientific input; Zemin Ren, Timothé Denaes and H el ene Duparc for help with mouse experiments and Sally-Ann Clark for help with sorting. We also thank Dr. Cacho-Soblechero for help with computational analysis. This work was funded by a Medical Research Council (MRC) Senior Clinical Fellowship (A.J.M.; MR/L006340/1), a CRUK Senior Cancer Research Fellowship (A.J.M.; C42639/A26988.), a Cancer Research  
1115 UK (CRUK) DPhil Prize Studentship (C5255/A20936) to A.R-M, a Sir Henry Wellcome Postdoctoral Fellowship from the Wellcome Trust (222800/Z/21/Z) to A.R-M, a British Spanish Society Scholarship to A.R-M., a MRC Confidence in Concept/MLSTF Grant to A.R-M. and A.J.M (MC\_PC\_19049), the MRC Molecular Haematology Unit core award (A.J.M and S.E.J. Eirik; MC\_UU\_12009/5), Emergence Canc erop ole Ile de  
1120 France 2017 (I.A-D.), Association pour la Recherche contre le Cancer 2018 (I.A-D.), Siric-Socrate 2019 (I.A-D.), INCA-PLBIO 2020 (I.A-D.). A.L.C. was supported by Paris University (MENRT grant), J.R.C by a CRUK Senior Cancer Research Fellowship (RCCSCF-Nov21\100004) and S.E.W.J, by a Swedish Research Council and Knut and Alice Wallenberg Foundation. The authors would like to acknowledge the flow

1125 cytometry facility at the MRC Weatherall Institute of Molecular Medicine (WIMM) which  
is supported by the MRC Human Immunology Unit; MRC Molecular Haematology Unit  
(MC\_UU\_12009); National Institute for Health Research (NIHR), Oxford Biomedical  
Research Centre (BRC); Kay Kendall Leukaemia Fund (KKL1057), John Fell Fund  
(131/030 and 101/517), the EPA fund (CF182 and CF170) and by the MRC WIMM  
1130 Strategic Alliance awards G0902418 and MC\_UU\_12025. The authors acknowledge  
the contributions of Dr. Neil Ashley at the MRC Weatherall Institute of Molecular  
Medicine (WIMM) Single Cell Facility and MRC-funded Oxford Consortium for Single-  
Cell Biology (MR/M00919X/1). The authors would also like to acknowledge the  
contribution of the WIMM Sequencing Facility, supported by the MRC Human  
1135 Immunology Unit and by the EPA fund (CF268), the Gustave Roussy flow cytometry  
platform and mouse facility. We also thank the Oxford Genomics Centre at the  
Wellcome Centre for Human Genetics (funded by Wellcome Trust grant reference  
203141/Z/16/Z) for the generation and initial processing of the OmniExpress SNP array  
data. The results published here are in whole or part based upon data generated by  
1140 the TCGA Research Network (<https://www.cancer.gov/tcga>) and the BeatAML team.  
The views expressed are those of the authors and not necessarily those of the National  
Health Service (NHS), the NIHR or the Department of Health.

### **Author contributions**

1145 A.R.M. conceived the project, designed and performed experiments, performed  
computational analysis, analysed data and wrote the manuscript. R.N., A.L.C., H.R.,  
J.O.S., E.L., A.P., J.C. and B.W. designed, performed experiments and analysed data.  
S.W., G.W. and W.W.K. performed computational analysis. J.E.M. collected primary  
samples and clinical and bibliographic data. C.D. provided clinical data. C.B. and M.B.  
1150 analysed SNP-Array data. J.O.S., C.B., N.S., F.G., F.P., I.P., M.D., C.H. provided  
patients samples, clinical data and scientific input. C.M., H.G. analysed and provided  
patients and PDX biological data (NGS and SNP-array). A.H. performed and analysed  
patient's targeted sequencing. D.M. and L.S.M. performed M-FISH experiments and  
analysed M-FISH data. J.R.C., S.E.W.J, B.P. and S.T provided scientific input and  
1155 conceptualization. S.T. supervised computational analysis. I.A-D. conceived and  
supervised the project, analysed data and wrote the manuscript. A.J.M. conceived and  
supervised the project, provided clinical care and wrote the manuscript.

1160 **Competing Interests statement**

A patent relating to the TARGET-seq technique is licensed to Alethiomics Ltd, a spin out company from the University of Oxford with equity owned by B.P. and A.J.M. The other authors declare no competing interests.

1165 **Materials & Correspondence.** Requests for material(s) should be addressed and will be fulfilled by corresponding authors: Alba Rodriguez-Meira ([albarmeira@gmail.com](mailto:albarmeira@gmail.com)), Iléana Antony-Debré ([ileana.antony-debre@gustaveroussy.fr](mailto:ileana.antony-debre@gustaveroussy.fr)) and Adam J. Mead ([adam.mead@imm.ox.ac.uk](mailto:adam.mead@imm.ox.ac.uk)).

1170 **Supplementary Tables**

**TableS1.** Clinical and genetic details from healthy donors and patients included in the study.

1175 **TableS2.** Differentially expressed genes between *TP53* multi-hit HSPCs and *TP53*-WT cells.

**TableS3.** Genesets used to calculate gene expression signature scores in TARGET-seq and publicly available bulk-transcriptomic datasets.

