

1 **Conserved patterns of somatic mutations in human peripheral blood**
2 **cells**

3 **L. Alexander Liggett^{1,2}, Anchal Sharma³, Subhajyoti De³, and *James**
4 **DeGregori^{1,2,4,5,6}**

5 **¹Department of Biochemistry and Molecular Genetics, ³Linda Crnic Institute for**
6 **Down Syndrome, ⁴Integrated Department of Immunology, ⁵Department of**
7 **Pediatrics, ⁶Department of Medicine, Section of Hematology, University of**
8 **Colorado School of Medicine, Aurora, CO 80045**

9 **³Rutgers Cancer Institute, New Brunswick, NJ 08901,**

10 ***Corresponding Author**

11 **Corresponding Author:**

12 **James DeGregori, Ph.D.**

13 **James.DeGregori@ucdenver.edu**

14 **Summary**

15 Mutation accumulation varies across a genome by chromosomal location, nucleotide
16 identity, surrounding sequence, and chromatin context¹⁻⁵. Nevertheless, while

17 mutagens, replication machinery, and repair processes exhibit identifiable mutation
18 signatures, at the tissue scale the aggregate manifestation of these processes has been
19 difficult to measure. The challenge in observing tissue-wide somatic mutation patterns is
20 that prior to clonal expansion, most mutations are relatively rare⁶⁻⁹. This challenge has
21 meant that somatic mutation detection in humans has largely been limited to in vitro
22 expanded stem cells¹⁰⁻¹³ or clonal expansions that occur in vivo¹⁴⁻¹⁷. Here we describe a
23 new method called FERMI (Fast Extremely Rare Mutation Identification), which
24 comprehensively captures and quantifies rare mutations at single DNA molecule
25 resolution, that exist at frequencies as rare as 10^{-4} . Using this method, we observed that
26 mutations are highly prevalent in human peripheral blood cells, with virtually every
27 position mutated across fewer than 10^5 cells. Our results revealed an unanticipated
28 degree of similarity in somatic mutation patterns across individuals, where most
29 assayed substitutions are found to occur at conserved frequencies across nearly all
30 individuals spanning a nine-decade age range. We observe substantial bias in changes
31 for many positions, including substitution to only a single base across all assayed
32 individuals. These observed mutational patterns existed both within non-conserved,
33 non-coding and non-repetitive regions of the genome and within the coding regions of
34 oncogenes implicated in hematopoietic malignancies. Finally, we identify individuals
35 who deviate from typical mutational patterns in a reproducible manner that resembles a
36 mild mismatch repair deficiency, suggesting that variance from typical somatic mutation
37 rates may be relatively common. This study provides an unprecedented characterization
38 of mutations in terminally differentiated somatic cells and demonstrates that somatic

39 mutations in such cells are significantly more frequent and deterministic than previously
40 believed.

41 Measuring somatic mutations has been technically challenging because
42 mutations occur within individual cells that do not necessarily clonally expand to
43 detectable representation. While these challenges have been somewhat overcome by
44 increasing the depth of sequencing, using clever methods of barcoding⁶ or by
45 performing paired strand collapsing¹⁸, it remains difficult to get enough sequencing
46 depth and breadth while sufficiently limiting false positive noise^{6,9,19}. To overcome these
47 sequencing limitations, we created FERMI, in which we adapted the amplicon
48 sequencing method of Illumina's TrueSeq Custom Amplicon platform to target only 32 x
49 150bp genomic regions, spanning AML-associated oncogenic mutations and the Tier III
50 regions of the human genome (non-conserved, non-protein coding and non-repetitive).
51 We further improved upon Illumina's capture efficiency to achieve approximately 1.2
52 million unique captures from 500ng-1µg of genomic DNA (gDNA) (see Methods). We
53 designed the targeting probes used in gDNA capture with a 16bp index of sequence
54 unique to each individual and a 12bp unique molecular identifier (UMI) of random DNA
55 unique to each capture (Fig. 1a). Sequencing reads were sorted by sample index and
56 UMI, producing bins of single cell sequencing which were collapsed to produce
57 relatively error-free consensus reads. Captures were only considered if supported by at
58 least 5 reads, and variants were only included if identified in both paired-end
59 sequences, and detected in at least 55% percent of supporting reads (Fig. 1a and
60 Methods; see also Extended Data Figure 1).

61 While all probed regions were successfully captured and amplified, capture
62 efficiency varied by 2-3-fold dependent on probe identity (Fig 1b). To understand assay

63 sensitivity, log-series dilutions of human heterozygous single nucleotide polymorphisms
64 (SNPs) were prepared and assayed by FERMI. Using these dilutions, we observed
65 robust quantification of diluted SNPs as rare as 10^{-4} (Fig. 1c). Even more accurate
66 quantifications of SNP frequency can be made when using strand information to follow
67 dilutions of multiple SNPs located on the same allele (Fig. 1d). For more description of
68 the methods used to maximize the accuracy of FERMI, see *Elimination of false positive*
69 *signal* in Methods and Extended Data Figure 1.

70 Using FERMI, we captured and sequenced gDNA from the peripheral blood of 22
71 apparently healthy donors ranging in age from 0 (cord blood) to 89 years of age
72 (Extended Data Table 1). Surprisingly, within each of the probed regions, nearly every
73 position is mutated in at least one individual, including all probed oncogenic mutations,
74 independent of segment location or individual age, indicating a mutation burden of
75 greater than 50 per megabase (See *Estimation of mutation burden* in Methods). While
76 FERMI could correctly identify individual-specific unique germline SNPs (Extended Data
77 Figure 2a), rare somatic variants are found at remarkably similar allele frequencies
78 across all sampled ages. The rare allele frequencies are similar enough between most
79 individuals that comparisons of the variant allele frequencies for each unique
80 substitution falls along a $y=x$ line (Fig. 2a). FERMI of biopsies taken 1 month apart from
81 the same individuals revealed the same germline SNPs (Extended Data Figure 2b), but
82 detected rare variants are not significantly more similar to each other than to other
83 individuals (Extended Data Figure 2c). Variant allele frequencies (VAFs) were averaged
84 across 22 sampled blood donors and used as a comparison to individuals, which

85 appear age-independent and still adhere to a $y=x$ line (R^2 Range = 0.426-0.631, Mean =
86 0.558) (Fig. 2b), and are similar across experiments (Extended Data Figure 3a-d shows
87 data from an additional 11 individuals). Variants with frequencies above 0.001 were
88 found in nearly all samples, while more rare variants were missed with a probability
89 inversely proportional to their allele frequencies. Furthermore, most variants likely
90 represent multiple independent events rather than clonal expansions, as they are found
91 at similar frequencies on both alleles (Extended Data Figure 3e). It thus appears that
92 instead of being semi-random, the aggregate effect of all DNA damage and
93 maintenance generates somatic mutations at predictable rates throughout the genome
94 independent of age. We suspect that such mutations primarily arise during the
95 generation of terminally differentiated blood cell types in a sequence context-dependent
96 manner, with minimal impact of selection, such that it reflects the basal DNA damage
97 and repair errors in hematopoietic cells.

98 We observed that the overall probability of a substitution occurring is biased by
99 nucleotide identity, with C>T substitutions being the most common and T>G
100 substitutions being the least common (Fig. 2c). These biases were largely expected, as
101 similar patterns have been observed both in other healthy tissues and in
102 cancers^{10,14,17,20,21}. There were notable differences, especially for C>N changes which we
103 observe as underrepresented within a CpG context (Fig. 2d). Regardless of functional
104 or oncogenic potential, each site tends to undergo the same substitutions across
105 individuals (Fig. 2e). These conserved substitution rates appear to be deterministic, and
106 cannot be explained by undersampling (Extended Data Figure 4) or known base change

107 biases (Extended Data Figure 5). It therefore appears that the combined sources of
108 external and internal DNA mutation result in systematic substitutions at frequencies that
109 are often predictable by location and sequence context. Suggestive of differences
110 during cancer evolution and normal somatic mutation, the integrated exome sequencing
111 pan cancer somatic mutation data from the TCGA exhibits different substitution patterns
112 from those that we find in healthy donor blood (Extended Data Figure 6a). Using the
113 trinucleotide contexts of the substitutions, 7 out of 30 previously identified mutations
114 signatures were identified, and these signatures did not differ significantly across
115 sampled genomic segments (Extended Data Figure 6b-c).

116 While we observe variants at conserved frequencies across many individuals,
117 previous studies have described clonal expansions bearing AML-associated oncogenic
118 changes that are largely restricted to old age^{14–16,22}. While we observe each queried
119 oncogenic change in every biopsied individual independent of age, we do not observe
120 significant age-related changes in the allele frequencies of either oncogenic or
121 non-oncogenic mutations within proto-oncogenes (Fig. 2f and Extended Data Figure 7).
122 This inability to observe any clonal expansions with age is most likely due to the fact
123 that the average age of the individuals within our cohort is 49 years, with only 5 donors
124 older than 70 years.

125 To explore the ability of FERMI to distinguish perturbations of somatic mutation
126 patterns, gDNA from mismatch repair deficient HCT116 cells (MMR^{MT}; hemizygous for
127 MLH1) was compared to MMR proficient parental cell line gDNA. Substantiating our
128 method, there was a substantial increase in VAFs within the MMR^{MT} gDNA when

129 compared to parental gDNA (Fig. 3a-b). Unexpectedly, while the VAFs for most
130 peripheral blood samples closely resemble those in other individuals, samples from two
131 individuals (2 and 19), contained a subset of variants that deviated from the population
132 averages with approximately a twofold increase in prevalence (Fig. 3c, 3d, and
133 Extended Data Figure 9). While the magnitude of deviation from mean VAFs was
134 different, the identities of the deviating variants were the same, such that a comparison
135 of VAFs between these two individuals correlate more closely to a $y=x$ line than to the
136 overall population average (Fig. 3e). This consistent deviation in VAFs for these two
137 individuals from the averaged population suggests that the mechanisms governing
138 mutation levels can be systematically perturbed. Surprisingly, the VAF changes in these
139 two individuals resemble those altered in the MMR^{MT} HCT116 cells, though the
140 magnitude of these changes are greater in the latter (Fig. 3f). Finally, the deviating
141 variants found within individuals 2 and 19 are not enriched for either oncogenic variants
142 or for other variants within coding regions (Fig. 3g), indicating that deviations from the
143 typical variant pattern are not likely the result of selection.

144 As expected from previous studies²³, the HCT116 MMR^{MT} gDNA showed an
145 increased prevalence of T>C and T>A substitutions when compared to parental gDNA
146 (Extended Data Figure 8). The samples from individuals #2 and #19 also exhibited
147 these increased rates of T>C and T>A substitutions, with less extensive increases at C
148 positions, compared with the average of the 22 individuals (Fig. 3h-j and Extended Data
149 Figure 9), mirroring the changes observed in MMR^{MT} HCT116 cells. Thus, these two
150 individuals appear to present with a mild MMR-like substitution pattern. In support of the

151 results, individual #2 shows the same increased rates of substitutions across multiple
152 experiments, with strong reproducibility in mutation patterns (Extended Data Figure
153 9h-j). Of note, the systematic variance from the typical mutational pattern for these two
154 individuals and the MMR^{MT} HCT116 cells serves as validation of the specificity of FERMI
155 to accurately detect variants. More importantly, this finding of two individuals with
156 deviating mutational patterns out of a sample size of only 22 individuals may indicate
157 that individuals with significant deviation from typical mutational profiles may be
158 relatively common in the human population.

159 **Conclusion**

160 These studies reveal an unprecedented degree of similarity in somatic mutational
161 patterns across most individuals, that almost all genomic positions are mutated within
162 less than a hundred-thousand leukocytes, and how mutational spectra can be
163 systematically disrupted in some individuals. Strikingly, we observed extremely
164 reproducible biases at *each particular* nucleotide position in terms of the frequency of
165 changes and the base to which it is changed. These strong position-dependent
166 substitution biases will restrict phenotypic diversity upon which somatic evolution can
167 act. It appears that mutation incidence, both non-oncogenic and oncogenic, are
168 relatively well tolerated, highlighting the importance of evolved tumor suppressive and
169 tissue maintenance mechanisms.

