

Combined landscape of single-nucleotide variants and copy-number alterations in clonal hematopoiesis

Ryunosuke Saiki¹, Yukihide Momozawa², Yasuhito Nannya¹, Masahiro M Nakagawa^{1,3}, Yotaro Ochi¹,
Tetsuichi Yoshizato¹, Chikashi Terao⁴, Yutaka Kuroda⁵, Yuichi Shiraishi⁶, Kenichi Chiba⁶, Hiroko Tanaka⁷,
Atsushi Niida⁸, Seiya Imoto⁹, Koichi Matsuda¹⁰, Takayuki Morisaki¹¹, Yoshinori Murakami¹¹, Yoichiro Kamatani^{4,10},
Shuichi Matsuda⁵, Michiaki Kubo¹², Satoru Miyano⁷, Hideki Makishima¹, Seishi Ogawa^{1,3,13}

¹Department of Pathology and Tumor Biology, Graduate School of Medicine, Kyoto University, Kyoto, Japan

²Laboratory for Genotyping Development, RIKEN Center for Integrative Medical Sciences, Yokohama, Japan

³Institute for the Advanced Study of Human Biology (WPI-ASHBi), Kyoto University, Kyoto, Japan

⁴Laboratory for Statistical and Translational Genetics, RIKEN Center for Integrative Medical Sciences, Yokohama, Japan

⁵Department of Orthopaedic Surgery, Kyoto University Graduate School of Medicine, Kyoto, Japan.

⁶Division of Cellular Signaling, National Cancer Center Research Institute, Tokyo, Japan

⁷Department of Integrated Data Science, M&D Data Science Center, Tokyo Medical and Dental University, Tokyo, Japan

⁸Laboratory of Molecular Medicine, Human Genome Center, The Institute of Medical Science, The University of Tokyo, Tokyo, Japan

⁹Division of Health Medical Data Science, Health Intelligence Center, Institute of Medical Science, University of Tokyo, Tokyo, Japan

¹⁰Department of Computational Biology and Medical Sciences, Graduate school of Frontier Sciences, The University of Tokyo, Tokyo, Japan

¹¹Division of Molecular Pathology, Institute of Medical Science, The University of Tokyo, Tokyo, Japan

¹²RIKEN Center for Integrative Medical Sciences, Yokohama, Japan

¹³Department of Medicine, Centre for Haematology and Regenerative Medicine, Karolinska Institute, Stockholm, Sweden

Correspondence should be addressed to:

Seishi Ogawa (sogawa-tky@umin.ac.jp).

Conflict of interest disclosure:

The authors declare no conflict of interest.

Text word count: 4566 words

Number of figures: 6; Number of references: 50

1 **Abstract**

2 Implicated in the development of hematological malignancies (HM) and cardiovascular mortality, clonal
3 hematopoiesis (CH) in apparently healthy individuals has been investigated by detecting either single-
4 nucleotide variants and indels (SNVs/indels) or copy number alterations (CNAs), but not both. Here by
5 combining targeted sequencing of 23 CH-related genes and array-based CNA detection of blood-derived DNA,
6 we have delineated the landscape of CH-related SNVs/indels and CNAs in a general population of 11,234
7 individuals, including 672 with subsequent HM development. Both CH-related lesions significantly co-occurred,
8 which combined, affected blood count, hypertension, and the mortality from HM and cardiovascular diseases
9 depending on the total number of both lesions, highlighting the importance of detecting both lesions in the
10 evaluation of CH.

11

12 **Introduction**

13 The presence of clonal components in an apparently normal hematopoietic compartment, or clonal
14 hematopoiesis (CH), has been drawing an increasing attention of recent years^{1,2}. Although suggested only
15 indirectly by skewed chromosome X inactivation in early studies³⁻⁷, CH has recently been demonstrated by
16 detecting copy number alterations (CNAs) in the peripheral blood samples from large cohorts of individuals
17 without blood cancers using single-nucleotide polymorphism (SNP) array data from genome-wide association
18 studies (GWAS)⁸⁻¹¹. Showing a substantial overlap to those characteristic of hematological malignancies (HM),
19 CNAs were shown to be associated with an elevated risk of developing HM^{8,9}. More recently, CH has also been
20 detected by the presence of somatic single-nucleotide variants and indels (SNVs/indels) in the peripheral blood
21 of apparently healthy individuals¹²⁻¹⁵ and cancer patients^{16,17} using next generation sequencing. In addition to
22 its link to HM, CH as detected by SNVs/indels has been highlighted by its unexpected association with a
23 significantly increased risk for cardiovascular diseases (CVD)^{12,13,18,19}.

24 Regardless of the type of genetic lesions by which it is detected, CH is strongly age-related with an
25 increasing frequency in the elderly⁸⁻¹³. With substantially improved technologies to identify CNAs and somatic
26 SNVs/indels, a complete registry of CNAs and SNVs/indels associated with CH has been elucidated, which are
27 thought to involve virtually every individual in the extreme elderly^{20,21}. However, to date, no studies have
28 evaluated both CNAs and SNVs/indels together at a comparable sensitivity in a large cohort of a general
29 population, although they have recently been investigated in a cancer population, where many had been
30 treated with chemo/radiotherapy²². What is the landscape of CH recognized by combining both CNAs and
31 SNVs/indels in a general population? Are there any interactions between SNVs/indels and CNAs that shaped
32 the landscape of CH? How are hematological phenotypes affected by both CH-related lesions? How does it
33 affect HM and CVD risks? These are the key questions to be answered for better understanding of CH and its
34 implication in HM and CVD.

35 In the present study, for the purpose of delineating the combined landscape of common driver
36 SNVs/indels and CNAs in CH, we performed SNP array-based copy number analysis and targeted sequencing
37 of major CH-related genes on blood-derived DNA from the Biobank Japan (BBJ)²³, which had been SNP-typed
38 for GWAS studies for common diseases, including hypertension, diabetes, autoimmune diseases and several
39 solid cancers²³. We then investigated the combined effect of both CH-related lesions on clinical phenotypes
40 and outcomes, particularly that on the mortality from HM and CVD.

41

42 **Identification of CH-related SNVs/indels and CNAs**

43 We enrolled a total of 11,234 subjects from the BBJ cohort (n=179,417), in which SNP array analysis of
44 peripheral blood-derived DNA had been performed for large-scale GWAS studies for common diseases
45 (Supplementary Table 1,2) (https://biobankjp.org/info/pdf/sample_collection.pdf)²³. Among these 10,623
46 were randomly selected from 60,787 cases who were aged ≥ 60 years at the time of sample collection and were

47 confirmed not to have solid cancers as of March 2013. This randomly selected set included 61 cases who were
48 known to develop and/or die from HM as of March 2017. The remaining 611 consisted of all cases from the
49 entire BBJ cohort who were confirmed to develop and/or die from HM as of the same date but were not
50 included in the randomly selected 10,623 cases. In total, 672 cases were reported to have HM in the entire BBJ
51 cohort, which included 215 myeloid, 420 lymphoid, and 37 lineage-unknown tumors (Extended Data Fig. 1a).
52 For these 11,234 cases, SNVs/indels in blood were investigated using multiplex PCR-based amplification of
53 exons of 23 CH-related genes, followed by high-throughput sequencing (Online methods).²⁴ Sensitivity of SNV
54 detection according to *in silico* simulations using known SNPs was >94% for 3% variant allele frequency (VAF)
55 and >74% for 2% VAF, but <20% for 1% VAF with a mean depth of ~800x (Supplementary Fig. 1a-b).

56 In total, we called 4,056 SNVs/indels (2,750 SNVs and 1,306 indels) in 3,071 (27.3 %) subjects, of which
57 2,312 (20.6%) had one, 586 (5.2%) two, and 173 (1.5%) ≥ 3 SNV/indels (Fig. 1a). Their VAFs widely distributed
58 from 0.5% to 85.6% with a median of 3.0% (Supplementary Fig. 1c). Age-dependence of CH-related SNVs/indels
59 was evident (Fig. 1b). In accordance with previous reports, *DNMT3A* (13.5%), *TET2* (9.5%), *ASXL1* (2.2%), and
60 *PPM1D* (1.4%) were most frequently mutated (Extended Data Fig. 2a,c). Several combinations of genes,
61 including *TET2/DNMT3A*, *ASXL1/TET2*, *ASXL1/CBL*, *SRSF2/TET2*, and *SRSF2/ASXL1*, were more frequently co-
62 mutated than expected only by chance (OR: 1.53-6.53, $q < 0.05$) (Extended Data Fig. 2d). Of interest, many of
63 these combinations are also co-mutated in myeloid neoplasms with large VAF values²⁵⁻²⁷, suggesting the
64 presence of these combinations of SNVs/indels in the same cell fraction. This was also expected for some cases
65 having a large (>50%) sum of VAFs of relevant SNVs/indels (“pigeonhole principle”),²⁸ although it was not
66 determined whether or not these combinations of SNVs/indels affected the same cell populations in the vast
67 majority of cases (Extended Data Fig. 2e-i).

68 CNAs data were available from the previous study²¹, in which SNP array-based copy number detection
69 in blood-derived DNA was performed for a larger cohort of BBJ cases ($n=179,417$), including all the cases
70 enrolled in the current study ($n=11,234$). In total, 2,797 CNA-positive regions/segments were identified in
71 2,254 (20.1%) cases (Extended Data Fig. 3, Online methods), of which 413 (3.7%) had multiple CNAs (Fig. 1a).
72 Reflecting a higher age distribution of the current cohort, the frequency of CNAs was higher than that in the
73 entire BBJ cohort²¹, even though age-stratified frequencies were almost equivalent between both cohorts (Fig.
74 1b). Estimated mutant cell fractions (MCF) for CNAs were ranged from 0.2% to 93.2% with a median of 2.0%
75 with $FDR < 0.05$, where a substantial number ($n=461$) of CNAs were seen in a cell fraction of $\leq 1\%$, which was
76 below the limit of detection for SNVs/indels. Thus, smaller clones were detected through CNAs, particularly
77 copy-neutral loss-of-heterozygosity (CN-LOH) or uniparental-disomy (UPD), compared with through
78 SNVs/indels (Supplementary Fig. 1c).

79 We found 27 significantly recurrent CNAs, many of which are also commonly seen in HM, supporting a
80 pathogenic link between CH and leukemogenesis (Extended Data Fig. 4a-c). In accordance with previous
81 reports⁸⁻¹¹, 14qUPD, +21q, del(20q), and +15q were among the most frequent CNA lesions (Extended Data Fig.

82 2b,c), while del(20q), 16pUPD, and 17pUPD showed the largest mean clone size (Supplementary Fig. 2). Several
83 CNAs, such as 14qUPD and +21, showed higher frequencies than reported in western populations, which is
84 likely due to a higher sensitivity for detecting CNAs in this study compared with that in previous studies in
85 western populations⁸⁻¹¹; when confined to lesions with $\geq 5\%$ cell fractions, the difference across studies
86 becomes less conspicuous for many CNA targets (Extended Data Fig. 4d,e). Nevertheless, even considering the
87 different sensitivities, several CNAs, including +15, del(14q), del(9q), del(20q) and del(13q), still showed a
88 different frequency across studies in both populations²¹, suggesting an ethnic difference in positive selection
89 of CH-related CNAs (Extended Data Fig. 4e), although the exact genetic basis of the ethnic difference is largely
90 unclear for most CNAs.

91 **Combined landscape of SNVs/indels and CNAs**

92 When SNVs/indels and CNAs were combined, CH was demonstrated in 4,242 (40%) of randomly selected
93 10,623 cases who were ≥ 60 years of age with no reported cancer history and in 376 (56%) of 672 cases who
94 developed HM, where 38 of the 376 were < 60 years old. Combining both lesions, more subjects ($n=1,503$) had
95 two or more lesions than judged by SNVs/indels ($n=759$) or CNA alone ($n=413$) (Fig. 1a). The frequency of CH
96 and the total number of CH-related lesions, as well as the maximum estimate of clone size in CH(+) cases, were
97 significantly larger in individuals with abnormal blood counts, particularly those with cytopenias, compared
98 with those with completely normal blood counts, depending on the number of blood lineages involved (Fig.
99 1c,d). A similar landscape of combined CH-lesions was observed in an independent cohort of 8,023 solid cancer
100 patients from The Cancer Genome Atlas (TCGA), although the sensitivity of CH-lesions, particularly CNAs, was
101 substantially lower than the current study due to a lower coverage of exome sequencing and a less accurate
102 haplotype phasing required for sensitive CNA detection (Extended Data Fig. 5a,b,c).

103 Accounting for 7% of the total cohort and 16% of all CH(+) cases, 740 individuals harbored both types
104 of lesions, which were significantly more frequent than expected only by chance (Extended Data Fig. 2j), even
105 after their age was adjusted (odds ratio [OR]=1.3; $P=0.0003$, age-stratified permutation test) (Online methods).
106 SNVs/indels in *TP53*, *TET2*, *JAK2*, *SF3B1*, and *U2AF1*, and less significantly in *DNMT3A*, *CBL*, and *SRSF2*, was
107 accompanied by significantly more CNAs (Supplementary Fig. 3). The number of cases with multiple CH-related
108 lesions was also significantly larger than expected from the number of all CH-related lesions ($P=0.0067$). The
109 significantly higher frequency of cases with both SNVs/indels and CNAs ($P<0.0001$) and those with multiple
110 lesions ($P<0.0001$) were confirmed in the TCGA cohort. These observations raise a possibility that it might be
111 the total number of lesions, rather than the combination of SNVs/indels and CNAs, that is relevant to the
112 positive selection in CH, in which multiple CH-related lesions in the same cell contributed to positive selection
113 in a substantial number of cases with multiple CH-lesions. In support of this, the maximum clone size in CH(+) cases
114 significantly correlated with the total number of CH-related SNVs/indels and CNAs, but not their
115 combinations per se (Fig. 1e).

116 Co-occurring multiple lesions were judged to be present in the same cell in 73 cases on the basis of their
117 large (>1.0) clone size sum²⁸, of which 8 were combinations between SNVs/indels and CNAs (Extended Data
118 Fig. 2k). In the vast majority of cases, we could not determine the cellular compartment of multiple lesions
119 due to small clone size of both lesions, which would be better addressed using single cell-based sequencing. A
120 representative case was shown in Supplementary Fig. 4, in which the presence of both del(13q) and a *TET2*-
121 involving SNV in the same cell compartment of myeloid lineages was demonstrated using single-cell
122 sequencing (Supplementary Fig. 4a-d). Some combinations of SNVs/indels and CNAs were significantly more
123 frequently observed than expected only by chance (Fig. 2a). Of particular interest among these were co-
124 occurring SNVs/indels and CNAs affecting the same gene/locus. Overall, we found 88 cases having co-occurring
125 SNVs/indels and CNAs affecting 8 genes/loci (Extended Data Fig. 6a), of which most frequently involved were
126 *TP53* (with 17pLOH) (n=24, OR=60.6, $q<0.001$), *TET2* (with 4qLOH) (n=22, OR=10.8, $q<0.001$), *JAK2* (with
127 9pLOH/gain) (n=18, OR: 414, $q<0.001$), and *DNMT3A* (with 2pLOH) (n=16, OR=4.02, $q=0.001$), which were also
128 found in the TCGA cases (Fig. 2a-e, Extended Data Fig. 5d). In reality, more cases are expected to have these
129 combinations, because there were many 'isolated' LOH lesions or allelic imbalances affecting these loci that
130 lacked accompanying SNVs/indels (n=64) (Fig. 2b-e), which were thought to escape from detection due to
131 lower sensitivity of detecting SNVs/indels than CNAs (Supplementary Fig. 1a,c). In fact, using highly sensitive
132 ddPCR assay targeting mutational hotspots, SNVs in *JAK2* and *TP53* were confirmed in 8 out of 48 and 24 out
133 of 41 samples with isolated LOH at 9p and 17p, respectively (Supplementary Fig. 5). Representing well-known
134 mechanisms of biallelic alterations of the relevant driver genes in myeloid malignancies, these combinations
135 of lesions in CH are predicted to affect the same cell, being involved even in very early stages of positive
136 selection in myeloid leukemogenesis²⁹⁻³¹. SNVs/indels were most frequently associated with LOH when they
137 affected *TP53* and *JAK2* in both myeloid malignancies^{32,33} and CH (Extended Data Fig. 6b), also supporting their
138 role in the mechanism of biallelic alterations. Unfortunately, none of these cases satisfied the pigeonhole
139 principle or no samples were available for single cell-sequencing analysis to directly confirm this at a single cell
140 level. However, in the case of SNVs/indels associated with UPD, their presence in the same cell compartments
141 in many cases was supported by a highly skewed distribution of mutant cell fractions of both lesions
142 (Supplementary Fig. 6, Online methods).

143 Besides SNVs/indels and CNAs affecting the same gene/locus, we also detected a significant
144 combination between SNVs/indels in *TET2* and microdeletions of the *TCRA* (14q11.2 involving the) locus (n=7,
145 OR=3.53, $q=0.059$), of which one case was reported to develop T-cell lymphoma (Fig. 2a and Extended Data
146 Fig. 2l). This combination is of potential interest, given that *TET2* is frequently mutated in mature T-cell
147 lymphomas³⁴, particularly in follicular-helper T-cell-derived lymphomas, such as angio-immunoblastic T-cell
148 lymphoma (AITL), which are also seen in *Tet2* knockdown mice³⁵. Other potentially relevant combinations
149 included *SF3B1*/14qUPD, *TET2*/14qUPD, *ASXL1*/1pUPD, *TP53*/1pUPD, and *TP53*/del(5q) (Fig. 2a), whose
150 biological significance, however, is largely unclear except for the interplay between del(5q) and mutated-*TP53*

151 intensively studied in myelodysplastic syndromes (MDS)^{36,37}.