1180 **Table S4.** Differentially upregulated genes in *TP53* multi-hit cells (globally or LSCs) and genes selected by lasso regression to derive p53-all-cells and p53-LSC signatures.

**TableS5.** Gene signatures from *TP53* mutant heterozygous HSPCs from CP-*TP53*-MPN and pre-*TP53*-AML patients.

**TableS6.** Primers used throughout the experiments presented in the manuscript.

**TableS7.** Antibodies used for all experiments presented throughout the manuscript.

1185 **TableS8.** Summary of mutation-specific homozygous status statistical testing for *TP53*-sAML and CP *TP53*-MPN patients. Related to Fig.1b-f; Fig.4b; Extended Data Fig.2b-o; Extended Data Fig.9e-m.

1190 **References**

1. Sill, H., Zebisch, A. & Haase, D. Acute Myeloid Leukemia and Myelodysplastic Syndromes with TP53 Aberrations - A Distinct Stem Cell Disorder. *Clin Cancer Res* **26**, 5304-5309 (2020).
- 1195 2. Bernard, E. *et al.* Implications of TP53 allelic state for genome stability, clinical presentation and outcomes in myelodysplastic syndromes. *Nat Med* **26**, 1549-1556 (2020).
3. Kasthuber, E.R. & Lowe, S.W. Putting p53 in Context. *Cell* **170**, 1062-1078 (2017).
- 1200 4. Lindsley, R.C. *et al.* Acute myeloid leukemia ontogeny is defined by distinct somatic mutations. *Blood* **125**, 1367-76 (2015).



