

2018

On Statistical Modeling of Sequencing Noise in High Depth Data to Assess Tumor Evolution

R. Rabadan

G. Bhanot

S. Marsilio

N. Chiorazzi

Zucker School of Medicine at Hofstra/Northwell

L. Pasqualucci

See next page for additional authors

Follow this and additional works at: <https://academicworks.medicine.hofstra.edu/publications>



Part of the [Medical Molecular Biology Commons](#)

Recommended Citation

Rabadan R, Bhanot G, Marsilio S, Chiorazzi N, Pasqualucci L, Khiabani H. On Statistical Modeling of Sequencing Noise in High Depth Data to Assess Tumor Evolution. . 2018 Jan 01; 172(1):Article 3812 [p.]. Available from: <https://academicworks.medicine.hofstra.edu/publications/3812>. Free full text article.

This Article is brought to you for free and open access by Donald and Barbara Zucker School of Medicine Academic Works. It has been accepted for inclusion in Journal Articles by an authorized administrator of Donald and Barbara Zucker School of Medicine Academic Works. For more information, please contact academicworks@hofstra.edu.

Authors

R. Rabadan, G. Bhanot, S. Marsilio, N. Chiorazzi, L. Pasqualucci, and H. Khiabani



HHS Public Access

Author manuscript

J Stat Phys. Author manuscript; available in PMC 2019 July 01.

Published in final edited form as:

J Stat Phys. 2018 July ; 172(1): 143–155. doi:10.1007/s10955-017-1945-1.

On statistical modeling of sequencing noise in high depth data to assess tumor evolution

Raul Rabadan,

Department of Systems Biology, Columbia University, New York, NY

Gyan Bhanot,

Department of Physics and Astronomy, Rutgers University, Piscataway, NJ

Sonia Marsilio,

The Feinstein Institute for Medical Research, Northwell Health, Manhasset, NY

Nicholas Chiorazzi,

The Feinstein Institute for Medical Research, Northwell Health, Manhasset, NY

Laura Pasqualucci, and

Institute for Cancer Genetics, Columbia University, New York, NY

Hossein Khiabani

Rutgers Cancer Institute of New Jersey, Rutgers University, New Brunswick, NJ

Abstract

One cause of cancer mortality is tumor evolution to therapy-resistant disease. First line therapy often targets the dominant clone, and drug resistance can emerge from preexisting clones that gain fitness through therapy-induced natural selection. Such mutations may be identified using targeted sequencing assays by analysis of noise in high-depth data. Here, we develop a comprehensive, unbiased model for sequencing error background. We find that noise in sufficiently deep DNA sequencing data can be approximated by aggregating negative binomial distributions. Mutations with frequencies above noise may have prognostic value. We evaluate our model with simulated exponentially expanded populations as well as data from cell line and patient sample dilution experiments, demonstrating its utility in prognosticating tumor progression. Our results may have the potential to identify significant mutations that can cause recurrence. These results are relevant in the pre-treatment clinical setting to determine appropriate therapy and prepare for potential recurrence pretreatment.

Keywords

87.23.Kg; 02.50.-r; 87.18.Tt

1 Introduction

Every extant organism is the result of over three billion years of evolution. Complex organisms consist of cells whose functions are regulated by a large number of interconnected pathways that ensure cellular, tissue, and organ homeostasis. Cancer is a result of the breakdown of this process in a single cell, which results in its unregulated growth. In most cases, the immune system is able to detect and eliminate such aberrant cells. Sometimes, however, a clone escapes this surveillance and manifests as clinically detectable disease [47]. Consequently, most clinically diagnosable tumors are clonal, i.e. they grow clonally from a single cell that finds a path to circumvent the body's defense mechanisms. The growing tumor accumulates mutations, most of which have low or no fitness and therefore are found at low frequencies, outcompeted by the dominant clone [30].

The clonal expansion process, which underlies genomic diversification within a tumor, was first studied by Salvador Luria and Max Delbrück. They designed a simple system of single-cell organisms to investigate patterns of mutation accumulation. Their rigorous quantitative methodology led them to discover that mutations arise randomly and their numbers follow a distinct probability distribution [28]. As the cell population in the tumor diversifies, it is able to explore the fitness landscape. Studying the dynamics of this genomic heterogeneity can yield insight into when the clonal expansion started, how fast the population evolved, and whether specific genomic alterations were selected in a particular host or under a treatment regimen.

The principal biochemical mechanisms in cancer are often recurrent across tumors in different tissues. For example, aberrations leading to unregulated cell growth or inactivation of the apoptotic pathway (cell suicide) are common to almost all tumors. Given the limits within which cells are regulated, the growing tumor has access to only a finite number of pathways that it can alter. As a result, tumors arising from different cells of origin often harbor identical genetic mutations, which alter the same pathways, and often have similar prognostic consequences [5].

First line therapy drugs target a tumor's dominant, fastest growing clone. Drug resistance often emerges from the rise of preexisting clones that harbor potential driver mutations that gain evolutionary fitness via therapy-induced natural selection. It has been shown that the presence of drug-resistant sub-clones in the primary tumor prior to therapy may be a strong predictor of poor survival, with direct implications for disease management [41,51,44,35]. As cancer therapy moves towards individualized treatment, it is important to identify and understand the role of such mutations, some of which may have prognostic value. Such potentially prognostic mutations are commonly identified using targeted deep sequencing of the tumor DNA in clinical settings, and their sensitive detection relies on the accurate analysis of background noise, specifically DNA sequencing errors.

Studying the evolution of chronic lymphocytic leukemia (CLL) under therapy is an illuminating example of these approaches [23,24]. CLL is the most common leukemia in adults and its clinical course ranges from asymptomatic disease that never requires therapy to rapidly progressive disease that requires intensive treatment. Genomic alterations in CLL

follow a time ordered process [52]. Patients who harbor genomic defects in the *TP53* gene, which regulates many pathways including the cell suicide or apoptotic pathway, are considered at high risk of failing conventional therapies [42]. Such patients are good candidates for stem cell transplant or new gene-specific therapeutics [46,2]. The presence of such secondary mutations in genes such as *TP53* is often assessed using traditional Sanger sequencing that only provides sufficient power to detect mutations present in at least 20% of leukemia cells [39]. To assess the presence of *TP53* prognostic mutations at lower abundances in newly diagnosed CLL patients, we used deep sequencing and evaluated thousands of leukemia cells and identified small *TP53* mutations that were missed by traditional methods such as Sanger sequencing [41]. We found that *TP53* mutated sub-clones identified before treatment became the predominant population at the time of CLL relapse, as a result of therapy induced selection pressure. These results suggest that tumors harboring small *TP53* mutations have the same clinical phenotype and risk of failing therapy as those with *TP53* defects in the dominant clone [41,34], and their early detection is essential for the identification and management of high-risk CLL patients [12].

