KEYWORDS

Sequencing

Genetics

Cancer

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

FERMI: A Novel Method for Sensitive Detection of **Rare Mutations in Somatic Tissue**

L. Alexander Liggett,*,[†] Anchal Sharma,[‡] Subhajyoti De,[‡] and James DeGregori^{*,†,§,**,††,1}

*Department of Biochemistry and Molecular Genetics, University of Colorado School of Medicine, Aurora, CO 80045, [†]Linda Crnic Institute for Down Syndrome, University of Colorado School of Medicine, Aurora, CO 80045, [‡]Rutgers Cancer Institute, New Brunswick, NJ 08901, [§]Integrated Department of Immunology, University of Colorado School of Medicine, Aurora, CO 80045, **Department of Pediatrics, University of Colorado School of Medicine, Aurora, CO 80045, and ⁺⁺Department of Medicine, Section of Hematology, University of Colorado School of Medicine, Aurora, CO 80045

ABSTRACT With growing interest in monitoring mutational processes in normal tissues, tumor heteroge-

neity, and cancer evolution under therapy, the ability to accurately and economically detect ultra-rare mutations is becoming increasingly important. However, this capability has often been compromised by significant sequencing, PCR and DNA preparation error rates. Here, we describe FERMI (Fast Extremely Rare Mutation Identification) - a novel method designed to eliminate the majority of these sequencing and library-preparation errors in order to significantly improve rare somatic mutation detection. This method leverages barcoded targeting probes to capture and sequence DNA of interest with single copy resolution. The variant calls from the barcoded sequencing data are then further filtered in a position-dependent fashion against an adaptive, context-aware null model in order to distinguish true variants. As a proof of principle, we employ FERMI to probe bone marrow biopsies from leukemia patients, and show that rare mutations and clonal evolution can be tracked throughout cancer treatment, including during historically 21 intractable periods like minimum residual disease. Importantly, FERMI is able to accurately detect nascent Ill clonal expansions within leukemias in a manner that may facilitate the early detection and characterization 6,5 of cancer relapse.

The simultaneous growth in accuracy and reduction in cost of DNA sequencing has encouraged its use throughout many diverse areas of biology. Accompanying this explosion of applications for sequencing has been a natural demand for increasingly sensitive sequencing methods. While the detection of high frequency variants like germline SNPs is not particularly challenging by most sequencing technologies, sequencer and library-preparation error rates are typically high enough to mask most rare or somatic variants. What is perhaps most challenging about library preparation is that the very isolation of DNA exposes it to oxidation that can change base identities (Shibutani et al. 1991; Cheng

- 52
 - Copyright © 2019 Liggett et al.
 - doi: https://doi.org/10.1534/g3.119.400438
- Manuscript received June 12, 2019; accepted for publication July 14, 2019; 55 published Early Online September 1, 2019.
- This is an open-access article distributed under the terms of the Creative 57 Commons Attribution 4.0 International License (http://creativecommons.org/ licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction 58 in any medium, provided the original work is properly cited.
- Supplemental material available at FigShare: https://doi.org/10.25387/ a3 9037457
- ¹Corresponding author. E-mail: james.degregori@ucdenver.edu

Genes Genes Genomes Genetics

et al. 1992), and high-temperature exposure can thermally alter nucleotide identities (Lindahl and Karlstrom 1973; Lindahl and Nyberg 1974).

Because of these sequencing and library-preparation limitations, quantitative PCR (qPCR) and multiparameter flow cytometry (MFC) have remained common methods of rare variant detection (Terwijn et al. 2013). More recent technologies such as high-throughput digital droplet PCR (Sykes et al. 1992; Vogelstein and Kinzler 1999; Hindson et al. 2011), COLD-PCR (Li et al. 2008; Milbury et al. 2012), and BEAMing (Dressman et al. 2003) have shown promise for rare mutation detection, but are often limited to variant allele frequencies (VAFs) greater than 1% or are restricted to assaying only a few chosen mutations at a time.

A number of studies have sought improvements in sequencing technology accuracies by targeting and labeling small regions of genomic DNA such as sMIPs (Hiatt et al. 2013), by paired strand collapsing (Kennedy et al. 2014) and through other targeting methods (Flaherty et al. 2012; Kim et al. 2013; Albitar et al. 2017; Onecha et al. 2019; Mansukhani et al. 2018; Thol et al. 2018). Some groups have also in- 7 corporated error-correction methods to eliminate sequencing and PCR errors, like PELE-Seq (Preston et al. 2016), and error-correcting



11 12

13 14 15

16

17

18

19

⁵¹

105 Here, we describe a novel integrated genomic method that utilizes 106 single molecule tagging and position specific background correction to 107 push the limit of detection to variants existing in as little as 0.01% of a 108 sample. Initial detection improvements come from the quantitative 109 tracking ability of molecular barcodes that facilitate the elimination of 110 the vast majority of sequencer and PCR amplification errors. Combined 111 with paired-end sequence collapsing, consensus reads are produced that 112 contain reduced numbers of false variants.

113 In a similar manner to previous methods (Newman et al. 2016; 114 Young et al. 2016), we then experimentally derive a background of 115 expected errors for each position within the consensus reads. As we 116 know that sequence context impacts nucleotide stability (Benzer 1961; 117 Nachman and Crowell 2000; Lercher et al. 2001; Hwang and Green 118 2004; Gaffney and Keightley 2005), we use this background to correct 119 our consensus reads based on probability density functions created for 120 each assayed nucleotide position. We build on previous methods by 121 then extensively characterizing the background error probabilities that 122 generally occur in our sequencing library preparations. These charac-123 terizations were sufficiently comprehensive when applying different 124 degrees of bottlenecks to cell lines, it appeared that no variants were 125 detected when none existed within the 1/10,000 detectable range, sup-126 porting the possibility that FERMI often eliminates all background 127 mutations.

128 One recent application of especially sensitive sequencing technolo-129 gies is assaying and understanding clonal evolution within cancerous 130 tissues (Greaves and Maley 2012). The rarity of somatic mutations, 131 even within the clonally expanding pool of cells that exists within a 132 tumor, has limited the observation of changes that can occur. Such an 133 understanding would be valuable, as cancer therapies often leave be-134 hind a small number of cells that can frequently lead to relapse. In 135 leukemias, the state during which these small numbers of cells remain 136 after initial treatment is referred to as minimal residual disease (MRD). 137 During this MRD stage, residual leukemia cells continue to evolve, and 138 successful detection of relapsing leukemia at early stages would facili-139 tate improved prognosis and treatment strategies (Krönke et al. 2011; 140 Ivey et al. 2016).

141As a proof of principle, we directly sample leukocyte genomic DNA142and demonstrate the ability of FERMI to detect oncogenic changes143during the MRD state, and monitor clonal changes with time. We also144show that by concurrently sampling a diverse panel of oncogenic regions,145we can detect the expansion of new oncogenic variants during MRD.146Such observations could be critically important in predicting relapse in147patients.

