

(51) International Patent Classification:
G06F 19/22 (2011.01)(21) International Application Number:
PCT/US2018/038342(22) International Filing Date:
19 June 2018 (19.06.2018)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:
62/522,605 20 June 2017 (20.06.2017) US(71) Applicant: **ILLUMINA, INC.** [US/US]; 5200 Illumina Way, San Diego, California 92122 (US).(72) Inventors: **LI, Yong**; 5200 Illumina Way., San Diego, California 92122 (US). **BRUAND, Jocelyne**; 5200 Illumina Way, San Diego, California 92122 (US). **KELLEY, Ryan**; 5200 Illumina Way, San Diego, California 92122 (US). **LEE, Chih**; 5200 Illumina Way, San Diego, California 92122 (US). **SCHEFFLER, Konrad**; 5200 Illumina Way, San Diego, California 92122 (US).(74) Agent: **OSTLER, Trenton B.**; 5200 Illumina Way, San Diego, California 92122 (US).

(81) Designated States (unless otherwise indicated, for every kind of national protection available): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JO, JP, KE, KG, KH, KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.

(84) Designated States (unless otherwise indicated, for every kind of regional protection available): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, RU, TJ, TM), European (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

(54) Title: METHODS AND SYSTEMS FOR DECOMPOSITION AND QUANTIFICATION OF DNA MIXTURES FROM MULTIPLE CONTRIBUTORS OF KNOWN OR UNKNOWN GENOTYPES

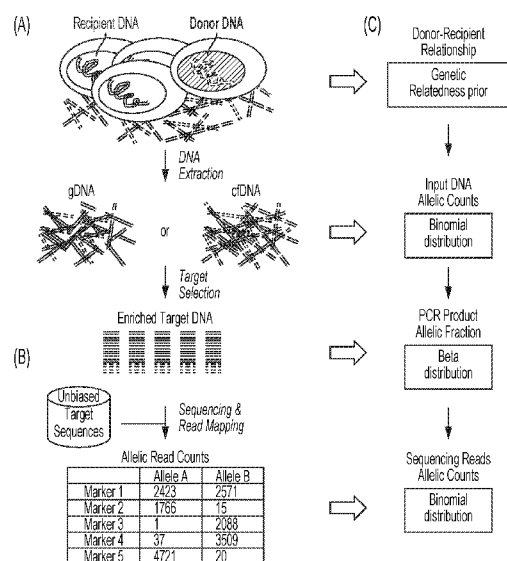


Figure 1

(57) Abstract: Methods and systems are provided for quantifying and deconvolving nucleic acid mixture samples including nucleic acid of one or more contributors having known or unknown genomes. The methods and systems provided herein implement processes that use Bayesian probabilistic modeling techniques to determine the abundances and confidence intervals of genetically distinct contributors in a chimerism sample, thereby improving specificity, accuracy and sensitivity, and greatly expanded the application scope over conventional methods.

Declarations under Rule 4.17:

- *as to applicant's entitlement to apply for and be granted a patent (Rule 4.17(ii))*

Published:

- *with international search report (Art. 21(3))*

**METHODS AND SYSTEMS FOR DECOMPOSITION AND QUANTIFICATION OF
DNA MIXTURES FROM MULTIPLE CONTRIBUTORS OF KNOWN OR UNKNOWN
GENOTYPES**

CROSS REFERENCE TO RELATED APPLICATIONS

[0001] This application claims benefits under 35 U.S.C. § 119(e) to U.S. Provisional Patent Application No. 62/522,605, entitled: METHODS FOR ACCURATE COMPUTATIONAL DECOMPOSITION OF DNA MIXTURES FROM CONTRIBUTORS OF UNKNOWN GENOTYPES, filed June 20, 2017, which is herein incorporated by reference in its entirety for all purposes.

BACKGROUND

[0002] Sequencing data from a nucleic acid (e.g., DNA or RNA) mixture of closely related genomes is frequently found in research as well as clinical settings, and quantifying the mixing contributors has been a challenge when the original genomes are unknown. For example, in the context of microbiology and metagenomics, researchers and clinicians may need to quantify closely related bacterial strains of the same species in an environmental sample. In the setting of forensics, law enforcement personnel may need to quantify as well as identify human individuals from a blood sample containing DNA of multiple individuals. In the setting of biomedical research, scientists may need to determine the purity or extend of contamination in a cell or DNA sample.

[0003] Another application is Next Generation Sequencing (NGS) coupled liquid biopsy. NGS-coupled liquid biopsy is an emerging diagnosis strategy with potential applications in various clinical settings. In the context of organ or tissue transplant, NGS-coupled liquid biopsy provides a non-invasive approach for monitoring the health of allogeneic graft by quantifying the amount of allogeneic DNA in recipient blood. In some applications, the donor and recipient genomes are unknown or partially unknown.

[0004] The term chimera has been used in modern medicine to describe individuals containing cell populations originated from different individuals. The state of chimerism may occur spontaneously through inheritance, but is more frequently produced artificially via transplantation, transfusion, or sample contamination.

[0005] Chimerism leaves informative signals in different DNA types depending on the type of transplant. For bone marrow and hematopoietic stem cell transplants, blood genomic DNA (gDNA) collected post-transplant will have varying levels of chimerism depending on the engraftment state of the transplant. For solid organ transplants, chimerism signals can be seen in the blood cell-free DNA (cfDNA). Such signals can be extracted through non-invasive liquid biopsy, as contrast to the invasive tissue biopsy procedure that is the current standard of care for organ transplant monitoring.

[0006] Reproducible and accurate determinations of the relative contributions of donor genomes to a chimerism DNA sample would provide an informative tool for transplant monitoring, allowing researchers and clinicians to non-invasively and objectively measure the changes in dynamics among donor and recipient cells, which reflect the health status of the donor cells and organs. This disclosure introduces novel and improved methods for quantifying the relative contribution of each genome to a chimerism sample.

SUMMARY

[0007] Some implementations presented herein provide computer-implemented methods and systems for quantification and deconvolution of nucleic acid mixture samples including nucleic acid of two or more contributors of unknown genotypes. One aspect of the disclosure relates to methods for quantifying nucleic acid fractions in nucleic acid samples including nucleic acid (e.g., DNA or RNA) of two or more contributors having different genomes. In some implementations, the nucleic acid mixture samples include biological tissues, cells, peripheral blood, saliva, urine, and other biological fluid, as described below. In some applications, the nucleic acid sample includes the nucleic acid of only a single contributor, and the implementations described herein can determine that the single contributor's nucleic acid accounts for 100% of the nucleic acid in the sample. So although the description hereinafter refers to the nucleic acid sample as a nucleic acid mixture sample in some implementations, it is understood that the sample can include a single contributor's nucleic acid, with the contributor's fraction being 100% or 1. Of course, the methods can also be used to quantify a sample including nucleic acid of two or more contributors.

[0008] Because various methods and systems provided herein implement strategies and processes that use probabilistic mixture models and Bayesian inference techniques, the

embodiments provide technological improvements over conventional methods in quantification and deconvolution of nucleic acid (e.g., DNA or RNA) mixture samples. Some implementations provide improved analytical sensitivity and specificity, providing more accurate deconvolution and quantification of nucleic acid mixture samples.

[0009] Some implementations allow accurate quantification of nucleic acid mixture samples with nucleic acid quantities that are too low for conventional methods to accurately quantify. Some implementations allow accurate quantification of 3-10ng of cell free DNA (cfDNA) mixture samples, which cannot be accurately quantified by conventional methods. Some implementations allow application to mixture samples with 3 or more contributors, which conventional methods cannot handle. Some implementations allow applications to mixtures with one or more unknown genomes, which conventional methods cannot handle. Some implementations described herein refer to a DNA sample, but it is understood that the implementations are also applicable to analyzing RNA samples.

[0010] In some embodiments, the method is implemented at a computer system that includes one or more processors and system memory configured to deconvolve and quantify a nucleic acid mixture sample including nucleic acid of two or more contributors.

[0011] Some embodiments provide a method for quantifying a fraction of nucleic acid of a contributor in a nucleic acid mixture sample comprising nucleic acid of the contributor and at least one other contributor. The method involves: (a) extracting nucleic acid molecules from the nucleic acid sample; (b) amplifying the extracted nucleic acid molecules; (c) sequencing the amplified nucleic acid molecules using a nucleic acid sequencer to produce nucleic acid sequence reads; (d) mapping, by the one or more processors, the nucleic acid sequence reads to one or more polymorphism loci on a reference sequence; (e) determining, using the mapped nucleic acid sequence reads and by the one or more processors, allele counts of nucleic acid sequence reads for one or more alleles at the one or more polymorphism loci; and (f) quantifying, using a probabilistic mixture model and by the one or more processors, one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample, wherein using the probabilistic mixture model includes applying a probabilistic mixture model to the allele counts of nucleic acid sequence reads, and wherein the probabilistic mixture model uses probability distributions to model the allele counts of nucleic acid sequence reads at the one or

more polymorphism loci, the probability distributions accounting for errors in the nucleic acid sequence reads.

[0012] In some implementations, the mapping of (d) includes mapping using computer hashing or computer dynamic programming. In some implementations, the quantifying of (f) comprises quantifying using a novel optimization method combining a multi-iteration grid searching and a Broyden–Fletcher–Goldfarb–Shanno (BFGS) - quasi-Newton method. In some implementations, the quantifying of (f) comprises quantifying using an iterative weighted linear regression. These features require computers to perform and are rooted in computer technology.

[0013] In some implementations, the method further includes, determining, using the probabilistic mixture model and by the one or more processors, one or more genotypes of the one or more contributors at the one or more polymorphism loci.

[0014] In some implementations, the method further includes, determining, using the one or more fractions of nucleic acid of the one or more contributors, a risk of one contributor (a donee) rejecting a tissue or an organ transplanted from another contributor (a donor).

[0015] In some implementations, the one or more contributors include two or more contributors.

[0016] In some implementations, the nucleic acid molecules include DNA molecules or RNA molecules.

[0017] In some implementations, the nucleic acid sample includes nucleic acid from zero, one, or more contaminant genomes and one genome of interest.

[0018] In some implementations, the one or more contributors include zero, one, or more donors of a transplant and a donee of the transplant, and wherein the nucleic acid sample includes a sample obtained from the donee.

[0019] In some implementations, the transplant includes an allogeneic or xenogeneic transplant.

[0020] In some implementations, the nucleic acid sample includes a biological sample obtained from the donee.

[0021] In some implementations, the nucleic acid sample includes a biological sample obtained from a cell culture.

[0022] In some implementations, the extracted nucleic acid molecules include cell-free nucleic acid.

[0023] In some implementations, the extracted nucleic acid molecules include cellular DNA.

[0024] In some implementations, the one or more polymorphism loci include one or more biallelic polymorphism loci.

[0025] In some implementations, the one or more alleles at the one or more polymorphism loci include one or more single nucleotide polymorphism (SNP) alleles.

[0026] In some implementations, the probabilistic mixture model uses a single-locus likelihood function to model allele counts at a single polymorphism locus. The single-locus likelihood function includes:

[0027] $M(n_{1i}, n_{2i} | p_{1i}, \theta)$

[0028] n_{1i} is the allele count of allele 1 at locus i , n_{2i} is the allele count of allele 2 at locus i , p_{1i} is an expected fraction of allele 1 at locus i , and θ includes one or more model parameters.

[0029] In some implementations, p_{1i} is modeled as a function of: (i) genotypes of the contributors at locus i , or $g_i = (g_{1i}, \dots, g_{Di})$, which is a vector of copy number of allele 1 at locus i in contributors 1... D ; (ii) read count errors resulting from the sequencing operation in (c), or λ ; and (iii) fractions of nucleic acid of contributors in the nucleic acid sample, or $\beta = (\beta_1, \dots, \beta_D)$, wherein D is the number of contributors. In some implementations, the contributors include two or more contributors, and $p_{1i} = p(g_i, \lambda, \beta) \leftarrow [(1 - \lambda) g_i + \lambda (2 - g_i)] / 2 \cdot \beta$, where \cdot is vector dot product operator.

[0030] In some implementations, the contributors include two contributors, and p_{1i} is obtained using the p_{1i}' values in Table 3.

[0031] In some implementations, zero, one or more genotypes of the contributors are unknown. In some implementations, (f) includes marginalizing over a plurality of possible combinations of genotypes to enumerate the probability parameter p_{1i} . In some implementations, the method further includes determining a genotype configuration at each of the one or more polymorphism loci, the genotype configuration including two alleles for each of the one or more contributors. In some implementations, the single-locus likelihood function include a first binomial distribution. In some implementations, the first binomial distribution is expressed as follows:

[0032] $n_{1i} \sim BN(n_i, p_{1i})$

[0033] n_{1i} is an allele count of nucleic acid sequence reads for allele 1 at locus i ; and n_i is a total read count at locus i , which equals to a total genome copy numbers n . In some implementations, (f) includes maximizing a multiple-loci likelihood function calculated from a plurality of single-locus likelihood functions.

[0034] In some implementations, (f) includes: calculating a plurality of multiple-loci likelihood values using a plurality of potential fraction values and a multiple-loci likelihood function of the allele counts of nucleic acid sequence reads determined in (e); identifying one or more potential fraction values associated with a maximum multiple-loci likelihood value; and quantifying the one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample as the identified potential fraction value.

[0035] In some implementations, the multiple-loci likelihood function includes:

[0036]
$$L(\beta, \theta, \lambda, \pi; n_1, n_2) = \prod_i [\sum g_i M(n_{1i}, n_{2i} | p(g_i, \lambda, \beta), \theta) \cdot P(g_i | \pi)]$$

[0037] $L(\beta, \theta, \lambda, \pi; n_1, n_2)$ is the likelihood of observing allele count vectors n_1 and n_2 for alleles 1 and 2; $p(g_i, \lambda, \beta)$ is the expected fraction or probability of observing allele 1 at locus i based on the contributors' genotypes g_i at locus i ; $P(g_i | \pi)$ is the prior probability of observing the genotypes g_i at locus i given a population allele frequency (π); and $\sum g_i$ denotes summing over a plurality of possible combinations of genotypes of the contributors.

[0038] In some implementations, the multiple-loci likelihood function includes:

[0039]
$$L(\beta, \lambda, \pi; n_1, n_2) = \prod_i [\sum g_i BN(n_{1i} | n_i, p(g_i, \lambda, \beta)) \cdot P(g_i | \pi)]$$

[0040] In some implementations, the contributors include two contributors and the likelihood function includes:

[0041]
$$L(\beta, \lambda, \pi; n_1, n_2) = \prod_i \sum_{g_{1i} g_{2i}} BN(n_{1i} | n_i, p_{1i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)$$

[0042] $L(\beta, \lambda, \pi; n_1, n_2)$ is the likelihood of observing allele count vectors n_1 to n_2 for alleles 1 and 2 given parameters β and π ; $p_{1i}(g_{1i}, g_{2i}, \lambda, \beta)$ is a probability parameter, taken as p_1' from Table 3, indicating a probability of allele 1 at locus i based on the two contributors' genotypes (g_{1i}, g_{2i}); and $P(g_{1i}, g_{2i} | \pi)$ is a prior joint probability of observing the two contributors' genotypes given a population allele frequency (π).

[0043] In some implementations, the prior joint probability is calculated using marginal distributions $P(g_{1i} | \pi)$ and $P(g_{2i} | \pi)$ that satisfy the Hardy-Weinberg equilibrium.

[0044] In some implementations, the prior joint probability is calculated using genetic relationship between the two contributors.

[0045] In some implementations, the probabilistic mixture model accounts for nucleic acid molecule copy number errors resulting from extracting the nucleic acid molecules performed in (a), as well as the read count errors resulting from the sequencing operation in (c). In some implementations, the probabilistic mixture model uses a second binomial distribution to model allele counts of the extracted nucleic acid molecules for alleles at the one or more polymorphism loci. In some implementations, the second binomial distribution is expressed as follows:

[0046] $n_{li}'' \sim BN(n_i'', p_{li})$

[0047] n_{li}'' is an allele count of extracted nucleic acid molecules for allele l at locus i ; n_i'' is a total nucleic acid molecule count at locus i ; and p_{li} is a probability parameter indicating the probability of allele l at locus i .

[0048] In some implementations, the first binomial distribution is conditioned on an allele fraction n_{li}''/n_i'' . In some implementations, the first binomial distribution is re-parameterized as follows:

[0049] $n_{li} \sim BN(n_i, n_{li}''/n_i'')$

[0050] n_{li} is an allele count of nucleic acid sequence reads for allele l at locus i ; n_i'' is a total number of nucleic acid molecules at locus i , which equals to a total genome copy numbers n'' ; n_i is a total read count at locus i ; and n_{li}'' is a number of extracted nucleic acid molecules for allele l at locus i .

[0051] In some implementations, the probabilistic mixture model uses a first beta distribution to approximate a distribution of n_{li}''/n'' . In some implementations, the first beta distribution has a mean and a variance that match a mean and a variance of the second binomial distribution. In some implementations, locus i is modeled as biallelic and the first beta distribution is expressed as follows:

[0052] $n_{li}''/n'' \sim Beta((n''-1)p_{1i}, (n''-1)p_{2i})$

[0053] p_{1i} is a probability parameter indicating the probability of a first allele at locus i ; and p_{2i} is a probability parameter indicating the probability of a second allele at locus i .

[0054] In some implementations, (f) includes combining the first binomial distribution, modeling sequencing read counts, and the first beta distribution, modeling extracted nucleic acid molecule number, to obtain the single-locus likelihood function of n_{li} that follows a first beta-binomial distribution. In some implementations, the first beta-binomial distribution has the form:

$n_{1i} \sim BB(n_i, (n''-1) \cdot p_{1i}, (n''-1) \cdot p_{2i})$, or an alternative approximation: $n_{1i} \sim BB(n_i, n'' \cdot p_{1i}, n'' \cdot p_{2i})$. In some implementations, the multiple-loci likelihood function includes:

[0055] $L(\beta, n'', \lambda, \pi; n_1, n_2) = \prod_i [\sum_{g_i} BB(n_{1i} | n_i, (n''-1) \cdot p_{1i}, (n''-1) \cdot p_{2i}) \cdot P(g_i | \pi)]$

[0056] $L(\beta, n'', \lambda, \pi; n_1, n_2)$ is the likelihood of observing allele count vectors n_1 and n_2 for alleles 1 and 2 at all loci, and $p_{1i} = p(g_i, \lambda, \beta)$, $p_{2i} = 1 - p_{1i}$.

[0057] In some implementations, the contributors include two contributors, and the multiple-loci likelihood function includes:

[0058] $L(\beta, n'', \lambda, \pi; n_1, n_2) = \prod_i \sum_{g_{1i} g_{2i}} BB(n_{1i}, n_{2i} | n_i, (n''-1) \cdot p_{1i}(g_{1i}, g_{2i}, \lambda, \beta), (n''-1) \cdot p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)$.

[0059] $L(\beta, n'', \lambda, \pi; n_1, n_2)$ is the likelihood of observing an allele count vector for the first allele of all loci (n_1) and an allele count vector for the second allele of all loci (n_2) given parameters β, n'', λ , and π ; $p_{1i}(g_{1i}, g_{2i}, \lambda, \beta)$ is a probability parameter, taken as p_1' from Table 3, indicating a probability of allele 1 at locus i based on the two contributors' genotypes (g_{1i}, g_{2i}); $p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)$ is a probability parameter, taken as p_2' from Table 3, indicating a probability of allele 2 at locus i based on the two contributors' genotypes (g_{1i}, g_{2i}); and $P(g_{1i}, g_{2i} | \pi)$ is a prior joint probability of observing the first contributor's genotype for the first allele (g_{1i}) and the second contributor's genotype for the first allele (g_{2i}) at locus i given a population allele frequency (π).

[0060] In some implementations, (f) includes estimating the total extracted genome copy number n'' from a mass of the extracted nucleic acid molecules. In some implementations, the estimated total extracted genome copy number n'' is adjusted according to fragment size of the extracted nucleic acid molecules.

[0061] In some implementations, the probabilistic mixture model accounts for nucleic acid molecule number errors resulting from amplifying the nucleic acid molecules performed in (b), as well as the read count errors resulting from the sequencing operation in (c). In some implementations, the amplification process of (b) is modeled as follows:

[0062] $x_{t+1} = x_t + y_{t+1}$

[0063] x_{t+1} is the nucleic acid copies of a given allele after cycle $t+1$ of amplification; x_t is the nucleic acid copies of a given allele after cycle t of amplification; y_{t+1} is the new copies generated at cycle $t+1$, and it follows a binomial distribution $y_{t+1} \sim BN(x_t, r_{t+1})$; and r_{t+1} is the amplification rate for cycle $t+1$.

[0064] In some implementations, the probabilistic mixture model uses a second beta distribution to model allele fractions of the amplified nucleic acid molecules for alleles at the one or more polymorphism loci.

[0065] In some implementations, locus i is biallelic and the second beta distribution is expressed as follows:

$$[0066] \quad n_{1i}' / (n_{1i}' + n_{2i}') \sim \text{Beta}(n'' \cdot \rho_i \cdot p_{1i}, n'' \cdot \rho_i \cdot p_{2i})$$

[0067] n_{1i}' is an allele count of amplified nucleic acid molecules for a first allele at locus i ; n_{2i}' is an allele count of amplified nucleic acid molecules for a second allele at locus i ; n'' is a total nucleic acid molecule count at any locus; ρ_i is a constant related to an average amplification rate r ; p_{1i} is the probability of the first allele at locus i ; and p_{2i} is the probability of the second allele at locus i . In some implementations, ρ_i is $(1+r)/(1-r) / [1-(1+r)^{-i}]$, and r is the average amplification rate per cycle. In some implementations, ρ_i is approximated as $(1+r)/(1-r)$.

[0068] In some implementations, (f) includes combining the first binomial distribution and the second beta distribution to obtain the single-locus likelihood function for n_{1i} that follows a second beta-binomial distribution. In some implementations, the second beta-binomial distribution has the form:

$$[0069] \quad n_{1i} \sim \text{BB}(n_i, n'' \cdot \rho_i \cdot p_{1i}, n'' \cdot \rho_i \cdot p_{2i})$$

[0070] n_{1i} is an allele count of nucleic acid sequence reads for the first allele at locus i ; p_{1i} is a probability parameter indicating the probability of a first allele at locus i ; and p_{2i} is a probability parameter indicating the probability of a second allele at locus i .

[0071] In some implementations, (f) includes, by assuming the one or more polymorphism loci have a same amplification rate, re-parameterizing the second beta-binomial distribution as: $n_{1i} \sim \text{BB}(n_i, n'' \cdot (1+r)/(1-r) \cdot p_{1i}, n'' \cdot (1+r)/(1-r) \cdot p_{2i})$, where r is an amplification rate. In some implementations, the multiple-loci likelihood function includes:

$$[0072] \quad L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i [\sum_{g_i} \text{BB}(n_{1i} | n_i, n'' \cdot (1+r)/(1-r) \cdot p_{1i}, n'' \cdot (1+r)/(1-r) \cdot p_{2i}) \cdot P(g_i | \pi)]$$

[0073] In some implementations, the contributors include two contributors and the multiple-loci likelihood function includes:

$$[0074] \quad L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i \sum_{g_{1i} g_{2i}} [\text{BB}(n_{1i} | n_i, n'' \cdot (1+r)/(1-r) \cdot p_{1i}(g_{1i}, g_{2i}, \lambda, \beta), n'' \cdot (1+r)/(1-r) \cdot p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)]$$

[0075] $L(\beta, n'', r, \lambda, \pi; n_1, n_2)$ is the likelihood of observing an allele count vector for the first allele of all loci (n_1) and an allele count vector for the second allele of all loci (n_2) given parameters β, n'', r, λ , and π .

[0076] In some implementations, (f) includes, by defining a relative amplification rate of each polymorphism locus to be proportional to a total reads of the locus, re-parameterizing the second beta-binomial distribution as:

[0077] $n_{1i} \sim BB(n_i, c' \cdot n_i \cdot p_{1i}, c' \cdot n_i \cdot p_{2i})$

[0078] c' is a parameter to be optimized; and n_i is the total reads at locus i .

[0079] In some implementations, the multiple-loci likelihood function includes:

[0080] $L(\beta, n'', c', \lambda, \pi; n_1, n_2) = \prod_i [\sum g_i BB(n_{1i} | n_i, c' \cdot n_i \cdot p_{1i}, c' \cdot n_i \cdot p_{2i}) \cdot P(g_i | \pi)]$

[0081] In some implementations, the probabilistic mixture model accounts for nucleic acid molecule number errors resulting from extracting the nucleic acid molecules performed in (a) and amplifying the nucleic acid molecules performed in (b), as well as the read count errors resulting from the sequencing operation in (c). In some implementations, the probabilistic mixture model uses a third beta distribution to model allele fractions of the amplified nucleic acid molecules for alleles at the one or more polymorphism loci, accounting for the sampling errors resulting from extracting the nucleic acid molecules performed in (a) and amplifying the nucleic acid molecules performed in (b). In some implementations, locus i is biallelic and the third beta distribution has the form of:

[0082] $n_{1i}' / (n_{1i}' + n_{2i}') \sim \text{Beta}(n'' \cdot (1 + r_i) / 2 \cdot p_{1i}, n'' \cdot (1 + r_i) / 2 \cdot p_{2i})$

[0083] n_{1i}' is an allele count of amplified nucleic acid molecules for a first allele at locus i ; n_{2i}' is an allele count of amplified nucleic acid molecules for a second allele at locus i ; n'' is a total nucleic acid molecule count; r_i is the average amplification rate for locus i ; p_{1i} is the probability of the first allele at locus i ; and p_{2i} is a probability of the second allele at locus i . In some implementations, (f) includes combining the first binomial distribution and the third beta distribution to obtain the single-locus likelihood function of n_{1i} that follows a third beta-binomial distribution.

[0084] In some implementations, the third beta-binomial distribution has the form:

[0085] $n_{1i} \sim BB(n_i, n'' \cdot (1 + r_i) / 2 \cdot p_{1i}, n'' \cdot (1 + r_i) / 2 \cdot p_{2i})$

[0086] r_i is an amplification rate.

[0087] In some implementations, the multiple-loci likelihood function includes:

[0088] $L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i [\sum_{g_i} BB(n_{1i} | n_i, n'' \cdot (1+r)/2 \cdot p_{1i}, n'' \cdot (1+r)/2 \cdot p_{2i}) \cdot P(g_i | \pi)]$

[0089] In some implementations, the contributors include two contributors, and wherein the multiple-loci likelihood function includes:

[0090] $L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i \sum_{g_{1i} g_{2i}} BB(n_{1i} | n_i, n'' \cdot (1+r)/2 \cdot p_{1i}(g_{1i}, g_{2i}, \lambda, \beta), n'' \cdot (1+r)/2 \cdot p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)$

[0091] $L(n_1, n_2 | \beta, n'', r, \lambda, \pi)$ is the likelihood of observing allele counts for the first allele vector n_1 and an allele count for the second allele vector n_2 given parameters β, n'', r, λ , and π .

[0092] In some implementations, the method further includes: (g) estimating one or more confidence intervals of the one or more fractions of nucleic acid of the one or more contributors using the hessian matrix of the log-likelihood using numerical differentiation.

[0093] In some implementations, the mapping of (d) includes identifying, by the one or more processors using computer hashing and computer dynamic programming, reads among the nucleic acid sequence reads matching any sequence of a plurality of unbiased target sequences, wherein the plurality of unbiased target sequences includes sub-sequences of the reference sequence and sequences that differ from the subsequences by a single nucleotide. In some implementations, the plurality of unbiased target sequences includes five categories of sequences encompassing each polymorphic site of a plurality of polymorphic sites: (i) a reference target sequence that is a sub-sequence of the reference sequence, the reference target sequence having a reference allele with a reference nucleotide at the polymorphic site; (ii) alternative target sequences each having an alternative allele with an alternative nucleotide at the polymorphic site, the alternative nucleotide being different from the reference nucleotide; (iii) mutated reference target sequences including all possible sequences that each differ from the reference target sequence by only one nucleotide at a site that is not the polymorphic site; (iv) mutated alternative target sequences including all possible sequences that each differ from an alternative target sequence by only one nucleotide at a site that is not the polymorphic site; and (v) unexpected allele target sequences each having an unexpected allele different from the reference allele and the alternative allele, and each having a sequence different from the previous four categories of sequences.

[0094] In some implementations, the method further includes estimating a sequencing error rate λ at the variant site base on a frequency of observing the unexpected allele target sequences of (v). In some implementations, (e) includes using the identified reads and their matching unbiased target sequences to determine allele counts of the nucleic acid sequence reads for the alleles at the one or more polymorphism loci. In some implementations, the plurality of unbiased target sequences includes sequences that are truncated to have the same length as the nucleic acid sequence reads. In some implementations, the plurality of unbiased target sequences includes sequences stored in one or more hash tables, and the reads are identified using the hash tables.

[0095] The disclosed embodiments also provide a computer program product including a non-transitory computer readable medium on which is provided program instructions for performing the recited operations and other computational operations described herein.

[0096] Some embodiments provide a system for quantifying a fraction of nucleic acid of a contributor in a nucleic acid mixture sample comprising nucleic acid of the contributor and at least one other contributor. The system includes a sequencer for receiving nucleic acids from the test sample providing nucleic acid sequence information from the sample, a processor; and one or more computer-readable storage media having stored thereon instructions for execution on the processor to deconvolve and quantify DNA mixture samples using the method recited herein.

[0097] Another aspect of the disclosure provides a system quantifying a nucleic acid sample including nucleic acid of one or more contributors. The system includes: (a) a sequencer configured to (i) receive nucleic acid molecules extracted from the nucleic acid sample, (ii) amplify the extracted nucleic acid molecules, and (iii) sequence the amplified nucleic acid molecules under conditions that produce nucleic acid sequence reads; and (b) a computer including one or more processors configured to: map the nucleic acid sequence reads to one or more polymorphism loci on a reference sequence; determine, using the mapped nucleic acid sequence reads, allele counts of nucleic acid sequence reads for one or more alleles at the one or more polymorphism loci; and quantify, using a probabilistic mixture model, one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample. Using the probabilistic mixture model includes applying a probabilistic mixture model to the allele counts of nucleic acid sequence reads, and the probabilistic mixture model uses probability distributions

to model the allele counts of nucleic acid sequence reads at the one or more polymorphism loci, the probability distributions accounting for errors in the nucleic acid sequence reads.

[0098] In some implementations, the system includes a tool for extracting nucleic acid molecules from the nucleic acid sample. In some implementations, the probability distributions include a first binomial distribution as follows:

[0099] $n_{1i} \sim BN(n_i, p_{1i})$.