These results are also pertinent to other hematological malignancies where the presence of leukemia-associated mutations in remission is associated with significantly increased risk of relapse and poor survival [38,44]. These data lead to the conclusion that it is imperative to identify alterations that induce therapeutic resistance in leukemia patients in the early stages of disease in order to properly guide individualized therapy with the goal of preventing disease relapse. However, the detection of mutations at low allele frequencies (e.g., 1 mutation in 10,000 cells) is hindered by the lack of a precise model of noise in diagnostic sequencing assays.

Targeted sequencing is the most commonly used method to track prognostic markers in both clinical and basic research applications [10]. However, finding such mutations in sequencing reads is often confounded by misreading a base in the sequencing instrument or misincorporation of DNA bases (nucleotides) during library enrichment by polymerase chain reaction (PCR) amplification cycles. More accurate sequencing protocols, which perform overlapping reads of the same genomic DNA region, allows the merging of such reads for improved accuracy. This facilitates correcting errors accumulated in the sequencer, while leaving uncorrected PCR errors that arise during library preparation steps [4,53,17].

The challenge in identifying potentially functional sub-dominant mutations is to determine the sensitivity thresholds of sequencing platforms, i.e. the depths above which PCR errors happen with a probability below a statistical cut-off. Such thresholds can be estimated by hypothesizing that all variants are due to errors and using deviations from this null hypothesis to indicate the presence of true variants. This can sometimes be confounded by the fact that different sequencing errors occur at different rates [6,3,11], as the mechanism of nucleotide misincorporations during PCR amplification by polymerase molecules is the same is the rise of spontaneous mutations that drive Darwinian variation. Hence a single threshold cannot comprehensively test the significance of all variants. As a result, more sophisticated statistical modeling of the background error distribution is necessary.

To model background error one may use different types of error distributions: (i) a single or a linear combination of Luria-Delbrück distributions, characterizing the expected number of spontaneous nucleotide misincorporations, when PCR error rate is assumed to be constant [20]; (ii) the negative binomial distribution, describing the depth distribution of clones after PCR amplification through a Poisson-Gamma mixture model [36]; and iii) the beta-binomial distribution, suitable for Bayesian models, where error rates are assumed to follow the Beta distribution [25].

The Luria-Delbrück distribution, which has been demonstrated to accurately describe mutagenesis by bacterial polymerases using high-depth error-corrected sequencing [17], is expected to better model the long tail of the error depths. However, empirical analysis has shown that the negative binomial distribution gives the best fit to the observed error depths based on goodness-of-fit log-likelihood [41]. The beta-binomial distribution, in conjunction with multiple filtering criteria based on normal control DNA samples, has also been proposed for somatic mutation detection from cancer genomes [27,8,43,9]. Finally, empirical Bayes methods that establish prior distributions based on observations in the data have also been used for detection and genotyping variants and comparing allele frequencies across different samples for both tumor and viral populations [33,48,21,14].

In this manuscript, we revisit this problem and provide a comprehensive model that illustrates how aggregate negative binomial distributions describe PCR error depths in ultra-deep targeted sequencing. We test our model with *in silico* as well as cell line and patient dilution experiments, and propose a highly sensitive, mutation-specific approach to detect true mutations, without the need for control data from un-mutated (wild type) normal tissue DNA.

2 Methods

Derivation of the error depth distribution

Here, we will only be discussing the distribution of low frequency errors in deep DNA sequencing analysis of tumor samples. Let us assume an experiment in which S independent tissue samples are subjected to ultra-deep sequencing. DNA sequencing of tumor samples produces strings of nucleotides (A, C, G, and T) of 100–200 base-pair length that correspond to the DNA sequences of different sections of the genome in the tumor sample. These sequences of DNA reads are mapped to a “reference” genome and deviations/mismatches are identified as potential mutations. Ideally, the reference sequence is the sequence from the patient’s “germ-line”, usually obtained from blood or some other tissue with normal cells. The sequencing read depth is the average number of reads that map to the same locus (section of the genome). At a nucleotide, three potential single base substitutions can occur: A (adenine) \rightarrow C, G, T, or C (cytosine) \rightarrow A, G, T, or G (guanine) \rightarrow A, C, T, or T (thymine) \rightarrow A, C, G. Alternately, there might be an insertion (addition of one or more A, C, G, T nucleotides) or a deletion (loss of A, C, G, T nucleotides). All of these will henceforth be referred to as variants. We want to derive the posterior probability distribution for these variants, assuming they are stochastic, i.e. they represent noise (statistical random errors).

Suppose that, at a genomic DNA locus, we see n_j such variant reads amongst N_j total reads. The distribution of n_j follows a binomial distribution, $\text{Bino}(n_j | N_j, \theta)$, where θ is the *a priori* probability of a variant's occurrence. Let $M = \sum_{i \neq j}^S N_i$ be the total number of reads across samples at that locus and $m = \sum_{i \neq j}^S n_i$ be the total number of variant (erroneous) reads across samples at that DNA locus. Then, the posterior predictive p value for having detected a true mutation in sample j , given $S - 1$ other samples, can be obtained from the posterior probability distribution:

$$\begin{aligned} P(n_j | N_j, \{n_i, N_i\}) &= \int_0^1 \frac{\text{Bino}(n_j | N_j, \theta) \prod_{i \neq j} \text{Bino}(n_i | N_i, \theta)}{\int_0^1 \prod_{i \neq j} \text{Bino}(n_i | N_i, \theta') d\theta'} d\theta \\ &= \binom{N_j}{n_j} \times \int_0^1 \frac{\theta^{n_j+m} (1-\theta)^{N_j-n_j+M-m}}{\int_0^1 \theta^m (1-\theta)^{M-m} d\theta} d\theta \\ &= \binom{N_j}{n_j} \times \frac{\text{Beta}(1+n_j+m, 1+N_j-n_j+M-m)}{\text{Beta}(1+m, 1+M-m)}, \end{aligned}$$

where Beta indicates the Beta function. Simplifying the algebra yields the beta-binomial distribution,

$$P(n_j | N_j, m, M) = \frac{1+M}{1+N_j+M} \frac{\binom{N_j}{n_j} \binom{M}{m}}{\binom{N_j+M}{n_j+m}}. \quad (1)$$