149 MATERIALS AND METHODS

150 151 Amplicon design

Amplicon probes for targeted annealing regions were created using the Illumina Custom Amplicon DesignStudio (https://designstudio.illumina.com/). UMIs were then added to the designed probe regions and generated by IDT using machine mixing for the randomized DNA. Probes were PAGE purified by IDT. All probes are listed in Table S2 along with binding locations and expected lengths of captured sequence.

159 160

148

Genomic DNA isolation

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

Amplicon capture

For amplicon capture from gDNA, we modified the Illumina protocol called "Preparing Libraries for Sequencing on the MiSeq" (Illumina Part #15039740 Revision D). DNA was quantified with a NanoDrop 2000c (ThermoFisher Catalog #ND-2000C). 500ng of input DNA in 15µl was used for each reaction instead of the recommended quantities. In place of 5µl of Illumina 'CAT' amplicons, 5µl of 4500ng/µl of our amplicons were used. During the hybridization reaction, after gDNA and amplicon reaction mixture was prepared, sealed, and centrifuged as instructed, gDNA was melted for 10 min at 95° in a heat block (SciGene Hybex Microsample Incubator Catalog #1057-30-O). Heat block temperature was then set to 60°, allowed to passively cool from 95° and incubated for 24hr. Following incubation, the heat block was set to 40° and allowed to passively cool for 1hr. The extension-ligation reaction was prepared using 90 µl of ELM4 master mix per sample and incubated at 37° for 24hr. PCR amplification was performed at recommended temperatures and times for 29 cycles. Successful amplification was confirmed immediately following PCR amplification using a Bioanalyzer (Agilent Genomics 2200 Tapestation Catalog #G2964-90002, High Sensitivity D1000 ScreenTape Catalog #5067-5584, High Sensitivity D1000 Reagents Catalog #5067-5585). PCR cleanup was then performed as described in Illumina's protocol using 45 µl of AMPure XP beads. Libraries were then normalized for sequencing using the Illumina KapaBiosystems qPCR kit (KapaBiosystems Reference # 07960336001).

Sequencing

Prepared libraries were pooled at a concentration of 5 nM. Libraries were sequenced on the Illumina HiSeq 4000 at a density of 12 samples per lane with 5% PhiX DNA included, or on the Illumina NovaSeq 6000, allocating approximately 30 million reads per sample.

Bioinformatics

The analysis pipeline used to process sequencing results can be found under FERMI here: http://software.laliggett.com/ or here: https:// github.com/liggettla/FERMI. For a detailed understanding of each function provided by the analysis pipeline, please refer directly to the software. The overall goal of the software built for this project is to analyze amplicon captured DNA that is tagged with equal length UMIs on the 5' and 3' ends of captures, and has been paired-end sequenced using dual indexes. Input fastq files are either automatically or

221

161

162

163

164

165

166

167

168

169

170

171

172

173 174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

manually combined with their paired-end sequencing partners into a

single fastq file. Paired reads are combined by eliminating any base that

does not match between Read1 and Read2, and concatenating this

consensus read with the 5' and 3' UMIs. A barcode is then created

for each consensus read from the 5' and 3' UMIs and the first five bases

at the 5' end of the consensus. All consensus sequences are then binned

together by their unique barcodes. The threshold for barcode mismatch

can be specified when running the software, and for all data shown in

this manuscript one mismatched base was allowed for a sequence to still

count as the same barcode. Bins are then collapsed into a single con-

sensus read by first removing the 5' and 3' UMIs. Following UMI

removal, consensus sequences are derived by incorporating the most

commonly observed nucleotide at each position, so long as the same

nucleotide is observed in at least a specified percent of supporting reads

(75% of reads was used for results in this manuscript) and there are least

some minimum number of reads supporting a capture (5 supporting

reads was used for results in this manuscript). Any nucleotide that does

not meet the minimum threshold for read support is not added to the

consensus read, and alignment is attempted with an unknown base at

that position. From this set of consensus reads, experimental quality

measurements are made, such as total captures, total sequencing reads,

average capture coverage, and estimated error rates. Typically we re-

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

probed position. These confidence intervals were created by using 263 FERMI to sequence control peripheral blood samples from the same 264 experiment as test samples. Following the logic described in Results, it 265 was assumed that low frequency variants that are detected across 266 multiple individuals (including in blood and sperm, where few variants 267 are expected) were not real signal but rather false positive background. 2.68All of the variants from these control samples were thus used to construct 269 a standard background. This background was calculated for each 270 position at which a variant was observed within the standard control 271 samples, and was uniquely calculated for each type of change. Often, in 272 the construction of the background, the highest frequency alleles were 273 eliminated in an effort to minimize the effect of true mutations on the 274 background. A student's t continuous random variable function was 275 used to create a probability density function that describes the back-276 ground distribution for each substitution type at every probed locus 277 (Oliphant 2007). By specifying a particular alpha fraction of the distri-278 bution, high and low VAF endpoints were derived that were then used 279 to determine if an experimental signal was significantly above back-280 ground. While the number of samples required to make a useful back-281 ground will certainly specific to a particular experiment, for the 282 analyses performed in this manuscript and associated work, 5-10

samples seemed to provide sufficient data for the construction of a
background. As outlier variants can also be eliminated at each nucle-
otide position, it is helpful to note that in some cases, samples can be
internally controlled without requiring separate samples.283
284285
286285

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

322

323

324

325

326

327

328

329

330

331

332

333

334

335

Elimination of false positive sequencing and library creation artifacts

A number of steps have been included within sample preparation and bioinformatics analysis specifically to reduce false background signal. Using the dilution series shown in Figures 1C-D, we can show sufficient sensitivity to identify signal diluted to levels as rare as 10^{-4} . While these dilutions show significantly improved sensitivity over many current sequencing methods, background error could still exist. The two largest sources of erroneous mutation when sequencing DNA will typically be from PCR amplification mutations (caused both by polymerase errors and exogenous insults like oxidative damage), and sequencing errors.

These are the steps taken to eliminate errors before final background derivation:

- Elimination of first round PCR amplification errors
- Elimination of subsequent PCR amplification errors
- Elimination of sequencing errors

Elimination of first round PCR amplification errors in consensus reads

The first round of PCR amplification performed during library-prep-308 aration causes mutations that are challenging to distinguish from those 309 that occurred endogenously. Since there is little difference between those 310 mutations that occur during the first round of PCR amplification and 311 those that occurred endogenously, we rely on probability to eliminate 312 these errors. Since we are performing sequencing of individually cap-313 tured alleles, we can ask whether requiring that a mutation be observed in 314 multiple captured alleles before it is called as a true positive signal alters 315 the frequency of variants identified. We expect about 400 first round 316 PCR amplification errors, and the probability that the identical mutation 317 will occur in multiple cells becomes exponentially unlikely. By requiring 318 a mutation be observed in just five captures before it is called as real 319 signal, theoretically, none of the first round PCR amplification errors 320 should make it into the final consensus reads. 321

Elimination of subsequent PCR amplification errors

Elimination of PCR amplification errors after the first round of PCR is done using UMI collapsing (Figure 1A). Each time a strand is amplified, the UMI will keep track of its identity. Any mutations that occur after the first round of PCR will be found on average in 25% of the reads (or fewer for subsequent rounds). This allows us to collapse each unique capture and eliminate any rarely observed variants (<75%) associated with a given UMI. Utilizing the UMI in this way allows us to essentially eliminate any PCR amplification errors that occurred after the first round of PCR. The method should also eliminate most errors resulting from DNA oxidation *in vitro*.