170 **Acknowledgments**

171 We would like to thank Ruth Hershberg of Technion University and Jay Hesselberth and
172 Robert Sclafani of the University of Colorado School of Medicine for useful suggestions
173 and for review of the manuscript. These studies were supported by grants from the
174 National Cancer Institute (R01CA180175 to J.D.), NIH/NCATS Colorado CTSI Grant
175 Number UL1TR001082CU (seed grant to J.D.), F31CA196231 (to L.A.L.), and the Linda
176 Crnic Institute for Down Syndrome (to J.D. and L.A.L.). The research utilized services of
177 the Cancer Center Genomics Shared Resource, which is supported in part by NIH grant
178 P30-CA46934.

179 **Contributions**

180 L.A.L. and J.D. developed the concept of this project, planned the experiments,
181 analyzed results, and wrote the manuscript. L.A.L. processed and prepared samples
182 from blood biopsy to sequencing, and wrote the bioinformatics software used for
183 analysis. A.S. and S.D. ran the analyses found in Extended Data Figure 6 and assisted
184 in checking the validity of results throughout the manuscript.

185 **Figure 1 | Amplicon sequencing accurately detects mutation allele frequencies as**
186 **rare as 1/10,000. a**, Graphical depiction of gDNA capture and analysis method. **b**,
187 Capture efficiencies vary in a probe dependent manner. **c**, Accurate detection of a
188 single heterozygous SNP in gDNA from one individual diluted into gDNA from another
189 (without this germline SNP) to frequencies as low as 1/10,000. **d**, Accurate detection of

three linked SNPs found within the same allele diluted as in c. For c and d, error shown is standard deviation.

Figure 2 | Mutations exist at conserved frequencies independently of age. a,

Comparison of VAFs of identified variants within a 34 year old (x-axis) and 62 year old

(y-axis); $R^2 = 0.408211$, $p=0.000$. R^2 values unless otherwise noted are calculated for all

points falling below VAFs of 0.003 which largely includes all variants but germline. **b,**

VAFs from a 34 year old (x-axis) compared to mean VAFs from individuals ranging in

190 ages from newborn to 89 years of age ($n=22$); $R^2 = 0.590412$, $p=0.000$. **c,** Relative

191 contribution rates of each base substitution to all substitutions identified. **d,** Relative

192 contribution rates of each base substitution segregated by surrounding 5' and 3'

193 nucleotide context. **e,** All identified base substitutions within a probed region are plotted

194 by their position and allele frequencies for individuals 7 and 15 (representative of all

195 other individuals, with greater deviation observed for individuals 2 and 19 as described

196 below), revealing highly reproducible patterns. **f,** Oncogenic VAFs plotted as a function

197 of donor age does not reveal evidence of clonal expansions.

198 **Figure 3 | Individuals Can Systematically Deviate from the Population Average. a,**

199 Comparing VAFs in HCT116 MMR+ vs MMR^{MT} cells reveals an increase in frequencies

200 for many of the observed variants in MMR^{MT} cells ($R^2 = 0.211479$). **b,** MMR^{MT} vs mean

201 VAFs from blood of the 22 individuals shows a similar pattern of increased VAFs as the

202 comparison with parental HCT116 cells ($R^2 = 0.120895$). **c,** blood from a 73 year old

203 person (individual #19) compared to the mean VAFs reveals a deviating population of
204 variants that exist at an increased frequency compared with average VAFs ($R^2 =$
205 0.387125). **d**, A cord blood sample (individual #2) also shows a subset of variants with
206 higher frequencies than in the average ($R^2 = 0.278250$). **e**, VAFs from individual #2 vs
207 individual #19 reveals that the deviating variants are at the same positions, causing the
208 comparison to fall close to the $y=x$ line ($R^2 = 0.613542$). **f**, Plotting the mean for VAFs
209 from individuals #2 and #19 versus VAFs from MMR^{MT} HCT116 cells reveals that the
210 variants within the blood are the same as those found within the MMR^{MT} cell line. While
211 variant frequencies are higher in the MMR^{MT} cell line, the proportional change for
212 different deviating variants are similar ($R^2 = 0.587474$). **g**, Variants detected in
213 individuals #2 and #19 are not enriched for oncogenic changes, indicated in blue. **h**,
214 Plot of only C>N/G>N variants shows relative similarity between individual #2 and the
215 average for all other individuals ($R^2 = 0.350623$). **i**, Plot of only T>N/A>N variants
216 reveals that the majority of deviating variants for individual #2 are substitutions affecting
217 T or A (R-Squared = 0.040712).

218 **Methods**

219 **Amplicon Design**

220 Amplicon probes for targeted annealing regions were created using the Illumina
221 Custom Amplicon DesignStudio (<https://designstudio.illumina.com/>). UMIs were then
222 added to the designed probe regions and generated by IDT using machine mixing for

223 the randomized DNA. Probes were PAGE purified by IDT. All probes are listed below
 224 along with binding locations and expected lengths of captured sequence.

Gene	Probe Up	Probe Down	Probe Start	Probe End	Length
JAK2	AGTTTACACTGACA CCTAGCTGTGATC	CCATAATTTAAAACC AAATGCTTGTGAGA 236 A	chr9:5073733	chr9:5073887	155
TP53-1	TCATCTTGGGCCTG TGTTATCTCCTA	ATCCTCACCATCAT CACACTGGAAGAC	chr17:7577504	chr17:7577635	132
TP53-2	CCCTCAACAAGATG TTTTGCCAACTG	ATGAGCGCTGCTCA GATAGCGATGGT	chr17:7578369	chr17:7578544	176
TP53-3	GGACAGGTAGGAC CTGATTTCCCTTACT	TGTCCTGGGAGAGA CCGGCCGACAGA	chr17:7577084	chr17:7577214	131
NRAS-1	CAATAGCATTGCAT TCCCTGTGGTTTT	GTACAGTGCCATGA GAGACCAATACAT	chr1:115256496	chr1:115256680	185
NRAS-2	GAAGGTCACACTAG GGTTTTCATTTCC	AAAAGCGCACTGAC AATCCAGCTA	chr1:115258713	chr1:115258897	185
HRAS	TCCTTGGCAGGTGG GGCAGGAGACCC	GCAAGAGTGCCTG ACCATCCA	chr11:534258	chr1:534385	128
KRAS-1	AGGTAAGTGGTGGAG TATTTGATAGTGT	CAAGAGTGCCTTGA CGATACAGCTAATT	chr12:25398247	chr12:25398415	169
KRAS-2	GACTGTGTTTCTCC CTTCTCAGGATTC	TACAGTGCAATGAG GGACCAGTACATG	chr12:25380242	chr12:25380368	127
TET2-1	CCATGTTTTGGCTC ATTATGCTCTTA	ACGGCCACTCCCC AATGTCAG	chr4:106197237	chr4:106197405	169
TET2-2	CTTTTGAAGAGTG CCACTGGTGTCT	GGTGATGGTATCAG GAATGGACTTAGTC	chr4:106155137	chr4:106155275	139
DNMT3A	TGTGTGGTTAGACG GCTTCCGGGCA	AGGCAGAGACTGCT GGGCCGGTCA	chr2:25457211	chr2:25457364	154
IDH1	CAAATGTGAAAATC ACCAAATGGCACC	TGGGGATCAAGTAA GTCATGTTGGCA	chr2:209113077	chr2:209113239	163
IDH2	GAAGAAGATGTGGA AAAGTCCAATGG	CATGGCGACCAGGT AGGCCAGG	chr15:90631809	chr15:90631969	161
GATA1	CTTCCAGCCATTTT TGAGATATCCTCA	CAGCTGCAGCGGT GGCTGTGCT	chrX:48649667	chrX:48649849	183
SF3B1	GTGAACATATTCTG CAGTTTGGCTGAA	ACCATCAGTGCTTT GGCCATTGC	chr2:198266803	chr2:198266967	165
TIIIA	CATCTATTCTGTGCT AGGCATTGTGTG	CAGACCTAGCATCT GTGCCAGAC	chr1:115227814	chr1:115227978	165
TIIIB	CAGTCTGGGTTTTG GAGCAATGATATC	GCAGTGAGCTCAGC CTTGATTTT	chr2:223190674	chr2:223190820	147
TIIIC	CCTGGTGCTTAGTC CTGTTCTGAAATT	AGTCTTCTATAATGC CACAACTGTAT	chr2:229041101	chr2:229041289	189
TIIID	GAACAGAACTTTG GTAGTTGACCATG	AGACAGGGAAGTGG CATGAAGAGTTT	chr4:110541172	chr4:110541302	131
TIIIE	GCCTAGAACAGGCA CCATACATTCAAT	AGATGGTGTGCTG TGCCGGATAGGAG	chr4:112997214	chr4:112997386	173
TIIIF	TGGCACTATGTGGA	GGATGTTGGTGCTA	chr4:121167756	chr4:121167884	129

	GATGTTAGTACAG	TCAGTAGCCATA			
TIIG	CTCTAGGCTTAGTG GTCAAGGAATGAA	AGAAGCAGGACTGT GCTTCCAAACAA	chr4:123547743	chr4:123547901	159
TIHH	CTTGGTGGTAGCCT AGGCAGTAATTAA	CACGTGGTTGGGAA GAGAAAGTG	chr4:124428637	chr4:124428767	131
TIHJ	TTCTATAGCACTGG TGACCAGGACACT	CTGGCCACAGTGCC TGGTTTCC	chr11:2126256	chr11:2126420	165
TIHK	AGACAGGAGGAAG GAGCAATTCAGAAG	CATGGAGATCTCGT CCCCTCAGA	chr11:2389983	chr11:2390171	189
TIHL	TAGGCCAGAAAACA CACAGTGTCCGGG	AACTCCGGTAAGTG GCGGGTGGGGGT	chr11:2593889	chr11:2594074	186
TIHM	ATCTGGGAACAGAC CTTCTCAGGCAT	GTTCTAAGTTACTCT GTGTACTTGACT	chr11:11486596	chr11:11486728	133
TIHN	AGCCTAGTTACCAT AGACGGATTCAAC	GAATATCTTCTAACT GGACTTAGAAAACC	chr15:92527052	chr15:92527176	125
TIHO	CCAACATGTTCTAA ATTCTGGCCACAG	TGGGTCTCAGCCAT CCCATTACTG	chr16:73379656	chr16:73379832	177
TIHP	CTAACATCTCACTTC TACCCTACGCTA	TAAGTGCCCACTAC CCCATCCTTAAT	chr16:82455026	chr16:82455164	139
TIHQ	TCATGACCCAGGCC TCCCAGAACTGAG	ATCTGTGAAGCCGG AGTGAAAACAAC	chr16:85949137	chr16:85949299	163

488 Genomic DNA Isolation

489 Human blood samples were purchased from the Bonfils Blood Center
 490 Headquarters of Denver Colorado. Our use of these samples was determined to be “Not
 491 Human Subjects” by our Institutional Review Board. Biopsies were collected as
 492 unfractionated whole blood from apparently healthy donors, though samples were not
 493 tested for infection. Samples were approximately 10 mL in volume, and collected in BD
 494 Vacutainer spray-coated EDTA tubes. Following collection, samples were stored at 4°C
 495 until processing, which occurred within 5 hours of donation. To remove plasma from the
 496 blood, samples were put in 50 mL conical tubes (Corning #430828) and centrifuged for
 497 10 minutes at 515 rcf. Following centrifugation, plasma was aspirated and 200 mL of
 498 4°C hemolytic buffer (8.3g NH₄Cl, 1.0g NaHCO₃, 0.04 Na₂ in 1L ddH₂O) was added to
 499 the samples and incubated at 4°C for 10 minutes. Hemolyzed cells were centrifuged at
 500 515 rcf for 10 minutes, supernatant was aspirated, and pellet was washed with 200 mL

501 of 4°C PBS. Washed cells were centrifuged for at 515rcf for 10 minutes, from which
502 gDNA was extracted using a DNeasy Blood & Tissue Kit (Qiagen REF 69504).

