

## **LFMD: detecting low-frequency mutations in high-depth genome sequencing data without molecular tags**

Rui Ye<sup>1,2,7\*</sup>, Jie Ruan<sup>2,7\*</sup>, Xuehan Zhuang<sup>3\*</sup>, Yanwei Qi<sup>2,7</sup>, Yitai An<sup>2,7</sup>, Jiaming Xu<sup>2,7</sup>, Timothy Mak<sup>4</sup>, Xiao Liu<sup>2,7</sup>, Xiuqing Zhang<sup>2,7</sup>, Huanming Yang<sup>2,6,7</sup>, Xun Xu<sup>2,7</sup>, Larry Baum<sup>1,4,5</sup>, Chao Nie<sup>2,7#</sup> & Pak Chung Sham<sup>1,4,5#</sup>

<sup>1</sup>Department of Psychiatry, Li Ka Shing Faculty of Medicine, The University of Hong Kong, Hong Kong, China;

<sup>2</sup>BGI-Shenzhen, Shenzhen 518083, China;

<sup>3</sup>Department of Surgery, Li Ka Shing Faculty of Medicine, The University of Hong Kong; Hong Kong, China;

<sup>4</sup>Center for Genomic Sciences, The University of Hong Kong, Hong Kong, China;

<sup>5</sup>State Key Laboratory of Brain and Cognitive Sciences, The University of Hong Kong, Hong Kong, China;

<sup>6</sup>James D. Watson Institute of Genome Sciences, Hangzhou, China;

<sup>7</sup>China National GeneBank, BGI-Shenzhen, Shenzhen 518120, China;

\*These authors contributed equally to this work.

#Correspondence should be addressed to C.N. ([niechao@genomics.cn](mailto:niechao@genomics.cn)) or P.C.S ([pcsham@hku.hk](mailto:pcsham@hku.hk)).

Rui Ye, [yerui@connect.hku.hk](mailto:yerui@connect.hku.hk)  
Jie Ruan, [ruanjie@genomics.cn](mailto:ruanjie@genomics.cn)  
Xuehan Zhuang, [sunnyzxh@connect.hku.hk](mailto:sunnyzxh@connect.hku.hk)  
Yanwei Qi, [ywqi.cu@gmail.com](mailto:ywqi.cu@gmail.com)  
Yitai An, [ayt\\_1988@126.com](mailto:ayt_1988@126.com)  
Jiaming Xu, [xujm@haplox.com](mailto:xujm@haplox.com)  
Timothy Mak, [timmak@yahoo.com](mailto:timmak@yahoo.com)  
Xiao Liu, [liuxiao@genomics.cn](mailto:liuxiao@genomics.cn)  
Xiuqing Zhang, [zhangxq@genomics.cn](mailto:zhangxq@genomics.cn)  
Huanming Yang, [yanghm@genomics.cn](mailto:yanghm@genomics.cn)  
Xun Xu, [xuxun@genomics.cn](mailto:xuxun@genomics.cn)  
Larry Baum, [lwbaum@hotmail.com](mailto:lwbaum@hotmail.com)  
Chao Nie, [niechao@genomics.cn](mailto:niechao@genomics.cn)  
Pak Chung Sham, [pcsham@hku.hk](mailto:pcsham@hku.hk)

V4.1 on 2019.10.31

## Abstract

As next-generation sequencing (NGS) and liquid biopsy become more prevalent in research and clinic, there is an increasing need for better methods to reduce cost and improve sensitivity and specificity of low-frequency mutation detection (Alternative Allele Frequency, AAF < 1%). Here we propose a likelihood-based approach, called Low-Frequency Mutation Detector (LFMD), which combines the advantages of duplex sequencing (DS) and bottleneck sequencing system (BotSeqS) to maximize the utilization of duplicate sequence reads. Compared with the existing state-of-the-art method DS, our method achieves higher sensitivity (improved by ~16%), comparable specificity ( $< 4 \times 10^{-10}$  errors per base sequenced) and lower cost (reduced by ~70%) without involving additional experimental steps, customized adapters or molecular tags. LFMD is useful in areas where high precision is required, such as drug resistance prediction and cancer screening. In addition, this method can also be used to improve power of other variant calling algorithms by making it unnecessary to remove polymerase chain reaction (PCR) duplicates.

## Keywords

Low-frequency mutation, liquid biopsy, drug resistance, cancer diagnosis, mitochondria heterogeneity, mosaic somatic mutation

## Introduction

Low-frequency mutations (LFMs) are defined as mutations with allele frequency lower than 1% in an individual's DNA or mixed DNA. It can indicate early stages of cancer[1-3], Alzheimer's Disease (AD)[4] and Parkinson's Disease (PD)[5], study aging[5, 6], identify disease-causing variants[7], predict potential drug resistance[8], diagnose mitochondrial disease before tri-parental *in vitro* fertilization[9], and track the mutational spectrum of viral genomes, malignant lesions, and somatic tissues[8, 10]. To effectively improve signal-to-noise ratio (SNR) and detect LFMs, researchers have developed methods with stringent thresholds, complex experimental procedures[4, 11], single cell sequencing[12-15], circle sequencing[16], o2n-seq[17], and more precise analytic models[6, 18-20]. The bottleneck sequencing system[21] (BotSeqS) and duplex sequencing[22, 23] (DS) utilize duplicate reads generated by polymerase chain reaction (PCR), which are discarded by other methods, to

achieve much higher accuracy. However, current methods still have some limitations in detecting LFMs.

#### *Disadvantages of single cell sequencing and circle sequencing*

For single cell sequencing, DNA extraction is laborious and exacting, with point mutations and copy number biases introduced during the amplification of small amounts of fragile DNA. To increase specificity, only variants shared by at least two cells are accepted as true variants[15]. At present, this method is not cost-efficient and cannot be used in large-scale clinical applications because a large number of single cells need to be sequenced to identify rare mutations.

Circle sequencing only utilizes a single strand of DNA, so its specificity is limited by the error rate of PCR. It controls errors to a rate as low as  $7.6 \times 10^{-6}$  per base sequenced[16] while DS can achieve  $4 \times 10^{-10}$  errors per base sequenced[22].