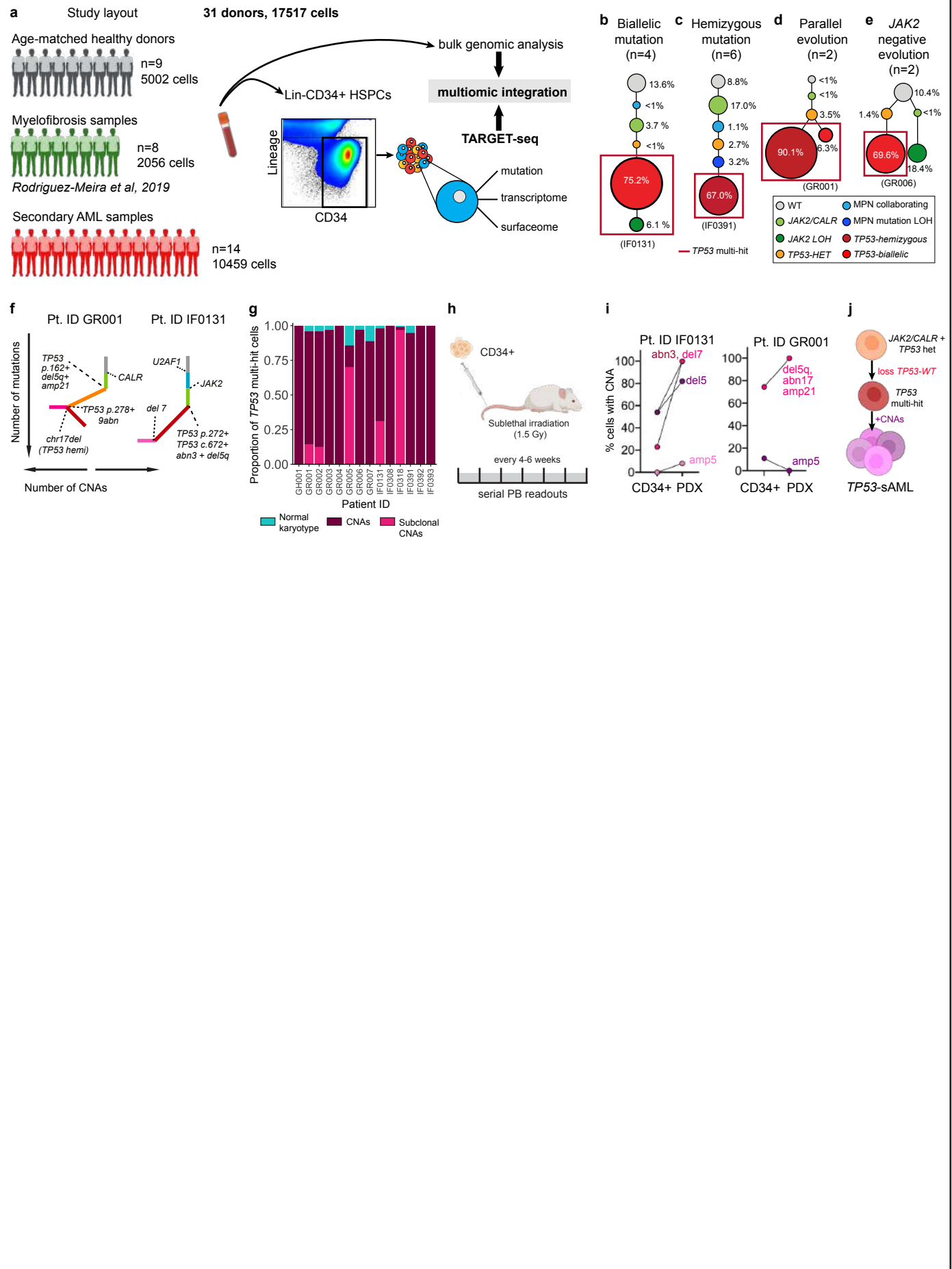
5. Granfeldt Østgård, L.S. *et al.* Epidemiology and Clinical Significance of Secondary and Therapy-Related Acute Myeloid Leukemia: A National Population-Based Cohort Study. *J Clin Oncol* **33**, 3641-9 (2015).
- 1205 6. Mead, A.J. & Mullally, A. Myeloproliferative neoplasm stem cells. *Blood* **129**, 1607-1616 (2017).
7. Celik, H. *et al.* A Humanized Animal Model Predicts Clonal Evolution and Therapeutic Vulnerabilities in Myeloproliferative Neoplasms. *Cancer Discov* (2021).
8. Dunbar, A.J., Rampal, R.K. & Levine, R. Leukemia secondary to myeloproliferative neoplasms. *Blood* **136**, 61-70 (2020).
- 1210 9. Lasho, T.L. *et al.* Targeted next-generation sequencing in blast phase myeloproliferative neoplasms. *Blood Adv* **2**, 370-380 (2018).
10. Luque Paz, D. *et al.* Leukemic evolution of polycythemia vera and essential thrombocythemia: genomic profiles predict time to transformation. *Blood Adv* **4**, 4887-4897 (2020).
- 1215 11. Rampal, R. *et al.* Genomic and functional analysis of leukemic transformation of myeloproliferative neoplasms. *Proc Natl Acad Sci U S A* **111**, E5401-10 (2014).
12. Marcellino, B.K. *et al.* Advanced forms of MPNs are accompanied by chromosomal abnormalities that lead to dysregulation of TP53. *Blood Adv* **2**, 3581-3589 (2018).
13. Courtier, F. *et al.* Genomic analysis of myeloproliferative neoplasms in chronic and acute phases. *Haematologica* **102**, e11-e14 (2017).
- 1220 14. Tsuruta-Kishino, T. *et al.* Loss of p53 induces leukemic transformation in a murine model of Jak2 V617F-driven polycythemia vera. *Oncogene* **36**, 3300-3311 (2017).
15. Kubsova, B. *et al.* Low-burden TP53 mutations in chronic phase of myeloproliferative neoplasms: association with age, hydroxyurea administration, disease type and JAK2 mutational status. *Leukemia* **32**, 450-461 (2018).
- 1225 16. Rodriguez-Meira, A. *et al.* Unravelling Intratumoral Heterogeneity through High-Sensitivity Single-Cell Mutational Analysis and Parallel RNA Sequencing. *Mol Cell* **73**, 1292-1305 e8 (2019).
17. Giustacchini, A. *et al.* Single-cell transcriptomics uncovers distinct molecular signatures of stem cells in chronic myeloid leukemia. *Nat Med* **23**, 692-702 (2017).
- 1230 18. Campbell, P.J. *et al.* Mutation of JAK2 in the myeloproliferative disorders: timing, clonality studies, cytogenetic associations, and role in leukemic transformation. *Blood* **108**, 3548-55 (2006).
19. Goardon, N. *et al.* Coexistence of LMPP-like and GMP-like leukemia stem cells in acute myeloid leukemia. *Cancer Cell* **19**, 138-52 (2011).
- 1235 20. Booth, C.A.G. *et al.* Ezh2 and Runx1 Mutations Collaborate to Initiate Lympho-Myeloid Leukemia in Early Thymic Progenitors. *Cancer Cell* **33**, 274-291.e8 (2018).
21. Ledergor, G. *et al.* Single cell dissection of plasma cell heterogeneity in symptomatic and asymptomatic myeloma. *Nat Med* **24**, 1867-1876 (2018).
- 1240 22. Mesa, R.A. *et al.* Leukemic transformation in myelofibrosis with myeloid metaplasia: a single-institution experience with 91 cases. *Blood* **105**, 973-7 (2005).
23. Passamonti, F. *et al.* Leukemic transformation of polycythemia vera: a single center study of 23 patients. *Cancer* **104**, 1032-6 (2005).
- 1245 24. Granja, J.M. *et al.* Single-cell multiomic analysis identifies regulatory programs in mixed-phenotype acute leukemia. *Nature Biotechnology* **37**, 1458-1465 (2019).
25. van Galen, P. *et al.* Single-Cell RNA-Seq Reveals AML Hierarchies Relevant to Disease Progression and Immunity. *Cell* **176**, 1265-1281 e24 (2019).
26. Tyner, J.W. *et al.* Functional genomic landscape of acute myeloid leukaemia. *Nature* **562**, 526-531 (2018).
- 1250 27. Ley, T.J. *et al.* Genomic and epigenomic landscapes of adult de novo acute myeloid leukemia. *N Engl J Med* **368**, 2059-74 (2013).
28. Boettcher, S. *et al.* A dominant-negative effect drives selection of TP53 missense mutations in myeloid malignancies. *Science* **365**, 599-604 (2019).