152 **Clinical associations with CH**

153 Next, we investigated common demographic factors that may influence CH-related SNVs/indels and CNAs and
154 the effect of both CH lesions on clinical features and outcomes. In addition to the large effect of age, several
155 factors impacted on CNAs and/or SNVs/indels were observed. Male gender and smoking were significantly
156 associated with SNVs/indels in *ASXL1*, *PPM1D*, splicing factors, and *TP53*, and with CNAs, particularly +15,
157 del(20q), and +21 (with male gender), and 14qUPD (with smoking), many of which remained significant in
158 multivariate analysis (Fig. 3a). The effect of alcohol consumption was less prominent and mostly confined to
159 an increased incidence of del(20q). Although none of the subjects in our cohort had been diagnosed with HM
160 at the time of sample collection, 1,314 cases had varying degrees of abnormal blood counts (Supplementary
161 Table 3). Even though the landscape of CH in these cytopenic individuals at a glance was largely similar to that
162 in non-cytopenic individuals (Extended Data Fig. 7a), cytopenic cases exhibited a significantly higher frequency
163 of CH, where the frequency significantly correlated with the severity of cytopenia (Fig. 1c). In particular,
164 individuals with abnormally high platelet counts had a higher frequency of *JAK2*-involving SNVs/indels and
165 9pUPD (OR=50.5, $q<0.001$ and OR=26.0, $q=0.0017$, respectively), while *U2AF1*-involving SNVs/indels, and
166 del(20q) were more common in those with cytopenia of any sort (OR=7.39, $q<0.001$, and OR=3.10, $q=0.015$,
167 respectively) (Extended Data Fig. 7b). Individuals with CH-related SNVs/indels had a higher frequency of
168 cytopenia and exhibited lower hemoglobin and mean corpuscular hemoglobin concentration (MCHC), while
169 CNAs was associated with lower WBC and platelet counts and larger mean corpuscular volume (MCV) (Fig. 3b).
170 The number of all co-occurring alterations, SNVs/indels or CNAs, and VAF of SNVs/indels predicted significantly
171 lower hemoglobin values, while MCF of CNAs predicted larger MCV and lower MCHC values (Fig. 3b,c,
172 Extended Data Fig. 7c). As for individual alterations, SNVs/indels in *JAK2* were significantly correlated with high
173 platelet counts ($q<0.001$) even when the analysis was confined to the individuals with normal blood counts.
174 Moreover, we found significant associations of lower hemoglobin values with SNVs/indels in *TP53*, *PPM1D*,
175 *SF3B1*, and *U2AF1*, and 4qUPD and del(20q), while SNVs/indels in *PPM1D*, *U2AF1*, 6pUPD, and del(20q) were
176 associated with lower platelet counts and SNVs/indels in *TP53*, and *SF3B1*, and 11qUPD correlated with larger
177 MCV (Fig. 3b, Extended Data Fig. 7c). VAF or cell fractions of SNVs/indels and CNAs were also predictive of the
178 changes in hemoglobin, platelet counts, or MCV (Fig. 3b, Extended Data Fig. 7d). SNVs/indels in *TET2* alone
179 were not associated with a reduced hemoglobin value (Fig. 3b). However, interestingly, we observed a
180 significant association of lower hemoglobin values with multiple SNVs/indels in *TET2* and any allelic imbalance
181 affecting 4q, which is most likely attributable to biallelic *TET2* alterations (Fig. 3b,d). We also tested the
182 relationships between CH and values of other blood tests to reveal a negative correlation between *GNB1*-
183 involving SNVs and uric acid achieved FDR<0.1 (Supplementary Fig. 7).

184 **Effect of SNVs/indels and CNAs on HM mortality**

185 Among the major interests in the current study is the effect of SNVs/indels and CNAs on the risk of HM,
186 particularly the combined effect of both CH-related lesions. To see this, we investigated the effect of CH on the
187 cumulative mortality from HM using the Fine and Gray regression modeling in a case-cohort design³⁸, where
188 7,937 of the 10,623 cases were regarded as a subcohort that were randomly selected from 43,662 cases who
189 had been followed up for survival and cause of deaths on the basis of the vital statistics of Japan³⁹ (Extended
190 Data Fig. 1b). The median follow-up of these cases was 10.4 years (range, 0.01-13.5), during which 401 HM
191 deaths were confirmed (Extended Data Fig. 1b). Age, sex, and versions of SNP array were adjusted and deaths
192 from any causes other than HM were analyzed as competing risks.

193 In accordance with previous reports²², both SNVs/indels and CNAs were significantly associated with a
194 higher mortality from HM than observed in CH(-) cases with an estimated cumulative 10-year mortality of
195 1.28% and 1.32%, respectively (Fig. 4a). Although lymphoid neoplasms accounted for two-thirds of all HM
196 mortality in the cohort of 43,662 elderly cases, attributable mortality in CH(+) vs. CH(-) cases was ~two times
197 higher from myeloid neoplasms (0.39%) than lymphoid (0.21%) neoplasms (Fig. 4b) and the hazard ratio
198 between CH(+) and CH(-) cases was >2.5 times larger for myeloid (3.64) than lymphoid (1.36) neoplasms (Fig.
199 5b). This suggests the predominant effects of CH on myeloid neoplasms, which is in line with the fact that most
200 CH-related lesions targeted driver genes in myeloid neoplasms. The number of SNVs/indels and CNAs and the
201 total number of CH-related lesions all significantly correlated with higher HM mortality (Fig. 4c, Extended Data
202 Fig. 8a,b). While the maximum clone size of CH-related lesions correlated with the number of CH-related
203 lesions (Fig. 1e), the former was also significantly associated with a higher HM mortality independently of the
204 latter (Fig. 4d, Fig. 5a), which was in line with a previous observation that SNVs/indels correlated with
205 development of HM only when they exhibited sufficiently large VAFs⁴⁰. In univariate analysis, the largest risk
206 of HM mortality was conferred by SNVs/indels of *U2AF1*, *EZH2*, *RUNX1*, *SRSF2* and *TP53*, +1q^{11,14,21} (Fig. 5c-d,
207 Supplementary Fig. 8,9). As expected from a ~2 times larger attributable mortality for myeloid than lymphoid
208 malignancy, HRs and ORs were higher in myeloid than lymphoid HM for most of the lesions, with an exception
209 of trisomy 12, which was associated with lymphoid, but not myeloid, neoplasms (Extended Data Fig. 9a). The
210 impact of CH on HM mortality was more prominent when it was present in combination with abnormal blood
211 counts, particularly cytopenia. A significantly higher HM mortality associated with CH was observed in subjects
212 with abnormality in blood counts than in those without (Fig. 4e), depending on the number of CH-related
213 lesions and on the severity of cytopenia; as large as 3.4% 10-year HM mortality was observed for those with
214 multi-lineage cytopenia and multiple CH-related SNVs/indels and CNAs, compared with 0.46% for those with
215 normal blood count lacking CH-related lesions.

216 The presence of both SNVs/indels and CNAs was associated with a significantly increased HM mortality
217 compared with that of SNVs/indels (HR=2.84, 95%CI:2.14-3.78) or CNA (HR=2.64, 95%CI:1.94-3.60) alone (Fig.
218 4f). It was observed even when subjects were stratified according to the number of SNVs/indels (Extended
219 Data Fig. 8c-e). However, the combined effect seems to be explained in large part by an increased total number

220 of alterations, rather than the type of lesions co-occurred, i.e., SNVs/indels vs. CNAs. In fact, the HM mortality
221 significantly correlated with the total number of CH-related lesions and the co-occurrence of both lesions did
222 not significantly affect the mortality of individuals having the same number of lesions (Extended Data Fig. 8f-
223 h). Many of SNVs/indels conferring a higher HM mortality, including those affecting *U2AF1*, *SRSR2*, *TP53*, and
224 *JAK2*, tended to have a higher total number of CH-related lesions, compared with other SNVs/indels (Extended
225 Data Fig. 2m). Nevertheless, the effect on HM mortality was not uniform across different combinations of
226 SNVs/indels and CNAs, regardless of the total number of lesions. In particular, those involving the same
227 gene/locus were associated with a higher HM mortality, compared with other combinations of SNVs/indels
228 and CNAs (Extended Data Fig. 6c). The increase mortality was largely explained by those affecting *TP53*.
229 However, even excluding *TP53*-involving SNVs/indels and CNAs, the combinations of lesions affecting the same
230 locus showed a higher HM mortality than other SNVs/indels and CNAs combinations. Of interest, *TP53*-
231 involving SNVs/indels also exhibited significant associations with del(5q) and multiple (≥ 3) CNAs mimicking a
232 complex karyotype (Fig. 2a, Extended Data Fig. 6f), which together with 17pLOH, are among the most common
233 lesions associated with *TP53* alterations in a variety of myeloid neoplasms with a very poor prognosis,
234 particularly in MDS^{25,33,41}. In agreement with this, these combinations involving *TP53* alterations were
235 significantly associated with a higher mortality from MDS, compared with *TP53*-involving SNVs/indels alone
236 (Extended Data Fig. 6g-h).

237 An almost identical risk estimation for HM was obtained in a case-control setting including all 672 cases
238 who developed HM (Extended Data Fig. 1a, 9a). A small number of cases in which the onset of HM was
239 recorded due to incomplete follow-up and exclusion of MDS and myeloproliferative neoplasms from the follow-
240 up prevented powered analyses of the effect of CH on cumulative incidence of HM, although a similar trend of
241 the effect of CH was observed with regard to the risk of HM that were seen in the analysis using mortality as
242 an endpoint (Extended Data Fig. 9b-f).

243 **Effect of SNVs/indels and CNAs on cardiovascular mortality**

244 Finally, we investigated the combined effect of SNVs/indels and CNAs on cardiovascular mortality in the cohort
245 of 10,623 individuals using multivariate models to take into account known risk factors other than CH: age,
246 gender, body-mass index, comorbidities (diabetes mellitus, hypertension, and dyslipidemia), history of
247 smoking/drinking, and versions of SNP array. In accordance with the previous reports¹³, the presence of CH-
248 related SNVs/indels with large clone size (VAFs $\geq 5\%$) were associated with an elevated cardiovascular and all-
249 cause mortality (HR=1.36, 95%CI:1.09-1.71 for cardiovascular mortality; HR=1.41, 95%CI:1.24-1.60 for all-
250 cause mortality) (Fig. 6a, Extended Data Fig. 10a). In support of this, we observed significant association of
251 SNVs/indels with hypertension (Fig. 6b), which was independent of known risk factors for hypertension,
252 including older age, a higher BMI, and diabetes. By contrast, regardless of their clone size, CNAs alone did not
253 seem to affect cardiovascular or all-cause mortality (Fig. 6c, Extended Data Fig. 10b). However, CNAs in

254 combination with SNVs/indels with $\geq 5\%$ VAFs were significantly associated with elevated cardiovascular
255 mortality and all-cause mortality, compared with CNAs alone, SNVs/indels alone and either SNVs/indels or
256 CNAs (Fig. 6d and Extended Data Fig. 10c), although there was no significant difference in cardiovascular
257 mortality or overall survival depending on whether or not they involved the same locus (Extended Data Fig.
258 6d,e). In multivariate analysis, the combined effect of both lesions was independent of the number of
259 cooccurring SNVs/indels (HR=1.77, $P=0.012$) (Extended Data Fig. 10d-f). Given no impact of CNAs alone, the
260 combined effect on cardiovascular and all-cause mortality does not seem to be explained by an increased total
261 number of CH-related lesions. In fact, the total number of CH-related lesions did not correlate with
262 cardiovascular and all-cause mortality, except for a significantly higher mortality for ≥ 3 CH-related lesions
263 (Extended Data Fig. 10g), likely involving both SNVs/indels and CNAs. Collectively, these observations
264 suggested that the presence of both SNVs/indels and CNAs increased the cardiovascular and all-cause mortality,
265 compared with either of both lesions.

266

267 Discussion

268 Combining targeted deep sequencing of major CH-related genes and SNP array-based copy number analysis of
269 blood-derived DNA from >10,000 individuals aged ≥ 60 years, we have delineated a comprehensive registry of
270 CH in a general population of elderly individuals in terms of both SNV/indel and CNA. A case-cohort study
271 design enabled an accurate estimation of CH-associated cumulative HM mortality in a large general cohort of
272 elderly individuals (>43,000) including >400 cases who developed HM, substantially saving the cost and effort
273 of sequencing, where only ~ 8300 ($\sim 18\%$) individuals/subcohort were fully genotyped. It should be noted that
274 with a much larger number of cases with HM mortality ($n=401$) compared with previous cohort studies (16
275 and 37 cases/cohort)^{12,13}, the estimation of HM mortality in individuals with CH-related SNV/indels was
276 substantially more accurate with a much smaller confidence interval for both myeloid and lymphoid
277 malignancies, where the mortality attributable to CH was mostly explained by myeloid malignancies regardless
278 of type of CH-related lesions. Estimation of odds ratios for CH(+) vs. CH(-) cases were even more accurate with
279 a total of 672 HM events in a case-control study setting.

280 Including both types of lesions, CH was found in as many as 40% of a general population of ≥ 60 years of
281 age, of which 11% had $\geq 10\%$ clone size. As a whole, SNVs/indels and CNAs co-occurred more frequently than
282 expected only by chance. In particular, as repeatedly highlighted in myeloid neoplasms^{29,33,42}, SNVs/indels in
283 *DNTM3A*, *TET2*, *JAK2*, and *TP53*, significantly co-occurred with LOH at each locus in CH, suggesting the role of
284 biallelic alterations of these genes even in an early stage during leukemogenic evolution. Co-occurrence of
285 *TET2*-involving SNVs/indels and deletions involving the *TCRA* locus that are suggestive of evolution of *TET2*-
286 mutated T-cell clones is also of interest. However, even excluding the subjects having these combinations
287 affecting the same gene, SNVs/indels and CNAs significantly co-occurred ($P=0.0042$). Given that most of the
288 CNAs in CH are recurrently seen in myeloid neoplasms, this suggests the presence of functional interactions

289 between CH-related SNVs/indels and CNAs for positive selection, although we cannot exclude a possibility that
290 CNAs might just represent chromosomal instability induced by one or more CH-related SNVs/indels.

291 Compared with those having SNVs/indels or CNAs alone, CH(+) individuals with both lesions showed a
292 higher clone size, more abnormal blood counts, and a higher mortality from HM, particularly of myeloid
293 lineages. The combined effect of SNVs/indels and CNAs⁴⁰, is typically exemplified by biallelic alterations in
294 *DNTM3A*, *TET2*, *JAK2*, and *TP53*, caused by LOH affecting the mutated locus. However, the effect of combined
295 SNVs/indels and CNAs is largely explained by an increased total number of CH-related lesions. Given that the
296 size of CH clones correlated with the number of CH-related lesions, the increasing number of mutations is
297 thought to promote expansion of clones, contributing to an earlier onset and progression of HM. This
298 underscores the importance of measuring both lesions for accurate estimation of HM mortality, which is
299 expected to increase the number of CH-related lesions evaluated only for SNVs/indels and CNAs alone by 0.25
300 and 0.36 on average, revising 10-year expected HM mortality by 0.14% and 0.19%, respectively. The combined
301 effect of both SNVs/indels and CNAs was also observed for cardiovascular and all-cause mortality. Of interest,
302 the effect was seen despite that CNAs alone did not affect the mortality. Because the effect of SNVs/indels on
303 cardiovascular mortality depended on their VAFs, which increased with the presence of CNAs, the combined
304 effect seems to be mediated in part by an increased size of clones having SNVs/indels, although CNA still
305 remained significant after the effects of clone size was adjusted.

306 Potential caveats in the current study include a limited number of CH-related genes analyzed (n=23), a
307 compromised sensitivity of detecting focal CNAs, and the study population exclusively including individuals
308 over 60 years of age. However, these 23 genes, which are estimated to capture ~90% of CH-related
309 SNVs/indels^{12,13}, were analyzed using deep sequencing to sensitively detect lesions in very small fractions (~1%),
310 which would not have been possible with a more unbiased sequencing with a larger target size. In addition,
311 CH and related HM and CVD are highly enriched in and mostly confined to this age group, respectively. Thus,
312 the limited number of genes and age group might not necessarily be the limitations, but rather contributed to
313 efficient analyses of comprehensive analysis of CH-related alterations in a large number of cases to investigate
314 their effects on clinical outcomes at an acceptable cost. However, clearly more comprehensive studies with
315 unbiased sequencing and improved copy number detection including all age groups should be warranted to
316 elucidate the full spectrum of CH-related alterations in future studies.

317

318 **Acknowledgement**

319 This work was supported by the Japan Agency for Medical Research and Development (AMED)
320 (JP15cm0106056h0005, JP19cm0106501h0004, JP16ck0106073h0003, JP19ck0106250h0003 to S.O.;
321 JP17km0405110h0005 and JP19ck0106470h0001 to H.M.; JP19ck0106353h0003 to Y.N.) and the Core
322 Research for Evolutional Science and Technology (CREST) (JP19gm1110011 to S.O.); the Ministry of Education,
323 Culture, Sports, Science and Technology of Japan; the High Performance Computing Infrastructure System
324 Research Project (hp160219, hp170227, hp180198 and hp190158 to S.O. and S.M.) (this research used
325 computational resources of the K computer provided by the RIKEN Advanced Institute for Computational
326 Science through the HPCI System Research project); the Japan Society for the Promotion of Science (JSPS);
327 Scientific Research on Innovative Areas (JP15H05909 to S.O. and S.M.; JP15H05912 to S.M.) and KAKENHI
328 (JP26221308 and JP19H05656 to S.O.; JP16H05338 and JP19H01053 to H.M.; JP15H05707 to S.M.); the Takeda
329 Science Foundation (S.O., H.M. and T.Y.). S.O. is a recipient of the JSPS Core-to-Core Program A: Advanced
330 Research Networks. DNA samples and subjects' clinical data were provided by Biobank Japan, the Institute of
331 Medical Science, the University of Tokyo. The super-computing resource was provided by Human Genome
332 Center, the Institute of Medical Science, the University of Tokyo.