503 **Amplicon Capture**

504 For amplicon capture from gDNA, we modified the Illumina protocol called
505 “Preparing Libraries for Sequencing on the MiSeq” (Illumina Part #15039740 Revision
506 D). DNA was quantified with a NanoDrop 2000c (ThermoFisher Catalog #ND-2000C).
507 500ng of input DNA in 15µl was used for each reaction instead of the recommended
508 quantities. In place of 5µl of Illumina ‘CAT’ amplicons, 5µl of 4500ng/µl of our amplicons
509 were used. During the hybridization reaction, after gDNA and amplicon reaction mixture
510 was prepared, sealed, and centrifuged as instructed, gDNA was melted for 10 minutes
511 at 95°C in a heat block (SciGene Hybex Microsample Incubator Catalog #1057-30-O).
512 Heat block temperature was then set to 60°C, allowed to passively cool from 95°C and
513 incubated for 24hr. Following incubation, the heat block was set to 40°C and allowed to
514 passively cool for 1hr. The extension-ligation reaction was prepared using 90 µl of ELM4
515 master mix per sample and incubated at 37°C for 24hr. PCR amplification was
516 performed at recommended temperatures and times for 29 cycles. Successful
517 amplification was confirmed immediately following PCR amplification using a
518 Bioanalyzer (Agilent Genomics 2200 TapeStation Catalog #G2964-90002, High
519 Sensitivity D1000 ScreenTape Catalog #5067-5584, High Sensitivity D1000 Reagents
520 Catalog #5067-5585). PCR cleanup was then performed as described in Illumina’s
521 protocol using 45 µl of AMPure XP beads. Libraries were then normalized for

522 sequencing using the Illumina KapaBiosystems qPCR kit (KapaBiosystems Reference #
523 07960336001).

524 **Sequencing**

525 Prepared libraries were pooled at a concentration of 5 nM and mixed with PhiX
526 sequencing control at 5%. Libraries were sequenced on the Illumina HiSeq 4000 at a
527 density of 12 samples per lane.

528 **Bioinformatics**

529 The analysis pipeline used to process sequencing results can be found under
530 FERMI here: <http://software.laliggett.com/>. For a detailed understanding of each
531 function provided by the analysis pipeline, refer directly to the software. The overall goal
532 of the software built for this project is to analyze amplicon captured DNA that is tagged
533 with equal length UMIs on the 5' and 3' ends of captures, and has been paired-end
534 sequenced using dual indexes. Input fastq files are either automatically or manually
535 combined with their paired-end sequencing partners into a single fastq file. Paired reads
536 are combined by eliminating any base that does not match between Read1 and Read2,
537 and concatenating this consensus read with the 5' and 3' UMIs. A barcode is then
538 created for each consensus read from the 5' and 3' UMIs and the first five bases at the
539 5' end of the consensus. All consensus sequences are then binned together by their
540 unique barcodes. The threshold for barcode mismatch can be specified when running
541 the software, and for all data shown in this manuscript one mismatched base was

542 allowed for a sequence to still count as the same barcode. Bins are then collapsed into
543 a single consensus read by first removing the 5' and 3' UMIs. Following UMI removal,
544 consensus sequences are derived by incorporating the most commonly observed
545 nucleotide at each position, so long as the same nucleotide is observed in at least a
546 specified percent of supporting reads (55% of reads was used for results in this
547 manuscript) and there are least some minimum number of reads supporting a capture
548 (5 supporting reads was used for results in this manuscript). Any nucleotide that does
549 not meet the minimum threshold for read support is not added to the consensus read,
550 and alignment is attempted with an unknown base at that position. From this set of
551 consensus reads, experimental quality measurements are made, such as total captures,
552 total sequencing reads, average capture coverage, and estimated error rates.

553 Derived consensus reads are then aligned to the specified reference genome
554 using Burrows-Wheeler²⁴, and indexed using SAMtools²⁵. For this manuscript
555 consensus reads were aligned to the human reference genome hg19^{26,27} (though the
556 software should be compatible with other reference genomes). Sequencing alignments
557 are then used to call variants using the Bayesian haplotype-based variant detector,
558 FreeBayes²⁸. Identified variants are then decomposed and block decomposed using the
559 variant toolset vt²⁹. Variants are then filtered to eliminate any that have been identified
560 outside of probed genomic regions. If necessary variants can also be eliminated if below
561 certain coverage or observation thresholds such that variants must be independently
562 observed multiple times in different captures to be included. For this manuscript, we

563 included all variants that passed previous filters and did not eliminate those that were
564 observed only within a single capture, unless otherwise indicated.

565 **Elimination of false positive signal**

566 A number of steps have been included within sample preparation and
567 bioinformatics analysis specifically to distinguish between true positive signal and false
568 positive signal. Using the dilution series shown in Figs. 1c-d we can show sufficient
569 sensitivity to identify signal diluted to levels as rare as 10^{-4} . While these dilutions show
570 significantly improved sensitivity over many current sequencing methods, they do not
571 address our background error rate. Unfortunately, because both endogenous and
572 exogenous DNA synthesis is error prone, it is challenging to find negative controls that
573 can be used to estimate background error rates with a method of mutation detection as
574 putatively sensitive as FERMI. Nevertheless, we have a number of steps that should
575 eliminate most sources of false signal. The two largest sources of erroneous mutation
576 when sequencing DNA will typically be from PCR amplification mutations (caused both
577 by polymerase errors and exogenous insults like oxidative damage), and sequencing
578 errors.

579 The steps are the following:

- 580 ● *Elimination of first round PCR amplification errors*
- 581 ● *Elimination of subsequent PCR amplification errors*
- 582 ● *Elimination of sequencing errors*

583 *Elimination of first round PCR amplification errors*

584 The first round of PCR amplification performed during library preparation causes
585 mutations that are challenging to distinguish from those that occurred endogenously.
586 Since there is little difference between those mutations that occur during the first round
587 of PCR amplification and those that occurred endogenously, we rely on probability to
588 eliminate these errors. Since we are performing single-cell sequencing, we can require
589 that a mutation be observed in multiple cells before it is called as a true positive signal.
590 We expect about 400 first round PCR amplification errors, and the probability that the
591 identical mutation will occur in multiple cells becomes exponentially unlikely (Extended
592 Data Figure 1). By requiring a mutation be observed in just three cells before it is called
593 as real signal, only about 1-2 first round PCR amplification errors should make it into the
594 final data. In contrast, when we process our data requiring up to 5 independent
595 observations of a mutation, the overall mutation spectrum does not change, apart from
596 a loss of the most rarely observed variants. This observation led us to include all
597 variants that were observed even once.

598 *Elimination of subsequent PCR amplification errors*

599 Elimination of PCR amplification errors after the first round of PCR is done using
600 UMI collapsing (Fig. 1a). Each time a strand is amplified, the UMI will keep track of its
601 identity. Any mutations that occur after the first round of PCR will be found on average in
602 25% of the reads (or fewer for subsequent rounds). This allows us to collapse each

603 unique capture and eliminate any rarely observed variants associated with a given UMI.
604 Utilizing the UMI in this way allows us to essentially eliminate any PCR amplification
605 errors that occurred after the first round of PCR.

606 *Elimination of sequencing errors*

607 Sequencing errors are eliminated in two ways. This first method is by using
608 paired-end sequencing to read the same fragment of DNA twice (Fig. 1a). The
609 sequence of these reads (Read1 and Read2) should match in lieu of sequencing errors.
610 For an error to escape elimination it would need to occur at the same position (changing
611 to the same new base) within both Read1 and Read2. Therefore, when the base call
612 differs at a position on Reads 1 and 2, these changes are eliminated from the final
613 sequence. This collapsing should eliminate most sequencing errors, although
614 sequencing errors of the same identity occurring at the same position will escape.
615 These errors should be removed when collapsing into single cell bins (Fig. 1a). As with
616 the logic when eliminating subsequent PCR amplification errors, most sequences
617 associated with each UMI pair should be identical. Therefore, sequencing errors
618 passing through Read1 and Read2 will be very unlikely to match other sequenced
619 strands from the same capture event, and are eliminated during consensus sequence
620 derivation.

621 **Mutation signature analysis**

622 Twenty somatic mutation signatures were previously identified²⁰ by analyzing
623 trinucleotide mutation context of cancer genomes using non-negative matrix
624 factorization (NMF) and principal component analysis (PCA). Here, we used
625 deconstructSig³⁰ to identify the relative presence of those mutation signatures within the
626 somatic mutations detected blood using somaticSignatures³¹. Codon triplet biases were
627 analyzed using the MutationalPatterns R package³².

628 Estimation of mutation burden

629 It is difficult to understand the somatic lineage development that gave rise to the
630 number of cells that are assayed from each blood biopsy. Therefore, estimating a
631 somatic mutation rate is challenging. Nevertheless, we can derive estimates of somatic
632 mutation burden.

633 An upper bound for the somatic mutation burden observed by FERMI analysis
634 can be estimated by using the number of captures and total observed variants, and
635 assume that all of these are de-novo mutations. In our data from Cohort 1, we observe
636 on average 1,232,458 unique captures per analyzed blood sample. These captures are
637 relatively uniformly spread across each of our 32 different probes, which span a total of
638 4838bp. From this, the total probed DNA, D_T , can be estimated as:

$$639 D_T = \frac{1232458 \text{ captures} * 4838 \text{ bp}}{32 \text{ probes}}$$
$$640 D_T = 186332243.9 \text{ bp}$$

641 The total number of observed variants within each blood sample is on average
642 168,940, from which the aggregate mutation burden, M , can be estimated as:

$$643 M = \frac{168940 \text{ mutations}}{186332243.9 \text{ bp}}$$
$$644 M = 9 * 10^{-4} \text{ mut/bp}$$

645
$$M = 900 \text{ mut/Mb}$$

646 A lower estimate can be made by assuming that mutations are not all unique
647 occurrences but might be the result of clonal expansions creating many copies of each
648 mutation. This mutation burden, M , can be roughly estimated by the approximately
649 40,000 captures per each of the 32 probes that captured roughly 6000 variants across a
650 conservative 100bp sized capture for each probe (probe region is realistically smaller
651 than 150bp because of collapsing conditions). Given that all variants for which allelic
652 information could be discerned were present on both alleles, we can realistically
653 conclude each of the ~3000 base positions queried was mutated at least twice (hence
654 the estimate of 6000 variants).

655
$$M = \frac{6000 \text{ variants/sample}}{40000 \text{ captures} * 32 \text{ probes} * 100 \text{ bp/probe}}$$

656
$$M = 5 * 10^{-5} \text{ mut/bp}$$

657
$$M = 50 \text{ mut/Mb}$$

- 658
1. Benzer, S. ON THE TOPOGRAPHY OF THE GENETIC FINE STRUCTURE. *Proc. Natl. Acad. Sci. U. S. A.* **47**, 403–415 (1961).
659
 2. Gaffney, D. J. & Keightley, P. D. The scale of mutational variation in the murid genome. *Genome Res.* **15**, 1086–1094 (2005).
660
 3. Lercher, M. J., Williams, E. J. B. & Hurst, L. D. Local similarity in evolutionary rates extends over whole chromosomes in human-rodent and mouse-rat comparisons: implications for understanding the mechanistic basis of the male mutation bias.
661
662
663
664
665
Mol. Biol. Evol. **18**, 2032–2039 (2001).
 4. Nachman, M. W. & Crowell, S. L. Estimate of the mutation rate per nucleotide in humans. *Genetics* **156**, 297–304 (2000).
666
667
 5. Hwang, D. G. & Green, P. Bayesian Markov chain Monte Carlo sequence analysis reveals varying neutral substitution patterns in mammalian evolution. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 13994–14001 (2004).
668
669
670
 6. Hiatt, J. B., Pritchard, C. C., Salipante, S. J., O’Roak, B. J. & Shendure, J. Single molecule molecular inversion probes for targeted, high-accuracy detection of
671
672
673
low-frequency variation. *Genome Res.* **23**, 843–854 (2013).
 7. Preston, J. L. *et al.* High-specificity detection of rare alleles with Paired-End Low Error Sequencing (PELE-Seq). *BMC Genomics* **17**, 464 (2016).
674
675
 8. Zhang, T.-H., Wu, N. C. & Sun, R. A benchmark study on error-correction by read-pairing and tag-clustering in amplicon-based deep sequencing. *BMC Genomics* 1–9 (2016).
676
677
678
 9. Schmitt, M. W. *et al.* Sequencing small genomic targets with high efficiency and
679