Variations of equation (1) have been previously derived for sequencing depths $> 100\times$ [8,43,9]. Today, it is possible to do ultra-deep sequencing, where $N_j > 5,000\times$. In such cases, for low frequency variants, we can assume that $n_i \ll N_i$. Therefore, we can use Stirling's

approximation, and estimate $\binom{N_i}{n_i} \approx \frac{N_i^{n_i}}{n_i!}$. Equation (1) can then be approximated by

$$P(n_j | N_j, m, M) = \binom{n_j+m}{n_j} \left(\frac{N_j}{N_j+M}\right)^{n_j} \left(\frac{M}{N_j+M}\right)^{m+1}, \quad (2)$$

which equals $\text{NB}(n_j | 1 + m, \frac{N_j}{N_j + M})$, with NB indicating the negative binomial distribution, and where $1 + m$ and $\frac{N_j}{N_j + M}$ are its two parameters, which we can interpret as the number of detected errors and the *a priori* probability of success in detecting an error, respectively.

Exponential expansions at varying error rates

An exponentially expanded population is generated through c PCR amplification cycles, where each cycle doubles the DNA population. If errors accumulate independently at a rate of μ substitutions per site per cycle, the average error depth (i.e. the average number of reads harboring errors) is $2^c \mu$. For S such populations, the error depth distribution is described by equation (1), or is approximated by a negative binomial distribution, $\text{NB}(1 + (S - 1)2^c \mu, \frac{1}{S})$, as derived above in equation (2).

It is well known that different types of PCR errors occur at different rates. For example, transitions, that exchange two-ring purines (A and G) or one-ring pyrimidines (C and T) are more common than transversions, which replace an A or G with one of C or T. Assuming R independent rates, the observed number of variants with error depth v , $D(v)$, is given by,

$$D(v) = \sum_{r=1}^R X_r P(v | 2^c, (S-1)2^c \mu_r, (S-1)2^c) \quad (3)$$

$$\approx \sum_{r=1}^R X_r \text{NB}(v | 1 + (S-1)2^c \mu_r, \frac{1}{S}),$$

where X_r represents the number of variants that occur with rate μ_r . Since error rates are often unknown and sequence context dependent, we can alternatively bin the variants based on their average error depth across samples and write $D(v)$ as

$$D(v) = \sum_{b=1}^B X_b P(v | \langle N \rangle, (S-1)\langle v \rangle_b, (S-1)\langle N \rangle) \quad (4)$$

$$\approx \sum_{b=1}^B X_b \text{NB}(v | 1 + (S-1)\langle v \rangle_b, \frac{1}{S}),$$

where B is the number of bins, X_b is the number of variants in each bin, and $\langle N \rangle$ is the average sequencing depth across S samples. It has been shown that the sum of negative binomial distributions with equal success probabilities is also a negative binomial distribution, though with a random parameter [7, 50]. Thus, the approximation of $D(v)$ in equations (3) and (4) with sums of negative binomial distributions that have success probability of $\frac{1}{S}$, suggests empirical observations [41].

The MATLAB implementation for simulating exponentially expanding populations as well as calculating cumulative P using equations (1) and (2) are available at software.khiabanian-lab.org.

3 Data

In the first experiment, a series of dilutions was generated using the SU-DHL-6 cell line (Diffuse Large B-Cell Lymphoma), which carries a heterozygous (one allele altered) *TP53*-Y234C missense mutation (one that changes an amino acid in a protein sequence) [32]. The cells were serially diluted at (1:10, 1:10², 1:10³, 5:10⁴, 1:10⁴, 5:10⁵, and 1:10⁵) by mixing the cell line DNA with *TP53* wild-type genomic DNA from a healthy donor. The *TP53* mutation locus was sequenced at depths of 10,000× (10K×), 100,000× (100K×), and 1,000,000× (1M×) [29].

In the second experiment, samples from undiluted cancer cells from a CLL patient, harboring a heterozygous *SF3B1*-K700E missense transition substitution were analyzed. A diluted sample (1:10³) was also generated by mixing this patient's CLL DNA with wild-type genomic DNA from a healthy donor. The mutated *SF3B1* locus in these samples, in addition to un-mutated genomic DNA from 18 healthy, volunteered individuals were sequenced at a mean depth of 620,000×.

For both experiments, each cell line dilution and patient sample was bar-coded and targeted with amplicon multiplexed sequencing using the Illumina MiSeq (2 × 150 bp) (Genewiz, South Plainfield, NJ). The number of raw sequence reads per samples depended on the utilized instrument. With sufficient DNA (in the form of PCR products or un-amplified genomic DNA), Illumina MiSeq could produce eight million paired-end 150 bp reads [15]. The primers were designed so that the paired-end reads substantially overlapped with each other and each read pair was merged to correct sequencing errors. The merged reads were mapped to the human reference genome (hg19) using the Burrows-Wheeler Aligner (BWA) alignment tool [26], and all variable sites were identified using an inclusive variant caller [11]. Raw sequence data are available at the Sequence Read Archive under Bio Project PRJNA421179.

4 Results

Simulated data

We generated a set of *in silico* experiments with exponentially expanded populations starting from a single, homogenous, 100 base-long sequence of binary bases. Each population was aggregated from four expansions that followed error rates of 10⁻³, 10⁻⁴, 10⁻⁵, and 10⁻⁶ substitutions per site per cycle. The number 12, 14, and 18 of cycles were chosen to produce populations with 16,384, 65,536 and 1,048,576 total reads respectively. Each experiment contained 50 independent populations ($S = 50$) and for each experiment, $D(v)$, the expected number of variants with depth v was calculated using equations (3). This experiment was repeated 100 times. Figure 1 shows the results, as well as statistically significant χ^2 p values indicating high accuracy of the estimates from both the beta-binomial model and its NB approximation.

Dilution experiments

We removed the real diluted *TP53* mutation from cell line sequencing data, and arranged the erroneous variants based on their depth in $5\times$ -sized bins. We then counted the number of variants X_b in each bin, and calculated $D(v)$ using equation (4). Figures 2, 3, and 4 show the results for sequencing depths of $10K\times$, $100K\times$, and $1M\times$, indicating statistically significant χ^2 p values that show a strong concordance between estimates from the beta-binomial model, its NB approximation, and ultra-deep sequencing data. Distinguishing transitions and transversions further clarified the importance of classifying variants using sequencing depth as a proxy for the error rates. We obtain similar results from modeling the ultra-deep sequencing data from the *SF3B1* locus (Figure 5).