- 1255 29. Wagner, K. *et al.* Absence of the transcription factor CCAAT enhancer binding protein alpha results in loss of myeloid identity in bcr/abl-induced malignancy. *Proc Natl Acad Sci U S A* **103**, 6338-43 (2006).
30. Bereshchenko, O. *et al.* Hematopoietic stem cell expansion precedes the generation of committed myeloid leukemia-initiating cells in C/EBPalpha mutant AML. *Cancer Cell* **16**, 390-400 (2009).
- 1260 31. Ng, S.W. *et al.* A 17-gene stemness score for rapid determination of risk in acute leukaemia. *Nature* **540**, 433-437 (2016).
32. Bryder, D. *et al.* Self-renewal of multipotent long-term repopulating hematopoietic stem cells is negatively regulated by Fas and tumor necrosis factor receptor activation. *J Exp Med* **194**, 941-52 (2001).
- 1265 33. Jacobsen, F.W., Stokke, T. & Jacobsen, S.E. Transforming growth factor-beta potently inhibits the viability-promoting activity of stem cell factor and other cytokines and induces apoptosis of primitive murine hematopoietic progenitor cells. *Blood* **86**, 2957-66 (1995).
- 1270 34. Demerdash, Y., Kain, B., Essers, M.A.G. & King, K.Y. Yin and Yang: The dual effects of interferons on hematopoiesis. *Exp Hematol* **96**, 1-12 (2021).
35. Trapp, S. *et al.* Double-stranded RNA analog poly(I:C) inhibits human immunodeficiency virus amplification in dendritic cells via type I interferon-mediated activation of APOBEC3G. *J Virol* **83**, 884-95 (2009).
- 1275 36. Cacemiro, M.D.C. *et al.* Philadelphia-negative myeloproliferative neoplasms as disorders marked by cytokine modulation. *Hematol Transfus Cell Ther* **40**, 120-131 (2018).
37. Ngkelo, A., Meja, K., Yeadon, M., Adcock, I. & Kirkham, P.A. LPS induced inflammatory responses in human peripheral blood mononuclear cells is mediated through NOX4 and G $\alpha$  dependent PI-3kinase signalling. *Journal of Inflammation* **9**, 1 (2012).
- 1280 38. Libregts, S.F. *et al.* Chronic IFN-gamma production in mice induces anemia by reducing erythrocyte life span and inhibiting erythropoiesis through an IRF-1/PU.1 axis. *Blood* **118**, 2578-88 (2011).
39. Essers, M.A. *et al.* IFNalpha activates dormant haematopoietic stem cells in vivo. *Nature* **458**, 904-8 (2009).
- 1285 40. Pietras, E.M. *et al.* Re-entry into quiescence protects hematopoietic stem cells from the killing effect of chronic exposure to type I interferons. *J Exp Med* **211**, 245-62 (2014).
41. Walter, D. *et al.* Exit from dormancy provokes DNA-damage-induced attrition in haematopoietic stem cells. *Nature* **520**, 549-52 (2015).
- 1290 42. Loizou, E. *et al.* A Gain-of-Function p53-Mutant Oncogene Promotes Cell Fate Plasticity and Myeloid Leukemia through the Pluripotency Factor FOXH1. *Cancer Discov* **9**, 962-979 (2019).
43. Boddu, P. *et al.* Erythroleukemia-historical perspectives and recent advances in diagnosis and management. *Blood Rev* **32**, 96-105 (2018).
- 1295 44. Iacobucci, I. *et al.* Genomic subtyping and therapeutic targeting of acute erythroleukemia. *Nature Genetics* **51**, 694-704 (2019).
45. Trainor, C.D., Mas, C., Archambault, P., Di Lello, P. & Omichinski, J.G. GATA-1 associates with and inhibits p53. *Blood* **114**, 165-73 (2009).
46. Enver, T. & Jacobsen, S.E. Developmental biology: Instructions writ in blood. *Nature* **461**, 183-4 (2009).
- 1300 47. Caiado, F., Pietras, E.M. & Manz, M.G. Inflammation as a regulator of hematopoietic stem cell function in disease, aging, and clonal selection. *J Exp Med* **218**(2021).
48. Hormaechea-Agulla, D. *et al.* Chronic infection drives Dnmt3a-loss-of-function clonal hematopoiesis via IFN $\gamma$  signaling. *Cell Stem Cell* **28**, 1428-1442.e6 (2021).
- 1305 49. Avagyan, S. *et al.* Resistance to inflammation underlies enhanced fitness in clonal hematopoiesis. *Science* **374**, 768-772 (2021).
50. Lussana, F. & Rambaldi, A. Inflammation and myeloproliferative neoplasms. *J Autoimmun* **85**, 58-63 (2017).