#### *Disadvantages of BotSeqS*

In contrast, BotSeqS uses endogenous molecular tags to group reads from the same DNA template and construct double-strand consensus reads. As a result, it can detect very rare mutations ( $<10^{-6}$ ) while being cheap enough to sequence the whole human genome[21]. However, it requires highly diluted DNA templates before PCR amplification to reduce endogenous tag conflicts and ensure sufficient sequencing of each DNA template. Thus, it has a high specificity with poor sensitivity. Also, it discards clonal variants and small insertions/deletions (InDels) to limit false positives.

#### *Disadvantages of DS*

Another promising method to eliminate tag conflicts is Duplex sequencing (DS), which ligates exogenous random molecular tags (also known as barcode, unique molecular identifier, UID or UMI) to both ends of each DNA template before PCR amplification. Although sensitive and accurate, much sequencing data is wasted on sequence tags, fixed sequences and a large proportion of read families that contain only one read pair, which arose from a sequencing error on a tag. Since random molecular tags are synthesized with

customized adapters, batch effects might occur during DNA library construction.

Additionally, DS only works on targeted small genome regions[6, 22] rather than on the whole genome.

### *Disadvantages of Tag clustering*

To solve the problem induced by the errors on tags, multiple methods have been developed to cluster similar tags[24-26], where one or two mismatches are allowed in merging two tags. Although tag clustering does improve the sensitivity, it is still not a straightforward way to solve the problem. Inappropriate tag clustering might occur because unstable synthesis of random tags can result in distinct but similar tags.

### *A new approach*

To avoid the aforementioned problems, we present here a new, efficient approach that combines the advantages of BotSeqS and DS. The method uses a likelihood-based model[6, 18, 27] to dramatically reduce endogenous tag conflicts. Then it groups reads into read families and constructs double-strand consensus reads to detect ultra-rare mutations accurately while maximizing the utilization of duplicate read pairs. This simplifies the DNA sequencing procedure, saves data and cost, achieves higher sensitivity and specificity, and can be used in whole genome sequencing. In addition, our new method offers a statistical solution to the problem of PCR duplication in the basic analysis pipeline of NGS data and can improve sensitivity and specificity of other variant calling algorithms without specific experimental designs. As the price of sequencing is falling, the depth and the rate of PCR duplication are rising. The method we present here might help deal with such high-depth data more accurately and efficiently.

## **Results**

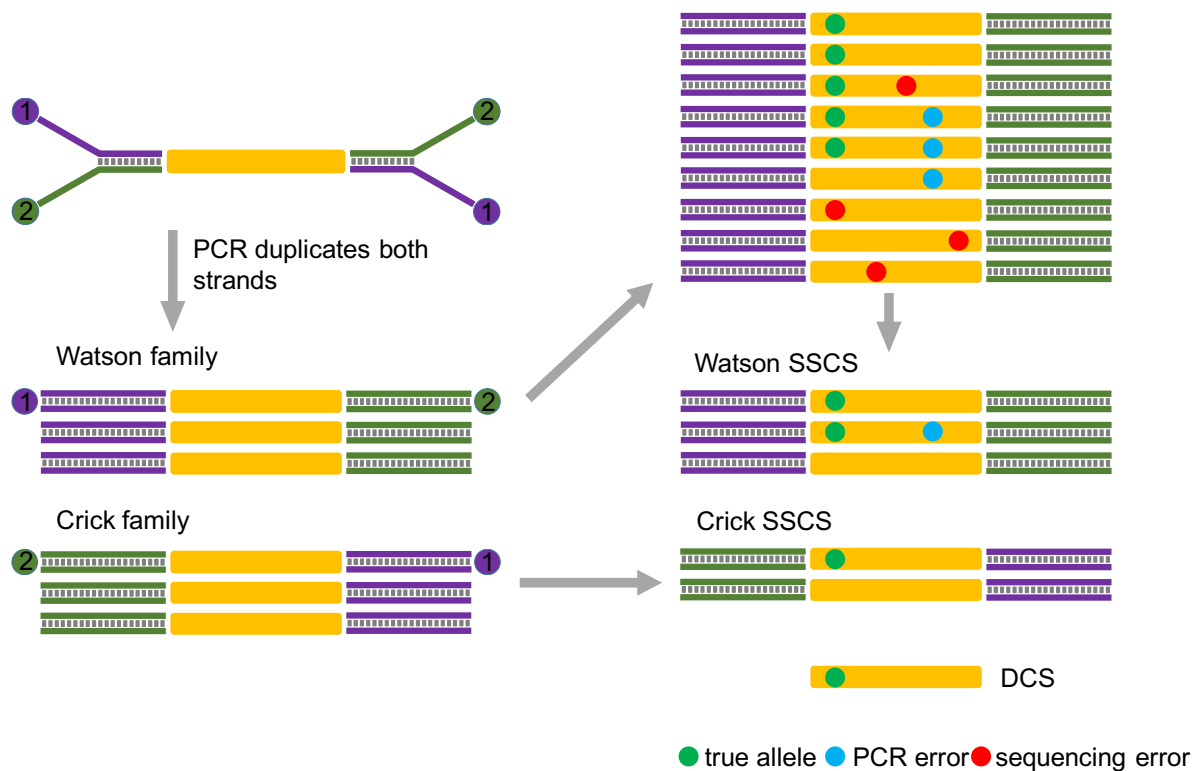
Intuitively, to distinguish LFM (signal) from background PCR and sequencing errors (noise), we need to increase the SNR. To increase SNR, we need to either increase the frequency of mutations or reduce sequencing errors. Single cell sequencing increases the frequency of mutations by isolating single cells from the bulk population, while BotSeqS and

DS reduce sequencing errors by identifying the major allele at each site of multiple reads from the same DNA template. In this paper, we only focus on the latter strategy.

To group reads from the same DNA template, the simplest idea is to group properly mapped reads with the same coordinates (i.e., chromosome, start position, and end position) because random shearing of DNA molecules can provide natural differences, called endogenous tags, between templates. A group of reads is called a read family. However, as the length of DNA template is approximately determined, random shearing cannot provide enough differences to distinguish each DNA template. Thus, it is common that two original DNA templates share the same coordinates. If two or more DNA templates shared the same coordinates, and their reads are grouped into a single read family, it is difficult to determine, using only their frequencies as a guide, whether an allele is a potential error or a mutation. Thus, BotSeqS introduced a strategy of dilution before PCR amplification to dramatically reduce the number of DNA templates in order to reduce the probability of endogenous tag conflicts. And DS introduced exogenous molecular tags before PCR amplification to dramatically increase the differences between templates. Thus, BotSeqS sacrifices sensitivity and DS sequences extra data: the tags.