333

334 **Author contributions**

335 R.S., H.M., and S.O. designed the study. K.M., Y.K., T.M., and Y.M. provided DNA samples and clinical data. Y.K.
336 and S.M. provided bone marrow samples. C.T., and Y.K. performed copy-number analysis. Y.M. and M.K.
337 performed sequencing. M.M.N. performed cell sorting and single-cell analysis. R.S., M.M.N., Y.O., T.Y., Y.S, K.C.,
338 H.T., A.N., S.I., and S.M. performed bioinformatics analysis. R.S., Y.N., M.M.N., Y.O., T.Y., H.M., and S.O. prepared
339 the manuscript. All authors participated in discussions and interpretation of the data and results.

340

341 **Online methods**

342

343 **Sample ascertainment**

344 All subjects in this study were derived from BioBank Japan (BBJ) project, a multi-hospital-based-registry²³. BBJ
345 project enrolled approximately 200,000 individuals with at least one of 47 target diseases between fiscal years
346 2003 and 2007. From 179,417 participants of BBJ project in which SNP array analysis of peripheral blood-
347 derived DNA had been performed, we enrolled a total of 11,234 subjects. Among these, 10,623 were randomly
348 selected from 60,787 cases who were aged ≥ 60 years at the time of sample collection and were confirmed not
349 to have solid cancers as of March 2013. Out of the randomly selected 10,623 cases, 61 were recorded to
350 develop or die from HM. The remaining 611 subjects, all of whom were recorded to have HM events, were
351 additionally enrolled to maximize the statistical power in survival analysis. In total, we enrolled 672 subjects
352 with any HM events, 138 and 589 of which were recorded to develop and die from HM, respectively. Subjects'
353 demographic summary was presented in Supplementary Table 1. The numbers of subjects with individual
354 targeted diseases were listed in Supplementary Table 2. The protocols for this study were approved by the
355 ethics committees at all the involved institutions; written informed consent had been obtained from all
356 participants.

357

358 **Multiplex PCR-based targeted sequencing**

359 To detect CH-associated driver mutations, we performed multiplex PCR-based targeted sequencing, as
360 previously described.²⁴ Primers were designed to cover coding regions of 23 driver genes commonly mutated
361 in clonal hematopoiesis or myeloid neoplasms: *ASXL1*, *CBL*, *CEBPA*, *DDX41*, *DNMT3A*, *ETV6*, *EZH2*, *GATA2*, *GNAS*,
362 *GNB1*, *IDH1*, *IDH2*, *JAK2*, *KRAS*, *MYD88*, *NRAS*, *PPM1D*, *RUNX1*, *SF3B1*, *SRSF2*, *TET2*, *TP53*, and *U2AF1*. PCR
363 product sizes were designed to be 180-300 bp to cover the amplicon by the sequencing reads. We added
364 CGCTCTCCGATCTCTG to the 5' end of the forward primers and CGCTCTCCGATCTGAC to the 5' end of the
365 reverse primers to perform second PCR.^{43,44} We performed multiplex PCR using different primer pools to cover
366 all coding regions of the 23 genes. Then we performed second PCR with primer sequences 5'-
367 AATGATACGGCGACCACCGAGATCTACACxxxxxxxACACTCTTCCCTACACGACGCTCTCCGATCTCTG-3' and 5'-
368 CAAGCAGAAGACGGCATAACGAGATxxxxxxxGTGACTGGAGTTCAGACGTGTGCTCTCCGATCTGAC-3', where
369 xxxxxxxx represents 8-bp barcodes. All second PCR products were pooled for one sequencing run. After each
370 library was purified using Agencourt AMPure XP (Beckman Coulter), we obtained 2x150-bp paired-end reads
371 with dual 8-bp barcode sequences on a HiSeq2500 instrument.

372

373 **Calling CH-related SNVs/indels**

374 Sequencing reads were aligned to the human genome reference (hg19) using Burrows-Wheeler Aligner, version
375 0.7.8, with default parameter settings. Mutation calling was performed through our established pipeline, as

376 previously reported^{33,45-47}, using the following parameters.

377 First, we adopted variants fulfilling the following criteria:

- 378 (i) Number of variant reads ≥ 10 (≥ 5 for TCGA dataset) †
- 379 (ii) Variant allele frequency (VAF) $\geq 0.5\%$ †
- 380 (iii) Non-synonymous variants within coding-sequence or splice-site variants

381 († For calculation of read counts and VAFs, we only counted base calls fulfilling Mapping Quality score ≥ 40 ,
382 and Base Quality score ≥ 20 .)

383

384 To further exclude false positive calls due to sequencing artifacts, we modeled site-specific error rates
385 as beta-binomial distribution. Parameters for beta-binomial distribution were determined by maximum
386 likelihood method⁴⁸, based on the read counts in all samples. Mutation calls whose VAFs were significantly
387 deviated from background-error distribution ($P_{\text{beta-binomial}} \leq 10^{-6}$) were regarded as true mutations.

388 Additionally, variants always appeared within similar ranges of VAFs (especially $<1\%$, or $>40\%$) were
389 likely to be sequencing artefacts or germline polymorphisms, rather than true somatic mutations. Based on
390 this assumption, we excluded candidates fulfilling both of the following criteria from the remaining candidates:

391

- 392 (i) Candidates observed in ≥ 5 samples
- 393 (ii) Mean VAF $<1\%$, or $>40\%$, or coefficient of variation of VAFs < 0.5 .

394

395 The candidates fulfilling the quality filter noted above were included in the subsequent analyses if they
396 fulfil one of the following criteria for driver mutations³³:

397

- 398 (i) Candidates resulting in amino-acid substitutions which were registered in the Catalogue of Somatic
399 Mutations in Cancer (COSMIC) v91 databases for ≥ 5 counts
- 400 (ii) Candidates which fulfill the Criteria 1 and at least one of the Criteria 2

401

402 Criteria 1

403 Candidates which were not registered in public databases, including dbSNP138, the 1000 genomes project as
404 of 2014 Oct, Human Genome Variation Database, and The Exome Aggregation Consortium (ExAC).

405

406 Criteria 2

- 407 a) Candidates located on the non-repeat region with VAFs $\geq 4\% < 40\%$ or $\geq 60\% < 96\%$
- 408 b) Nonsense, frameshift, or splice-site candidates
- 409 c) Candidates which were computationally predicted to have negative consequences by SIFT (score < 0.05),
410 PolyPhen-2 (damaging), and MutationAssessor (high or medium)

411

412 Finally, the resulting set of driver mutations were manually reviewed in Integrated Genome Viewer
413 (<http://software.broadinstitute.org/software/igv/>).

414

415 ***In silico* simulation of mutation calling**

416 To benchmark the performance in detection of low-*VAF* mutations, we performed *in silico* simulation. Mixing
417 2 bam files with variable proportions, we diluted 750 heterozygous SNPs and artificially created low-*VAF*
418 mutations (ranging from 0.5% to 5%). Each diluted SNPs were classified into 6 bins according to sequencing
419 depths (x100-x300, x300-x500, x500-x750, x750-x1000, x1000-1500, and x1500-), and sensitivities were
420 calculated separately for the 6 bins. We calculated sensitivity as a fraction of detected variants within all
421 simulated variants:

$$422 \quad SN_{\text{VAF} = x\%} = TP_{\text{VAF} = x\%} / (TP_{\text{VAF} = x\%} + FN_{\text{VAF} = x\%})$$

423 (SN_{VAF = x%}: sensitivity for variants with x% *VAFs*,

424 TP_{VAF = x%}: number of detected SNPs whose *VAFs* were diluted to x%,

425 FN_{VAF = x%}: number of missed SNPs whose *VAFs* were diluted to x%).

426 Together with sensitivity, we calculated specificity by sampling genomic positions without known SNPs (n =
427 5000/simulation). We counted mutation calls on these positions as false positives, and calculated the
428 specificity as follows:

$$429 \quad SP = 1 - FP / N$$

430 (SP: specificity,

431 FP: number of false-positive mutation calls,

432 N: number of sampled genomic positions).

433 To draw receiver operator characteristic (ROC) curves, we calculated sensitivities and specificities for 9 different
434 cutoffs on beta-binomial *P* values (10⁻², 10⁻³, 10⁻⁴, 10⁻⁵, 10⁻⁶, 10⁻⁷, 10⁻⁸, 10⁻⁹, and 10⁻¹⁰).

435

436 **Copy-number analysis**

437 Our analysis pertaining CNAs are based on the result in previous publication²¹, in which blood derived DNA
438 samples from the 11,234 subjects were examined by either of three different versions of microarrays:
439 Illumina Infinium OmniExpress (n=708), Infinium OmniExpressExome v.1.0 (n=3,152), or v.1.2 (n=7,374). For
440 detection of CNAs, we analyzed allele-specific hybridization intensities for the polymorphisms examined by
441 all versions of arrays (n = 515,355). Haplotype phasing and calculation of log R ratio (LRR) and B-allele
442 frequency (BAF) were performed as previously described²¹. Based on long-range haplotype information and
443 LRR/BAF values, we detected allelic imbalances and classified them into duplications, deletions, and UPDs,
444 with false discovery rate around 5%.^{10,21} Because the power to detect allelic imbalances exceeded the power
445 to distinguish UPD from copy-number gain or loss, CNAs were designated as “unclassifiable” when we could

446 not assign them into specific types of CNAs. In the analyses where the exact discrimination between UPD,
447 duplication, or deletion (e.g., lesion-specific analysis in Fig 2a, 3) was relevant, we excluded unclassifiable
448 CNAs from the analysis. Although we cannot calculate precise cell fractions for unclassifiable CNAs, their cell
449 fractions are basically expected to be quite small. Therefore, when we classified CNAs by their cell fractions
450 (e.g., Fig. 4d, 5b, and 6c), unclassifiable CNAs were regarded to be smaller than the thresholds. When we
451 analyzed CNAs in terms of their cell fractions (e.g., Fig. 1d,e, 5a), unclassifiable CNAs were excluded.
452 Otherwise, we did include those unclassifiable CNAs in the analysis (e.g., Fig 1a-c, 2b-e, 3c-d, 4, 6b,d). Based
453 on the detected CNAs, we determined chromosomal regions significantly affected with CNAs by PART
454 (parametric aberration recurrence test)⁴⁹ (<https://www.hgc.jp/~aniida/PART/manual.html>).

455

456 **Definition of abnormalities in blood counts**

457 Subjects fulfilling at least one of the following criteria were considered to have abnormalities in blood counts.

- 458 (i) White blood cells (/μL): ≥10000, or <3000
- 459 (ii) Hemoglobin (g/dL): ≥16.5 (male), ≥16 (female), or <10
- 460 (iii) Hematocrit (%): ≥50
- 461 (iv) Platelet (10000/μL): ≥50, or <10

462 These cutoffs on blood counts were adopted from diagnostic criteria for MDS or myeloproliferative
463 neoplasms⁵⁰. Out of subjects with available counts for all of WBC, hemoglobin, hematocrit, and platelet (n =
464 8,345), 7,031 subjects (84.3%) had normal blood cell counts.

465

466 **Analysis of lineage-sorted samples**

467 Bone marrow Frozen bone marrow was thawed in Dulbecco's Modified Eagle Medium (Sigma-Aldrich)
468 containing 10% of foetal bovine serum (FBS, biosera) and 1% of Penicillin-Streptomycin solution
469 (ThermoFisher). After the cell pellets were washed with PBS containing 2% FBS, the cells were stained with an
470 antibody mix for 20 min, followed by washing with PBS containing 2% FBS and filtered with a 5 mL Round
471 Bottom Polystyrene Test Tube with Cell Strainer Snap Cap (ThermoFisher). We mixed 500 μL of the filtered cell
472 suspension in PBS containing 2% FBS was mixed with 5 μL of Propidium Iodide Staining Solution (BD Bioscience),
473 which was then sorted with the FACSAria III cell sorter (BD Bioscience). The antibodies used in flow cytometry
474 are listed in Supplementary Table 4. For digital droplet PCR (ddPCR) and amplicon-sequencing, we sorted
475 myeloid, erythroid, T cell, and B cell fractions and gDNA was extracted from sorted cells. To detect allelic
476 imbalances in the region of del(13q), amplicon sequencing was performed with custom primers targeting
477 heterozygous SNPs within the deleted region (ThermoFisher, Supplementary Table 5). To detect the A1153V
478 substitution in *TET2*, ddPCR was performed as described below. For single-cell analysis, CD34⁺ cells were sorted.
479 Cells were re-suspended in StemSpan Serum-Free Expansion Medium (STEMCELL Technologies) at 400–1,600
480 cells/μL, which was then applied into Fluidigm C1 platform for combined single-cell gene expression analysis

481 and SNV detection. Detailed methods for single cell analysis are in preparation for publication (shared upon
482 request, Masahiro M Nakagawa, Ryosaku Inagaki, et al.).

483

484 **ddPCR**

485 For ddPCR, predesigned probes were purchased from BioRad. We mixed 50 ng of gDNA with enzymes (ddPCR
486 Supermix for Probes (no dUTP), BioRad) and the probe mix, followed by droplet generation and PCR
487 amplification according to the manufacturer's protocol. Annealing temperatures was set at 55°C. We measured
488 amplified droplets using the QX200 system and QuantaSoft (version 1.7, BioRad). Catalogue numbers of probe
489 mix are shown in Supplemental Table 6.

490

491 **Statistical analysis**

492 All the statistical analyses were performed using the R statistical platform (<https://www.r-project.org/>) v.3.6.1.
493 All statistical tests were two-sided. Benjamini–Hochberg multiple testing correction was applied when
494 appropriate.

495

496 *Age-stratified permutation test for cooccurrences of CH-related alterations*

497 We tested the significance of cooccurrences between SNVs/indels and CNAs under the stratification by subjects'
498 age, because age-dependent frequencies of both CH-related alterations can confound their cooccurrences.
499 First, we stratified subjects into 41 bins according to their age (60, 61, ..., 100 years old) and calculated
500 frequencies of SNVs/indels, CNAs, and their cooccurrences within each bin. In single iteration of permutation,
501 we randomized the status of SNVs/indels and CNAs in all subjects while retaining their frequencies in each age
502 bin. Then, the number of cooccurrences were summed up across all age bins. By repeating this process, we
503 obtained null random distribution of the number of subjects with cooccurring SNVs/indels and CNAs.
504 Comparing the null distribution and the actual number of cooccurrences, we obtained *P* value for significance
505 of cooccurrences between SNVs/indels and CNAs. Significant cooccurrences of multiple CH-related alterations
506 was also demonstrated in a similar way, in which we counted the total number of CH-related alterations within
507 each age bin. In single iteration, these alterations were randomly re-assigned to the subjects retaining the total
508 number of alterations in each bin. Then, the number of subjects to whom multiple alterations were assigned
509 was counted across all bins. *P* value was calculated by comparing the actual number of cases with ≥ 2 alterations
510 and null distribution generated by repeating the process above.

511

512 *Simulation test for cell-level coexistence of SNVs/indels and CNAs involving the same genes*

513 Regarding the combinations of SNVs/indels and UPDs involving the same genes (*DNMT3A*, *TET2*, *TP53*, and
514 *JAK2*), we observed higher VAFs of SNVs/indels than cell fractions of CNAs in 49 of the 51 cases, which
515 suggested they were likely to be acquired in the same cells and resulted in biallelic alterations (Supplementary

516 Figure 6a,b). To examine how many of the 49 cases should be explained by cell-level coexistence of SNVs/indels
517 and UPDs, we performed random simulation on their clone sizes putting a null hypothesis, $H_0(x)$: SNVs/indels
518 and UPDs were independently acquired in at least x cases ($x=3,4,\dots,51$). P value for $H_0(x)$ was calculated
519 assuming VAFs of SNVs/indels and cell fractions of UPDs follows independent distributions (Supplementary
520 Figure 6c-e). We searched for the maximum x with which P value for $H_0(x)$ was below 0.05 to obtain minimum
521 estimate of the number of cases in which cell-level coexistence of SNVs/indels and UPDs was expected
522 (Supplementary Figure 6f).

523

524 *Risk factors for CH*

525 To extract risk factors for CH, we examined correlations between genetic alterations in CH and baseline
526 characteristics of subjects (age, sex, history of smoking and drinking). Information regarding the history of
527 smoking and drinking were based on self-report questionnaires at DNA sampling. First, we performed
528 univariate logistic regressions for presence of genetic alterations. Based on factors significantly correlated with
529 genetic alterations ($q < 0.1$), we then performed multivariate logistic regressions to extract independent risk
530 factors ($P < 0.05$).

531

532 *Effect of CH on blood cell counts*

533 To elucidate effects of genetic alterations on blood cell counts, we examined correlations between genetic
534 alterations and blood cell counts. After Cox-Box transformation of blood counts, linear regressions were
535 performed. To correct for confounding effects, all regressions were performed in multivariate models including
536 age, gender, and versions of SNP array as covariates, in comparison with subjects without detectable CH.

537

538 *Prediction models for hypertension*

539 To elucidate the relationships between CH and hypertension, we performed multivariate logistic regression.
540 Optimal sets of variables were selected by stepwise method from known risk factors and blood test values
541 available for $\geq 70\%$ of the subjects: presence of SNVs/indels and CNAs, age (+10 years), gender, BMI (+5), history
542 of smoking and drinking (based on self-report questionnaires), white blood cell count, red blood cell count,
543 hemoglobin, hematocrit, MCHC, platelet, aspartate aminotransferase (AST), alanine aminotransferase (ALT),
544 lactate dehydrogenase (LDH), creatinine, blood urea nitrogen, total cholesterol, and glucose.

