

1 Conserved patterns of somatic mutations in human peripheral blood

2 cells

³ L. Alex Liggett¹, Anchal Sharma², Subhajyoti De², and *James DeGregori^{1,3,4,5}

⁴ **¹Department of Biochemistry and Molecular Genetics, ²Rutgers Cancer Institute,
⁵ New Brunswick, NJ 08901, ³Integrated Department of Immunology, ⁴Department of
⁶ Pediatrics, ⁵Department of Medicine, Section of Hematology, University of
⁷ Colorado School of Medicine, Aurora, CO 80045**

8 *Corresponding Author

9 Corresponding Author:

10 James DeGregori, Ph.D.

11 James.DeGregori@ucdenver.edu

12 **Summary**

Mutation accumulation varies across a genome by chromosomal location, nucleotide identity, surrounding sequence, and chromatin context¹⁻⁵. Nevertheless, while mutagens, replication machinery, and repair processes exhibit identifiable mutation signatures, at the tissue or organismal scale the aggregate somatic manifestation of these processes has been difficult to measure, and often appears to be semi-random.

18 This randomness is often believed to contribute to the stochasticity of diseases like
19 cancer⁶ and physiological decline during aging⁷. The challenge in observing any
20 tissue-wide somatic mutation patterns is that prior to clonal expansion, most mutations
21 are rare in healthy tissue^{8–11}. Here we describe a new method called FERMI (Fast
22 Extremely Rare Mutation Identification), which comprehensively captures and quantifies
23 rare mutations at single DNA molecule resolution that exist at frequencies as rare as
24 10^{-4} in human peripheral blood. Using this method, we observed an unanticipated
25 degree of ubiquity and similarity between the somatic mutation loads of different
26 individuals, where most assayed substitutions are found to occur at conserved
27 frequencies across nearly all individuals spanning a nine decade age range. These
28 observed mutational patterns existed both within non-conserved, non-coding and
29 non-repetitive regions of the genome and within the coding regions of oncogenes
30 implicated in hematopoietic malignancies. Furthermore, we find that nucleotides
31 preferentially mutate to particular bases in a manner that is specified by nucleotide
32 identity, position, and sequence context. Finally, we identify individuals who deviate from
33 typical mutational patterns in a reproducible manner that resembles a mild mismatch
34 repair deficiency, suggesting that variation in somatic mutation rates may be relatively
35 common. This study provides one of the first estimates of mutation burden in terminally
36 differentiated somatic cells and demonstrates that somatic mutations in such cells are
37 significantly more frequent and deterministic than previously believed, and are governed
38 by mechanisms that when perturbed, result in predictable outcomes.

39 Measuring somatic mutations has been technically challenging because
40 mutations occur within individual cells that do not necessarily clonally expand to
41 detectable representation. While these challenges have been somewhat overcome by
42 increasing the depth of sequencing, using clever methods of barcoding⁸ or by
43 performing paired strand collapsing¹², it remains difficult to get enough sequencing
44 depth and breadth while sufficiently limiting false positive noise^{8,11,13}. To overcome these
45 sequencing limitations, we created FERMI, in which we adapted the amplicon
46 sequencing method of Illumina's TrueSeq Custom Amplicon platform to target only 32 x
47 150bp genomic regions, spanning AML-associated oncogenic mutations and the Tier III
48 regions of the human genome (non-conserved, non-protein coding and non-repetitive).
49 We further improved upon Illumina's capture efficiency to achieve approximately 1.2
50 million unique captures from 500ng-1µg of genomic DNA (gDNA) (See Methods). The
51 targeting probes used in gDNA capture were designed with 16bp index of sequence
52 unique to each individual and a 12bp unique molecular identifier (UMI) of random DNA,
53 ideally unique to each capture (Fig. 1a). Sequencing reads were sorted both by sample
54 index and UMI, producing bins of single cell sequencing which were collapsed to
55 produce relatively error-free consensus reads. Captures were only considered if
56 supported by at least 5 reads, and variants were only included if identified in both
57 paired-end sequences, and detected in at least 55% percent of supporting reads (Fig.
58 1a and Methods; see also Extended Data Figure 1).

59 While all probed regions were successfully captured and amplified, capture
60 efficiency was variable and dependent on probe identity (Fig 1b). To understand assay
61 sensitivity, log-series dilutions of human heterozygous single nucleotide polymorphisms
62 (SNPs) were prepared and assayed by FERMI. Using these dilutions, we observed
63 robust quantification of diluted SNPs as rare as 10^{-4} (Fig. 1c). Even more accurate
64 quantifications of SNP frequency can be made when using strand information to follow
65 dilutions of multiple SNPs located on the same allele (Fig. 1d). For more description of
66 the accuracy of FERMI, see Elimination of false positive signal in Methods and
67 Extended Data Figure 1.

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

81 (Extended Data Figure 2c). Variant allele frequencies (VAFs) were averaged across 22
82 sampled individuals and used as a comparison to individuals, which appear
83 age-independent and still adhere to a $y=x$ line (R^2 Range = 0.426-0.631, Mean = 0.558)
84 (Fig. 2b), and are similar across experiments (Extended Data Figure 3a-d shows data
85 from an additional 11 individuals). Variants with frequencies above 0.001 were found in
86 nearly all samples, while more rare variants were missed with a probability inversely
87 proportional to their allele frequencies. Furthermore, most variants likely represent
88 multiple independent events rather than clonal expansions, as they are found at similar
89 frequencies on both alleles (Extended Data Figure 3e). It thus appears that instead of
90 being semi-random, the aggregate effect of all DNA damage and maintenance
91 generates somatic mutations at predictable rates throughout the genome independent
92 of age. We suspect that such mutations primarily arise in terminally differentiated blood
93 cell types in a sequence context-dependent manner, without constraints imposed by
94 developmental lineages and selection, such that it reflects the basal DNA damage and
95 repair errors in hematopoietic cells.

96 We observed that the overall probability of a substitution occurring is biased by
97 nucleotide identity, with C>T substitutions being the most common and T>G
98 substitutions being the least common (Fig. 2c). These biases were largely expected, as
99 similar patterns have been observed both in other healthy tissues and in cancers¹⁴⁻¹⁸.
100 There were notable differences, especially for C>N changes which we observe as
101 underrepresented within a CpG context (Fig. 2d). Regardless of functional or oncogenic
102 potential, each site tends to undergo the same substitutions across individuals (Fig. 2e).

103 These conserved substitution rates appear to be deterministic, and cannot be explained
104 by undersampling (Extended Data Figure 4) or known base change biases (Extended
105 Data Figure 5). It therefore appears that the combined sources of external and internal
106 DNA mutation result in systematic substitutions at frequencies that are often predictable
107 by location and sequence context. Suggestive of differences during cancer evolution
108 and normal somatic mutation, the integrated exome sequencing pan cancer somatic
109 mutation data from the TCGA exhibits different substitution patterns from those that we
110 find in healthy donor blood (Extended Data Figure 6a). Using the trinucleotide contexts
111 of the substitutions, 7 out of 30 previously identified mutations signatures were
112 identified, and these signatures did not differ significantly across sampled genomic
113 segments (Extended Data Figure 6b-c).

114 While we observe variants at conserved frequencies across many individuals,
115 previous studies have shown that selection for oncogenic changes can increase certain
116 mutation frequencies with age¹⁷. While we observe each queried oncogenic change in
117 every biopsied individual regardless of age, we do not observe significant age-related
118 changes in the allele frequencies of either oncogenic or non-oncogenic mutations within
119 proto-oncogenes (Fig. 2f and Extended Data Figure 7). This inability to observe any
120 clonal expansions with age is most likely due to the fact that the average age of the
121 individuals within the cohort is 49 years, with only 5 donors older than 70 years.

122 To explore the ability of FERMI to distinguish perturbations of somatic mutation
123 patterns, gDNA from mismatch repair deficient HCT116 cells (MMR^{MT}; hemizygous for
124 MLH1) was compared to MMR proficient parental cell line gDNA. Substantiating our

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

141 As expected from previous studies¹⁹, the HCT116 MMR^{MT} gDNA showed an
142 increased prevalence of T>C and T>A substitutions when compared to parental gDNA
143 (Extended Data Figure 8). The samples from individuals #2 and #19 also exhibited
144 these increased rates of T>C and T>A substitutions, with less extensive increases at C
145 positions, compared with the average of the 22 individuals (Fig. 3h-i and Extended Data
146 Figure 9), mirroring the changes observed in MMR^{MT} HCT116 cells. Thus, these two

147 individuals appear to present with a mild MMR-like substitution pattern. In support of the
148 results, individuals #2 and #19 show the same increased rates of substitutions across
149 multiple experiments (Extended Data Figure 9h-i). Of note, the systematic variance from
150 the typical mutational pattern for these two individuals and the MMR^{MT} HCT116 cells
151 serves as validation of the specificity of FERMI to accurately detect variants. More
152 importantly, this finding of two individuals with deviating mutational patterns out of a
153 sample size of only 22 individuals may be indicative of a broad spectrum of mutational
154 profiles that exist across the human population.

155 **Conclusion**

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

165 **Acknowledgments**

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

174 **Figure 1 | Amplicon sequencing accurately detects mutation allele frequencies as**
175 **rare as 1/10,000. a,** Graphical depiction of gDNA capture and analysis method. **b,**
176 Capture efficiencies vary in a probe dependent manner. **c,** Accurate detection of a
177 single heterozygous SNP in gDNA from one individual diluted into gDNA from another
178 (without this germline SNP) to frequencies as low as 1/10,000. **d,** Accurate detection of
179 three linked SNPs found within the same allele diluted as in c. For c and d, error shown
180 is standard deviation. **e,** ddPCR showing detection of R882H chr2:25457242 C>T
181 mutation at approximately the expected frequencies in normal human blood.

182 **Figure 2 | Mutations exist at conserved frequencies independently of age. a,**
183 Comparison of VAFs of identified variants within a 34 year old (x-axis) and 62 year old
184 (y-axis); $R^2 = 0.408211$, $p=0.000$. **b,** VAFs from a 34 year old (x-axis) compared to mean
185 VAFs from individuals ranging in ages from newborn to 89 years of age ($n=22$); $R^2 =$

186 0.590412, p=0.000. **c**, Relative contribution rates of each base substitution to all
187 substitutions identified. **d**, Relative contribution rates of each base substitution
188 segregated by surrounding 5' and 3' nucleotide context. **e**, All identified base
189 substitutions within a probed region are plotted by their position and VAFs for individuals
190 7 and 15 (representative of all other individuals, with greater deviation observed for
191 individuals 2 and 19 as described below), revealing highly reproducible patterns. **f**,
192 Oncogenic VAFs plotted as a function of donor age does not reveal evidence of clonal
193 expansions.