[00100] n_{1i} is an allele count of nucleic acid sequence reads for allele 1 at locus i ; n_i is a total read count at locus i , which equals to a total genome copy numbers n ; and p_{1i} is a probability parameter indicating the probability of allele 1 at locus i .

[00101] An additional aspect of the disclosure provides a computer program product including a non-transitory machine readable medium storing program code that, when executed by one or more processors of a computer system, causes the computer system to implement a method of quantifying a nucleic acid sample including nucleic acid of one or more contributors, said program code including: code for mapping the nucleic acid sequence reads to one or more polymorphism loci on a reference sequence; code for determining, using the mapped nucleic acid sequence reads, allele counts of nucleic acid sequence reads for one or more alleles at the one or more polymorphism loci; and code for quantifying, using a probabilistic mixture model, one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample. Using the probabilistic mixture model includes applying a probabilistic mixture model to the allele counts of nucleic acid sequence reads, and the probabilistic mixture model uses probability distributions to model the allele counts of nucleic acid sequence reads at the one or more polymorphism loci, the probability distributions accounting for errors in the nucleic acid sequence reads.

[00102] Yet another aspect of the disclosure provides a method, implemented at a computer system that includes one or more processors and system memory, of quantifying a nucleic acid sample including nucleic acid of one or more contributors. The method includes: (a) receiving, by the one or more processors, nucleic acid sequence reads obtained from the nucleic acid sample; (b) mapping, by the one or more processors, using computer hashing and computer dynamic programming, the nucleic acid sequence reads to one or more polymorphism loci on a reference sequence; (c) determining, using the mapped nucleic acid sequence reads and by the

one or more processors, allele counts of nucleic acid sequence reads for one or more alleles at the one or more polymorphism loci; and (d) quantifying, using a probabilistic mixture model and by the one or more processors, one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample and confidence of the fractions. Using the probabilistic mixture model includes applying a probabilistic mixture model to the allele counts of nucleic acid sequence reads. The probabilistic mixture model uses probability distributions to model the allele counts of nucleic acid sequence reads at the one or more polymorphism loci, the probability distributions accounting for errors in the mapped nucleic acid sequence reads. The quantifying employs (i) a computer optimization method combining multi-iteration grid searching and a BFGS - quasi-Newton method, or an iterative weighted linear regression, and (ii) a numerical differentiation method.

[00103] Although the examples herein concern humans and the language is primarily directed to human concerns, the concepts described herein are applicable to genomes from any plant or animal. These and other objects and features of the present disclosure will become more fully apparent from the following description and appended claims, or may be learned by the practice of the disclosure as set forth hereinafter.

INCORPORATION BY REFERENCE

[00104] All patents, patent applications, and other publications, including all sequences disclosed within these references, referred to herein are expressly incorporated herein by reference, to the same extent as if each individual publication, patent or patent application was specifically and individually indicated to be incorporated by reference. All documents cited are, in relevant part, incorporated herein by reference in their entireties for the purposes indicated by the context of their citation herein. However, the citation of any document is not to be construed as an admission that it is prior art with respect to the present disclosure.

BRIEF DESCRIPTION OF THE DRAWINGS

[00105] Figures 1A-1C show an overview of a method and statistical model designed for contributor DNA quantification.

[00106] Figure 2A shows a block diagram illustrating a process for quantifying one or more fractions of nucleic acid (e.g., DNA or RNA) of one or more contributors in the nucleic acid sample.

[00107] Figure 2B shows a block diagram illustrating various components of a probabilistic mixture model.

[00108] Figure 2C schematically illustrates sequencing errors that convert one allele to another allele and true alleles to unexpected alleles.

[00109] Figure 3 shows a block diagram illustrating a process for evaluating a nucleic acid sample including nucleic acid of one or more contributors.

[00110] Figure 4 shows block diagram of a typical computer system that can serve as a computational apparatus according to certain embodiments.

[00111] Figure 5 shows one implementation of a dispersed system for producing a call or diagnosis from a test sample.

[00112] Figure 6 shows options for performing various operations of some implementations at distinct locations.

[00113] Figure 7 shows the performance of disclosed and baseline methods each under different choices of cfDNA length parameter.

[00114] Figure 8 shows analytical accuracy of some implementations in another format.

[00115] Figure 9 shows the coefficient of variance (CV) of 16 conditions for determining .limit of quantification (LOQ) for some implementations.

DETAILED DESCRIPTION

Definitions

[00116] Unless otherwise indicated, the practice of the method and system disclosed herein involves conventional techniques and apparatus commonly used in molecular biology, microbiology, protein purification, protein engineering, protein and DNA sequencing, and recombinant DNA fields, which are within the skill of the art. Such techniques and apparatus are known to those of skill in the art and are described in numerous texts and reference works (See

e.g., Sambrook et al., “Molecular Cloning: A Laboratory Manual,” Third Edition (Cold Spring Harbor), [2001]); and Ausubel et al., “Current Protocols in Molecular Biology” [1987]).

[00117] Numeric ranges are inclusive of the numbers defining the range. It is intended that every maximum numerical limitation given throughout this specification includes every lower numerical limitation, as if such lower numerical limitations were expressly written herein. Every minimum numerical limitation given throughout this specification will include every higher numerical limitation, as if such higher numerical limitations were expressly written herein. Every numerical range given throughout this specification will include every narrower numerical range that falls within such broader numerical range, as if such narrower numerical ranges were all expressly written herein.

[00118] The headings provided herein are not intended to limit the disclosure.

[00119] Unless defined otherwise herein, all technical and scientific terms used herein have the same meaning as commonly understood by one of ordinary skill in the art. Various scientific dictionaries that include the terms included herein are well known and available to those in the art. Although any methods and materials similar or equivalent to those described herein find use in the practice or testing of the embodiments disclosed herein, some methods and materials are described.

[00120] The terms defined immediately below are more fully described by reference to the Specification as a whole. It is to be understood that this disclosure is not limited to the particular methodology, protocols, and reagents described, as these may vary, depending upon the context they are used by those of skill in the art. As used herein, the singular terms “a,” “an,” and “the” include the plural reference unless the context clearly indicates otherwise.

[00121] Unless otherwise indicated, nucleic acids are written left to right in 5' to 3' orientation and amino acid sequences are written left to right in amino to carboxy orientation, respectively.

[00122] The term “chimerism sample” is used herein to refer to a sample believed to contain DNA of two or more genomes. Chimerism analysis is used herein to refer to the biological and chemical processing of a chimerism sample and/or the quantification of the nucleic acid of two or more organisms in the chimera sample. In some implementations, a

chimerism analysis also determines some or all of the sequence information of the genomes of the two or more organisms.

[00123] The term donor DNA (dDNA) refers to DNA molecules originating from cells of a donor of a transplant. In various implementations, the dDNA is found in a sample obtained from a donee who received a transplanted tissue/organ from the donor.

[00124] Circulating cell-free DNA or simply cell-free DNA (cfDNA) are DNA fragments that are not confined within cells and are freely circulating in the bloodstream or other bodily fluids. It is known that cfDNA have different origins, in some cases from donor tissue DNA circulating in a donee's blood, in some cases from tumor cells or tumor affected cells, in other cases from fetal DNA circulating in maternal blood. In general, cfDNA are fragmented and include only a small portion of a genome, which may be different from the genome of the individual from which the cfDNA is obtained.

[00125] The term non-circulating genomic DNA (gDNA) or cellular DNA are used to refer to DNA molecules that are confined in cells and often include a complete genome.

[00126] The term "allele count" refers to the count or number of sequence reads of a particular allele. In some implementations, it can be determined by mapping reads to a location in a reference genome, and counting the reads that include an allele sequence and are mapped to the reference genome.

[00127] A beta distribution is a family of continuous probability distributions defined on the interval $[0, 1]$ parameterized by two positive shape parameters, denoted by, e.g., α and β , that appear as exponents of the random variable and control the shape of the distribution. The beta distribution has been applied to model the behavior of random variables limited to intervals of finite length in a wide variety of disciplines. In Bayesian inference, the beta distribution is the conjugate prior probability distribution for the Bernoulli, binomial, negative binomial and geometric distributions. For example, the beta distribution can be used in Bayesian analysis to describe initial knowledge concerning probability of success. If the random variable X follows the beta distribution, the random variable X is written as $X \sim \text{Beta}(\alpha, \beta)$.

[00128] A binomial distribution is a discrete probability distribution of the number of successes in a sequence of n independent experiments, each asking a yes-no question, and each

with its own Boolean-valued outcome: a random variable containing single bit of information: positive (with probability p) or negative (with probability $q = 1 - p$). For a single trial, i.e., $n = 1$, the binomial distribution is a Bernoulli distribution. The binomial distribution is frequently used to model the number of successes in a sample of size n drawn with replacement from a population of size N . If the random variable X follows the binomial distribution with parameters $n \in \mathbb{N}$ and $p \in [0,1]$, the random variable X is written as $X \sim B(n, p)$.

[00129] Poisson distribution, denoted as $\text{Pois}()$ herein, is a discrete probability distribution that expresses the probability of a given number of events occurring in a fixed interval of time and/or space if these events occur with a known average rate and independently of the time since the last event. The Poisson distribution can also be used for the number of events in other specified intervals such as distance, area or volume. The probability of observing k events in an interval according to a Poisson distribution is given by the equation:

$$P(k \text{ events in interval}) = \frac{\lambda^k e^{-\lambda}}{k!}$$

where λ is the average number of events in an interval or an event rate, also called the rate parameter e is 2.71828, Euler's number, or the base of the natural logarithms, k takes values 0, 1, 2, ..., and $k!$ is the factorial of k .

[00130] Gamma distribution is a two-parameter family of continuous probability distributions. There are three different parametrizations in common use: with a shape parameter k and a scale parameter θ ; with a shape parameter $\alpha = k$ and an inverse scale parameter $\beta = 1/\theta$, called a rate parameter; or with a shape parameter k and a mean parameter $\mu = k/\beta$. In each of these three forms, both parameters are positive real numbers. The gamma distribution is the maximum entropy probability distribution for a random variable X for which $E[X] = k\theta = \alpha/\beta$ is fixed and greater than zero, and $E[\ln(X)] = \psi(k) + \ln(\theta) = \psi(\alpha) - \ln(\beta)$ is fixed (ψ is the digamma function).

[00131] Polymorphism and genetic polymorphism are used interchangeably herein to refer to the occurrence in the same population of two or more alleles at one genomic locus, each with appreciable frequency.

[00132] Polymorphism site and polymorphic site are used interchangeably herein to refer to a locus on a genome at which two or more alleles reside. In some implementations, it is used to refer to a single nucleotide variation with two alleles of different bases.

[00133] Allele frequency or gene frequency is the frequency of an allele of a gene (or a variant of the gene) relative to other alleles of the gene, which can be expressed as a fraction or percentage. An allele frequency is often associated with a particular genomic locus, because a gene is often located at with one or more locus. However, an allele frequency as used herein can also be associated with a size-based bin of DNA fragments. In this sense, DNA fragments such as cfDNA containing an allele are assigned to different size-based bins. The frequency of the allele in a size-based bin relative to the frequency of other alleles is an allele frequency.

[00134] The term “parameter” herein refers to a numerical value that characterizes a property of a system such as a physical feature whose value or other characteristic has an impact on a relevant condition such as a sample or DNA fragments. In some cases, the term parameter is used with reference to a variable that affects the output of a mathematical relation or model, which variable may be an independent variable (i.e., an input to the model) or an intermediate variable based on one or more independent variables. Depending on the scope of a model, an output of one model may become an input of another model, thereby becoming a parameter to the other model.

[00135] The term “plurality” refers to more than one element.

[00136] The term “paired end reads” refers to reads from paired end sequencing that obtains one read from each end of a nucleic acid fragment. Paired end sequencing may involve fragmenting strands of polynucleotides into short sequences called inserts. Fragmentation is optional or unnecessary for relatively short polynucleotides such as cell free DNA molecules.

[00137] The terms “polynucleotide,” “nucleic acid” and “nucleic acid molecules” are used interchangeably and refer to a covalently linked sequence of nucleotides (i.e., ribonucleotides for RNA and deoxyribonucleotides for DNA) in which the 3' position of the pentose of one nucleotide is joined by a phosphodiester group to the 5' position of the pentose of the next. The nucleotides include sequences of any form of nucleic acid, including, but not limited to RNA and DNA molecules such as cfDNA or cellular DNA molecules. The term “polynucleotide” includes, without limitation, single- and double-stranded polynucleotide.

[00138] The term “test sample” herein refers to a sample typically derived from a biological fluid, cell, tissue, organ, or organism, comprising a nucleic acid or a mixture of nucleic acids. Such samples include, but are not limited to sputum/oral fluid, amniotic fluid, blood, a blood fraction, or fine needle biopsy samples (e.g., surgical biopsy, fine needle biopsy, etc.), urine, peritoneal fluid, pleural fluid, and the like. Although the sample is often taken from a human subject (e.g., patient), the assays can be used in samples from any mammal, including, but not limited to dogs, cats, horses, goats, sheep, cattle, pigs, etc. The sample may be used directly as obtained from the biological source or following a pretreatment to modify the character of the sample. For example, such pretreatment may include preparing plasma from blood, diluting viscous fluids and so forth. Methods of pretreatment may also involve, but are not limited to, filtration, precipitation, dilution, distillation, mixing, centrifugation, freezing, lyophilization, concentration, amplification, nucleic acid fragmentation, inactivation of interfering components, the addition of reagents, lysing, etc. If such methods of pretreatment are employed with respect to the sample, such pretreatment methods are typically such that the nucleic acid(s) of interest remain in the test sample, sometimes at a concentration proportional to that in an untreated test sample (e.g., namely, a sample that is not subjected to any such pretreatment method(s)). Such “treated” or “processed” samples are still considered to be biological “test” samples with respect to the methods described herein.

[00139] The term “Next Generation Sequencing (NGS)” herein refers to sequencing methods that allow for massively parallel sequencing of clonally amplified molecules and of single nucleic acid molecules. Non-limiting examples of NGS include sequencing-by-synthesis using reversible dye terminators, and sequencing-by-ligation.

[00140] The term “read” refers to a sequence obtained from a portion of a nucleic acid sample. Typically, though not necessarily, a read represents a short sequence of contiguous base pairs in the sample. The read may be represented symbolically by the base pair sequence (in A, T, C, or G) of the sample portion. It may be stored in a memory device and processed as appropriate to determine whether it matches a reference sequence or meets other criteria. A read may be obtained directly from a sequencing apparatus or indirectly from stored sequence information concerning the sample. In some cases, a read is a DNA sequence of sufficient length (e.g., at least about 25 bp) that can be used to identify a larger sequence or region, e.g., that can be aligned and specifically assigned to a chromosome or genomic region or gene.

[00141] The term “genomic read” is used in reference to a read of any segments in the entire genome of an individual.

[00142] As used herein, the terms “aligned,” “alignment,” or “aligning” refer to the process of comparing a read or tag to a reference sequence and thereby determining whether the reference sequence contains the read sequence. If the reference sequence contains the read, the read may be mapped to the reference sequence or, in certain embodiments, to a particular location in the reference sequence. In some cases, alignment simply tells whether or not a read is a member of a particular reference sequence (i.e., whether the read is present or absent in the reference sequence). For example, the alignment of a read to the reference sequence for human chromosome 13 will tell whether the read is present in the reference sequence for chromosome 13. A tool that provides this information may be called a set membership tester. In some cases, an alignment additionally indicates a location in the reference sequence where the read or tag maps to. For example, if the reference sequence is the whole human genome sequence, an alignment may indicate that a read is present on chromosome 13, and may further indicate that the read is on a particular strand and/or site of chromosome 13.

[00143] Aligned reads or tags are one or more sequences that are identified as a match in terms of the order of their nucleic acid molecules to a known sequence from a reference genome. Alignment can be done manually, although it is typically implemented by a computer program, as it would be impossible to align reads in a reasonable time period for implementing the methods disclosed herein. One example of a program for aligning sequences is the Efficient Local Alignment of Nucleotide Data (ELAND) computer program distributed as part of the Illumina Genomics Analysis pipeline. Alternatively, a Bloom filter or similar set membership tester may be employed to align reads to reference genomes. See US Patent Application No. 61/552,374 filed October 27, 2011 which is incorporated herein by reference in its entirety. The matching of a sequence read in aligning can be a 100% sequence match or less than 100% (non-perfect match).

[00144] The term “mapping” used herein refers to specifically assigning a sequence read to a larger sequence, e.g., a reference genome, a subsequence of the larger sequence using alignment or membership assignment.

[00145] As used herein, the term “reference genome” or “reference sequence” refers to any particular known genome sequence, whether partial or complete, of any organism or virus which may be used to reference identified sequences from a subject. For example, a reference genome used for human subjects as well as many other organisms is found at the National Center for Biotechnology Information at ncbi.nlm.nih.gov. A “genome” refers to the complete genetic information of an organism or virus, expressed in nucleic acid sequences.

[00146] In various embodiments, the reference sequence is significantly larger than the reads that are aligned to it. For example, it may be at least about 100 times larger, or at least about 1000 times larger, or at least about 10,000 times larger, or at least about 10^5 times larger, or at least about 10^6 times larger, or at least about 10^7 times larger.

[00147] In one example, the reference sequence is that of a full length human genome. Such sequences may be referred to as genomic reference sequences. In another example, the reference sequence is limited to a specific human chromosome such as chromosome 13. In some embodiments, a reference Y chromosome is the Y chromosome sequence from human genome version hg19. Such sequences may be referred to as chromosome reference sequences. Other examples of reference sequences include genomes of other species, as well as chromosomes, sub-chromosomal regions (such as strands), etc., of any species.

[00148] In various embodiments, the reference sequence is a consensus sequence or other combination derived from multiple individuals. However, in certain applications, the reference sequence may be taken from a particular individual.

[00149] The term “derived” when used in the context of a nucleic acid or a mixture of nucleic acids, herein refers to the means whereby the nucleic acid(s) are obtained from the source from which they originate. For example, in one embodiment, a mixture of nucleic acids that is derived from two different genomes means that the nucleic acids, e.g., cfDNA, were naturally released by cells through naturally occurring processes such as necrosis or apoptosis. In another embodiment, a mixture of nucleic acids that is derived from two different genomes means that the nucleic acids were extracted from two different types of cells from a subject. For instance, a mixture of nucleic acids includes nucleic acids originating from donor cells and donee cells obtained from an organ transplant subject. In some implementations, a mixture of nucleic acids comprise biological materials of two or more contributor individuals. For example, a forensic

sample including biological materials of two or more individuals includes DNA of the two or more individuals.

[00150] The term “based on” when used in the context of obtaining a specific quantitative value, herein refers to using another quantity as input to calculate the specific quantitative value as an output.

[00151] The term “biological fluid” herein refers to a liquid taken from a biological source and includes, for example, blood, serum, plasma, sputum, lavage fluid, cerebrospinal fluid, urine, semen, sweat, tears, saliva, and the like. As used herein, the terms “blood,” “plasma” and “serum” expressly encompass fractions or processed portions thereof. Similarly, where a sample is taken from a biopsy, swab, smear, etc., the “sample” expressly encompasses a processed fraction or portion derived from the biopsy, swab, smear, etc.

[00152] As used herein, the term “corresponding to” sometimes refers to a nucleic acid sequence, e.g., a gene or a chromosome, that is present in the genome of different subjects, and which does not necessarily have the same sequence in all genomes, but serves to provide the identity rather than the genetic information of a sequence of interest, e.g., a gene or chromosome.

[00153] The term “contributor” herein refers to a human contributor as well as a non-human contributor such as a mammal, an invertebrate, a vertebrate, a fungus, a yeast, a bacterium, and a virus. Although the examples herein concern humans and the language is primarily directed to human concerns, the concepts disclosed herein are applicable to genomes from any plant or animal, and are useful in the fields of veterinary medicine, animal sciences, research laboratories and such.

[00154] The term “sensitivity” as used herein refers to the probability that a test result will be positive when the condition of interest is present. It may be calculated as the number of true positives divided by the sum of true positives and false negatives.

[00155] The term “specificity” as used herein refers to the probability that a test result will be negative when the condition of interest is absent. It may be calculated as the number of true negatives divided by the sum of true negatives and false positives.

[00156] The term “primer,” as used herein refers to an isolated oligonucleotide that is capable of acting as a point of initiation of synthesis when placed under conditions inductive to

synthesis of an extension product (e.g., the conditions include nucleotides, an inducing agent such as DNA polymerase, and a suitable temperature and pH). The primer is preferably single stranded for maximum efficiency in amplification, but may alternatively be double stranded. If double stranded, the primer is first treated to separate its strands before being used to prepare extension products. Preferably, the primer is an oligodeoxyribonucleotide. The primer must be sufficiently long to prime the synthesis of extension products in the presence of the inducing agent. The exact lengths of the primers will depend on many factors, including temperature, source of primer, use of the method, and the parameters used for primer design.

Introduction

[00157] This disclosure provides methods and systems for quantification and deconvolution of nucleic acid mixture samples including nucleic acid of two or more contributors of unknown genotypes, providing various advantages and technological improvements. For instance, some implementations apply probabilistic mixture modeling, Bayesian inference techniques, and numerical optimization methods to quantify contributor DNA in a mixture without knowing contributor's genotypes.

[00158] Sequencing data from a nucleic acid (e.g., DNA or RNA) mixture of closely related genomes is frequently found in research as well as clinical settings, and quantifying the mixing contributors has been a challenge when the original genomes are unknown.

[00159] Conventional methods of chimerism analysis (for bone marrow and blood stem cell transplants only) utilize capillary electrophoresis (CE) fragment analysis or quantitative polymerase chain reaction (qPCR) analysis of short tandem repeats (STRs) or small insertions and deletions (Indels). These methods tend to have poor limit of quantification, dynamic range, or reproducibility. They have limited number of targets, complicated workflow, and time-consuming and inaccurate manual input for analysis. The conventional methods tend to comprise among these different metrics. CE approach has a LOQ ranging from 1%-5%, and suffers from low reproducibility. These limitations can be significant in clinical use. For example, an actual chimerism result of 99% will be reported as 100%. The qPCR approach can achieve low LOQ of 0.1% but that requires 66ng or more chimerism DNA not considering the DNA required for pure baseline samples. Neither 66ng nor 10ng is possible for routine cfDNA analysis for solid organ

transplant. In addition, the dynamic range of qPCR-based chimerism suffers, and chimerism predictions when the minor contributor is greater than 30% are not reliable.

[00160] Given the high input DNA requirement, CE and qPCR approaches are only applicable to bone marrow or blood stem cell transplant. Neither approach works for solid organ transplant monitoring, for which the cfDNA amount from a typical blood draw is much less than 10ng. In addition, even at the same amount, cfDNA is not as effective as gDNA as PCR template.

[00161] Besides high DNA input requirement, both CE and qPCR approaches require pure pre-transplant baseline samples to be available. They are also associated with complex assay and require manual intervention in selecting the appropriate markers before quantification.

[00162] In addition to these, there are two fundamental challenges in chimerism quantification that our methods systematically addressed, while existing methods do not work.

[00163] The first challenge is to quantify chimerism sample with more than two contributors, corresponding to transplant with more than one donor. Multi-donor transplants are common for bone marrow and blood stem cell transplants. It also occurs in solid organ transplant, for example for 2nd kidney transplant following the failure of previous kidney, or when solid organ transplant coincide with blood transfusion from another donor.

[00164] The second challenge is to quantify chimerism samples when one of the contributors is unknown. This occurs frequently in clinical setting, for example 1) when the donor genome is not available, 2) when in the multi-donor cases, when an old organ's donor genome is not available, or 3) when solid organ transplant recipient also received blood transfusion from unknown donors.

[00165] While conventional methods do not address these challenges, the methods disclose herein can accurately quantify chimerism samples when there is an unknown donor. When there is only one donor neither the donor or recipient genome are required using the disclosed methods. Further, the disclosed methods can work with arbitrary number of donors. Some empirical studies have validated the performance of the disclosed method for 4 donors and achieved an LOQ of less than 0.35% at 10ng total gDNA input.

[00166] In some implementations, the disclosed methods achieve 0.1% to 0.2% LOQ with as low as 3ng cfDNA input, and achieve a broad dynamic range from 0.1% - 99.9%. Some implementations do not require baseline genomes to be known, although knowing the baseline can improve performance. The disclosed methods can work with chimerism samples of arbitrary number of donors, and have been experimentally validated for samples with 0-4 donors, which covers nearly all clinically relevant cases for solid organ transplant, bone marrow transplant, and hematopoietic stem cell transplants. In addition, the disclosed methods do not require any manual intervention in selecting genetic markers, allowing digitization and automation of quantification of nucleic acids.

[00167] Some implementations provide methods and systems for quantifying contributor DNA from multi-marker targeted-resequencing data of blood cfDNA or gDNA samples. Some implementations provide methods and systems for quantifying contributor DNA from multi-marker targeted-resequencing data of blood cfDNA or gDNA samples using novel probabilistic models and numerical optimization methods. Some implementations provide methods and systems for quantifying contributor DNA for genetically related donor and recipient of unknown genotypes using Bayesian modeling with prior distributions that encode genetic-relationship. By using genetic-relationship information to provide prior information in a Bayesian framework, quantification of DNA mixture can be improved compared to methods that do not use the genetic-relationship information.

[00168] Some implementations provide methods and systems for estimating the confidence interval of DNA quantification by numerically computing the Cramer-Rao bound from the estimated Hessian matrices of log-likelihood functions.

[00169] Allelic bias in short sequencing read mapping confounds DNA quantification. In some implementations, we reduce the confounding effect through an unbiased mapping strategy of reads spanning variant sites.

[00170] Implementations described herein can accurately estimate the contributor DNA fraction even though the genotypes for the contributing genomes are totally unknown. The allele fraction of a marker site after PCR amplification can be reliably modeled with a beta-distribution.

[00171] Using the unbiased reference DNA sequence database containing both reference and alternate allele, one can remove read mapping biases towards the reference alleles, and reliably estimate the allele counts and sequencing error at the variant sites.

[00172] Implementations described herein can estimate the confidence interval of the predicted contributor DNA fractions with a single sequencing run of a mixture DNA sample.

[00173] Formally, the problem of contributor DNA quantification (CDQ) is stated as following: Given the sequencing data of a DNA sample comprised of one or more contributors, determine the fraction of each contributor in the sample. When the genotypes of the contributor genomes are unknown, the CDQ problem is referred to as blind contributor DNA quantification (blind-CDQ); the opposite is referred to as non-blind-CDQ. Some descriptions regarding some implementations refer to the two contributors as the donor and the recipient, but they do not limit the applications of the methods to the organ donation setting. In some description hereinafter regarding some implementations, a contributor is equivalent to a donor, and the other contributor is equivalent to a donee.

[00174] Blind-CDQ is a harder problem compared to non-blind CDQ, but it is of wider application to all scenarios where only a single sequencing experiment of the mixture sample is achievable, while the non-blind-CDQ requires prior sequencing experiments to determine genotypes of the contributors (e.g. organ donors and recipients).

[00175] The computational methods described in this document address both the blind-CDQ and the non-blind-CDQ problems with single, two, or multiple contributors.

[00176] Figures 1A-1C show an overview of methods and statistical model designed for contributor DNA quantification. Figure 1A shows an experimental pipeline for sequencing based allogeneic DNA detection. Figure 1B shows an unbiased read mapping workflow for allele counting. Figure 1C shows a hierarchical, probabilistic mixture model for allelic counts per marker locus.

[00177] Some implementations apply experimental pipeline as depicted in Figure 1A. This generic experimental pipeline has the following steps.

[00178] 1) A blood or other type of sample is obtained containing DNA from multiple genetic origins.

[00179] 2) The appropriate type of DNA is extracted, e.g. cellular DNA (also referred to as gDNA) or cell free DNA (cfDNA), depending on the application.

[00180] 3) Specific variant sites or polymorphism sites of the genome is targeted and enriched by approaches such as PCR amplification or hybridization. The variant sites are prior-selected to be variable among diverse populations of human (or another organism of interest). Alternatively, untargeted (whole genome) sequencing can be done, and all variant sites will be covered.

[00181] 4) The DNA sample is sequenced by NGS or other DNA sequencing techniques such as some of the ones described hereinafter to obtain sequencing reads that cover variant sites of interest.

[00182] The computational method for CDQ has three main components:

[00183] 1) Allele Counting: an computer program based on hashing and dynamic programming for unbiased counting of sequencing reads from each allele for each target marker site (Figure 1B), and

[00184] 2) Contributor DNA Quantification: a hierarchical probabilistic model and novel combination of multi-iteration grid search strategy with BFGS - quasi-Newton method, or in some implementations an iterative weighted linear regression, for quantifying the contributor DNA fraction (Figure 1C).

[00185] 3) Confidence interval (uncertainly) determination: around the quantified mixture fractions, variances are determined based on the hessian matrix of the log likelihood function base on information inequality.

[00186] The totality of these components for chimerism quantification is impossible to execute manually by human experts or be carried out in their heads. They require computers and are computer-implemented technology. These computational components allow the disclosed methods to achieve unparalleled quantification sensitivity, dynamic range, and reproducibility. They also enable the disclosed methods to reliably quantify diverse set of chimerism samples, including cfDNA or gDNA, 3-10ng or more input DNA, 0 to 4 or more donors, and genetically related or unrelated donor with known or unknown genomes.

[00187] Although some implementations only address “relative quantification” here, meaning that the implementations estimate the percentage or fraction of the DNA sample that is originated from the contributor sources, rather than the absolute amount (in terms of mass or copy numbers). Additional steps can be taken to convert the relative abundance to absolute abundance if the total amount of input DNA is measured or known.

Overview of Processes for Quantifying Contributor Fractions in a Nucleic Acid Sample

[00188] Figure 2A shows a block diagram illustrating a process 200 for quantifying one or more fractions of nucleic acid (e.g., DNA or RNA) of one or more contributors in the nucleic acid sample. The method is implemented on a computer system that includes one or more processors and system memory such as the systems described hereinafter. Descriptions herein refer to DNA in some implementations and applications, but one skilled in the art appreciates that other forms of nucleic acids can also be analyzed using the implementations described herein. The various implementations described herein can be used to analyze a nucleic acid sample containing nucleic acid from one or more contributors. In some implementations, methods and systems are provided to quantify one or more fractions of nucleic acid of the one or more contributors. In some descriptions herein, the nucleic acid sample is referred to as a mixture sample because the sample can include nucleic acid from two more contributors. However, it is understood that the use of the term “mixture” indicates the possibility that the sample includes two or more contributors’ nucleic acid, and it does not exclude the possibility that the sample includes nucleic acid from only a single contributor. In the latter case, a fraction of 1 or a percentage of 100% (or values within a margin of error) may be determined for the one contributor.