545

546 *Survival analysis*

547 We evaluated the effects of CH-related mutations, CNAs, and their combinations on mortality from HM, all-
548 cause mortality, and cardiovascular mortality. To define mortality from hematologic malignancies, we included
549 diagnoses within ICD10 code groups C81–C96 (malignant neoplasms of lymphoid, hematopoietic and related
550 tissue), D45 (polycythemia vera), D46 (MDS), D47 (other neoplasms of uncertain behavior of lymphoid,

551 hematopoietic and related tissue), and D7581 (myelofibrosis). For CVD, we included I20-I25 (ischemic heart
552 diseases), I48-49 (arrhythmia), I50 (heart failure), I60-I67, I69 (cerebrovascular diseases), I70-I72 (aortic
553 atherosclerosis, aortic aneurysm, aortic dissection), and I74 (peripheral artery diseases). In analysis on all-
554 cause mortality, we performed Cox proportional hazards regression using the R package, “survival”
555 (<http://cran.r-project.org/web/packages/survival/index.html>). In analysis of HM events (mortality or
556 development) or mortality from CVD, we performed competing risk regression based on fine-gray model. In
557 the analysis of events of HM (mortality and development), we applied a case-cohort design to maximize the
558 statistical power as previously described³⁸ (Extended Data Fig. 1b, 9b), including all subjects with HM events
559 within the target cohort. Meanwhile, cardiovascular mortality and overall survival were analyzed in a cohort
560 of the randomly selected 10,623 subjects. To correct for confounding effects, we included subjects’ age, gender
561 and version of SNP array in the multivariate models for events of HM, while age, gender, BMI, presence of
562 diabetes mellitus, hyperlipidemia, and hypertension, history of tobacco smoking and alcohol drinking, and
563 version of SNP array were included in the models for all-cause and cardiovascular mortalities.

564

565 **Data availability**

566 A table of detected somatic mutations is available at https://github.com/RSaikiRSaiki/CH_2020. Clinical data
567 and a list of CNAs can be provided by the BBJ project upon request (<https://biobankjp.org/english/index.html>).
568 All other data will be made available upon request to the corresponding author.

569

570 **Code availability**

571 All computational codes are available upon request to the corresponding author.

572

Fig. 1

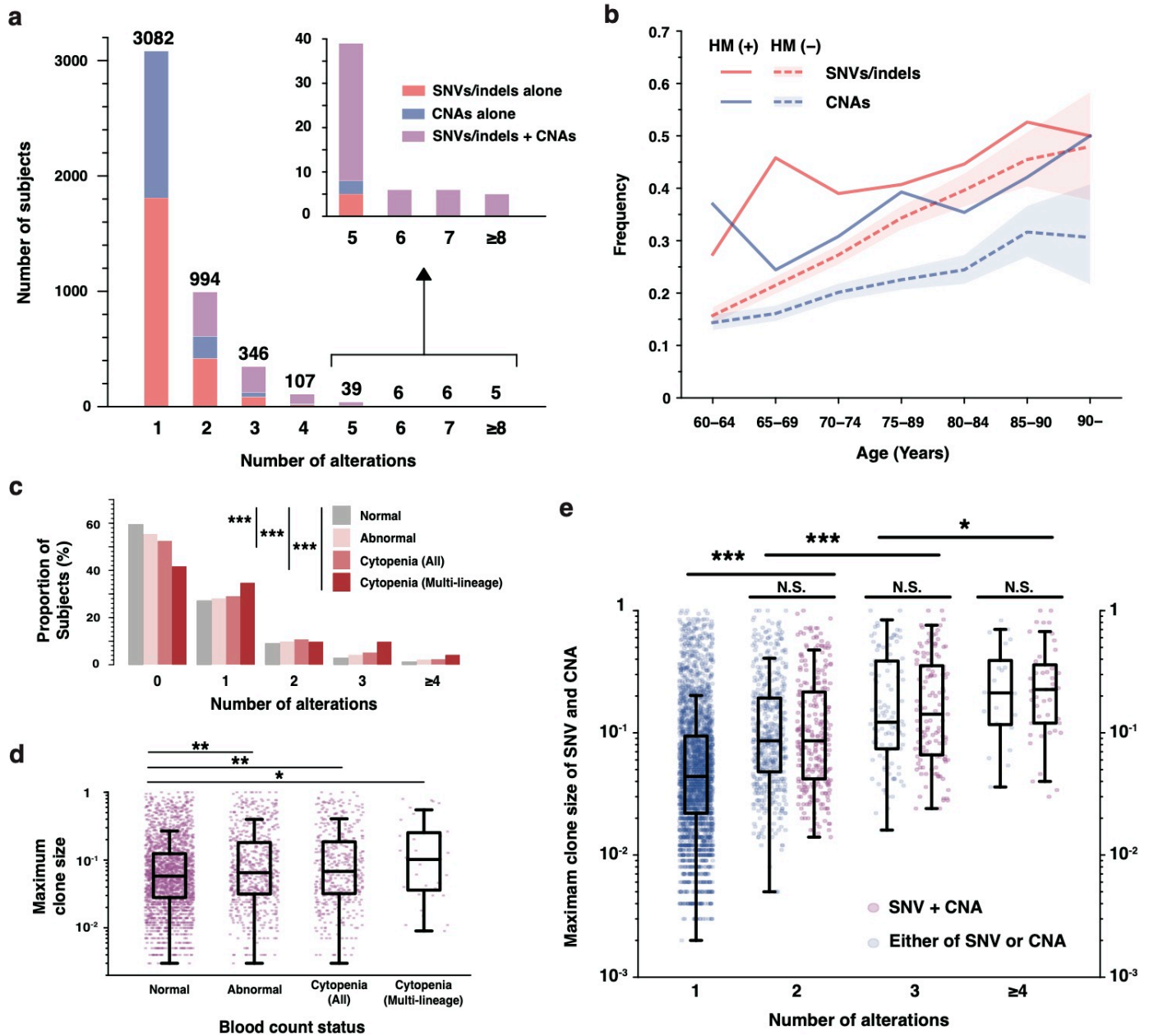
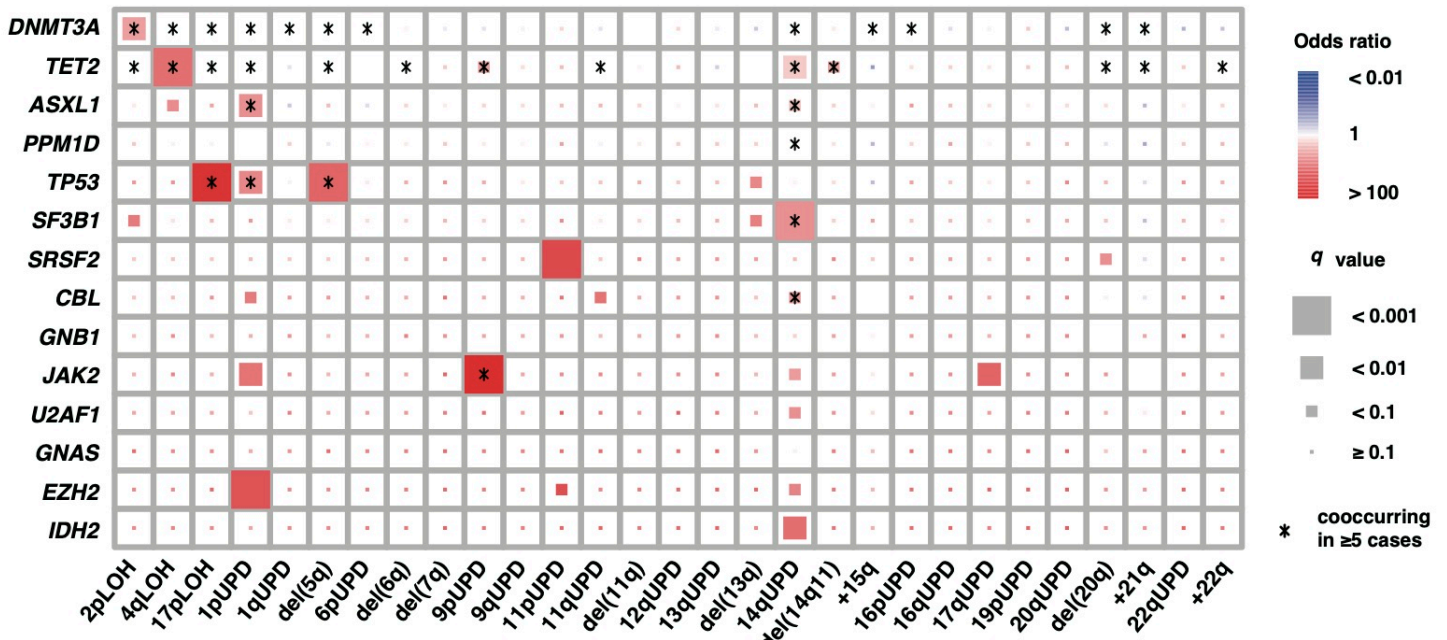


Fig. 1 | Landscape of SNVs/indels and CNAs in clonal hematopoiesis.

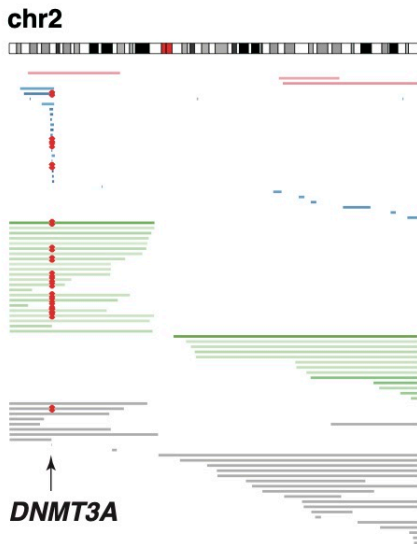
a, Distribution of the number of genetic alterations in each subject. Subjects with SNV/indels alone, with CNAs alone, or with both of them are illustrated by different colors. b, The prevalence of CH-related SNVs/indels and CNAs in subjects with or without HM events, according to age. Colored bands represent the 95% confidence intervals. c, Number of cooccurring alterations in those with subjects with abnormalities in blood cell counts, or cytopenia. d, Maximum cell fraction of CH-related alterations in subjects with abnormalities in blood cell counts, and in those with cytopenia. e, Dotplot of maximum cell fractions of SNVs/indels or CNAs across different numbers of cooccurring alterations. Cell fractions of SNVs/indels are defined as 2 times VAF. Those with both of SNVs/indels and CNAs are shown in purple, while those with either are shown in blue. In panel (d,e), unclassifiable CNAs were excluded because we cannot calculate their precise cell fractions. The box plots indicate the median, first and third quartiles (Q1 and Q3) and whiskers extend to the furthest value between Q1 – 1.5×the interquartile range (IQR) and Q3 + 1.5×IQR. *P* values were calculated by wilcoxon rank-sum test. N.S., not significant; *, *P*<0.05; **, *P*<0.001; ***, *P*<0.0001.

Fig. 2

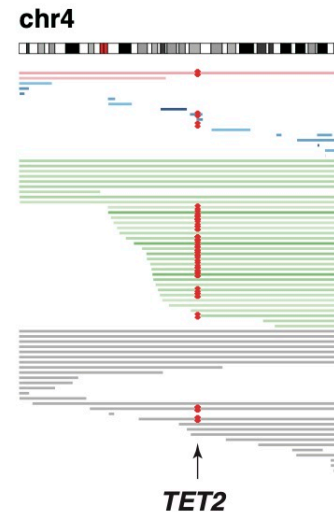
a



b



c



d



e

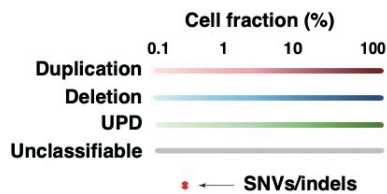


Fig. 2 | Cooccurrences of SNVs/indels and CNAs in clonal hematopoiesis.

a, The correlations between individual SNVs/indels and CNAs. The size of rectangles indicate the significance of correlations. Red rectangles represent positive correlations while blue rectangles represent negative correlations. Combinations of SNVs/indels and CNAs seen in 5 or more subjects are indicated by asterisks. b-e, The distributions of CNAs on chromosome 2 (b), 4 (c), 9 (d), and 17 (e). Horizontal bars represent CNAs, and cooccurring SNVs/indels in *DNMT3A*, *TET2*, *JAK2*, and *TP53* are indicated by red asterisks. Colors of horizontal bars represent the types and cell fractions of CNAs. Allele imbalances which cannot be classified into any of UPD, deletion, or duplication are indicated as unclassifiable CNAs (gray).

Fig. 3

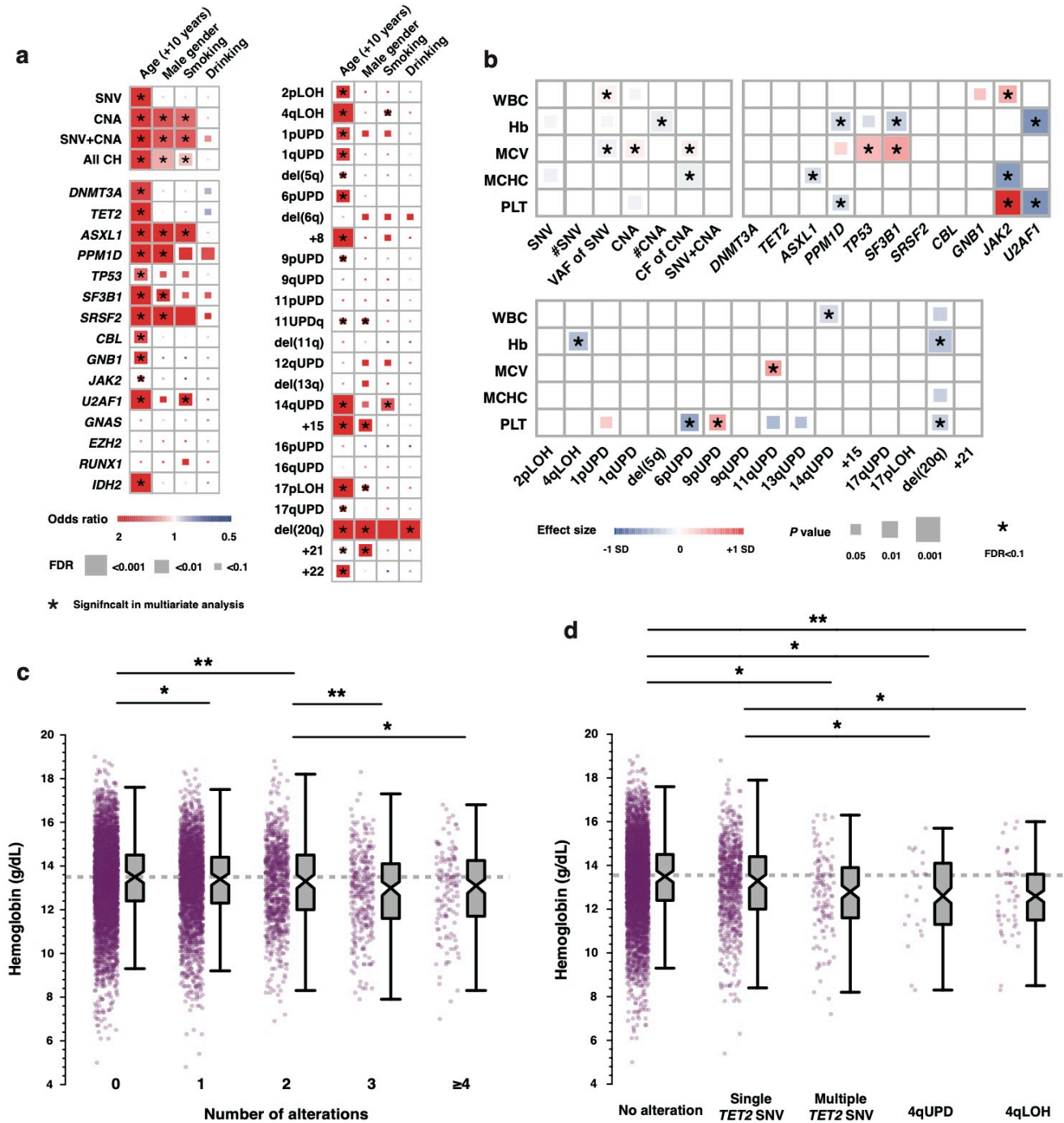


Fig. 3 | Risk factors for CH and effects on blood counts.

a, Correlations of genetic alterations with age, male gender, history of smoking and drinking. Sizes and colors of rectangles represent the significance and effect size, respectively. Asterisks indicate the clinical factors significantly correlated with each alteration in multivariate logistic regression ($P < 0.05$). c, Correlations between genetic alterations and blood counts. The sizes and colors of rectangles indicate the significance, and effect size of correlations. Correlations significant after correction for multiple testing (FDR < 0.1) are indicated by asterisks. WBC: white blood cell, Hb: hemoglobin, MCV: mean corpuscular volume, MCHC: mean corpuscular hemoglobin concentration, Plt: Platelet. c, Distributions of hemoglobin in subjects with different number of alterations. d, Distributions of hemoglobin in subjects with no alterations, with single SNV/indels in *TET2* (Single *TET2* SNV), multiple SNVs/indels (Multiple *TET2* SNV) in *TET2*, with 4qUPD, or with any allelic imbalances in 4q (4qAI) are illustrated in dot plots and boxplots. In all box plots, the median, first and third quartiles (Q1 and Q3) are indicated and whiskers extend to the furthest value between $Q1 - 1.5 \times \text{IQR}$ and $Q3 + 1.5 \times \text{IQR}$. * $P < 0.05$; ** $P < 0.01$.