- 680 extreme accuracy. *Nat. Methods* 1–4 (2015).
- 681 10. Blokzijl, F. *et al.* Tissue-specific mutation accumulation in human adult stem cells
682 during life. *Nature* **538**, 260–264 (2016).
- 683 11. Welch, J. S. *et al.* The Origin and Evolution of Mutations in Acute Myeloid
684 Leukemia. *Cell* **150**, 264–278 (2012).
- 685 12. Vijg, J., Dong, X. & Zhang, L. A high-fidelity method for genomic sequencing of
686 single somatic cells reveals a very high mutational burden. *Exp. Biol. Med.* **242**,
687 1318–1324 (2017).
- 688 13. Saini, N. *et al.* The Impact of Environmental and Endogenous Damage on Somatic
689 Mutation Load in Human Skin Fibroblasts. *PLoS Genet.* **12**, e1006385 (2016).
- 690 14. Jaiswal, S. *et al.* Age-Related Clonal Hematopoiesis Associated with Adverse
691 Outcomes. *N. Engl. J. Med.* 1–11 (2014).
- 692 15. Genovese, G. *et al.* Clonal hematopoiesis and blood-cancer risk inferred from blood
693 DNA sequence. *N. Engl. J. Med.* **371**, 2477–2487 (2014).
- 694 16. Xie, M. *et al.* Age-related mutations associated with clonal hematopoietic expansion
695 and malignancies. *Nat. Med.* **20**, 1472–1478 (2014).
- 696 17. Martincorena, I. *et al.* Tumor evolution. High burden and pervasive positive
697 selection of somatic mutations in normal human skin. *Science* **348**, 880–886
698 (2015).
- 699 18. Kennedy, S. R. *et al.* Detecting ultralow-frequency mutations by Duplex
700 Sequencing. *Nat. Protoc.* **9**, 2586–2606 (2014).
- 701 19. Chen, L., Liu, P., Evans, T. C., Jr & Ettwiller, L. M. DNA damage is a pervasive

- 702 cause of sequencing errors, directly confounding variant identification. *Science* **355**,
703 752–756 (2017).
- 704 20. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer.
705 *Nature* **500**, 415–421 (2013).
- 706 21. Alexandrov, L. B. *et al.* Mutational signatures associated with tobacco smoking in
707 human cancer. *Science* **354**, 618–622 (2016).
- 708 22. McKerrell, T. *et al.* Leukemia-Associated Somatic Mutations Drive Distinct Patterns
709 of Age-Related Clonal Hemopoiesis. *Cell Rep.* **10**, 1239–1245 (2015).
- 710 23. Zhao, H. *et al.* Mismatch repair deficiency endows tumors with a unique mutation
711 signature and sensitivity to DNA double-strand breaks. *eLife Sciences* **3**, e02725
712 (2014).
- 713 24. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler
714 transform. *Bioinformatics* **25**, 1754–1760 (2009).
- 715 25. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
716 2078–2079 (2009).
- 717 26. Lander, E. S. *et al.* Initial sequencing and analysis of the human genome. *Nature*
718 **409**, 860–921 (2001).
- 719 27. Fujita, P. A. *et al.* The UCSC genome browser database: update 2011. *Nucleic*
720 *Acids Res.* **39**, D876–D882 (2010).
- 721 28. Garrison, E. & Marth, G. Haplotype-based variant detection from short-read
722 sequencing. *arXiv [q-bio.GN]* (2012).
- 723 29. Tan, A., Abecasis, G. R. & Kang, H. M. Unified representation of genetic variants.

- 724 *Bioinformatics* **31**, 2202–2204 (2015).
- 725 30. Rosenthal, R., McGranahan, N., Herrero, J., Taylor, B. S. & Swanton, C.
726 DeconstructSigs: delineating mutational processes in single tumors distinguishes
727 DNA repair deficiencies and patterns of carcinoma evolution. *Genome Biol.* **17**, 31
728 (2016).
- 729 31. Gehrung, J. S., Fischer, B., Lawrence, M. & Huber, W. SomaticSignatures: inferring
730 mutational signatures from single-nucleotide variants. *Bioinformatics* **31**,
731 3673–3675 (2015).
- 732 32. Blokzijl, F., Janssen, R., Van Boxtel, R. & Cuppen, E. MutationalPatterns: an
733 integrative R package for studying patterns in base substitution catalogues. *bioRxiv*
734 071761 (2016). doi:10.1101/071761

Figure 1

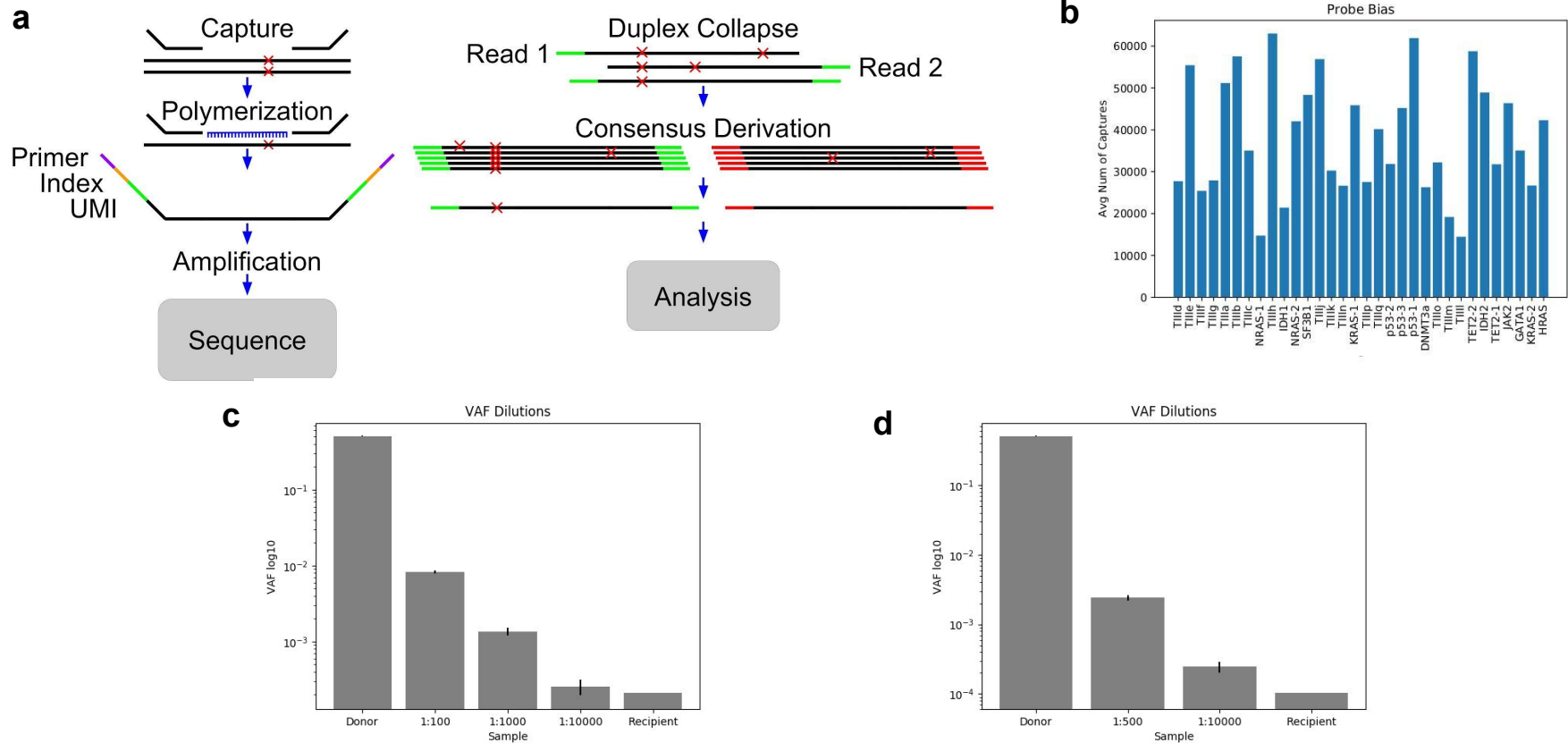


Figure 1 | Amplicon sequencing accurately detects mutation allele frequencies as rare as 1/10,000. **a**, Graphical depiction of gDNA capture and analysis method. **b**, Capture efficiencies vary in a probe dependent manner. **c**, Accurate detection of a single heterozygous SNP in gDNA from one individual diluted into gDNA from another (without this germline SNP) to frequencies as low as 1/10,000. **d**, Accurate detection of three linked SNPs found within the same allele diluted as in **c**. Error shown is standard deviation.

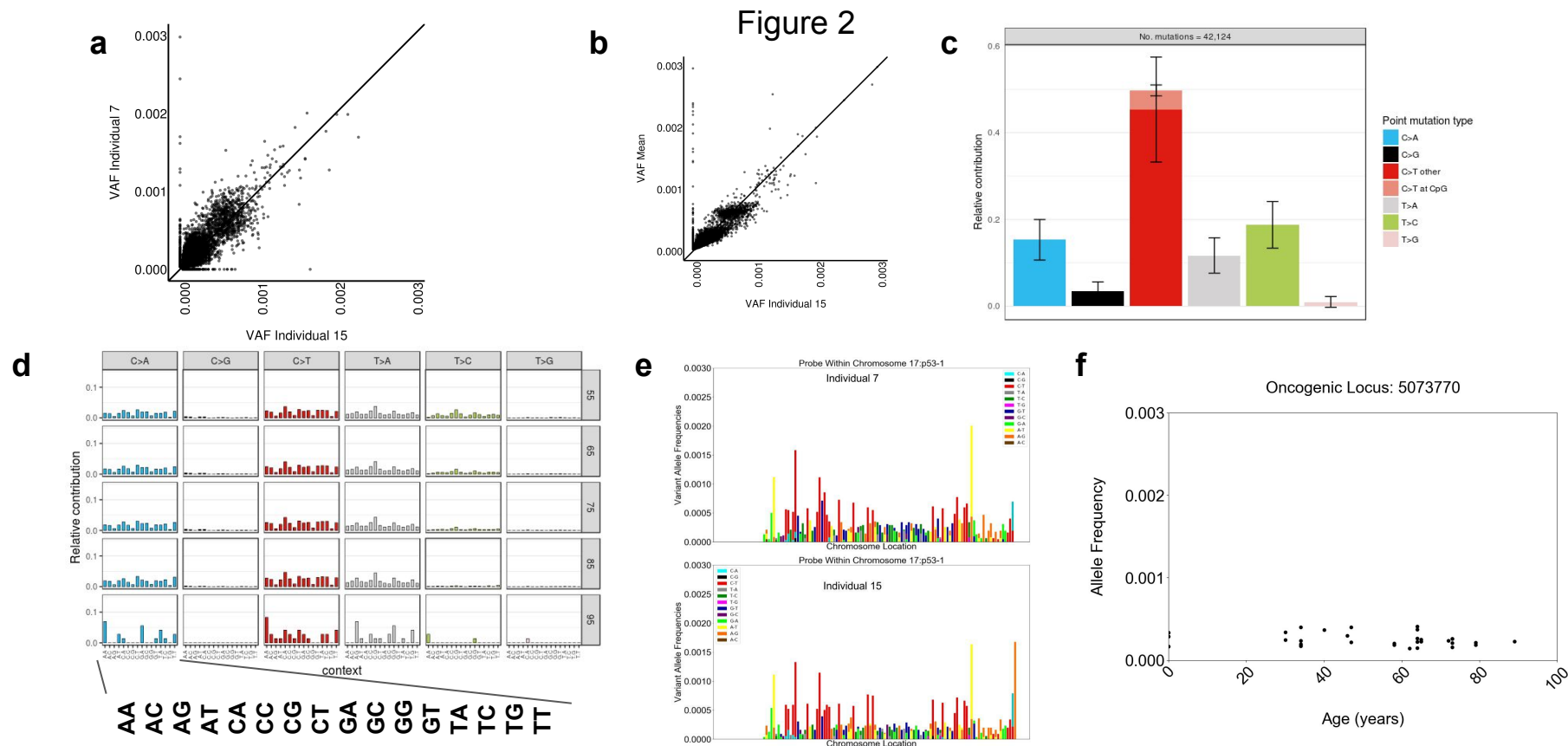


Figure 2 | Mutations exist at conserved frequencies independently of age. **a**, Comparison of VAFs of identified variants within a 34 year old (x-axis) and 62 year old (y-axis); $R^2 = 0.408211$, $p=0.000$. R^2 values unless otherwise noted are calculated for all points falling below VAFs of 0.003 which largely includes all variants but germline. **b**, VAFs from a 34 year old (x-axis) compared to mean VAFs from individuals ranging in ages from newborn to 89 years of age ($n=22$); R -Squared = 0.590412, $p=0.000$. **c**, Relative contribution rates of each base substitution to all substitutions identified. **d**, Relative contribution rates of each base substitution identified by surrounding 5' and 3' nucleotide context. **e**, All identified base substitutions within a probed region are plotted by their position and VAFs for individuals 7 and 15 (representative of most other individuals), revealing highly reproducible patterns. **f**, Oncogenic VAFs plotted as a function of donor age show little evidence of clonal expansion.