Elimination of sequencing errors

Sequencing errors are eliminated in two ways. This first method is by 336 337 using paired-end sequencing to read each strand of a DNA fragment 338 (Figure 1A). The sequence of these reads (Read1 and Read2) should 339 match if no sequencing errors have been made. For an error to escape elimination it would need to occur at the same position (changing to 340 the same new base) within both Read1 and Read2. Therefore, when the 341 base call differs at a position on Reads 1 and 2, these changes are 342 343 eliminated from the final sequence. This collapsing should eliminate



most sequencing errors, although sequencing errors of the same iden-tity occurring at the same position will escape. These errors should be removed when collapsing into single capture bins (Figure 1A). As with the logic when eliminating subsequent PCR amplification errors, most sequences associated with each UMI pair should be identical. There-fore, sequencing errors passing through Read1 and Read2 will be very unlikely to match other sequenced strands from the same capture event, and are eliminated during consensus sequence derivation.

Data availability

4 | L. A. Liggett et al.

The raw sequence data produced for this study are available in the Sequence Read Archive. The data are available as raw fastq files which have been prepared and sequenced as described in this manuscript. These fastq files can be analyzed using the FERMI software provided on https://github.com/liggettla/FERMI. The fastq files can be found in the BioProject repository under BioProject ID PRJNA525088. This is the direct link to the hosted files: https://www.ncbi.nlm.nih.gov/bioproject/ PRJNA525088. Sequence data are available at BioProject with the ac-cession number: PRJNA525088. The code used to generate and analyze the data can be found at https://github.com/liggettla/fermi The authors affirm that all other data necessary for confirming the conclusions of the article are present within the article, figures, and tables. Supplemen-tal material available at FigShare: https://doi.org/10.25387/g3.9037457.

RESULTS

Method overview

We devised FERMI as a method to overcome current sequencing challenges facing rare mutation detection. FERMI is based on Illumina's TrueSeq Custom Amplicon and AmpliSeq Myeloid protocols, which are designed for mutation detection across selected genomic regions. In FERMI, sequences found within human genomic DNA (gDNA) are captured by targeted oligomer probes, which are then sequenced and analyzed for the presence of any existing mutations. We adapted the AmpliSeq process to target a much smaller number of regions of the genome (32 vs. 1500 regions) in order to achieve a greater sequencing depth per location with a reduced sequencing cost. We designed DNA probes to our 32 selected regions, each approximately 150bp in length, that span either AML-associated oncogenic mutations or Tier III (non-conserved, non-protein coding and non-repetitive sequence) regions of the human genome. The exonic regions were selected to efficiently cover loci that are commonly oncogenically mutated by substitutions in leukemias based on COSMIC classifications. The gDNA used for capture and sequencing was purified either from blood, cancer cell line or sperm cells, though most of our work focused on peripheral blood cells. The method should be adaptable to any species.

Barcode-guided single molecule sequencing

Capture of gDNA, including any existing variants, begins by incubating double-stranded gDNA together with oligomer probes designed to bind specified regions of the genome (Figure 1A). These probes span regions of approximately 150bp in length and contain two identifying indexes. The first index is a 16-bp sequence specific to each sample being

Figure 1 Method overview. A) Schematic representing the steps involved in identifying mutations with FERMI. B) The average number of unique captures varies by probe location. Error is standard deviation across 20 samples.

processed, and the second is a 12-bp unique molecular identifier (UMI)
of randomized DNA that should be unique to each captured strand of
gDNA. Double-stranded gDNA is melted apart to allow these targeting
probes to bind the resulting single-stranded DNA. Probe annealing is
then achieved by slowly cooling the samples to allow for efficient
targeting.

472 Following hybridization of the probes and gDNA, DNA polymerase 473 is used to copy the template, and DNA ligase joins the strands together 474 into a single contiguous amplicon. Using the sample indexes and the 475 capture-specific UMIs to ensure each capture is tracked, amplicons are 476 amplified by polymerase chain reaction (PCR) and pooled together for 477 sequencing. Samples were sequenced using paired-end 150 bp sequenc-478 ing, allocating approximately 30 million Illumina HiSeq or NovoSeq 479 reads per sample. This coverage encompasses on average about 480 1,000,000 capture reactions per sample, resulting in about 30X sequenc-481 ing coverage for each capture (given an average of 30,000 captures per 482 probed region). Though capture efficiency was not uniform for the 483 different probes, which exhibit a fivefold range in the numbers of 484 successful unique captures, we show sufficient coverage at each probed 485 location to capture mutations at least as rare as 0.01% (Figure 1B). 486

487 Assessment of background error profile

488 Following sequencing, reads are distributed into sample-specific bins by 489 their sample index. Within these sample-specific bins, paired-end reads 490 are combined into single consensus reads by marking all mismatched 491 base calls as an unknown identity. This approach yielded better results 492 than elimination of pairs with some threshold of mismatches, as it 493 retained substantially more sequencing information. These paired-end 494 consensus reads are then sorted into capture-specific bins by their UMI 495 sequences. These capture-specific bins are then collapsed into final 496 consensus reads. In order to qualify for this final UMI-based consensus 497 derivation, a UMI-specified capture is required to have at least 5 sup-498 porting sequencing reads, and the base at each position is only called if 499 75% of supporting reads agree with its identity. The final consensus reads 500 are then compared against an experimentally determined background to 501 distinguish true-positive variants from false positive signal, as described 502 below.

503 Though UMI barcode collapsing of sequencing probes is an impor-504 tant technique by which sequencing sensitivity and accuracy can be 505 increased (Hiatt et al. 2013), we find that UMI-collapsed data still 506 retains a significant amount of false-positive variant signal. Using leu-507 kocyte gDNA purified from putatively healthy blood donors, we find 508 approximately 5000 unique variants within our UMI collapsed consen-509 sus sequences in each individual. To estimate how much of this signal 510 might be false-positive background, consensus sequences were compu-511 tationally binned by the presence or absence of heterozygous SNPs 512 found within our probed individuals. This sorting created bins of se-513 quencing that should have originated from only a single allele. Theo-514 retically, if rare variants were indicative of mutations that existed 515 in-vivo, by their very nature of being rare, the mutations should exhibit 516 an associative bias with only one of the two alleles. When we call 517 variants within these two allele-specific bins however, we find that 518 the variants associate quite uniformly across both alleles, suggesting 519 that much of the variant signal found within our final consensus reads 520 is erroneous (Figure 2A).