Detecting true mutations

We propose two comprehensive approaches to assess the presence of true mutations at very low abundance relative to background. Our methodology does not require matched normal samples or extensive filtering based on variant annotation resources.

First, having established an accurate model to describe the sequencing error distribution, a threshold is determined above which sequencing errors happen with a probability below an established statistical cut-off. These thresholds can be derived from all variants or a subset of variants, for example, only transitions or transversions. Figure 6 shows such thresholds for detecting the *TP53*-Y234C transition mutation in dilution experiments, where we are able to identify the mutation in abundances as low as $5:10^4$ at $10K\times$ and $100K\times$, and $1:10^4$ at $1M\times$, without any false positive calls. As shown in Figures 2, 3, and 4, there is better sensitivity for detecting a transversion substitution.

In the absence of matched normal samples, this approach is especially practical for identifying mutations that may exist in more than one tumor sample. Its application to 309 newly diagnosed CLL patients identified small sub-clonal prognostic mutations in four frequently mutated drivers of this neoplasm, present in 2 out of 1,000 wild-type alleles. These mutations were missed by traditional Sanger sequencing, but were validated by independent deep sequencing and allele-specific PCR [41,40].

Second, we tested an individual mutation in each sample against all other sequenced samples and calculated the cumulative P using equation (1). After correcting for multiple hypotheses using the Benjamini and Hochberg method [1], we generated a list of variants that satisfied a pre-determined false discovery rate. This approach is particularly powerful in identifying patient-specific mutations. We assess the presence of the *SF3B1*-K700E mutation in patient samples, and find the probability of observing the mutation in $1:10^3$ CLL dilution to be extremely significant compared to controls (Table 1). This approach can accurately identify sample-specific mutations by comparing multiple samples at the same exact mutated base.

In comparison of our method to other published variant calling algorithms, one comparable unbiased method is EBCall, whose implementation is based on beta-binomial distributions and establishing priors from normal sequencing data [43]. EBCall requires normal samples; therefore, we removed the reads harboring the diluted mutations in the EBCall analysis to simulate matched normal data. EBCall, with a sensitivity-adjusted configuration,

successfully identified the *SF3B1*-K700E mutation in 1:10³ CLL dilution sample, as well as the *TP53*-Y234C mutation in the least diluted samples at all sequencing depths (i.e. 1:10 in 10K×, 1:10² in 100K×, and 1:10³ in 1M×). However, it failed to detect the mutation at higher dilution levels, and also resulted in four false positive calls at 1M×.

5 Conclusion

Therapeutic resistance, one of the main causes of eventual disease relapse and mortality in cancer patients, is often associated with natural selection of preexisting resistant clones under treatment [13,41]. The detection of such low frequency sub-clones is hindered by a lack of precision-tested diagnostic assays.

Allele-specific, real-time PCR assays have been proposed to identify prognostic variants [31,18,49]. These approaches only target known mutations, and their adaptation to situations with large numbers of variants requires extensive primer calibration. In contrast, high-throughput sequencing provides an unbiased view of tumor heterogeneity and its genomic profile. Various techniques based on unique molecular identifiers have been proposed to correct both polymerase and sequencing errors [22,19,17,37] that facilitate distinguishing real mutations from mistakes that arise during amplification. However, the main hurdle in clinical utilization of these approaches is the requirement for generating very large numbers of sequencing reads to assemble the genome of a single DNA molecule with high confidence at depth > 2,000×.

Here, we addressed this important problem in cancer therapy by introducing a highly sensitive method to model sequencing noise, which allows the detection of prognostic markers of disease recurrence using ultra-deep targeted sequencing. Our approach is based on interrogating data from multiple tumor samples at identical genomic regions and provides an accurate assessment of the error rate at a given position without relying on normal samples. Instead of establishing a fixed detection threshold for all variants, we directly calculate mutation-specific sensitivities. Overall, since ultra-deep sequencing methods are now routinely implemented in the clinic, we believe that the application of our comprehensive model to tumor samples will increase the speed with which patients can be evaluated during disease surveillance. Our method opens up the possibility of exploring the dynamics of cancer clones after treatment, timing the rise of resistance to therapy, and determining the clinical importance of minimal residual disease assessed from liquid biopsy samples for precise disease management [16,45].

Acknowledgments

The authors gratefully acknowledge the constructive feedback of Mohammad Hadigol and Alexandra Jacunski. R.R. acknowledges funding from the NIH (U54CA193313, R01CA185486, and R01CA179044). H.K. acknowledges support from the ACS (IRG-15-168-01), Rutgers Cancer Institute (P30CA072720), and Rutgers Office of Advanced Research Computing (NIH 1S10OD012346-01A1).

References

1. Benjamini Y, Hochberg Y. Controlling the false discovery rate - a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society Series B-Methodological*. 1995; 57(1):289–300.