Here we introduce a third strategy to eliminate tag conflicts. It is a likelihood-based approach based on an intuitive hypothesis: that if reads of two or more DNA templates group together due to endogenous tags, a true allele's frequency in this read family is high enough to distinguish the allele from background sequencing errors. The overview of the LFMD pipeline is shown in Figure 1.

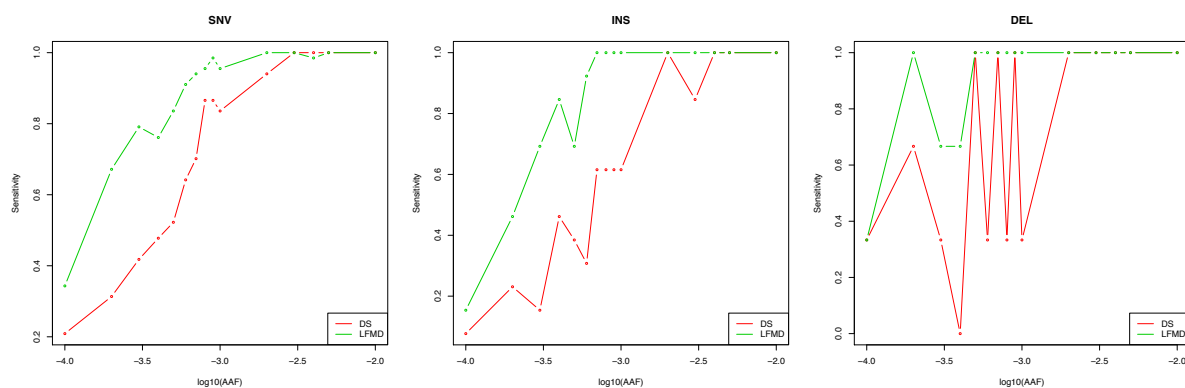
**Figure 1.** Overview of the LFMD pipeline. The Y-shaped adapters determine read 1 (purple bar) and 2 (green bar). The directions of reads determine +/- strands. So after the first cycle of the PCR amplification, the Watson and Crick families are well defined. Then within a read family, true alleles (green dots) and accumulated PCR errors (blue dots) are detected via the likelihood-base model and given a combined error rate. Sequencing errors (red dots) are eliminated. Combining single-strand consensus sequences (SSCSs) of paired read families, high-quality double-strand consensus sequence (DCSs) with estimated error rates are generated and used in the downstream analysis.



### Simulated data

We used Python scripts developed by the Du Novo[28] team to simulate mixed double-strand sequencing data. Although the simulation may not be entirely realistic, the results are still useful to evaluate the performance and the potential drawbacks of LFMD and DS. Because the underlying true mutations are known, the true positives and false positives can be calculated. We found that DS made two false positive rare variant calls due to mapping errors, which were not made by LFMD because it skipped the realignment step. In the meanwhile, LFMD is much more sensitive than DS according to Figure 2 and Table S1. The specificity of LFMD is slightly higher than that of DS shown in Table S2. The simulated mutations are listed in Table S7.

**Figure 2.** The sensitivity of DS and LFMD. SNV, single nucleotide variant; INS, small insertion; DEL, small deletion; AAF, alternative allele frequency.



### Mouse mtDNA

To evaluate the performance of LFMD in real data, we compared LFMD with DS on a DS dataset of mouse mtDNA: SRR1613972. A comparison of DS and LFMD pipelines is shown in Figure S4. We controlled almost all parameters to be the same in DS and LFMD and then compared the results although LFMD can achieve much higher sensitivity if we set the minimum number of supporting reads in each read family as 2 (DS suggests the parameter as 3 to balance sensitivity and specificity). We found that mapping quality influenced the performance of both methods. To reduce the influence of mapping quality, we only use unique proper mapped read pairs to call mutations. The results are shown in Table 1.

**Table 1.** The number of mutations identified by DS and LFMD. The results are based on all unique proper mapped reads of the mouse mtDNA. “DS\_only” and “LFMD\_only” represent mutations only identified by DS or LFMD respectively. “Overlap” is shown for mutations found by both DS and LFMD.

	DS_only	Overlap	LFMD_only	DS_only/Overlap	LFMD_only/Overlap
A>C	0	8	4	0.00%	50.00%
A>G	0	72	8	0.00%	11.11%
A>T	0	28	4	0.00%	14.29%
A>del	0	2	0	0.00%	0.00%
A>ins	0	3	1	0.00%	33.33%
C>A	0	8	3	0.00%	37.50%

C>G	0	2	0	0.00%	0.00%
C>T	0	57	8	0.00%	14.04%
C>del	0	2	2	0.00%	100.00%
C>ins	0	4	0	0.00%	0.00%
G>A	1	19	6	5.26%	31.58%
G>C	0	4	0	0.00%	0.00%
G>T	0	2	3	0.00%	150.00%
G>del	0	1	1	0.00%	100.00%
G>ins	0	1	1	0.00%	100.00%
T>A	0	11	2	0.00%	18.18%
T>C	0	82	10	0.00%	12.20%
T>G	1	10	0	10.00%	0.00%
T>del	1	7	0	14.29%	0.00%
T>ins	0	1	0	0.00%	0.00%
<b>Total</b>	<b>3</b>	<b>324</b>	<b>53</b>	<b>0.93%</b>	<b>16.36%</b>

We investigated the discordant mutations one by one. The DS\_only mutations are all false positives due to read mapping errors (MT:7G>A and MT:3418T>del) and low sequencing quality (MT:5462T>G). The mapping errors are from the following: 1) DS assigns highest base quality for all bases of DCS[22] even though the base is 'N' and 2) DS introduces more 'N' in DCS reads because the minimum proportion of “true” bases in a read family is set to 0.7 by default (it can be set to 0.5 to improve sensitivity).

Among the 53 LFMD\_only mutations, 2 insertions and 3 deletions are missed by DS due to mapping errors of DCS, and 48 SNVs are not detected by DS due to three technical reasons: 1) sequencing and PCR errors on tags lead to smaller read families in DS, decreasing its sensitivity; 2) DS discards some low complexity tags; 3) DS assumes “true” mutations should occupy most reads (proportion  $\geq 0.7$ ) in the reads family, although this assumption is unreasonable considering PCR errors in the first PCR cycle.