521 Further suggestive of a significant false-positive presence within 522 consensus reads, we show that when the rare variants found within 523 the blood of any two individuals are compared, the same variants are 524 found in each sample at nearly the same allele frequencies (Figure 2B). 525 This similarity is not limited to inter-blood sample comparisons, as 526

527 sperm gDNA shows similar patterns (Supplemental Figure 1). Finally, 528 being similar to that in blood, the somatic mutation load we observe in sperm cells is well above previous estimates of less than 100 mutations 529 per genome (Lynch 2016). Furthermore, when blood from healthy 530 individuals was compared, mutations were no more similar in repeats 531 532 from the same individual than between individuals (Supplemental Fig-533 ure 4). Combined, these observations suggest that a false-positive background exists relatively uniformly across samples and sample types, 534 and invites the possibility for a correction algorithm to distinguish real 535 from false signal in order to significantly improve sequencing detection 536 limits. 537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

While over 90% of detected background variants were substitutions, occasionally insertions and deletions (indels) were observed. As many of these indels were observed in multiple captures for the same regions, we thought they might represent real mutations. However, as shown in Figure 2C, individual samples contain roughly 500 insertions or deletions, and about 250 of these are conserved across all samples. Furthermore, when a group of 20 individuals was pooled, only one insertion was not found at least twice within the pool. As these indels are often found in multiple captures, the repetitive occurrence between individual samples suggests that some mutagenic mechanism during sample processing is responsible for indel occurrence. Because of this recurrent observation, without modification to the protocol it seems detecting indels would be quite challenging.

Within our final consensus reads, single-nucleotide substitutions account for the majority of falsely identified variants, and within this group of variants, there is significant identity bias. We find that C > Tsubstitutions account for nearly 50% of the variants present in our final consensus reads, while other changes like C > G and T > G are far more rare indicative of differences in nucleotide stability or mutability (Figure 2D). Breaking down these substitutions into their trinucleotide contexts by including the bases located 5' and 3' of each change, we find that sequence context significantly impacts the probability of a false variant being identified (Figure 2E). Among the trinucleotide contexts, false variants within CpG sites are overrepresented within our final consensus reads.

Importantly, the patterns we identify in the trinucleotide contextindependent and context-dependent substitutions mirror those identified in other studies of both normal tissues and cancers (Alexandrov *et al.* 2013; Martincorena *et al.* 2015; Blokzijl *et al.* 2016). The similarity of these patterns provides a cautionary note for mutation detection, as obedience to known patterns does not necessarily provide confidence in the accuracy of calls. It is possible that in-vivo mechanisms of mutation generation are similar to those experienced by template DNA ex-vivo, and therefore results in similar patterns within the background.

Nucleotide context insufficiently explains background signal

575 In search of common patterns within our false positive background, we 576 looked for surrounding sequence contexts that play a role in the 577 prevalence of a false variant. While trinucleotide context does impact the probability that a substitution is found within our final consensus 578 579 read pool, it often incompletely predicts the resulting variant allele 580 frequency (VAF). We observe that many of the background substitu-581 tions found within our final consensus reads such as C > A within the 582 contexts of CCA and ACA, exist within two relatively distinct VAF 583 groups (Figure 3A). This separation indicates that within a given tri-584 nucleotide context, a substitution such as C > A will occur with either a 585 high or a low frequency. Alternatively, some substitutions such as C >586 A within the CCG context largely occur with a low frequency, while 587



616 Figure 2 The background of false positive variants is similar across individuals. A) Using heterozygous SNPs to identify different alleles in human 617 leukocyte gDNA, rare variants within consensus reads equally associate with each of the alleles suggesting they are occurring randomly ex-vivo. 618 Axes are VAFs of variants found on corresponding alleles. B) Using FERMI to measure somatic mutation loads within two different samples shows 619 that background signal is similar within leukocyte gDNA. Axes are VAFs of variants found within each individual. Points represent specific 620 substitutions at each locus. C) The insertions/deletions found within three different individuals were identified. Indel counts are shown on the y-axis. Black dots represent sample groups, where the indel counts are those found within all samples indicated (either each sample alone, or 621 commonly found across the indicated samples; for example, the next to the last vertical bar reflects indels found in both samples 3 and 2). 622 Horizontal bars quantify the total number of indels found in a sample. D) Relative prevalence of observed substitutions within background signal 623 found in leukocyte gDNA consensus reads. Complementary changes such as C > T and G > A are combined. Error is standard deviation across 624 20 individuals. E) Relative prevalence of observed substitutions classified by the neighboring upstream and downstream nucleotides (trinucleotide 625 context). Error is standard deviation across 20 individuals. 626

628 other changes such as T > G in the context of CTA almost never exist 629 at a high frequency. Both results suggest that trinucleotide context is 630 not sufficient to predict background substitution rates at a given locus 631 (Supplemental Table 1), consistent with recent reports that broader 632 (epi)genomic contexts play key roles in replication errors, DNA dam-633 age, and repair (Coleman and De 2018). We do find, that regardless of 634 the substitution identity or the trinucleotide context, a substitution 635 defined only by trinucleotide context never exclusively occurs at high 636 frequency. 637

If the background variants are separated by their presence in either 638 the upper or the lower VAF population, we find that for some changes 639 such as C > A, both the 5' and the 3' nucleotides of the trinucleotide 640 context significantly impact the VAF of the change (Figure 3B). This 641 impact of the trinucleotide context is however not present for all 642 changes, such as C > T substitutions, which show minimal bias of 643 any of the possible trinucleotide contexts. We further searched for 644 patterns within the 10bp upstream and downstream of a given change 645 and find that only the triplet context showed any meaningful impact on 646 mutation rate (Supplemental Figure 2). 647

Algorithmic background subtraction eliminates most false positive signal

Although the trinucleotide context alone does not provide a sufficient amount of contextual information to determine the frequency with which a background variant is observed, nucleotide position strongly impacts the VAF of a substitution within the final consensus sequences. Throughout the probed regions, each nucleotide locus shows a unique background signal pattern that is relatively conserved across individuals (Figure 4A; similar conservation of background signal is observed across all other segments, data not shown). Some of the mutational patterns we observe within our background signal are similar to those found in The Cancer Genome Atlas (TCGA) (Supplemental Figure 3). Notably, enrichment for previously defined signatures are evident in background variants, representing artifacts of damage to isolated gDNA.

Within our observed backgrounds, some nucleotide loci exhibit a705strong bias toward a particular base change, showing only one type of706substitution across all tested individuals. This effect is most commonly707observed at nucleotides that exhibit a C > T substitution, where it is708often the only observed change at that locus. Other nucleotide positions709

627

648

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703



Figure 3 Trinucleotide context is insufficient to predict VAF. A) Mean VAF of background variants calculated from two different groups of 10 individuals probed with FERMI. Variants are either classified by substitution identity alone or within a particular trinucleotide context. B) Relative substitution rates for different substitutions classified by triplet context across all probed regions. Error is standard deviation across 20 individuals.

exhibit multiple different background substitutions, some changing to all three possible other bases. We noticed that while the background signal was surprisingly conserved across samples, that variability did exist (Figure 4B). It is often the case that a particular nucleotide locus will exhibit the same types of variants across different samples, but the allele frequencies vary.

Because each nucleotide locus tends to show a similar background across all tested individuals (Figure 2B), it was possible to derive a governing probability distribution for each observed substitution at every probed position. To create this distribution, a probability density function was created using a student's *t* continuous random variable function. This probability density function was then used to calculate the high and low VAF endpoints of a confidence range by using a specified alpha fraction of the distribution.