194 **Figure 3 | Individuals Can Systematically Deviate from the Population Average. a**,
195 Comparing VAFs in HCT116 MMR+ vs MMR^{MT} cells reveals an increase in frequencies
196 for many of the observed variants in MMR^{MT} cells ($R^2 = 0.211479$). **b**, MMR^{MT} vs mean
197 VAFs from blood of the 22 individuals shows a similar pattern of increased VAFs as the
198 comparison with parental HCT116 cells ($R^2 = 0.120895$). **c**, blood from a 73 yr old
199 person (individual #19) compared to the mean VAFs reveals a deviating population of
200 variants that exist at an increased frequency compared with average VAFs ($R^2 =$
201 0.387125). **d**, A cord blood sample (individual #2) also shows a subset of variants with
202 higher frequencies than in the average ($R^2 = 0.278250$). **e**, VAFs from individual #2 vs
203 individual #19 reveals that the deviating variants are at the same positions, causing the
204 comparison to fall close to the y=x line ($R^2 = 0.613542$). **f**, Plotting the mean for VAFs
205 from individuals #2 and #19 versus VAFs from MMR^{MT} HCT116 cells reveals that the
206 variants within the blood are the same as those found within the MMR^{MT} cell line. While

207 variant frequencies are higher in the MMR^{MT} cell line, the proportional change for
208 different deviating variants are similar ($R^2 = 0.587474$). **g**, Variants detected in
209 individuals #2 and #19 are not enriched for oncogenic changes, indicated in blue. **h**,
210 Plot of only C>N/G>N variants shows relative similarity between individual #2 and the
211 average for all other individuals ($R^2 = 0.350623$). **i**, Plot of only T>N/A>N variants
212 reveals that the majority of deviating variants for individual #2 are substitutions affecting
213 T or A (R -Squared = 0.040712).

214 Methods

215 Amplicon Design

216 Amplicon probes for targeted annealing regions were created using the Illumina
217 Custom Amplicon DesignStudio (<https://designstudio.illumina.com/>). UMIs were then
218 added to the designed probe regions and generated by IDT using machine mixing for
219 the randomized DNA. Probes were PAGE purified by IDT. All probes are listed below
220 along with binding locations and expected lengths of captured sequence.

Gene	Probe Up	Probe Down	Probe Start	Probe End	Length
JAK2	AGTTTACACTGACA CCTAGCTGTGATC	CCATAATTAAAACC AAATGCTTGAGA 232 A	chr9:5073733	chr9:5073887	155
TP53-1	TCATCTGGGCCTG TGTTATCTCTTA	ATCCTCACCATCAT CACACTGGAAAGAC	chr17:7577504	chr17:7577635	132
TP53-2	CCCTCAACAAAGATG TTTGGCCAACTG	ATGAGCGCTGCTCA GATAGCGATGGT	chr17:7578369	chr17:7578544	176
TP53-3	GGACAGGTAGGAC CTGATTCCTTACT	TGTCCTGGGAGAGA CCGGCGCACAGA	chr17:7577084	chr17:7577214	131
NRAS-1	CAATAGCATTGCAT TCCCTGTGGTTT	GTACAGTGCCATGA GAGACCAATACAT	chr1:115256496	chr1:115256680	185
NRAS-2	GAAGGTCACACTAG GGTTTCATTCC	AAAAGCGCACTGAC AATCCAGCTA	chr1:115258713	chr1:115258897	185
HRAS	TCCTTGGCAGGTGG	GCAAGAGTGCCTG	chr11:534258	chr1:534385	128

	GGCAGGAGACCC	ACCATCCA			
KRAS-1	AGGTACTGGTGGAG TATTGTAGTGT	CAAGAGTGCCTG CGATACTAGCTAATT	chr12:25398247	chr12:25398415	169
KRAS-2	GACTGTGTTCTCC CTTCTCAGGATT	TACAGTGCAATGAG GGACCAGTACATG	chr12:25380242	chr12:25380368	127
TET2-1	CCATGTTGGCTC ATTATGCTCTTA	ACGCCCACTCCCC AATGTCAG	chr4:106197237	chr4:106197405	169
TET2-2	CTTTGAAAGAGTG CCACTTGGTGTCT	GGTGATGGTATCAG GAATGGACTTAGTC	chr4:106155137	chr4:106155275	139
DNMT3A	TGTGTGGTTAGACG GCTTCCGGGCA	AGGCAGAGACTGCT GGGCCGGTCA	chr2:25457211	chr2:25457364	154
IDH1	CAAATGTGAAATC ACCAAATGGCACC	TGGGGATCAAGTAA GTCATGTTGGCA	chr2:209113077	chr2:209113239	163
IDH2	GAAGAAGATGTGGA AAAGTCCAATGG	CATGGCGACCAGGT AGGCCAGG	chr15:90631809	chr15:90631969	161
GATA1	CTTCCAGCCATTTC TGAGATATCCTCA	CAGCTGCAGCGGT GGCTGTGCT	chrX:48649667	chrX:48649849	183
SF3B1	GTGAACATATTCTG CAGTTGGCTGAA	ACCATCAGTGCTTT GGCCATTGC	chr2:198266803	chr2:198266967	165
TIIIA	CATCTATTCTGTGCT AGGCATTGTGTG	CAGACCTAGCATCT GTGCCAGAC	chr1:115227814	chr1:115227978	165
TIIB	CAGTCTGGTTTTG GAGCAATGATATC	GCAGTGAGCTCAGC CTTGATTT	chr2:223190674	chr2:223190820	147
TIIIC	CCTGGTGCTTAGTC CTGTTCTGAAATT	AGTCTTCTATAATGC CACAAACCTGTAT	chr2:229041101	chr2:229041289	189
TIIID	GAACAGAACACTTG GTAGTTGACCATG	AGACAGGGAACTGG CATGAAGAGTTT	chr4:110541172	chr4:110541302	131
TIIIE	GCCTAGAACAGGCA CCATACATTCAAT	AGATGGTGTGCTG TGCCGGATAGGAG	chr4:112997214	chr4:112997386	173
TIIIF	TGGCACTATGTGGA GATGTTAGTACAG	GGATGTTGGTGCTA TCAGTAGCCATA	chr4:121167756	chr4:121167884	129
TIIIG	CTCTAGGCTTAGTG GTCAAGGAATGAA	AGAACGAGGACTGT GCTTCCAAACAA	chr4:123547743	chr4:123547901	159
TIIIH	CTTGGTGGTAGCCT AGGCAGTAATTAA	CACGTGGTTGGGAA GAGAAAGTG	chr4:124428637	chr4:124428767	131
TIIIJ	TTCTATAGCACTGG TGACCAGGACACT	CTGGCCACAGTGCC TGGTTTCC	chr11:2126256	chr11:2126420	165
TIIIK	AGACAGGAGGAAG GAGCAATTAGAAG	CATGGAGATCTCGT CCCCTCAGA	chr11:2389983	chr11:2390171	189
TIIIL	TAGGCCAGAAAACA CACAGTGTGCGGG	AACTCCGGTAAGTG CGGGGTGGGGGT	chr11:2593889	chr11:2594074	186
TIIIM	ATCTGGAACAGAC CTTCTCAGGCAT	GTTCTAAGTTACTCT GTGTACTTGACT	chr11:11486596	chr11:11486728	133
TIIIN	AGCCTAGTTACCAT AGACGGATTCAAC	GAATATCTTCTAACT GGACTTAGAAAACC	chr15:92527052	chr15:92527176	125
TIIIO	CCAACATGTTCTAA ATTCTGGCACAG	TGGGTCTCAGCCAT CCCATTACTG	chr16:73379656	chr16:73379832	177
TIIIP	CTAACATCTCACTTC TACCTCTACGCTA	TAAGTGCCTCACTAC CCCATCCTTAAT	chr16:82455026	chr16:82455164	139
TIIIQ	TCATGACCCAGGCC TCCCGAGAACTGAG	ATCTGTGAAGCCGG AGTGAAAACAAC	chr16:85949137	chr16:85949299	163

484 **Genomic DNA Isolation**

485 Human blood samples were purchased from the Bonfils Blood Center
486 Headquarters of Denver Colorado. Our use of these samples was determined to be “Not
487 Human Subjects” by our Institutional Review Board. Biopsies were collected as
488 unfractionated whole blood from apparently healthy donors, though samples were not
489 tested for infection. Samples were approximately 10 mL in volume, and collected in BD
490 Vacutainer spray-coated EDTA tubes. Following collection, samples were stored at 4°C
491 until processing, which occurred within 5 hours of donation. To remove plasma from the
492 blood, samples were put in 50 mL conical tubes (Corning #430828) and centrifuged for
493 10 minutes at 515 rcf. Following centrifugation, plasma was aspirated and 200 mL of
494 4°C hemolytic buffer (8.3g NH₄Cl, 1.0g NaHCO₃, 0.04 Na₂ in 1L ddH₂O) was added to
495 the samples and incubated at 4°C for 10 minutes. Hemolyzed cells were centrifuged at
496 515 rcf for 10 minutes, supernatant was aspirated, and pellet was washed with 200 mL
497 of 4°C PBS. Washed cells were centrifuged for at 515rcf for 10 minutes, from which
498 gDNA was extracted using a DNeasy Blood & Tissue Kit (Qiagen REF 69504).

499 **Amplicon Capture**

500 For amplicon capture from gDNA, we modified the Illumina protocol called
501 “Preparing Libraries for Sequencing on the MiSeq” (Illumina Part #15039740 Revision
502 D). DNA was quantified with a NanoDrop 2000c (ThermoFisher Catalog #ND-2000C).
503 500ng of input DNA in 15μl was used for each reaction instead of the recommended
504 quantities. In place of 5μl of Illumina ‘CAT’ amplicons, 5μl of 4500ng/μl of our amplicons

505 were used. During the hybridization reaction, after gDNA and amplicon reaction mixture
506 was prepared, sealed, and centrifuged as instructed, gDNA was melted for 10 minutes
507 at 95°C in a heat block (SciGene Hybex Microsample Incubator Catalog #1057-30-O).
508 Heat block temperature was then set to 60°C, allowed to passively cool from 95°C and
509 incubated for 24hr. Following incubation, the heat block was set to 40°C and allowed to
510 passively cool for 1hr. The extension-ligation reaction was prepared using 90 µl of ELM4
511 master mix per sample and incubated at 37°C for 24hr. PCR amplification was
512 performed at recommended temperatures and times for 29 cycles. Successful
513 amplification was confirmed immediately following PCR amplification using a
514 Bioanalyzer (Agilent Genomics 2200 Tapestation Catalog #G2964-90002, High
515 Sensitivity D1000 ScreenTape Catalog #5067-5584, High Sensitivity D1000 Reagents
516 Catalog #5067-5585). PCR cleanup was then performed as described in Illumina's
517 protocol using 45 µl of AMPure XP beads. Libraries were then normalized for
518 sequencing using the Illumina KapaBiosystems qPCR kit (KapaBiosystems Reference #
519 07960336001).