Figure 3

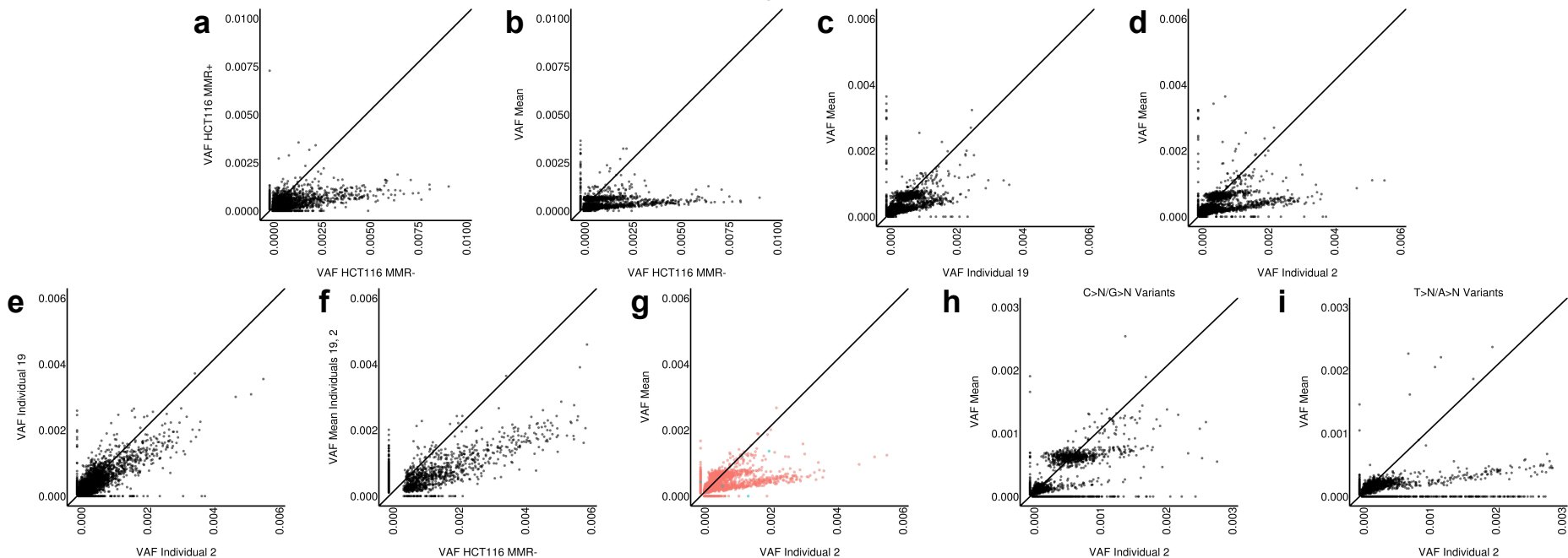


Figure 3 | Individuals Can Systematically Deviate from Population Average. **a**, Comparing VAFs in HCT116 MMR+ vs MMR^{MT} cells reveals an increase in frequencies for many of the observed variants in MMR^{MT} cells (R-Squared = 0.211479). **b**, MMR^{MT} vs mean VAFs from blood of the 22 individuals shows a similar pattern of increased VAFs as the comparison with parental (R-Squared = 0.120895). **c**, blood from a 73 yr old person (individual #19) compared to the mean VAFs reveals a deviating population of variants that exist at an increased frequency compared with average VAFs (R-Squared = 0.387125). **d**, A cord blood sample (individual #2) also shows a subset of variants with higher frequencies than in the average (R-Squared = 0.278250). **e**, VAFs from individual #2 vs individual #19 reveals that the deviating variants are at the same positions causing the comparison to fall close to the y=x line (R-Squared = 0.613542). **f**, Plotting the mean for VAFs from individuals #2 and #19 versus VAFs from MMR^{MT} HCT116 cells reveals that the variants within the blood are the same as those found within the MMR^{MT} cell line. While variant frequencies are higher in the MMR^{MT} cell line, the identities of the deviating variants are the same (R-Squared = 0.587474). **g**, Variants detected in individuals #2 and #19 are not enriched for oncogenic changes, indicated in blue **h**, Plot of only C>N/G>N variants shows relative similarity between MMR- and parental cells (R-Squared = 0.350623). **i**, Plot of only T>N/A>N variants reveals that the majority of deviating variants between MMR^{MT} and parental cells are substitutions affecting T or A.

Extended Data Figure 1

a

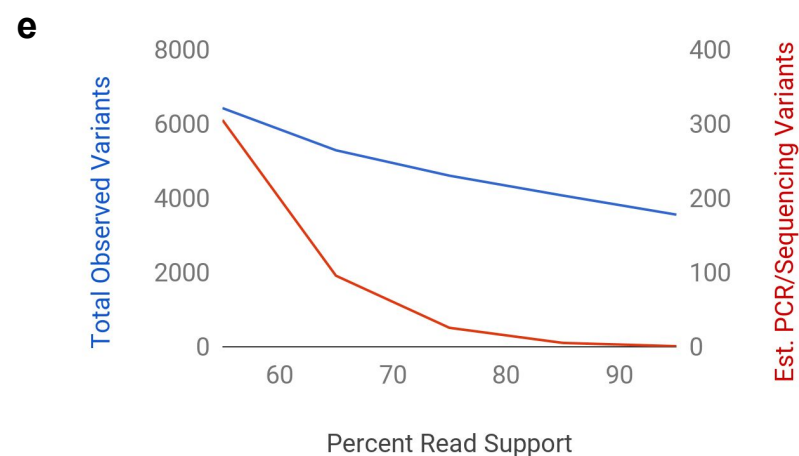
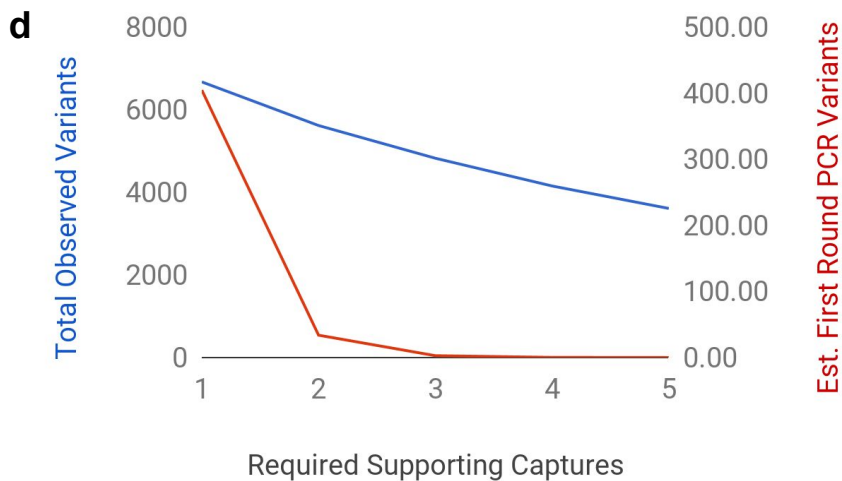
Supporting Captures	Duplex	Mock-Duplex	% Vars Eliminated
4	4240	4264	0.56285
3	4912	4928	0.32468
2	5704	5734	0.52319
1	6760	6794	0.50044

b

Enzyme	Error Rate (mut/base)	Unique UMIs	Captures per UMI	Total Amplicon Size	# Bases In First Amplification	Total Errors
Phusion HF Buffer	0.00000044	2818388	88075	4838	426105036	187
Phusion GC Buffer	0.00000095	2818388	88075	4838	426105036	405

c

Supporting Captures	1	2	3	4	5
	187.49	7.27	0.28	0.01	0.00
	404.80	33.87	2.83	0.24	0.02



Extended Data Table 1

Cohort 1

a

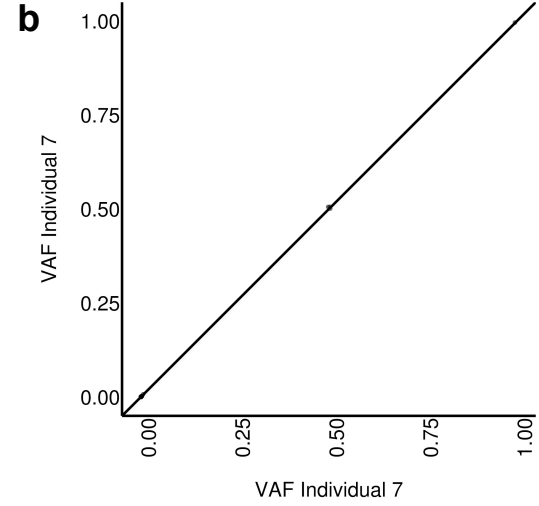
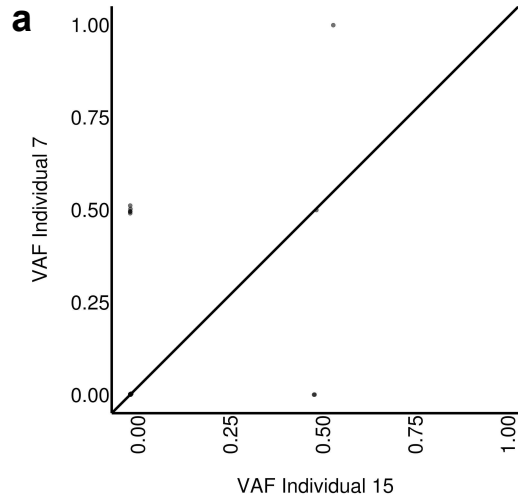
Individual	Age (years)
1	0
2	0
3	0
4	34
5	34
6	30
7	34
8	46
9	47
10	40
11	59
12	59
13	58
14	62
15	65
16	64
17	64
18	73
19	73
20	72
21	79
22	89

Cohort 2

b

Individual	Age (years)
25	0
26	34
27	44
28	43
29	46
30	44
31	46
32	49
33	41
34	57
35	62

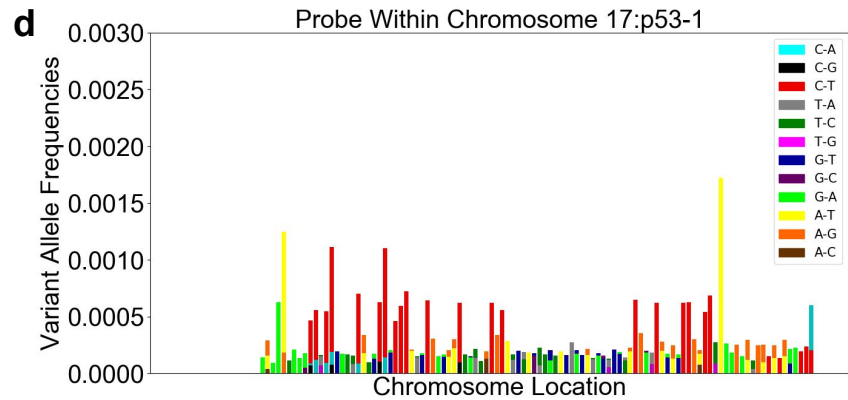
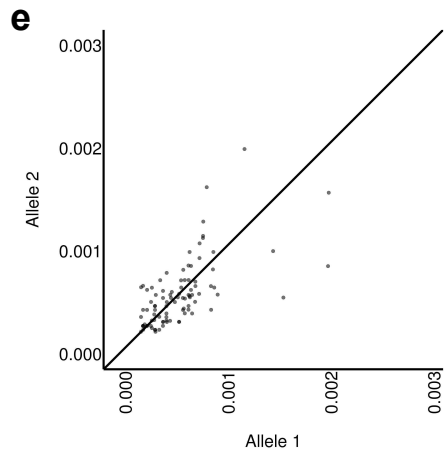
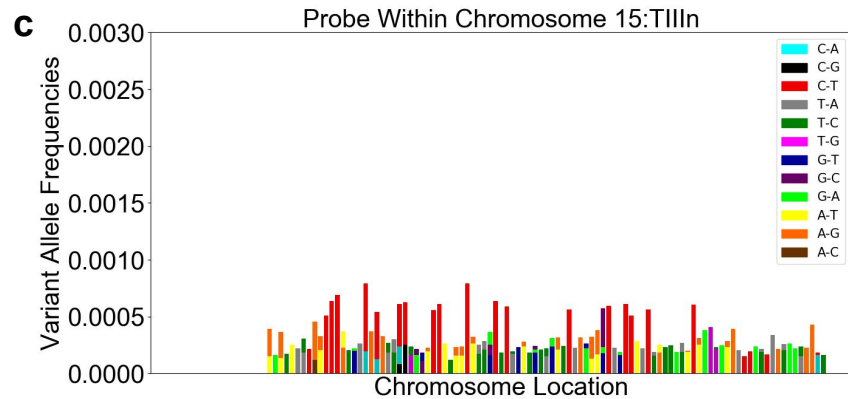
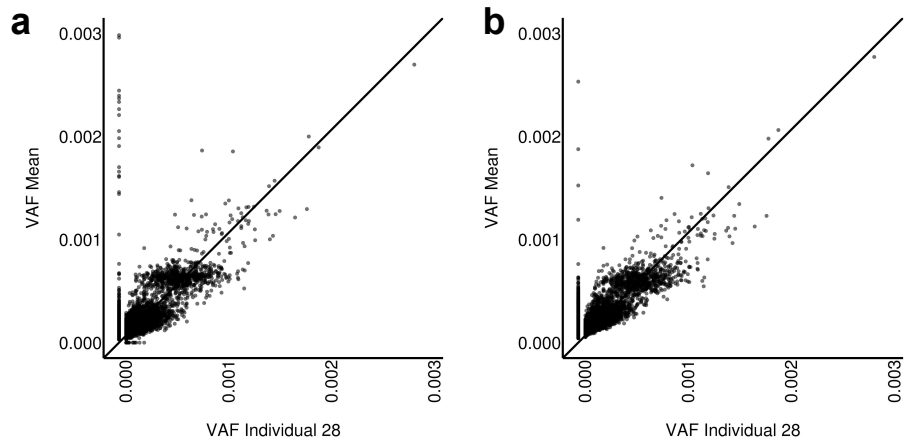
Extended Data Figure 2



c

Individual	0mo vs 1mo
Individual A	0.460348
Individual B	0.538478
Individual C	0.436766
Individual D	0.522387
Individual E	0.519219
Individual F	0.482805