We compared a number of different types of backgrounds to understand which best allowed us to eliminate false variants (Figure 4C). Initially a generic background was created by deriving a probability density function for each type of nucleotide change based on its neighbors (the three nucleotide patterns shown in Figure 2D), independent of genomic position. This approach was modestly effective at eliminating background variants from samples, as it reduced the total variant calls by about 80% (Figure 4C, "Generic"). By incorporating positional information and deriving a density function for each observed substitution at all nucleotide positions, about 99.9% of variant calls were eliminated, providing very clean sequencing data. Importantly, experimental variability seems to play a significant role in the accuracy of the background. While the same sample sequenced across multiple experiments generally shows a very similar background, experimental variability does appear (Supplemental Figure 4). While a background derived from samples taken from a different experiment ("External") will result in about 10 variants being called as real in a given peripheral blood sample, a background created from samples run in the same experiment ("Internal") will result in about 1-3 variants being called as real (Figure 4C). As expected variants are always retained, but the

total number of significant variants is minimized when using an internal background created from samples of the same experiment, internal backgrounds are used for all subsequent analyses. To understand what alpha fraction of the probability density functions should be called as *bona fide* mutations, we used 10 healthy blood samples to derive a confidence interval range for each observed substitution across all probed nucleotide positions. These confidence intervals were then compiled into a comprehensive false-positive background against which experimental samples were then compared.

For 5 healthy blood samples, variants were called from their derived consensus sequences, and then compared against the comprehensive background. Within these 5 samples, variants were called as confidently above background if their VAFs were high enough to fall within the specific alpha of their governing probability density function. As expected, as the confidence interval alpha fraction was increased from called as confidently above background exponentially decreases (Figure 4D). This method eliminates nearly all background signal by confidence interval alpha fractions in the range of ten 9's. Furthermore, at higher confidence intervals even germline variants are often eliminated for being too close to background, indicating excessive stringency. Peripheral samples taken from leukemic patients at different points during therapy were also tested, and show similar exponential decreases in confident variant calls, though the overall numbers of variants are higher than in healthy blood (Figure 4E).

Assessing FERMI sensitivity and specificity

To help understand the specificity of FERMI in detecting only true mutations, Molm13 acute myeloid leukemia cells were expanded *in vitro* after passing them through bottlenecks of 1, 100, or 1 million cells. We show that many mutations can be observed within the cell cultures started from 100 cells, given that the 100-cell bottleneck should create clones at approximately 1% frequency each with occasional variants in our probed region (Figure 5A). Heterozygous mutations are expected at



Figure 4 Background confidence intervals eliminate most variants. A) Observed substitution VAFs for a Tier III probe, illustrating the varying 853 mutation presence at each nucleotide locus. Error is standard deviation across 20 individuals. B) Subsets of the IDH2 probe region from two 854 different groups of 10 individuals illustrates the degree of similarity and differences that are commonly observed between samples. Error is 855 standard deviation across the individuals. C) Total number of variants deemed significantly above background when only triplet context was used 856 to generate expected background substitution rates (Generic), or when position-specific substitution rates generated from a different experiment 857 (External) or the same experiment (Internal) are used. D) The total numbers of variants called as significantly above background for 5 individuals 858 (labeled as sample number) at confidence intervals from 0.9 - 0.999999999999999. where the number of trailing 9's is indicated by x-axis value. E) 859 The total numbers of variants called as significantly above background for two leukemic patients using sample 4 from MRD1 and sample 3 from MRD2 (See Figure 6), which are both points that had followed treatment, and at which leukemic burden was low. 860

allele frequencies of 0.005 at 2N loci, and at lower VAFs if a variant falls
in a region with greater ploidy. Indeed, most variants fall within this
range. As expected, very few mutations are detected in gDNA isolated

861

from the cell cultures started from 1 million cells, as most mutations will exist at rare allele frequencies (below our limit of detection). Similarly, we observed no mutations within the cultures initiated with



Figure 5 Assessing FERMI specificity and accuracy. A) MOLM-13 cells grown from 1, 100, or 1 million initial cells were expanded to a pool of 1 million cells, and then probed with FERMI. Samples illustrate how clonality impacts mutation detection by FERMI by altering the VAFs of somatic variants. B) Observed frequencies of a serially diluted blood gDNA sample with a heterozygous germline SNP show successful detection at allele frequencies from 1/2 to 1/10,000 (legend indicates the dilutions not the allele frequencies, where a 1/5,000 dilution of a heterozygous mutation should result in a 1/10,000 allele frequency). Background signal mean and standard deviation shown in red and purple respectively, calculated from 12 samples. C) Limit of detection improvements observed when multiple mutations in linkage disequilibrium are leveraged to eliminate erroneous reads. Background signal calculated from 12 samples.

914

915

916

917

918

919

920

921

922

923

924

925

926

927 928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945



Figure 6 Oncogenic mutation detection during leukemia treatments. A) Oncogenic driver detection using FERMI on bone marrow biopsies taken at 6 different timepoints starting with clinical presentation of the patient and ending with relapse (x-axis is in chronological order of leukemia samplings). Background signal mean and standard deviation shown in red and purple respectively, and is derived from 20 samples. B) Oncogenic driver detection throughout leukemia treatment in a case where relapse was not driven by a mutation within our panel. Background derived from 20 samples. 27. C) Example of a leukemic relapse in which the detected driver was not unknown prior to blood sampling. Background derived from 8 samples. D,E,F) Corresponding blast counts as percentage of bone marrow biopsies.

986single cells, consistent with the low odds that a mutation would occur in987our probed region during the ~ 14 cell divisions required to generate988the 10,000 cell limit of detection. The absence of mutations detected in989the cultures initiated with 1 cell each, but the presence of mutations990within the cultures initiated with 100 cells, indicates that we have991sufficiently limited false positive variants, but retained true mutations.

984

985

992 To assess the sensitivity and limit of detection of FERMI, gDNA from 993 human blood containing known heterozygous SNPs was serially diluted 994 into blood gDNA lacking these SNPs. We find that the detection limits of 995 tracking single dilutions to be variable as the level of background noise is 996 position specific, but diluted germline variants were detected at fre-997 quencies at least as rare as 1:10,000 at the expected VAFs (Figure 5B). In 998 other positions, where the background can be much higher, a dilution 999 series would not be detected as low as 1:10,000.