520 **Sequencing**

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

524 **Bioinformatics**

525 The analysis pipeline used to process sequencing results can be found under
526 FERMI here: <http://software.laliggett.com/>. For a detailed understanding of each
527 function provided by the analysis pipeline, refer directly to the software. The overall goal
528 of the software built for this project is to analyze amplicon captured DNA that is tagged
529 with equal length UMIs on the 5' and 3' ends of captures, and has been paired-end
530 sequenced using dual indexes. Input fastq files are either automatically or manually
531 combined with their paired-end sequencing partners into a single fastq file. Paired reads
532 are combined by eliminating any base that does not match between Read1 and Read2,
533 and concatenating this consensus read with the 5' and 3' UMIs. A barcode is then
534 created for each consensus read from the 5' and 3' UMIs and the first five bases at the
535 5' end of the consensus. All consensus sequences are then binned together by their
536 unique barcodes. The threshold for barcode mismatch can be specified when running
537 the software, and for all data shown in this manuscript one mismatched base was
538 allowed for a sequence to still count as the same barcode. Bins are then collapsed into
539 a single consensus read by first removing the 5' and 3' UMIs. Following UMI removal,
540 consensus sequences are derived by incorporating the most commonly observed
541 nucleotide at each position so long as the same nucleotide is observed in at least a
542 specified percent of supporting reads (55% of reads was used for results in this
543 manuscript) and there are least some minimum number of reads supporting a capture
544 (5 supporting reads was used for results in this manuscript). Any nucleotide that does
545 not meet the minimum threshold for read support is not added to the consensus read,
546 and alignment is attempted with an unknown base at that position. From this set of

547 consensus reads, experimental quality measurements are made, such as total captures,
548 total sequencing reads, average capture coverage, and estimated error rates. Derived
549 consensus reads are then aligned to the specified reference genome using
550 Burrows-Wheeler²⁰, and indexed using SAMtools²¹. For this manuscript consensus
551 reads were aligned to the human reference genome hg19^{22,23} (though the software
552 should be compatible with other reference genomes). Sequencing alignments are then
553 used to call variants using the Bayesian haplotype-based variant detector, FreeBayes²⁴.
554 Identified variants are then decomposed and block decomposed using the variant
555 toolset vt²⁵. Variants are then filtered to eliminate any that have been identified outside
556 of probed genomic regions. If necessary variants can also be eliminated if below certain
557 coverage or observation thresholds such that variants must be independently observed
558 multiple times in different captures to be included. For this manuscript, we included all
559 variants that passed previous filters and did not eliminate those that were observed only
560 within a single capture, unless otherwise indicated.

561 **Elimination of false positive signal**

562 A number of steps have been included within sample preparation and
563 bioinformatics analysis specifically to distinguish between true positive signal and false
564 positive signal. Using the dilution series shown in Figs. 1c-d we can show sufficient
565 sensitivity to identify signal diluted to levels as rare as 10^{-4} . While these dilutions show
566 significantly improved sensitivity over many current sequencing methods, they do not
567 address our background error rate. Unfortunately, because both endogenous and

568 exogenous DNA synthesis is error prone, it is challenging to find negative controls that
569 can be used to estimate background error rates with a method of mutation detection as
570 putatively sensitive as FERMI. Nevertheless, we have a number of steps that should
571 eliminate most sources of false signal. The two largest sources of erroneous mutation
572 when sequencing DNA will typically be from PCR amplification mutations (caused both
573 by polymerase errors and exogenous insults like oxidative damage), and sequencing
574 errors.

575 The steps are the following:

- 576 ● *Elimination of first round PCR amplification errors*
577 ● *Elimination of subsequent PCR amplification errors*
578 ● *Elimination of sequencing errors*

579 *Elimination of first round PCR amplification errors*

580 The first round of PCR amplification performed during library preparation causes
581 mutations that are challenging to distinguish from those that occurred endogenously.
582 Since there is little difference between those mutations that occur during the first round
583 of PCR amplification and those that occurred endogenously, we rely on probability to
584 eliminate these errors. Since we are performing single-cell sequencing, we can require
585 that a mutation be observed in multiple cells before it is called as a true positive signal.
586 As we expect about 400 first round PCR amplification errors, the probability that the
587 identical mutation will occur in multiple cells becomes exponentially unlikely (Extended

588 Data Figure 1). By requiring a mutation be observed in just three cells before it is called
589 as real signal, only about 1-2 first round PCR amplification errors should ever make it
590 into the final data. When we process our data requiring up to 5 independent
591 observations of a mutation, the overall mutation spectrum does not change, apart from
592 a loss of the most rarely observed variants. This observation led us to include all
593 variants that were observed even once. Our logic is that while about 400 variants will be
594 the result of the first round of PCR amplification, these same variants are already
595 occurring endogenously, meaning that absolute variant allele frequency accuracy will
596 be affected, but not the identities of variants.

597 *Elimination of subsequent PCR amplification errors*

598 Elimination of PCR amplification errors after the first round of PCR is done using
599 UMI collapsing (Fig. 1a). Each time a strand is amplified, the UMI will keep track of its
600 identity. Any mutations that occur after the first round of PCR will be found in 25% of the
601 reads or fewer. This allows us to collapse each unique capture and eliminate any rarely
602 observed variants associated with a given UMI. Utilizing the UMI in this way allows us to
603 essentially eliminate any PCR amplification errors that occurred after the first round of
604 PCR.

605 *Elimination of sequencing errors*

606 Sequencing errors are eliminated in two ways. This first method is by using
607 paired-end sequencing to sequence the same fragment of DNA twice (Fig. 1a). The

608 information found within each of these reads (Read1 and Read2) should theoretically be
609 the same, meaning that unless the same sequencing error is made at the same locus
610 within both Read1 and Read2, the two strands will differ. Any differences are simply
611 eliminated from the data as it is unknown which base call is correct, and the rest of the
612 data is included in the analysis. This collapsing should eliminate most sequencing
613 errors, though it will sequencing errors of the same identity occurring at the same locus.
614 These errors are removed when collapsing into single cell bins (Fig. 1a). As with the
615 logic when eliminating subsequent PCR amplification errors, all sequence associated
616 with each UMI should be identical. Therefore, sequencing errors passing through Read1
617 and Read2 will not match other sequenced strands from the same capture event, and
618 are eliminated during consensus sequence derivation.

619 **Mutation signature analysis**

620 20 somatic mutation signatures were previously identified¹⁵ by analyzing
621 trinucleotide mutation context of cancer genomes using non-negative matrix
622 factorization (NMF) and principal component analysis (PCA). Here, we used
623 deconstructSig²⁶ to identify the relative presence of those mutation signatures within the
624 somatic mutations detected blood using somaticSignatures²⁷. Codon triplet biases were
625 analyzed using the MutationalPatterns R package²⁸.

626 **Estimation of mutation rate**

627 It is difficult to understand the somatic lineage development that gave rise to the
628 number of cells that are assayed from each blood biopsy. So estimating a somatic
629 mutation rate is challenging. Nevertheless we can estimate what somatic mutation rates
630 might be assuming the number of cells that were assayed contained all unique variants
631 or not.

632 An upper bound for the somatic mutation rate observed by FERMI analysis can
633 be estimated by using the number of captures and total observed variants, and assume
634 that all of these are de-novo mutations. In our data we observe about 1232458 unique
635 captures per analyzed blood sample. These captures are relatively uniformly spread
636 across each of our 32 different probes which span a total of 4838bp. From this the total
637 probed DNA, D_T , can be estimated as:

638
$$D_T = \frac{1232458 \text{ captures} * 4838 \text{ bp}}{32 \text{ probes}}$$

639
$$D_T = 186332243.9 \text{ bp}$$

640 The total number of observed observed variants within each blood sample is
641 approximately 168940, from which the aggregate mutation rate, M, can be estimated as:

642
$$M = \frac{168940 \text{ mutations}}{186332243.9 \text{ bp}}$$

643
$$M = 9 * 10^{-4} \text{ mut/bp}$$

644
$$M = 900 \text{ mut/Mb}$$

645 A lower estimate can be made by assuming that mutations are not all unique
646 occurrences but might be the result of clonal expansions creating many copies of each
647 mutation. This mutation rate, M, can be roughly estimated by the approximately 40000
648 captures per each of the 32 probes that captured roughly 6000 variants across a

649 conservative 100bp sized capture for each probe (probe region is realistically smaller
650 than 150bp because of collapsing conditions).

651

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

652

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

653

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

- 654 1. Benzer, S. ON THE TOPOGRAPHY OF THE GENETIC FINE STRUCTURE. *Proc.*
- 655 *Natl. Acad. Sci. U. S. A.* **47**, 403–415 (1961).
- 656 2. Gaffney, D. J. & Keightley, P. D. The scale of mutational variation in the murid
- 657 genome. *Genome Res.* **15**, 1086–1094 (2005).
- 658 3. Lercher, M. J., Williams, E. J. B. & Hurst, L. D. Local similarity in evolutionary rates
- 659 extends over whole chromosomes in human-rodent and mouse-rat comparisons:
- 660 implications for understanding the mechanistic basis of the male mutation bias. *Mol.*
- 661 *Biol. Evol.* **18**, 2032–2039 (2001).
- 662 4. Nachman, M. W. & Crowell, S. L. Estimate of the mutation rate per nucleotide in
- 663 humans. *Genetics* **156**, 297–304 (2000).
- 664 5. Hwang, D. G. & Green, P. Bayesian Markov chain Monte Carlo sequence analysis
- 665 reveals varying neutral substitution patterns in mammalian evolution. *Proc. Natl.*
- 666 *Acad. Sci. U. S. A.* **101**, 13994–14001 (2004).
- 667 6. Tomasetti, C., Li, L. & Vogelstein, B. Stem cell divisions, somatic mutations, cancer
- 668 etiology, and cancer prevention. *Science* **355**, 1330–1334 (2017).
- 669 7. Kujoth, G. C. *et al.* Mitochondrial DNA mutations, oxidative stress, and apoptosis in
- 670 mammalian aging. *Science* **309**, 481–484 (2005).
- 671 8. Hiatt, J. B., Pritchard, C. C., Salipante, S. J., O’Roak, B. J. & Shendure, J. Single
- 672 molecule molecular inversion probes for targeted, high-accuracy detection of
- 673 low-frequency variation. *Genome Res.* **23**, 843–854 (2013).
- 674 9. Preston, J. L. *et al.* High-specificity detection of rare alleles with Paired-End Low
- 675 Error Sequencing (PELE-Seq). *BMC Genomics* **17**, 464 (2016).