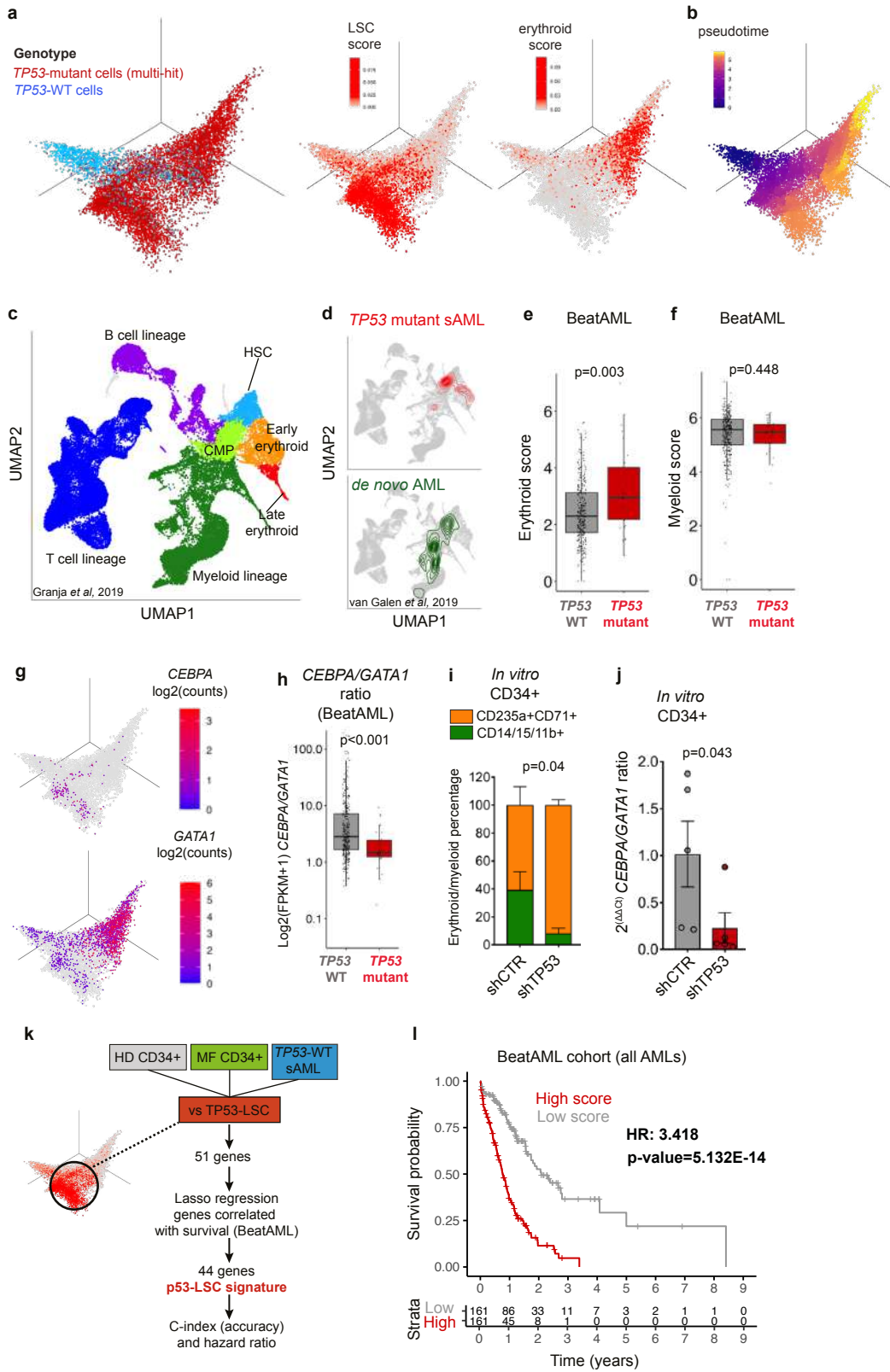
51. Kleppe, M. *et al.* Dual Targeting of Oncogenic Activation and Inflammatory Signaling Increases Therapeutic Efficacy in Myeloproliferative Neoplasms. *Cancer Cell* **33**, 29-43 e7 (2018).
52. Pan-cancer analysis of whole genomes. *Nature* **578**, 82-93 (2020).
53. Hamblin, A. *et al.* Development and Evaluation of the Clinical Utility of a Next Generation Sequencing (NGS) Tool for Myeloid Disorders. *Blood* **124**, 2373 (2014).
54. Zhou, X. *et al.* Exploring genomic alteration in pediatric cancer using ProteinPaint. *Nat Genet* **48**, 4-6 (2016).
55. Papaemmanuil, E. *et al.* Genomic Classification and Prognosis in Acute Myeloid Leukemia. *N Engl J Med* **374**, 2209-2221 (2016).
56. Coombs, C.C. *et al.* Therapy-Related Clonal Hematopoiesis in Patients with Non-hematologic Cancers Is Common and Associated with Adverse Clinical Outcomes. *Cell Stem Cell* **21**, 374-382.e4 (2017).
57. Desai, P. *et al.* Somatic mutations precede acute myeloid leukemia years before diagnosis. *Nat Med* **24**, 1015-1023 (2018).
58. Young, A.L., Challen, G.A., Birmann, B.M. & Druley, T.E. Clonal haematopoiesis harbouring AML-associated mutations is ubiquitous in healthy adults. *Nat Commun* **7**, 12484 (2016).
59. Loh, P.R. *et al.* Insights into clonal haematopoiesis from 8,342 mosaic chromosomal alterations. *Nature* **559**, 350-355 (2018).
60. Loh, P.R., Genovese, G. & McCarroll, S.A. Monogenic and polygenic inheritance become instruments for clonal selection. *Nature* **584**, 136-141 (2020).
61. Olshen, A.B., Venkatraman, E.S., Lucito, R. & Wigler, M. Circular binary segmentation for the analysis of array-based DNA copy number data. *Biostatistics* **5**, 557-72 (2004).
62. Venkatraman, E.S. & Olshen, A.B. A faster circular binary segmentation algorithm for the analysis of array CGH data. *Bioinformatics* **23**, 657-63 (2007).
63. Lawrence, M. *et al.* Software for computing and annotating genomic ranges. *PLoS Comput Biol* **9**, e1003118 (2013).
64. Bashton, M. *et al.* Concordance of copy number abnormality detection using SNP arrays and Multiplex Ligation-dependent Probe Amplification (MLPA) in acute lymphoblastic leukaemia. *Sci Rep* **10**, 45 (2020).
65. Mermel, C.H. *et al.* GISTIC2.0 facilitates sensitive and confident localization of the targets of focal somatic copy-number alteration in human cancers. *Genome Biol* **12**, R41 (2011).
66. Rodriguez-Meira, A., O'Sullivan, J., Rahman, H. & Mead, A.J. TARGET-Seq: A Protocol for High-Sensitivity Single-Cell Mutational Analysis and Parallel RNA Sequencing. *STAR Protoc* **1**, 100125 (2020).
67. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21 (2013).
68. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-9 (2009).
69. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**, 1297-303 (2010).
70. Schischlik, F. *et al.* Mutational landscape of the transcriptome offers putative targets for immunotherapy of myeloproliferative neoplasms. *Blood* **134**, 199-210 (2019).
71. Quinlan, A.R. & Hall, I.M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841-2 (2010).
72. Morita, K. *et al.* Clonal evolution of acute myeloid leukemia revealed by high-throughput single-cell genomics. *Nature Communications* **11**, 5327 (2020).
73. Jahn, K., Kuipers, J. & Beerenwinkel, N. Tree inference for single-cell data. *Genome Biol* **17**, 86 (2016).
74. Roy, A. *et al.* Transitions in lineage specification and gene regulatory networks in hematopoietic stem/progenitor cells over human development. *Cell Reports* **36**, 109698 (2021).