1000 One of the samples tested in the dilution series contained three 1001 heterozygous SNPs within the same probe region, on the same allele, 1002 allowing for an extra level of error correction. Within this sample, it was 1003 assumed that only those consensus reads with all three SNPs or those 1004 without any of the SNPs were correct, and all other reads were 1005 eliminated. This analysis significantly reduced the background error 1006 rates at the SNP positions, and allowed detection of diluted mutations at 1007 least as low as 1:10,000 (Figure 5C). 1008

1009 Oncogenic driver detection in leukemias

1010 To test the ability of FERMI to detect and follow mutations throughout
1011 leukemia treatments, patient biopsies were collected at disease inception,
1012 throughout treatment, and during relapse when possible. Using FERMI,
1013 we are able to detect these mutations when they are present and observe
1014

clonal evolution as it occurs. In some cases, we detect the principle oncogenic driver and watch it fluctuate in frequency in response to treatment without ever disappearing below background (Figure 6A). In another case, a JAK2 mutation is initially observed at high frequency, but treatment eliminates the clone. As relapse occurs, blast counts increase (Figure 6, D–F), but the initial JAK2 clone does not increase in frequency, as the genetics of the leukemia has clearly changed with treatment (Figure 6B). In a third sample, we observe a patient relapse with a previously undetected JAK2 mutation (Figure 6C). While this time point was taken at relapse, we detect it at a frequency significantly below that of most sequencing method sensitivities, requiring only 5 ml of peripheral blood. The early detection of such a clone could allow treatment with a kinase inhibitor before overt disease relapse.

DISCUSSION

In this study, we designed a sensitive sequencing method that enables the accurate detection of rare variants and clonal evolution within primary samples. In leveraging the quantitative power of capture-unique UMI barcodes, we achieve single-allele sequencing resolution from gDNA, and by then combining these sequencing results with a comprehensive analysis of expected background signal, we achieve exceptional sequencing fidelity.

While capture-specific barcoding has been effectively used in the
past, an inability to achieve sufficient capture numbers and high
background have often held the theoretical limit of detection to variants
existing at a frequency of >1/1,000. By probing roughly 30,000 different
unique captures for each region of interest per sample, we pushed our
theoretical limit of detection to at least 1/10,000. This limit is governed1069
10701071
1072
10731071
1073

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1076 by the magnitude of false positive background observed at a particular 1077 location, as variants are difficult to identify when they exist at allele 1078 frequencies below the background signal.

1079 Paired-end collapsing has been successfully used to reduce the 1080 number of sequencing errors within sequencing data. Unfortunately, 1081 errors also occur during library-preparation, and while the molecular 1082 barcodes assist with the elimination of these library-preparation errors, 1083 mistakes made before or during the first round of PCR will typically 1084 appear indistinguishable from a heterozygous variant, such that neither 1085 paired-end collapsing nor molecular barcode collapsing will be capable 1086 of eliminating them. This understanding prompted the development of 1087 an expected false positive background that could be used to filter out common mistakes that occur during library preparation. 1088

1089 Our experimentally derived backgrounds proved vitally important in 1090 determining whether or not a variant found in a sample was sufficiently 1091 elevated in allele frequency that it could be classified as truly existing in 1092 the in-vivo gDNA from which it originated. Similar to background 1093 subtraction employed by other groups (Newman et al. 2016; Chaudhuri 1094 et al. 2017), our generalized background allowed us to not only detect 1095 expected mutations, but also discover new mutations within samples.

1096 It is interesting to note that within our correction background, we 1097 observe similar mutation patterns to those observed for other studies, 1098 and even similar mutational signatures (Alexandrov et al. 2013; Behjati 1099 et al. 2014). This reproducibility may indicate a surprising degree of 1100 similarity between the intrinsic mutagenic processes in-vivo and error-1101 causing processes involved in sample preparation. It is possible that this 1102 similarity is the result of the conserved behavior of error-inducing 1103 machinery like DNA polymerase and DNA ligase both in-vivo and 1104 in-vitro, or even similar mutagenic exposures such as oxidative damage. 1105 These observed similarities suggest caution against using mutation 1106 signatures as validation of the accuracy of sequencing data when 1107 attempting to identify rare variants.

1108 The early detection of clonal evolution within cancer samples has 1109 held the promise of more comprehensive diagnoses and improved 1110 treatment strategies for patients. While deep sequencing has been 1111 applied to patient leukemias in the past, mutation discovery accuracies 1112 have typically limited these approaches to more of a validation role (Thol 1113 et al. 2018). We obtained a number of leukemic patient biopsies, taken 1114 at initial clinical presentation, and throughout treatment, and used 1115 FERMI to search for somatic mutations. We find that while the de-1116 tection limit of FERMI is quite low, the greatest improvements are 1117 made through its accurate mutation detection ability. Because we elim-1118 inate nearly all background variants, we can accurately detect unex-1119 pected relapse mutations and drivers of clonal expansions.

1120 In requiring only around 5ml of blood, FERMI could be easily used in 1121 a clinical setting to quickly, cheaply and easily identify important driver 1122 mutations and clonal evolution within patient's cancers. If relapse mu-1123 tations were caught by FERMI when they are still rare, targeted ther-1124 apies could be used to prevent them from clonally expanding to fixation 1125 and driving leukemic relapse. Additional applications of FERMI could 1126 include analyses of mutational patterns in normal epithelial tissues, 1127 premalignancies and carcinomas obtained through direct biopsies or 1128 via detection in blood. 1129

ACKNOWLEDGMENTS

1130

1131 We would like to thank Ruth Hershberg of Technion University and 1132 Jay Hesselberth and Robert Sclafani of the University of Colorado 1133 School of Medicine for useful suggestions and for review of the 1134 manuscript. We would like to thank Craig Jordan and Amanda 1135 Winters of the University of Colorado School of Medicine for the 1136

primary leukemia samples. We would like to thank the Illumina 1137 Concierge team who assisted with the initial probe design. We would 1138 also like to thank Joe Hiatt, Beth Martin, and Jay Shendure of the 1139 Shendure lab, who assisted in the initial design of the probing 1140 technique. These studies were supported by grants from the National 1141 Cancer Institute (R01CA180175 to J.D.), NIH/NCATS Colorado CTSI 1142 Grant Number UL1TR001082CU (seed grant to J.D.), a SCOR grant 1143 from the Leukemia and Lymphoma Society (to Craig Jordan), 1144 F31CA196231 (to L.A.L.), the Linda Crnic Institute for Down 1145 Syndrome (to J.D. and L.A.L.), P30-CA072720 (to A.S. and S.D.), 1146 and R01-GM129066, P30-CA072720, and Robert Wood Johnson 1147 Foundation (to A.S. and S.D.). The research utilized services of the 1148 Cancer Center Genomics Shared Resource, which is supported in part 1149 by NIH grant P30-CA46934. L.A.L. and J.D. developed the concept of 1150 this project, planned the experiments, analyzed results, and wrote the 1151 1152 manuscript. L.A.L. processed and prepared samples from blood biopsy 1153 to sequencing, and wrote the software used for analysis. A.S. and S.D. 1154 analyzed results, and contributed to writing of the manuscript.