In summary, the discordant mutations in this dataset are mainly false positives and false negatives of DS. All LFMD\_only mutations are manually checked and well supported by



read families under the same criteria of DS considering 1 or 2 sequencing or PCR errors on tags. Therefore, LFMD achieves higher sensitivity and specificity in this dataset.

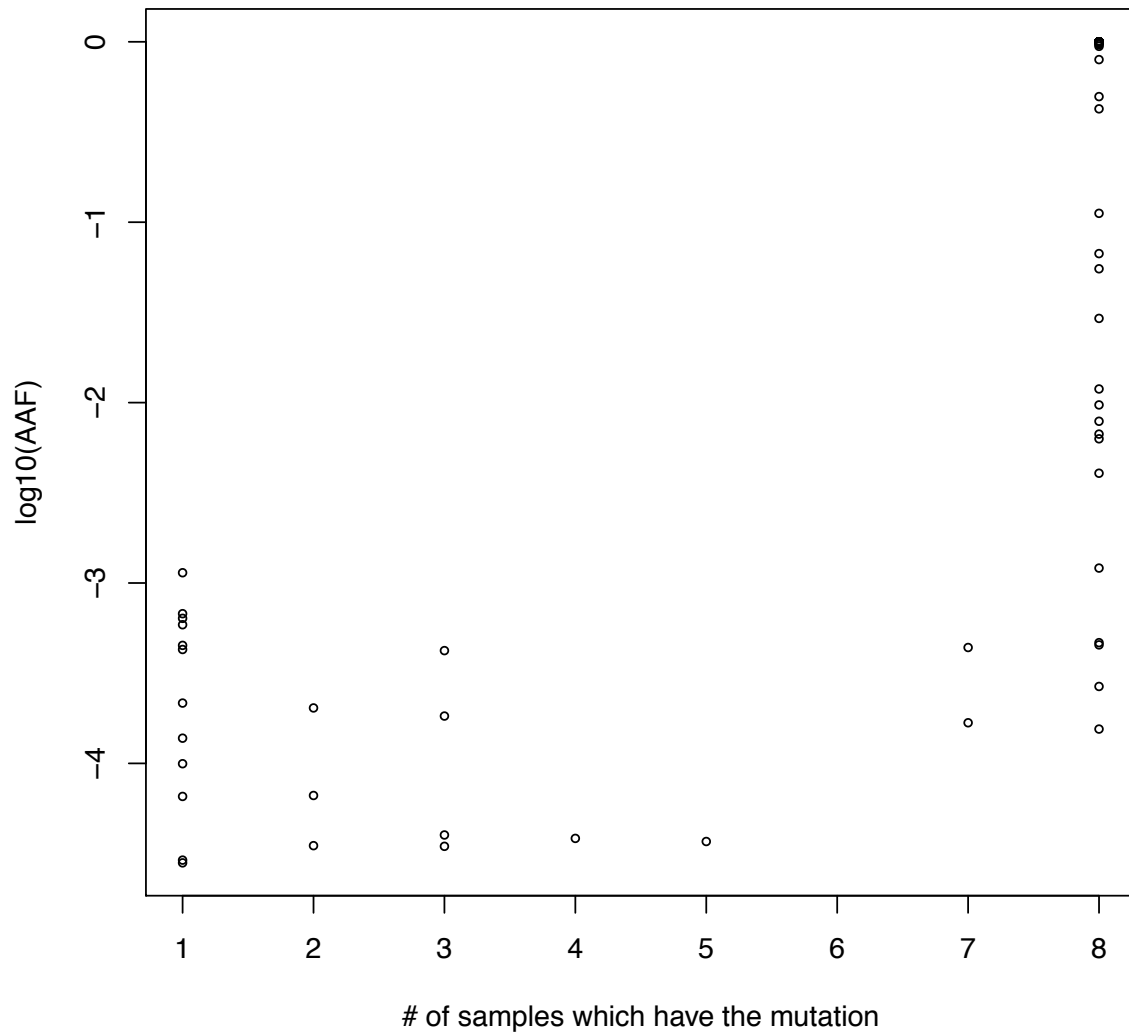
*Twenty-six human mtDNA samples from Prof. Kennedy's laboratory[4]*

We compared the performance of DS and LFMD on 26 samples from Prof. Scott R. Kennedy's laboratory. Only unique proper mapped reads were used to detect LFMs. The majority of LFMs were detected by both tools. Almost all LFMs detected only by DS are false positives due to alignment errors of DCS, while LFMD outputs BAM files directly and avoids alignment errors (Supplemental Material, Figure S4). LFMs only detected by LFMD are well supported by raw reads if considering PCR and sequencing errors on molecular tags. As a result, LFMD is much more sensitive and accurate than DS. The improvement of sensitivity is about 16% according to Table S3.

*YH cell line*

We sequenced the YH cell line (passage 19) in 8 independent experiments to evaluate the stability of LFMD. The experimental details can be found in the Materials part. The results are highly consistent in terms of numbers of mutations detected (range 61~68) shown in Figure 3 and Supplemental Materials, Table S4. Under the hypothesis that true mutations should be identified from at least two samples, we detected 68 "true" mutations and the mean true positive rate (TPR) and false discovery rate (FDR) are around 91.36% and 2.36%, respectively. Due to the high specificity and relative low sensitivity of LFMD, all mutations are almost impossible to be false positive, and the above TPR and FDR only indicate the consistency of parallel experiments.

**Figure 3.** Distribution of mutations found in mtDNA of YH cell lines compared with human Revised Cambridge Reference Sequence (rCRS).



### *ABL1* data

Using the duplex sequencing method in 2015, Schmitt et al. analysed an individual with chronic myeloid leukaemia who relapsed after targeted therapy with the drug, Imatinib (the Short Read Archive under accession SRR1799908). We analyzed this individual and found 5 extra LFMs. Two of them are in the coding region of the *ABL1* gene and change amino acids: E255G and V256G. In the drug resistance database of COSMIC[29], we found that 1) E255VDK, change of the 255<sup>th</sup> amino acid, is associated with resistance to the drugs Dasatinib, Imatinib, and Nilotinib; and 2) V256L is related to resistance to the drug Imatinib. Although the directions of amino acid changes, in this case, are not the same as those in the database, these two additional LFMs still infer potential resistance to the drug Imatinib and

provide an additional explanation for the clinical relapse of leukemia. In the meanwhile, the prediction of potential drug resistance of Dasatinib and Nilotinib is of certain guiding significance for clinical application. The annotation results of 5 LFMs are shown in Supplemental Materials, Table S5.