75. Patel, A.P. *et al.* Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. *Science* **344**, 1396-401 (2014).
- 1365 76. Qiu, X. *et al.* Reversed graph embedding resolves complex single-cell trajectories. *Nature Methods* **14**, 979-982 (2017).
77. Fischer, M. Census and evaluation of p53 target genes. *Oncogene* **36**, 3943-3956 (2017).
78. Granja, J.M. *et al.* Single-cell multiomic analysis identifies regulatory programs in mixed-phenotype acute leukemia. *Nat Biotechnol* **37**, 1458-1465 (2019).
- 1370 79. Psaila, B. *et al.* Single-Cell Analyses Reveal Megakaryocyte-Biased Hematopoiesis in Myelofibrosis and Identify Mutant Clone-Specific Targets. *Mol Cell* **78**, 477-492 e8 (2020).
80. La Manno, G. *et al.* RNA velocity of single cells. *Nature* **560**, 494-498 (2018).
- 1375 81. Satija, R., Farrell, J.A., Gennert, D., Schier, A.F. & Regev, A. Spatial reconstruction of single-cell gene expression data. *Nat Biotechnol* **33**, 495-502 (2015).
82. Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biol* **20**, 296 (2019).
83. Barile, M. *et al.* Coordinated changes in gene expression kinetics underlie both mouse and human erythroid maturation. *Genome Biology* **22**, 197 (2021).
- 1380 84. Bergen, V., Lange, M., Peidli, S., Wolf, F.A. & Theis, F.J. Generalizing RNA velocity to transient cell states through dynamical modeling. *Nat Biotechnol* **38**, 1408-1414 (2020).
85. Heath, A.P. *et al.* The NCI Genomic Data Commons. *Nat Genet* **53**, 257-262 (2021).
86. Cerami, E. *et al.* The cBio cancer genomics portal: an open platform for exploring multidimensional cancer genomics data. *Cancer Discov* **2**, 401-4 (2012).
- 1385 87. Anande, G. *et al.* RNA Splicing Alterations Induce a Cellular Stress Response Associated with Poor Prognosis in Acute Myeloid Leukemia. *Clin Cancer Res* **26**, 3597-3607 (2020).
88. Mahfoudhi, E. *et al.* P53 activation inhibits all types of hematopoietic progenitors and all stages of megakaryopoiesis. *Oncotarget* **7**, 31980-92 (2016).
- 1390 89. Olive, K.P. *et al.* Mutant p53 gain of function in two mouse models of Li-Fraumeni syndrome. *Cell* **119**, 847-60 (2004).
90. de Boer, J. *et al.* Transgenic mice with hematopoietic and lymphoid specific expression of Cre. *Eur J Immunol* **33**, 314-25 (2003).
- 1395 91. Göthert, J.R. *et al.* In vivo fate-tracing studies using the Scl stem cell enhancer: embryonic hematopoietic stem cells significantly contribute to adult hematopoiesis. *Blood* **105**, 2724-32 (2005).