2. Burger JA, Tedeschi A, Barr PM, Robak T, Owen C, Ghia P, Bairey O, Hillmen P, Bartlett NL, Li J, Simpson D, Grosicki S, Devereux S, McCarthy H, Coutre S, Quach H, Gaidano G, Maslyak Z, Stevens DA, Janssens A, Offner F, Mayer J, ODwyer M, Hellmann A, Schuh A, Siddiqi T, Polliack A, Tam CS, Suri D, Cheng M, Clow F, Styles L, James DF, Kipps TJ. Ibrutinib as initial therapy for patients with chronic lymphocytic leukemia. *New England Journal of Medicine*. 2015; 373(25): 2425–2437. [PubMed: 26639149]
3. Chen L, Liu P, Evans TC, Ettwiller LM. Dna damage is a pervasive cause of sequencing errors, directly confounding variant identification. *Science*. 2017; 355(6326):752–756. [PubMed: 28209900]
4. Chen-Harris H, Borucki MK, Torres C, Slezak TR, Allen JE. Ultra-deep mutant spectrum profiling: improving sequencing accuracy using overlapping read pairs. *BMC Genomics*. 2013; 14(1):96. [PubMed: 23402258]
5. Ciriello G, Miller ML, Aksoy BA, Senbabaoglu Y, Schultz N, Sander C. Emerging landscape of oncogenic signatures across human cancers. *Nat Genet*. 2013; 45(10):1127–33. [PubMed: 24071851]
6. Costello M, Pugh TJ, Fennell TJ, Stewart C, Lichtenstein L, Meldrim JC, Fostel JL, Friedrich DC, Perrin D, Dionne D, Kim S, Gabriel SB, Lander ES, Fisher S, Getz G. Discovery and characterization of artifactual mutations in deep coverage targeted capture sequencing data due to oxidative dna damage during sample preparation. *Nucleic Acids Res*. 2013; 41(6):e67. [PubMed: 23303777]
7. Furman E. On the convolution of the negative binomial random variables. *Statistics and Probability Letters*. 2007; 77(2):169–172.
8. Gerstung M, Beisel C, Rechsteiner M, Wild P, Schraml P, Moch H, Beerenwinkel N. Reliable detection of subclonal single-nucleotide variants in tumour cell populations. *Nature Communications*. 2012; 3:811 EP.
9. Gerstung M, Papaemmanuil E, Campbell PJ. Subclonal variant calling with multiple samples and prior knowledge. *Bioinformatics*. 2014; 30(9):1198–1204. [PubMed: 24443148]
10. Grossmann V, Roller A, Klein HU, Weissmann S, Kern W, Haferlach C, Dugas M, Haferlach T, Schnittger S, Kohlmann A. Robustness of amplicon deep sequencing underlines its utility in clinical applications. *J Mol Diagn*. 2013; 15(4):473–84. [PubMed: 23680131]
11. Hadigol M, , Khiabani H. Merit: a mutation error rate identification toolkit for ultra-deep sequencing applications. *bioRxiv* 2017
12. Hallek M. Chronic lymphocytic leukemia: 2015 update on diagnosis, risk stratification, and treatment. *American Journal of Hematology*. 2015; 90(5):446–460. [PubMed: 25908509]
13. Hata AN, Niederst MJ, Archibald HL, Gomez-Caraballo M, Siddiqui FM, Mulvey HE, Maruvka YE, Ji F, Bhang HE, Krishnamurthy Radhakrishna V, Siravegna G, Hu H, Raouf S, Lockerman E, Kalsy A, Lee D, Keating CL, Ruddy DA, Damon LJ, Crystal AS, Costa C, Piotrowska Z, Bardelli A, Iafrate AJ, Sadreyev RI, Stegmeier F, Getz G, Sequist LV, Faber AC, Engelman JA. Tumor cells can follow distinct evolutionary paths to become resistant to epidermal growth factor receptor inhibition. *Nat Med*. 2016; 22(3):262–9. [PubMed: 26828195]
14. Huang G, Wang S, Wang X, You N. An empirical bayes method for genotyping and snp detection using multi-sample next-generation sequencing data. *Bioinformatics*. 2016; 32(21):3240–3245. [PubMed: 27378300]
15. Illumina. [Accessed: 2017-12-05] Specifications for the miseq system 2017 <https://www.illumina.com/systems/sequencing-platforms/miseq/specifications.html>
16. Ivey A, Hills RK, Simpson MA, Jovanovic JV, Gilkes A, Grech A, Patel Y, Bhudia N, Farah H, Mason J, Wall K, Akiki S, Griffiths M, Solomon E, McCaughan F, Linch DC, Gale RE, Vyas P, Freeman SD, Russell N, Burnett AK, Grimwade D. Group U.K.N.C.R.I.A.W. Assessment of minimal residual disease in standard-risk aml. *N Engl J Med*. 2016; 374(5):422–33. [PubMed: 26789727]
17. Jee J, Rasouly A, Shamovsky I, Akivis Y, Steinman RS, Mishra B, Nudler E. Rates and mechanisms of bacterial mutagenesis from maximum-depth sequencing. *Nature*. 2016; 534(7609): 693–696. [PubMed: 27338792]

18. Jia Y, Sanchez JA, Wangh LJ. Kinetic hairpin oligonucleotide blockers for selective amplification of rare mutations. *Sci Rep.* 2014; 4:5921. [PubMed: 25082368]
19. 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 Protocols.* 2014; 9(11):2586–2606. [PubMed: 25299156]
20. Kessler DA, Levine H. Large population solution of the stochastic luria-delbruck evolution model. *Proc Natl Acad Sci U S A.* 2013; 110(29):11,682–7.
21. Khiabani H, Carpenter Z, Kugelman J, Chan J, Trifonov V, Nagle E, Warren T, Iversen P, Bavari S, Palacios G, Rabadan R. Viral diversity and clonal evolution from unphased genomic data. *BMC Genomics.* 2014; 15(6):S17.
22. Kinde I, Wu J, Papadopoulos N, Kinzler KW, Vogelstein B. Detection and quantification of rare mutations with massively parallel sequencing. *Proceedings of the National Academy of Sciences.* 2011; 108(23):9530–9535.
23. Kipps TJ, , Stevenson FK, , Wu CJ, , Croce CM, , Packham G, , Wierda WG, , O'Brien S, , Gribben J, , Rai K. Chronic lymphocytic leukaemia. *Nature Reviews Disease Primers.* 3, 16,096 EP. 2017
24. Lazarian G, Guieze R, Wu CJ. Clinical implications of novel genomic discoveries in chronic lymphocytic leukemia. *Journal of Clinical Oncology.* 2017; 35(9):984–993. [PubMed: 28297623]
25. Lee JC, Sabavala DJ. Bayesian estimation and prediction for the beta-binomial model. *Journal of Business and Economic Statistics.* 1987; 5(3):357–367.
26. Li H, Durbin R. Fast and accurate short read alignment with burrows-wheeler transform. *Bioinformatics.* 2009; 25(14):1754–60. [PubMed: 19451168]
27. Li M, Stoneking M. A new approach for detecting low-level mutations in next-generation sequence data. *Genome Biology.* 2012; 13(5):R34–R34. [PubMed: 22621726]
28. Luria SE, Delbrück M. Mutations of bacteria from virus sensitivity to virus resistance. *Genetics.* 1943; 28(6):491–511. [PubMed: 17247100]
29. Marsilio S, , Khiabani H, , Fabbri G, , Vergani S, , Scuoppo C, , Montserrat E, , Shpall EJ, , Hadigol M, , Marin P, , Rai KR, , Rabadan R, , Devereux S, , Pasqualucci L, , Chiorazzi N. Somatic cll mutations occur at multiple distinct hematopoietic maturation stages: Documentation and cautionary note regarding cell fraction purity. *Leukemia* 2017 EP –URL
30. Martincorena I, Raine KM, Gerstung M, Dawson KJ, Haase K, Van Loo P, Davies H, Stratton MR, Campbell PJ. Universal patterns of selection in cancer and somatic tissues. *Cell.* 2017; 171(5): 1029–1041.e21. DOI: 10.1016/j.cell.2017.09.042 [PubMed: 29056346]
31. Milbury CA, Li J, Makrigiorgos GM. Pcr-based methods for the enrichment of minority alleles and mutations. *Clin Chem.* 2009; 55(4):632–640. [PubMed: 19201784]
32. Morin RD, Mungall K, Pleasance E, Mungall AJ, Goya R, Huff RD, Scott DW, Ding J, Roth A, Chiu R, Corbett RD, Chan FC, Mendez-Lago M, Trinh DL, Bolger-Munro M, Taylor G, Hadj Khodabakhshi A, Ben-Neriah S, Pon J, Meissner B, Woolcock B, Farnoud N, Rogic S, Lim EL, Johnson NA, Shah S, Jones S, Steidl C, Holt R, Birol I, Moore R, Connors JM, Gascoyne RD, Marra MA. Mutational and structural analysis of diffuse large b-cell lymphoma using whole-genome sequencing. *Blood.* 2013; 122(7):1256–1265. [PubMed: 23699601]
33. Muralidharan O, Natsoulis G, Bell J, Ji H, Zhang NR. Detecting mutations in mixed sample sequencing data using empirical bayes. *Ann Appl Stat.* 2012; 6(3):1047–1067. DOI: 10.1214/12-AOAS538
34. Nadeu F, Delgado J, Royo C, Baumann T, Stankovic T, Pinyol M, Jares P, Navarro A, Martín-García D, Beà S, Salaverria I, Oldreive C, Aymerich M, Suárez-Cisneros H, Rozman M, Villamor N, Colomer D, López-Guillermo A, González M, Alcoceba M, Terol MJ, Colado E, Puente XS, López-Otín C, Enjuanes A, Campo E. Clinical impact of clonal and subclonal tp53, sf3b1, birc3, notch1, and atm mutations in chronic lymphocytic leukemia. *Blood.* 2016; 127(17):2122– 2130. [PubMed: 26837699]
35. Naxerova K, Reiter JG, Brachtel E, Lennerz JK, van de Wetering M, Rowan A, Cai T, Clevers H, Swanton C, Nowak MA, Elledge SJ, Jain RK. Origins of lymphatic and distant metastases in human colorectal cancer. *Science.* 2017; 357(6346):55–60. [PubMed: 28684519]