1155

1197

LITERATURE CITED

LITERATURE CITED		1156
Albitar, A. Z., W. Ma, and M. Albitar, 2017 Wild-type Blocking PCR		1157
Combined with Direct Sequencing as a Highly Sensitive Method for		1158
Detection of Low-Frequency Somatic Mutations. J. Vis. Exp. (121).		1159
https://doi.org/10.3791/55130	8	1160
Alexandrov, L. B., S. Nik-Zainal, D. C. Wedge, S. A. J. R. Aparicio, S. Behjati		1161
et al., 2013 Signatures of mutational processes in human cancer. Nature	_	1162
500: 415–421. https://doi.org/10.1038/nature12477	9	1163
Benjati, S., M. Huch, R. van Boxtel, W. Karthaus, D. C. Wedge <i>et al.</i> ,		1164
2014 Genome sequencing of normal cells reveals developmental line-		1165
10 1038/paturel 3448		1166
Benzer, S. 1961 on the topography of the genetic fine structure Proc. Natl.		1167
Acad. Sci. USA 47: 403–415. https://doi.org/10.1073/pnas.47.3.403		1168
Blokzijl, F., J. de Ligt, M. Jager, V. Sasselli, S. Roerink et al., 2016 Tissue-		1169
specific mutation accumulation in human adult stem cells during life.		1170
Nature 538: 260-264. https://doi.org/10.1038/nature19768		1171
Chaudhuri, A. A., J. J. Chabon, A. F. Lovejoy, A. M. Newman, H. Stehr et al.,		1172
2017 Early Detection of Molecular Residual Disease in Localized Lung		1172
Cancer by Circulating Tumor DNA Profiling. Cancer Discov. 7: 1394-		1174
1403. https://doi.org/10.1158/2159-8290.CD-17-0716		1175
Cheng, K. C., D. S. Cahill, H. Kasai, S. Nishimura, and L. A. Loeb, 1992 8-		1176
Hydroxyguanine, an abundant form of oxidative DINA damage, causes		1177
Coleman N and S De 2018 Mutation Signatures Depend on Engenomic		1178
Contexts Trends Cancer Res 4: 659–661 https://doi.org/10.1016/		1170
i.trecan.2018.08.001		119
Dressman, D., H. Yan, G. Traverso, K. W. Kinzler, and B. Vogelstein,		1100
2003 Transforming single DNA molecules into fluorescent magnetic		1101
particles for detection and enumeration of genetic variations. Proc. Natl.		1182
Acad. Sci. USA 100: 8817-8822. https://doi.org/10.1073/pnas.1133470100		1103
Flaherty, P., G. Natsoulis, O. Muralidharan, M. Winters, J. Buenrostro et al.,		1184
2012 Ultrasensitive detection of rare mutations using next-generation		1185
targeted resequencing. Nucleic Acids Res. 40: e2. https://doi.org/10.1093/		1186
nar/gkr861		1187
Fujita, P. A., B. Rhead, A. S. Zweig, A. S. Hinrichs, D. Karolchik <i>et al.</i> ,		1188
2010 The UCSC genome browser database: update 2011. Nucleic Acids		1189
Caffney D L and P D Keightley 2005 The scale of mutational variation		1190
in the murid genome Genome Res 15: 1086–1094 https://doi.org/		1191
10.1101/gr 3895005		1192
Garrison, E., and G. Marth, 2012 Haplotype-based variant detection from		1193
short-read sequencing. arXiv [q-bio.GN].		1194
Greaves, M., and C. C. Maley, 2012 Clonal evolution in cancer. Nature 481:		1195
306-313. https://doi.org/10.1038/nature10762		1196

- Hiatt, J. B., C. C. Pritchard, S. J. Salipante, B. J. O'Roak, and J. Shendure,
 2013 Single molecule molecular inversion probes for targeted, highaccuracy detection of low-frequency variation. Genome Res. 23: 843–854.
 https://doi.org/10.1101/gr.147686.112
- Hindson, B. J., K. D. Ness, D. A. Masquelier, P. Belgrader, N. J. Heredia *et al.*,
 2011 High-throughput droplet digital PCR system for absolute quantitation of DNA copy number. Anal. Chem. 83: 8604–8610. https://doi.org/10.1021/ac202028g
- Hwang, D. G., and P. Green, 2004 Bayesian Markov chain Monte Carlo sequence analysis reveals varying neutral substitution patterns in mammalian evolution. Proc. Natl. Acad. Sci. USA 101: 13994–14001. https://doi.org/10.1073/pnas.0404142101
- Ivey, A., R. K. Hills, M. A. Simpson, J. V. Jovanovic, A. Gilkes *et al.*,
 2016 Assessment of Minimal Residual Disease in Standard-Risk AML.
 N. Engl. J. Med. 374: 422–433. https://doi.org/10.1056/NEJMoa1507471
- Kennedy, S. R., M. W. Schmitt, E. J. Fox, B. F. Kohrn, J. J. Salk *et al.*,
 2014 Detecting ultralow-frequency mutations by Duplex Sequencing. Nat. Protoc. 9: 2586–2606. https://doi.org/10.1038/nprot.2014.170
- Kim, S., K. Jeong, K. Bhutani, J. Lee, A. Patel *et al.*, 2013 Virmid: accurate detection of somatic mutations with sample impurity inference. Genome Biol. 14: R90. https://doi.org/10.1186/gb-2013-14-8-r90
- 1217 Krönke, J., R. F. Schlenk, K.-O. Jensen, F. Tschürtz, A. Corbacioglu *et al.*,
 1218 2011 Monitoring of minimal residual disease in NPM1-mutated acute
 1219 myeloid leukemia: a study from the German-Austrian acute myeloid
 1220 leukemia study group. J. Clin. Oncol. 29: 2709–2716. https://doi.org/
 10.1200/JCO.2011.35.0371
- Lander, E. S., L. M. Linton, B. Birren, C. Nusbaum, M. C. Zody *et al.*,
 2001 Initial sequencing and analysis of the human genome. Nature 409: 860–921. https://doi.org/10.1038/35057062
- Lercher, M. J., E. J. B. Williams, and L. D. Hurst, 2001 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. Mol. Biol. Evol. 18: 2032–2039. https://doi.org/10.1093/oxfordjournals.molbev.a003744
- Li, H., and R. Durbin, 2009 Fast and accurate short read alignment with
 Burrows-Wheeler transform. Bioinformatics 25: 1754–1760. https://
 doi.org/10.1093/bioinformatics/btp324
- Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan *et al.*, 2009 The Sequence Alignment/Map format and SAMtools. Bioinformatics 25: 2078–2079. https://doi.org/10.1093/bioinformatics/btp352
- Lindahl, T., and B. Karlstrom, 1973 Heat-induced depyrimidination of deoxyribonucleic acid in neutral solution. Biochemistry 25: 5151–5154.
 https://doi.org/10.1021/bi00749a020
- Lindahl, T., and B. Nyberg, 1974 Heat-induced deamination of cytosine residues in deoxyribonucleic acid. Biochemistry 13: 3405–3410. https:// doi.org/10.1021/bi00713a035
- Li, J., L. Wang, H. Mamon, M. H. Kulke, R. Berbeco *et al.*, 2008 Replacing PCR with COLD-PCR enriches variant DNA sequences and redefines the sensitivity of genetic testing. Nat. Med. 14: 579–584. https://doi.org/10.1038/nm1708
- Lynch, M., 2016 Mutation and Human Exceptionalism: Our Future Genetic Load. Genetics 202: 869–875. https://doi.org/10.1534/genetics.115.180471
 Mangukhani S. L. L. Barber, D. Kleftogiannis, S. Y. Moorgraft, M. Davideon
- Mansukhani, S., L. J. Barber, D. Kleftogiannis, S. Y. Moorcraft, M. Davidson *et al.*, 2018 Ultra-Sensitive Mutation Detection and Genome-Wide
 DNA Copy Number Reconstruction by Error-Corrected Circulating
- 1249