- 676 10. Zhang, T.-H., Wu, N. C. & Sun, R. A benchmark study on error-correction by
677 read-pairing and tag-clustering in amplicon-based deep sequencing. *BMC Genomics*
678 1–9 (2016).
- 679 11. Schmitt, M. W. *et al.* Sequencing small genomic targets with high efficiency and
680 extreme accuracy. *Nat. Methods* 1–4 (2015).
- 681 12. Kennedy, S. R. *et al.* Detecting ultralow-frequency mutations by Duplex Sequencing.
682 *Nat. Protoc.* **9**, 2586–2606 (2014).
- 683 13. Chen, L., Liu, P., Evans, T. C., Jr & Ettwiller, L. M. DNA damage is a pervasive cause
684 of sequencing errors, directly confounding variant identification. *Science* **355**,
685 752–756 (2017).
- 686 14. Blokzijl, F. *et al.* Tissue-specific mutation accumulation in human adult stem cells
687 during life. *Nature* **538**, 260–264 (2016).
- 688 15. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature*
689 **500**, 415–421 (2013).
- 690 16. Alexandrov, L. B. *et al.* Mutational signatures associated with tobacco smoking in
691 human cancer. *Science* **354**, 618–622 (2016).
- 692 17. Jaiswal, S. *et al.* Age-Related Clonal Hematopoiesis Associated with Adverse
693 Outcomes. *N. Engl. J. Med.* 1–11 (2014).
- 694 18. Martincorena, I. *et al.* Tumor evolution. High burden and pervasive positive selection
695 of somatic mutations in normal human skin. *Science* **348**, 880–886 (2015).
- 696 19. Zhao, H. *et al.* Mismatch repair deficiency endows tumors with a unique mutation
697 signature and sensitivity to DNA double-strand breaks. *eLife Sciences* **3**, e02725

- 698 (2014).
- 699 20. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler
700 transform. *Bioinformatics* **25**, 1754–1760 (2009).
- 701 21. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
702 2078–2079 (2009).
- 703 22. Lander, E. S. *et al.* Initial sequencing and analysis of the human genome. *Nature*
704 **409**, 860–921 (2001).
- 705 23. Fujita, P. A. *et al.* The UCSC genome browser database: update 2011. *Nucleic Acids
706 Res.* **39**, D876–D882 (2010).
- 707 24. Garrison, E. & Marth, G. Haplotype-based variant detection from short-read
708 sequencing. *arXiv [q-bio.GN]* (2012).
- 709 25. Tan, A., Abecasis, G. R. & Kang, H. M. Unified representation of genetic variants.
710 *Bioinformatics* **31**, 2202–2204 (2015).
- 711 26. Rosenthal, R., McGranahan, N., Herrero, J., Taylor, B. S. & Swanton, C.
712 DeconstructSigs: delineating mutational processes in single tumors distinguishes
713 DNA repair deficiencies and patterns of carcinoma evolution. *Genome Biol.* **17**, 31
714 (2016).
- 715 27. Gehring, J. S., Fischer, B., Lawrence, M. & Huber, W. SomaticSignatures: inferring
716 mutational signatures from single-nucleotide variants. *Bioinformatics* **31**, 3673–3675
717 (2015).
- 718 28. Blokzijl, F., Janssen, R., Van Boxtel, R. & Cuppen, E. MutationalPatterns: an
719 integrative R package for studying patterns in base substitution catalogues. *bioRxiv*

720

071761 (2016). doi:10.1101/071761

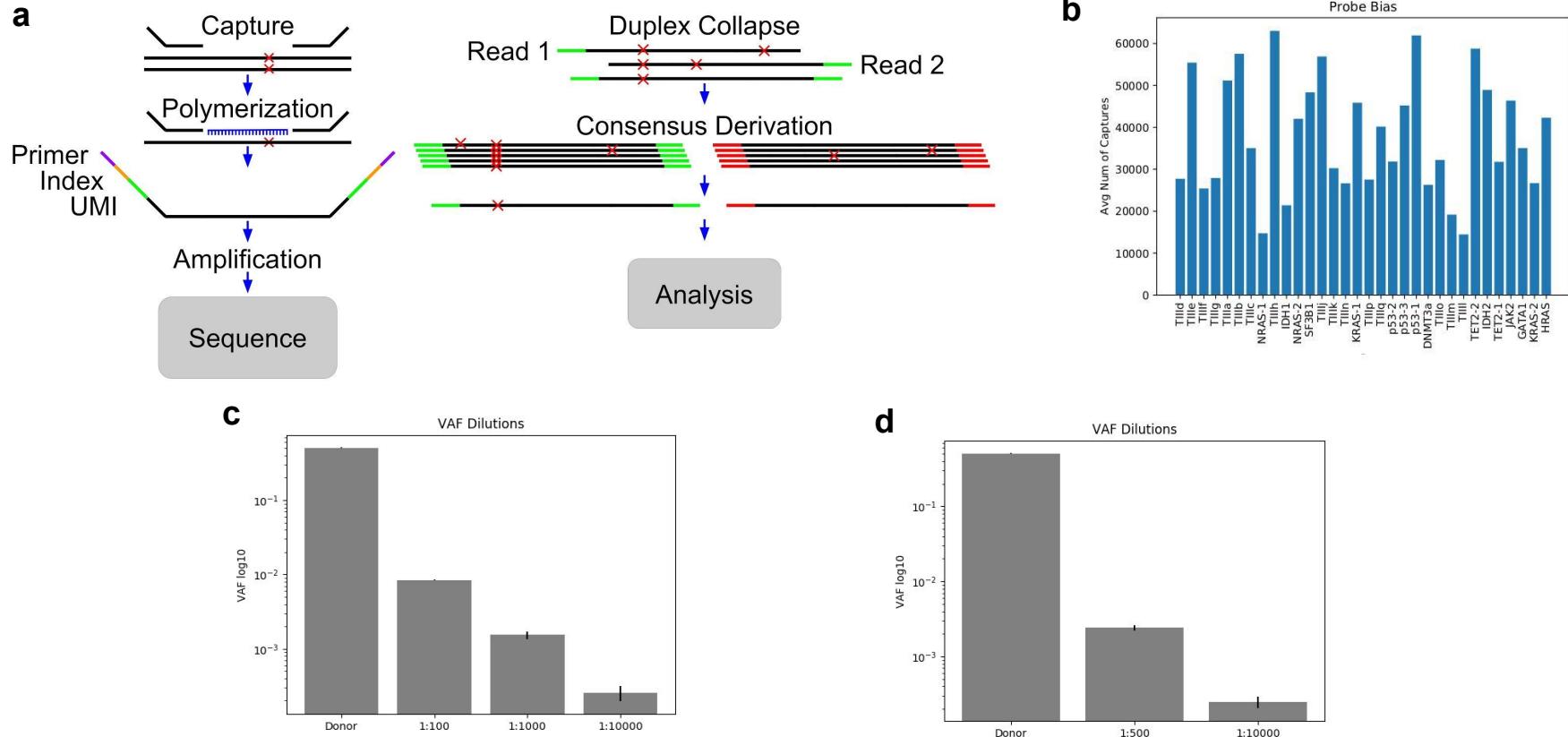


Figure 1 | Amplicon sequencing accurately detects mutation allele frequencies as rare as 1/10,000. **a**, Graphical depiction of gDNA capture and analysis method. **b**, Capture efficiencies vary in a probe dependent manner. **c**, Accurate detection of a single heterozygous SNP in gDNA from one individual 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.

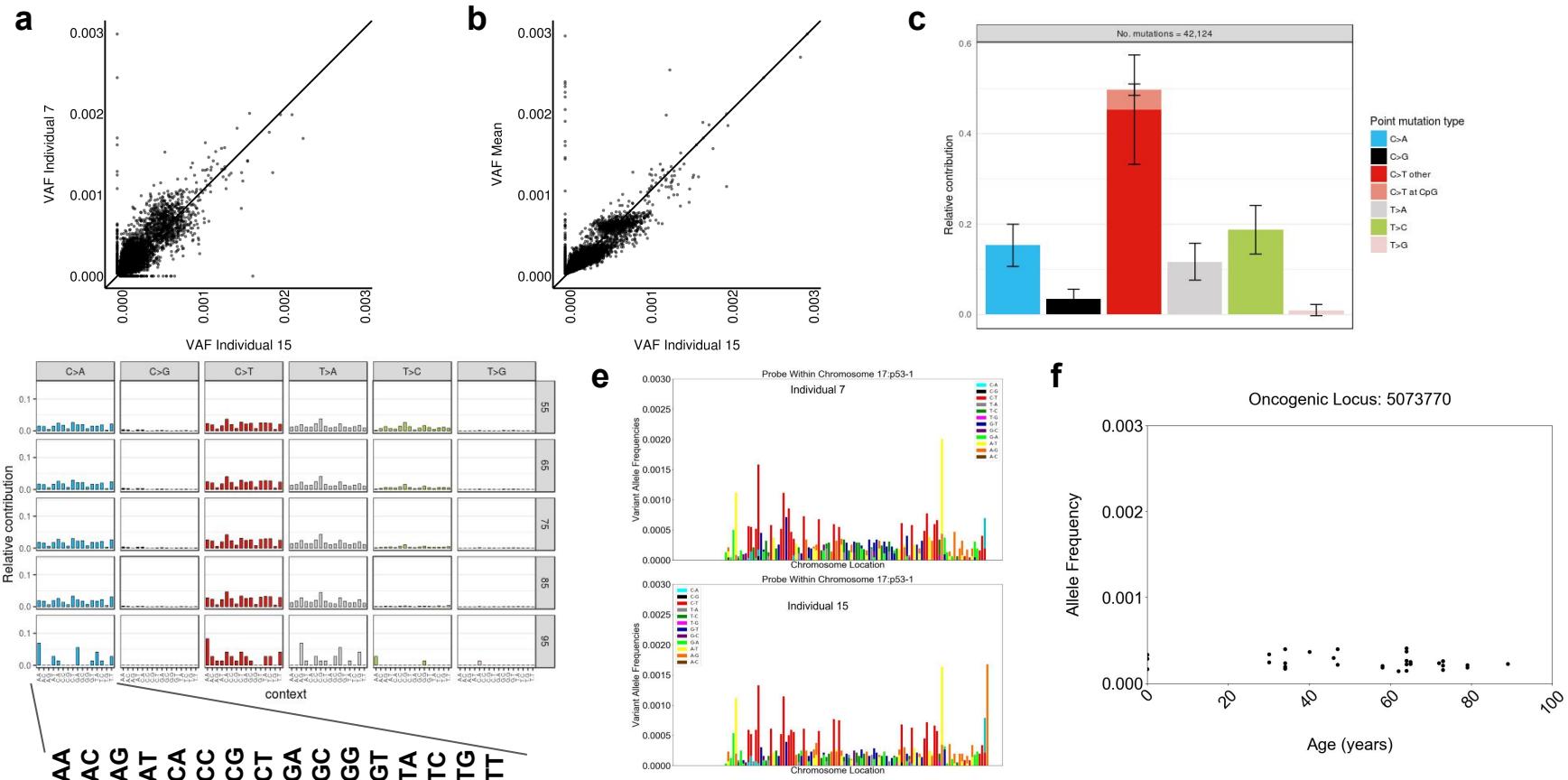


Figure 2 | Mutations exist at conserved frequencies independently of age. **a**, Comparison of VAFs of identified variants within a 34 year old (x-axis) and 62 year old (y-axis); $R^2 = 0.408211$, $p=0.000$. **b**, VAFs from a 34 year old (x-axis) compared to mean VAFs from individuals ranging in ages from newborn to 89 years of age ($n=22$); $R^2 = 0.590412$, $p=0.000$. **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.