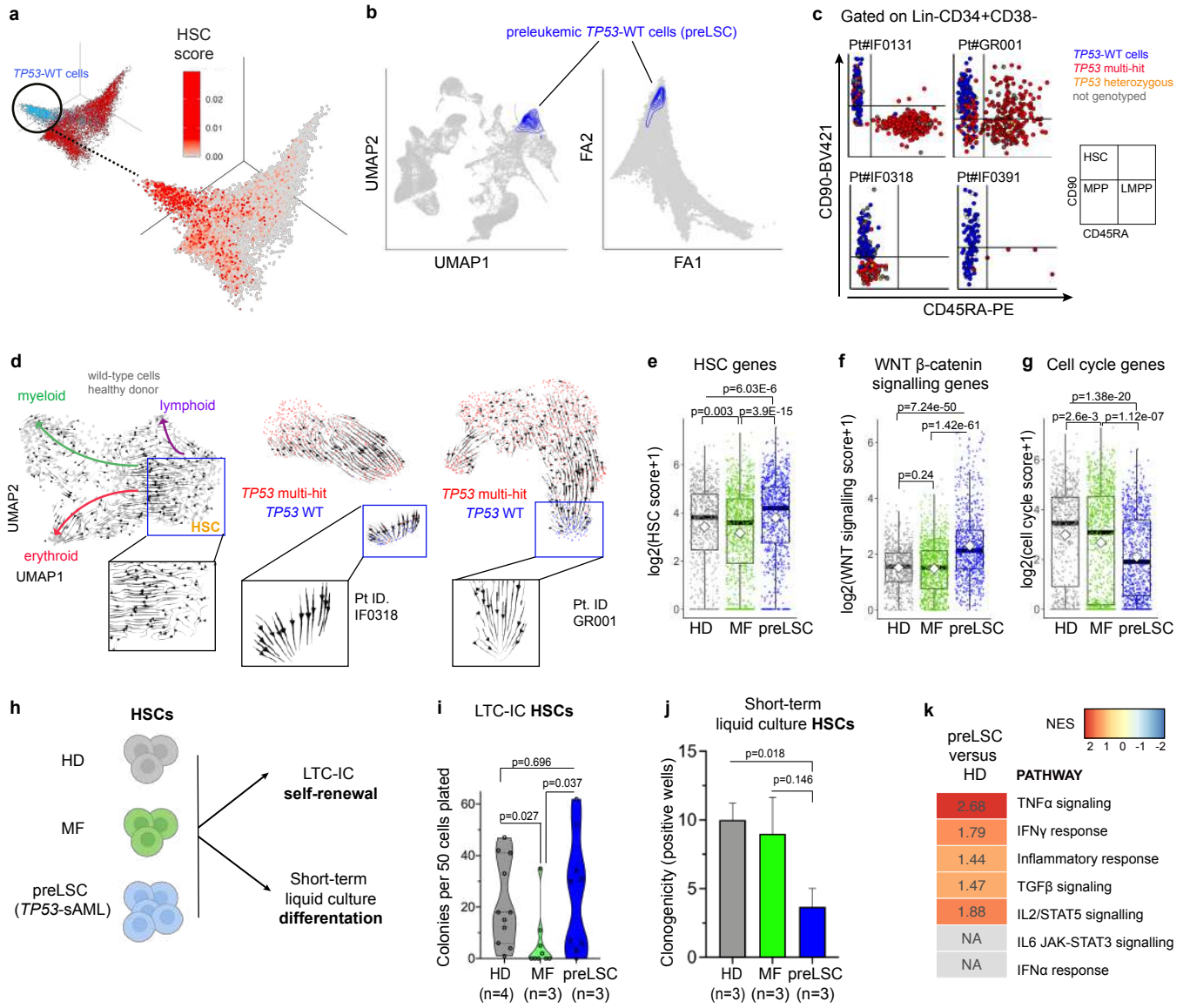
# Figure 1



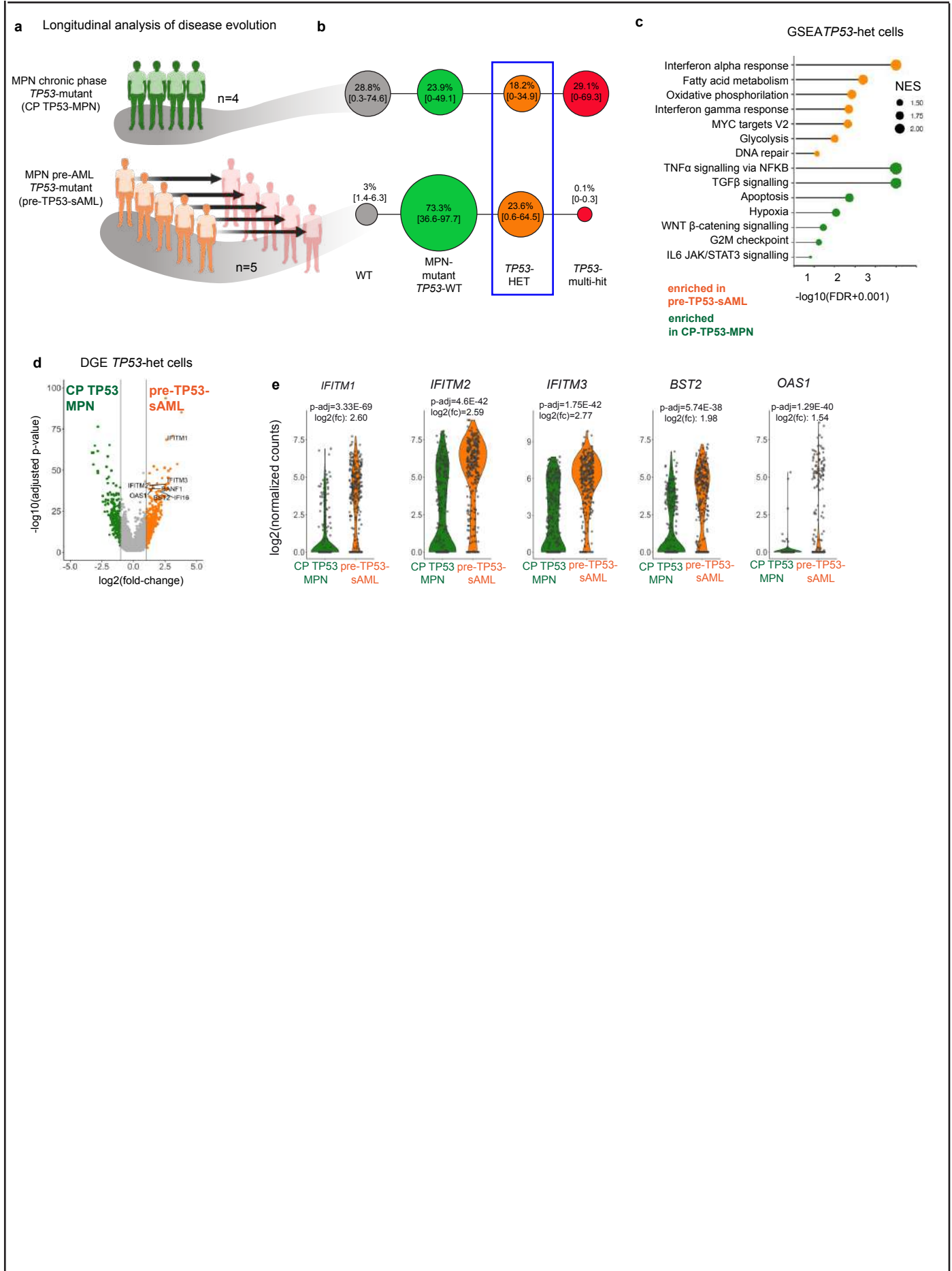