## Discussion

LFMD is still expensive for target regions >2 Mbp in size because of the need for high sequencing depth. However, as the cost of sequencing continues to fall, it will become more and more practical. In order to sequence a larger region with lower cost, the dilution step of BotSeqS can be introduced into the LFMD pipeline. Because LFMD can deal with tag conflicts, the dilution level can be decreased 1 or 2 magnitudes to increase the sensitivity. Additional experiments will be done soon.

Only accepting random sheared DNA fragments, not working on short amplicon sequencing data, and only working on pair-end sequencing data are known limitations of LFMD. Moreover, LFMD's precision is limited by the accuracy of the alignment software. Although tags are excluded in the model of LFMD, LFMD still has the potential to utilize tags and deal with amplicon sequencing data. The basic assumption in our model, that error rates are the same in all three directions, is not close enough to reality according to experimental data at present. These remain issues may be solved in the next version of LFMD.

To estimate the theoretical limit of LFMD, let read length be 100bp and the standard deviation (SD) of insert size be 20bp. Furthermore, let N represent the number of read families across one point. Then,  $N = (2 * 100) * (2 * 3 * 20) = 24000$  if only considering  $\pm 3$  SD. As the sheering of DNA is not random in the real world, it is safe to set N as 20,000. Ideally, the likelihood ratio test can detect mutations whose frequency is greater than 0.2% in a read family with Q30 bases. Thus, the theoretical limit of minor allele frequency is around  $1e-7$  ( $= 0.002 / 20000$ ).

LFMD reduces the cost dramatically mainly because it discards tags. First, for a typical 100bp read, the lengths of the tags and the fixed sequences between the tag and the true sequence are 12bp and 5bp respectively. So  $(12+5) / 100 = 17\%$  of data are saved if we discard tags directly. Second, the efficiency of targeted capture decreases by about 10 to 20%

because of the tags and the fixed sequences, according to in-house experiments. Third, LFMD can work on short read data (PE50) of BGISEQ and then 30 to 40% of the cost can be saved because of the cheaper sequencing platform. Totally, the cost can be reduced by about 70%.

Based on the concept of molecular tags, Du novo[28], UMI-tools[26] and others have made significant improvements since 2012. In contrast, LFMD shows a completely different approach that discards tags in the first place. It is essential to report an alternative path which shows potential advantages. As DS data and high-depth conventional whole genome sequencing data accumulated in public databases, LFMD provides opportunities for further mining of existing datasets.

## **Conclusion**

To eliminate endogenous tag conflicts, we use a likelihood-based model to separate the read family of the minor allele from that of the major allele. Without additional experimental steps and the customized adapters of DS, LFMD achieves higher sensitivity and specificity with lower cost. LFMD is useful in areas that require ultra-high precision, such as prediction of drug resistance and early screening of cancer. It is a general method that can be used in several cutting-edge areas and its mathematical methodology can be generalized to increase the power of other NGS tools.

## **Methods and Materials**

### *Subject recruitment and sampling*

A lymphoblastoid cell line (YH cell line) established from the first Asian genome donor[30] was used. Total DNA was extracted with the MagPure Buffy Coat DNA Midi KF Kit (MAGEN). The DNA concentration was quantified by Qubit (Invitrogen). DNA integrity was examined by agarose gel electrophoresis. The extracted DNA was kept frozen at -80°C until further processing.

### *Mitochondrial genome DNA isolation*

Mitochondrial DNA (mtDNA) was isolated and enriched by double/single primer set amplifying the complete mitochondrial genome. The samples were isolated using a single primer set (LR-PCR4) by ultra-high-fidelity Q5 DNA polymerase following the protocol of the manufacturer (NEB) (Table S6).

#### *Library construction and mitochondrial whole genome DNA sequencing*

For the BGISEq-500 sequencing platform, mtDNA PCR products were fragmented directly by Covaris E220 (Covaris, Brighton, UK) without purification. Sheared DNA ranging from 150bp to 500bp without size selection was purified with an Axygen™ AxyPrep™ Mag PCR Clean-Up Kit. 100 ng of sheared mtDNA was used for library construction. End-repairing and A-tailing was carried out in a reaction containing 0.5 U Klenow Fragment (ENZYMATICS™ P706-500), 6 U T4 DNA polymerase (ENZYMATICS™ P708-1500), 10 U T4 polynucleotide kinase (ENZYMATICS™ Y904-1500), 1 U rTaq DNA polymerase (TAKARA™ R500Z), 5 pmol dNTPs (ENZYMATICS™ N205L), 40 pmol dATPs (ENZYMATICS™ N2010-A-L), 1 X PNK buffer (ENZYMATICS™ B904) and water with a total reaction volume of 50 µl. The reaction mixture was placed in a thermocycler running at 37°C for 30 minutes and heat-denatured at 65°C for 15 minutes with the heated lid at 5°C above the running temperature. Adaptors with 10bp tags (Ad153-2B) were ligated to the DNA fragments by T4 DNA ligase (ENZYMATICS™ L603-HC-1500) at 25°C. The ligation products were PCR amplified. Twenty to twenty-four purified PCR products were pooled together in equal amounts and then denatured at 95°C and ligated by T4 DNA ligase (ENZYMATICS™ L603-HC-1500) at 37°C to generate a single-strand circular DNA library. Pooled libraries were made into DNA Nanoballs (DNB). Each DNB was loaded into one lane for sequencing.

Sequencing was performed according to the BGISEq-500 protocol (SOP AO) employing the PE100 mode. For reproducibility analyses, YH cell line mtDNA was processed four times following the same protocol as described above to serve as library replicates, and one of the DNBs from the same cell line was sequenced twice as sequencing replicates. A total of 8 datasets were generated using the BGISEQ-500 platform. MtDNA sequencing was performed on the BGISEq-500 with 100bp paired-end reads. The libraries were processed for high-throughput sequencing with a mean depth of ~60000x.

The data that support the findings of this study have been deposited in the CNSA (<https://db.cngb.org/cnsa/>) of CNGBdb with accession code CNP0000297. Analysis codes used in this paper can be accessed at <https://github.com/RainyEricYe/LFMD>.