[00189] In some implementations, the one or more contributors of the nucleic acid sample include a donor of a transplant and a donee of the transplant. In some implementations, the transplant includes an allogeneic or a xenogeneic transplant. In some implementations, the nucleic acid sample is a biological sample obtained from the donee. In some implementations, the nucleic acid sample includes cell-free nucleic acid. In some implementations, the sample includes cellular DNA. In some implementations, the nucleic acid sample includes nucleic acid from zero, one, or more contaminant genomes and one genome of interest. In some implementations, the nucleic acid sample includes a biological sample obtained from a cell

culture, which can be a mixture of multiple cell lines of different genetic origins in some implementations.

[00190] Process 200 involves extracting nucleic acid molecules from the nucleic acid sample using techniques such as those described herein. See block 202.

[00191] Process 200 further involves amplifying or enriching the extracted nucleic acid molecules. See block 204. Various amplification or enrichment techniques such as those described herein may be used. In some implementations, PCR are used to amplify the extracted nucleic acid molecules. In some implementations, the amplification targets specific polymorphisms, which amplification is also referred to as targeted enrichment. In other implementations, whole genome amplification may be performed, and allele data for specific polymorphism sites may be obtained by sequencing.

[00192] Process 200 also involves sequencing the amplified or enriched nucleic acid molecules using a nucleic acid sequencer to produce nucleic acid sequence reads. See block 206. Various sequencing techniques and devices are further described hereinafter, which may be applied in operation 206.

[00193] Process 200 further involves mapping the nucleic acid sequence reads to one or more polymorphism loci on a reference sequence. In some implementations, alignment techniques may be used to map the nucleic acid sequence reads to one or more polymorphism loci. In other implementations, an unbiased mapping techniques may be used to match the nucleic acid sequence reads to the polymorphism loci. See block 208. In some implementations, the nucleic acid sequence reads are mapped to specific alleles at the polymorphism loci. The unbiased mapping technique is further described hereinafter. In some implementations, the one or more polymorphism loci (or polymorphic loci) include biallelic loci. In some implementations, the alleles at the one or more polymorphism loci include single nucleotide polymorphism (SNP) alleles.

[00194] In some implementations, unique molecular indexes (UMIs) are attached to the extracted nucleic acid molecules, which are then amplified, sequenced, and mapped to the polymorphism loci or alleles. The unique molecular indices provide mechanisms to reduce the errors that can occur in the sample processing and analysis steps. For instance, different reads sharing a same unique molecular index (UMI) can be combined or collapsed to determine a

sequence from which the reads are derived, effectively removing errors that have occurred during amplification and sequencing.

[00195] Process 200 further involves determining, using the method nucleic acid sequence reads, allele counts of nucleic acid sequence reads for alleles at the one or more polymorphism loci. See block 210.

[00196] Process 200 also involves applying the probabilistic mixture model to the allele counts of nucleic acid sequence reads. The probabilistic mixture model uses probability distributions to model allele count of nucleic acid sequence reads at the one or more polymorphism loci. The probability distributions account for errors and noises in the nucleic acid sequence reads. The probabilistic mixture model treats each allele count of nucleic acid sequence reads as a random sample from a probability distribution.

[00197] In the equations hereinafter, the notations below are used.

[00198] d : indicator for donors, $d = 1, 2, \dots, D$, where D is the total number of contributors. D can be any natural number. In some implementations, D is 5 or smaller. In some implementations, D is 9 or smaller.

[00199] a : indicator for alleles. In some implementations, the alleles include biallelic SNPs, and $a = 1$ or 2.

[00200] i : indicator for marker loci, $i = 1 \dots I$, where I is the total number of markers, e.g. 300.

[00201] g_{dai} : genotypes of contributor d allele type a for marker i . It takes value 0, 1, or 2, representing the number of copies of allele a for this locus in this contributor.

[00202] n_{ai} , n_{ai}' , n_{ai}'' : copies of reads, DNA molecules after amplification, and DNA molecules before amplification, of allele type a and marker locus i .

[00203] n_i , n_i' , n_i'' : total copies of reads, nucleic acid molecules after amplification, and DNA molecules before amplification, for marker locus i .

[00204] r_{ai} : fractions of read counts for allele type a and marker locus i .

[00205] p_{ai} : probability of seeing a read of allele type a at a given marker locus i .

[00206] Note that for g_{dai} , n_{ai} , n_{ai}' , n_{ai}'' , n_i , n_i' , n_i'' , r_{ai} , and p_{ai} , the subscript i is sometimes omitted when the implementations are focused on a single locus.

[00207] β_d : fraction of nucleic acids from contributor d that contribute to a mixture sample.

[00208] λ : Sequencing error rate.

[00209] Bold letters represent vectors or matrices:

[00210] $\mathbf{g} = [g_{dli}]_{i=1\dots I, d=1\dots D}$: genotype matrix with reference allele counts in all contributors and all loci.

[00211] $\mathbf{g}_i = [g_{dli}]_{d=1\dots D}$: genotype vector with reference allele counts for all contributors and given locus i .

[00212] $\mathbf{r} = [r_{li}]_{i=1\dots I}$: allele fraction vector with fractions of allele 1 reads for every loci.

[00213] $\mathbf{n} = [n_i]_{i=1\dots I}$: read count vector with read count for every loci.

[00214] $\mathbf{p} = [p_{li}]_{i=1\dots I}$: vector with expected allele 1 fraction for every loci.

[00215] $\boldsymbol{\beta} = [\beta_d]_{d=1\dots D}$: contributor fraction vector with relative fraction of each contributor contributing to the nucleic acid sample.

[00216] In some implementations, the probabilistic mixture model uses a single-locus likelihood function to model allele counts at a single polymorphism locus, the single-locus likelihood function can be expressed as:

[00217] $M(n_{1i}, n_{2i} | p_{1i}, \theta)$, where n_{1i} is the allele count of allele 1 at locus i , n_{2i} is the allele count of allele 2 at locus i , p_{1i} is an expected fraction of allele 1 at locus i , and θ includes one or more model parameters.

[00218] In some implementations, p_{li} is modeled as a function $p(\mathbf{g}_i, \lambda, \boldsymbol{\beta})$ of: (i) genotypes of the contributors at locus i , or $\mathbf{g}_i = (g_{1li}, \dots, g_{Dli})$, which is a vector of copy number of allele 1 at locus i in contributors $1\dots D$; (ii) read count errors, or λ , resulting from the sequencing; and (iii) fractions of nucleic acid of contributors in the nucleic acid sample, or $\boldsymbol{\beta} = (\beta_1, \dots, \beta_D)$, where D is the number of contributors.

[00219] In some implementations, p_{li} is calculated as $p_{li} = p(\mathbf{g}_i, \lambda, \boldsymbol{\beta}) \leftarrow [(1 - \lambda) \mathbf{g}_i + \lambda (2 - \mathbf{g}_i)] / 2 \cdot \boldsymbol{\beta}$, where \cdot is vector dot product operator.

[00220] In some implementations, the contributors include two contributors, and p_{li} is obtained using the p_{li}' values in Table 3 described hereinafter.

[00221] In some implementations (method S), the single-locus likelihood function is a probability distribution that includes a first binomial distribution. In some implementations, the first binomial distribution includes a quantity parameter indicating the total allele count at a locus and a probability parameter indicating a probability of the first allele at the locus. In some implementations, the first binomial distribution is expressed as follows:

[00222] $n_{li} \sim BN(n_i, p_{li})$

[00223] where n_{li} is an allele count of nucleic acid sequence reads for allele 1 at locus i ; n_i is a total read count at locus i ; and p_{li} is a probability parameter indicating the probability of allele 1 at locus i .

[00224] In some implementations, the probability parameter p is a function of a fraction of nucleic acid of a contributor, or β . The probability parameter is also a function of genotypes of the one or more contributors \mathbf{g} . The probability parameter is also a function of errors resulting from the sequencing operation of 206, or λ . In some implementations, zero, one or more genotypes of the contributors were unknown. In some implementations the probabilistic mixture model includes various probability distributions as shown in Figure 2B.

[00225] Returning to Figure 2A, process 200 involves quantifying, using the probabilistic mixture model, one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample. See block 214. In some implementations, the quantifying includes marginalizing over a plurality of possible combinations of genotypes to enumerate the probability parameter p . In some implementations, the quantifying includes determining β , the fractions of nucleic acid of the contributors, using a multiple-loci likelihood function of the allele counts of nucleic acid sequence reads determined in operation 210 conditioned on parameters of the probabilistic mixture model.

[00226] In some implementations, the quantification includes calculating a plurality of likelihood values using a plurality of potential fraction values and a multiple-loci likelihood function of the allele counts of nucleic acid sequence reads. The quantification also involves identifying a potential fraction value that is associated with a likelihood value that is the maximum value among the plurality of likelihood values. In some implementations, the plurality of likelihood values is obtained for a plurality of parameters and the values thereof in a multi-dimensional grid. The quantification also involves quantifying the fraction of nucleic acid of the contributor(s) in the nucleic acid sample at the identified potential fraction value having the maximum likelihood. In some implementations, the multiple-loci likelihood function includes a plurality of marginal distributions for the one or more polymorphism loci.

[00227] In some implementations, the multiple-loci likelihood function of the one or more contributors, with known, unknown, or partially known genotypes, is computed as following:

$$[00228] \quad L(\beta, \theta, \lambda, \pi; n_1, n_2) = \prod_i [\sum_{g_i} M(n_{1i}, n_{2i} | p(g_i, \lambda, \beta), \theta) \cdot P(g_i | \pi)]$$

[00229] where $L(\beta, \theta, \lambda, \pi; n_1, n_2)$ is the likelihood of observing allele count vectors n_1 and n_2 for alleles 1 and 2; $p(g_i, \lambda, \beta)$ is the expected fraction or probability of observing allele 1 at locus i based on the contributors' genotypes g_i at locus i ; $P(g_i | \pi)$ is the prior probability of observing the genotypes g_i at locus i given a population allele frequency (π); and, \sum_{g_i} denotes summing over a plurality of possible combinations of genotypes of the contributors, subjecting to constraints of the known genotypes for some or all the contributors.

[00230] In some implementations, the prior joint probability is calculated using marginal distributions $P(g_{1i} | \pi)$ and $P(g_{2i} | \pi)$ that satisfy the Hardy-Weinberg equilibrium.

[00231] In some implementations, all genotypes are known, and the multi-loci likelihood function is computed using the genotype vector g_i representing the known genotype combination for the contributors: $L(\beta, \theta, \lambda, \pi; n_1, n_2) = \prod_i [M(n_{1i}, n_{2i} | p(g_i, \lambda, \beta), \theta) \cdot P(g_i | \pi)]$.

[00232] In some implementations, the probabilistic mixture model accounts for nucleic acid molecule number errors resulting from extracting the nucleic acid molecules performed in 202, as well as the read count errors resulting from the sequencing operation in 206.

[00233] In some implementations, the probabilistic mixture model uses a second binomial distribution to model allele counts of the extracted nucleic acid molecules for alleles at the one or

more polymorphism loci. In some implementations, the second binomial distribution is expressed as follows:

[00234] $n_{li}'' \sim BN(n_i'', p_{li})$

[00235] where n_{li}'' is an allele count of extracted nucleic acid molecules for allele l at locus i ; n_i'' is a total nucleic acid molecule count at locus i , which equals to a total genome copy numbers n'' ; and p_{li} is a probability parameter indicating the probability of allele l at locus i .

[00236] In some implementations, the first binomial distribution is conditioned on an allele fraction n_{li}''/n_i'' . In some implementations, the first binomial distribution is re-parameterized as follows:

[00237] $n_{li} \sim BN(n_i, n_{li}''/n_i'')$

[00238] where n_{li} is an allele count of nucleic acid sequence reads for allele l at locus i .

[00239] In some implementations, the probabilistic mixture model uses a first beta distribution to approximate a distribution of n_{li}''/n_i'' . In some implementations, the first beta distribution has a mean and a variance that match a mean and a variance of the second binomial distribution.

[00240] In some implementations, locus i is modeled as biallelic and the first beta distribution is expressed as follows:

[00241] $n_{li}''/n_i'' \sim \text{Beta}((n''-1)p_{1i}, (n''-1)p_{2i})$

[00242] where p_{1i} is a probability parameter indicating the probability of a first allele at locus i ; and p_{2i} is a probability parameter indicating the probability of a second allele at locus i .

[00243] In some implementations, the process includes combining the first binomial distribution, modeling sequencing read counts, and the first beta distribution, modeling extracted nucleic acid molecule number, to obtain the single-locus likelihood function of n_{li} that follows a first beta-binomial distribution.

[00244] In some implementations, the first beta-binomial distribution has the form:

[00245] $n_{li} \sim BB(n_i, (n''-1) \cdot p_{1i}, (n''-1) \cdot p_{2i}),$

[00246] or an alternative approximation:

[00247] $n_{1i} \sim BB(n_i, n'' \cdot p_{1i}, n'' \cdot p_{2i})$.

[00248] In some implementations, the multiple-loci likelihood function can be expressed as:

[00249] $L(\beta, n'', \lambda, \pi; n_1, n_2) = \prod_i [\sum_{g_i} BB(n_{1i} | n_i, (n''-1) \cdot p_{1i}, (n''-1) \cdot p_{2i}) \cdot P(g_i | \pi)]$

[00250] where $L(\beta, n'', \lambda, \pi; n_1, n_2)$ is the likelihood of observing allele count vectors n_1 and n_2 for alleles 1 and 2 at all loci, and $p_{1i} = p(g_i, \lambda, \beta)$, $p_{2i} = 1 - p_{1i}$.

[00251] In some implementations, the contributors include two contributors, and the multiple-loci likelihood function is expressed as:

[00252] $L(\beta, n'', \lambda, \pi; n_1, n_2) = \prod_i \sum_{g_{1i} g_{2i}} BB(n_{1i}, n_{2i} | n_i, (n''-1) \cdot p_{1i}(g_{1i}, g_{2i}, \lambda, \beta), (n''-1) \cdot p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)$

[00253] where $L(\beta, n'', \lambda, \pi; n_1, n_2)$ is the likelihood of observing an allele count vector for the first allele of all loci (n_1) and an allele count vector for the second allele of all loci (n_2) given parameters β, n'', λ , and π ; $p_{1i}(g_{1i}, g_{2i}, \lambda, \beta)$ is a probability parameter, taken as p_1' from Table 3, indicating a probability of allele 1 at locus i based on the two contributors' genotypes (g_{1i}, g_{2i}); $p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)$ is a probability parameter, taken as p_2' from Table 3, indicating a probability of allele 2 at locus i based on the two contributors' genotypes (g_{1i}, g_{2i}); and $P(g_{1i}, g_{2i} | \pi)$ is a prior joint probability of observing the first contributor's genotype for the first allele (g_{1i}) and the second contributor's genotype for the first allele (g_{2i}) at locus i given a population allele frequency (π).

[00254] In some implementations, operation 214 includes estimating the total extracted genome copy number n'' from a mass of the extracted nucleic acid molecules. In some implementations, the estimated total extracted genome copy number n'' is adjusted according to fragment size of the extracted nucleic acid molecules as further described hereinafter.

[00255] In some implementations, the probabilistic mixture model accounts for nucleic acid molecule number errors resulting from amplifying the nucleic acid molecules performed in 204, as well as the read count errors resulting from the sequencing operation in 206. In some implementations, the nucleic acid amplification process is modeled as follows:

$$x_{t+1} = x_t + y_{t+1}$$

wherein x_{t+1} is the nucleic acid copies of a given allele after cycle $t+1$ of amplification; x_t is the nucleic acid copies of a given allele after cycle t of amplification; y_{t+1} is the new copies generated at cycle $t+1$, and it follows a binomial distribution $y_{t+1} \sim BN(x_t, r_{t+1})$; and r_{t+1} is the amplification rate for cycle $t+1$.

[00256] In some implementations, the probabilistic mixture model uses a second beta distribution to model allele fractions of the amplified nucleic acid molecules for alleles at the one or more polymorphism loci. In some implementations, locus i is modeled as biallelic and the second beta distribution is expressed as follows:

$$\mathbf{[00257]} \quad n_{1i}' / (n_{1i}' + n_{2i}') \sim \text{Beta}(n'' \cdot \rho_i \cdot p_{1i}, n'' \cdot \rho_i \cdot p_{2i})$$

[00258] where n_{1i}' is an allele count of amplified nucleic acid molecules for a first allele at locus i ; n_{2i}' is an allele count of amplified nucleic acid molecules for a second allele at locus i ; n'' is a total nucleic acid molecule count at any locus; ρ_i is a constant related to an average amplification rate r_i over all amplification cycles; p_{1i} is the probability of the first allele at locus i ; and p_{2i} is the probability of the second allele at locus i . In some implementations, ρ_i is $(1+r_i)/(1-r_i) / [1-(1+r_i)^{-t}]$. In some implementations, ρ_i is approximated as $(1+r_i)/(1-r_i)$.

[00259] In some implementations, operation 214 includes combining the first binomial distribution and the second beta distribution to obtain the single-locus likelihood function for n_{1i} that that follows a second beta-binomial distribution. In some implementations, the second beta-binomial distribution has the form:

$n_{1i} \sim BB(n_i, n'' \cdot \rho_i \cdot p_{1i}, n'' \cdot \rho_i \cdot p_{2i})$, wherein n_{1i} is an allele count of nucleic acid sequence reads for the first allele at locus i ; p_{1i} is a probability parameter indicating the probability of a first allele at locus i ; and p_{2i} is a probability parameter indicating the probability of a second allele at locus i .

[00260] In some implementations, operations 214 includes, by assuming the one or more polymorphism loci have a same amplification rate, re-parameterizing the second beta-binomial distribution as:

$$n_{1i} \sim BB(n_i, n'' \cdot (1+r)/(1-r) \cdot p_{1i}, n'' \cdot (1+r)/(1-r) \cdot p_{2i}), \text{ wherein } r \text{ is an amplification rate.}$$

[00261] In some implementations, operation 214 includes quantifying the one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample using a

multiple-loci likelihood function obtained using the second beta-binomial distribution, the multiple-loci likelihood function is as follows:

$$L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i [\sum g_i BB(n_{1i} | n_i, n'' \cdot (1+r)/(1-r) \cdot p_{1i}, n'' \cdot (1+r)/(1-r) \cdot p_{2i}) \cdot P(g_i | \pi)]$$

[00262] In some implementations, the contributors include two contributors and the multiple-loci likelihood function comprises:

$$[00263] \quad L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i \sum_{g_{1i} g_{2i}} [BB(n_{1i} | n_i, n'' \cdot (1+r)/(1-r) \cdot p_{1i}(g_{1i}, g_{2i}, \lambda, \beta), n'' \cdot (1+r)/(1-r) \cdot p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)]$$

[00264] wherein $L(\beta, n'', r, \lambda, \pi; n_1, n_2)$ is the likelihood of observing an allele count vector for the first allele of all loci (n_1) and an allele count vector for the second allele of all loci (n_2) given parameters β, n'', r, λ , and π .

[00265] In some implementations, operation 214 includes, by defining a relative amplification rate of each polymorphism locus to be proportional to a total reads per locus, re-parameterizing the second beta-binomial distribution as:

$$[00266] \quad n_{1i} \sim BB(n_i, c' \cdot n_i \cdot p_{1i}, c' \cdot n_i \cdot p_{2i}), \text{ wherein } c' \text{ is a parameter to be optimized.}$$

[00267] In some implementations, operation 214 includes quantifying the one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample using a multiple-loci likelihood function obtained using the second beta-binomial distribution, the multiple-loci likelihood function follows:

$$[00268] \quad L(\beta, n'', c', \lambda, \pi; n_1, n_2) = \prod_i [\sum g_i BB(n_{1i} | n_i, c' \cdot n_i \cdot p_{1i}, c' \cdot n_i \cdot p_{2i}) \cdot P(g_i | \pi)]$$

[00269] In some implementations, the probabilistic mixture model accounts for nucleic acid molecule number errors resulting from extracting the nucleic acid molecules performed in 202 and amplifying the nucleic acid molecules performed in 204, as well as the read count errors resulting from the sequencing operation in 206.

[00270] In some implementations, the probabilistic mixture model uses a third beta distribution to model allele fractions of the amplified nucleic acid molecules for alleles at the one or more polymorphism loci, accounting for the sampling errors resulting from extracting the nucleic acid molecules performed in 202 and amplifying the nucleic acid molecules performed in

204. In some implementations, locus i is modeled as biallelic and the third beta distribution has the form of:

[00271]
$$n_{1i}' / (n_{1i}' + n_{2i}') \sim \text{Beta}(n'' \cdot (1 + r_i) / 2 \cdot p_{1i}, n'' \cdot (1 + r_i) / 2 \cdot p_{2i})$$

[00272] where n_{1i}' is an allele count of amplified nucleic acid molecules for a first allele at locus i ; n_{2i}' is an allele count of amplified nucleic acid molecules for a second allele at locus i ; n'' is a total nucleic acid molecule count; r_i is the average amplification rate for locus i ; p_{1i} is the probability of the first allele at locus i ; and p_{2i} is a probability of the second allele at locus i .

[00273] In some implementations, operation 214 includes combining the first binomial distribution and the third beta distribution to obtain the single-locus likelihood function of n_{1i} that follows a third beta-binomial distribution. In some implementations, the third beta-binomial distribution has the form:

[00274]
$$n_{1i} \sim \text{BB}(n_i, n'' \cdot (1 + r_i) / 2 \cdot p_{1i}, n'' \cdot (1 + r_i) / 2 \cdot p_{2i})$$

[00275] where r_i is an amplification rate.

[00276] In some implementations, the multiple-loci likelihood function is:

[00277]
$$L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i [\sum_{g_i} \text{BB}(n_{1i} | n_i, n'' \cdot (1 + r) / 2 \cdot p_{1i}, n'' \cdot (1 + r) / 2 \cdot p_{2i}) \cdot P(g_i | \pi)]$$

[00278] where r is an amplification rate assumed to be equal for all loci.

[00279] In some implementations, the contributors include two contributors, and the multiple-loci likelihood function is:

[00280]
$$L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i \sum_{g_{1i} g_{2i}} \text{BB}(n_{1i} | n_i, n'' \cdot (1 + r) / 2 \cdot p_{1i}(g_{1i}, g_{2i}, \lambda, \beta), n'' \cdot (1 + r) / 2 \cdot p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)$$

[00281] where $L(n_1, n_2 | \beta, n'', r, \lambda, \pi)$ is the likelihood of observing allele counts for the first allele vector n_1 and an allele count for the second allele vector n_2 given parameters β, n'', r, λ , and π .

[00282] In some implementations, process 200 further includes estimating, using the Cramer-Rao inequality, one or more confidence intervals of the one or more fractions of nucleic acid of the one or more contributors.

[00283] In some implementations, the mapping operation of 208 includes identifying reads among the nucleic acid sequence reads matching any sequence of a plurality of unbiased target sequences, wherein the plurality of unbiased target sequences includes sub-sequences of the reference sequence and sequences that differ from the subsequences by a single nucleotide.

[00284] In some implementations, the plurality of unbiased target sequences comprises five categories of sequences encompassing each polymorphic site of a plurality of polymorphic sites: (i) a reference target sequence that is a sub-sequence of the reference sequence, the reference target sequence having a reference allele with a reference nucleotide at the polymorphic site; (ii) alternative target sequences each having an alternative allele with an alternative nucleotide at the polymorphic site, the alternative nucleotide being different from the reference nucleotide; (iii) mutated reference target sequences comprising all possible sequences that each differ from the reference target sequence by only one nucleotide at a site that is not the polymorphic site; (iv) mutated alternative target sequences comprising all possible sequences that each differ from an alternative target sequence by only one nucleotide at a site that is not the polymorphic site; and (v) one or more unexpected allele target sequences each having an unexpected allele different from the reference allele and the alternative allele, and each having a sequence different from the previous four categories of sequences. In some implementations, the five categories of sequences have the same length and are located at the same region of a genome.

[00285] In some implementations, operation 208 includes using the identified reads and their matching unbiased target sequences to determine allele counts of the nucleic acid sequence reads for the alleles at the one or more polymorphism loci. In some implementations, the plurality of unbiased target sequences includes sequences that are truncated to have the same length as the nucleic acid sequence reads. In some implementations, the plurality of unbiased target sequences includes sequences stored in one or more hash tables, and the reads are subsequently identified using the hash tables.

[00286] In some implementations, the process 200 further includes a procedure to determine if a contributor of known genotype is a true contributor to a mixture sample by comparing two versions of maximized multi-loci likelihood values, one version using a genotype

matrix containing the known genotype for the contributor, another version using a genotyping matrix with unknown genotype for the contributor.

[00287] In some implementations, the process further includes determining one or more genotypes of the one or more contributors at the one or more polymorphism loci. In some implementations, the process includes determining, using the one or more fractions of nucleic acid of the one or more contributors, a risk of one contributor (a donee) rejecting a tissue or an organ transplanted from another contributor (a donor). In many applications, the risk is not or cannot be based solely on the estimated contributor fractions. Instead, the contributor fractions are used as an intermediate parameter or intermediate result for determining the risk. In various implementations, other parameters obtained from other methods are combined with the contributor fractions to determine the risk. Such other methods include, and are not limited to, tissue biopsy, serum creatinine measurement, HLA-DSA (donor specific antibody) analysis.

[00288] Figure 3 shows a block diagram illustrating process 300 for evaluating a nucleic acid sample including nucleic acid of one or more contributors. Process 300 starts by receiving nucleic acid sequence reads of one or more alleles at one or more polymorphism loci obtained from the nucleic acid sample. See block 302. In some implementations, the nucleic acid sequence reads were obtained by sequencing the nucleic acid in the nucleic acid sample using various techniques described herein.

[00289] In some implementations, unique molecular indexes (UMIs) are attached to the extracted nucleic acid molecules, which are then amplified, sequenced, and mapped to the polymorphism loci or alleles. The unique molecular indices provide mechanisms to reduce the errors that can occur in the sample processing and analysis steps. For instance, different reads sharing a same unique molecular index (UMI) can be combined or collapsed to determine a sequence from which the reads are derived, effectively removing errors that have occurred during sample processing and sequencing. US Patent Application No. 15/130,668, filed April 16, 2016, and US Patent Application No. 15/863,737, filed January 5, 2018 describe various methods and systems for sequencing nucleic acids using unique molecular indexes, which are incorporated by reference by their entireties for all purposes.

[00290] When UMI is used in an assay, the redundant DNA molecules resulting from PCR amplification of a template nucleotide acid are collapsed into a single read. For such

experimental procedure, a preferred model for single locus read counts is the first beta-binomial distribution, which combined the first binomial distribution, modeling sequencing read counts, and the first beta distribution, modeling extracted nucleic acid molecule number.

[00291] When UMI is not used in an assay, nucleic acid extraction, amplification, and sequencing all contribute to the statistical variability in read counts. For such experimental procedure, a preferred model for single locus read counts is the third beta-binomial distribution, which combined the first binomial distribution, modeling sequencing read counts, the third beta distribution, modeling allele fractions of the amplified nucleic acid molecules, and the first beta distribution, modeling allele fractions in extracted nucleic acid molecule.

[00292] Process 300 further involves determining, using the nucleic acid sequence reads, allele counts for the one or more alleles at the one or more polymorphism loci.

[00293] Process 300 also involves applying the probabilistic mixture model to the allele counts. The probabilistic model uses probabilistic distributions to model allele counts of alleles at the one or more polymorphism loci. The probabilistic distributions count for errors in the allele data. The errors include errors originating from nucleic acid extraction, sample processing, and sequencing operations.

[00294] In some implementations, the probabilistic distributions include a first binomial distribution. In some implementations, the first binomial distribution includes a parameter indicating the total allele count at a locus and a probability parameter indicating the probability of the first allele at the locus. In some implementations, the probability parameter is a function of the fractions of nucleic acid of the one or more contributors in the nucleic acid sample. The probability parameter is also a function of genotypes of the one or more contributors, or G , and a function of errors in the nucleic acid sequence read data, or θ . In some implementations, the errors in the read data include errors originating from nucleic acid extraction, sample processing, and sequencing operations.

[00295] Process 300 also involves obtaining likelihood values of observing the allele data given model parameters and potential nucleic acid fraction values. See block 308.

[00296] In some implementations, process 300 involves quantifying, using the likelihood values, fractions of nucleic acid of the one or more contributors in the nucleic acid sample. See block 310.

[00297] In some implementations, process 300 further involves determining, using the likelihood values, at least one genotype for at least one of the contributors. See block 312.

[00298] In some implementations, genotypes of the contributors were unknown prior to process 300.

[00299] In some implementations, the probabilistic mixture model uses a beta distribution to model the errors in the allele data. In some implementations, the beta distribution is defined by a mean parameter and a concentration parameter. In some implementations, the concentration parameter has discrete prior representing different noise conditions. The concentration parameter varies across loci.

[00300] In some implementations, the quantification of operation 310 includes combining the first binomial distribution and the beta distribution to obtain a marginal distribution that follows a beta-binomial distribution.

[00301] In some implementations, the quantification of 310 includes quantifying the fractions of nucleic acid of the one or more contributors in the nucleic acid sample using a multiple-loci likelihood function of the allele data. In some implementations, the quantification involves calculating a plurality of likelihood values using a plurality of potential fraction values and a multiple-loci likelihood function of the allele counts. The quantification also involves identifying a potential fraction vector associated with the maximum likelihood value, and quantifying the fractions of nucleic acid of the one or more contributors in the nucleic acid sample using the identified potential fraction vector.

[00302] In some implementations, the multiple-loci likelihood function depends on $P(G|\pi)$, which is a prior probability of the genotype of the one or more contributors given a population allele frequency (π). In some implementations, the prior probability is calculated considering a dummy allele with a fixed prior probability representing mechanistic drop-out.

[00303] In some implementations, the one or more contributors include two or more contributors. In some implementations, process 300 includes an operation of determining a total

number of contributors in the one or more contributors. In some implementations, one or more genotypes of the one or more contributors were unknown, and process 300 includes an operation of determining an allele configuration at each of the one or more polymorphism loci, the allele configuration comprising an allele for each of the one or more contributors. In some implementations, process 300 includes an operation of determining an estimated probability for the allele configuration.

[00304] In some implementations, process 300 further includes obtaining a posterior probability that a specific contributor among the one or more contributors has a specific genotype. In some implementations, process 300 further includes calling, based on the posterior probability, that the nucleic acid sample includes nucleic acid from the specific contributor. In some implementations, obtaining the posterior probability that a specific contributor among the one or more contributors has a specific genotype includes: (i) multiplying prior probabilities of genotype configurations by likelihoods of the genotype configurations; (ii) normalizing a product of (i) by a sum over genotype space; and (iii) summing over genotype configurations containing the specific genotype to obtain the posterior probability.