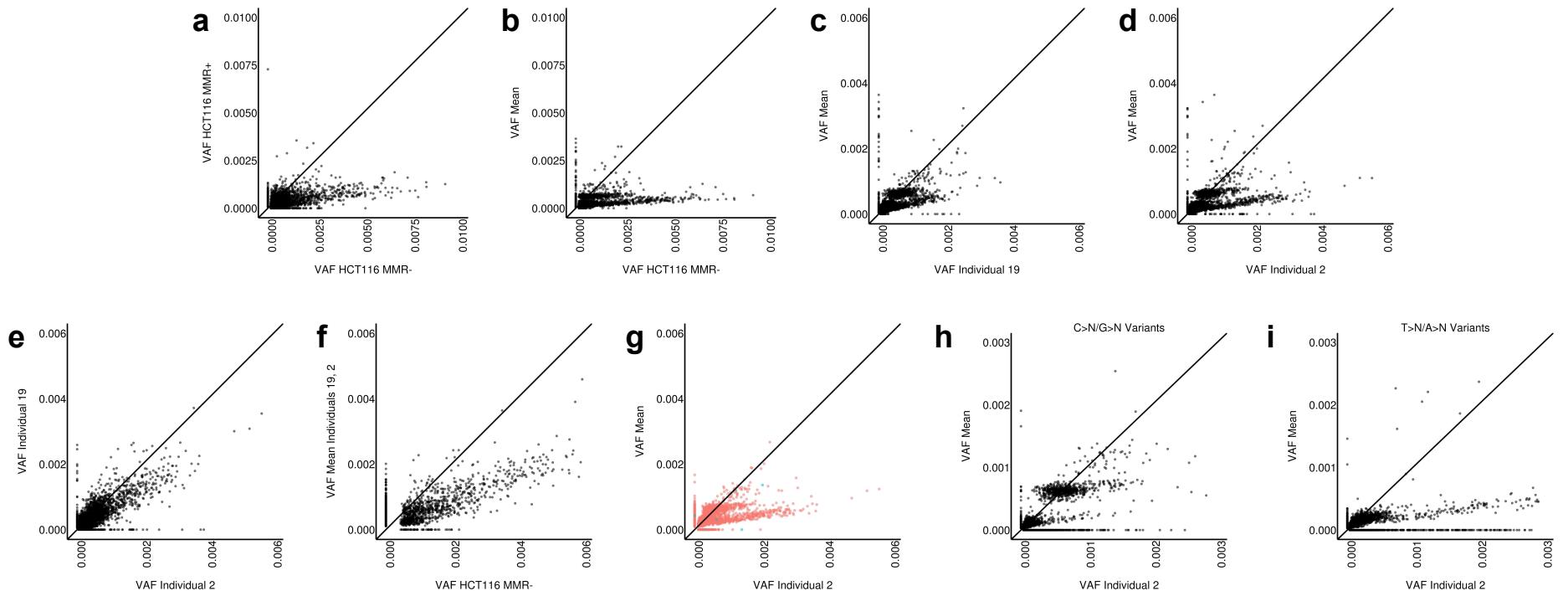


Figure 3 | Individuals Can Systematically Deviate from Population Average. **a**, Comparing VAFs in HCT116 MMR+ vs MMR^{MT} cells reveals an increase in frequencies for many of the observed variants in MMR^{MT} cells (R-Squared = 0.211479). **b**, MMR^{MT} vs mean VAFs from blood of the 22 individuals 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

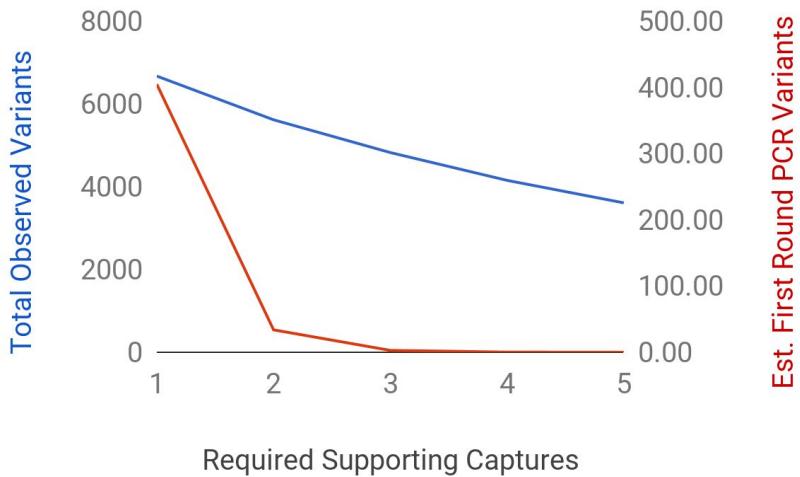
b

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

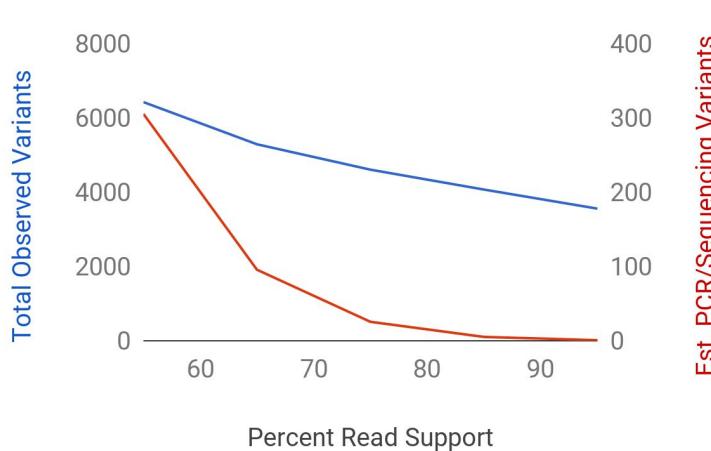
c

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

d



e



Extended Data Table 1

Cohort 1

a

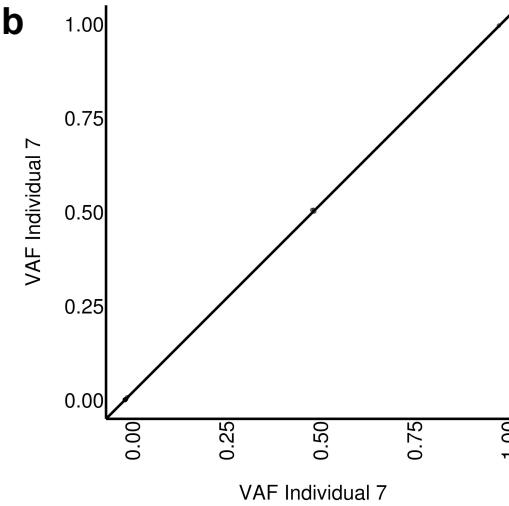
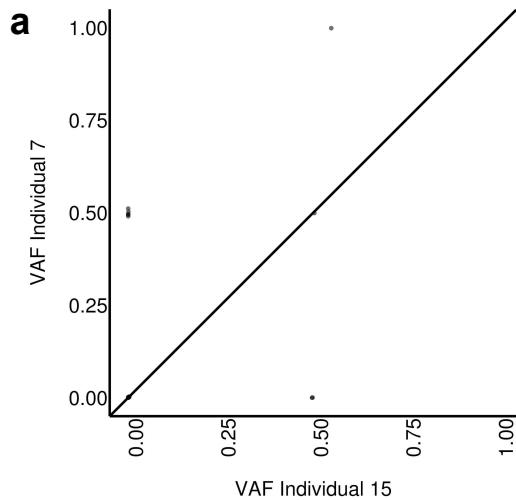
Individual	Age (years)
1	0
2	0
3	0
4	34
5	34
6	30
7	34
8	46
9	47
10	40
11	59
12	59
13	58
14	62
15	65
16	64
17	64
18	73
19	73
20	72
21	79
22	89