10.1272/dincham 2018 280620		1051
10.15/5/clinchem.2018.289029 Martinagana I. A. Dashan, M. Canstung, D. Ellis, D. Van Las, et al.		1251
2015 Turn on evolution High hunder and nervoive negitive coloritor of		1252
2015 Tumor evolution. High burden and pervasive positive selection of		1253
somatic mutations in normal numan skin. Science 548: 880–886. https://		1254
dol.org/10.1126/science.aaab806		1255
Milbury, C. A., M. Correll, J. Quackenbush, K. Rubio, and G. M. Makri-		1256
giorgos, 2012 COLD-PCR enrichment of rare cancer mutations prior to		1257
targeted amplicon resequencing. Clin. Chem. 58: 580–589. https://		1257
doi.org/10.1373/clinchem.2011.176198		1250
Nachman, M. W., and S. L. Crowell, 2000 Estimate of the mutation rate per		1259
nucleotide in humans. Genetics 156: 297–304.		1260
Newman, A. M., A. F. Lovejoy, D. M. Klass, D. M. Kurtz, J. J. Chabon et al.,		1261
2016 Integrated digital error suppression for improved detection of		1262
circulating tumor DNA. Nat. Biotechnol. 34: 547–555. https://doi.org/		1263
10.1038/nbt.3520		1264
Oliphant, I. E., 2007 Python for Scientific Computing. Comput. Sci. Eng. 9:		1265
10-20. https://doi.org/10.1109/MCSE.2007.58		1266
Onecha, E., M. Linares, I. Rapado, Y. Ruiz-Heredia, P. Martinez-Sanchez		1267
et al., 2019 A Novel deep targeted sequencing method for minimal		1269
residual disease monitoring in acute myeloid leukemia. Haematologica		1200
104: 288–296. https://doi.org/10.3324/haematol.2018.194/12	12	1209
Preston, J. L., A. E. Royall, M. A. Randel, K. L. Sikkink, P. C. Phillips et al.,		12/0
2016 High-specificity detection of rare alleles with Paired-End Low		1271
Error Sequencing (PELE-Seq). BMC Genomics 17: 464. https://doi.org/		1272
10.1186/S12864-016-2669-3		1273
Schmitt, M. W., E. J. Fox, M. J. Prindle, K. S. Reid-Bayliss, L. D. True <i>et al.</i> ,		1274
2015 Sequencing small genomic targets with high efficiency and ex-		1275
treme accuracy. Nat. Methods 12: 423–425. https://doi.org/10.1038/	43	1276
nmetn.5551 Skilastani G. M. Talaakita and A. D. Caallanan 1001. Jacontina of marifa	13	1277
Snibutani, S., M. Takesnita, and A. P. Grollman, 1991 Insertion of specific		1278
bases during DNA synthesis past the oxidation-damaged base 8-oxodG.		1270
Nature 349: 431–434. https://doi.org/10.1038/349431a0		12/9
Sykes, P. J., S. H. Neon, M. J. Brisco, E. Hugnes, J. Condon <i>et al.</i> ,		1280
1992 Quantitation of targets for PCK by use of limiting dilution. Bio-		1281
Ten A C D Abassaia and H M Kang 2015 Unified representation of		1282
ran, A., G. K. Abecasis, and H. M. Kang, 2013 Unined representation of		1283
bioinformatics (http://doi.org/10.1093/		1284
Torruin M. W. I. I. van Duttan, A. Kaldar, V. H. I. van dar Valdan, D. A.		1285
Programme et al. 2012 High programmet of flow		1286
cytometric minimal residual disease detection in acute myeloid leukemia:		1287
data from the HOVON/SAKK AMI 42A study I Clin Oncol 31: 3889-		1288
3807 https://doi.org/10.1200/ICO.2012.45.9628		1289
Thol F. R. Gabdoulline A Liebich P. Klement I. Schiller et al.		1200
2018 Measurable residual disease monitoring by NGS before allogeneic		1290
hematonoietic cell transplantation in AMI_Blood 132: 1703-1713		1291
https://doi.org/10.1182/blood-2018-02-829911	14	1292
Vogelstein B and K W Kinzler 1999 Digital PCB Proc Natl Acad Sci	14	1293
USA 96: 9236-9241 https://doi.org/10.1073/ppas.96.16.9236		1294
Young, A. L. G. A. Challen, B. M. Birmann, and T. E. Druley 2016 Clonal		1295
haematopoiesis harbouring AML-associated mutations is ubiquitous in		1296
healthy adults. Nat. Commun. 7: 12484. https://doi.org/10.1038/		1297
ncomms12484		1298
		1299
Communicating editor: M. Boutros		1300
0		1301
		1.001

Tumor DNA Sequencing. Clin. Chem. 64: 1626-1635. https://doi.org/

GGG September (2019 Liggett et al.) Author query sheet Liggett (GGG_400438)

QA1 If you or your coauthors would like to include an ORCID ID in this article, please provide your respective ORCID IDs along with your corrections.

Note: If you do not yet have an ORCID ID and would like one, you may register for this unique digital identifier at https://orcid.org/register.

- I Please confirm that your Data Availability statement is accurate, or if not, update to include the details of where your data can be found. Details are available in the Materials and Methods instructions at http://www.genetics.org/content/prep-manuscript#text.
- 2 Please supply a mailing address for the corresponding author.
- 3 Any alternations between capitalization and/or italics in genetic and taxonomic nomenclature have been retained per the original manuscript. *G3* style is for genes and alleles to be italicized; please confirm that all nomenclature has been formatted properly throughout. Uppercase Greek letters should remain roman per journal style even when appearing in a term where the overall style is italic (e.g., a gene name such as $kap108\Delta$). Note that headings are set all roman or all italics based on journal style and should not be changed.
- 4 Please confirm or update any and all URLs in your article.
- 5 Please check all figure legends carefully to confirm that any and all labels, designators, directionals, colors, etc. are represented accurately in comparison with the figure images.
- 6 Please verify styling of Greek and math symbols in text and equations throughout article. Check carefully for correct use of boldface, italics, operators, qualifiers, spacing, superscripts, and subscripts. Journal style includes Greek letters set roman (not italic), math variables set in italic, and variable modifiers set in roman.
- [7] In-text citation "Onecha *et al.* (2018)" has been changed to "Onecha *et al.* (2019)" to match reference with "Literature Cited" section. Kindly check.
- 8 Please provide page range and volume for the "Albitar et al. (2017)." Please check.
- 9 Please add erratum for "Alexandrov *et al.* (2013)."
- 10 Please add erratum for "Kennedy et al. (2014)."
- 11 Reference has been updated. Please check.
- 12 Reference has been updated as per PubMed. Please check.
- 13 Reference has been updated. Please check.
- 14 Reference has been updated as per PubMed. Please check.