[00305] In some implementations, the specific genotype includes a multiple-locus genotype, and the method further includes: summing, over all contributors, a posterior probability that a contributor has the specific genotype at all loci; and determining, based on the summed probability, the specified multiple-locus genotype appears in any contributor.

[00306] In some implementations, the nucleic acid sample is a forensic sample and the data of the multiple-locus genotype is obtained from a person of interest. The process further includes determining that the person of interest is a contributor of the nucleic acid sample.

[00307] In some implementations, the probabilistic mixture model uses a second binomial distribution to model stutter errors in the allele data. In some implementations, the second binomial distribution is expressed as follows:

$$s_{ik} \sim BN(n_{i(k+I)}, r_i)$$

where s_{ik} is a stutter allele count at locus i of a stutter allele that appears to be allele k but actually results from a stutter error of allele $k+I$; $n_{i(k+I)}$ is an original allele count of allele $k+I$ at locus i ; and r_i is a stutter rate for locus i .

[00308] In some implementations, the stutter rate r varies across loci and has a prior representing different noise conditions, the prior being shared across loci.

[00309] In some implementations, operation 310 includes quantifying fractions of nucleic acid of the one or more contributors in the nucleic acid sample using a multiple-loci likelihood function including a product of likelihoods of non-stutter allele counts and likelihoods of stutter allele counts.

[00310] In some implementations, applying the probabilistic mixture model includes adding a fixed number of molecules to an allele count assigned to allele $k+1$ when determining a number of molecules from which stutter can potentially originate.

[00311] In some implementations, the probabilistic mixture model uses a dummy out-of-sample allele to model natural drop-out. In some implementations, the prior of the dummy out-of-sample allele is proportional to a number of unobserved alleles. In some implementations, the number of unobserved alleles is estimated by: interpolating all integers between the shortest and longest observed integer-valued alleles, adding any observed non-integer-valued alleles, and returning the maximum of the resulting value and a criterion value.

[00312] In some implementations, applying the probabilistic mixture model involves pruning genotype configurations from data used to quantify the fractions of nucleic acid of the one or more contributors in the nucleic acid sample. In some implementations, pruning genotype configurations involves: limiting genotype configurations that are plausible by constructing a list of required alleles and excluding loci with not enough contributors to explain all required alleles. In some implementations, the list of required alleles consists essentially of alleles having allele counts above a threshold and too high to be plausible due to stutter drop-in. In some implementations, the threshold is a sum of (i) a maximum non-stutter allele count, and (ii) a value multiplied by a count of potential stutter donor alleles. In some implementations, pruning genotype configurations involves removing genotype configurations that have poor matches between the allele data and expected allele counts. In some implementations, the genotype configurations that have poor matches have root mean squared error (RMSE) values larger than one or more thresholds.

[00313] In some implementations, the alleles at the one or more polymorphism loci include single nucleotide polymorphism (SNP) alleles and/or short tandem repeat (STR) alleles.

Method for Unbiased Mapping of Reads to Marker Sites

[00314] Conventional computational methods for mapping nucleic acid (e.g., DNA or RNA) sequencing reads to the genome can be biased by the reference genome used. Since only one allele (the reference allele) for each variant site is present in the reference genome, mismatches between the reads and references are treated as sequencing errors in existing read mapping strategies. The problem is that when reads containing the non-reference alleles are treated as containing sequencing errors, the alignment confidence (score) is decreased, and hence they are less likely to be retained as confidently mapped reads in subsequent filtering steps. This mapping bias will skew the allele counts (Figure 1B), and subsequently compromise the estimation of contributor DNA fractions.

[00315] To address the mapping bias issue and enable optimal CDQ, some implementations provide a novel workflow for mapping reads to variant sites. The new read mapping approach enables unbiased counting of alleles and estimation of sequencing error on variant sites and non-variant sites.

[00316] The read mapping workflow is as follows. The workflow first generates five types of sequences (see Table 1) based on 1) the reference sequences and 2) the known alleles of the variant sites. If more than one single mutation is allowed per sequence, more types of sequences will be generated. The five types of sequences are referred to as ref, alt, ref.mut, alt.mut, and snp.mut respectively. For example, for each biallelic SNP marker site covered by a target sequence of length L , there are one ref, one alt, $[L - 1] \times 3$ ref.mut, $[L - 1] \times 3$ alt.mut, and 2 snp.mut sequences. All five types of sequences are then included in the database of “unbiased target sequences” (Figure 1B). Depending on the length of the reads from the sequencer, the unbiased target sequences are then truncated into two versions. Let r be the read length. Version 1 of the truncated target sequences comprises the r 5' bases of all unbiased target sequences, while version 2 of the truncated target sequences comprises the reverse complement of the r 3' bases of all unbiased target sequences. Redundant sequences in truncated target sequences are then removed. The unique sequences in the two truncated sequence databases are then recorded into two hash tables. Next, sequencing reads are counted using the hash tables. For pair end sequencing strategies, R1 reads and R2 reads are counted using the first and second hash tables respectively. For non-pair end sequencing, all reads are counted using the first hash table.

Finally, for each marker site, the counts are aggregated into the five types defined above depending on which type the truncated unbiased target sequences corresponds to in Table 1.

[00317] A similar strategy can be implemented when sequence alignment tools are used instead of using hash table for the mapping. For each marker site, the ref and alt types of sequences are generated to form the unbiased sequence database. Each sequencing read is then aligned to this database with up to a predefined number of sequencing errors. The mapped reads are then categorized based on Table 1. For SNP markers only the bi-allelic scenario is presented here, but the method extends to multi-allelic loci.

Table 1. Definition of five types of target sequences to be generated from the reference sequence around a variant site.

Type	Definition
ref	SNP site taking reference allele
alt	SNP site taking alternative allele
ref.mut	Single mutation on non SNP site when the SNP site is ref
alt.mut	Single mutation on non SNP site when the SNP site is alt
snp.mut	SNP site taking neither reference nor alternative alleles

[00318] The proposed read mapping workflow addresses the read mapping bias issue when tested using real data. With the workflow, the observed error rates of the reference to alternative errors and the alternative to reference errors are identical. The sequencing error rate on the non-variant sites on the reference DNA copy and that on the alternative DNA copy are also identical.

Linking Contributor DNA Fraction with Allele Fractions

Sequencing Error-Free Scenario

[00319] We denote n_1 as the number of contributor 1 (e.g. organ recipient) cells and n_2 as the number of contributor 2 (e.g. organ donor) cells that supplied DNA to the sample. Based on these cells, the implementations define the contributor 2 fraction as $\beta_2 = n_2/(n_1 + n_2)$. For two-contributor scenario, we denote β_2 as β for short. Depending on the genotypes of the two

contributors at each specific locus, the two alleles have different fractions (see Table 2 for details), and the generic formula for calculating them is $p_1 = [g_{11}(1-\beta) + g_{21}\beta] / 2$ and $p_2 = [g_{12}(1-\beta) + g_{22}\beta] / 2$. Note that g_{11} and g_{12} are the contributor 1 (recipient) genotype, i.e. copies of allele 1 and 2 in the recipient genome; g_{21} and g_{22} are contributor 2 (donor) genotype, i.e. copies of allele 1 and 2 in the donor genome.

[00320] In matrix notation, the relationship for multiple contributor cases is generally implemented as $\mathbf{p} \leftarrow \mathbf{g}/2 \cdot \boldsymbol{\beta}$, where \mathbf{p} is a vector of expected allele 1 fraction for all loci, \mathbf{g} is a matrix of genotype of all loci in all contributor, and $\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_D]$ is the vector of nucleic acid fractions for all contributors. The implementation is generally applied to single-, two-, and multi-contributor scenarios.

Table 2: The binomial model parameters expected allele 1 and allele 2 fractions p_1 and p_2 for the 9 possible genotype combinations between a contributor 1 and contributor 2 pair for a given variant site

g_{11}	g_{21}	p_1	p_2
0	0	0	1
0	1	$\beta/2$	$1-\beta/2$
0	2	β	$1-\beta$
1	0	$(1-\beta)/2$	$(1+\beta)/2$
1	1	$1/2$	$1/2$
1	2	$(1+\beta)/2$	$(1-\beta)/2$
2	0	$1-\beta$	β
2	1	$1-\beta/2$	$\beta/2$
2	2	1	0

General Scenario with Sequencing Error

[00321] When there are two known alleles at a variant site, sequencing errors will convert one allele to another in addition to converting the two known alleles to the two remaining nucleotides at this locus. The consequence is that the allele fractions in the sequenced reads will deviate from the allele fractions in the NGS input sample.

[00322] Figure 2C schematically illustrates sequencing errors that convert one allele to another allele and true alleles to unexpected alleles. Panel (A) shows nucleotide-dependent sequencing error, and panel (B) shows uniform sequencing error.

[00323] Let N_1, N_2 be the allele 1 and allele 2 nucleotides. Let p_1', p_2' be the probability of observing allele 1 and allele 2 reads respectively, whether it is real or due to sequencing error; and $p_0' = 1 - p_1' - p_2'$ be the probability of observing the two unexpected alleles due to sequencing error. Let $\lambda_{N_1 N_2}$ be the mutation rate (probability) from N_1 to N_2 , where N_1 and N_2 are unique to each SNP site, and

[00324] $\lambda_{N_1 \#}$: mutation probability from N_1 to any of the 3 nucleotide non- N_1 nucleotides.

[00325] The transition diagram among the 4 nucleotide of a SNP site is shown in Figure 2C. Based on this, the implementations obtain the following equations for converting from true allele fractions p_1, p_2 to observed allele fractions p_1', p_2' , and p_0' :

[00326]
$$p_1' = p_1 - p_1 \cdot \lambda_{N_1 \#} + p_2 \cdot \lambda_{N_2 N_1}$$

[00327]
$$p_2' = p_2 - p_2 \lambda_{N_2 \#} + p_1 \cdot \lambda_{N_1 N_2}$$

[00328]
$$p_0' = p_1 \cdot (\lambda_{N_1 \#} - \lambda_{N_1 N_2}) + p_2 \cdot (\lambda_{N_2 \#} - \lambda_{N_2 N_1}).$$

[00329] When the implementations assume uniform sequencing error rate that is independent to the nucleotide identity, the implementations have,

[00330]
$$p_1' = p_1 \cdot (1 - 3 \cdot \lambda) + p_2 \cdot \lambda$$

[00331]
$$p_2' = p_2 \cdot (1 - 3 \cdot \lambda) + p_1 \cdot \lambda$$

[00332]
$$p_0' = 2\lambda.$$

[00333] When the implementations ignore the unexpected alleles

[00334]
$$p_1' = (p_1 \cdot (1 - 3 \cdot \lambda) + p_2 \cdot \lambda) / (1 - 2\lambda)$$

[00335]
$$p_2' = (p_2 \cdot (1 - 3 \cdot \lambda) + p_1 \cdot \lambda) / (1 - 2\lambda),$$

[00336] with $o(\lambda^2)$ approximation error, these are rewritten as

[00337]
$$p_1' = p_1 \cdot (1 - \lambda) + p_2 \cdot \lambda$$

[00338]
$$p_2' = p_2 \cdot (1 - \lambda) + p_1 \cdot \lambda$$

[00339] Or for locus i and substituting g and β for p :

[00340]
$$p_{li}' \leftarrow \Sigma_d [(g_{dli} \cdot (1 - \lambda) + g_{d2i} \cdot \lambda) \cdot \beta_d] / 2$$

[00341] $p_{2i}' \leftarrow \sum_d [(g_{d2i} \cdot (1 - \lambda) + g_{d1i} \cdot \lambda) \cdot \beta_d] / 2$

[00342] which is referred to as an error-adjusted-genotype weighted mixing coefficients.

[00343] The formula linking contributor 2 fraction β with the observed allele fraction p_1' in two contributor scenario is listed in Table 3.

Table 3: Expected probabilities of observing alleles 1 and 2 allowing for sequencing errors, conditioned on each donor/recipient genotype combination in a two-contributor setting. Here a uniform sequencing error rate $\lambda_{N_1N_2} = \lambda$ is used for all nucleotide pairs N_1 and N_2 . Since mutation rate λ is small, a first order approximation is used.

g_{11}	g_{21}	p_1'	p_2'
0	0	λ	$1 - \lambda$
0	1	$\beta/2 + \lambda - \beta\lambda$	$1 - \beta/2 - \lambda + \beta\lambda$
0	2	$\beta + \lambda - 2\beta\lambda$	$1 - \beta - \lambda + 2\beta\lambda$
1	0	$(1-\beta)/2 + \beta\lambda$	$(1+\beta)/2 - \beta\lambda$
1	1	$1/2$	$1/2$
1	2	$(1+\beta)/2 - \beta\lambda$	$(1-\beta)/2 + \beta\lambda$
2	0	$1 - \beta - \lambda + 2\beta\lambda$	$\beta + \lambda - 2\beta\lambda$
2	1	$1 - \beta/2 - \lambda + \beta\lambda$	$\beta/2 + \lambda - \beta\lambda$
2	2	$1 - \lambda$	λ

[00344] In matrix format, error-adjusted-genotype for allele 1 accounting for sequencing error λ is implemented as: $G \leftarrow [(1 - \lambda) g + \lambda (2 - g)] / 2$

[00345] For general cases with more than two contributors, the expected mixing fraction vector for allele 1 is computed as: $p \leftarrow G \cdot \beta$, which is implemented for nucleic acid mixtures with single, two, or multiple contributors.

[00346] When $\lambda = 0$, the implementation has the special case: $p \leftarrow g/2 \cdot \beta$

Overview of the DNA Extraction, PCR (Amplification), and Sequencing Models

[00347] Three probabilistic models (Figure 1C) are provided to model the three major components in the generic experimental pipeline (Figure 1A): 1) DNA/RNA extraction; 2) DNA/RNA amplification (e.g., PCR) as an approach for enriching target DNA/RNA; 3) sequencing (e.g., NGS sequencing). These and other modeling components are then integrated to implement the single-locus model and compute the single-locus likelihood function $M(n_{1i}, n_{2i} | p_{1i}, \theta)$.

[00348] The following notations are used in the mathematical models detailed in Table 4 and the remaining of this section.

[00349] $B()$: beta function

[00350] $Beta()$, $BN()$, $Pois()$, $Gamma()$: beta distribution, binomial distribution, Poisson distribution, and Gamma distribution

[00351] $NB()$ denotes a negative binomial distribution, which is a discrete probability distribution of the number of successes in a sequence of independent and identically distributed Bernoulli trials before a specified (non-random) number of failures (denoted r) occurs.

Table 4: Statistical models for the three major components in the generic experiment pipeline (Figure 1). The model for each component is conditioned on the previous component. The models are per each locus and locus index i is omitted.

gDNA or cfDNA extracted from a blood sample (Model E)	Copies of allele 1: $n_1'' \sim \text{Pois}(c \cdot p_1)$	Copies of allele 2: $n_2'' \sim \text{Pois}(c \cdot p_2)$
	Copies of allele 1 given the total copies of the locus: $n_1'' n'' \sim \text{BN}(n'', p_1)$, where $n'' = n_1'' + n_2''$.	
PCR amplified DNA (Model P)	Copies of allele 1: $n_1' \sim \text{Gamma}(n_1'' \cdot p, \theta)$	Copies of allele 2: $n_2' \sim \text{Gamma}(n_2'' \cdot p, \theta)$
	Fraction of allele 1 of the locus in the PCR product, conditioning on allele 1 and allele 2 copies in the extracted DNA: $n_1'/n' n_1'', n_2'' \sim \text{Beta}(n_1'' \cdot p, n_2'' \cdot p)$, where $n' = n_1' + n_2'$. Ignoring DNA sampling variation (hence $n_1'' = n'' \cdot p_1$, $n_2'' = n'' \cdot p_2$): $n_1'/n' \sim \text{Beta}(n'' \cdot p \cdot p_1, n'' \cdot p \cdot p_2)$	
Copies of sequenced reads mapped to the loci (without sequencing error) (Model S)	Copies of allele 1 given the total copies of the locus, conditioning on fraction of allele 1 of a locus in the PCR product: $n_1 n, n_1'/n' \sim \text{BN}(n, p = n_1'/n')$, where $n = n_1 + n_2$.	

DNA Extraction Model: Model E

[00352] When cfDNA or cellular DNA is extracted from a blood sample, the obtained DNA is a small sample from the large pool of DNA, and hence the implementations model the counts of two alleles at each locus as two Poisson distributions. Hence the DNA copy (n_1'') for

allele 1 at a locus conditioned on the total counts n follows the binomial distribution: $n_1 \sim \text{BN}(n, p_1)$, with mean $\mu_0 = n \cdot p_1$ and variance $\delta_0^2 = n \cdot p_1 \cdot p_2$.

[00353] When gDNA is extracted from a sample, the resulting gDNA amount for each locus can again be variable due to extraction losses. Viewing p_1 as the fraction of allele 1 in the input sample, the amount of allele 1 in the extracted DNA can again be modeled by a binomial distribution: $n_1 \sim \text{BN}(n, p_1)$.

PCR Amplification Model: Model P

[00354] We model the PCR amplification process as a stochastic process in order to obtain a probabilistic distribution of allele 1 counts in the PCR product. Let x_t be the DNA copies of a given allele after cycle t of PCR amplification, let r_t be the amplification rate for cycle t , and let y_t be the new copies generated at cycle t . By assuming each piece of DNA has a probability r_t of getting amplified and added to the DNA pool, the implementations have the following model for amplification:

[00355] $x_{t+1} = x_t + y_{t+1}$, where $y_{t+1} \sim \text{BN}(x_t, r_{t+1})$ follows a binomial distribution with x_t and r_{t+1} as parameters.

[00356] Based on this model, the implementations assume that the DNA copy number for a locus in the PCR product follows the Gamma distribution approximately. Below is the justification.

[00357] Step 1: Using Yule process (a continuous time stochastic process) to approximate PCR (a discrete time stochastic process).

[00358] The PCR process $x_{t+1} = x_t + y_{t+1}$, where $y_{t+1} \sim \text{BN}(x_t, r_{t+1})$ is a discrete time pure-birth process: in a given cycle of time t , each copy of DNA "gives birth" independently at some rate r_t . The continuous time version of the pure-birth process is well-known as the Yule-Furry Process. For the continuous time birth process, the final copy number for a locus at a given time t is known to follow a negative binomial distribution. The implementations can use the same distribution to approximate the discrete time birth process, when the total number of PCR cycles is not close to 1.

[00359] Step 2: Using Gamma distribution (a continuous distribution) to approximate negative binomial distribution (a discrete distribution).

[00360] A negative binomial random variable can be written as a sum of independent and identically distributed (i.i.d.) geometric random variables. The exponential distribution is known to be the continuous version of the geometric distribution. Hence, the sum of i.i.d. exponential random variables, which follows the Gamma distribution, is the continuous version of the sum of binomial random variable, which is negative binomial.

[00361] Below the implementations that estimate the parameters of the Gamma distributions of the allele counts in the PCR products.

[00362] Based on the law of total variance $\text{var}(x_{t+1}) = \text{var}(E(x_{t+1}|x_t)) + E(\text{var}(x_{t+1}|x_t))$, the implementations can derive the mean and variance of x_t recursively as follows:

$$[00363] \quad \mu_{t+1} = \mu_t \cdot (1+r_{t+1})$$

$$[00364] \quad \delta_{t+1}^2 = \mu_t \cdot r_{t+1} \cdot (1-r_{t+1}) + \delta_t^2 \cdot (1+r_{t+1})^2,$$

$$[00365] \quad \text{where } \mu_t = E(x_t), \delta_t^2 = \text{var}(x_t).$$

[00366] Assuming an average amplification rate per PCR cycle $r_{t+1} = r$, the implementations have

$$[00367] \quad \mu_t = \mu_0 \cdot (1+r)^t$$

$$[00368] \quad \delta_t^2 = \mu_0 \cdot (1+r)^t \cdot [(1+r)^t - 1] \cdot (1-r)/(1+r) + \delta_0^2 \cdot (1+r)^{2t}$$

[00369] Notice that μ_0 and δ_0^2 are the mean and variance of DNA allele counts in the PCR amplification input, and they can be computed based on the DNA extraction model (model E) described above. Alternatively, if the implementations do not treat cfDNA/cellular DNA allele counts as random variables, the implementations have $\mu_0 = n_1$ or n_2 , and $\delta_0^2 = 0$.

[00370] The corresponding gamma distribution $G(x_t | k, \theta) = x^{k-1} e^{-x/\theta} / [\theta^k \cdot \Gamma(k)]$ that matches this mean and variance has parameters:

$$[00371] \quad \theta = [(1+r)^t - 1] \cdot (1-r)/(1+r) + \delta_0^2 / \mu_0 \cdot (1+r)^t$$

$$[00372] \quad k = \mu_0 \cdot (1+r)^t / [(1+r)^t - 1] \cdot (1-r)/(1+r) + \delta_0^2 / \mu_0 \cdot (1+r)^t.$$

[00373] For a given locus with two alleles and two initial copies (n_1 , n_2), assuming identical amplification rate $r_1 = r_2 = r$ for two alleles for each locus, the two corresponding gamma distributions $G(n_1' | k_1, \theta_1)$ and $G(n_2' | k_2, \theta_2)$ have the following parameters:

$$[00374] \quad \theta_1 = [(1+r)^t - 1] \cdot (1-r)/(1+r) + p_2 \cdot (1+r)^t$$

$$[00375] \quad \theta_2 = [(1+r)^t - 1] \cdot (1-r)/(1+r) + p_1 \cdot (1+r)^t$$

$$[00376] \quad k_1 = n'' p_1 / [(1 - (1+r)^{-t}) \cdot (1-r)/(1+r) + p_2]$$

$$[00377] \quad k_2 = n'' p_2 / [(1 - (1+r)^{-t}) \cdot (1-r)/(1+r) + p_1].$$

[00378] When the implementations condition the PCR model on the DNA extraction model, s.t. $\mu_0 = n_1''$ or n_2'' and $\delta_0^2 = 0$, the implementations then have

$$[00379] \quad \theta_1 = [(1+r)^t - 1] \cdot (1-r)/(1+r)$$

$$[00380] \quad \theta_2 = [(1+r)^t - 1] \cdot (1-r)/(1+r)$$

$$[00381] \quad k_1 = n_1'' \cdot (1+r)/(1-r) / [1 - (1+r)^{-t}]$$

$$[00382] \quad k_2 = n_2'' \cdot (1+r)/(1-r) / [1 - (1+r)^{-t}].$$

[00383] Hence the allele copies n_1' and n_2' in the PCR product follow two Gamma distributions with identical scale parameters θ_1 and θ_2 , which are only dependent on the PCR process (the number of cycles and amplification rate). Therefore,

$$[00384] \quad n_1' / (n_1' + n_2') \sim \text{Beta}(n_1'' \cdot \rho, n_2'' \cdot \rho),$$

[00385] where ρ is a constant related to the amplification rate r , which is only dependent on the PCR process: $\rho = (1+r)/(1-r) / [1 - (1+r)^{-t}]$, or approximately $\rho = (1+r)/(1-r)$ when the number of cycles t is large. For a specific locus, this is written as $n_{1i}' / (n_{1i}' + n_{2i}') \sim \text{Beta}(n_{1i}'' \cdot \rho_i, n_{2i}'' \cdot \rho_i)$, to capture the locus specific PCR amplification rate.

[00386] If the implementations ignore DNA sampling and assume all loci have the same total DNA copy number $n_i'' = n''$, then $n_{1i}'' = n'' \cdot p_{1i}$ and $n_{2i}'' = n'' \cdot p_{2i}$. The allele fraction for a locus in the PCR product follows:

$$[00387] \quad n_{1i}' / (n_{1i}' + n_{2i}') \sim \text{Beta}(n'' \cdot \rho_i \cdot p_{1i}, n'' \cdot \rho_i \cdot p_{2i}).$$

[00388] Note that without the Gamma distribution approximation, the allele counts of PCR products have $n_1' \sim \text{NB}(r_1, p)$ and $n_2' \sim \text{NB}(r_2, p)$, and the ratio $n_1'/(n_1' + n_2')$ has no closed form distribution. With the Gamma distribution approximation, $n_1' \sim \text{Gamma}(n_1'' \cdot \rho, \theta)$ and $n_2' \sim \text{Gamma}(n_2'' \cdot \rho, \theta)$, and $n_1'/(n_1' + n_2')$ follows the beta distribution.

Sequencing Read Count Model: Model S

[00389] NGS sequencing is a process that samples from the pool of DNA molecules supplied to the sequencer and reads out the sequences of these molecules. The fraction of allele 1 for a locus i in the PCR product is $n_{1i}' / (n_{1i}' + n_{2i}')$. This fraction determines the probability that allele 1 reads occur in the sequencing results. Conditioning on n_i , the total number of reads per locus, the distribution of n_{1i} , the allele 1 read count of a locus, is then modeled as a binomial distribution $n_{1i} \sim \text{BN}(n_i, n_{1i}' / (n_{1i}' + n_{2i}'))$.

Modeling the Genetic Relationship between the Contributors as a Prior Distribution

[00390] If the contributor genotypes are completely known, they can be directly incorporated (using Table 2 or Table 3) as parameters of the component models described above. However, when the genotypes are unknown, the implementations make use of the genetic-relationship information between the donor and recipient in a two-contributor setting to achieve accurate mixture quantification. Genetic relationship is commonly available in clinical applications such as organ transplant. Here we present the implementation for two-contributor scenario, but this “genetic prior” approach can be generalized to any number of contributors.

[00391] We formulate different types of donor-recipient relationships as distinct prior distributions on the space of possible genotype combinations of the donor (contributor 2) and recipient (contributor 1). Assuming Hardy-Weinberg equilibrium, the genotype distribution for a given loci for a single individual is $P(g = [0, 1, 2]) = [(1-\pi)^2, 2\pi(1-\pi), \pi^2]$, where π is the population frequency of allele 1, g is the allele 1 copy number. Notice that all genetic relationships are the results of parent-child relationships. Based on the genetic-relationships between parent and child for a give biallelic marker site (Table 5), the implementations can compute the joint distribution for any genetic relationship among two or multiple contributors.

Table 5: Probability distribution of child genotype given parents' genotypes (father genotype g_{Father} and mother genotype g_{Mother}) for a given locus, as well as the joint distribution between father and mother assuming they are not relatives.

g_{Father}	g_{Mother}	Child probability for genotype $[0, 1, 2]$ respectively conditioned on parent genotypes	$P(g_{\text{Father}}, g_{\text{Mother}})$
0	0	$[1, 0, 0]$	$(1-\pi)^4$
0	1	$[1/2, 1/2, 0]$	$2\pi(1-\pi)^3$

g _{Father}	g _{Mother}	Child probability for genotype [0, 1, 2] respectively conditioned on parent genotypes	P(g _{Father} , g _{Mother})
0	2	[0, 1, 0]	$\pi^2(1-\pi)^2$
1	0	[1/2, 1/2, 0]	$2\pi(1-\pi)^3$
1	1	[1/4, 1/2, 1/4]	$4\pi^2(1-\pi)^2$
1	2	[0, 1/2, 1/2]	$2\pi^3(1-\pi)$
2	0	[0, 1, 0]	$\pi^2(1-\pi)^2$
2	1	[0, 1/2, 1/2]	$2\pi^3(1-\pi)$
2	2	[0, 0, 1]	π^4

[00392] The prior distributions for various types of genetic-relationship between two contributors are further provided below.

Joint Distribution between Father and Child Genotypes

[00393] As an example, the Father-Child donor-recipient genotype (GT) joint distribution is computed using the following formula:

[00394] $P(\text{Recipient} = \text{Me GT}, \text{Donor} = \text{Father GT}) = \sum_{\text{mother GT}} [P(\text{Me GT} | \text{Father GT}, \text{Mother GT}) \cdot P(\text{Father GT}, \text{Mother GT})],$

[00395] where values of $P(\text{Me GT} | \text{Father GT}, \text{Mother GT})$ and $P(\text{Father GT}, \text{Mother GT})$ are taken from the Table 5 columns 3 and 4 respectively.

Joint Distribution between Sibling Genotypes

[00396] As an example, the Me-Sibling donor-recipient genotype joint distribution is computed using the following formula, based on the conditional independence of two sibling genotypes given parents genomes:

[00397] $P(\text{Recipient} = \text{Me GT}, \text{Donor} = \text{Sibling GT}) = \sum_{\text{mother GT}} \sum_{\text{father GT}} [P(\text{Me GT} | \text{Father GT}, \text{Mother GT}) \cdot P(\text{Sibling GT} | \text{Father GT}, \text{Mother GT}) \cdot P(\text{Father GT}, \text{Mother GT})],$

[00398] Where values of $P(\text{Me GT} | \text{Father GT}, \text{Mother GT})$, $P(\text{Sibling GT} | \text{Father GT}, \text{Mother GT})$, and $P(\text{Father GT}, \text{Mother GT})$ are taken from the Table 5 columns 3, column 3, and column 4 respectively.

Joint Distribution between Uncle-Nephew Genotypes

[00399] As an example, the Uncle/Aunt-Nephew/Niece donor-recipient genotype joint distribution is computed using the following formula:

[00400] $P(\text{Recipient} = \text{Me GT}, \text{Donor} = \text{Uncle GT})$

[00401] $= \sum_{\text{grandmother GT}} \sum_{\text{grandfather GT}} \sum_{\text{mother GT}} \sum_{\text{father GT}} [P(\text{Me GT} | \text{Father GT}, \text{Mother GT}) \cdot P(\text{Mother GT}) \cdot P(\text{Father GT} | \text{GrandFather GT}, \text{GrandMother GT}) \cdot P(\text{Uncle GT} | \text{GrandFather GT}, \text{GrandMother GT}) \cdot P(\text{GrandFather GT}, \text{GrandMother GT})]$

[00402] $= \sum_{\text{mother GT}} \sum_{\text{father GT}} P(\text{Me GT} | \text{Father GT}, \text{Mother GT}) \cdot P(\text{Mother GT}) \cdot P(\text{Father GT}, \text{Uncle GT}),$

[00403] where values of $P(\text{Me GT} | \text{Father GT}, \text{Mother GT})$ is taken from column 3 of table 5, and $P(\text{Father GT}, \text{Uncle GT})$ is the same as $P(\text{Recipient} = \text{Me GT}, \text{Donor} = \text{Sibling GT})$.