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

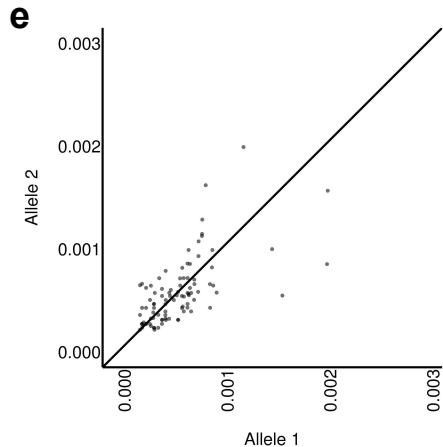
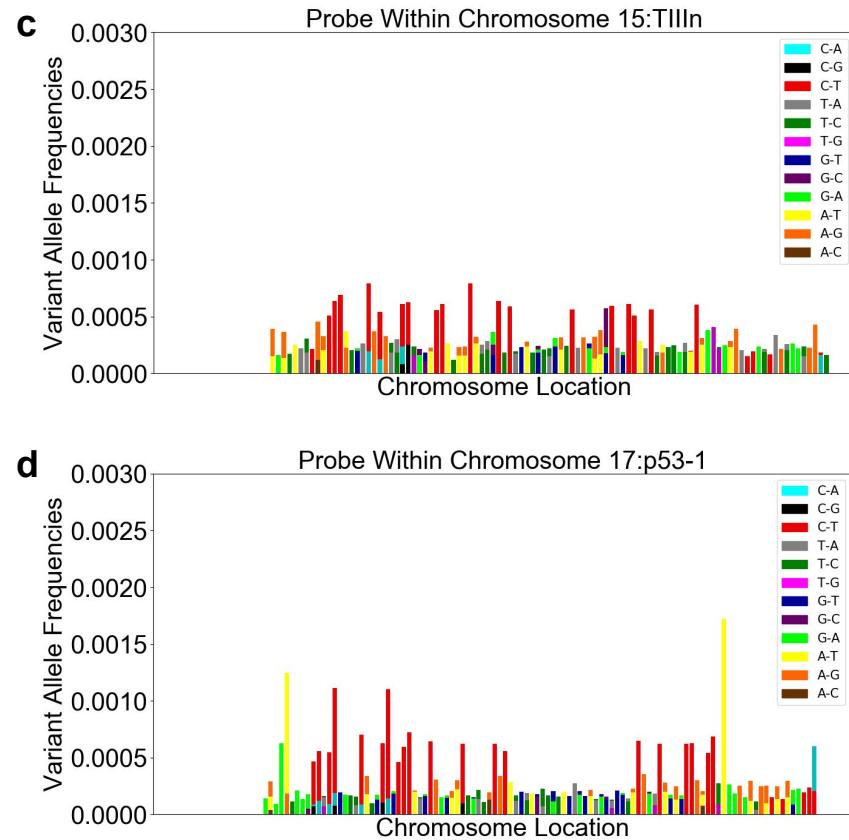
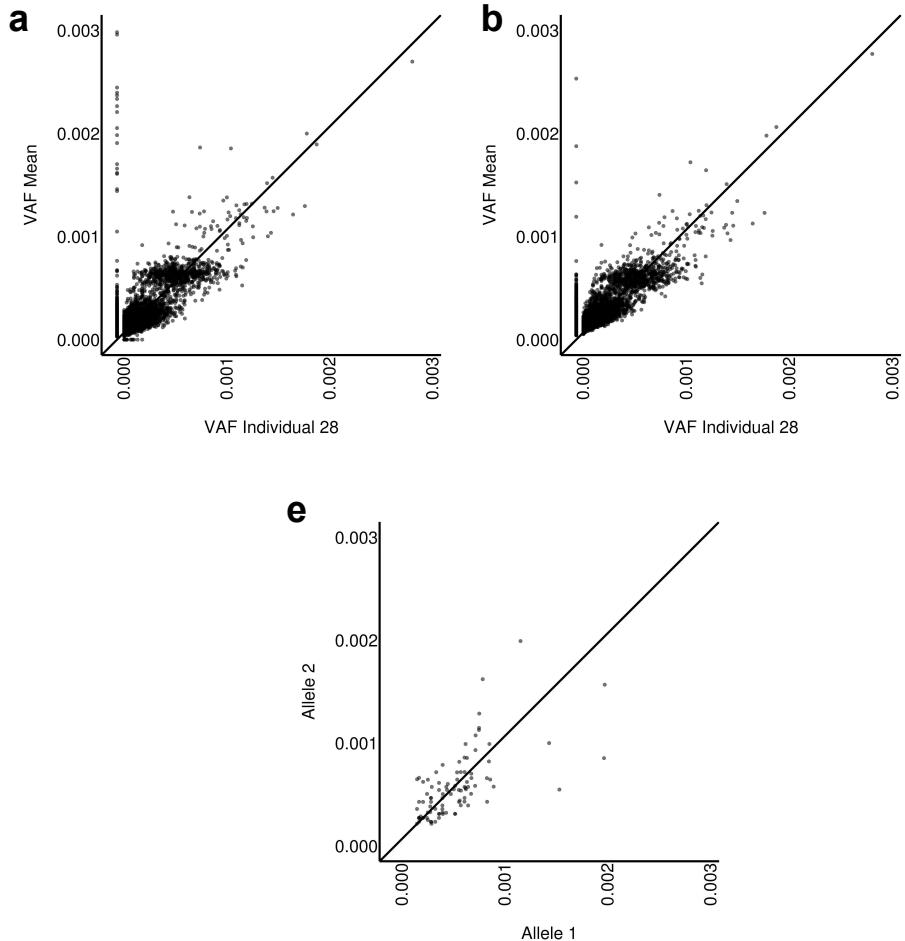
Extended Data Figure 2



c

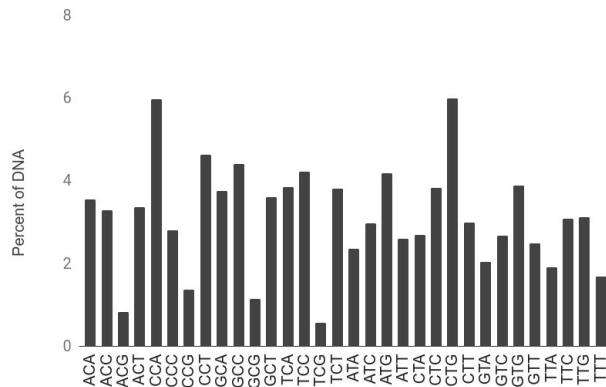
Individual	0mo vs 1mo
Individual A	0.460348
Individual B	0.538478
Individual C	0.436766
Individual D	0.522387
Individual E	0.519219
Individual F	0.482805

Extended Data Figure 3

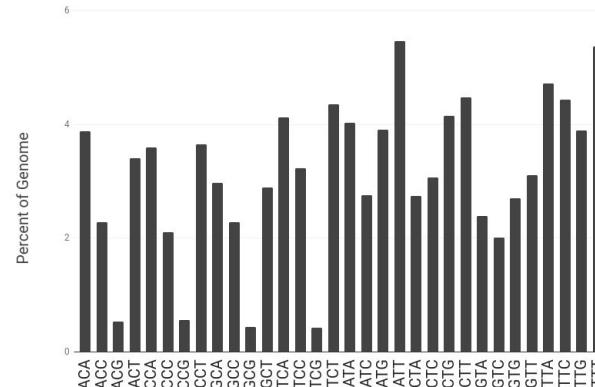


Extended Data Figure 4

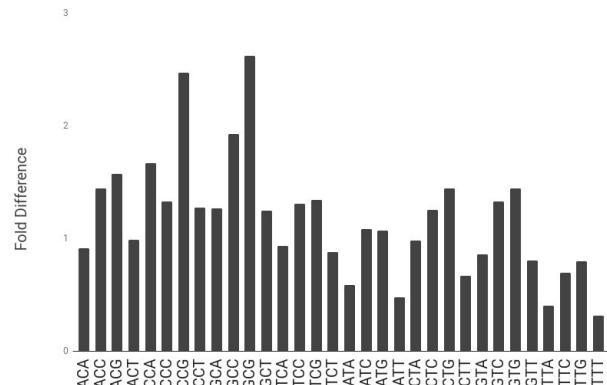
a Trinucleotide Representation Probed Region



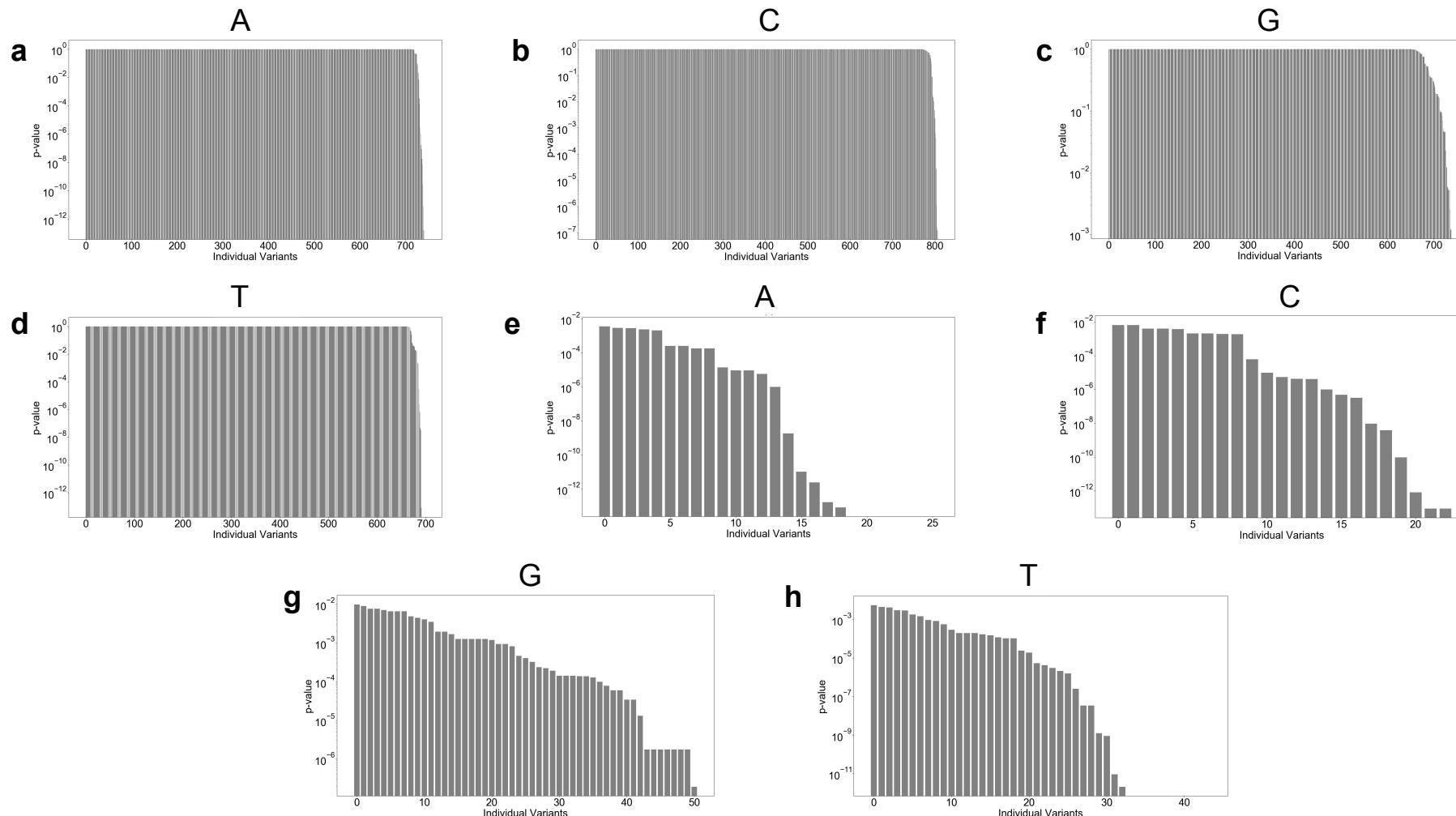
b Trinucleotide Representation Human Genome