Fig. 4

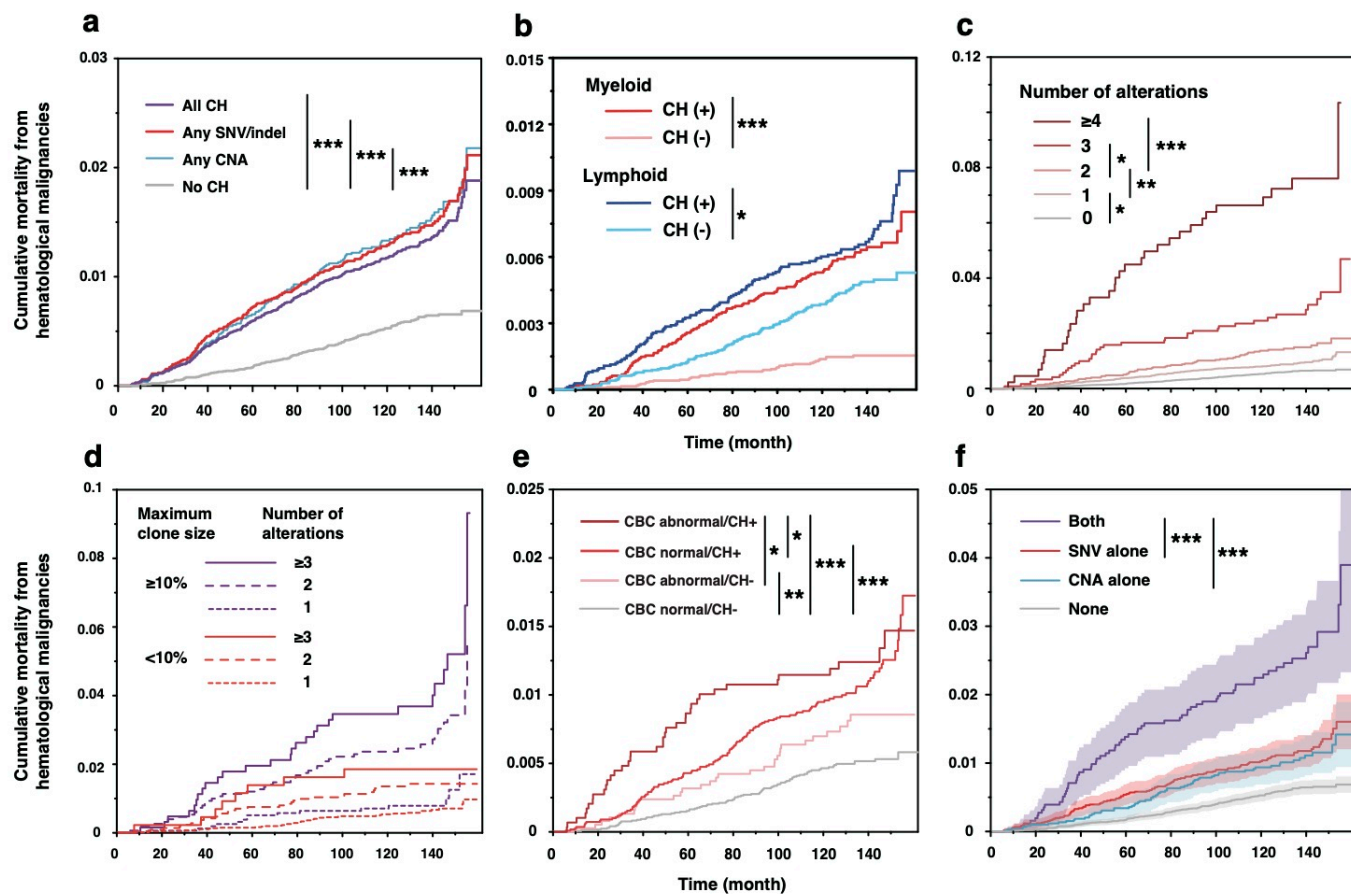


Fig. 4 | Impact of CH on mortality from hematological malignancies.

a, Cumulative mortality from HM in subjects with any CH, any SNV/indel, any CNA, or without CH are shown. b, Cumulative mortality from myeloid and lymphoid malignancies in subjects with or without CH are shown. c, Cumulative mortality from HM in subjects with different numbers of CH-related alterations. d, Cumulative mortality from HM in subjects with different numbers of cooccurring alterations and maximum clone sizes (<10% or ≥10%). Cell fractions of unclassifiable CNAs were regarded to be smaller than 10%. e, Cumulative mortality from HM in subjects with or without abnormalities in complete blood counts (CBC) and/or CH. f, Cumulative mortality from HM in subjects with both SNV/indels and CNA, SNV/indels alone, CNAs alone, and without any alterations are shown. Colored bands indicate 95% confidence intervals. * $P < 0.05$, ** $P < 0.001$, *** $P < 0.0001$.

Fig. 5

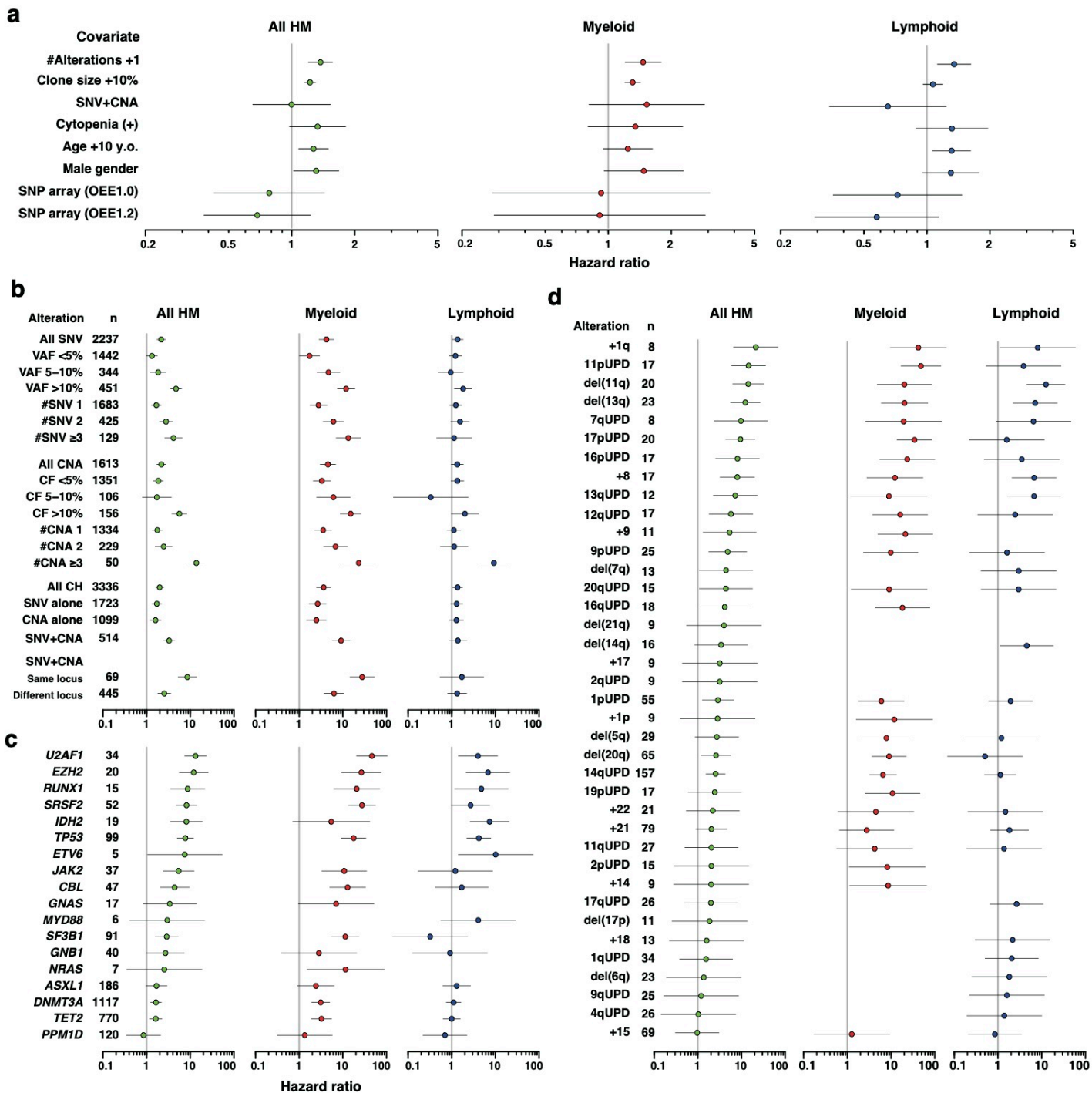


Fig. 5 | Impact of CH-related alterations on mortality from HM.

a-d, Hazard ratios for mortality from All hematological malignancies (All HM), myeloid neoplasms, and lymphoid neoplasms. Error bars indicate 95% confidence intervals. In panel (a), hazard ratios of the indicated covariates calculated by multivariate Fine-Gray regression are shown. In panel (b-d), hazard ratios of the indicated alterations are calculated in comparison with subjects without any alteration. Hazard ratios are not shown for alterations without any event. Cell fractions of unclassifiable CNAs were regarded to be zero in panel (a), and smaller than 5% in panel (b). n, number of cases with the indicated alterations; #Alteration, additional one alteration; Clone size +10%, 10% increase in cell fraction; SNV+CNA, cooccurrence of both SNVs/indels and CNAs; #SNV, number of SNVs/indels; CF, cell fraction of CNAs; #CNA, number of CNAs.

Fig. 6

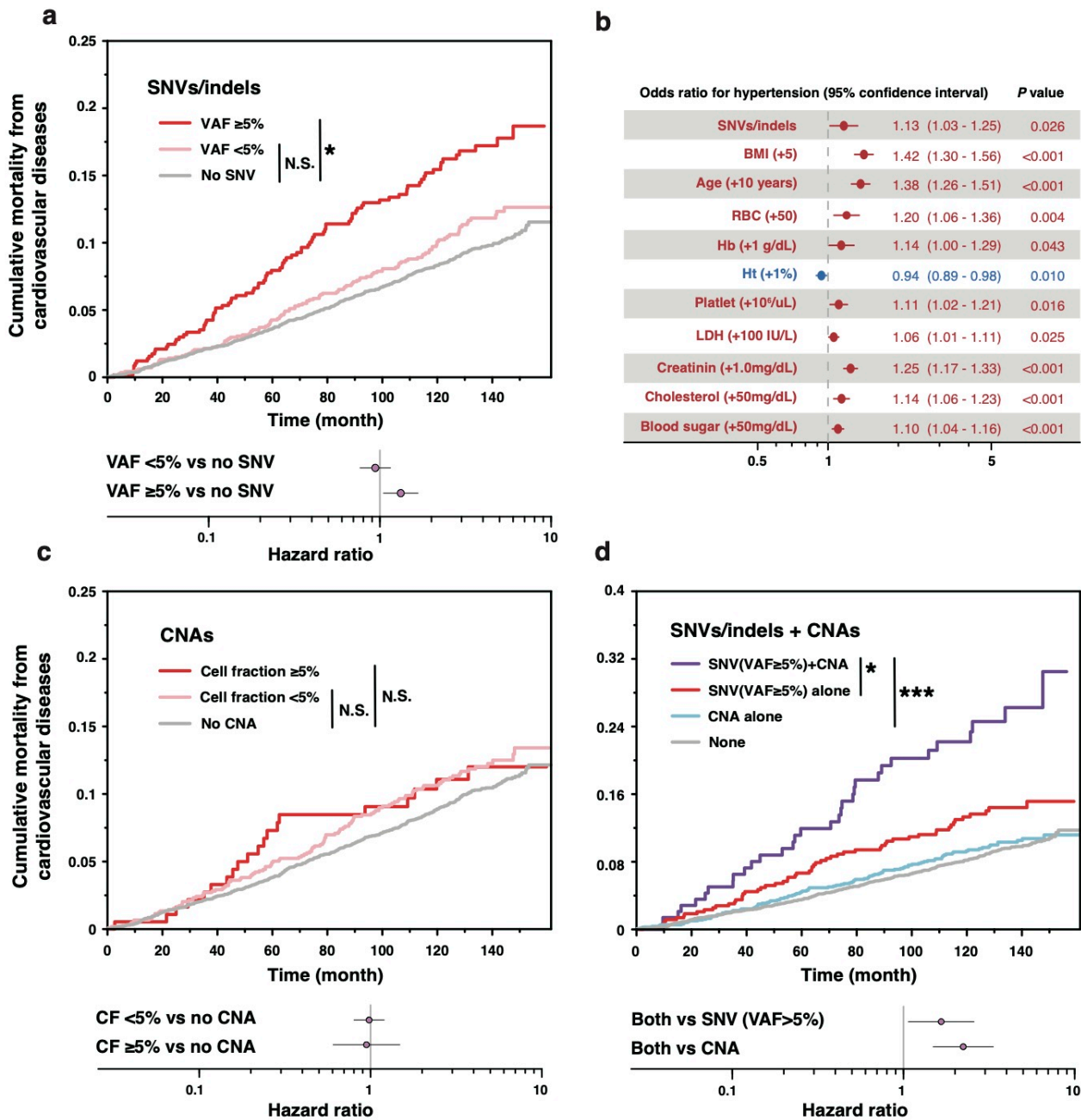


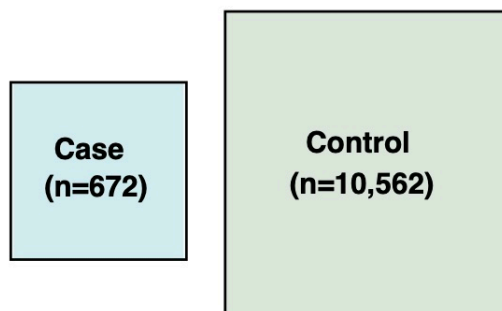
Fig. 6 | Effect of SNV/indels and CNAs on cardiovascular mortality.

a, Cardiovascular mortality in subjects with SNV/indels (VAF $\geq 5\%$ or $< 5\%$), and those without SNV/indels. Hazard ratios and P values are calculated in comparison with those without SNV/indels. b, Results of multivariate logistic regressions for the presence of hypertension. Explanatory variables were selected by stepwise method from following factors: presence of SNV/indels, CNAs, age (+10 years), male gender, BMI (+5), history of drinking and smoking, presence of diabetes mellitus, hyperlipidemia, hypertension, and 13 blood test values available in $\geq 70\%$ of the subjects. Only remaining variables are shown. c, Cardiovascular mortality in subjects with CNAs (cell fraction $\geq 5\%$ or $< 5\%$), and those without CNAs. Hazard ratios and P values for subjects are calculated in comparison with those without CNAs. d, Cardiovascular mortality in subjects with both SNV/indels (VAF $\geq 5\%$) and CNAs (purple), with SNV/indels (VAF $\geq 5\%$) alone (red), with CNAs alone (blue), and without SNV/indels (VAF $\geq 5\%$) or CNAs (gray). Hazard ratios and P values were calculated by comparing those with both SNV/indels (VAF $\geq 5\%$) and CNAs with those with SNV/indels (VAF $\geq 5\%$) alone, or CNAs alone. In (a), (b) and (d), all comparisons were performed in multivariate models including age, gender, body-mass index, comorbidities (diabetes mellitus, hypertension, and dyslipidemia), history of smoking/drinking, and the versions of SNP array.

Extended Data Fig. 1

a

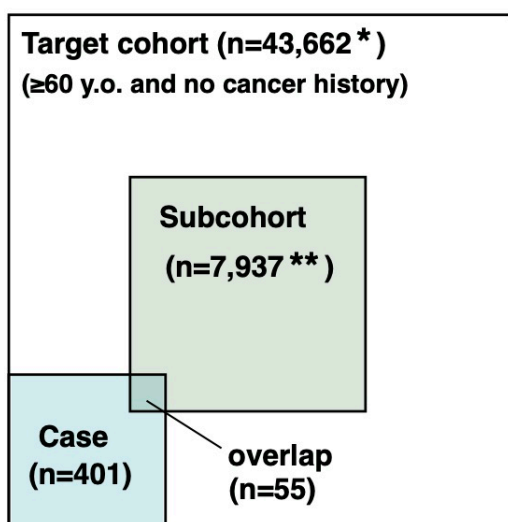
Case-control study for all HM



| | CH(+) | | | | CH(-) | Total |
|----------------|-----------|-----------|------|-----------|-------|--------|
| | SNV alone | CNA alone | Both | All CH(+) | | |
| Case | 154 | 115 | 107 | 376 | 296 | 672 |
| Myeloid | 53 | 41 | 66 | 160 | 55 | 215 |
| AML | 32 | 12 | 19 | 63 | 27 | 90 |
| MDS | 16 | 25 | 34 | 75 | 25 | 100 |
| MPN | 1 | 1 | 2 | 4 | 1 | 5 |
| CML | 1 | 1 | 5 | 7 | 2 | 9 |
| Others | 3 | 2 | 6 | 11 | 0 | 11 |
| Lymphoid | 90 | 69 | 32 | 191 | 229 | 420 |
| B-NHL | 61 | 44 | 18 | 123 | 143 | 266 |
| T-NHL | 4 | 7 | 4 | 15 | 17 | 32 |
| CLL | 3 | 2 | 2 | 7 | 0 | 7 |
| ALL | 4 | 3 | 0 | 7 | 12 | 19 |
| MM/PCT | 17 | 12 | 7 | 36 | 53 | 89 |
| Others | 1 | 1 | 1 | 3 | 4 | 7 |
| Linage Unknown | 11 | 5 | 9 | 25 | 12 | 37 |
| Control | 2,177 | 1,399 | 633 | 4,209 | 6,353 | 10,562 |
| Total | 2,331 | 1,514 | 740 | 4,585 | 6,649 | 11,234 |

b

Case-cohort study for HM death



- * Among 60,787 cases aged ≥ 60 years and confirmed not to have solid cancers as of March 2013, 43,662 had the follow up data for survival.
- ** Among 10,623 cases randomly selected from the 60,787 cases, 7,937 had the follow up data for survival.