36. Ndifon W, Gal H, Shifrut E, Aharoni R, Yissachar N, Waysbort N, Reich-Zeliger S, Arnon R, Friedman N. Chromatin conformation governs t-cell receptor jbeta gene segment usage. *Proc Natl Acad Sci U S A*. 2012; 109(39):15,865–70. [PubMed: 22219319]
37. Newman AM, Lovejoy AF, Klass DM, Kurtz DM, Chabon JJ, Scherer F, Stehr H, Liu CL, Bratman SV, Say C, Zhou L, Carter JN, West RB, Sledge GW Jr, Shrager JB, Loo BW Jr, Neal JW, Wakelee HA, Diehn M, Alizadeh AA. Integrated digital error suppression for improved detection of circulating tumor dna. *Nat Biotech*. 2016; 34(5):547–555.
38. Oshima K, Khiabani H, da Silva-Almeida AC, Tzoneva G, Abate F, Ambesi-Impiombato A, Sanchez-Martin M, Carpenter Z, Penson A, Perez-Garcia A, Eckert C, Nicolas C, Balbin M, Sulis ML, Kato M, Koh K, Paganin M, Basso G, Gastier-Foster JM, Devidas M, Loh ML, Kirschner-Schwabe R, Palomero T, Rabadan R, Ferrando AA. Mutational landscape, clonal evolution patterns, and role of ras mutations in relapsed acute lymphoblastic leukemia. *Proc Natl Acad Sci U S A*. 2016
39. Pospisilova S, Gonzalez D, Malcikova J, Trbusek M, Rossi D, Kater AP, Cymbalista F, Eichhorst B, Hallek M, Dohner H, Hillmen P, van Oers M, Gribben J, Ghia P, Montserrat E, Stilgenbauer S, Zenz T. Eric recommendations on tp53 mutation analysis in chronic lymphocytic leukemia. *Leukemia*. 2012; 26(7):1458–1461. [PubMed: 22297721]
40. Rasi S, Khiabani H, Ciardullo C, Terzi-di Bergamo L, Monti S, Spina V, Brusca A, Cerri M, Deambrogi C, Martuscelli L, Biasi A, Spaccarotella E, De Paoli L, Gattei V, Foa R, Rabadan R, Gaidano G, Rossi D. Clinical impact of small subclones harboring notch1, sf3b1 or birc3 mutations in chronic lymphocytic leukemia. *Haematologica*. 2016; 101(4):e135–8. [PubMed: 26819056]
41. Rossi D, Khiabani H, Spina V, Ciardullo C, Brusca A, Fama R, Rasi S, Monti S, Deambrogi C, De Paoli L, Wang J, Gattei V, Guarini A, Foa R, Rabadan R, Gaidano G. Clinical impact of small tp53 mutated subclones in chronic lymphocytic leukemia. *Blood*. 2014; 123(14):2139–47. [PubMed: 24501221]
42. Rossi D, Rasi S, Spina V, Brusca A, Monti S, Ciardullo C, Deambrogi C, Khiabani H, Serra R, Bertoni F, Forconi F, Laurenti L, Marasca R, Dal-Bo M, Rossi FM, Bulian P, Nomdedeu J, Del Poeta G, Gattei V, Pasqualucci L, Rabadan R, Foà R, Dalla-Favera R, Gaidano G. Integrated mutational and cytogenetic analysis identifies new prognostic subgroups in chronic lymphocytic leukemia. *Blood*. 2013; 121(8):1403–1412. [PubMed: 23243274]
43. Shiraishi Y, Sato Y, Chiba K, Okuno Y, Nagata Y, Yoshida K, Shiba N, Hayashi Y, Kume H, Homma Y, Sanada M, Ogawa S, Miyano S. An empirical bayesian framework for somatic mutation detection from cancer genome sequencing data. *Nucleic Acids Res*. 2013; 41(7):e89. [PubMed: 23471004]
44. Shlush LI, Mitchell A, Heisler L, Abelson S, Ng SWK, Trotman-Grant A, Medeiros JF, Rao-Bhatia A, Jaciw-Zurakowsky I, Marke R, McLeod JL, Doedens M, Bader G, Voisin V, Xu C, McPherson JD, Hudson TJ, Wang JCY, Minden MD, Dick JE. Tracing the origins of relapse in acute myeloid leukaemia to stem cells. *Nature*. 2017; 547(7661):104–108. [PubMed: 28658204]
45. Siravegna G, Marsoni S, Siena S, Bardelli A. Integrating liquid biopsies into the management of cancer. *Nat Rev Clin Oncol*. 2017 advance online publication, --.
46. Souers AJ, Levenson JD, Boghaert ER, Ackler SL, Catron ND, Chen J, Dayton BD, Ding H, Enschede SH, Fairbrother WJ, Huang DCS, Hymowitz SG, Jin S, Khaw SL, Kovar PJ, Lam LT, Lee J, Maecker HL, Marsh KC, Mason KD, Mitten MJ, Nimmer PM, Oleksijew A, Park CH, Park CM, Phillips DC, Roberts AW, Sampath D, Seymour JF, Smith ML, Sullivan GM, Tahir SK, Tse C, Wendt MD, Xiao Y, Xue JC, Zhang H, Humerickhouse RA, Rosenberg SH, Elmore SW. Abt-199, a potent and selective bcl-2 inhibitor, achieves antitumor activity while sparing platelets. *Nat Med*. 2013; 19(2):202–208. [PubMed: 23291630]
47. Stewart TJ, Abrams SI. How tumours escape mass destruction. *Oncogene*. 2008; 27(45):5894–5903. [PubMed: 18836470]
48. Trifonov V, Pasqualucci L, Tiacci E, Falini B, Rabadan R. Savi: a statistical algorithm for variant frequency identification. *BMC Syst Biol*. 2013; 7(Suppl 2):S2.
49. Vargas DY, Kramer FR, Tyagi S, Marras SAE. Multiplex real-time pcr assays that measure the abundance of extremely rare mutations associated with cancer. *PLoS One*. 2016; 11(5):e0156, 546.