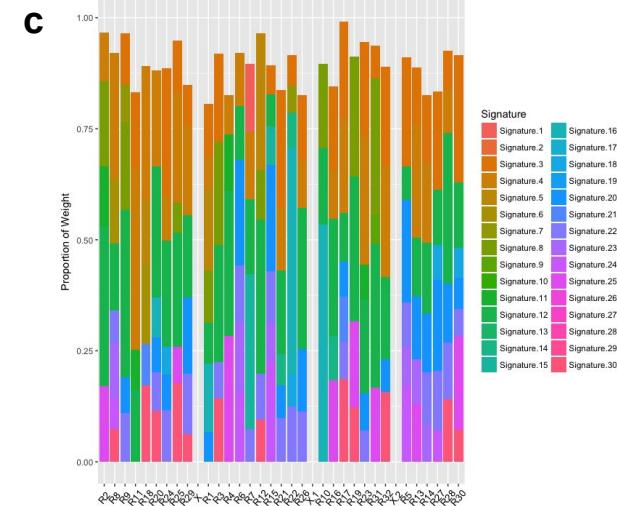
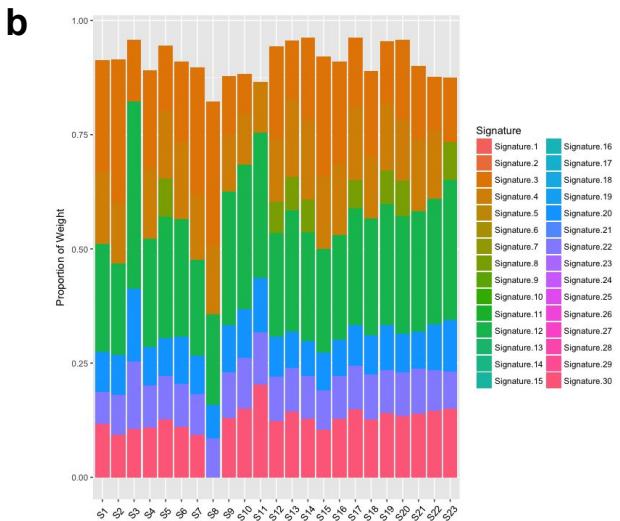
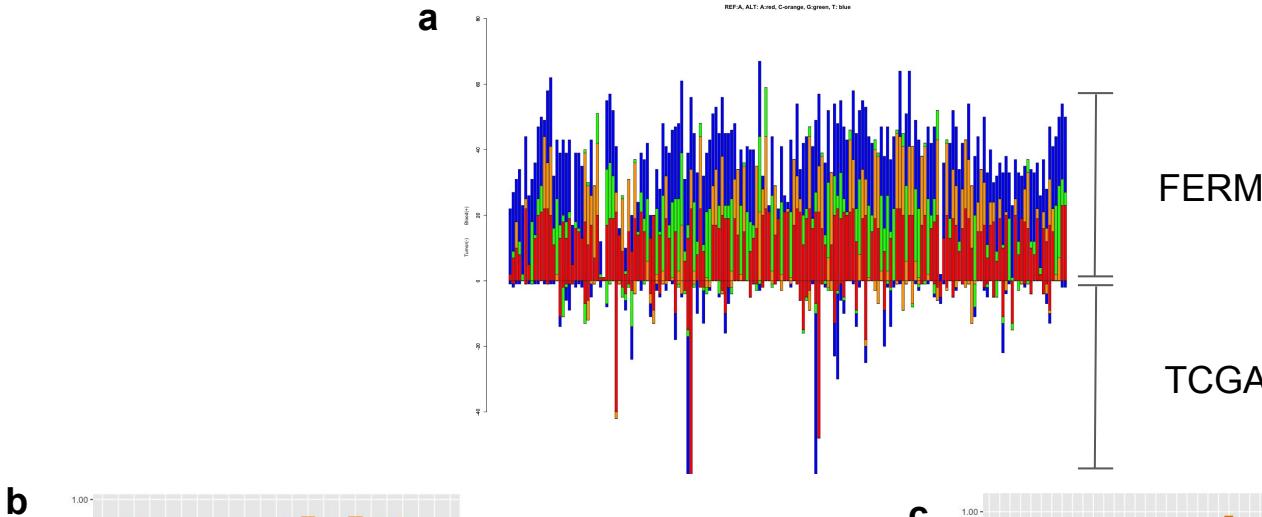
c Trinucleotide Representation (Probed Region/hg19)