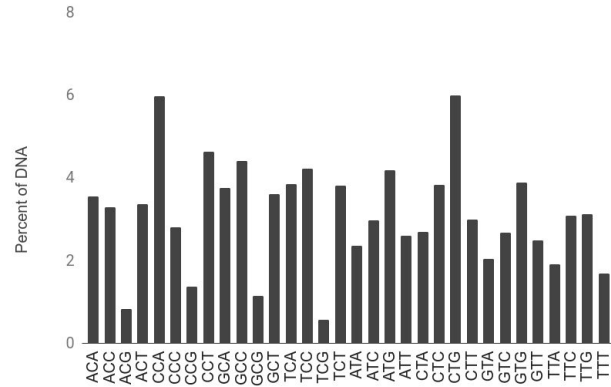
Extended Data Figure 3



Extended Data Figure 4

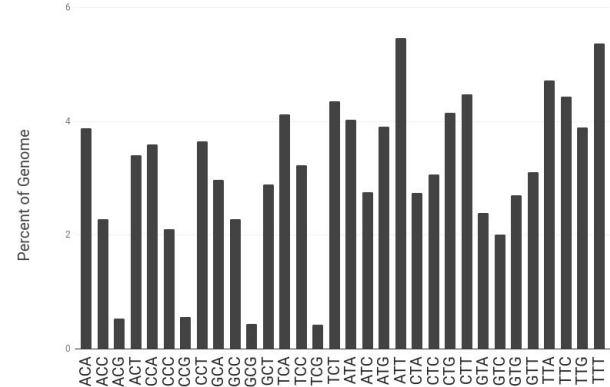
a

Trinucleotide Representation Probed Region



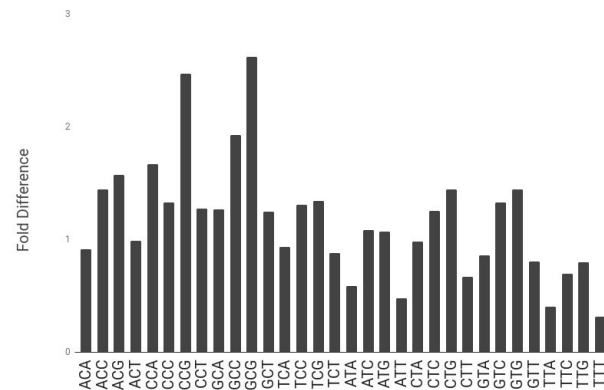
b

Trinucleotide Representation Human Genome



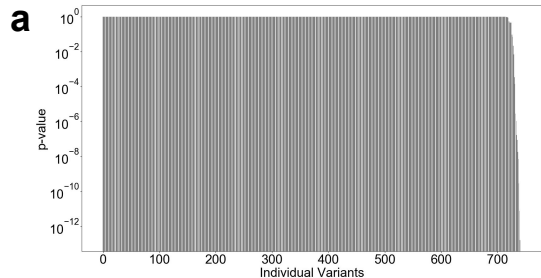
c

Trinucleotide Representation (Probed Region/hg19)

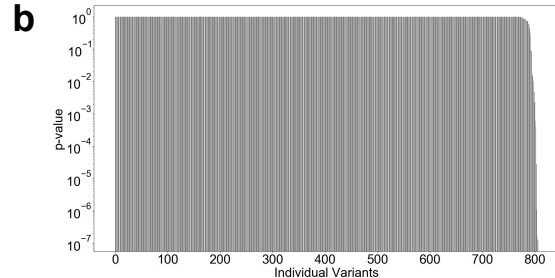


Extended Data Figure 5

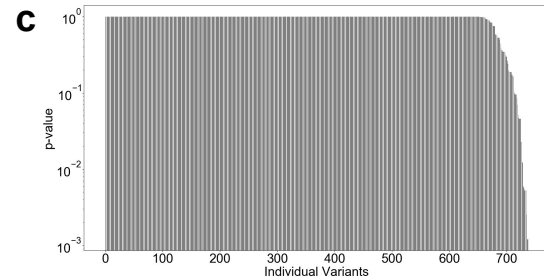
A



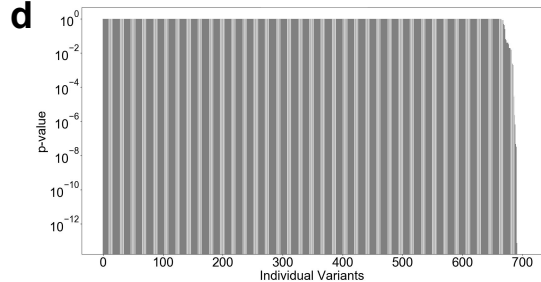
C



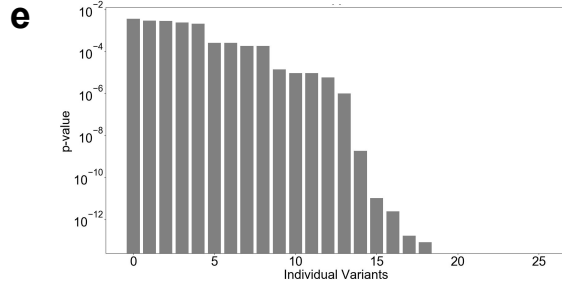
G



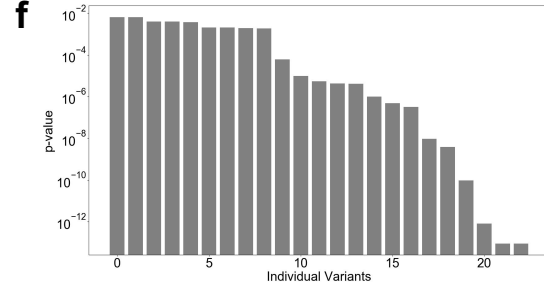
T



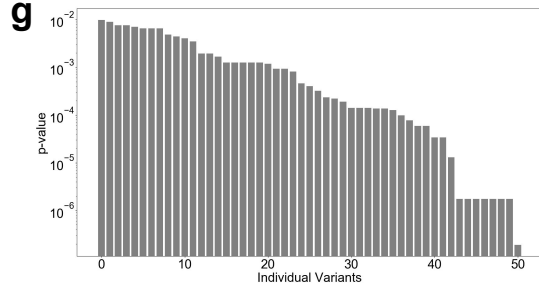
A



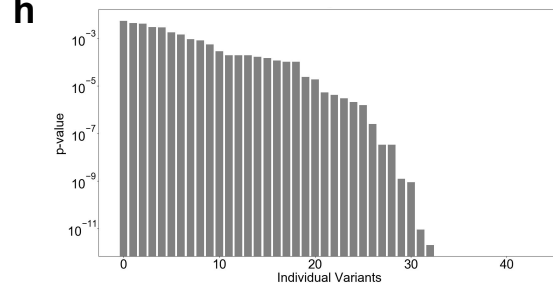
C



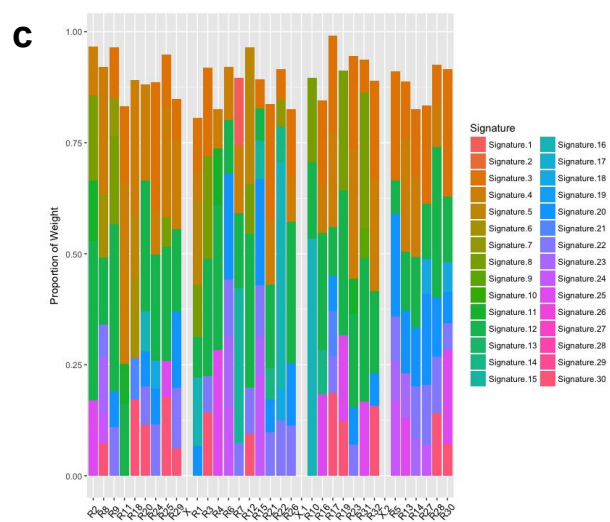
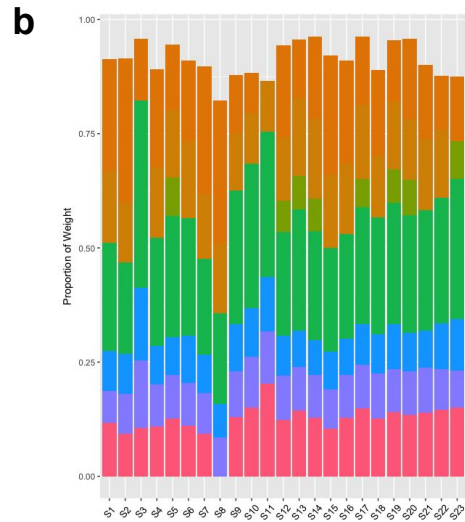
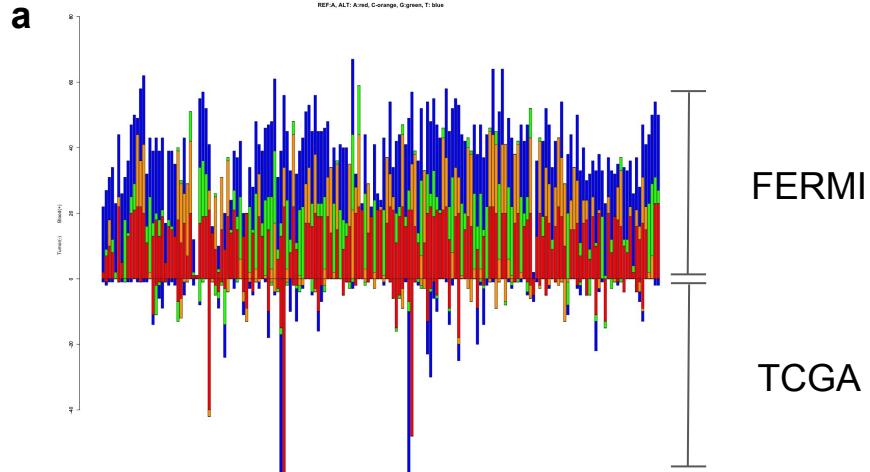
G