50. Vellaisamy P, Upadhye NS. On the sums of compound negative binomial and gamma random variables. *Journal of Applied Probability*. 2009; 46(1):272–283.
51. Wang J, Cazzato E, Ladewig E, Frattini V, Rosenbloom DIS, Zairis S, Abate F, Liu Z, Elliott O, Shin YJ, Lee JK, Lee IH, Park WY, Eoli M, Blumberg AJ, Lasorella A, Nam DH, Finocchiaro G, Iavarone A, Rabadan R. Clonal evolution of glioblastoma under therapy. *Nat Genet*. 2016; 48(7): 768–776. [PubMed: 27270107]
52. Wang J, Khiabani H, Rossi D, Fabbri G, Gattei V, Forconi F, Laurenti L, Marasca R, Del Poeta G, Fo R, Pasqualucci L, Gaidano G, Rabadan R. Tumor evolutionary directed graphs and the history of chronic lymphocytic leukemia. *eLife*. 2014; 3:e02,869.
53. Zhang J, Kobert K, Flouri T, Stamatakis A. Pear: a fast and accurate illumina paired-end read merger. *Bioinformatics*. 2014; 30(5):614–620. [PubMed: 24142950]

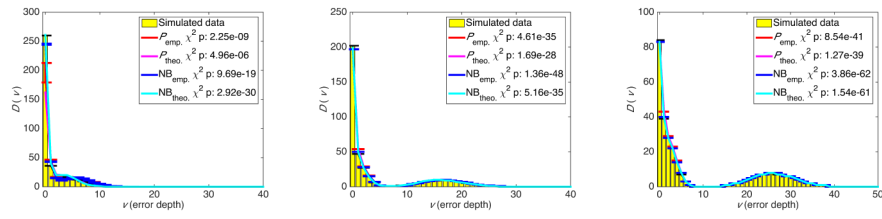


Fig. 1. Number of variants with error depth of v from aggregated simulated cycles of PCR amplification at four error rates: 12 cycles (left), 14 cycles (middle), and 18 cycles (right). P_{theo} . and NB_{theo} . are calculated using equation (3), and P_{emp} . and NB_{emp} . are calculates using equation (4). The χ^2 test was used to compare the distributions.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

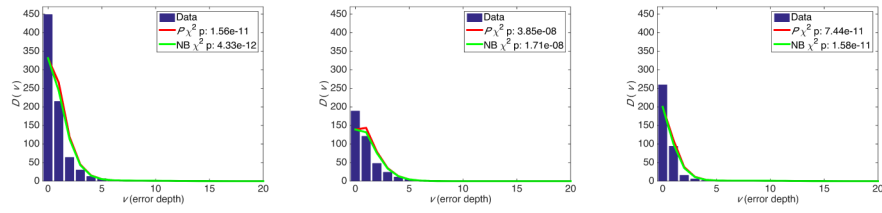


Fig. 2. Error depth distribution in ultra-deep sequencing of a *TP53* locus at 10,000× for all variants (left), transitions (middle), and transversions (right).

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

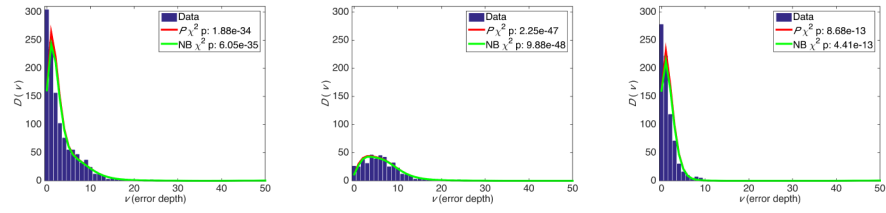


Fig. 3. Error depth distribution in ultra-deep sequencing of a *TP53* locus at 100,000× for all variants (left), transitions (middle), and transversions (right).

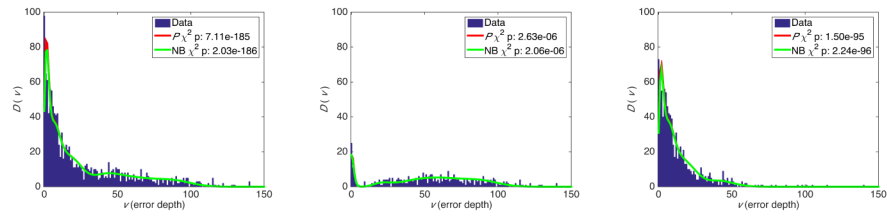


Fig. 4. Error depth distribution in ultra-deep sequencing of a $TP53$ locus at 1,000,000 \times for all variants (left), transitions (middle), and transversions (right).