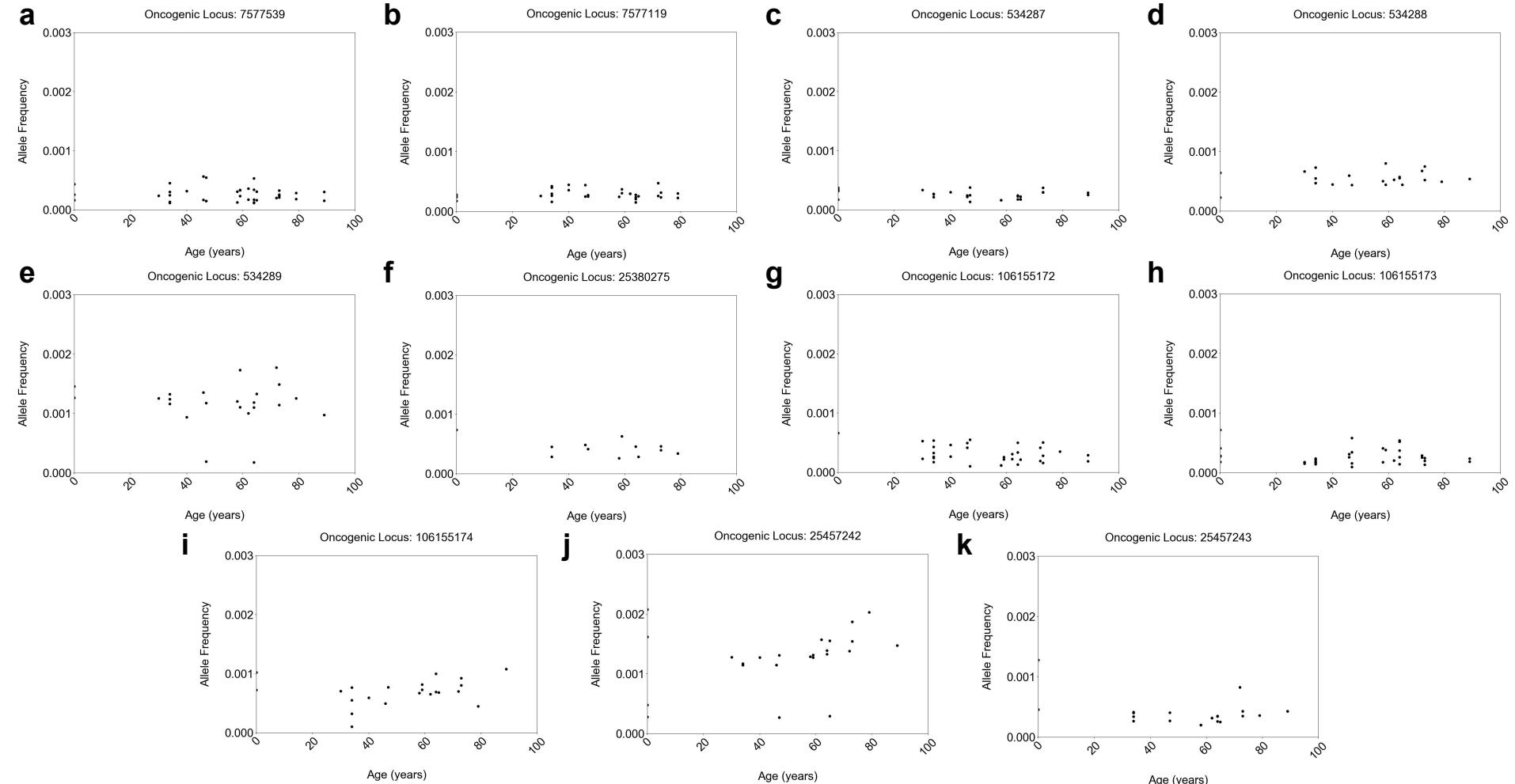
Extended Data Figure 5



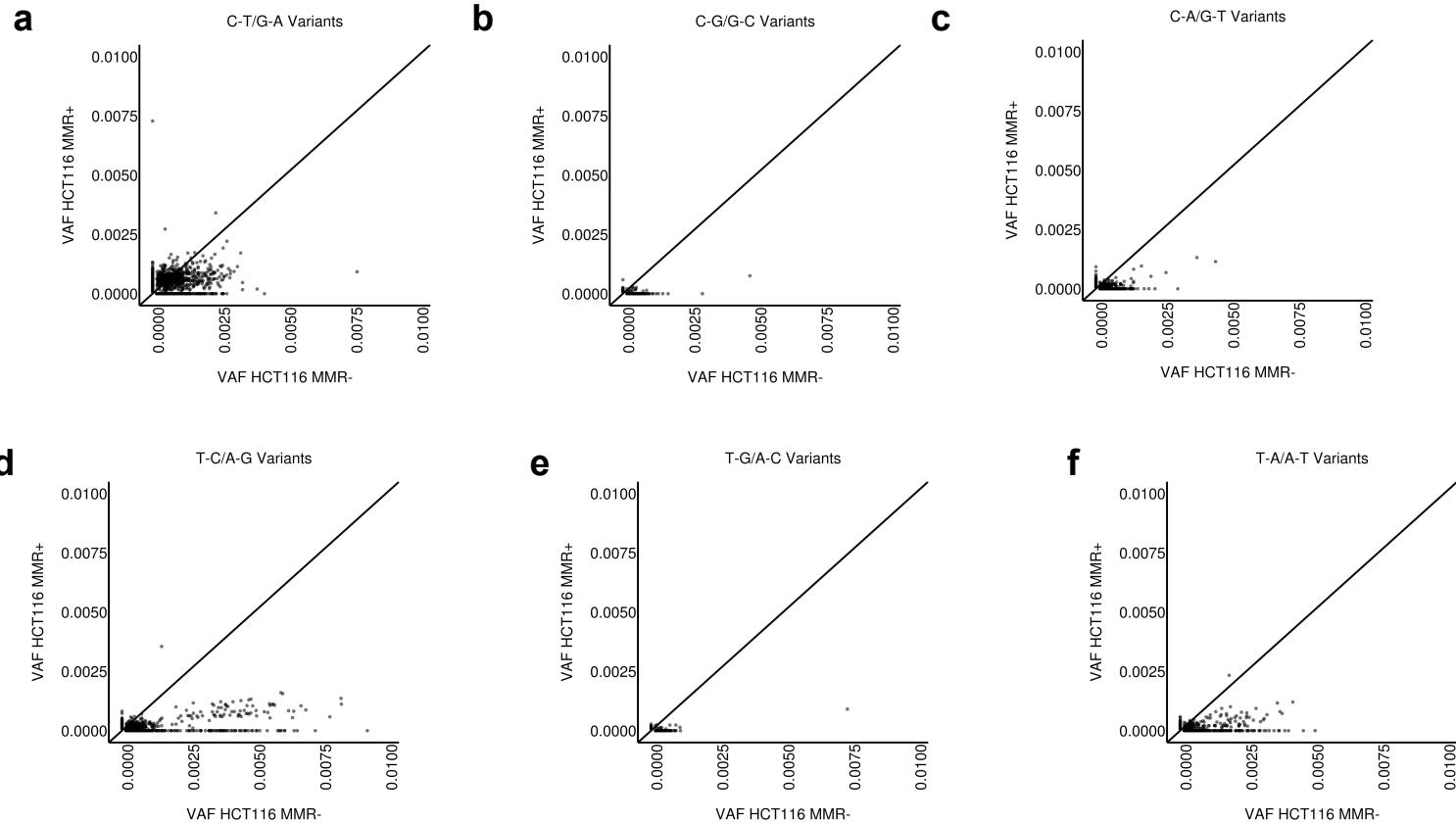
Extended Data Figure 6



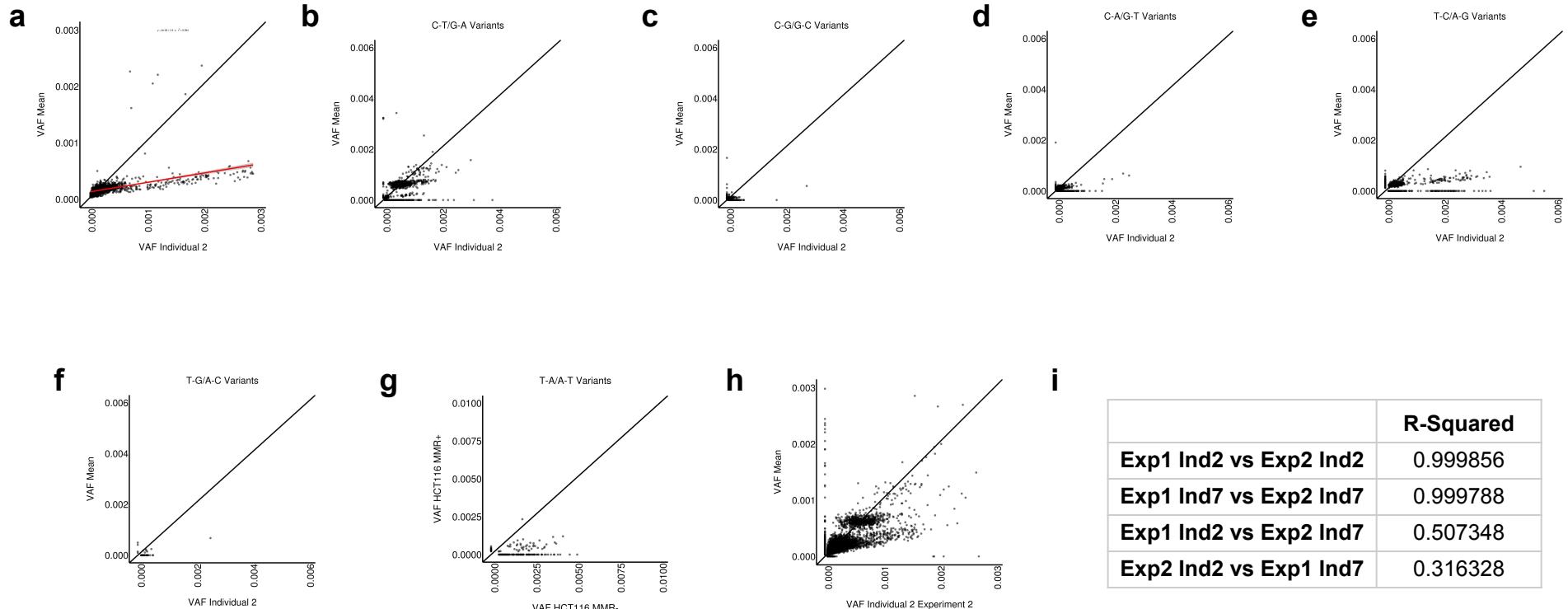
Extended Data Figure 7



Extended Data Figure 8



Extended Data Figure 9



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

from amplification using Phusion polymerase and are shown in the table (two estimations are included because Illumina's reaction mixtures are proprietary and we do not know the exact reaction conditions). **c**, When only requiring one supporting capture, 3-6% of variants should be derived from first round PCR errors, although more than half of these will be eliminated by the requirement that 55% of reads for a capture support the variant (errors from subsequent PCR rounds will be even more efficiently eliminated by the 55% cutoff). If we require that the same variant be present at the same location across multiple captures before it is included in the final results, it becomes exponentially more unlikely that a first round PCR error would get included. In contrast, increased capture number requirements have a much more modest effect on variants called. **d**, While increasing the number of required supporting captures eliminates rare variants as well as first round PCR errors, the numbers of identified variants only decreases modestly for all individuals (blue line, left y-axis). In contrast, the number of variants expected to be identified as a result of first round PCR amplification errors exponentially decreases with each extra capture requirement (red line, right y-axis). When compared to the number of variants that pass all filters and processing, the first round PCR errors appear to have minimal effect even when only a single capture is required. Expectedly, as we increase the number of required captures supporting a variant, the total number of variants also decreases, and after two required captures should essentially not include mutations created by PCR amplification. Throughout most of this paper, a single capture is used, so as to not bias results by variant 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 $E_{tot} = 0.09 * 150 \text{ bp} * 13.5 \text{ reads}$

57 $E_{tot} = 182.25$

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

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

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

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

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

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

68 The number of expected PCR amplification errors to pass all cutoffs is then estimated
69 using a Gaussian distribution. The logic is that the first round of PCR amplification will
70 create errors that will be at an allele frequency near 50 percent as an error will be
71 created in one of two strands of a captured sequence. Using a Gaussian distribution
72 with a mean at 50, the number of all PCR amplification errors expected to pass the 1
73 supporting capture and 55-95 percent of sequencing reads criteria can be calculated by
74 integrating under the Gaussian distribution. Since we expected about 405 first round
75 PCR amplification errors, and subsequent errors will exist at much smaller allele
76 frequencies, the expected number of variants expected to pass criteria is calculated as
77 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.

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

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

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

98 **a**, Low frequency variants tend to exist close to a $y=x$ line, while high frequency SNPs
99 differ across individuals. As expected, such SNPs cluster around frequencies of 0.5 and
100 1 ($R^2=0.243364$). **b**, When samples are re-sequenced, they show a high degree
101 of similarity, both among SNPs and more rare variants ($R^2=0.568749$). **c**,
102 Though repeat sequencing of individuals typically results in close matches of VAF,
103 repeats do not more closely each other than they match the VAF population mean or
104 any other typical sample. This suggests that the differences observed between samples
105 is likely due to sampling differences than to real differences in individual mutation loads.

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

108 For consistency, all samples used in the main analysis derive from a single bulk library
109 preparation and sequencing run. To ensure that the observed trends are not the result
110 of some bias specific to this single preparation, the entire process was independently
111 repeated, with eleven different blood biopsies (Cohort 2). **a**, Cohort 2 samples closely
112 resembled averaged allele frequencies from the Cohort 1 ($R^2 = 0.455316$,
113 $p\text{-value} = 0.000000$). **b**, Comparing Cohort 2 samples against the VAF mean created
114 from Cohort 2 samples produces a similar pattern to the same comparison using the
115 Cohort 1 data ($R^2 = 0.615327$, $p\text{-value} = 0.000000$). **c-d**, Similar mutation
116 patterns along captured regions were observed for Cohort 2 as for cohort #1 (Fig. 2e).
117 **e**, To understand if observed variant frequencies are the result of clonal expansions or
118 independent events, heterozygous variants were separated by allele. The logic behind
119 this analysis is that if independently captured variants result from the same original
120 event (i.e. a clone), then these variants should be found on the same allele.
121 Alternatively, if variants result from independent events, then such variants should be
122 frequently found on both alleles. By following linkage between variants and
123 heterozygous SNPs, the two alleles can be distinguished. Shown here are the allele
124 frequencies of variants found on either Allele 1 along the x-axis or Allele 2 along the
125 y-axis (analyses are restricted to genomic segments from individuals containing
126 heterozygous SNPs). As the variants adhere to a $y=x$ line, they appear randomly
127 distributed between both alleles, suggesting that variants detected represent multiple
128 independent events rather than clonal expansions.

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

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

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

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

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

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

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

171 identified signatures, our data showed overrepresentation of only 7 of them (Signatures
172 3, 4, 8, 12, 20, 22 and 30) across different samples. Out of seven signatures, Signature
173 12, 3 and 4 had maximum contributions. Signature 3 and 4 are known to be associated
174 with failure of DNA double stranded break repair by homologous repair mechanism and
175 tobacco mutagens respectively, whereas the aetiology of Signature 12 remains
176 unknown. **c**, There was no systematic difference in mutation signatures between
177 amplicons when grouped by their genomic context, and they also showed similar
178 pattern of enrichment of few signatures as compared to others, with signature 12, 3 and
179 4 having maximum contributions. Signature 12 and 4 exhibits transcriptional strand bias
180 for T>C and C>A substitutions respectively, whereas signature 3 is associated with
181 increased numbers of large InDels.

182 **Extended Data Figure 7: Oncogenic mutations do not show evidence of selection.**
183 As shown in Fig. 2f, known oncogenic mutations within probed regions do not show
184 evidence of positive selection. Shown here are additional probed oncogenic loci
185 according the their observed VAFs across donor ages, which also do not show an
186 increase in variant allele frequency in older ages.

187 **Extended Data Figure 8: MMR^{MT} VAFs are elevated over parental frequencies.**
188 When compared to MMR sufficient HCT116 parental cell line genomic DNA, MMR
189 deficient HCT116 cell DNA ($R^2 = 0.066023$) contains substitution mutations at
190 significantly elevated frequencies, as expected with DNA repair deficiencies (Fig. 3a-b).

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

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

197 As for comparisons of MMR^{MT} and HCT116 parental cell lines, a cord blood donor
198 showed a variant population that significantly deviated from expected VAFs (Fig. 3d). **a**,
199 The mutation spectrum found within individual 2 fits to a linear regression line of
200 $y=1.9x+0.00004$, from which it can be seen that variants are approximately twofold
201 more prevalent than in the overall population average. Similar to the data in Extended
202 Figure 8, base substitutions altering **b-d**) C or G nucleotides did not show elevated
203 frequencies. As in the in the MMR^{MT} cells, **e-g**) T or A changes appear at elevated
204 frequencies. Data from individual 19 looked similar to the data shown here, but is not
205 shown. **h**, To ensure that the increased frequencies of variants are not the result of
206 some experimental anomaly, the DNA from individuals #19 (not shown) and #2 was
207 used in a second experiment. In the experimental repeat, the samples showed nearly
208 identical mutational spectra, with similarly elevated levels of T or A changes. **i**,
209 Indicative of experimental repeatability, when samples were freshly captured and
210 sequenced using FERMI, the same individual was highly similar across experiments,
211 and different individuals were less similar.