[00404] In matrix notation, this can be computed using the parent/child prior matrix, the sibling prior matrix, and the single genome prior vector

[00405] $= [P(\text{Me GT}, \text{Father GT})]_{\text{Me}, \text{Father}} \cdot \text{diag}(1 / [P(\text{Father GT})]_{\text{Father}}) \cdot [P(\text{Father GT}, \text{Uncle GT})]_{\text{Father}, \text{Uncle}}$

Joint Distribution between Cousin Genotypes

[00406] Assuming cousin is genetically linked by their fathers, who are brother, and the mothers are genetically unrelated, then,

[00407] $P(\text{Recipient} = \text{Me GT}, \text{Donor} = \text{Cousin GT})$

[00408] $= \sum_{\text{aunt GT}} \sum_{\text{uncle GT}} \sum_{\text{mother GT}} \sum_{\text{father GT}} P(\text{Me GT} | \text{Father GT}, \text{Mother GT}) \cdot P(\text{Mother GT}) \cdot P(\text{Father GT}, \text{Uncle GT}) \cdot P(\text{Aunt GT}) \cdot P(\text{Cousin GT} | \text{Uncle GT}, \text{Aunt GT})$

[00409] $= \sum_{\text{aunt GT}} \sum_{\text{uncle GT}} P(\text{Me GT}, \text{Uncle GT}) \cdot P(\text{Aunt GT}) \cdot P(\text{Cousin GT} | \text{Uncle GT}, \text{Aunt GT})$

[00410] $= \sum_{\text{uncle GT}} P(\text{Me GT}, \text{Uncle GT}) \cdot P(\text{Cousin GT}, \text{Uncle GT}) / P(\text{Uncle GT})$

[00411] In matrix notation, this can be computed using the uncle/niece prior matrix, the parent/child prior matrix, and the single genome prior vector

[00412] $= [P(\text{Me GT, Uncle GT})]_{\text{Me, Uncle}} \cdot \text{diag}(1 / [P(\text{Uncle GT})]_{\text{Uncle}}) \cdot [P(\text{Cousin GT, Uncle GT})]_{\text{Uncle, Cousin}}$

[00413] Notice that $P(\text{Cousin GT, Uncle GT})$ is the same as parent-child relationship.

Joint Distribution between Half Sibling Genotypes

[00414] Assuming half sibling is linked by a single mother, and the two fathers are unrelated:

[00415] $P(\text{Recipient} = \text{Me GT, Donor} = \text{HafSib GT})$

[00416] $= \sum_{\text{Father GT}} \sum_{\text{Mother GT}} \sum_{\text{StepFather GT}} P(\text{Me GT} | \text{Father GT, Mother GT}) \cdot P(\text{HalfSib GT} | \text{StepFather GT, Mother GT}) \cdot P(\text{Mother GT}) \cdot P(\text{Father GT}) \cdot P(\text{StepFather GT})$

[00417] $= \sum_{\text{Mother GT}} P(\text{Me GT, Mother GT}) \cdot P(\text{HalfSib GT, Mother GT}) / P(\text{Mother GT})$

[00418] In matrix notation, this can be computed using the two parent child prior matrix, and the single genome prior vector

[00419] $= [P(\text{Me GT, Mother GT})]_{\text{Me, Mother}} \cdot \text{diag}(1 / [P(\text{Mother GT})]_{\text{Mother}}) \cdot [P(\text{HalfSib GT, Mother GT})]_{\text{HalfSib, Mother}}$

[00420] Note that under HardyWeinburg equilibrium, half sibling relationship follows the same distribution as the uncle/aunt/nephew/niece relationship. This may not be true without HardyWeinburg equilibrium.

Summary

[00421] The results from the above derivations is summarized in Table 6, and the specific instances given population SNP allele frequency $\pi = 0.5$ is provided in Table 7. Additional relationships, such as grandparent-grandchild relationship or multi-contributor relationship, can be derived based on the same underlying principle.

Table 6: Prior distributions $P(g_{11}, g_{21})$ of related or unrelated genomes. Assuming all SNPs are from autosomes, all married couples are genetically-unrelated, and in Hardy Weinberg equilibrium. g_{11} is the recipient genome, g_{21} is the donor genome.

g_{11}	g_{21}	Donor's relationship to Recipient					
		Parent or	Sibling	Half Siblings	Uncle/Aunt or Nephew/Niece	Cousin	Unrelated

		Child					
0	0	$(1-\pi)^3$	$(1-\pi)^2(1-\pi/2)^2$	$(1-\pi)^3(1-\pi/2)$	$(1-\pi)^3(1-\pi/2)$	$(1-\pi)^3(1-3\pi/4)$	$(1-\pi)^4$
0	1	$\pi(1-\pi)^2$	$\pi(1-\pi)^2(1-\pi/2)$	$\pi(1-\pi)^2(3/2-\pi)$	$\pi(1-\pi)^2(3/2-\pi)$	$\pi(1-\pi)^2(7/4-3\pi/2)$	$2\pi(1-\pi)^3$
0	2	0	$\pi^2(1-\pi)^2/4$	$\pi^2(1-\pi)^2/2$	$\pi^2(1-\pi)^2/2$	$3/4\pi^2(1-\pi)^2$	$\pi^2(1-\pi)^2$
1	0	$\pi(1-\pi)^2$	$\pi(1-\pi)^2(1-\pi/2)$	$\pi(1-\pi)^2(3/2-\pi)$	$\pi(1-\pi)^2(3/2-\pi)$	$\pi(1-\pi)^2(7/4-3\pi/2)$	$2\pi(1-\pi)^3$
1	1	$\pi(1-\pi)$	$\pi(1-\pi)[1+\pi-\pi^2]$	$\pi(1-\pi)[1/2+2\pi-2\pi^2]$	$\pi(1-\pi)[1/2+2\pi-2\pi^2]$	$\pi(1-\pi)[1/4+3\pi-3\pi^2]$	$4\pi^2(1-\pi)^2$
1	2	$\pi^2(1-\pi)$	$\pi^2(1-\pi)(1/2+\pi/2)$	$\pi^2(1-\pi)(1/2+\pi)$	$\pi^2(1-\pi)(1/2+\pi)$	$\pi^2(1-\pi)(1/4+3\pi/2)$	$2\pi^3(1-\pi)$
2	0	0	$\pi^2(1-\pi)^2/4$	$\pi^2(1-\pi)^2/2$	$\pi^2(1-\pi)^2/2$	$3/4\pi^2(1-\pi)^2$	$\pi^2(1-\pi)^2$
2	1	$\pi^2(1-\pi)$	$\pi^2(1-\pi)(1/2+\pi/2)$	$\pi^2(1-\pi)(1/2+\pi)$	$\pi^2(1-\pi)(1/2+\pi)$	$\pi^2(1-\pi)(1/4+3\pi/2)$	$2\pi^3(1-\pi)$
2	2	π^3	$\pi^2(1/2+\pi/2)^2$	$\pi^3(1/2+\pi/2)$	$\pi^3(1/2+\pi/2)$	$\pi^3(1/4+3\pi/4)$	π^4

Table 7: Prior distributions $P(g_{11}, g_{21})$ of related or unrelated genomes given SNP population allele frequency $\pi = 0.5$.

g_{11}	g_{21}	Donor's relationship to Recipient					
		Parent or Child	Sibling	Half Sibling	Uncle/Aunt or Nephew/Niece	Cousin	Unrelated
0	0	8/64	9/64	6/64	6/64	5/64	4/64
0	1	8/64	6/64	8/64	8/64	8/64	8/64
0	2	0	1/64	2/64	2/64	3/64	4/64
1	0	8/64	6/64	8/64	8/64	8/64	8/64
1	1	16/64	20/64	16/64	16/64	16/64	16/64
1	2	8/64	6/64	8/64	8/64	8/64	8/64
2	0	0	1/64	2/64	2/64	3/64	4/64
2	1	8/64	6/64	8/64	8/64	8/64	8/64
2	2	8/64	9/64	6/64	6/64	5/64	4/64

[00422] The distributions for parent-child and sibling relationship are quite different from unrelated, while uncle/aunt-nephew/niece are close to unrelated. In the case when the donor genotype is unknown, the implementations can infer the genetic relationship by evaluating the likelihood function of fitted models of each of the above genetic relationships. Alternatively, the

implementations can allow multiple free parameters in the genetic priors distribution (with additional constraints that the marginal distributions should follow Hardy-Weinberg equilibrium), and estimate these parameters together with the estimation of donor fraction.

Adjustment of DNA Copy Numbers based on DNA Length

[00423] For an amplicon-based assay that involves PCR DNA amplification, the DNA length impacts the effectiveness of the DNA as PCR template. In the extreme, when DNA fragments are shorter than the intended amplicon length, they are 0% effective as PCR template. To correct for this effect, we used the following procedure to adjust the DNA copy numbers using the average DNA length, which varies depending on the type of the input DNA. Some implementations adjust the effective input DNA molecule number based on the average length of input DNA template. In some implementation, the effective input DNA molecule number is adjusted according to the equation below:

[00424]
$$n'' = w/w_0 \cdot (L - L_a + 1)/L,$$

[00425] where n'' is the effective input DNA molecule number (haploid), w is the input DNA amount, w_0 (3.59×10^3 ng/copy) is the weight of haploid human genome, L is the average length of input DNA template, and L_a is the average amplicon length (110bp for our amplicon design).

[00426] DNA template efficiency is defined as $e = (L - L_a + 1)/L$, which is defined for $L \geq L_a$. Table 8 shows example DNA types and their efficiency as PCR templates.

Table 8. DNA type and their efficiency as PCR template

DNA type	DNA length parameter (L)	Template Efficiency (e)
genomic DNA (gDNA)	100,000	0.9989
cell free DNA (cfDNA)	165	0.3394
mock cfDNA (mcfDNA)	160	0.3188

Integration of the Modeling Components

[00427] The components of the probabilistic mixture model are integrated to provide a solution to the contributor DNA quantification (CDQ) problem. The population allele frequency π for each SNP site can be obtained from public databases such as dbSNP. If one selects the most informative SNP markers, i.e. SNPs with $\pi = 0.5$, in an experiment design, one can set $\pi = 0.5$ for all loci and let $P(g_{11}, g_{21})$ be the genetic-relationship prior distribution as described in the previous section.

[00428] On a schematic level, Figure 2B shows a block diagram illustrating various components of the probabilistic mixture model 250. Some components are optional in some implementations. The probabilistic mixture model 250 includes a binomial distribution 258 for modeling allelic counts of sequencing reads. In some implementations, the probabilistic mixture model also includes a component for modeling donor-donee (or recipient) relationship using a genetic relationship prior distribution 252. In some implementations, the probabilistic mixture model also includes a binomial distribution 254 for modeling DNA extraction allelic counts. In some implementations, the probabilistic mixture model 250 also includes a beta distribution 256 for modeling PCR product or amplification product allelic fraction. See block 256.

[00429] In some implementations, the mixture model combines the binomial distribution 208 with binomial distribution 254 to model both the DNA extraction errors and sequencing errors. In such implementations, the mixture model uses a beta-binomial distribution 260 to model the allelic counts of sequencing reads while capturing variability in the allelic counts due to DNA extraction.

[00430] In some implementations, the probabilistic mixture model 250 combines beta distribution 256 and binomial distribution 258, and uses a beta-binomial distribution 262 to model both errors in the PCR or amplification process and errors of sequencing process.

[00431] In some implementations, the probabilistic mixture model 250 combines binomial distribution 254, beta distribution 256, and binomial distribution 258 to account for variance resulting from DNA extraction, amplification process, and sequencing process, respectively. In such implementations, probabilistic mixture model 200 first uses a beta distribution 264 to approximate the effects of binomial distribution 254 and beta distribution 256. The probabilistic mixture model 250 then combines beta distribution 264 and binomial distribution 258 using beta-binomial distribution 256.

The Sequencing Model: Model S

[00432] A basic version of the full model ignores the DNA extraction model and the PCR model, and only considers the sequencing model. For each locus, the sequencing read count for the reference allele is modeled by a binomial distribution (Figure 1C), $n_{1i} \sim \text{BN}(n_i, p_{1i})$, where the value of parameter $p_{1i}(g_{1i}, g_{2i}, \lambda, \beta)$ is a function on the donor-recipient genotype combination for the loci (Table 2 and Table 3). Given that the genotypes are unknown, the implementations marginalize over the 9 possible genotype combinations for each locus with $P(g_{1i}, g_{2i} | \pi)$ as prior distribution (Table 6 and Table 7). The complete likelihood function across all loci is the product of the marginal distributions for all loci:

[00433] $L(\beta, \lambda, \pi ; n_1, n_2) = \prod_i \sum_{g_{1i}g_{2i}} \text{BN}(n_{1i} | n_i, p_{1i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)$, where $L(\beta, \lambda, \pi ; n_1, n_2)$ is the likelihood of observing allele count vectors n_1 to n_2 for alleles 1 and 2 given parameters β and π ; $p_{1i}(g_{1i}, g_{2i}, \lambda, \beta)$ is a probability parameter, taken as p_1' from Table 3, indicating a probability of allele 1 at locus i based on the two contributors' genotypes (g_{1i}, g_{2i}); and $P(g_{1i}, g_{2i} | \pi)$ is a prior joint probability of observing the two contributors' genotypes given a population allele frequency (π).

[00434] Expanding it to multiple contributors, the likelihood function can be expressed as:

[00435] $L(\beta, \lambda, \pi ; n_1, n_2) = \prod_i [\sum_{g_i} \text{BN}(n_{1i} | n_i, p(g_i, \lambda, \beta)) \cdot P(g_i | \pi)]$

The Extraction-Seq Compound Model: Model ES

[00436] A more advanced model combines the DNA extraction model as well as the Sequencing model. The implementations ignore the PCR step (i.e. assume that, for each locus, the allele fraction in the PCR product is the same as the allele fraction in the DNA sample), and only model DNA sampling and sequencing steps. For each locus, there is a binomial distribution for the allele counts in the input DNA sample. This captures the locus-to-locus variability of the allele fractions in the input DNA provided to the NGS sequencing.

[00437] For the DNA extraction model, the implementations have $n_{1i}'' \sim \text{BN}(n'', p_{1i})$, while conditioning on the DNA extraction model, the sequencing model is $n_{1i} | n_{1i}'', n'' \sim \text{BN}(n_i, n_{1i}''/n'')$, where $n_i'' = n''$ is the copies of haploid genomes the input DNA correspond to. Unfortunately, the marginal distribution of n_{1i} has no closed form formula. The implementations approximate the distribution of n_{1i}''/n'' with a beta distribution $\text{Beta}(a, b)$, and the best Beta distribution is selected

by matching the mean and variance of n_{li}''/n'' with those derived from the binomial model $n_{li}'' \sim \text{BN}(n'', p_{li})$:

[00438] $p_{li} = a/(a+b)$

[00439] $p_{li} \cdot (1-p_{li})/n'' = ab/(a+b)^2/(a+b+1).$

[00440] Solving the equations gives the beta distribution $\text{Beta}((n''-1)p_{li}, (n''-1)p_{2i})$ as the best approximation. With this approximation to the DNA extraction model, the marginal distribution of n_{li} then follows a beta-binomial distribution of the form:

[00441] $n_{li} \sim \text{BB}(n_i, (n''-1) \cdot p_{li}, (n''-1) \cdot p_{2i}).$

[00442] Or in an alternative approximation:

[00443] $n_{li} \sim \text{BB}(n_i, n'' \cdot p_{li}, n'' \cdot p_{2i}).$

[00444] The corresponding full likelihood function considering the genetic-relationship prior is then:

[00445] $L(\beta, n'', \lambda, \pi; n_1, n_2) = \prod_i [\sum g_i \text{BB}(n_{li} | n_i, (n''-1) \cdot p_{li}, (n''-1) \cdot p_{2i}) \cdot P(g_i | \pi)]$

[00446] wherein $L(\beta, n'', \lambda, \pi; n_1, n_2)$ is the likelihood of observing allele count vectors n_1 and n_2 for alleles 1 and 2 at all loci, and $p_{li} = p(g_i, \lambda, \beta)$, $p_{2i} = 1 - p_{li}$.

[00447] Notice that both n'' and $\pi = 0.5$ are known parameters, and the final full likelihood function has only a single unknown parameter β , the donor DNA fraction.

[00448] The input DNA (haploid) copy numbers n'' can be derived from the input DNA mass. When input DNA amount is 8ng, $n'' = 8 \text{ ng} / [3.59 \times 10^{-3} \text{ ng/copy}] = 2228.412$.

PCR-Seq Compound Model: Model PS

[00449] Ignoring the DNA extraction model, and assuming a known genotype combination for a given locus, then the PCR model: $n_{li}' / (n_{li}' + n_{2i}') \sim \text{Beta}(n'' \cdot p_i \cdot p_{li}, n'' \cdot p_i \cdot p_{2i})$ and Sequencing model $n_{li} \sim \text{BN}(n_i, n_{li}' / (n_{li}' + n_{2i}'))$ can be combined into the beta-binomial distribution: $\text{BB}(n_i, n'' \cdot p_i \cdot p_{li}, n'' \cdot p_i \cdot p_{2i})$. Notice that both the underlying loci specific PCR amplification rates p_i are unknown. If the implementations assume all loci have the same inherent amplification rate, then the implementations have, $\text{BB}(n_i, c \cdot p_{li}(g_{11}, g_{21}, \beta), c \cdot p_{2i}(g_{11}, g_{21}, \beta))$.

[00450] The complete likelihood model across all loci is then: $L(\beta, n'', c, \lambda, \pi; n_1, n_2) = \prod_i [\sum_{g_i} BB(n_{1i} | n_i, c \cdot p_{1i}, c \cdot p_{2i}) \cdot P(g_i | \pi)]$, where c and β are two parameters to be estimated.

[00451] Alternatively, the implementations can define the relative amplification rate of each locus to be proportional to the total reads per locus, and re-parameterize the beta-binomial as $n_{1i} \sim BB(n_i, c' \cdot n_i \cdot p_{1i}, c' \cdot n_i \cdot p_{2i})$, where c' is a parameter to be optimized; and n_i is the total reads at locus i .

[00452] The complete likelihood model across all loci is then: $L(\beta, n'', c', \lambda, \pi; n_1, n_2) = \prod_i [\sum_{g_i} BB(n_{1i} | n_i, c' \cdot n_i \cdot p_{1i}, c' \cdot n_i \cdot p_{2i}) \cdot P(g_i | \pi)]$, where c' and β are two parameters to be estimated

Extraction-PCR-Seq Compound Model: Model EPS

[00453] All three components in the Extraction-PCR-sequencing generic experimental pipeline can be modeled together by a beta-binomial if the implementations combine DNA extraction and PCR models into one model and approximate it by a single beta distribution. Intuitively, although the expected value of allele 1 fraction in the PCR product (n_{1i}/n' , see Table 4) remains p_{1i} , the uncertainty (variance) of n_{1i}/n' originates from both the DNA extraction and the PCR steps. To obtain a beta distribution $\text{beta}(a, b)$ to model DNA extraction and PCR together, the implementations compute the unconditional mean and variance of n_{1i}/n' based on the following laws: $E(n_{1i}/n') = E(E(n_{1i}/n_i' | n_{1i}''/n''))$, and $\text{var}(n_{1i}/n') = \text{var}(E(n_{1i}/n_i' | n_{1i}''/n'')) + E(\text{var}(n_{1i}/n_i' | n_{1i}''/n''))$. This gives: $E(n_{1i}/n') = p_{1i}$, and $\text{var}(n_{1i}/n') = p_{1i}p_{2i}/n'' + p_{1i}p_{2i}/(n'' \cdot \rho_i + 1) - p_{1i}p_{2i}/[n'' \cdot (n'' \cdot \rho_i + 1)]$, where $\rho_i = (1+r_i)/(1-r_i) > 1$ is the constant related to the amplification rate r_i . Since n'' is large, the implementations have the following approximation $\text{var}(n_{1i}/n') = p_{1i}p_{2i}/[n'' \cdot (1+r_i)/2]$. The best beta distribution that models DNA extraction and PCR is then $\text{Beta}([n'' \cdot (1+r_i)/2 - 1]p_{1i}, [n'' \cdot (1+r_i)/2 - 1]p_{2i})$. Notice this is close to the beta distribution for cfDNA/gDNA extraction $\text{Beta}((n''-1)p_{1i}, (n''-1)p_{2i})$, yet the variance is now larger. For a typical PCR reaction with $r_i = 0.8$ to 0.95 , the implementations have $n'' \cdot (1+r_i)/2 = 0.9 \cdot n''$ to $0.975 \cdot n''$.

[00454] The full multiple-loci likelihood function for cfDNA-PCR-Seq model is:

$$L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i \sum_{g_{1i} g_{2i}} BB(n_{1i} | n_i, n'' \cdot (1+r_i)/2 \cdot p_{1i}(g_{1i}, g_{2i}, \lambda, \beta), n'' \cdot (1+r_i)/2 \cdot p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)$$

Baseline Method: NaiveLM or KGT.NaiveLM

[00455] A conventional method for quantifying DNA fractions of contributors uses a basic linear regression formulation, which does not use the same probabilistic model or cost functions described above. Instead, its cost function is expressed as:

[00456] $E = [r - p]^T \cdot [r - p]$, where r is the allele fraction vector, $p = g/2 \cdot \beta$ is the expected allele fraction vector, g is the genotype matrix, and β is the contributor DNA fraction vector. The naïve method is only applicable when all base lines are known.

Method for Estimating Contributor Nucleic Acid Fractions and Their Confidence Intervals

Numerical Optimization for Estimating Contributor DNA Fractions

[00457] The contributor DNA fraction β is estimated as the value that maximize the full likelihood function $L(n_1, n_2 | \beta)$. As mentioned above, although DNA is referred to in this and other examples, RNA and other nucleic acid molecules may be processed and analyzed similarly. Also, although the examples refer to nucleic acid mixture samples, the sample may include only a single contributor's nucleic acid, in which case the contributor fraction would be estimated as 1 or within a margin of error from 1.

[00458] During the calculation of $L(n_1, n_2 | \beta)$, multiple small probabilities values are multiplied. To avoid numerical underflowing when multiplying small probabilities, the implementations perform all summation and multiplications on log scale. The sum of small probability on log scale is performed as following: 1) obtain the max of the log probabilities as x_{max} ; 2) subtract all the log probabilities by the max; 3) exponentiate and then sum the resulting values; 4) log transform the resulting sum; 5) add back the max of the log probabilities. $\log(\exp(x_1 - x_{max}) + \exp(x_2 - x_{max}) + \dots + \exp(x_n - x_{max})) + x_{max}$.

[00459] To ensure positive contributor fractions within 0 to 1, the logit transformation $\beta = 1/(1+e^{-\eta})$ is used.

[00460] A novel numerical optimization computer strategy that seamlessly integrating iterative grid search with Broyden–Fletcher–Goldfarb–Shanno (BFGS) quasi-Newton method is implemented as described below.

[00461] Step 1: A grid initialization method generates even grids in $N-1$ dimensional space, where N is the number of contributors. In applications with only two contributors, to ensure global optimization and avoiding local optimums, the full likelihood function is initialized with $\beta_0 = 1/(1+e^{-\eta_0})$, where η_0 is the value among $-10, -9.9, -9.8, \dots, -0.1, 0$ that maximizes $L(n_1, n_2 | \beta_0 = 1/(1+e^{-\eta_0}))$ for two contributor cases. In applications with for multi-contributor cases, β is transformed using softmax, and then initialized over a high dimensional grid.

[00462] Step 2: An exhaustive search on the grid is performed to identified mixture fractions that minimizes $-\log_2(L)$.

[00463] Step 3: Initializing using the identified mixture fractions, numerical optimization of η is then performed using Broyden–Fletcher–Goldfarb–Shanno (BFGS) quasi-Newton method to minimize $-\log_2(L)$. Record the optimized mixture fraction as well as the convergence.

[00464] Step 4: Hessian matrix of $-\log_2(L)$ is computed using numerical differentiation on the identified mixture fractions.

[00465] Step 5: Errors and confidence interval around computed mixture fractions is determined based on the inverse of the hessian matrix. Meanwhile, determine if the hessian matrix is positive semi-definite.

[00466] Step 6: If BFGS optimization did not converge or if the hessian matrix is not positive semi-definite, then the procedure is configured for a next iteration of optimization. Otherwise optimization complete.

[00467] Step 7: When next iteration of optimization is to be performed, a finer $N-1$ dimensional grid is constructed covering 2^{N-1} original grids around the previously determined η , which corresponds to the estimated mixture fractions. The procedure then loops back to step 2 for next iteration of grid search and BFGS optimization.

[00468] The totality of these steps cannot be performed by human experts manually or in their heads. Instead, one or more computers are required to perform these steps.

Iterative Strategy for Model S with Known Genotypes (KGT.IterLM)

[00469] In some implementations, the single-locus likelihood function comprises a binomial distribution and the multiple-loci likelihood function is as follows: $L(\beta, \lambda, \pi; n_1, n_2) = \prod_i [\sum_{g_i} BN(n_{1i} | n_i, p(g_i, \lambda, \beta)) \cdot P(g_i | \pi)]$

[00470] In some implementations, the contributors include two contributors and the likelihood function is: $L(\beta, \lambda, \pi; n_1, n_2) = \prod_i \sum_{g_{1i}, g_{2i}} BN(n_{1i} | n_i, p_{1i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)$

[00471] where $L(\beta, \lambda, \pi; n_1, n_2)$ is the likelihood of observing allele count vectors n_1 to n_2 for alleles 1 and 2 given parameters β and π ; $p_{1i}(g_{1i}, g_{2i}, \lambda, \beta)$ is a probability parameter, taken as p_{1i}' from Table 3, indicating a probability of allele 1 at locus i based on the two contributors' genotypes (g_{1i}, g_{2i}); and $P(g_{1i}, g_{2i} | \pi)$ is a prior joint probability of observing the two contributors' genotypes given a population allele frequency (π).

[00472] In some implementations, the genotypes of all contributors are known, and the likelihood functions are expressed as $L(\beta, \lambda; n_1, n_2) = \prod_i n_i! / (n_{1i}! n_{2i}!) p_{1i}^{n_{1i}} p_{2i}^{n_{2i}}$, where $p_{ai} = \sum_{d=0 \dots D-1} g_{dai} \cdot \beta_d / [\sum_{d=0 \dots D-1} \beta_d \cdot (\sum_{a=1,2} g_{dai})]$. If all markers are on somatic chromosomes, then $p_{ai} = 1/2 \cdot \sum_{d=0 \dots D-1} g_{dai} \cdot \beta_d$. In matrix notation, this is $\mathbf{p} \leftarrow \mathbf{g}/2 \cdot \boldsymbol{\beta}$.

[00473] The iterative weighted linear regression method is developed by constructing a cost function that has the same gradient as that of $\log[L(\beta; n_1, n_2)]$ in each iteration when $\boldsymbol{\beta} = \boldsymbol{\beta}_0$:

[00474] $E = 1/2 \cdot \sum_i n_i / [p_{1i}(\boldsymbol{\beta}_0) \cdot (1 - p_{1i}(\boldsymbol{\beta}_0))] \cdot (n_{1i}/n_i - p_{1i}(\boldsymbol{\beta}))^2$.

[00475] In matrix notation, this is $E = 1/2 \cdot (\mathbf{r} - \mathbf{p})^T \cdot \mathbf{W}^2 \cdot (\mathbf{r} - \mathbf{p})$, where $\mathbf{W} = \text{diag}[(\mathbf{n} / (\mathbf{p}_0 \cdot (1 - \mathbf{p}_0)))^{1/2}]$ is a diagonal matrix, and $\mathbf{p}_0 = \mathbf{g}/2 \cdot \boldsymbol{\beta}_0$.

[00476] Iterative weighted linear regression is carried out by executing the following steps, given inputs: \mathbf{r} , \mathbf{n} , \mathbf{g} , and λ

[00477] Step 1. Initialize $\boldsymbol{\beta}$ as a uniform length D probability vector $\boldsymbol{\beta} \leftarrow [1/D]_D$

[00478] Step 2. Compute error correction of genotype matrix \mathbf{g} : $\mathbf{G} \leftarrow [(1 - \lambda) \mathbf{g} + \lambda (2 - \mathbf{g})]/2$

[00479] Step 3: Repeat Step a - Step e until convergence

- [00480] Step a:. Update expected allele 1 fraction using the previous computed contributor fraction: $p \leftarrow G \cdot \beta$
- [00481] Step b. Compute the weights for weighted regression $W \leftarrow \text{diag}([n / (p \cdot (1 - p))]^{1/2})$
- [00482] Step c. Solve the weighted linear regression: $\beta \leftarrow (W \cdot G)^{-1} \cdot (W \cdot r)$
- [00483] Step d. Ensure non-negativity: for each contributor i, $\beta_i \leftarrow \max(\beta_i, 0)$
- [00484] Step e. Normalization to probability vector: $\beta \leftarrow \beta / \sum_i \beta_i$ — normalization

Estimate the Confidence Interval

[00485] The lower bound of the confidence interval of the estimates are determined based on the Cramer-Rao inequality: $\text{var}(\theta_{\text{ML}}) \geq 1/I(\theta_{\text{ML}})$, where θ_{ML} is the maximum likelihood estimate of parameter θ , and $I(\theta_{\text{ML}})$ is fisher's information at θ_{ML} . Based on this, one can estimate the variance of β and c in the above described likelihood functions. The standard error is estimated as $\text{sqrt}(1/H)$ following the Cramér–Rao bound, where H is the Hessian matrix which can be approximated and is estimated in the BFGS - quasi-Newton method.

[00486] We use the following reparameterizations during the numerical optimization to estimate β and c ,

[00487] $\beta = 1/(1+e^{-\eta})$,

[00488] $c = e^{\kappa}$.

[00489] Let $I(\eta)$ and $I(\kappa)$ be the Fisher's information under parameterization η and κ , then the Fisher's information of the original parameters are

[00490] $I(\beta) = I(\eta) (1/(\beta(1-\beta)))^2$

[00491] $I(c) = I(\kappa) (1/c)^2$.

[00492] Hence the implementations have the following transformation on top of the numerical optimization method for estimating stand deviations,

[00493] $\text{std}(\beta) = \text{std}(\eta) \cdot \beta \cdot (1-\beta)$

[00494] $\text{std}(\beta) = \text{std}(\eta) \cdot c.$

Samples

[00495] Samples used herein contain nucleic acids that are “cell-free” (e.g., cfDNA) or cell-bound (e.g., cellular DNA). Cell-free nucleic acids, including cell-free DNA, can be obtained by various methods known in the art from biological samples including but not limited to plasma, serum, and urine (see, e.g., Fan et al., *Proc Natl Acad Sci* 105:16266-16271 [2008]; Koide et al., *Prenatal Diagnosis* 25:604-607 [2005]; Chen et al., *Nature Med.* 2: 1033-1035 [1996]; Lo et al., *Lancet* 350: 485-487 [1997]; Botezatu et al., *Clin Chem.* 46: 1078-1084, 2000; and Su et al., *J Mol. Diagn.* 6: 101-107 [2004]). To separate cell-free DNA from cells in a sample, various methods including, but not limited to fractionation, centrifugation (e.g., density gradient centrifugation), DNA-specific precipitation, or high-throughput cell sorting and/or other separation methods can be used. Commercially available kits for manual and automated separation of cfDNA are available (Roche Diagnostics, Indianapolis, IN, Qiagen, Valencia, CA, Macherey-Nagel, Duren, DE). Biological samples comprising cfDNA have been used in assays to determine the presence or absence of chromosomal abnormalities, e.g., trisomy 21, by sequencing assays that can detect chromosomal aneuploidies and/or various polymorphisms.