Subcohort

| | CH(+) | | | | CH(-) | Total |
|-------------------------------------|-----------|-----------|------|-----------|-------|-------|
| | SNV alone | CNA alone | Both | All CH(+) | | |
| Hematological malignancy (+) | 14 | 11 | 7 | 32 | 23 | 55 |
| Myeloid | 5 | 5 | 4 | 14 | 5 | 19 |
| AML | 1 | 0 | 1 | 2 | 3 | 5 |
| MDS | 3 | 4 | 3 | 10 | 1 | 11 |
| MPN | 0 | 0 | 0 | 0 | 0 | 0 |
| CML | 1 | 0 | 0 | 1 | 1 | 2 |
| Others | 0 | 1 | 0 | 1 | 0 | 1 |
| Lymphoid | 8 | 6 | 3 | 17 | 18 | 35 |
| B-NHL | 6 | 4 | 3 | 13 | 14 | 27 |
| T-NHL | 1 | 0 | 0 | 1 | 3 | 4 |
| CLL | 0 | 0 | 0 | 0 | 0 | 0 |
| ALL | 0 | 0 | 0 | 0 | 0 | 0 |
| MM/PCT | 1 | 2 | 0 | 3 | 1 | 4 |
| Others | 0 | 0 | 0 | 0 | 0 | 0 |
| Linage Unknown | 1 | 0 | 0 | 1 | 0 | 1 |
| Hematological malignancy (-) | 1,614 | 1,036 | 447 | 3,097 | 4,785 | 7,882 |
| Total | 1,628 | 1,047 | 454 | 3,129 | 4,808 | 7,937 |

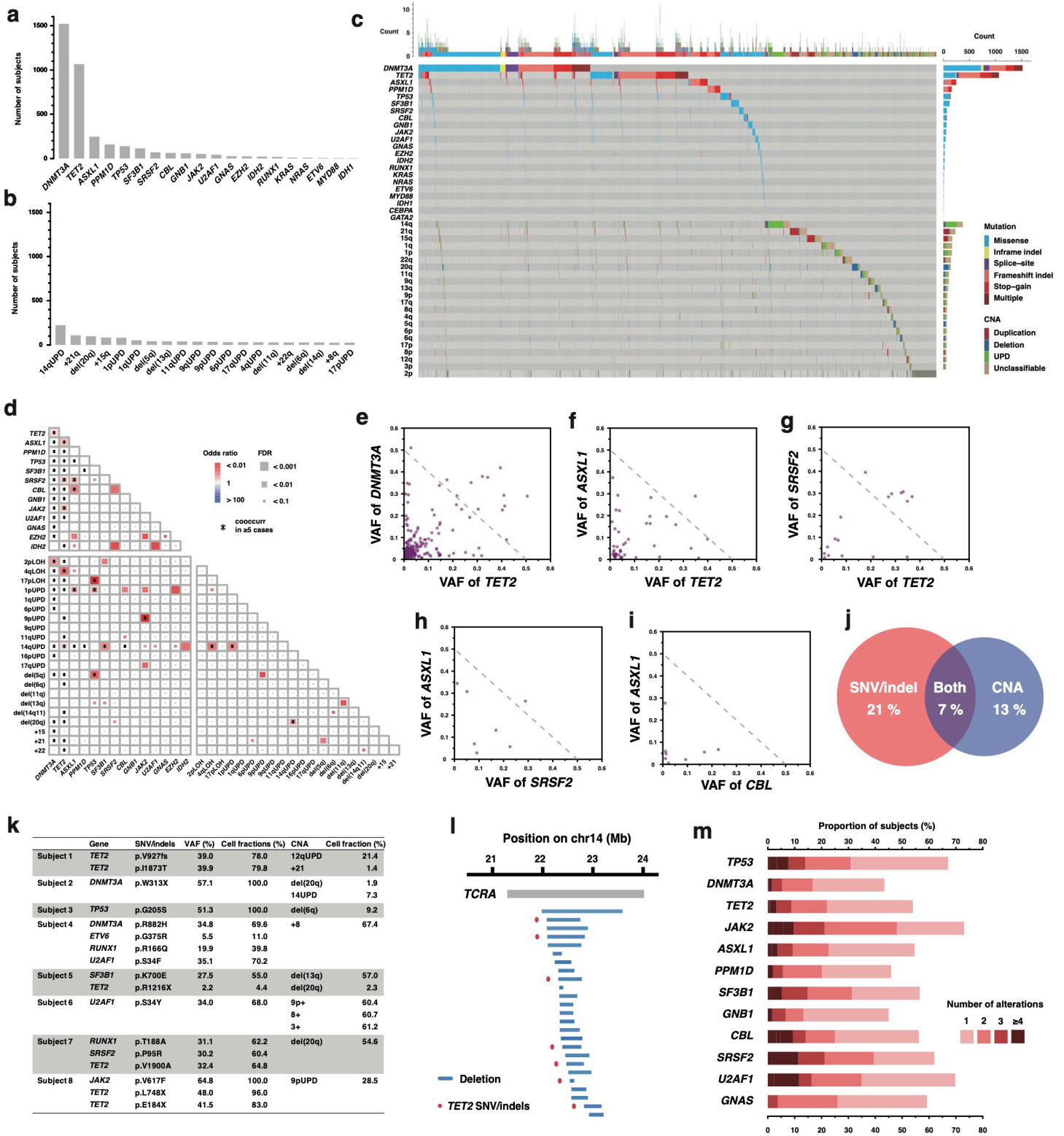
Case (Death from HM)

| | CH(+) | | | | CH(-) | Total |
|-------------------------------------|-----------|-----------|------|-----------|-------|-------|
| | SNV alone | CNA alone | Both | All CH(+) | | |
| Hematological malignancy (+) | 109 | 63 | 67 | 239 | 162 | 401 |
| Myeloid | 41 | 24 | 42 | 107 | 39 | 146 |
| AML | 24 | 8 | 8 | 40 | 20 | 60 |
| MDS | 14 | 13 | 23 | 50 | 17 | 67 |
| MPN | 0 | 1 | 2 | 3 | 0 | 3 |
| CML | 1 | 1 | 5 | 7 | 2 | 9 |
| Others | 2 | 1 | 4 | 7 | 0 | 7 |
| Lymphoid | 62 | 38 | 22 | 122 | 122 | 244 |
| B-NHL | 38 | 25 | 11 | 74 | 74 | 148 |
| T-NHL | 3 | 3 | 3 | 9 | 12 | 21 |
| CLL | 3 | 1 | 2 | 6 | 0 | 6 |
| ALL | 3 | 1 | 0 | 4 | 5 | 9 |
| MM/PCT | 13 | 8 | 4 | 25 | 28 | 53 |
| Others | 2 | 0 | 2 | 4 | 3 | 7 |
| Linage Unknown | 6 | 1 | 3 | 10 | 1 | 11 |
| Hematological malignancy (-) | 0 | 0 | 0 | 0 | 0 | 0 |
| Total | 109 | 63 | 67 | 239 | 162 | 401 |

Extended Data Fig. 1 | Design of case-control and case-cohort study.

a, Design of case-control study (Left). Diagnosis of hematological malignancies (HM) in subjects with or without CH enrolled in the case-control study (Right). b, Design of case-cohort study for death from HM (Left). Diagnosis of HM in subjects with or without CH enrolled in the case-cohort study (Right). AML, acute myeloid leukemia; MDS, myelodysplastic syndromes; MPN, myeloproliferative neoplasms; CML, chronic myeloid leukemia; B-NHL, B-cell non-Hodgkin lymphoma; T-NHL, T-cell non-Hodgkin lymphoma; CLL, chronic lymphoid leukemia; ALL, acute lymphoblastic leukemia; MM, multiple myeloma; PCT, plasma cell tumor.

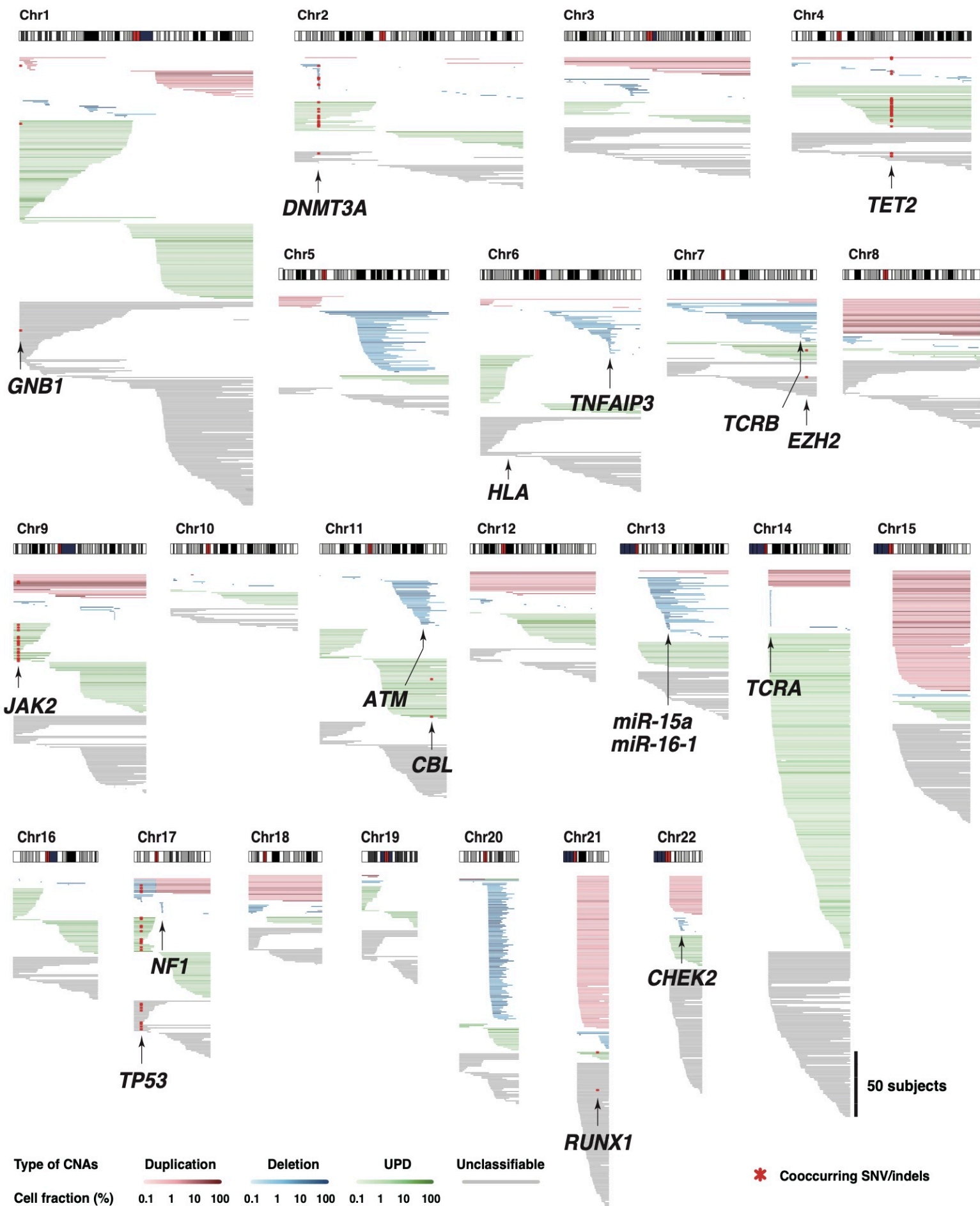
Extended Data Fig. 2



Extended Data Fig. 2 | Landscape of genetic alterations in CH.

a-b, The number of subjects with individual SNVs/indels (a) and CNAs (b). The vertical axis represents the number of subjects with indicated alterations. Unclassifiable CNAs are not included in (b). c, Landscape of SNVs/indels and CNAs in 11,234 subjects. Those without CH-related alterations are omitted. d, The correlations between individual genetic alterations. Combinations seen in 5 or more cases are indicated by asterisks. e-i, VAF of cooccurring SNVs/indels in diagonal plot. Dots above the dashed line fulfill "pegeonhall principle". j, Venn diagram illustrating the overlap between subjects with SNVs/indels and those with CNAs. Frequencies within all subjects in whom SNVs/indels and CNAs were examined (n=11,234) are indicated. k, Subjects in whom cooccurring SNVs/indels and CNAs were suspected to coexist in the same cells on the basis of "pegeonhall principle." l, A magnified illustration of microdeletions around *TCRA* locus (14q11.2). A gray bar represents gene body of *TCRA*. Blue horizontal bars represent microdeletions. Cooccurring *TET2* SNVs are indicated by red dots. Genomic coordinates in hg19 are indicated above. m, Proportions of subjects with different number of cooccurring alterations within those who harbor SNVs/indels in the indicated genes. The proportions of subjects with 1, 2, 3, and ≥ 4 CNAs are depicted by different colors.

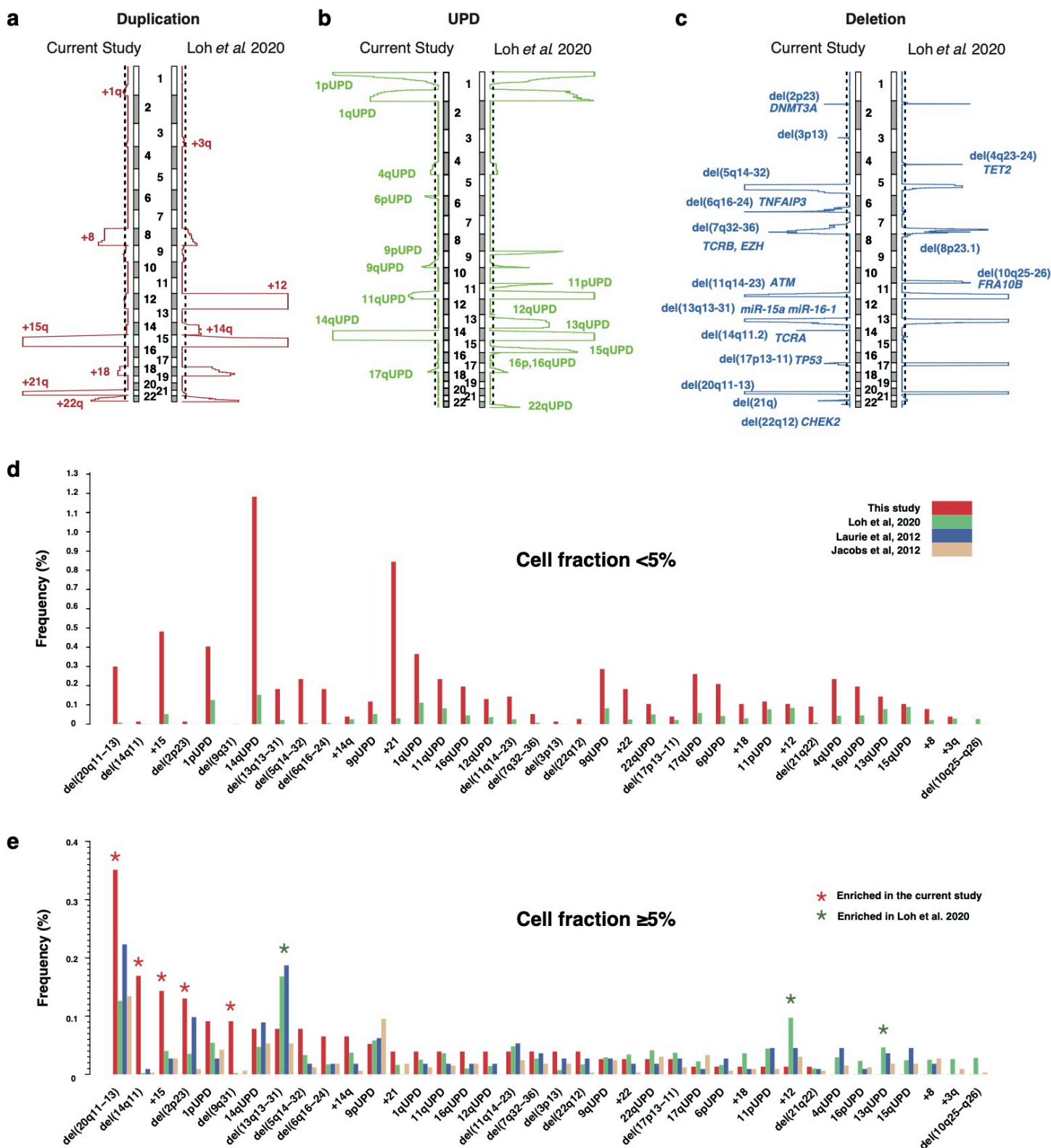
Extended Data Fig. 3



Extended Data Fig. 3 | Distribution of CNAs in all chromosomes.

Distributions of CNAs on all chromosomes are illustrated. Loci of known driver genes are indicated by arrows. Each horizontal bar represents one CNA. Cooccurring SNV/indels are indicated by red dots. Types of CNAs are depicted by different colors as indicated in the annotations.