### *Likelihood-based model*

We aim to identify alleles at each potential heterozygous position in a read family (grouped according to endogenous tags). Then based on those heterozygous sites, we split the mixed read family into smaller ones, and compress each one into a consensus read. Finally, we detect mutations based on all consensus reads, which have much lower error rates than 0.1%.

First, we define a Watson strand as a read pair for which read 1 is the plus strand while read 2 is the minus strand. A Crick strand is defined as a read pair for which read 1 is the minus strand while read 2 is the plus strand. The plus and minus strands are also known as the forward and reverse strands according to the reference genome. Read 1 and 2 are derived from raw pair-end fastq files. Thus a read family which contains Watson and Crick strand reads simultaneously is an ideal read family because it is supported by both strands of the original DNA template. Second, we select potential heterozygous sites which meet the following criteria: 1) the minor allele is supported by both Watson and Crick reads; 2) minor allele frequencies in both Watson and Crick read family are greater than approximately the average sequencing error rate, often 1% or 0.1%; 3) low-quality bases ( $<Q20$ ) and low quality alignments ( $<Q30$ ) are excluded. Finally, we calculate the genotype likelihood in the Watson and Crick family independently in order to eliminate PCR errors during the first PCR cycle.

At each position of a Watson or Crick read family, let  $X$  denote the sequenced base and  $\theta$  the allele frequencies. Let  $P(x|\theta)$  be the probability mass function of the random variable  $X$ , indexed by the parameter  $\theta = (\theta_A, \theta_C, \theta_G, \theta_T)^T$ , where  $\theta$  belongs to a parameter space  $\Omega$ . Let  $g \in \{A, C, G, T\}$ , and  $\theta_g$  represent the frequency of allele  $g$  at this position. Obviously, we have boundary constraints for  $\theta$ :  $\theta_g \in [0, 1]$  and  $\sum \theta_g = 1$ .

Assuming  $N$  reads cover this position,  $x_i$  represents the base on read  $i \in \{1, 2, \dots, N\}$ , and  $e_i$  denotes sequencing error of the base, we get

$$\begin{aligned}
 P(x_i = g|\theta) &= P(\text{no sequencing error} \mid \text{the base is } g) \cdot P(\text{the base is } g) \\
 &\quad + P(\text{sequencing error with specific direction} \mid \text{the base is not } g) \\
 &\quad \cdot P(\text{the base is not } g) \\
 &= (1 - e_i) \theta_g + \frac{e_i}{3} (1 - \theta_g)
 \end{aligned}$$

So the log-likelihood function can be written as

$$\ell(\theta) = \sum_{i=1}^N \log P(x_i|\theta) = \sum_{i=1}^N \log \left( (1 - e_i) \theta_g + \frac{e_i}{3} (1 - \theta_g) \right), \quad g = x_i$$

Thus, for each candidate allele  $g$ , under the null hypothesis  $H_0: \theta_g = 0, \theta \in \Omega$ , and the alternative hypothesis  $H_1: \theta_g \neq 0, \theta \in \Omega$ , the likelihood ratio test is

$$t_g = -2\{\ell_0(\theta) - \ell_1(\theta)\} \sim \chi_1^2$$

However, as  $\theta_g = 0$  lies on the boundary of the parameter space, the general likelihood ratio test needs an adjustment to fit  $\chi_1^2$ . Because the adjustment is related to calculation of a tangent cone[31] in constrained 3-dimensional parameter space, and the computation is too complicated and time-consuming for large scale NGS data, here we use a simplified, straightforward adjustment[32] presented by Chen et al in 2017. (Details in Supplemental materials)

Interestingly, we finally arrive at a general conclusion that the further adjustment of  $\chi_1^2$  is not helpful in similar cases although the asymptotic distribution we use is not perfect when  $N$  is small (e.g.,  $N < 5$ ). Alternative approaches might be derived in the future. We also compared theoretical P-values with empirical P-values from Monte Carlo procedures (Supplemental Material, Figure S1), explored the power of our model under truly and uniformly distributed sequencing errors (Supplemental Material, Figure S2), and evaluated the accuracy of allele frequency (Supplemental Material, Figure S3). The simulation results support the theoretical conclusion sufficiently when P-value is equal or less than  $10^{-3}$ .

Because the null and alternative hypotheses have two and three free variables respectively, the Chi-square distribution has 1 degree of freedom. The P-value of the allele  $g$  can then be given

$$P_g = 1 - \text{cdf}(t_g)$$

where  $\text{cdf}(x)$  is the cumulative density function of the  $\chi_1^2$  distribution. If  $P_g$  is less than a given threshold  $\alpha$ , the null hypothesis is rejected and the allele  $g$  is treated as a candidate allele of the read family.

Although  $P_g$  cannot be interpreted as the probability that  $H_{0,g}$  is true and allele  $g$  is an error, it is a proper approximation of the error rate of allele  $g$  when  $P_g \leq 10^{-3}$ . We only reserve alleles with  $P_g \leq \alpha$  in both Watson and Crick families and substitute others with “N”. Then Watson and Crick families are compressed into several single-strand consensus sequences (SSCSs). The SSCSs might contain haplotype information if more than one heterozygous sites were detected. Finally, SSCSs which are consistent in both Watson and Crick families are claimed as double-strand consensus sequences (DCSs).

For each allele on a DCS, let  $P_w(g)$  and  $P_c(g)$  represent the relative error rates of the given allele  $g$  in the Watson and Crick family respectively, and let  $P_{wc}(g)$  denote the error rate of the allele  $g$  on the DCS. Thus,

$$P_{wc}(g) = P_w(g) + P_c(g) - P_w(g)P_c(g)$$

For a read family which proliferated from  $n$  original templates, a coalescent model can be used to model the PCR procedure[33]. The exact coalescent PCR error rate is too complicated to be calculated quickly, so we tried to give a rough estimate. According to the model, a PCR error proliferates and its frequency decreases exponentially with the number of PCR cycles,  $m$ . For example, an error that occurs in the first PCR cycle would occupy half of the PCR products, an error that occurs in the second cycle occupies a quarter, the third only 1/8, and so on. As we only need to consider PCR errors which are detectable, the coalescent PCR error rate is defined as the probability to detect a PCR error whose frequency  $\geq 2^{-m}/n$ , and it is equal to or less than

$$1 - (1 - \text{error rate per cycle})^{2^m - 1}$$

Let  $e_{pcr}(g)$  denote the coalescent PCR error rate and  $P_{pcr}(g)$  the united PCR error rate of the double strand consensus allele. Then we get

$$P_{pcr}(g) \approx n * e_{pcr}(g)^2$$



Because the PCR fidelity ranges from  $10^{-5}$  to  $10^{-7}$ , we get  $P_{wc}(g)P_{pcr}(g) \approx 0$ , then the combined base quality of the allele on the DCS is

$$Q(g) = -10 \log_{10} (P_{wc}(g) + P_{pcr}(g))$$

Then  $Q(g)$  is transferred to an ASCII character, and a series of characters make a base quality sequence for the DCS. Finally, we generate a BAM file with DCSs and their quality sequences.