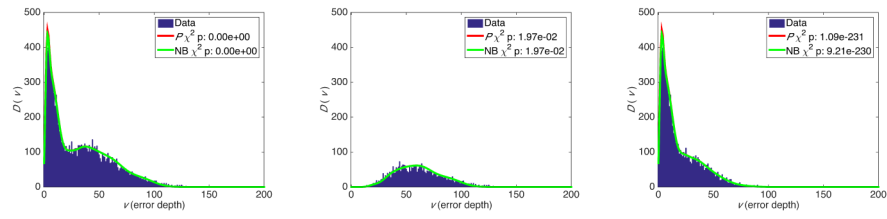


Fig. 5. Error depth distribution in ultra-deep sequencing of a *SF3B1* locus at mean 620,000 \times for all variants (left), transitions (middle), and transversions (right).

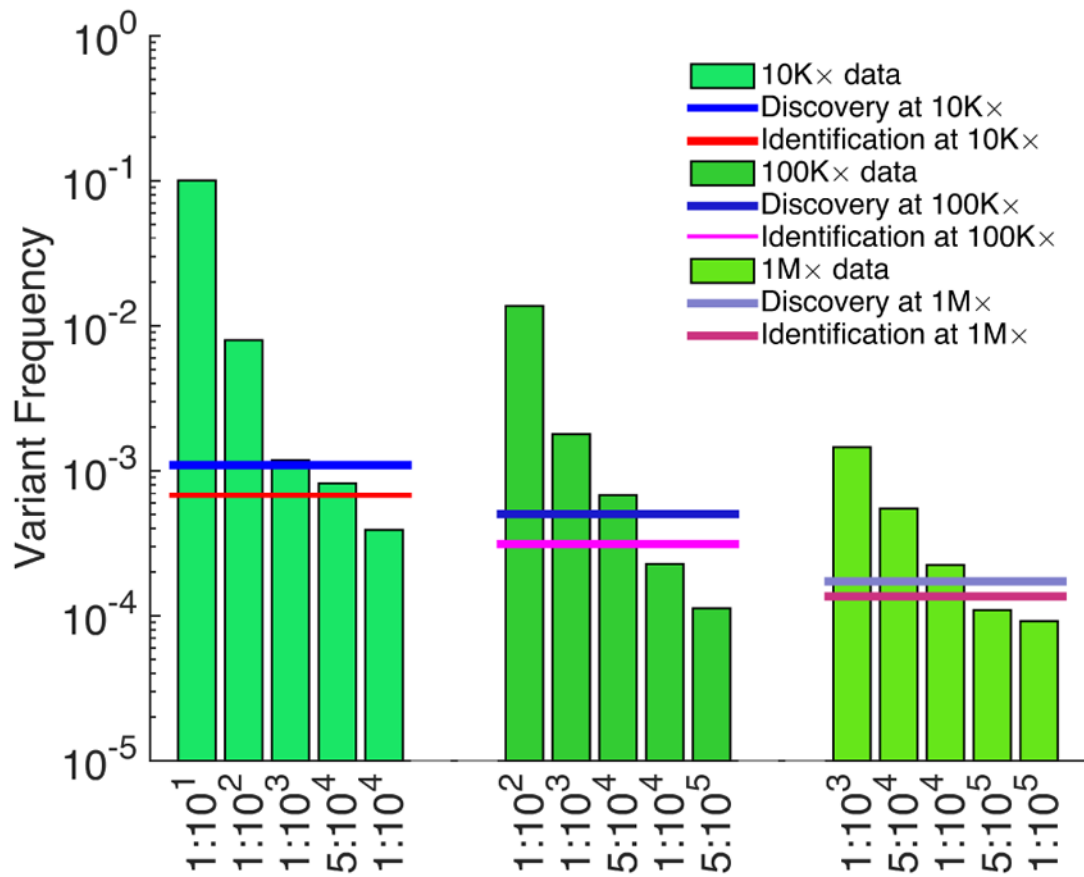


Fig. 6. Sensitivity of detecting *TP53*-Y234C mutation dilutions. Assessing the presence of a variant requires correcting for multiple hypotheses based on the number of sequenced genomic positions (Bonferroni correction). Testing the presence of a discovered variant does not require such a correction; here, significance is set at 0.01.

Table 1

Presence of the *SF3B1*-K700E mutation in undiluted and diluted patient samples are tested against 18 samples that harbor wild-type allele.

Sample	Variant Depth (v)	Total Depth	Variant Frequency	Cumulative P	FDR	Cumulative NB
Control	64	711703	0.00009	9.48E-01	9.48E-01	8.66E-01
Control	62	642586	0.00010	8.45E-01	9.48E-01	6.85E-01
Control	74	717154	0.00010	7.02E-01	9.37E-01	5.09E-01
Control	94	630510	0.00015	2.95E-03	9.43E-03	5.00E-04
Control	56	505857	0.00011	4.68E-01	7.89E-01	2.60E-01
Control	61	509147	0.00012	2.49E-01	5.69E-01	1.07E-01
Control	88	699082	0.00013	1.12E-01	2.98E-01	4.12E-02
Control	75	749932	0.00010	7.91E-01	9.48E-01	6.22E-01
Control	62	657036	0.00009	8.84E-01	9.48E-01	7.47E-01
Control	56	581178	0.00010	8.34E-01	9.48E-01	6.63E-01
Control	81	731934	0.00011	4.75E-01	7.89E-01	2.85E-01
Control	70	636485	0.00011	4.93E-01	7.89E-01	2.89E-01
Control	40	452271	0.00009	9.15E-01	9.48E-01	7.92E-01
Control	59	511932	0.00012	3.51E-01	7.03E-01	1.72E-01
Control	46	518211	0.00009	9.27E-01	9.48E-01	8.15E-01
Control	80	714670	0.00011	4.35E-01	7.89E-01	2.50E-01
Control	85	736865	0.00012	3.33E-01	7.03E-01	1.75E-01
Control	74	691495	0.00011	5.87E-01	8.53E-01	3.82E-01
CLL	281058	630750	0.44559	0.00E+00	0.00E+00	0.00E+00
CLL 1:1000 dilution	2678	440301	0.00608	0.00E+00	0.00E+00	0.00E+00