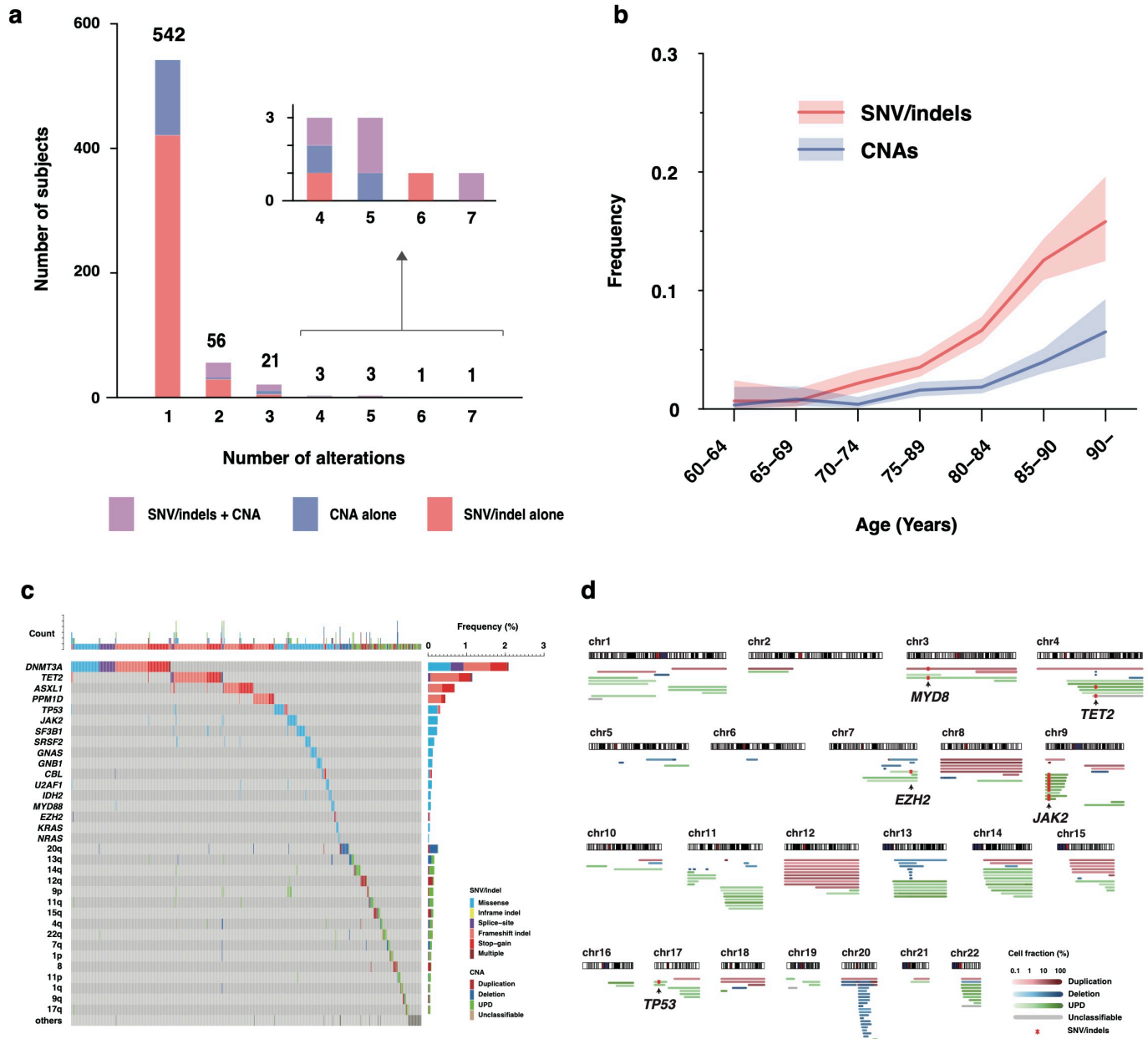
Extended Data Fig. 4



Extended Data Fig. 4 | Chromosomal regions significantly affected by CNAs.

a-c, Chromosomal regions significantly affected by duplications (a), UPDs (b), and deletions (c) in Japanese cohort (current study) and in British cohort¹¹. Statistical significance for recurrence of CNAs were evaluated by PART⁴⁴. Dashed lines indicate thresholds for statistical significance ($q = 0.25$). d-e, Comparison of frequencies of individual CNAs between the current and previous studies¹¹. Comparisons were performed in those aged 60-75 years. In (d) or (e), CNAs in <5% or ≥5% cell fractions were taken into account, respectively. CNAs significantly enriched in either cohort were indicated by asterisks ($q < 0.1$) in (e).

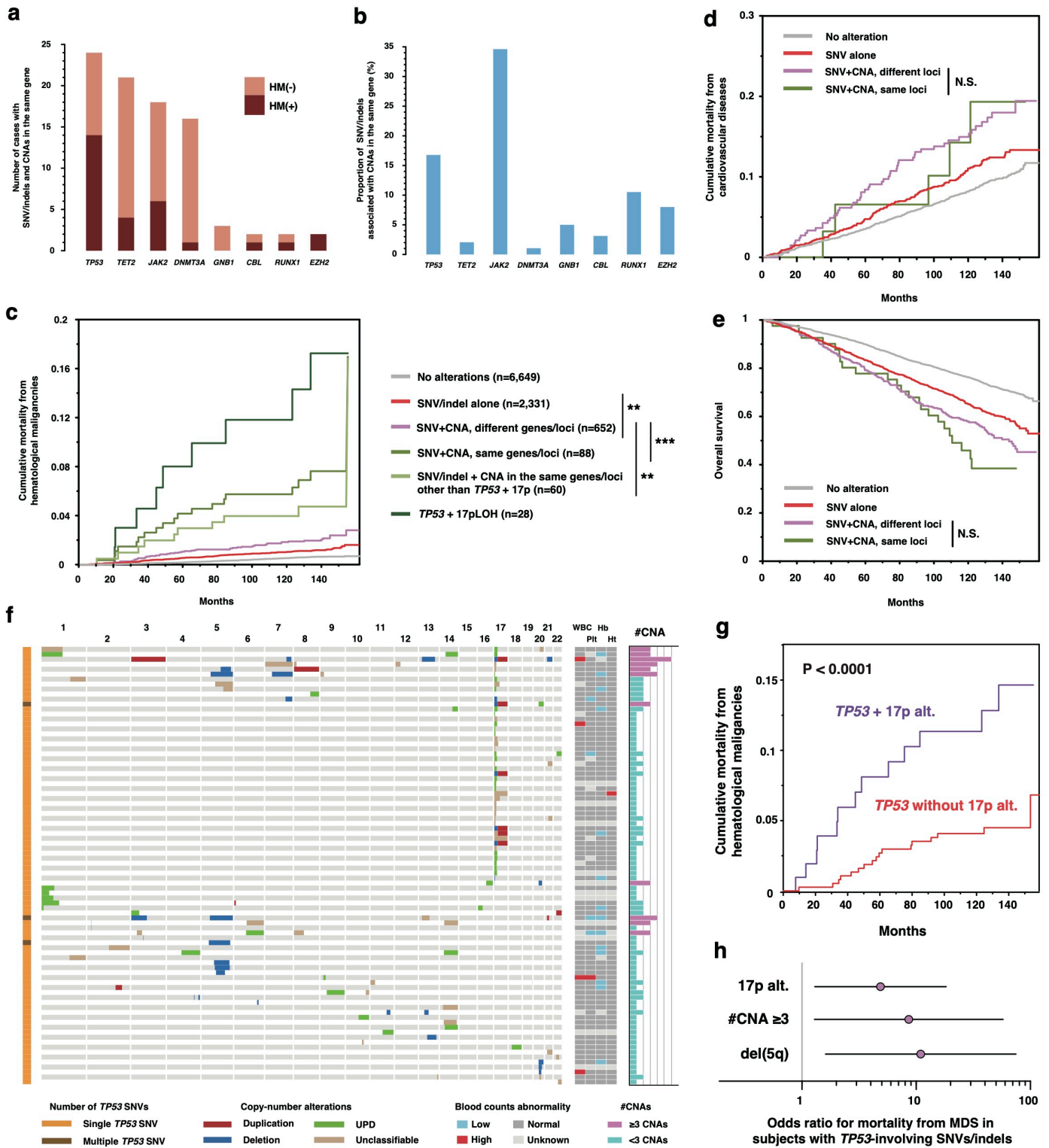
Extended Data Fig. 5



Extended Data Fig. 5 | Analysis of SNVs/indels and CNAs in peripheral blood samples in TCGA cohort.

a, Distribution of the number of genetic alterations in each subject. Subjects with SNVs/indels alone, with CNAs alone, or with both of them are illustrated by different colors. b, The prevalence of CH-related SNVs/indels and CNAs, according to age. Colored bands represent the 95% confidence intervals. c, The landscape of CH-related SNVs/indels and CNAs. Each row represents genetic alterations or affected chromosomal arms, and each column represents subjects. Subjects without any alterations are omitted. Types of SNVs/indels and CNAs are depicted by different colors. d, Distributions of CNAs on all chromosomes are illustrated. Loci of cooccurring SNVs/indels are indicated by arrows. Each horizontal bar represents one CNA. Cooccurring SNVs/indels are indicated by red asterisks. Types of CNAs are depicted by different colors.

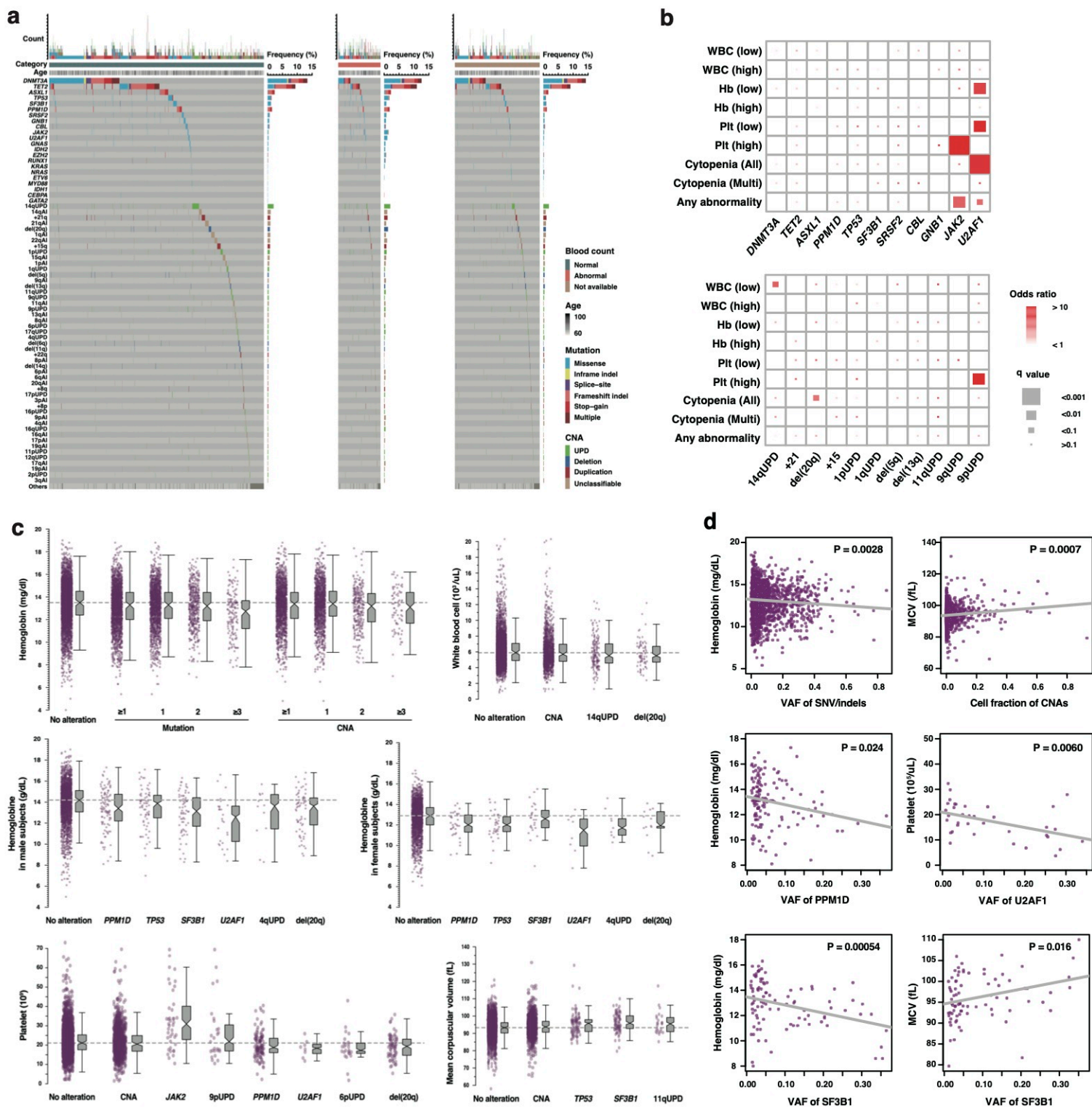
Extended Data Fig. 6



Extended Data Fig. 6 | Interplay between SNVs/indels and CNAs

a, Number of subjects with SNVs/indels and CNAs involving the same genes/loci. b, Proportion of SNVs/indels associated with CNAs in the same genes/loci. c, Cumulative mortality from hematological malignancies. d, Cumulative mortality from cardiovascular diseases. e, Survival curves for overall survival. f, Profiles of CNAs in subjects with SNV/indels in *TP53*. Abnormally high or low blood counts (WBC, Platelet, Hemoglobin, and Hematocrit) are indicated by red or blue, respectively. Numbers of cooccurring CNAs are indicated on the right side (#CNA), where subjects with ≥ 3 CNAs were highlighted by purple. Subjects without any CNA were abbreviated. g, Mortality from hematological malignancies in *TP53*-mutated cases with or without CNAs in 17p. h, Odds ratio for mortality from MDS calculated by multivariate logistic regression in subjects with *TP53*-involving SNVs/indels. We included unclassifiable CNAs involving 17p in 17p alterations (17p alt.) in panel (g-h) because they are most likely to be LOH (UPDs or deletions). *TP53*-involving SNVs/indels in panel (f-h) included those detected by ddPCR (Supplementary Fig. 3). N.S., not significant; **, $P < 0.001$; ***, $P < 0.0001$.

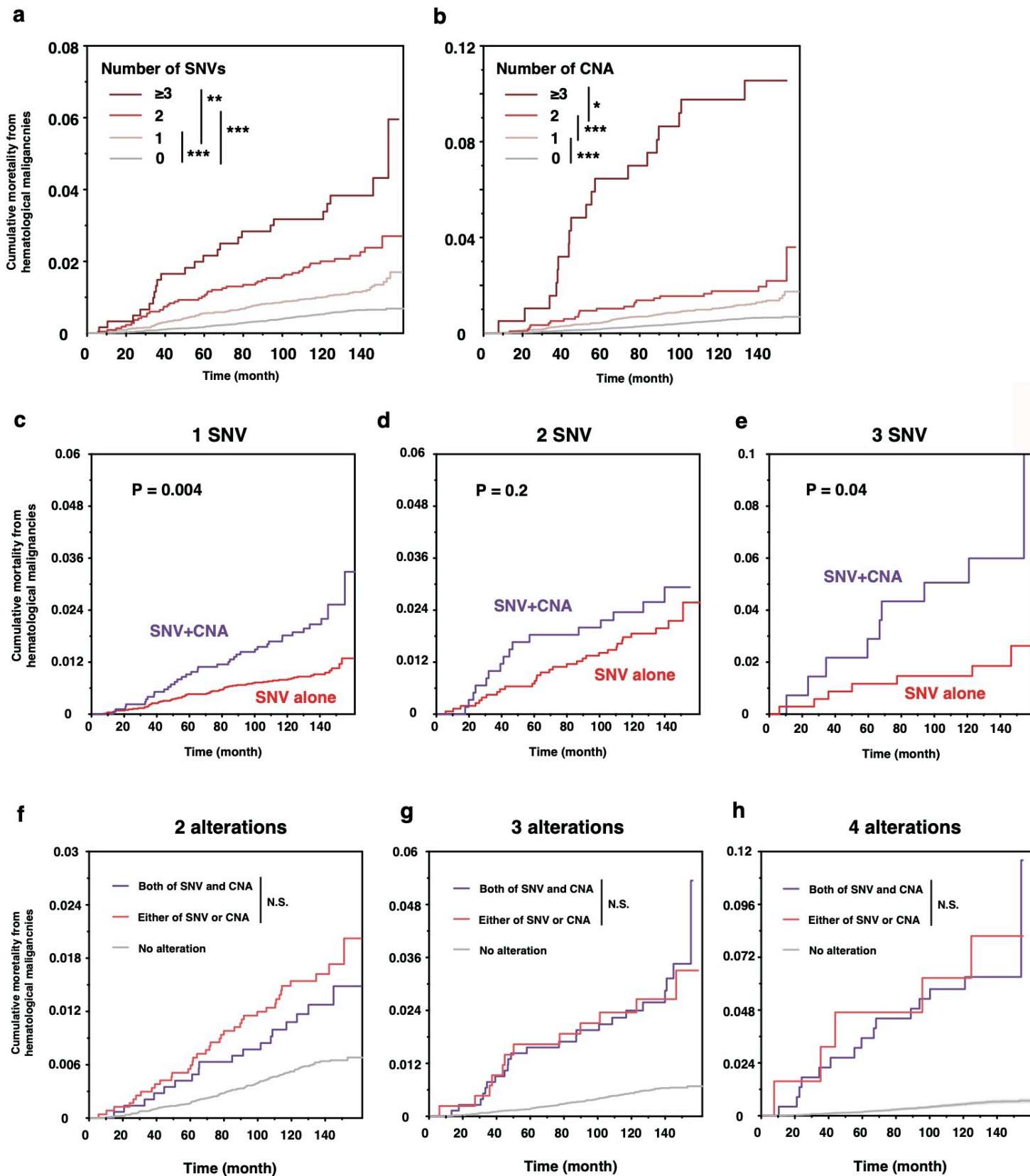
Extended Data Fig. 7



Extended Data Fig. 7 | Genetic alterations in CH and abnormalities in blood counts.

a, Landscape of SNVs/indels and CNAs in subjects without abnormalities in blood counts (left), in those with any abnormalities in blood counts (middle), and in those with no available blood counts (right). Each row represents a genetic alteration while each column represents a subject. Subjects without any alteration are omitted. Different types of mutations and CNAs are depicted by different colors. b, Enrichment of genetic alterations in subjects with abnormalities in blood counts. Sizes of rectangles indicate significance of enrichment. Colors of rectangles indicate odds ratios. The enrichment of alterations were examined by Fisher exact test. Cytopenia (All), subjects with cytopenia in at least one lineage; Cytopenia (Multi), subjects with cytopenia in ≥ 2 lineage. WBC, white blood cell; Hb, hemoglobin; Plt, platelet. c, Distribution of blood cell counts in subjects with different CH-related alterations. d, Relationships between blood cell counts and VAF of SNVs/indels.