With the BAM file which contains all the high-quality DCS reads, the same approach is used to give each allele a P-value at each genomic position which is covered by DCS reads. Adjusted P-values (q-values) are given via the Benjamin-Hochberg procedure. The threshold of q-values is selected according to the total number of tests conducted and false discovery rate (FDR) which can be accepted.

A similar mathematical model was described in detail in previous papers by Ding et al[6] and Guo et al[18]. Ding et al. used this model to reliably call mutations with frequency  $> 4\%$ . In contrast, we use this model to deal with read families rather than non-duplicate reads. In a mixed read family, most of the minor allele frequencies are larger than 4%, so the power of the model meets our expectation.

For those reads containing InDels, the CIGAR strings in BAM files contain I or D. It is obvious that reads with different CIGAR strings cannot fit into one read family. Thus, CIGAR strings can also be used as part of endogenous tags. In contrast, the soft-clipped part of CIGAR strings cannot be ignored when considering start and end positions because low-quality parts of reads tend to be clipped, and the coordinates after clipping are not a proper endogenous tag for the original DNA template.

### **Acknowledgements**

We thank Prof. Scott R. Kennedy for kindly sharing their valuable data. Thank Jeremiah Wala and Rameen Beroukhim for developing a wonderful C++ API, SeqLib[34], and thank S. Bochkanov, V. Bystritsky, and other contributors for the development and maintenance of ALGLIB[35], which accelerated the coding progress of LFMD significantly. We also thank

Heng Li, Richard Durbin and other developers for creating bwa[36] and samtools[37]. We are so grateful to McKenna et al. and Broad Institute for the introduction and maintenance of the GATK[38] framework. And we do appreciate the open source scripts of Nicholas Stoler et al. to simulate tagged sequences[28]. Without the open source spirit and knowledge sharing, our work could not be done.

### **Authors' contributions**

P.C.S. and C.N. supervised the entire work. R.Y. wrote the manuscript. R.Y. and X.H.Z wrote the code and evaluated the performance of LFMD. Y.W.Q assisted the analysis of mitochondrial data. R.J. led all wet-lab experiments and tested the LR-PCR on mitochondria. R.J., Y.T.A and J.M.X prepared the samples and did the experiments. T.M. and P.C.S. helped the mathematical part. P.C.S, C.N., L.B., X.X., H.M.Y, X.Q.Z, X.L., T.M. contributed to critical revision of the manuscript. All authors read and approved the final manuscript.

### **Funding**

This work was supported by the Shenzhen Municipal Government of China [JCYJ20170412153100794]. The funders had no role in study design, data collection, analysis, decision to publish, or preparation of the manuscript.

### **Conflict of Interest**

None declared.

### **References**

1. Newman AM, Lovejoy AF, Klass DM, Kurtz DM, Chabon JJ, Scherer F, Stehr H, Liu CL, Bratman SV, Say C: **Integrated digital error suppression for improved detection of circulating tumor DNA.** *Nature biotechnology* 2016, **34**:547.
2. Phallen J, Sausen M, Adleff V, Leal A, Hruban C, White J, Anagnostou V, Fiksel J, Cristiano S, Papp E: **Direct detection of early-stage cancers using circulating tumor DNA.** *Science translational medicine* 2017, **9**:eaan2415.
3. Newman AM, Bratman SV, To J, Wynne JF, Eclov NC, Modlin LA, Liu CL, Neal JW, Wakelee HA, Merritt RE: **An ultrasensitive method for quantitating circulating tumor DNA with broad patient coverage.** *Nature medicine* 2014, **20**:548.
4. Hoekstra JG, Hipp MJ, Montine TJ, Kennedy SR: **Mitochondrial DNA mutations increase in early stage Alzheimer disease and are inconsistent with oxidative damage.** *Annals of neurology* 2016, **80**:301-306.