# Figure 2



# Figure 3

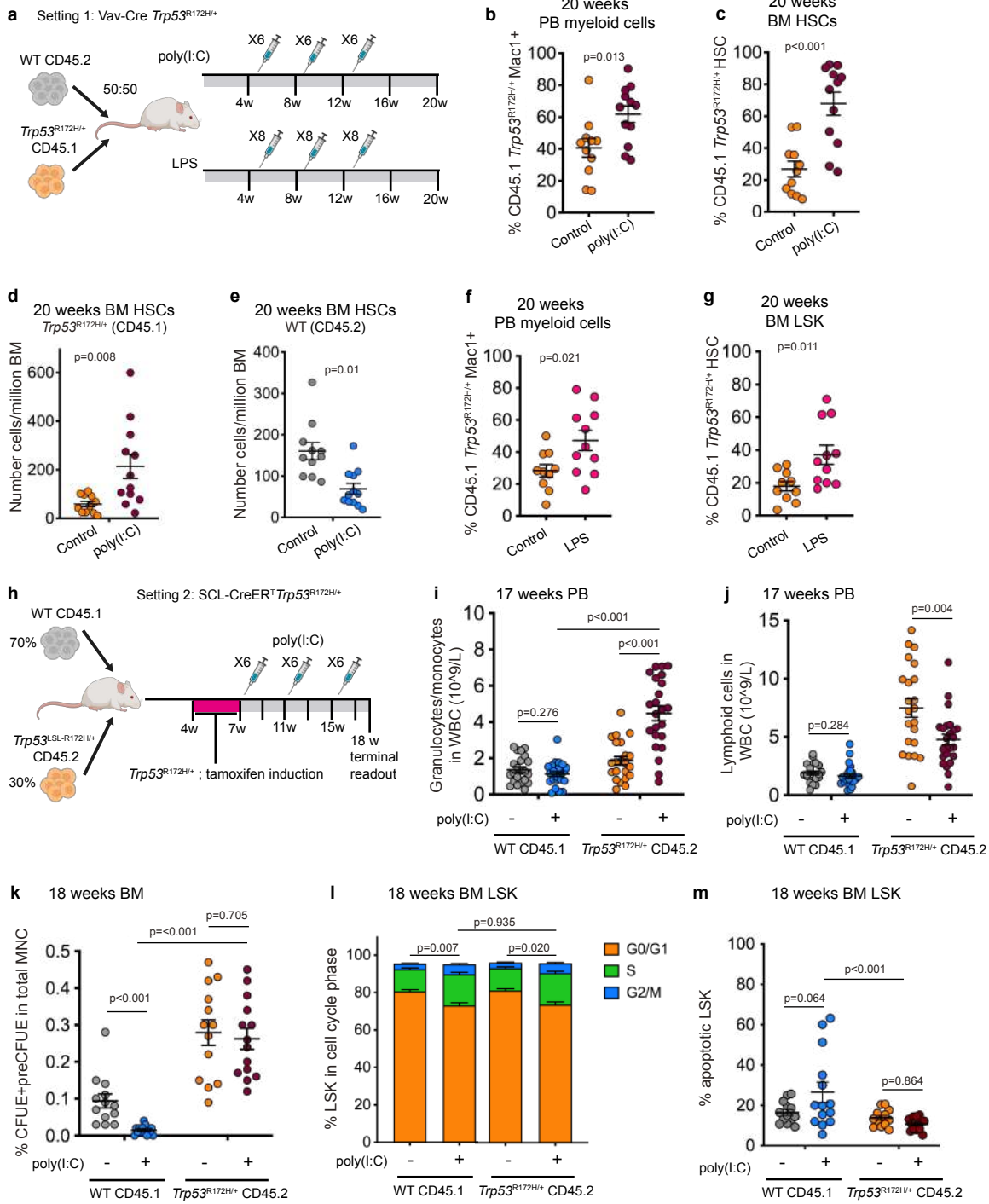


# Figure 4





# Figure 5



# Figure 6