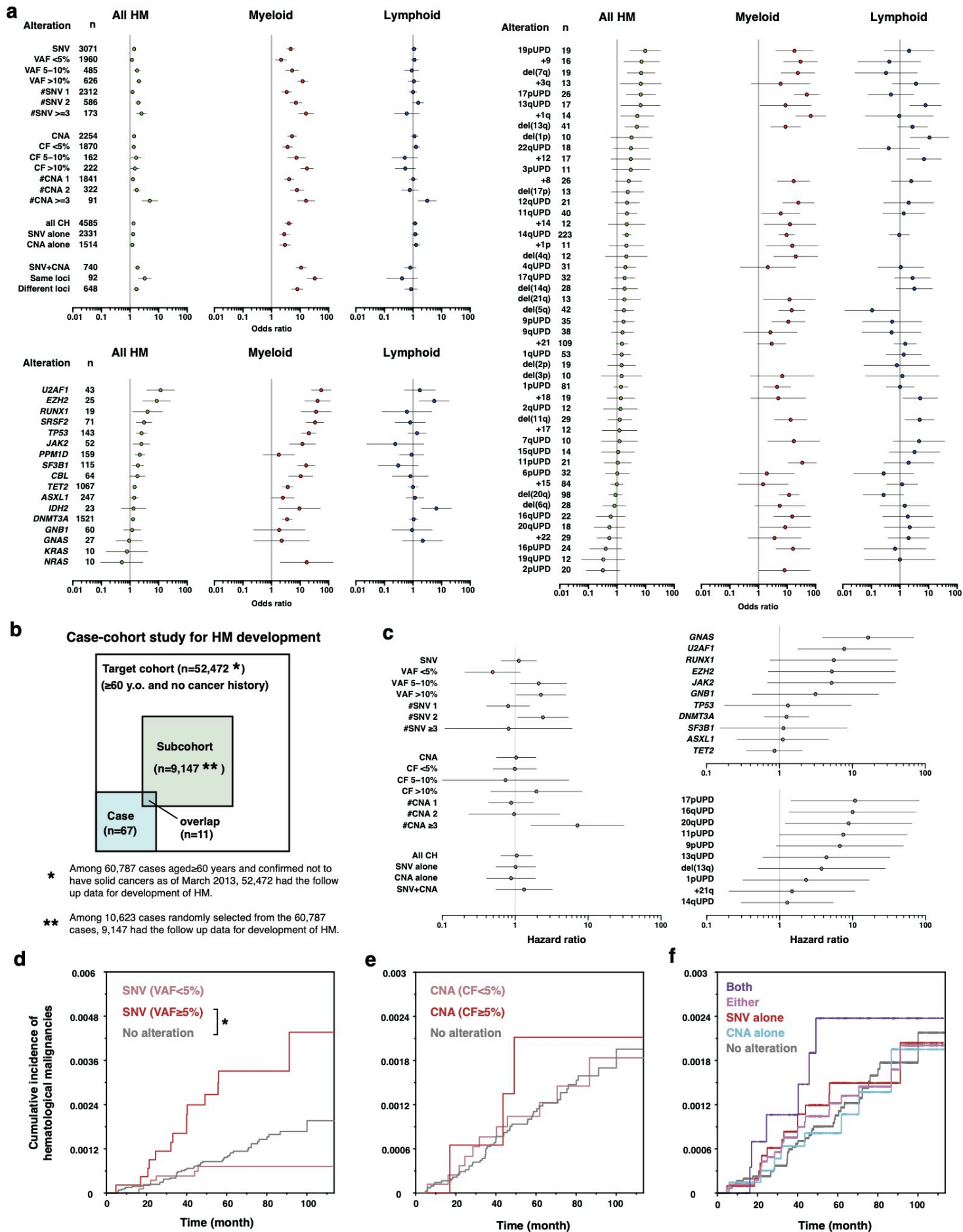
Extended Data Fig. 8



Extended Data Fig. 8 | Impact of CH on mortality from HM stratified by number of alterations..

a-b, Cumulative mortality from hematological malignancies in subjects with different number of SNVs/indels (a), or CNAs (b). c-e, Cumulative mortality from hematological malignancies in subjects with both SNVs/indels and CNAs or in those with SNVs/indels alone. Subjects with 1 (c), 2 (d), or ≥ 3 alterations are separately shown. f-h, Cumulative mortality from hematological malignancies in subjects with both SNV/indels and CNAs or in those with either of them. Subjects with 2 (f), 3 (g), or 4 alterations (h) are separately shown.

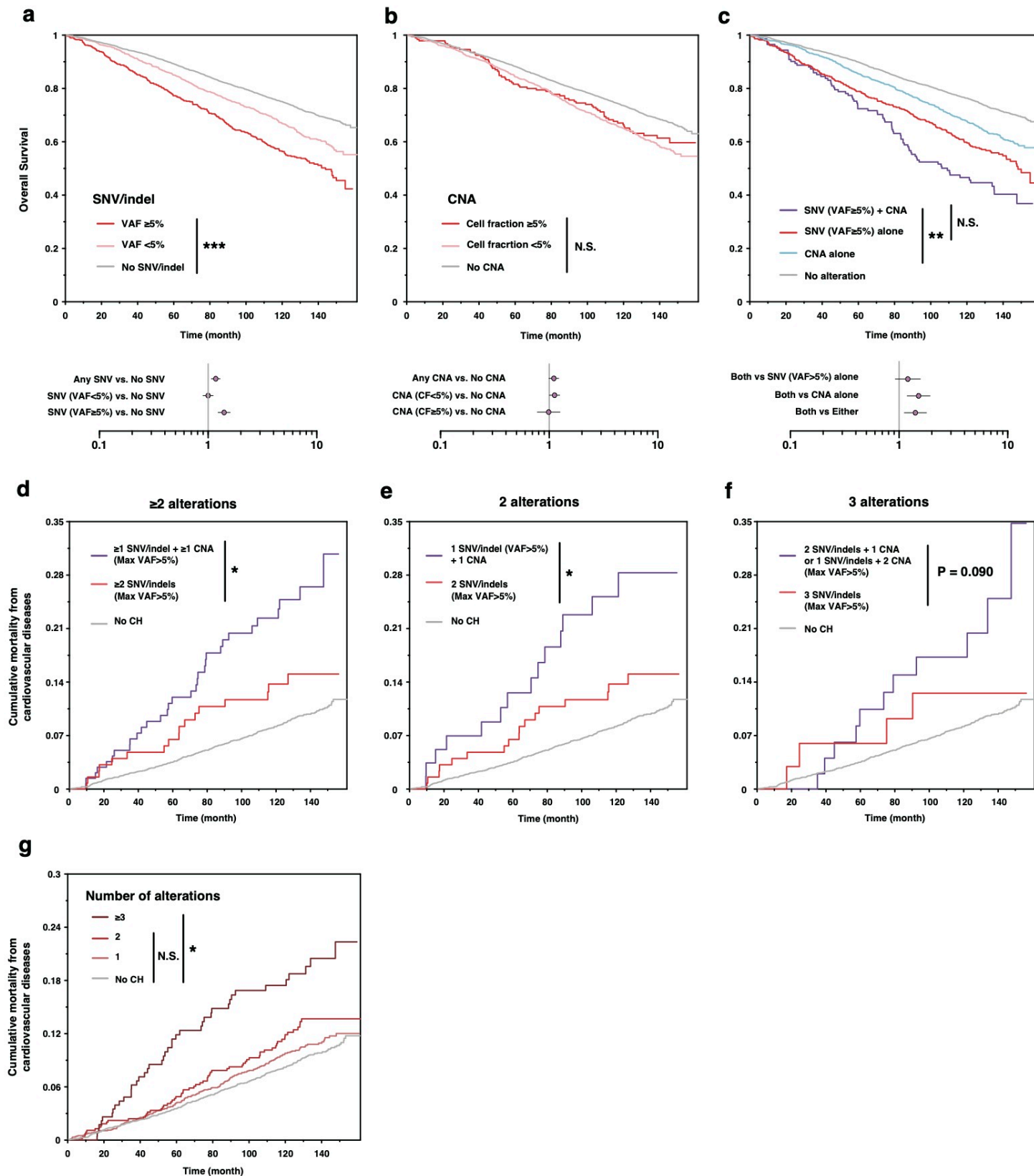
Extended Data Fig. 9



Extended Data Fig. 9 | Interplay between SNVs/indels in *TP53* and CNAs

a, Odds ratios for the events (death and/or development) of hematological malignancies in case-control study (Extended Data Fig. 1a). b, Design of case-cohort study for development of hematological malignancies. c, Hazard ratios for development of hematological malignancies. d-f, Effect of SNVs/indels (d), CNAs (e), and combined SNVs/indels and CNAs (f) on the cumulative incidence of development of hematological malignancies. n, number of cases with the indicated alterations; SNV+CNA, cocurrence of both SNVs/indels and CNAs; #SNV, number of SNVs/indels; CF, cell fraction of CNAs; #CNA, number of CNAs.

Extended Data Fig. 10



Extended Data Fig. 10 | Interplay between SNVs/indels in *TP53* and CNAs

a-c, Effect of SNVs/indels(a), CNAs(b), or combined SNVs/indels and CNAs (c) on overall survivals. d-f, Cumulative mortality from cardiovascular diseases in subjects with SNVs/indels (Max VAF $> 5\%$) alone and those with both of SNVs/indels (Max VAF $> 5\%$) and CNAs. Subject with ≥ 2 (d), 2 (e), and 3 (f) alterations are separately shown. g, Cumulative mortality from cardiovascular diseases in subjects with different number of CH-related alterations. N.S., not significant; *, $P < 0.05$; **, $P < 0.001$; ***, $P < 0.0001$.

References

1. Steensma, D.P., *et al.* Clonal hematopoiesis of indeterminate potential and its distinction from myelodysplastic syndromes. *Blood* **126**, 9-16 (2015).
2. Shlush, L.I. Age-related clonal hematopoiesis. *Blood* **131**, 496-504 (2018).
3. Busque, L., *et al.* Skewing of X-inactivation ratios in blood cells of aging women is confirmed by independent methodologies. *Blood* **113**, 3472-3474 (2009).
4. Gale, R.E., Wheadon, H. & Linch, D.C. X-chromosome inactivation patterns using HPRT and PGK polymorphisms in haematologically normal and post-chemotherapy females. *Br J Haematol* **79**, 193-197 (1991).
5. Fey, M.F., *et al.* Clonality and X-inactivation patterns in hematopoietic cell populations detected by the highly informative M27 beta DNA probe. *Blood* **83**, 931-938 (1994).
6. Champion, K.M., Gilbert, J.G., Asimakopoulos, F.A., Hinshelwood, S. & Green, A.R. Clonal haemopoiesis in normal elderly women: implications for the myeloproliferative disorders and myelodysplastic syndromes. *Br J Haematol* **97**, 920-926 (1997).
7. Busque, L., *et al.* Nonrandom X-inactivation patterns in normal females: lyonization ratios vary with age. *Blood* **88**, 59-65 (1996).
8. Jacobs, K.B., *et al.* Detectable clonal mosaicism and its relationship to aging and cancer. *Nat Genet* **44**, 651-658 (2012).
9. Laurie, C.C., *et al.* Detectable clonal mosaicism from birth to old age and its relationship to cancer. *Nat Genet* **44**, 642-650 (2012).
10. Loh, P.R., *et al.* Insights into clonal haematopoiesis from 8,342 mosaic chromosomal alterations. *Nature* **559**, 350-355 (2018).
11. Loh, P.R., Genovese, G. & McCarroll, S.A. Monogenic and polygenic inheritance become instruments for clonal selection. *Nature* (2020).
12. Genovese, G., *et al.* Clonal hematopoiesis and blood-cancer risk inferred from blood DNA sequence. *N Engl J Med* **371**, 2477-2487 (2014).
13. Jaiswal, S., *et al.* Age-related clonal hematopoiesis associated with adverse outcomes. *N Engl J Med* **371**, 2488-2498 (2014).
14. Abelson, S., *et al.* Prediction of acute myeloid leukaemia risk in healthy individuals. *Nature* **559**, 400-404 (2018).
15. Desai, P., *et al.* Somatic mutations precede acute myeloid leukemia years before diagnosis. *Nat Med* **24**, 1015-1023 (2018).
16. 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 e374 (2017).
17. Bolton, K.L., *et al.* Cancer therapy shapes the fitness landscape of clonal hematopoiesis. *Nat Genet* **52**,

1219-1226 (2020).

18. Jaiswal, S., *et al.* Clonal Hematopoiesis and Risk of Atherosclerotic Cardiovascular Disease. *N Engl J Med* **377**, 111-121 (2017).

19. Fuster, J.J., *et al.* Clonal hematopoiesis associated with TET2 deficiency accelerates atherosclerosis development in mice. *Science* **355**, 842-847 (2017).

20. Young, A.L., Challen, G.A., Birman, B.M. & Druley, T.E. Clonal haematopoiesis harbouring AML-associated mutations is ubiquitous in healthy adults. *Nat Commun* **7**, 12484 (2016).

21. Terao, C., *et al.* Chromosomal alterations among age-related haematopoietic clones in Japan. *Nature* (2020).

22. Gao, T., *et al.* Interplay between chromosomal alterations and gene mutations shapes the evolutionary trajectory of clonal hematopoiesis. *Nat Commun* **12**, 338 (2021).

23. Nagai, A., *et al.* Overview of the BioBank Japan Project: Study design and profile. *J Epidemiol* **27**, S2-S8 (2017).

24. Momozawa, Y., *et al.* Low-frequency coding variants in CETP and CFB are associated with susceptibility of exudative age-related macular degeneration in the Japanese population. *Hum Mol Genet* **25**, 5027-5034 (2016).

25. Ogawa, S. Genetics of MDS. *Blood* **133**, 1049-1059 (2019).

26. Ochi, Y., *et al.* Combined Cohesin-RUNX1 Deficiency Synergistically Perturbs Chromatin Looping and Causes Myelodysplastic Syndromes. *Cancer Discov* **10**, 836-853 (2020).

27. Papaemmanuil, E., *et al.* Genomic Classification and Prognosis in Acute Myeloid Leukemia. *N Engl J Med* **374**, 2209-2221 (2016).

28. Nik-Zainal, S., *et al.* The life history of 21 breast cancers. *Cell* **149**, 994-1007 (2012).

29. Kralovics, R., *et al.* A gain-of-function mutation of JAK2 in myeloproliferative disorders. *N Engl J Med* **352**, 1779-1790 (2005).

30. Langemeijer, S.M., *et al.* Acquired mutations in TET2 are common in myelodysplastic syndromes. *Nat Genet* **41**, 838-842 (2009).

31. Jasek, M., *et al.* TP53 mutations in myeloid malignancies are either homozygous or hemizygous due to copy number-neutral loss of heterozygosity or deletion of 17p. *Leukemia* **24**, 216-219 (2010).

32. Thoennissen, N.H., *et al.* Prevalence and prognostic impact of allelic imbalances associated with leukemic transformation of Philadelphia chromosome-negative myeloproliferative neoplasms. *Blood* **115**, 2882-2890 (2010).

33. Yoshizato, T., *et al.* Genetic abnormalities in myelodysplasia and secondary acute myeloid leukemia: impact on outcome of stem cell transplantation. *Blood* **129**, 2347-2358 (2017).

34. Watatani, Y., *et al.* Molecular heterogeneity in peripheral T-cell lymphoma, not otherwise specified revealed by comprehensive genetic profiling. *Leukemia* **33**, 2867-2883 (2019).

35. Muto, H., *et al.* Reduced TET2 function leads to T-cell lymphoma with follicular helper T-cell-like features in mice. *Blood Cancer J* **4**, e264 (2014).
36. Schneider, R.K., *et al.* Rps14 haploinsufficiency causes a block in erythroid differentiation mediated by S100A8 and S100A9. *Nat Med* **22**, 288-297 (2016).
37. Stoddart, A., *et al.* Haploinsufficiency of del(5q) genes, Egr1 and Apc, cooperate with Tp53 loss to induce acute myeloid leukemia in mice. *Blood* **123**, 1069-1078 (2014).
38. Wolkewitz, M., Palomar-Martinez, M., Olaechea-Astigarraga, P., Alvarez-Lerma, F. & Schumacher, M. A full competing risk analysis of hospital-acquired infections can easily be performed by a case-cohort approach. *J Clin Epidemiol* **74**, 187-193 (2016).
39. Hirata, M., *et al.* Overview of BioBank Japan follow-up data in 32 diseases. *J Epidemiol* **27**, S22-S28 (2017).
40. Young, A.L., Tong, R.S., Birmann, B.M. & Druley, T.E. Clonal hematopoiesis and risk of acute myeloid leukemia. *Haematologica* **104**, 2410-2417 (2019).
41. 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).
42. Mutation in TET2 in Myeloid Cancers. (2009).
43. Harismendy, O., *et al.* Detection of low prevalence somatic mutations in solid tumors with ultra-deep targeted sequencing. *Genome Biol* **12**, R124 (2011).
44. Forshew, T., *et al.* Noninvasive identification and monitoring of cancer mutations by targeted deep sequencing of plasma DNA. *Sci Transl Med* **4**, 136ra168 (2012).
45. Yoshida, K., *et al.* Frequent pathway mutations of splicing machinery in myelodysplasia. *Nature* **478**, 64-69 (2011).
46. Haferlach, T., *et al.* Landscape of genetic lesions in 944 patients with myelodysplastic syndromes. *Leukemia* **28**, 241-247 (2014).
47. Suzuki, H., *et al.* Mutational landscape and clonal architecture in grade II and III gliomas. *Nat Genet* **47**, 458-468 (2015).
48. Shiraishi, Y., *et al.* An empirical Bayesian framework for somatic mutation detection from cancer genome sequencing data. *Nucleic Acids Res* **41**, e89 (2013).
49. Niida, A., Imoto, S., Shimamura, T. & Miyano, S. Statistical model-based testing to evaluate the recurrence of genomic aberrations. *Bioinformatics* **28**, i115-120 (2012).
50. Arber, D.A., *et al.* The 2016 revision to the World Health Organization classification of myeloid neoplasms and acute leukemia. *Blood* **127**, 2391-2405 (2016).