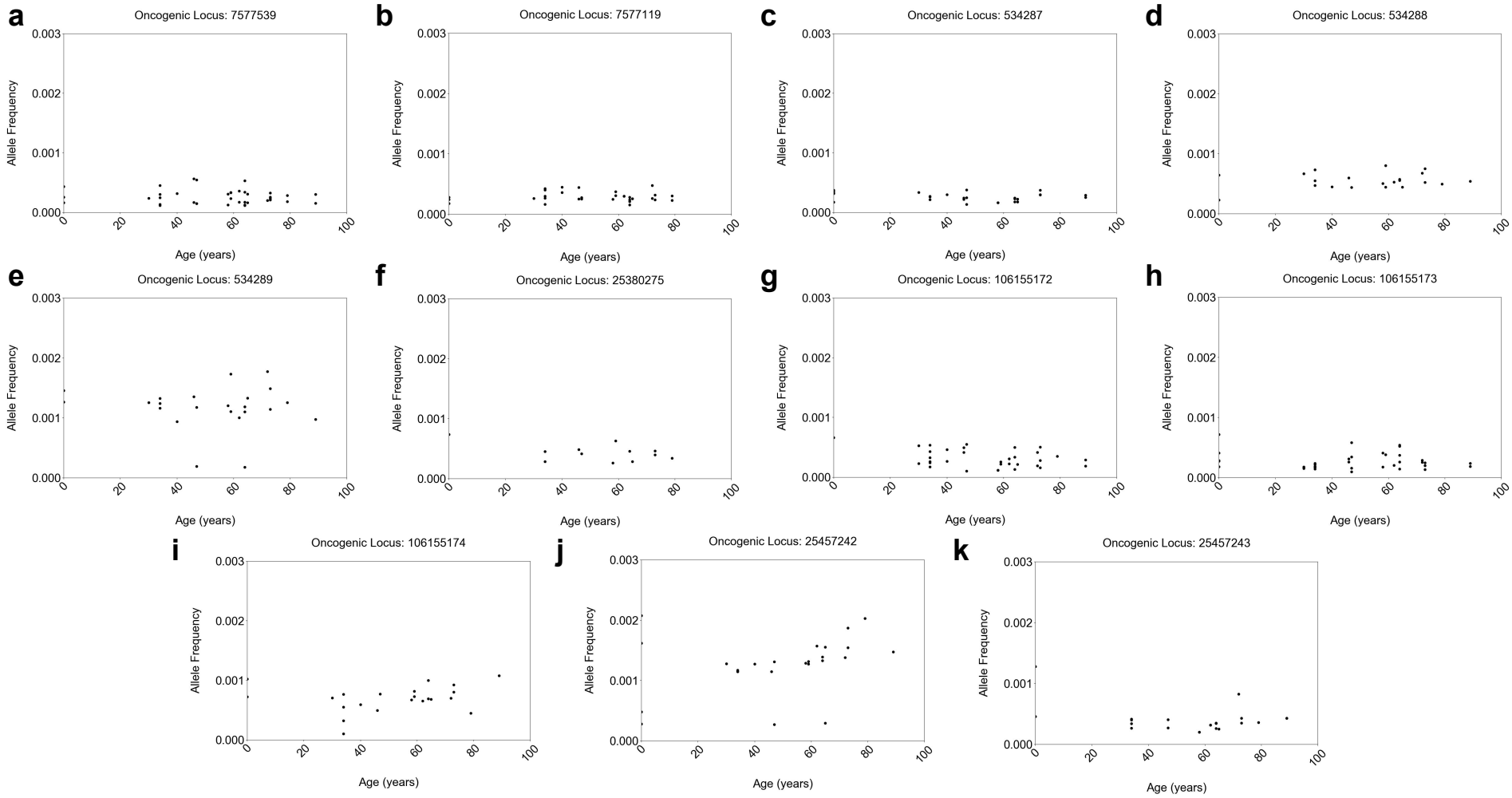
T



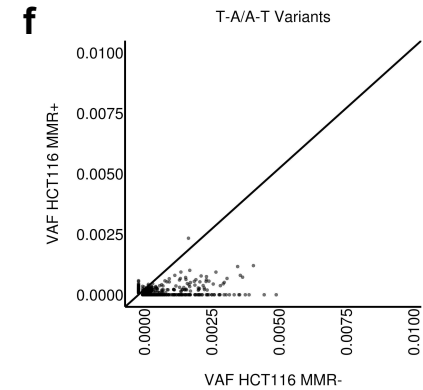
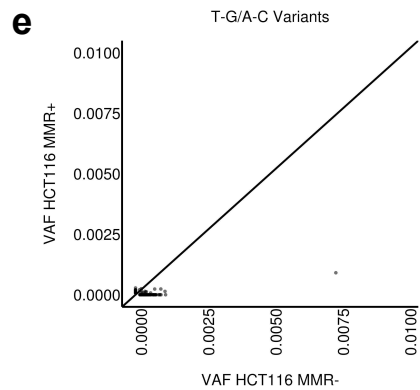
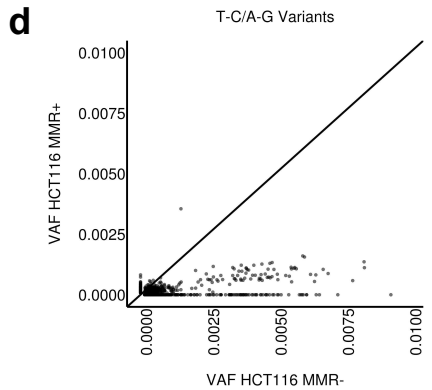
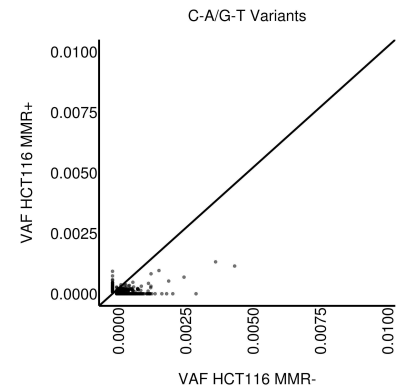
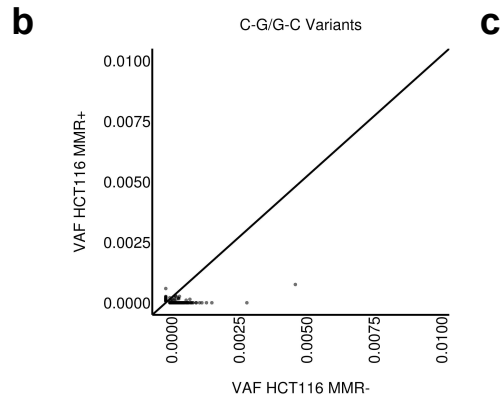
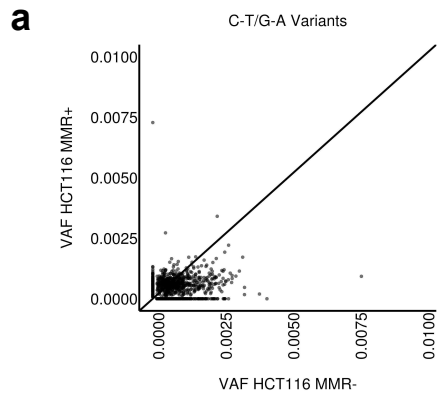
Extended Data Figure 6



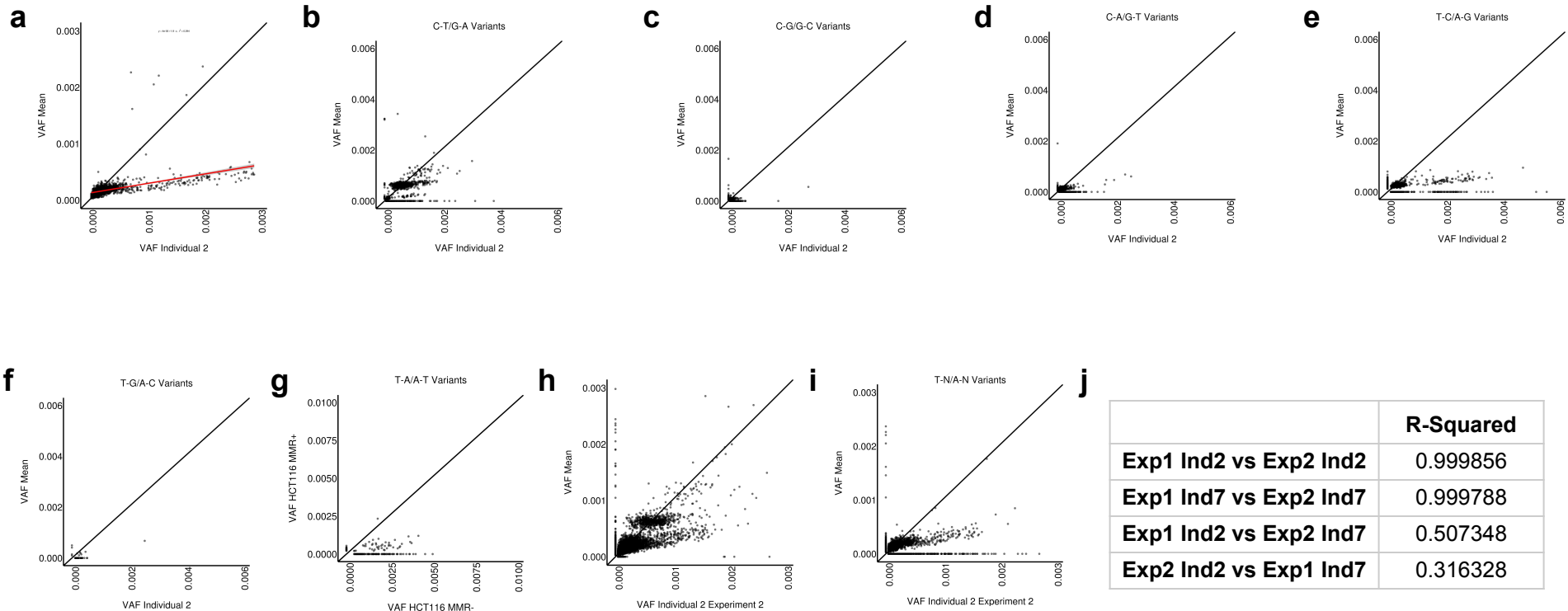
Extended Data Figure 7



Extended Data Figure 8



Extended Data Figure 9



	R-Squared
Exp1 Ind2 vs Exp2 Ind2	0.999856
Exp1 Ind7 vs Exp2 Ind7	0.999788
Exp1 Ind2 vs Exp2 Ind7	0.507348
Exp2 Ind2 vs Exp1 Ind7	0.316328

1 **Extended Data Figure 1: Estimation of false-positive rates due to sequencing and**
2 **PCR errors.**

3 **a**, The use of sequencing information found within Read 1 and Read 2 of paired-end
4 sequencing is often used to correct sequencing errors. We performed paired-end
5 collapsing prior to consensus read derivation (Fig. 1a), though the effect was
6 surprisingly mild. In this table, the number of identified variants are shown when duplex
7 collapsing is used or not in consensus read derivation (mock duplexing processes the
8 collapsing in the exact same way as duplex collapsing without eliminating variants for
9 not being in both reads). These variant counts are shown while also varying the number
10 of required independent supporting captures for a variant to pass filtering. The logic
11 behind this analysis is that the fewer captures in which a variant is found, the less
12 confidence we have that it represents true biological signal. Lower confidence variants
13 should be more likely to be eliminated by duplex collapsing reads, if other filters were
14 otherwise insufficient. We show that whether reads are first duplex collapsed or not,
15 there is little effect on the percent of variants that are eliminated, suggesting that our
16 other filtering parameters appear to adequately eliminate sequencing errors. **b**, While
17 the filters used for FERMI should eliminate the majority of errors introduced during PCR
18 amplification and those errors arising from sequencing mistakes, errors made in the first
19 round of PCR amplification could be identified as false positives. If there is a sufficient
20 number of PCR errors made within the first round of amplification, these errors could
21 create artificial patterns within the data. Using one supporting capture as the lower limit
22 for variants to be identified as true signal, the expected number of errors were estimated

23 from amplification using Phusion polymerase and are shown in the table (two
24 estimations are included because Illumina's reaction mixtures are proprietary and we do
25 not know the exact reaction conditions). **c**, When only requiring one supporting capture,
26 3-6% of variants should be derived from first round PCR errors, although more than half
27 of these will be eliminated by the requirement that 55% of reads for a capture support
28 the variant (errors from subsequent PCR rounds will be even more efficiently eliminated
29 by the 55% cutoff). If we require that the same variant be present at the same location
30 across multiple captures before it is included in the final results, it becomes
31 exponentially more unlikely that a first round PCR error would get included. In contrast,
32 increased capture number requirements have a much more modest effect on variants
33 called. **d**, While increasing the number of required supporting captures eliminates rare
34 variants as well as first round PCR errors, the numbers of identified variants only
35 decreases modestly for all individuals (blue line, left y-axis). In contrast, the number of
36 variants expected to be identified as a result of first round PCR amplification errors
37 exponentially decreases with each extra capture requirement (red line, right y-axis).
38 When compared to the number of variants that pass all filters and processing, the first
39 round PCR errors appear to have minimal effect even when only a single capture is
40 required. Expectedly, as we increase the number of required captures supporting a
41 variant, the total number of variants also decreases, and after two required captures
42 should essentially not include mutations created by PCR amplification. Throughout most
43 of this paper, a single capture is used, so as to not bias results by variant
44 representation. Nonetheless, the patterns of mutations identified look very similar when