5. Simon DK, Lin MT, Zheng L, Liu G-J, Ahn CH, Kim LM, Mauck WM, Twu F, Beal MF, Johns DR: **Somatic mitochondrial DNA mutations in cortex and substantia nigra in aging and Parkinson's disease.** *Neurobiology of aging* 2004, **25**:71-81.
6. Ding J, Sidore C, Butler TJ, Wing MK, Qian Y, Meirelles O, Busonero F, Tsoi LC, Maschio A, Angius A: **Assessing mitochondrial DNA variation and copy number in lymphocytes of ~ 2,000 Sardinians using tailored sequencing analysis tools.** *PLoS genetics* 2015, **11**:e1005306.
7. Wallace DC, Chalkia D: **Mitochondrial DNA genetics and the heteroplasmy conundrum in evolution and disease.** *Cold Spring Harbor perspectives in biology* 2013, **5**:a021220.
8. Schmitt MW, Fox EJ, Prindle MJ, Reid-Bayliss KS, True LD, Radich JP, Loeb LA: **Sequencing small genomic targets with high efficiency and extreme accuracy.** *Nature methods* 2015, **12**:423.
9. Dimond R: **Social and ethical issues in mitochondrial donation.** *British medical bulletin* 2015, **115**:173.
10. Jabara CB, Jones CD, Roach J, Anderson JA, Swanstrom R: **Accurate sampling and deep sequencing of the HIV-1 protease gene using a Primer ID.** *Proceedings of the National Academy of Sciences* 2011, **108**:20166-20171.
11. Marquis J, Lefebvre G, Kourmpetis YA, Kassam M, Ronga F, De Marchi U, Wiederkehr A, Descombes P: **MitoRS, a method for high throughput, sensitive, and accurate detection of mitochondrial DNA heteroplasmy.** *BMC genomics* 2017, **18**:326.
12. Kang E, Wang X, Tippner-Hedges R, Ma H, Folmes CD, Gutierrez NM, Lee Y, Van Dyken C, Ahmed R, Li Y: **Age-related accumulation of somatic mitochondrial DNA mutations in adult-derived human iPSCs.** *Cell Stem Cell* 2016, **18**:625-636.
13. Blandini F, Greenamyre JT, Nappi G: **The role of glutamate in the pathophysiology of Parkinson's disease.** *Functional neurology* 1996, **11**:3-15.
14. Navin N, Kendall J, Troge J, Andrews P, Rodgers L, McIndoo J, Cook K, Stepansky A, Levy D, Esposito D: **Tumour evolution inferred by single-cell sequencing.** *Nature* 2011, **472**:90.
15. Baslan T, Hicks J: **Single cell sequencing approaches for complex biological systems.** *Current opinion in genetics & development* 2014, **26**:59-65.
16. Lou DI, Hussmann JA, McBee RM, Acevedo A, Andino R, Press WH, Sawyer SL: **High-throughput DNA sequencing errors are reduced by orders of magnitude using circle sequencing.** *Proceedings of the National Academy of Sciences* 2013, **110**:19872-19877.
17. Wang K, Lai S, Yang X, Zhu T, Lu X, Wu C-I, Ruan J: **Ultrasensitive and high-efficiency screen of de novo low-frequency mutations by o2n-seq.** *Nature communications* 2017, **8**:15335.
18. Guo Y, Li J, Li C-I, Shyr Y, Samuels DC: **MitoSeek: extracting mitochondria information and performing high-throughput mitochondria sequencing analysis.** *Bioinformatics* 2013, **29**:1210-1211.
19. Kim J, Kim D, Lim JS, Maeng JH, Son H, Kang H-C, Nam H, Lee JH, Kim S: **The use of technical replication for detection of low-level somatic mutations in next-generation sequencing.** *Nature communications* 2019, **10**:1047.
20. Wang TT, Abelson S, Zou J, Li T, Zhao Z, Dick JE, Shlush LI, Pugh TJ, Bratman SV: **High efficiency error suppression for accurate detection of low-frequency variants.** *Nucleic acids research* 2019.

21. Hoang ML, Kinde I, Tomasetti C, McMahon KW, Rosenquist TA, Grollman AP, Kinzler KW, Vogelstein B, Papadopoulos N: **Genome-wide quantification of rare somatic mutations in normal human tissues using massively parallel sequencing.** *Proceedings of the National Academy of Sciences* 2016, **113**:9846-9851.
22. Schmitt MW, Kennedy SR, Salk JJ, Fox EJ, Hiatt JB, Loeb LA: **Detection of ultra-rare mutations by next-generation sequencing.** *Proceedings of the National Academy of Sciences* 2012, **109**:14508-14513.
23. Kennedy SR, Schmitt MW, Fox EJ, Kohn BF, Salk JJ, Ahn EH, Prindle MJ, Kuong KJ, Shen JC, Risques RA, Loeb LA: **Detecting ultralow-frequency mutations by Duplex Sequencing.** *Nat Protoc* 2014, **9**:2586-2606.
24. Peng Q, Satya RV, Lewis M, Randad P, Wang Y: **Reducing amplification artifacts in high multiplex amplicon sequencing by using molecular barcodes.** *BMC genomics* 2015, **16**:589.
25. Kou R, Lam H, Duan H, Ye L, Jongkam N, Chen W, Zhang S, Li S: **Benefits and challenges with applying unique molecular identifiers in next generation sequencing to detect low frequency mutations.** *PloS one* 2016, **11**:e0146638.
26. Smith T, Heger A, Sudbery I: **UMI-tools: modeling sequencing errors in Unique Molecular Identifiers to improve quantification accuracy.** *Genome Res* 2017, **27**:491-499.
27. Wilm A, Aw PPK, Bertrand D, Yeo GHT, Ong SH, Wong CH, Khor CC, Petric R, Hibberd ML, Nagarajan N: **LoFreq: a sequence-quality aware, ultra-sensitive variant caller for uncovering cell-population heterogeneity from high-throughput sequencing datasets.** *Nucleic acids research* 2012, **40**:11189-11201.
28. Stoler N, Arbeithuber B, Guiblet W, Makova KD, Nekrutenko A: **Streamlined analysis of duplex sequencing data with Du Novo.** *Genome biology* 2016, **17**:180.
29. Tate JG, Bamford S, Jubb HC, Sondka Z, Beare DM, Bindal N, Boutselakis H, Cole CG, Creatore C, Dawson E: **COSMIC: the catalogue of somatic mutations in cancer.** *Nucleic acids research* 2018, **47**:D941-D947.
30. Wang J, Wang W, Li R, Li Y, Tian G, Goodman L, Fan W, Zhang J, Li J, Zhang J, et al: **The diploid genome sequence of an Asian individual.** *Nature* 2008, **456**:60-65.
31. Drton M: **Likelihood ratio tests and singularities.** *The Annals of Statistics* 2009, **37**:979-1012.
32. Chen Y, Huang J, Ning Y, Liang K-Y, Lindsay BG: **A conditional composite likelihood ratio test with boundary constraints.** *Biometrika* 2017, **105**:225-232.
33. Weiss G, Von Haeseler A: **A coalescent approach to the polymerase chain reaction.** *Nucleic acids research* 1997, **25**:3082-3087.
34. Wala J, Beroukhir R: **SeqLib: a C++ API for rapid BAM manipulation, sequence alignment and sequence assembly.** *Bioinformatics* 2016, **33**:751-753.
35. Bochkanov S, Bystritsky V: **Alglib.** Available from: [www.alglib.net](http://www.alglib.net) 2013, **59**.
36. Li H, Durbin R: **Fast and accurate short read alignment with Burrows–Wheeler transform.** *bioinformatics* 2009, **25**:1754-1760.
37. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R: **The sequence alignment/map format and SAMtools.** *Bioinformatics* 2009, **25**:2078-2079.
38. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, Garimella K, Altshuler D, Gabriel S, Daly M: **The Genome Analysis Toolkit: a MapReduce**

**framework for analyzing next-generation DNA sequencing data.** *Genome research* 2010, **20**:1297-1303.