[00496] In various embodiments the DNA present in the sample can be enriched specifically or non-specifically prior to use (e.g., prior to preparing a sequencing library). Non-specific enrichment of sample DNA refers to the whole genome amplification of the genomic DNA fragments of the sample that can be used to increase the level of the sample DNA prior to preparing a DNA sequencing library. Non-specific enrichment can be the selective enrichment of one of the two genomes present in a sample that comprises more than one genome. For example, non-specific enrichment can be selective of the cancer genome in a plasma sample, which can be obtained by known methods to increase the relative proportion of cancer to normal DNA in a sample. Alternatively, non-specific enrichment can be the non-selective amplification of both genomes present in the sample. For example, non-specific amplification can be of cancer and normal DNA in a sample comprising a mixture of DNA from the cancer and normal genomes. Methods for whole genome amplification are known in the art. Degenerate oligonucleotide-primed PCR (DOP), primer extension PCR technique (PEP) and multiple

displacement amplification (MDA) are examples of whole genome amplification methods. In some embodiments, the sample comprising the mixture of cfDNA from different genomes is un-enriched for cfDNA of the genomes present in the mixture. In other embodiments, the sample comprising the mixture of cfDNA from different genomes is non-specifically enriched for any one of the genomes present in the sample.

[00497] The sample comprising the nucleic acid(s) to which the methods described herein are applied typically comprises a biological sample (“test sample”), e.g., as described above. In some embodiments, the nucleic acid(s) to be analyzed is purified or isolated by any of a number of well-known methods.

[00498] Accordingly, in certain embodiments the sample comprises or consists of a purified or isolated polynucleotide, or it can comprise samples such as a tissue sample, a biological fluid sample, a cell sample, and the like. Suitable biological fluid samples include, but are not limited to blood, plasma, serum, sweat, tears, sputum, urine, sputum, ear flow, lymph, saliva, cerebrospinal fluid, rashes, bone marrow suspension, vaginal flow, trans-cervical lavage, brain fluid, ascites, milk, secretions of the respiratory, intestinal and genitourinary tracts, amniotic fluid, milk, and leukopheresis samples. In some embodiments, the sample is a sample that is easily obtainable by non-invasive procedures, e.g., blood, plasma, serum, sweat, tears, sputum, urine, sputum, ear flow, saliva or feces. In certain embodiments the sample is a peripheral blood sample, or the plasma and/or serum fractions of a peripheral blood sample. In other embodiments, the biological sample is a swab or smear, a biopsy specimen, or a cell culture. In another embodiment, the sample is a mixture of two or more biological samples, e.g., a biological sample can comprise two or more of a biological fluid sample, a tissue sample, and a cell culture sample. As used herein, the terms “blood,” “plasma” and “serum” expressly encompass fractions or processed portions thereof. Similarly, where a sample is taken from a biopsy, swab, smear, etc., the “sample” expressly encompasses a processed fraction or portion derived from the biopsy, swab, smear, etc.

[00499] In certain embodiments, samples can be obtained from sources, including, but not limited to, samples from different individuals, samples from different developmental stages of the same or different individuals, samples from different diseased individuals (e.g., individuals with cancer or suspected of having a genetic disorder), normal individuals, samples obtained at

different stages of a disease in an individual, samples obtained from an individual subjected to different treatments for a disease, samples from individuals subjected to different environmental factors, samples from individuals with predisposition to a pathology, samples individuals with exposure to an infectious disease agent (e.g., HIV), and the like.

[00500] In one illustrative, but non-limiting embodiment, the sample is a donee sample that is obtained from a donee of an organ transplant, such as a plasma sample from a donee, which includes cfDNA originating from the donee and cfDNA originating from a tissue or organ transplanted from the donor. In this instance, the sample can be analyzed using the methods described herein to quantify donee and donor DNA portions. The donee sample can be a tissue sample, a biological fluid sample, or a cell sample. A biological fluid includes, as non-limiting examples, blood, plasma, serum, sweat, tears, sputum, urine, sputum, ear flow, lymph, saliva, cerebrospinal fluid, ravages, bone marrow suspension, vaginal flow, transcervical lavage, brain fluid, ascites, milk, secretions of the respiratory, intestinal and genitourinary tracts, and leukophoresis samples.

[00501] In another illustrative, but non-limiting embodiment, the donee sample is a mixture of two or more biological samples, e.g., the biological sample can comprise two or more of a biological fluid sample, a tissue sample, and a cell culture sample. In some embodiments, the sample is a sample that is easily obtainable by non-invasive procedures, e.g., blood, plasma, serum, sweat, tears, sputum, urine, milk, sputum, ear flow, saliva and feces. In some embodiments, the biological sample is a peripheral blood sample, and/or the plasma and serum fractions thereof. In other embodiments, the biological sample is a swab or smear, a biopsy specimen, or a sample of a cell culture. As disclosed above, the terms “blood,” “plasma” and “serum” expressly encompass fractions or processed portions thereof. Similarly, where a sample is taken from a biopsy, swab, smear, etc., the “sample” expressly encompasses a processed fraction or portion derived from the biopsy, swab, smear, etc.

[00502] In certain embodiments samples can also be obtained from in vitro cultured tissues, cells, or other polynucleotide-containing sources. The cultured samples can be taken from sources including, but not limited to, cultures (e.g., tissue or cells) maintained in different media and conditions (e.g., pH, pressure, or temperature), cultures (e.g., tissue or cells) maintained for different periods of length, cultures (e.g., tissue or cells) treated with different

factors or reagents (e.g., a drug candidate, or a modulator), or cultures of different types of tissue and/or cells.

[00503] Methods of isolating nucleic acids from biological sources are well known and will differ depending upon the nature of the source. One of skill in the art can readily isolate nucleic acid(s) from a source as needed for the method described herein. In some instances, it can be advantageous to fragment the nucleic acid molecules in the nucleic acid sample. Fragmentation can be random, or it can be specific, as achieved, for example, using restriction endonuclease digestion. Methods for random fragmentation are well known in the art, and include, for example, limited DNase digestion, alkali treatment and physical shearing. In one embodiment, sample nucleic acids are obtained from as cfDNA, which is not subjected to fragmentation.

Sequencing Library Preparation

[00504] In one embodiment, the methods described herein can utilize next generation sequencing technologies (NGS), that allow multiple samples to be sequenced individually as genomic molecules (i.e., singleplex sequencing) or as pooled samples comprising indexed genomic molecules (e.g., multiplex sequencing) on a single sequencing run. These methods can generate up to several hundred million reads of DNA sequences. In various embodiments the sequences of genomic nucleic acids, and/or of indexed genomic nucleic acids can be determined using, for example, the Next Generation Sequencing Technologies (NGS) described herein. In various embodiments analysis of the massive amount of sequence data obtained using NGS can be performed using one or more processors as described herein.

[00505] In various embodiments the use of such sequencing technologies does not involve the preparation of sequencing libraries.

[00506] However, in certain embodiments the sequencing methods contemplated herein involve the preparation of sequencing libraries. In one illustrative approach, sequencing library preparation involves the production of a random collection of adapter-modified DNA fragments (e.g., polynucleotides) that are ready to be sequenced. Sequencing libraries of polynucleotides can be prepared from DNA or RNA, including equivalents, analogs of either DNA or cDNA, for example, DNA or cDNA that is complementary or copy DNA produced from an RNA template,

by the action of reverse transcriptase. The polynucleotides may originate in double-stranded form (e.g., dsDNA such as genomic DNA fragments, cDNA, PCR amplification products, and the like) or, in certain embodiments, the polynucleotides may originate in single-stranded form (e.g., ssDNA, RNA, etc.) and have been converted to dsDNA form. By way of illustration, in certain embodiments, single stranded mRNA molecules may be copied into double-stranded cDNAs suitable for use in preparing a sequencing library. The precise sequence of the primary polynucleotide molecules is generally not material to the method of library preparation, and may be known or unknown. In one embodiment, the polynucleotide molecules are DNA molecules. More particularly, in certain embodiments, the polynucleotide molecules represent the entire genetic complement of an organism or substantially the entire genetic complement of an organism, and are genomic DNA molecules (e.g., cellular DNA, cell free DNA (cfDNA), etc.), that typically include both intron sequence and exon sequence (coding sequence), as well as non-coding regulatory sequences such as promoter and enhancer sequences. In certain embodiments, the primary polynucleotide molecules comprise human genomic DNA molecules, e.g., cfDNA molecules present in peripheral blood of a pregnant subject.

[00507] Preparation of sequencing libraries for some NGS sequencing platforms is facilitated by the use of polynucleotides comprising a specific range of fragment sizes. Preparation of such libraries typically involves the fragmentation of large polynucleotides (e.g. cellular genomic DNA) to obtain polynucleotides in the desired size range.

[00508] Fragmentation can be achieved by any of a number of methods known to those of skill in the art. For example, fragmentation can be achieved by mechanical means including, but not limited to nebulization, sonication and hydroshear. However mechanical fragmentation typically cleaves the DNA backbone at C-O, P-O and C-C bonds resulting in a heterogeneous mix of blunt and 3'- and 5'-overhanging ends with broken C-O, P-O and/ C-C bonds (see, e.g., Alnemri and Liwack, J Biol. Chem 265:17323-17333 [1990]; Richards and Boyer, J Mol Biol 11:327-240 [1965]) which may need to be repaired as they may lack the requisite 5'-phosphate for the subsequent enzymatic reactions, e.g., ligation of sequencing adaptors, that are required for preparing DNA for sequencing.

[00509] In contrast, cfDNA, typically exists as fragments of less than about 300 base pairs and consequently, fragmentation is not typically necessary for generating a sequencing library using cfDNA samples.

[00510] Typically, whether polynucleotides are forcibly fragmented (e.g., fragmented in vitro), or naturally exist as fragments, they are converted to blunt-ended DNA having 5'-phosphates and 3'-hydroxyl. Standard protocols, e.g., protocols for sequencing using, for example, the Illumina platform as described elsewhere herein, instruct users to end-repair sample DNA, to purify the end-repaired products prior to dA-tailing, and to purify the dA-tailing products prior to the adaptor-ligating steps of the library preparation.

[00511] Various embodiments of methods of sequence library preparation described herein obviate the need to perform one or more of the steps typically mandated by standard protocols to obtain a modified DNA product that can be sequenced by NGS. An abbreviated method (ABB method), a 1-step method, and a 2-step method are examples of methods for preparation of a sequencing library, which can be found in patent application 13/555,037 filed on July 20, 2012, which is incorporated by reference by its entirety.

Sequencing Methods

[00512] As indicated above, the prepared samples (e.g., Sequencing Libraries) are sequenced as part of the procedure for quantifying and deconvolving DNA mixture samples. Any of a number of sequencing technologies can be utilized.

[00513] Some sequencing technologies are available commercially, such as the sequencing-by-hybridization platform from Affymetrix Inc. (Sunnyvale, CA) and the sequencing-by-synthesis platforms from 454 Life Sciences (Bradford, CT), Illumina/Solexa (Hayward, CA) and Helicos Biosciences (Cambridge, MA), and the sequencing-by-ligation platform from Applied Biosystems (Foster City, CA), as described below. In addition to the single molecule sequencing performed using sequencing-by-synthesis of Helicos Biosciences, other single molecule sequencing technologies include, but are not limited to, the SMRT™ technology of Pacific Biosciences, the ION TORRENT™ technology, and nanopore sequencing developed for example, by Oxford Nanopore Technologies.

[00514] While the automated Sanger method is considered as a ‘first generation’ technology, Sanger sequencing including the automated Sanger sequencing, can also be employed in the methods described herein. Additional suitable sequencing methods include, but are not limited to nucleic acid imaging technologies, e.g., atomic force microscopy (AFM) or transmission electron microscopy (TEM). Illustrative sequencing technologies are described in greater detail below.

[00515] In one illustrative, but non-limiting, embodiment, the methods described herein comprise obtaining sequence information for the nucleic acids in a test sample, e.g., cfDNA in a donee sample including donor DNA and donee DNA, cfDNA or cellular DNA in a subject being screened for a cancer, and the like, using Illumina’s sequencing-by-synthesis and reversible terminator-based sequencing chemistry (e.g. as described in Bentley et al., *Nature* 6:53-59 [2009]). Template DNA can be genomic DNA, e.g., cellular DNA or cfDNA. In some embodiments, genomic DNA from isolated cells is used as the template, and it is fragmented into lengths of several hundred base pairs. In other embodiments, cfDNA is used as the template, and fragmentation is not required as cfDNA exists as short fragments. For example fetal cfDNA circulates in the bloodstream as fragments approximately 170 base pairs (bp) in length (Fan et al., *Clin Chem* 56:1279-1286 [2010]), and no fragmentation of the DNA is required prior to sequencing. Circulating tumor DNA also exist in short fragments, with a size distribution peaking at about 150-170bp. Illumina’s sequencing technology relies on the attachment of fragmented genomic DNA to a planar, optically transparent surface on which oligonucleotide anchors are bound. Template DNA is end-repaired to generate 5’-phosphorylated blunt ends, and the polymerase activity of Klenow fragment is used to add a single A base to the 3’ end of the blunt phosphorylated DNA fragments. This addition prepares the DNA fragments for ligation to oligonucleotide adapters, which have an overhang of a single T base at their 3’ end to increase ligation efficiency. The adapter oligonucleotides are complementary to the flow-cell anchor oligos (not to be confused with the anchor/anchored reads in the analysis of repeat expansion). Under limiting-dilution conditions, adapter-modified, single-stranded template DNA is added to the flow cell and immobilized by hybridization to the anchor oligos. Attached DNA fragments are extended and bridge amplified to create an ultra-high density sequencing flow cell with hundreds of millions of clusters, each containing about 1,000 copies of the same template. In one embodiment, the randomly fragmented genomic DNA is amplified using PCR

before it is subjected to cluster amplification. Alternatively, an amplification-free (e.g., PCR free) genomic library preparation is used, and the randomly fragmented genomic DNA is enriched using the cluster amplification alone (Kozarewa et al., Nature Methods 6:291-295 [2009]). The templates are sequenced using a robust four-color DNA sequencing-by-synthesis technology that employs reversible terminators with removable fluorescent dyes. High-sensitivity fluorescence detection is achieved using laser excitation and total internal reflection optics. Short sequence reads of about tens to a few hundred base pairs are aligned against a reference genome and unique mapping of the short sequence reads to the reference genome are identified using specially developed data analysis pipeline software. After completion of the first read, the templates can be regenerated in situ to enable a second read from the opposite end of the fragments. Thus, either single-end or paired end sequencing of the DNA fragments can be used.

[00516] Various embodiments of the disclosure may use sequencing by synthesis that allows paired end sequencing. In some embodiments, the sequencing by synthesis platform by Illumina involves clustering fragments. Clustering is a process in which each fragment molecule is isothermally amplified. In some embodiments, as the example described here, the fragment has two different adaptors attached to the two ends of the fragment, the adaptors allowing the fragment to hybridize with the two different oligos on the surface of a flow cell lane. The fragment further includes or is connected to two index sequences at two ends of the fragment, which index sequences provide labels to identify different samples in multiplex sequencing. In some sequencing platforms, a fragment to be sequenced is also referred to as an insert.

[00517] In some implementation, a flow cell for clustering in the Illumina platform is a glass slide with lanes. Each lane is a glass channel coated with a lawn of two types of oligos. Hybridization is enabled by the first of the two types of oligos on the surface. This oligo is complementary to a first adapter on one end of the fragment. A polymerase creates a complement strand of the hybridized fragment. The double-stranded molecule is denatured, and the original template strand is washed away. The remaining strand, in parallel with many other remaining strands, is clonally amplified through bridge application.

[00518] In bridge amplification, a strand folds over, and a second adapter region on a second end of the strand hybridizes with the second type of oligos on the flow cell surface. A

polymerase generates a complimentary strand, forming a double-stranded bridge molecule. This double-stranded molecule is denatured resulting in two single-stranded molecules tethered to the flow cell through two different oligos. The process is then repeated over and over, and occurs simultaneously for millions of clusters resulting in clonal amplification of all the fragments. After bridge amplification, the reverse strands are cleaved and washed off, leaving only the forward strands. The 3' ends are blocked to prevent unwanted priming.

[00519] After clustering, sequencing starts with extending a first sequencing primer to generate the first read. With each cycle, fluorescently tagged nucleotides compete for addition to the growing chain. Only one is incorporated based on the sequence of the template. After the addition of each nucleotide, the cluster is excited by a light source, and a characteristic fluorescent signal is emitted. The number of cycles determines the length of the read. The emission wavelength and the signal intensity determine the base call. For a given cluster all identical strands are read simultaneously. Hundreds of millions of clusters are sequenced in a massively parallel manner. At the completion of the first read, the read product is washed away.

[00520] In the next step of protocols involving two index primers, an index 1 primer is introduced and hybridized to an index 1 region on the template. Index regions provide identification of fragments, which is useful for de-multiplexing samples in a multiplex sequencing process. The index 1 read is generated similar to the first read. After completion of the index 1 read, the read product is washed away and the 3' end of the strand is de-protected. The template strand then folds over and binds to a second oligo on the flow cell. An index 2 sequence is read in the same manner as index 1. Then an index 2 read product is washed off at the completion of the step.

[00521] After reading two indices, read 2 initiates by using polymerases to extend the second flow cell oligos, forming a double-stranded bridge. This double-stranded DNA is denatured, and the 3' end is blocked. The original forward strand is cleaved off and washed away, leaving the reverse strand. Read 2 begins with the introduction of a read 2 sequencing primer. As with read 1, the sequencing steps are repeated until the desired length is achieved. The read 2 product is washed away. This entire process generates millions of reads, representing all the fragments. Sequences from pooled sample libraries are separated based on the unique indices introduced during sample preparation. For each sample, reads of similar stretches of base

calls are locally clustered. Forward and reversed reads are paired creating contiguous sequences. These contiguous sequences are aligned to the reference genome for variant identification.

[00522] The sequencing by synthesis example described above involves paired end reads, which is used in many of the embodiments of the disclosed methods. Paired end sequencing involves two reads from the two ends of a fragment. When a pair of reads are mapped to a reference sequence, the base-pair distance between the two reads can be determined, which distance can then be used to determine the length of the fragments from which the reads were obtained. In some instances, a fragment straddling two bins would have one of its pair-end read aligned to one bin, and another to an adjacent bin. This gets rarer as the bins get longer or the reads get shorter. Various methods may be used to account for the bin-membership of these fragments. For instance, they can be omitted in determining fragment size frequency of a bin; they can be counted for both of the adjacent bins; they can be assigned to the bin that encompasses the larger number of base pairs of the two bins; or they can be assigned to both bins with a weight related to portion of base pairs in each bin.

[00523] Paired end reads may use insert of different length (i.e., different fragment size to be sequenced). As the default meaning in this disclosure, paired end reads are used to refer to reads obtained from various insert lengths. In some instances, to distinguish short-insert paired end reads from long-inserts paired end reads, the latter is also referred to as mate pair reads. In some embodiments involving mate pair reads, two biotin junction adaptors first are attached to two ends of a relatively long insert (e.g., several kb). The biotin junction adaptors then link the two ends of the insert to form a circularized molecule. A sub-fragment encompassing the biotin junction adaptors can then be obtained by further fragmenting the circularized molecule. The sub-fragment including the two ends of the original fragment in opposite sequence order can then be sequenced by the same procedure as for short-insert paired end sequencing described above. Further details of mate pair sequencing using an Illumina platform is shown in an online publication at the following URL, which is incorporated by reference by its entirety: resl.lilluminal.com/documents/products/technotes/technote_nextera_matepair_data_processing. Additional information about paired end sequencing can be found in US Patent No. 7601499 and US Patent Publication No. 2012/0,053,063, which are incorporated by reference with regard to materials on paired end sequencing methods and apparatuses.

[00524] After sequencing of DNA fragments, sequence reads of predetermined length, e.g., 100 bp, are mapped or aligned to a known reference genome. The mapped or aligned reads and their corresponding locations on the reference sequence are also referred to as tags. In one embodiment, the reference genome sequence is the NCBI36/hg18 sequence, which is available on the world wide web at [genome dot ucsc dot edu/cgi-bin/hgGateway?org=Human&db=hg18&hgside=166260105](http://genome.ucsc.edu/cgi-bin/hgGateway?org=Human&db=hg18&hgside=166260105)). Alternatively, the reference genome sequence is the GRCh37/hg19, which is available on the World Wide Web at [genome dot ucsc dot edu/cgi-bin/hgGateway](http://genome.ucsc.edu/cgi-bin/hgGateway). Other sources of public sequence information include GenBank, dbEST, dbSTS, EMBL (the European Molecular Biology Laboratory), and the DDBJ (the DNA Databank of Japan). A number of computer programs are available for aligning sequences, including but not limited to BLAST (Altschul et al., 1990), BLITZ (MPsrch) (Sturrock & Collins, 1993), FASTA (Person & Lipman, 1988), BOWTIE (Langmead et al., Genome Biology 10:R25.1-R25.10 [2009]), or ELAND (Illumina, Inc., San Diego, CA, USA). In one embodiment, one end of the clonally expanded copies of the plasma cfDNA molecules is sequenced and processed by bioinformatics alignment analysis for the Illumina Genome Analyzer, which uses the Efficient Large-Scale Alignment of Nucleotide Databases (ELAND) software.

[00525] In one illustrative, but non-limiting, embodiment, the methods described herein comprise obtaining sequence information for the nucleic acids in a test sample, e.g., cfDNA in a donee sample including donee and donor DNA, cfDNA or cellular DNA in a subject being screened for a cancer, and the like, using single molecule sequencing technology of the Helicos True Single Molecule Sequencing (tSMS) technology (e.g. as described in Harris T.D. et al., Science 320:106-109 [2008]). In the tSMS technique, a DNA sample is cleaved into strands of approximately 100 to 200 nucleotides, and a polyA sequence is added to the 3' end of each DNA strand. Each strand is labeled by the addition of a fluorescently labeled adenosine nucleotide. The DNA strands are then hybridized to a flow cell, which contains millions of oligo-T capture sites that are immobilized to the flow cell surface. In certain embodiments the templates can be at a density of about 100 million templates/cm². The flow cell is then loaded into an instrument, e.g., HeliScope™ sequencer, and a laser illuminates the surface of the flow cell, revealing the position of each template. A CCD camera can map the position of the templates on the flow cell surface. The template fluorescent label is then cleaved and washed away. The sequencing

reaction begins by introducing a DNA polymerase and a fluorescently labeled nucleotide. The oligo-T nucleic acid serves as a primer. The polymerase incorporates the labeled nucleotides to the primer in a template directed manner. The polymerase and unincorporated nucleotides are removed. The templates that have directed incorporation of the fluorescently labeled nucleotide are discerned by imaging the flow cell surface. After imaging, a cleavage step removes the fluorescent label, and the process is repeated with other fluorescently labeled nucleotides until the desired read length is achieved. Sequence information is collected with each nucleotide addition step. Whole genome sequencing by single molecule sequencing technologies excludes or typically obviates PCR-based amplification in the preparation of the sequencing libraries, and the methods allow for direct measurement of the sample, rather than measurement of copies of that sample.

[00526] In another illustrative, but non-limiting embodiment, the methods described herein comprise obtaining sequence information for the nucleic acids in the test sample, e.g., cfDNA in a donee test sample including donee and donor DNA, cfDNA or cellular DNA in a subject being screened for a cancer, and the like, using the 454 sequencing (Roche) (e.g. as described in Margulies, M. et al. *Nature* 437:376-380 [2005]). 454 sequencing typically involves two steps. In the first step, DNA is sheared into fragments of approximately 300-800 base pairs, and the fragments are blunt-ended. Oligonucleotide adaptors are then ligated to the ends of the fragments. The adaptors serve as primers for amplification and sequencing of the fragments. The fragments can be attached to DNA capture beads, e.g., streptavidin-coated beads using, e.g., Adaptor B, which contains 5'-biotin tag. The fragments attached to the beads are PCR amplified within droplets of an oil-water emulsion. The result is multiple copies of clonally amplified DNA fragments on each bead. In the second step, the beads are captured in wells (e.g., picoliter-sized wells). Pyrosequencing is performed on each DNA fragment in parallel. Addition of one or more nucleotides generates a light signal that is recorded by a CCD camera in a sequencing instrument. The signal strength is proportional to the number of nucleotides incorporated. Pyrosequencing makes use of pyrophosphate (PPi), which is released upon nucleotide addition. PPi is converted to ATP by ATP sulfurylase in the presence of adenosine 5' phosphosulfate. Luciferase uses ATP to convert luciferin to oxyluciferin, and this reaction generates light that is measured and analyzed.

[00527] In another illustrative, but non-limiting, embodiment, the methods described herein comprises obtaining sequence information for the nucleic acids in the test sample, e.g., cfDNA in a donee test sample, cfDNA or cellular DNA in a subject being screened for a cancer, and the like, using the SOLiD™ technology (Applied Biosystems). In SOLiD™ sequencing-by-ligation, genomic DNA is sheared into fragments, and adaptors are attached to the 5' and 3' ends of the fragments to generate a fragment library. Alternatively, internal adaptors can be introduced by ligating adaptors to the 5' and 3' ends of the fragments, circularizing the fragments, digesting the circularized fragment to generate an internal adaptor, and attaching adaptors to the 5' and 3' ends of the resulting fragments to generate a mate-paired library. Next, clonal bead populations are prepared in microreactors containing beads, primers, template, and PCR components. Following PCR, the templates are denatured and beads are enriched to separate the beads with extended templates. Templates on the selected beads are subjected to a 3' modification that permits bonding to a glass slide. The sequence can be determined by sequential hybridization and ligation of partially random oligonucleotides with a central determined base (or pair of bases) that is identified by a specific fluorophore. After a color is recorded, the ligated oligonucleotide is cleaved and removed and the process is then repeated.

[00528] In another illustrative, but non-limiting, embodiment, the methods described herein comprise obtaining sequence information for the nucleic acids in the test sample, e.g., cfDNA in a donee test sample, cfDNA or cellular DNA in a subject being screened for a cancer, and the like, using the single molecule, real-time (SMRT™) sequencing technology of Pacific Biosciences. In SMRT sequencing, the continuous incorporation of dye-labeled nucleotides is imaged during DNA synthesis. Single DNA polymerase molecules are attached to the bottom surface of individual zero-mode wavelength detectors (ZMW detectors) that obtain sequence information while phospholinked nucleotides are being incorporated into the growing primer strand. A ZMW detector comprises a confinement structure that enables observation of incorporation of a single nucleotide by DNA polymerase against a background of fluorescent nucleotides that rapidly diffuse in and out of the ZMW (e.g., in microseconds). It typically takes several milliseconds to incorporate a nucleotide into a growing strand. During this time, the fluorescent label is excited and produces a fluorescent signal, and the fluorescent tag is cleaved off. Measurement of the corresponding fluorescence of the dye indicates which base was incorporated. The process is repeated to provide a sequence.

[00529] In another illustrative, but non-limiting embodiment, the methods described herein comprise obtaining sequence information for the nucleic acids in the test sample, e.g., cfDNA in a maternal or donee test sample, cfDNA or cellular DNA in a subject being screened for a cancer, and the like, using nanopore sequencing (e.g. as described in Soni GV and Meller A. Clin Chem 53: 1996-2001 [2007]). Nanopore sequencing DNA analysis techniques are developed by a number of companies, including, for example, Oxford Nanopore Technologies (Oxford, United Kingdom), Sequenom, NABsys, and the like. Nanopore sequencing is a single-molecule sequencing technology whereby a single molecule of DNA is sequenced directly as it passes through a nanopore. A nanopore is a small hole, typically of the order of 1 nanometer in diameter. Immersion of a nanopore in a conducting fluid and application of a potential (voltage) across it results in a slight electrical current due to conduction of ions through the nanopore. The amount of current that flows is sensitive to the size and shape of the nanopore. As a DNA molecule passes through a nanopore, each nucleotide on the DNA molecule obstructs the nanopore to a different degree, changing the magnitude of the current through the nanopore in different degrees. Thus, this change in the current as the DNA molecule passes through the nanopore provides a read of the DNA sequence.

[00530] In another illustrative, but non-limiting, embodiment, the methods described herein comprises obtaining sequence information for the nucleic acids in the test sample, e.g., cfDNA in a donee test sample, cfDNA or cellular DNA in a subject being screened for a cancer, and the like, using the chemical-sensitive field effect transistor (chemFET) array (e.g., as described in U.S. Patent Application Publication No. 2009/0026082). In one example of this technique, DNA molecules can be placed into reaction chambers, and the template molecules can be hybridized to a sequencing primer bound to a polymerase. Incorporation of one or more triphosphates into a new nucleic acid strand at the 3' end of the sequencing primer can be discerned as a change in current by a chemFET. An array can have multiple chemFET sensors. In another example, single nucleic acids can be attached to beads, and the nucleic acids can be amplified on the bead, and the individual beads can be transferred to individual reaction chambers on a chemFET array, with each chamber having a chemFET sensor, and the nucleic acids can be sequenced.

[00531] In another embodiment, the present method comprises obtaining sequence information for the nucleic acids in the test sample, e.g., cfDNA in a donee test sample, using

transmission electron microscopy (TEM). The method, termed Individual Molecule Placement Rapid Nano Transfer (IMPRNT), comprises utilizing single atom resolution transmission electron microscope imaging of high-molecular weight (150kb or greater) DNA selectively labeled with heavy atom markers and arranging these molecules on ultra-thin films in ultra-dense (3nm strand-to-strand) parallel arrays with consistent base-to-base spacing. The electron microscope is used to image the molecules on the films to determine the position of the heavy atom markers and to extract base sequence information from the DNA. The method is further described in PCT patent publication WO 2009/046445. The method allows for sequencing complete human genomes in less than ten minutes.