45 greater numbers of supporting captures are required. **e**, As shown in Fig. 1a, when
46 deriving consensus reads, variants are eliminated for being rarely observed across
47 reads supporting a given capture. The cutoff we use throughout most of this manuscript
48 is 55%, such that a given variant must be present in at least 55 percent of sequencing
49 reads supporting a capture or they are ignored. The logic behind this chosen cutoff is
50 that more stringent cutoffs largely do not alter the observed mutation spectra, but result
51 in a significant loss in putatively true positive signal. With this cutoff, the expected
52 number of sequencing errors can be estimated. We observe that 9 percent of bases are
53 mismatched within reads supporting a given capture. Each capture is approximately
54 150bp in length and is supported by an average 13.5 reads. This yields an average of
55 182.25 errors within each sequenced capture.

$$56 \quad E_{tot} = 0.09 * 150 \text{ bp} * 13.5 \text{ reads}$$

$$57 \quad E_{tot} = 182.25$$

58 Applying the requirements that 55-95 percent of reads must support a given variant
59 (shown as m), the number of false positive signals that pass filtering for each prepared
60 blood sample can be computed. Within each capture there are approximately 450 total
61 possible changes, and an average of 18 reads supporting each capture:

$$62 \quad E_{seq} = m * 18 \text{ reads/capture} * \frac{182.5 \text{ PCR err}}{450 \text{ bp}} * 1200000 \text{ captures/sample}$$

$$63 \quad m = 0.55 : E_{seq} = 155.95 \text{ errors/sample}$$

$$64 \quad m = 0.65 : E_{seq} = 31.48 \text{ errors/sample}$$

$$65 \quad m = 0.75 : E_{seq} = 6.19 \text{ errors/sample}$$

$$66 \quad m = 0.85 : E_{seq} = 1.22 \text{ errors/sample}$$

$$m = 0.95 : E_{seq} = 0.24 \text{ errors/sample}$$

The number of expected PCR amplification errors to pass all cutoffs is then estimated using a Gaussian distribution. The logic is that the first round of PCR amplification will create errors that will be at an allele frequency near 50 percent as an error will be created in one of two strands of a captured sequence. Using a Gaussian distribution with a mean at 50, the number of all PCR amplification errors expected to pass the 1 supporting capture and 55-95 percent of sequencing reads criteria can be calculated by integrating under the Gaussian distribution. Since we expected about 405 first round PCR amplification errors, and subsequent errors will exist at much smaller allele frequencies, the expected number of variants expected to pass criteria is calculated as follows:

$$E_{tot} = 405 * \int_c^{100} f(x) + m_c$$

Above we integrate from the support allele frequency c to 100 under the Gaussian distribution $f(x)$, multiply this by the expected total number of first round PCR amplification errors, and add to this the number of expected sequencing errors m as a function of the support frequency c . As shown here, when variants must be supported by at least one unique capture and at least 55 percent of supporting reads, we anticipate only about 150 total variants false variants to make through all FERMI analysis. We believed this to be an acceptable amount of noise given that we see about 6000 total variants from each sample and generated most of the data in this manuscript with these criteria.

67 **Extended Data Table 1: Cohort of sequenced individuals.**

68 **a**, This table contains the ages of the individuals used throughout the manuscript, and
69 their corresponding sample numbers. Those samples shown as age '0' are cord blood
70 samples that had been previously banked. All other samples were taken from
71 apparently healthy blood donors that passed the requirements to donate blood. **b**, This
72 table contains the ages of individuals used to ensure that the data generated by FERMI
73 was not experiment specific. These samples were used as the comparison to generate
74 Extended Data Figs. 3a-b.

75 **Extended Data Figure 2: Resequenced samples are not more similar to each other**
76 **than to other individuals.**

77 **a**, Low frequency variants tend to exist close to a $y=x$ line, while high frequency SNPs
78 differ across individuals. As expected, such SNPs cluster around frequencies of 0.5 and
79 1 (R-Squared=0.243364). **b**, When samples are re-sequenced, they show a high degree
80 of similarity, both among SNPs and more rare variants (R-squared=0.568749). **c**,
81 Though repeat sequencing of individuals typically results in close matches of VAF,
82 repeats do not more closely each other than they match the VAF population mean or
83 any other typical sample. This suggests that the differences observed between samples
84 is likely due to sampling differences than to real differences in individual mutation loads.

85 **Extended Data Figure 3: Variants detected represent multiple independent events**
86 **and reproduce across multiple experiments.**

87 For consistency, all samples used in the main analysis derive from a single bulk library
88 preparation and sequencing run. To ensure that the observed trends are not the result
89 of some bias specific to this single preparation, the entire process was independently
90 repeated, with eleven different blood biopsies (Cohort 2). **a**, Cohort 2 samples closely
91 resembled averaged allele frequencies from the Cohort 1 (R-squared = 0.455316,
92 p-value = 0.000000). **b**, Comparing Cohort 2 samples against the VAF mean created
93 from Cohort 2 samples produces a similar pattern to the same comparison using the
94 Cohort 1 data (R-Squared = 0.615327, p-value = 0.000000). **c-d**, Similar mutation
95 patterns along captured regions were observed for Cohort 2 as for cohort #1 (Fig. 2e).
96 **e**, To understand if observed variant frequencies are the result of clonal expansions or
97 independent events, heterozygous variants were separated by allele. The logic behind
98 this analysis is that if independently captured variants result from the same original
99 event (i.e. a clone), then these variants should be found on the same allele.
100 Alternatively, if variants result from independent events, then such variants should be
101 frequently found on both alleles. By following linkage between variants and
102 heterozygous SNPs, the two alleles can be distinguished. Shown here are the allele
103 frequencies of variants found on either Allele 1 along the x-axis or Allele 2 along the
104 y-axis (analyses are restricted to genomic segments from individuals containing
105 heterozygous SNPs). As the variants adhere to a $y=x$ line, they appear randomly
106 distributed between both alleles, suggesting that variants detected represent multiple
107 independent events rather than clonal expansions.

108 **Extended Data Figure 4: Triplet prevalence in probed regions does not sufficiently**
109 **explain base bias.**

110 To understand how representative our total captured region was of the overall human
111 genome, the trinucleotide sequence counts **a**, found within our 32 probes was
112 compared to **b**, the overall trinucleotide counts found within hg19. CpG sites were less
113 prevalently mutated in our samples than previously observed in other tissues and
114 cancers. The lower incidence numbers of CpG mutations does not appear to be due to
115 any effect of undersampling within our selected probe regions, as shown by **c**, the fold
116 difference in the number of triplets found in our probed region and in the hg19 reference
117 genome. Note that these analyses are of total sequence, not identified variants.

118 **Extended Data Figure 5: Multiple positions show nonrandom base bias.**

119 Not only is there significant conservation in the bases to which a position will change
120 across individuals, but many locations are only observed to mutate to a single base. To
121 understand the likelihood of this pattern arising due to random chance, every instance
122 of a given substitution was quantified for each probed site across all individuals. These
123 changes were used to derive an overall probability that each base would change to any
124 of the other 3 bases if mutated. Using a chi-squared algorithm to test goodness of fit,
125 individual probabilities were computed for the base substitution pattern observed at
126 each base locus. These probabilities were then multi-comparison corrected using
127 Bonferroni correction, separated by reference base, ordered in descending order, and
128 plotted here. When a variant was only observed in a small number of individuals, the

129 probability of this change exclusively occurring at a given location due to chance was
130 relatively high, resulting in a substantial number of non-significant loci (**a-d**; p values
131 ~ 1). Plotting only positions exhibiting significant bias reveals a substantial number of
132 bases that predictably mutate across individuals in a manner unlikely to be explained by
133 chance (**e-h**; p values that approach zero lack bars). The total number of variants
134 passing significance for each base are: A) 27 C) 23 G) 51 T) 44. This suggests that
135 sequence context and base location may both be playing significant roles in determining
136 the substitution probabilities for a number of base positions throughout the genome.

137 **Extended Data Figure 6: Blood shows previously identified signatures but is**
138 **different from cancers**

139 **a**, We focused on the amplicons in coding regions, and integrated Pan cancer somatic
140 mutation data from exome sequencing in the TCGA to analyze patterns of base
141 substitutions at genomic positions in the target regions which were mutated in both
142 blood and tumor genomes. Substitution frequency and substitution patterns were both
143 significantly different between blood and tumors, both at highly mutated sites (mutation
144 count > 10 ; Chi square test; FDR adjusted p-value < 0.05) and across all such sites
145 (Mantel test; p-value $< 1e-5$), with substitution patterns in tumor genomes being more
146 skewed. It is possible that selection during cancer evolution (as opposed to nearly
147 neutral evolution in terminally differentiated blood cells) contribute to the observed
148 patterns. **b**, Integrating trinucleotide contexts of the substitutions, we determined the
149 contributions of different mutation signatures previously identified. Out of 30 previously

150 identified signatures, our data showed overrepresentation of only 7 of them (Signatures
151 3, 4, 8, 12, 20, 22 and 30) across different samples. Out of seven signatures, Signature
152 12, 3 and 4 had maximum contributions. Signature 3 and 4 are known to be associated
153 with failure of DNA double stranded break repair by homologous repair mechanism and
154 tobacco mutagens respectively, whereas the aetiology of Signature 12 remains
155 unknown. **c**, There was no systematic difference in mutation signatures between
156 amplicons when grouped by their genomic context, and they also showed similar
157 pattern of enrichment of few signatures as compared to others, with signature 12, 3 and
158 4 having maximum contributions. Signature 12 and 4 exhibits transcriptional strand bias
159 for T>C and C>A substitutions respectively, whereas signature 3 is associated with
160 increased numbers of large InDels.

161 **Extended Data Figure 7: Oncogenic mutations do not show evidence of selection.**

162 As shown in Fig. 2f, known oncogenic mutations within probed regions do not show
163 evidence of positive selection. Shown here are additional probed oncogenic loci
164 according to their observed VAFs across donor ages, which also do not show an
165 increase in variant allele frequency in older ages.

166 **Extended Data Figure 8: MMR^{MT} VAFs are elevated over parental frequencies.**

167 When compared to MMR sufficient HCT116 parental cell line genomic DNA, MMR
168 deficient HCT116 cell DNA (R-Squared = 0.066023) contains substitution mutations at
169 significantly elevated frequencies, as expected with DNA repair deficiencies (Fig. 3a-b).

170 Although most VAFs appear elevated within MMR deficient cells, the magnitude of
171 increase was context dependent. Base substitutions altering **a-c)** C or G exhibited
172 elevated allele frequencies in MMR^{MT} cells, but substantially less compared to **d-f)** T or
173 A nucleotides, which exhibit much higher VAFs compared to parental.

174 **Extended Data Figure 9: Base bias for cord blood individual #2 resembles MMR^{MT}**
175 **Cells.**

176 As for comparisons of MMR^{MT} and HCT116 parental cell lines, a cord blood donor
177 showed a variant population that significantly deviated from expected VAFs (Fig. 3d). **a,**
178 The mutation spectrum found within individual 2 fits to a linear regression line of
179 $y=1.9x+0.00004$, from which it can be seen that variants are approximately twofold
180 more prevalent than in the overall population average. Similar to the data in Extended
181 Figure 8, base substitutions altering **b-d)** C or G nucleotides did not show elevated
182 frequencies. As in the in the MMR^{MT} cells, **e-g)** T or A changes appear at elevated
183 frequencies. Data from individual 19 looked similar to the data shown here, but is not
184 shown. **h,** To ensure that the increased frequencies of variants are not the result of
185 some experimental anomaly, the DNA from individuals #19 (not shown) and #2 was
186 used in a second experiment. In the experimental repeat, the samples showed nearly
187 identical mutational spectra, with similarly elevated levels of T or A changes. **i,** T or A
188 changes again appear at elevated frequencies in a similar manner to the first
189 experiment. The deviating population fits a regression line of $y=2.2x-9.6*10^{-5}$. **j,**
190 Indicative of experimental repeatability, when samples were freshly captured and

191 sequenced using FERMI, the same individual was highly similar across experiments,
192 and different individuals were less similar. R^2 values are calculated to include all
193 variants, including germline.