[00532] In another embodiment, the DNA sequencing technology is the Ion Torrent single molecule sequencing, which pairs semiconductor technology with a simple sequencing chemistry to directly translate chemically encoded information (A, C, G, T) into digital information (0, 1) on a semiconductor chip. In nature, when a nucleotide is incorporated into a strand of DNA by a polymerase, a hydrogen ion is released as a byproduct. Ion Torrent uses a high-density array of micro-machined wells to perform this biochemical process in a massively parallel way. Each well holds a different DNA molecule. Beneath the wells is an ion-sensitive layer and beneath that an ion sensor. When a nucleotide, for example a C, is added to a DNA template and is then incorporated into a strand of DNA, a hydrogen ion will be released. The charge from that ion will change the pH of the solution, which can be detected by Ion Torrent's ion sensor. The sequencer—essentially the world's smallest solid-state pH meter—calls the base, going directly from chemical information to digital information. The Ion personal Genome Machine (PGM™) sequencer then sequentially floods the chip with one nucleotide after another. If the next nucleotide that floods the chip is not a match. No voltage change will be recorded and no base will be called. If there are two identical bases on the DNA strand, the voltage will be double, and the chip will record two identical bases called. Direct detection allows recordation of nucleotide incorporation in seconds.

[00533] In another embodiment, the present method comprises obtaining sequence information for the nucleic acids in the test sample, e.g., cfDNA in a donee test sample, using sequencing by hybridization. Sequencing-by-hybridization comprises contacting the plurality of polynucleotide sequences with a plurality of polynucleotide probes, wherein each of the plurality of polynucleotide probes can be optionally tethered to a substrate. The substrate might be flat

surface comprising an array of known nucleotide sequences. The pattern of hybridization to the array can be used to determine the polynucleotide sequences present in the sample. In other embodiments, each probe is tethered to a bead, e.g., a magnetic bead or the like. Hybridization to the beads can be determined and used to identify the plurality of polynucleotide sequences within the sample.

[00534] In some embodiments of the methods described herein, the mapped sequence tags comprise sequence reads of about 20bp, about 25bp, about 30bp, about 35bp, about 40bp, about 45bp, about 50bp, about 55bp, about 60bp, about 65bp, about 70bp, about 75bp, about 80bp, about 85bp, about 90bp, about 95bp, about 100bp, about 110bp, about 120bp, about 130, about 140bp, about 150bp, about 200bp, about 250bp, about 300bp, about 350bp, about 400bp, about 450bp, or about 500bp. It is expected that technological advances will enable single-end reads of greater than 500bp enabling for reads of greater than about 1000bp when paired end reads are generated. In one embodiment, the mapped sequence tags comprise sequence reads that are 36bp. Mapping of the sequence tags is achieved by comparing the sequence of the tag with the sequence of the reference to determine the chromosomal origin of the sequenced nucleic acid (e.g. cfDNA) molecule, and specific genetic sequence information is not needed. A small degree of mismatch (0-2 mismatches per sequence tag) may be allowed to account for minor polymorphisms that may exist between the reference genome and the genomes in the mixed sample.

[00535] A plurality of sequence tags are typically obtained per sample. In some embodiments, at least about 1×10^5 sequence tags comprising between 75bp read are obtained from mapping the reads to the reference genome per sample.

[00536] The accuracy required for correctly quantifying DNA mixture samples, is predicated on the variation of the number of sequence tags that map to the reference genome among samples within a sequencing run (inter-run variability), and the variation of the number of sequence tags that map to the reference genome in different sequencing runs (inter-run variability). Other variations can result from using different protocols for the extraction and purification of the nucleic acids, the preparation of the sequencing libraries, and the use of different sequencing platforms.

Apparatus and System for Deconvolving and Quantifying Mixtures of Nucleic Acid from Multiple Sources

[00537] Analysis of the sequencing data and the diagnosis derived therefrom are typically performed using various computer programs. Therefore, certain embodiments employ processes involving data stored in or transferred through one or more computer systems or other processing systems. Embodiments disclosed herein also relate to apparatus for performing these operations. This apparatus may be specially constructed for the required purposes, or it may be a general-purpose computer (or a group of computers) selectively activated or reconfigured by a computer program and/or data structure stored in the computer. In some embodiments, a group of processors performs some or all of the recited analytical operations collaboratively (e.g., via a network or cloud computing) and/or in parallel. A processor or group of processors for performing the methods described herein may be of various types including microcontrollers and microprocessors such as programmable devices (e.g., CPLDs and FPGAs) and non-programmable devices such as gate array ASICs or general purpose microprocessors.

[00538] In addition, certain embodiments relate to tangible and/or non-transitory computer readable media or computer program products that include program instructions and/or data (including data structures) for performing various computer-implemented operations. Examples of computer-readable media include, but are not limited to, semiconductor memory devices, magnetic media such as disk drives, magnetic tape, optical media such as CDs, magneto-optical media, and hardware devices that are specially configured to store and perform program instructions, such as read-only memory devices (ROM) and random access memory (RAM). The computer readable media may be directly controlled by an end user or the media may be indirectly controlled by the end user. Examples of directly controlled media include the media located at a user facility and/or media that are not shared with other entities. Examples of indirectly controlled media include media that is indirectly accessible to the user via an external network and/or via a service providing shared resources such as the “cloud.” Examples of program instructions include both machine code, such as produced by a compiler, and files containing higher level code that may be executed by the computer using an interpreter.

[00539] In various embodiments, the data or information employed in the disclosed methods and apparatus is provided in an electronic format. Such data or information may

include reads and tags derived from a nucleic acid sample, counts or densities of such tags that align with particular regions of a reference sequence (e.g., that align to a chromosome or chromosome segment), reference sequences (including reference sequences providing solely or primarily polymorphisms), calls such as SNV or aneuploidy calls, counseling recommendations, diagnoses, and the like. As used herein, data or other information provided in electronic format is available for storage on a machine and transmission between machines. Conventionally, data in electronic format is provided digitally and may be stored as bits and/or bytes in various data structures, lists, databases, etc. The data may be embodied electronically, optically, etc.

[00540] One embodiment provides a computer program product for generating an output indicating the presence or absence of an SNV or aneuploidy associated with a cancer, in a test sample. The computer product may contain instructions for performing any one or more of the above-described methods for determining a chromosomal anomaly. As explained, the computer product may include a non-transitory and/or tangible computer readable medium having a computer executable or compilable logic (e.g., instructions) recorded thereon for enabling a processor to quantify DNA mixture samples. In one example, the computer product comprises a computer readable medium having a computer executable or compilable logic (e.g., instructions) recorded thereon for enabling a processor to quantify DNA mixture samples.

[00541] The sequence information from the sample under consideration may be mapped to chromosome reference sequences to identify a number of sequence tags for each of any one or more chromosomes of interest. In various embodiments, the reference sequences are stored in a database such as a relational or object database, for example.

[00542] It should be understood that it is not practical, or even possible in most cases, for an unaided human being to perform the computational operations of the methods disclosed herein. For example, mapping a single 30 bp read from a sample to any one of the human chromosomes might require years of effort without the assistance of a computational apparatus.

[00543] The methods disclosed herein can be performed using a system for quantifying DNA mixture samples. The system comprising: (a) a sequencer for receiving nucleic acids from the test sample providing nucleic acid sequence information from the sample; (b) a processor; and (c) one or more computer-readable storage media having stored thereon instructions for execution on said processor to carry out a method for quantifying DNA mixture samples.

[00544] In some embodiments, the methods are instructed by a computer-readable medium having stored thereon computer-readable instructions for carrying out a method for quantifying DNA mixture samples. Thus one embodiment provides a computer program product comprising one or more computer-readable non-transitory storage media having stored thereon computer-executable instructions that, when executed by one or more processors of a computer system, cause the computer system to implement a method for quantifying DNA mixture samples. The method includes: (a) extracting nucleic acid molecules from the nucleic acid sample; (b) amplifying the extracted nucleic acid molecules; (c) sequencing the amplified nucleic acid molecules using a nucleic acid sequencer to produce nucleic acid sequence reads; (d) mapping, by the one or more processors, the nucleic acid sequence reads to one or more polymorphism loci on a reference sequence; (e) determining, using the mapped nucleic acid sequence reads and by the one or more processors, allele counts of nucleic acid sequence reads for one or more alleles at the one or more polymorphism loci; and (f) quantifying, using a probabilistic mixture model and by the one or more processors, one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample, wherein using the probabilistic mixture model comprises applying a probabilistic mixture model to the allele counts of nucleic acid sequence reads, and wherein the probabilistic mixture model uses probability distributions to model the allele counts of nucleic acid sequence reads at the one or more polymorphism loci, the probability distributions accounting for errors in the nucleic acid sequence read sequences and counts.

[00545] In some embodiments, the instructions may further include automatically recording information pertinent to the method in a patient medical record for a human subject providing the donee test sample. The patient medical record may be maintained by, for example, a laboratory, physician's office, a hospital, a health maintenance organization, an insurance company, or a personal medical record website. Further, based on the results of the processor-implemented analysis, the method may further involve prescribing, initiating, and/or altering treatment of a human subject from whom the donee test sample was taken. This may involve performing one or more additional tests or analyses on additional samples taken from the subject.

[00546] Disclosed methods can also be performed using a computer processing system which is adapted or configured to perform a method for quantifying DNA mixture samples. One embodiment provides a computer processing system, which is adapted or configured to perform

a method as described herein. In one embodiment, the apparatus comprises a sequencing device adapted or configured for sequencing at least a portion of the nucleic acid molecules in a sample to obtain the type of sequence information described elsewhere herein. The apparatus may also include components for processing the sample. Such components are described elsewhere herein.

[00547] Sequence or other data, can be input into a computer or stored on a computer readable medium either directly or indirectly. In one embodiment, a computer system is directly coupled to a sequencing device that reads and/or analyzes sequences of nucleic acids from samples. Sequences or other information from such tools are provided via interface in the computer system. Alternatively, the sequences processed by system are provided from a sequence storage source such as a database or other repository. Once available to the processing apparatus, a memory device or mass storage device buffers or stores, at least temporarily, sequences of the nucleic acids. In addition, the memory device may store tag counts for various chromosomes or genomes, etc. The memory may also store various routines and/or programs for analyzing the presenting the sequence or mapped data. Such programs/routines may include programs for performing statistical analyses, etc.

[00548] In one example, a user provides a sample into a sequencing apparatus. Data is collected and/or analyzed by the sequencing apparatus, which is connected to a computer. Software on the computer allows for data collection and/or analysis. Data can be stored, displayed (via a monitor or other similar device), and/or sent to another location. The computer may be connected to the internet which is used to transmit data to a handheld device utilized by a remote user (e.g., a physician, scientist or analyst). It is understood that the data can be stored and/or analyzed prior to transmittal. In some embodiments, raw data is collected and sent to a remote user or apparatus that will analyze and/or store the data. Transmittal can occur via the internet, but can also occur via satellite or other connection. Alternately, data can be stored on a computer-readable medium and the medium can be shipped to an end user (e.g., via mail). The remote user can be in the same or a different geographical location including, but not limited to a building, city, state, country or continent.

[00549] In some embodiments, the methods also include collecting data regarding a plurality of polynucleotide sequences (e.g., reads, tags and/or reference chromosome sequences) and sending the data to a computer or other computational system. For example, the computer

can be connected to laboratory equipment, e.g., a sample collection apparatus, a nucleotide amplification apparatus, a nucleotide sequencing apparatus, or a hybridization apparatus. The computer can then collect applicable data gathered by the laboratory device. The data can be stored on a computer at any step, e.g., while collected in real time, prior to the sending, during or in conjunction with the sending, or following the sending. The data can be stored on a computer-readable medium that can be extracted from the computer. The data collected or stored can be transmitted from the computer to a remote location, e.g., via a local network or a wide area network such as the internet. At the remote location various operations can be performed on the transmitted data as described below.

[00550] Among the types of electronically formatted data that may be stored, transmitted, analyzed, and/or manipulated in systems, apparatus, and methods disclosed herein are the following:

Reads obtained by sequencing nucleic acids in a test sample

Tags obtained by aligning reads to a reference genome or other reference sequence or sequences

The reference genome or sequence

Allele counts - Counts or numbers of tags for each allele and regions of a reference genome or other reference sequences

Determined contributor nucleic acid fractions and the associated confidence intervals

Diagnoses (clinical condition associated with the calls)

Recommendations for further tests derived from the calls and/or diagnoses

Treatment and/or monitoring plans derived from the calls and/or diagnoses

[00551] These various types of data may be obtained, stored transmitted, analyzed, and/or manipulated at one or more locations using distinct apparatus. The processing options span a wide spectrum. At one end of the spectrum, all or much of this information is stored and used at the location where the test sample is processed, e.g., a doctor's office or other clinical setting. In other extreme, the sample is obtained at one location, it is processed and optionally sequenced at a different location, reads are aligned and calls are made at one or more different locations, and

diagnoses, recommendations, and/or plans are prepared at still another location (which may be a location where the sample was obtained).

[00552] In various embodiments, the reads are generated with the sequencing apparatus and then transmitted to a remote site where they are processed to produce calls. At this remote location, as an example, the reads are aligned to a reference sequence to produce tags, which are counted and assigned to chromosomes or segments of interest. Also at the remote location, the doses are used to generate calls.

[00553] Among the processing operations that may be employed at distinct locations are the following:

- Sample collection

- Sample processing preliminary to sequencing

- Sequencing

- Analyzing sequence data and quantifying DNA mixture samples

- Diagnosis

- Reporting a diagnosis and/or a call to patient or health care provider

- Developing a plan for further treatment, testing, and/or monitoring

- Executing the plan

- Counseling

[00554] Any one or more of these operations may be automated as described elsewhere herein. Typically, the sequencing and the analyzing of sequence data and quantifying DNA mixture samples will be performed computationally. The other operations may be performed manually or automatically.

[00555] Examples of locations where sample collection may be performed include health practitioners' offices, clinics, patients' homes (where a sample collection tool or kit is provided), and mobile health care vehicles. Examples of locations where sample processing prior to sequencing may be performed include health practitioners' offices, clinics, patients' homes (where a sample processing apparatus or kit is provided), mobile health care vehicles, and

facilities of DNA analysis providers. Examples of locations where sequencing may be performed include health practitioners' offices, clinics, health practitioners' offices, clinics, patients' homes (where a sample sequencing apparatus and/or kit is provided), mobile health care vehicles, and facilities of DNA analysis providers. The location where the sequencing takes place may be provided with a dedicated network connection for transmitting sequence data (typically reads) in an electronic format. Such connection may be wired or wireless and have and may be configured to send the data to a site where the data can be processed and/or aggregated prior to transmission to a processing site. Data aggregators can be maintained by health organizations such as Health Maintenance Organizations (HMOs).

[00556] The analyzing and/or deriving operations may be performed at any of the foregoing locations or alternatively at a further remote site dedicated to computation and/or the service of analyzing nucleic acid sequence data. Such locations include for example, clusters such as general purpose server farms, the facilities of a DNA analysis service business, and the like. In some embodiments, the computational apparatus employed to perform the analysis is leased or rented. The computational resources may be part of an internet accessible collection of processors such as processing resources colloquially known as the cloud. In some cases, the computations are performed by a parallel or massively parallel group of processors that are affiliated or unaffiliated with one another. The processing may be accomplished using distributed processing such as cluster computing, grid computing, and the like. In such embodiments, a cluster or grid of computational resources collective form a super virtual computer composed of multiple processors or computers acting together to perform the analysis and/or derivation described herein. These technologies as well as more conventional supercomputers may be employed to process sequence data as described herein. Each is a form of parallel computing that relies on processors or computers. In the case of grid computing these processors (often whole computers) are connected by a network (private, public, or the Internet) by a conventional network protocol such as Ethernet. By contrast, a supercomputer has many processors connected by a local high-speed computer bus.

[00557] In certain embodiments, the diagnosis is generated at the same location as the analyzing operation. In other embodiments, it is performed at a different location. In some examples, reporting the diagnosis is performed at the location where the sample was taken, although this need not be the case. Examples of locations where the diagnosis can be generated

or reported and/or where developing a plan is performed include health practitioners' offices, clinics, internet sites accessible by computers, and handheld devices such as cell phones, tablets, smart phones, etc. having a wired or wireless connection to a network. Examples of locations where counseling is performed include health practitioners' offices, clinics, internet sites accessible by computers, handheld devices, etc.

[00558] In some embodiments, the sample collection, sample processing, and sequencing operations are performed at a first location and the analyzing and deriving operation is performed at a second location. However, in some cases, the sample collection is collected at one location (e.g., a health practitioner's office or clinic) and the sample processing and sequencing is performed at a different location that is optionally the same location where the analyzing and deriving take place.

[00559] In various embodiments, a sequence of the above-listed operations may be triggered by a user or entity initiating sample collection, sample processing and/or sequencing. After one or more these operations have begun execution the other operations may naturally follow. For example, the sequencing operation may cause reads to be automatically collected and sent to a processing apparatus which then conducts, often automatically and possibly without further user intervention, the sequence analysis and quantifying DNA mixture samples. In some implementations, the result of this processing operation is then automatically delivered, possibly with reformatting as a diagnosis, to a system component or entity that processes reports the information to a health professional and/or patient. As explained such information can also be automatically processed to produce a treatment, testing, and/or monitoring plan, possibly along with counseling information. Thus, initiating an early stage operation can trigger an end to end sequence in which the health professional, patient or other concerned party is provided with a diagnosis, a plan, counseling and/or other information useful for acting on a physical condition. This is accomplished even though parts of the overall system are physically separated and possibly remote from the location of, e.g., the sample and sequence apparatus.

[00560] Figure 4 illustrates, in simple block format, a typical computer system that, when appropriately configured or designed, can serve as a computational apparatus according to certain embodiments. The computer system 2000 includes any number of processors 2002 (also referred to as central processing units, or CPUs) that are coupled to storage devices including

primary storage 2006 (typically a random access memory, or RAM), primary storage 2004 (typically a read only memory, or ROM). CPU 2002 may be of various types including microcontrollers and microprocessors such as programmable devices (e.g., CPLDs and FPGAs) and non-programmable devices such as gate array ASICs or general-purpose microprocessors. In the depicted embodiment, primary storage 2004 acts to transfer data and instructions uni-directionally to the CPU and primary storage 2006 is used typically to transfer data and instructions in a bi-directional manner. Both of these primary storage devices may include any suitable computer-readable media such as those described above. A mass storage device 2008 is also coupled bi-directionally to primary storage 2006 and provides additional data storage capacity and may include any of the computer-readable media described above. Mass storage device 2008 may be used to store programs, data and the like and is typically a secondary storage medium such as a hard disk. Frequently, such programs, data and the like are temporarily copied to primary memory 2006 for execution on CPU 2002. It will be appreciated that the information retained within the mass storage device 2008, may, in appropriate cases, be incorporated in standard fashion as part of primary storage 2004. A specific mass storage device such as a CD-ROM 2014 may also pass data uni-directionally to the CPU or primary storage.

[00561] CPU 2002 is also coupled to an interface 2010 that connects to one or more input/output devices such as such as a nucleic acid sequencer (2020), video monitors, track balls, mice, keyboards, microphones, touch-sensitive displays, transducer card readers, magnetic or paper tape readers, tablets, styluses, voice or handwriting recognition peripherals, USB ports, or other well-known input devices such as, of course, other computers. Finally, CPU 2002 optionally may be coupled to an external device such as a database or a computer or telecommunications network using an external connection as shown generally at 2012. With such a connection, it is contemplated that the CPU might receive information from the network, or might output information to the network in the course of performing the method steps described herein. In some implementations, a nucleic acid sequencer (2020) may be communicatively linked to the CPU 2002 via the network connection 2012 instead of or in addition to via the interface 2010.

[00562] In one embodiment, a system such as computer system 2000 is used as a data import, data correlation, and querying system capable of performing some or all of the tasks described herein. Information and programs, including data files can be provided via a network

connection 2012 for access or downloading by a researcher. Alternatively, such information, programs and files can be provided to the researcher on a storage device.

[00563] In a specific embodiment, the computer system 2000 is directly coupled to a data acquisition system such as a microarray, high-throughput screening system, or a nucleic acid sequencer (2020) that captures data from samples. Data from such systems are provided via interface 2010 for analysis by system 2000. Alternatively, the data processed by system 2000 are provided from a data storage source such as a database or other repository of relevant data. Once in apparatus 2000, a memory device such as primary storage 2006 or mass storage 2008 buffers or stores, at least temporarily, relevant data. The memory may also store various routines and/or programs for importing, analyzing and presenting the data, including sequence reads, UMIs, codes for determining sequence reads, collapsing sequence reads and correcting errors in reads, etc.

[00564] In certain embodiments, the computers used herein may include a user terminal, which may be any type of computer (e.g., desktop, laptop, tablet, etc.), media computing platforms (e.g., cable, satellite set top boxes, digital video recorders, etc.), handheld computing devices (e.g., PDAs, e-mail clients, etc.), cell phones or any other type of computing or communication platforms.

[00565] In certain embodiments, the computers used herein may also include a server system in communication with a user terminal, which server system may include a server device or decentralized server devices, and may include mainframe computers, mini computers, super computers, personal computers, or combinations thereof. A plurality of server systems may also be used without departing from the scope of the present invention. User terminals and a server system may communicate with each other through a network. The network may comprise, e.g., wired networks such as LANs (local area networks), WANs (wide area networks), MANs (metropolitan area networks), ISDNs (Integrated Service Digital Networks), etc. as well as wireless networks such as wireless LANs, CDMA, Bluetooth, and satellite communication networks, etc. without limiting the scope of the present invention.

[00566] Figure 5 shows one implementation of a dispersed system for producing a call or diagnosis from a test sample. A sample collection location 01 is used for obtaining a test sample from a patient such as a pregnant female or a putative cancer patient. The samples then provided

to a processing and sequencing location 03 where the test sample may be processed and sequenced as described above. Location 03 includes apparatus for processing the sample as well as apparatus for sequencing the processed sample. The result of the sequencing, as described elsewhere herein, is a collection of reads which are typically provided in an electronic format and provided to a network such as the Internet, which is indicated by reference number 05 in Figure 5.

[00567] The sequence data is provided to a remote location 07 where analysis and call generation are performed. This location may include one or more powerful computational devices such as computers or processors. After the computational resources at location 07 have completed their analysis and generated a call from the sequence information received, the call is relayed back to the network 05. In some implementations, not only is a call generated at location 07 but an associated diagnosis is also generated. The call and or diagnosis are then transmitted across the network and back to the sample collection location 01 as illustrated in Figure 5. As explained, this is simply one of many variations on how the various operations associated with generating a call or diagnosis may be divided among various locations. One common variant involves providing sample collection and processing and sequencing in a single location. Another variation involves providing processing and sequencing at the same location as analysis and call generation.

[00568] Figure 6 elaborates on the options for performing various operations at distinct locations. In the most granular sense depicted in Figure 6, each of the following operations is performed at a separate location: sample collection, sample processing, sequencing, read alignment, calling, diagnosis, and reporting and/or plan development.

[00569] In one embodiment that aggregates some of these operations, sample processing and sequencing are performed in one location and read alignment, calling, and diagnosis are performed at a separate location. See the portion of Figure 6 identified by reference character A. In another implementation, which is identified by character B in Figure 6, sample collection, sample processing, and sequencing are all performed at the same location. In this implementation, read alignment and calling are performed in a second location. Finally, diagnosis and reporting and/or plan development are performed in a third location. In the implementation depicted by character C in Figure 6, sample collection is performed at a first

location, sample processing, sequencing, read alignment, calling,, and diagnosis are all performed together at a second location, and reporting and/or plan development are performed at a third location. Finally, in the implementation labeled D in Figure 6, sample collection is performed at a first location, sample processing, sequencing, read alignment, and calling are all performed at a second location, and diagnosis and reporting and/or plan management are performed at a third location.

[00570] One embodiment provides a system for analyzing cell-free DNA (cfDNA) for simple nucleotide variants associated with tumors, the system including a sequencer for receiving a nucleic acid sample and providing nucleic acid sequence information from the nucleic acid sample; a processor; and a machine readable storage medium comprising instructions for execution on said processor, the instructions comprising: code for mapping the nucleic acid sequence reads to one or more polymorphism loci on a reference sequence; code for determining, using the mapped nucleic acid sequence reads, allele counts of nucleic acid sequence reads for one or more alleles at the one or more polymorphism loci; and code for quantifying, using a probabilistic mixture model, one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample, wherein using the probabilistic mixture model comprises applying a probabilistic mixture model to the allele counts of nucleic acid sequence reads, and the probabilistic mixture model uses probability distributions to model the allele counts of nucleic acid sequence reads at the one or more polymorphism loci, the probability distributions accounting for errors in the nucleic acid sequence reads.

[00571] In some embodiments of any of the systems provided herein, the sequencer is configured to perform next generation sequencing (NGS). In some embodiments, the sequencer is configured to perform massively parallel sequencing using sequencing-by-synthesis with reversible dye terminators. In other embodiments, the sequencer is configured to perform sequencing-by-ligation. In yet other embodiments, the sequencer is configured to perform single molecule sequencing.

EXPERIMENTAL

Example 1

[00572] This example uses data obtained from actual DNA mixture samples to illustrate that some implementations can provide higher accuracy and reliability, as well as lower empirical bias, in quantifying DNA mixture samples, than conventional technologies that do not use the probabilistic approaches disclosed herein.

[00573] The DNA mixture samples included two DNA from genomes (contributors), and the minor fractions are 0.1%, 0.2%, 0.4%, and 2% in different samples. Some samples included 3 ng of input DNA, and others included 10 ng. The samples were processed in two experimental procedures labeled as Nack or Nack2 to indicate two primer designs, where the numbers of target loci are different for the two designs. Some samples were processed using the MiSeq sequencing platform and some using the MiniSeq platform.

[00574] The sample data were analyzed using three different methods. Table 8 shows the average of coefficient of variance (CV, defined as $\text{standard_deviation_of_predictions} / \text{true_fraction}$) values over multiple mixture fractions and the average of coefficient of variation + bias (CVB, commonly denoted as CV(RMSD) and defined as $\text{RMSD}/\text{true_fraction}$) values over multiple mixture fractions for the three different methods using various samples and experimental procedures. The first method applies a probabilistic model including a binomial distribution for modeling sequencing errors. The first method corresponds to some implementations described as the Seq Model above. The data for the first method (Seq) are shown in the third row of Table 8. The second method applies a probabilistic mixture model including probability distributions accounting for DNA extraction errors, PCR amplification errors, and sequencing errors. The second method corresponds to some implementations described as the Extraction-PCR-Seq Model above. The data for the second method (EPS) are shown in the fourth row of Table 8.

[00575] The third method corresponds to the baseline method NaiveLM or also called KGT.NaiveLM as described above. It determines DNA fractions of contributors using a basic linear regression formulation. The data for the third method (NaiveLM) are shown in the fifth row of Table 8.

[00576] It is worth noting that the genotype information of the contributors was not used to quantify the contributor fractions in the Seq or EPS method, but it was used in the NaiveLM method. Despite the fact that the Seq method and the EPS method did not need to use the genotype information of the contributors, they produced more reliable results as indicated by the smaller coefficient of variation values than the NaiveLM method. Moreover, the Seq method and the EPS method had lower bias as indicated by the smaller CVB values than the NaiveLM method. The best results among the three methods are bolded in Table 8. In short, the two methods using probabilistic mixture models produced more reliable, accurate, and less biased results than the linear regression method.

Table 8. CV and CVB performance metrics for two of disclosed methods (Seq and EPS) compared to baseline method (NaiveLM) on four different datasets.

Method	Genotype	Nack (3ng, 10ng)		Nack2 (3ng, 10ng)		Nack2 MiniSeq (3ng, 10ng)		Nack2 Validation (3ng, 10ng)	
		CV	CVB	CV	CVB	CV	CVB	CV	CVB
Seq	Not used	0.151	0.796	0.214	0.688	0.175	0.414	0.21	0.488
EPS	Not used	0.139	0.207	0.126	0.193	0.125	0.216	0.165	0.253
NaiveLM	Used	0.771	1.117	0.889	1.388	3.818	1.03	8.83	1.407

Example 2

[00577] There are multiple free parameters, such as the average length of DNA template, average length of amplicon, human genome molecular weight, that are used together with the input DNA amount to estimate the effective input DNA amount and read counts. Justified adjustment of these parameters may ensure less biases and robust prediction performance. This example investigates how the average length of DNA template affects the performance of the various methods described above for quantifying DNA mixtures.

[00578] This example uses mock cfDNA (mcfDNA) to mimic real cfDNA. In order to obtain a proper correcting factor for real cfDNA, we need to 1) generate the similar standard mixtures using real cfDNA extracted from two individuals; 2) perform gDNA spike in experiments over real cfDNA mixtures.

Source Genomes

[00579] mcfDNA: mcfDNA from one of the tested cell lines, for which Nack4 target sites do not have CNV for the cell line.

[00580] cfDNA: cfDNA from a healthy person but not maternal cfDNA

[00581] gDNA: gDNA from one of the tested cell line or a normal cell line

Mixture Composition Design

[00582] Mixture 1: 75% cfDNA or mcfDNA, 25% gDNA

[00583] Mixture 2: 50% cfDNA or mcfDNA, 50% gDNA

[00584] Mixture 3: 25% cfDNA or mcfDNA, 75% gDNA

[00585] Mixture 4: 10% cfDNA or mcfDNA, 90% gDNA

[00586] Each with 3 replicates.

Mixing Strategy

[00587] 1. cfDNA and gDNA templates are quantified;

[00588] 2. cfDNA and gDNA templates are mixed at 3:1, 1:1, 1:3, 1:9 ratios;

[00589] 3. PCR over the mixed templates.

[00590] The resulting mixtures and their compositions are shown in Table 9.

Table 9. Mock and real DNA mixtures and estimated mixture fractions reflecting the impact of cfDNA/mcfDNA length on their relative effectiveness as PCR templates. There are three replicates for each type of mixture.

Mixing Fraction	mcfDNA + gDNA	cfDNA + gDNA
Mixture 1 (75%)	82.9%, 84.0% , 84.9%	68.6%, 68.8%, 69.1%
Mixture 3 (25%)	14.2%, 14.5%, 14.4%	20.5%, 20.2%, 20.0%
Mixture 4 (10%)	5.37%, 5.51%, 5.40%	8.94%, 8.86%, 9.09%

[00591] Figure 7 shows the CVB performance of various methods each under different choices of cfDNA length parameter. The following lengths: 120bp, 130bp, 140bp, 150bp, 160bp,

216bp, 300bp, 409bp, and 100k bp are evaluated. Different shades of bars indicate different mcfDNA lengths.

[00592] The different methods are labeled as follows.

[00593] S: probabilistic model accounting for errors due to sequencing. Not using baseline genomes as input. (without knowing D and R genome)

[00594] EPS: probabilistic model accounting for errors due to DNA extraction, PCR, and sequencing. Not using baseline genomes as input.

[00595] PUGT.EPS00: generic implementation of EPS model allowing both known, unknown, and partially known baselines. Not using baseline genomes as input.

[00596] PUGT.EPS: generic implementation of EPS model allowing both known, unknown, and partially known baselines. Using baseline genomes as input.

[00597] KGT.IterLM: Iterative Linear Model. Using baseline genomes as input.

[00598] KGT.Seq: probabilistic model accounting for errors due to sequencing. Using baseline genomes as input.

[00599] KGT.NaiveLM: Baseline method, the naïve linear model with known genotype. Using baseline genomes as input.

[00600] At default DNA length parameter of 160bp, the EPS models have the best performance (indicated by arrows), both when the baseline genomes are available and not available.

[00601] Moreover, the quantification performances of the EPS methods remain outstanding even when practitioners perturb the DNA length parameter from 160bp to 120bp or 216bp. This indicates robustness of the methods to the cfDNA length parameter. The range is comfortably wider than the parameters used in the implementations described above: 160bp for mcfDNA, and 165bp for cfDNA.

[00602] The performance ranking among different methods is:

[00603] PUGT.EPS (using baseline genomes) > KGT.seq or KGT.IterLM (using baseline genomes) > PUGT.EPS or EPS (not using baseline genomes) > S (not using baseline genomes) > KGT.Naive (using baseline genomes).

[00604] Notably, the three EPS methods have markedly lower CVB than the naïve linear model with known genotype, indicating that the EPS methods have improved accuracy and

reduced bias over conventional linear model methods. Note that conventional methods are only applicable to mixture samples with known baseline genomes.

[00605] Furthermore, under default DNA length parameters, methods described in this disclosure have lower limit of blank (LOB) and higher analytical sensitivity than the method using conventional linear model. As shown in Table 10, the limit of blank (LOB) is below 0.1% for four methods disclosed, but the conventional, naïve linear model method's LOB is at 0.42%.

Table 10. LOB of Different Methods

Method	S	PUGT.EPS00	PUGT.EPS	KGT.seq	KGT.NaiveLM
LOB	0.05%	0.08%	0.06%	0.03%	0.42%

Example 3

[00606] This example uses data obtained from mock cfDNA (mcfDNA) and actual genomic DNA (gDNA) to investigate the sensitivities of some of the disclosed methods, and compare them to a known method KIMERDx that uses a qPCR technique.

[00607] Table 11 shows the LOQ of two probabilistic models labeled as follows.

[00608] EPS: probabilistic model accounting for errors due to DNA extraction, PCR, and sequencing. Not using baseline genomes as input.

[00609] PUGT.EPS: generic implementation of EPS model allowing both known, unknown, and partially known baselines. Using baseline genomes as input.

[00610] LOQ, or limit of quantification, is a measure of quantification sensitivity. It is defined as the minimal donor fraction that can be determined at no greater than 20% coefficient of variation (CV).

[00611] Under mcfDNA conditions (top two rows of data in Table 11), which mimic cfDNA samples from solid organ transplant patients, DNA mixture samples of two contributors were generated. Each sample included 3 ng of DNA. The probabilistic methods PUGT.EPS (using baseline genotypes) and EPS (without using baseline genotypes from pre-transplant recipient and donor) were applied to 5 samples x 3 replicates. Both probabilistic methods

achieved LOQ of $\leq 0.2\%$ when using only 3 ng of input DNA, indicating high sensitivity for both disclosed methods.

[00612] Under a gDNA condition (third row of data in Table 11), which mimics blood gDNA samples from bone marrow transplant patients, DNA mixture samples of two contributors were generated. Each sample includes 10 ng of DNA. The PUGT.EPS method was used to analyze 5 samples x 3 replicates. The PUGT.EPS method achieved an LOQ of $\leq 0.1\%$ when using 10 ng of input DNA, which, as expected, is lower than the LOQ in the mcfDNA conditions using 3ng of input DNA.

[00613] Under another gDNA condition (four row of data in Table 11), DNA mixture samples of five contributors were generated. Each sample includes 10 ng total amount of DNA. The PUGT.EPS method was used to analyze 5 samples x 3 replicates. The PUGT.EPS method achieved an LOQ of $\leq 0.35\%$. Even for such a difficult condition with five contributors, the method achieved a great LOQ significantly lower than 1%.

Table 11. Sensitivity of Disclosed Methods

Sample Type	LOQ	Sample Size	Method
mcfDNA 3ng, 2 contributors	0.2%	5 samples x 3	PUGT.EPS
mcfDNA 3ng, 2 contributors	0.2%	5 samples x 3	EPS (no baseline)
gDNA 10ng, 2 contributors	0.1%	5 samples x 3	PUGT.EPS
gDNA 10ng, 5 contributors	0.35%	4 samples x 3	PUGT.EPS

[00614] Table 12 shows the sensitivity (LOQ) values of a KIMERDx method that uses a qPCR technique on mixture samples of only two contributors. The KIMERDx method was used to analyze different quantity of input gDNA. To achieve 0.1% of LOQ, it requires 66 ng of input gDNA. In comparison, the PUGT.EPS method only requires ≤ 10 ng of input DNA to achieve the same level of sensitivity. With 10ng input gDNA, KIMERDx would achieve an LOQ of 0.7% compared to $< 0.1\%$ for PUGT.EPS.

Table 12. Sensitivity of qPCR KIMERDx Method

LOQ	# Cells	Input DNA (ng)
0.05%	20000	132
0.1%	10000	66
1%	1000	7
2%	500	3
5%	200	1

[00615] Therefore, this example illustrates that the disclosed probabilistic methods required significantly less input DNA to achieve a same level of sensitivity compared to the state of art method. Conversely, the disclose methods achieves a significantly higher sensitivity at low input DNA amount. Due to their improved sensitivity, the methods may allow for faster sample processing, require less reagent and improve accuracy of DNA mixture quantification.

[00616] Existing chimerism assays do *not* work for solid organ transplant monitoring, which our methods are designed for. The disclosed methods improve the sensitivity of DNA mixture quantification, which is particularly beneficial in applications where the input DNA quantity is limited, which covers all solid organ transplant cases. Solid organ transplant monitoring using cfDNA is challenging because the amount of cfDNA extracted from a typical blood sample is typically < 10ng, much lower than the amount of extractable gDNA. Meanwhile, cfDNA is much less effective as PCR template compared to gDNA of the same amount.

[00617] Existing methods also do not work for transplant with more than one donor, for which the methods we disclosed still achieved high sensitivity. Transplants with more than one donor occur frequently for bone marrow transplants, and are also commonly seen in organ transplant with blood transfusion and in patients with prior organ transplants .

Example 4

[00618] Conventional methods of chimerism analysis utilize capillary electrophoresis (CE) fragment analysis or quantitative polymerase chain reaction (qPCR) analysis of short tandem repeats (STRs) or small insertions and deletions (Indels). There are a number of drawbacks associated with these methods including limit of quantification, dynamic range, number of targets, workflow, analysis, and reproducibility. An alternative approach to these conventional methods utilizes next-generation sequencing (NGS) targeting hundreds of SNPs to quantitatively assess chimerism with low limit of quantification, broad dynamic range, simple workflow, automated analysis, and robust reproducibility.

Conventional Chimerism Analysis Using CE**Targets: STRs**

[00619] STRs are loci found throughout the genome. They are comprised of short sequences, usually between 2 and 8 nucleotides and most commonly 4, that are repeated tandemly (e.g. *gata* tandemly repeated as *gatagatagatagatagata*). The number of repeats varies between 4 and 40 repeats making a typical STR less than 400 total nucleotides in length. The number of repeats is highly variable within the human population. These two characteristics of STRs, relatively short total length and high variability, have made them attractive targets for human identification in forensic science. The short length is important for poor quality forensic samples because amplification of larger regions is difficult with these types of samples. The high variability in the population is an attractive feature because a relatively small number are needed for positive identification. While more than 100 STRs have been well characterized in the human genome, most applications use less than 30.

Assay Design

[00620] PCR primers are designed in the conserved flanking regions surrounding the STR. Primers can be multiplexed with each of the four fluorophores containing 4 to 7 STRs of varying lengths. This means that the multiplexes support between 10 and 21 unique STRs. The CE system measures the relative fluorescence units and the elapse time to detection to generate an electropherogram for each STR. Most labs utilize the full multiplexes for generating pre-transplant baseline genotypes for the recipient and the donor. The pre-transplant genotypes are

compared to one another to select informative markers, markers in which the recipient and donor have unique alleles. The chimerism samples may be run with the entire multiplex or with individual singleplex assays for the informative STRs. Singleplex assays generally provide the highest level of sensitivity, but many labs prefer to run the multiplex assay.

Workflow

- DNA is extracted from peripheral blood, bone marrow, or cell lineages isolated with magnetic beads or by flow cytometry.
- PCR amplification of the target STRs is performed including fluorescent tagging.
- Separation and detection of the STR-PCR amplicons is performed with electrophoresis, most frequently a CE instrument. The CE system measures the relative fluorescence units and the elapse time to detection to generate electropherograms for each allele present in the sample.
- The person performing the analysis reviews the electropherograms for each informative marker to determine the relative frequency of the donor to the recipient. In cases with multiple informative markers, the average frequency is usually taken as the final measure of chimerism after taking into account variable performance of the different markers.

[00621] From extracted DNA to data analysis takes about seven hours with about 2 hours of that hands-on time. The analysis of the data is highly variable and takes from 15 minutes to two hours to analyze a single chimerism sample depending on the number of informative markers, the variability between the markers, and the complexity of the stutter peak subtraction.

Limitations

[00622] There are three primary limitations to CE analysis of STR regions for chimerism analysis.

[00623] First, the electropherogram peaks alone are often difficult to analyze and percent chimerism from multiple peaks within the same sample frequently vary 10-15%. As a consequence of this variability, analysis can often take hours for a single sample and the results are still semiquantitative.

[00624] Second, limit of quantification (LOQ), often referred to as limit of detection (LOD) or sensitivity, ranges from 1-5% with this methodology. This broad range exists because each STR will have its own LOQ depending on the PCR enzyme stutter or “slippage” on the STR and variable performance of the fluorophores.

[00625] Third, although more than 100 STR targets are well characterized in the genome, including more than 21 STRs in an assay has not been reliable. This is because multiplexing that many specific primer pools into a single assay is very difficult to make robust and reliable. Therefore, chimerism mixtures from closely related individuals may have difficulty identifying informative markers and cases with many donors may be very difficult to analyze.

[00626] These limitations can be significant in clinical use. For example, an actual chimerism result of 99% will be reported as 100%.

Conventional Chimerism Analysis Using qPCR

Targets: Indels

[00627] An indel is an insertion or deletion of 1 to 10,000 nucleotide bases. Millions of indels have been discovered in the human genome making it the second largest contributor to human genome variability after SNPs. Similar to STRs, many indels are short and can be easily amplified even from highly degraded DNA and small amounts of DNA. In addition, there is a wide variety of indels available in different lengths, different allele frequencies, and they are broadly distributed throughout the genome. These features of indels make them attractive targets for human identification and chimerism analysis.

Assay Design

[00628] PCR primers are designed to amplify the indel and are designed as singleplex, small multiplexes (~3 targets), or large multiplexes (30-40 targets). It has been shown that 30-40 appropriately selected indels are needed to distinguish individuals from one another. With the commercially available kits, pre-transplant donor and recipient baseline samples are run through 30 to 40 indel targets in either 3-indel multiplexes or individual indels laid out on a 96-well plate. This step identifies informative targets in which the donor and recipient have different alleles. A

minimum of two informative targets are then selected for each donor-recipient pair to be used for chimerism analysis.

[00629] Each indel is targeted by a set of fluorescently labeled primers that hybridize the DNA of interest. As the amplicon undergoes PCR cycling, the increasing fluorescence is proportional to the quantity of amplicon present. The quantification is determined by the number of PCR cycles required to reach the threshold cycle (Ct) value. The informative markers are usually selected to amplify the genome of the minor contributor, usually the recipient in the case of stem cell transplantation. The quantity is then determined by comparing the Ct values of the post-transplant sample, the matched pre-transplant baseline, and the reference control sample.

Workflow

- DNA is extracted from peripheral blood, bone marrow, or cell lineages isolated with magnetic beads or by flow cytometry.
- Purified DNA is quantified and diluted as needed to achieve target concentrations.
- Baseline genotyping is performed by testing both the donor and recipient pre-transplant samples for every target indel in the system. In the small multiplex system this includes 10 individual reactions of 2-3 indel targets per reaction. In the singleplex system, this requires 46 individual reactions with a single indel target in each reaction. Each baseline sample run must also include a positive control and a no template control. This means that the small multiplex system can fit 8 baseline samples on a 96-well plate and the singleplex system can fit 2 per plate.
- 10ng of baseline DNA is added to each reaction well (100ng total for small multiplex and 460ng for singleplex)
- PCR Master Mix is prepared and added to each reaction well.
- Amplification primers are added to the appropriate wells (8x10 for the small multiplex and 2x46 for the singleplex)
- Plates are sealed, vortexed, centrifuged, and loaded onto the qPCR instrument.
- Results are loaded into the application-specific software.
- Recipient and donor baselines are compared in the software and informative markers are selected for chimerism analysis. Usually two informative targets are selected for each transplant recipient/donor pair.

- For each target to be amplified, the pre-transplant baseline sample from the minor contributor must be run in triplicate, each post-transplant chimerism sample must be run in triplicate, a positive control for every two reaction wells, and a no template control for each target. In other words, to perform a single post-transplant chimerism analysis 60ng (6 wells) of reference DNA must be run, 60ng (6 wells) of pre-transplant baseline DNA must be run, and 60ng (6 wells) of post-transplant chimerism DNA must be run. This is a total of 21 wells to generate data from 2 targets.
- PCR Master Mix is prepared and added to each reaction well.
- Amplification primers are added to the appropriate wells (7 wells per sample – 3 pre-transplant, 3 post-transplant, and 1 no template control)
- Plates are sealed, vortexed, centrifuged, and loaded onto the qPCR instrument.
- Results are loaded into the application-specific software.

[00630] From extracted DNA to genotyping data for informative marker selection takes about 3 total hour with one and half hours hands-on. After selection of informative markers and DNA extraction from chimerism samples, an addition 3 hours and one and a half hours of hands-on time is needed for generation of the chimerism data.

Limitations

[00631] There are three primary limitations of qPCR-based chimerism analysis of indel targets.

[00632] First, each chimerism analysis requires 60ng of pre-transplant recipient baseline sample. This is in addition to the 100-500ng of baseline DNA required for the initial genotyping. For programs frequently performing chimerism analysis, the pre-transplant baseline samples may be depleted, limiting the ability to run this assay for long periods of times.

[00633] Second, the requirement to run the chimerism analyses as singleplex reactions complicates the overall system requiring dozens of unique assays to be held in inventory. In addition, the cost of each reaction usually limits the analysis to only two targets per donor-recipient pair and these targets are likely to be different for each donor-recipient pair, making the setup susceptible to error.

[00634] Third, while the LOQ for qPCR is very low, the dynamic range of qPCR-based chimerism suffers and chimerism predictions when the minor contributor is greater than 30% are not reliable.

Novel Chimerism Analysis by NGS

Targets: SNPs

[00635] SNPs are single nucleotide positions in which variation is present to a measurable degree within the human population or within specific populations. *dbSNP* is a database of SNPs managed by the National Center for Biotechnology Information (NCBI) and it currently contains more than 170 million human SNPs with nearly 25 million of them validated. This means that SNPs are responsible for the vast majority of variability within the human population averaging one SNP per 1,000 nucleotide bases. SNPs can be biallelic (two observed alleles), triallelic (three observed alleles), or tetra-allelic (four observed alleles). A single base variant can be considered a SNP when the minor allele has a frequency of at least 1% in a random set of individuals in a population. SNPs are excellent targets for chimerism analysis because of their low mutation rate, small amplicon size, and compatibility with high-throughput sequencing technology.

Assay Design

[00636] SNPs are selected to be biallelic with roughly 50/50 allele frequency within various populations around the world. In addition, SNPs having low mutation rates and no linkage disequilibrium with the SNP pool are selected. Finally, SNPs were assessed for designability, both in terms of minimizing primer-primer interaction and uniformity in PCR amplification and in sequencing coverage. The total number of SNPs is determined based on power to discriminate between first-degree relatives from all populations around the world.

[00637] A single PCR step amplifies the DNA, isolates the amplicons of interest, and incorporates flowcell adapters (inverse oligonucleotide sequences to those on the Illumina flowcells allowing the sample amplicons to bind to the flowcell), sequencing primers (oligonucleotide sequences that serve as initiation sites for the Illumina sequencing by synthesis (SBS) process), and index barcode sequences (oligonucleotide sequences that allow multiple samples to be run simultaneously).

[00638] The NGS system sequences each amplicon hundreds to thousands of times. In pre-transplant baseline samples, this information is used to genotype each contributor. In post-transplant chimerism samples, the reads counts for each nucleotide at a SNP location can be used with or without the baseline genotypes to accurately estimate the percent chimerism of each contributor, up to five total contributors.

Workflow

- DNA is extracted from peripheral blood, bone marrow, or cell lineages isolated with magnetic beads or by flow cytometry.
- Purified DNA is quantified and diluted as needed to achieve target concentrations.
- Unique index barcodes are added to each sample DNA.
- Master mix is added to every sample, mixed, sealed, and centrifuged.
- PCR amplification is performed.
- All samples are pooled into a single well and then a PCR clean-up is performed.
- The cleaned pool is quantified, diluted, and denatured.
- The final pool, also called a library, is loaded onto the sequencer and sequencing is initiated.
- Sequencing data is imported into the chimerism-specific analysis software for automated quality control and chimerism analysis.

[00639] From extracted DNA to loading of the sequencer takes less than 3 hours with less than 2 hours of hands-on time. The sequencing run requires 9 to 13 hours depending on the number of samples being run simultaneously. Once sequencing data are collected, the analysis of the data does not require manual intervention, allowing automation of the analysis and reducing human errors.

Limitations

[00640] There is one primary limitation of NGS-based chimerism analysis using SNPs: compared to CE and qPCR-based chimerism analysis, NGS-based sample processing and sequencing take longer, although the hands-on time is equivalent. The NGS-based library preparation can be completed in the afternoon with the sequencing completed overnight. This allows 24-hour turnaround for samples received in the morning. However, because sequencing

may be multiplexed, this method can combine multiple samples for sequencing, thereby improving the overall efficiency of sample processing.

Summary

[00641] NGS-based chimerism analysis using SNP targets is an efficient, accurate, and reliable method to overcome many of the limitations associated with conventional methods of chimerism analysis. The results are truly quantitative and can be automatically generated without the need for laborious human review of electropherograms and stutter subtractions. The NGS-based chimerism analysis has a broad dynamic range with low LOQ and no performance degradation at high levels of mixed chimerism. More than 200 SNP targets are used with the NGS system and they are multiplexed into a single reaction. This allows for utility with more than one donor and with donor-recipient pairs that are very closely related. The indexing capabilities and throughput of the NGS system allow for baseline and chimerism samples to be run simultaneously, only one assay and kit to be stored in inventory, and low potential for human error in the workflow.

Example 5

[00642] This example shows that some implementations improve over conventional methods because of the throughput of the NGS sequencer, the assay design with incredibly high uniformity, and the use of SNPs as targets. The disclosed methods can analyze far more targets than conventional methods, which are limited to <30 targets. The process allows multiplex many samples to boost efficiency. The methods are quantitative, and can all be done cost-effectively.

[00643] One experiment compares the performance of the methods in some implementations with baseline genomes known or unknown. Table 12 shows the DNA quantifications for four samples with different recipient portions for three baseline conditions (both baselines known, both baselines unknown, and recipient known but donor unknown). The results show that the methods can be performed with and without baselines with similar performance at different recipient portions. When baselines are known, the methods tend to produce results with smaller confidence intervals (and higher reliability).

Table 12. DNA Quantifications with Known and Unknown Baselines

	Both Baselines Known			Both Baselines Unknown			Recipient Known Donor Unknown		
	% recipient	Low 95% Ct	High 95% Ct	% recipient	Low 95% Ct	High 95% Ct	% recipient	Low 95% Ct	High 95% Ct
Sample 1	0.7%	0.6%	0.8%	0.7%	0.5%	0.9%	0.7%	0.6%	0.8%
Sample 2	5.8%	5.4%	6.2%	5.7%	5.2%	6.2%	5.8%	5.4%	6.2%
Sample 3	12.3%	12.0%	12.6%	12.3%	11.8%	12.7%	12.3%	12%	12.6%
Sample 4	38.6%	38.1%	39.0%	38.7%	38.0%	39.3%	38.6%	38.1%	39.0%

[00644] Figure 8 compares DNA portion determined by some implementations (Y axis) and actual DNA portions (X axis). The horizontal lines indicate the values of actual portions. The chimerism sample includes cfDNA mixtures that are mock cfDNA provided by Horizon Discovery (Catalog No. 12498714289). As the figure shows, the predicted minor contributor portion are quite close to the actual minor contributor portion at 0.1%, 0.2%, 0.4%, and 2%.

[00645] Figure 9 shows the coefficient of variance (CV) of 16 conditions for determining the limit of quantification (LOQ) for some implementations. LOQ is defined as the lowest concentration at which an analyte can be reliably detected at which the imprecision (CV) is less than 20%. This measurement takes into account both analytical sensitivity (i.e., limit of detection) and reproducibility (i.e., precision). The four different groups of bars represent different minor contributor fractions of 0.1%, 0.2%, 0.4%, and 2%. The four bars in a group represent, from left to right, four input DNA conditions: 10ng of gDNA, 3 ng of gDNA, 10ng of cfDNA, and 3 ng of cfDNA. At each minor contributor fraction, there is a consistent pattern as expected —samples of smaller amounts lead to higher CV, and cfDNA lead to higher CV.

[00646] All but one condition (0.1% minor contributor fraction, 3ng of cfDNA) can detect an analyte with an imprecision (CV) less than 20%. In other words, the one condition (3ng of cfDNA) has an LOQ of 0.2%, while the rest of the conditions have an LOQ of 0.1%.

[00647] Table 13 summarizes the data above. It clearly shows that all four input DNA conditions have LOQ values smaller than 0.2%, and all but the most challenging input condition (3 ng cfDNA) have LOQ of 0.1%.

Table 13. Limit of Quantification for some Implementations with Different DNA Input

Input DNA	Limit of Quantification
10 ng cfDNA	<0.1%
3 ng cfDNA	<0.2%
10 ng gDNA	<0.1%
3 ng gDNA	<0.1%

Discussions

[00648] Conventional chimerism methods using qPCR or CE technologies sacrifice ease of use, number of targets, sensitivity, or dynamic range.

[00649] The implemented methods multiplex samples in a single assay. This enables pre-transplant baseline samples and post-transplant chimerism samples to be run using the same assay and side by side on the same sequencing run. The methods can be performed with and without baselines with near identical performance.

[00650] qPCR and CE chimerism methods may provide some level of multiplexing for the pre-transplant baseline samples, but these methods have performance degradation for post-transplant chimerism quantification when targets are multiplexed. This means that baseline and chimerism samples must be run separately and the entire system may requires at least a dozen unique assays.

[00651] While qPCR is sensitive for microchimerism detection, it lacks the dynamic range to be reliable for mixed chimerism. CE-based chimerism analysis offers a broad dynamic range for mixed chimerism detection, but lacks the sensitivity for microchimerism. The disclosed methods provide both a reliable low limit of quantification (LOQ) and a broad dynamic range, enabling one solution that can cover all different types of chimerism.

CLAIMS

What is claimed is:

1. A method, implemented at a computer system that includes one or more processors and system memory, of quantifying a nucleic acid sample comprising nucleic acid of one or more contributors, the method comprising:
 - (a) extracting nucleic acid molecules from the nucleic acid sample;
 - (b) amplifying the extracted nucleic acid molecules;
 - (c) sequencing the amplified nucleic acid molecules using a nucleic acid sequencer to produce nucleic acid sequence reads;
 - (d) mapping, by the one or more processors, the nucleic acid sequence reads to one or more polymorphism loci on a reference sequence;
 - (e) determining, using the mapped nucleic acid sequence reads and by the one or more processors, allele counts of nucleic acid sequence reads for one or more alleles at the one or more polymorphism loci; and
 - (f) quantifying, using a probabilistic mixture model and by the one or more processors, one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample, wherein using the probabilistic mixture model comprises applying a probabilistic mixture model to the allele counts of nucleic acid sequence reads, and wherein the probabilistic mixture model uses probability distributions to model the allele counts of nucleic acid sequence reads at the one or more polymorphism loci, the probability distributions accounting for errors in the nucleic acid sequence reads.
2. The method of claim 1, further comprising, determining, using the probabilistic mixture model and by the one or more processors, one or more genotypes of the one or more contributors at the one or more polymorphism loci.
3. The method of claim 1, further comprising, determining, using the one or more fractions of nucleic acid of the one or more contributors, a risk of one contributor (a donee) rejecting a tissue or an organ transplanted from another contributor (a donor).

4. The method of claim 1, wherein the one or more contributors comprise two or more contributors.
5. The method of claim 1, wherein the nucleic acid molecules comprise DNA molecules or RNA molecules.
6. The method of claim 1, wherein the nucleic acid sample comprises nucleic acid from zero, one, or more contaminant genomes and one genome of interest.
7. The method of claim 1, wherein the one or more contributors comprise zero, one, or more donors of a transplant and a donee of the transplant, and wherein the nucleic acid sample comprises a sample obtained from the donee.
8. The method of claim 1, wherein the transplant comprises an allogeneic or xenogeneic transplant.
9. The method of claim 1, wherein the nucleic acid sample comprises a biological sample obtained from the donee.
10. The method of claim 1, wherein the nucleic acid sample comprises a biological sample obtained from a cell culture.
11. The method of claim 1, wherein the extracted nucleic acid molecules comprise cell-free nucleic acid.
12. The method of claim 1, wherein the extracted nucleic acid molecules comprise cellular DNA.
13. The method of claim 1, wherein the one or more polymorphism loci comprise one or more biallelic polymorphism loci.

14. The method of claim 1, wherein the one or more alleles at the one or more polymorphism loci comprise one or more single nucleotide polymorphism (SNP) alleles.

15. The method of claim 1, wherein the probabilistic mixture model uses a single-locus likelihood function to model allele counts at a single polymorphism locus, the single-locus likelihood function comprising

$$M(n_{1i}, n_{2i} | p_{1i}, \theta)$$

wherein

n_{1i} is the allele count of allele 1 at locus i ,

n_{2i} is the allele count of allele 2 at locus i ,

p_{1i} is an expected fraction of allele 1 at locus i , and

θ comprises one or more model parameters.

16. The method of claim 15, wherein p_{1i} is modeled as a function of:

(i) genotypes of the contributors at locus i , or $g_i = (g_{11i}, \dots, g_{Di})$, which is a vector of copy number of allele 1 at locus i in contributors 1...D;

(ii) read count errors resulting from the sequencing operation in (c), or λ ; and

(iii) fractions of nucleic acid of contributors in the nucleic acid sample, or $\beta = (\beta_1, \dots, \beta_D)$, wherein D is the number of contributors.

17. The method of claim 16, wherein the contributors comprise two or more contributors, and

$$p_{1i} = p(g_i, \lambda, \beta) \leftarrow [(1 - \lambda) g_i + \lambda (2 - g_i)] / 2 \bullet \beta, \text{ where } \bullet \text{ is vector dot product operator}$$

18. The method of claim 17, wherein the contributors comprise two contributors, and p_{1i} is obtained using the p_{1i}' values in Table 3.

19. The method of claim 16, wherein zero, one or more genotypes of the contributors are unknown.

20. The method of claim 19, wherein (f) comprises marginalizing over a plurality of possible combinations of genotypes to enumerate the probability parameter p_{1i} .

21. The method of claim 19, further comprising determining a genotype configuration at each of the one or more polymorphism loci, the genotype configuration comprising two alleles for each of the one or more contributors.

22. The method of claim 16, wherein the single-locus likelihood function comprise a first binomial distribution.

23. The method of claim 22, wherein the first binomial distribution is expressed as follows:

$$n_{1i} \sim BN(n_i, p_{1i})$$

wherein

n_{1i} is an allele count of nucleic acid sequence reads for allele 1 at locus i ; and

n_i is a total read count at locus i , which equals to a total genome copy numbers n .

24. The method of claim 23, wherein (f) comprises maximizing a multiple-loci likelihood function calculated from a plurality of single-locus likelihood functions.

25. The method of claim 24, wherein (f) comprises:

calculating a plurality of multiple-loci likelihood values using a plurality of potential fraction values and a multiple-loci likelihood function of the allele counts of nucleic acid sequence reads determined in (e);

identifying one or more potential fraction values associated with a maximum multiple-loci likelihood value; and

quantifying the one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample as the identified potential fraction value.

26. The method of claim 24, wherein the multiple-loci likelihood function comprises:

$$L(\boldsymbol{\beta}, \theta, \lambda, \pi; n_1, n_2) = \prod_i [\sum g_i M(n_{1i}, n_{2i} | \mathbf{p}(g_i, \lambda, \boldsymbol{\beta}), \theta) \cdot \mathbf{P}(g_i | \pi)]$$

wherein

$L(\boldsymbol{\beta}, \theta, \lambda, \pi; n_1, n_2)$ is the likelihood of observing allele count vectors n_1 and n_2 for alleles 1 and 2;

$p(g_i, \lambda, \beta)$ is the expected fraction or probability of observing allele 1 at locus i based on the contributors' genotypes g_i at locus i ;

$P(g_i|\pi)$ is the prior probability of observing the genotypes g_i at locus i given a population allele frequency (π); and,

$\sum g_i$ denotes summing over a plurality of possible combinations of genotypes of the contributors.

27. The method of claim 26, wherein the multiple-loci likelihood function comprises:

$$L(\beta, \lambda, \pi ; n_1, n_2) = \prod_i [\sum g_i BN(n_{1i} | n_i, p(g_i, \lambda, \beta)) \cdot P(g_i | \pi)]$$

28. The method of claim 27, wherein the contributors comprise two contributors and the likelihood function comprises:

$$L(\beta, \lambda, \pi ; n_1, n_2) = \prod_i \sum_{g_{1i}g_{2i}} BN(n_{1i} | n_i, p_{1i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)$$

wherein

$L(\beta, \lambda, \pi ; n_1, n_2)$ is the likelihood of observing allele count vectors n_1 to n_2 for alleles 1 and 2 given parameters β and π ;

$p_{1i}(g_{1i}, g_{2i}, \lambda, \beta)$ is a probability parameter, taken as p_{1i}' from Table 3, indicating a probability of allele 1 at locus i based on the two contributors' genotypes (g_{1i}, g_{2i}); and

$P(g_{1i}, g_{2i} | \pi)$ is a prior joint probability of observing the two contributors' genotypes given a population allele frequency (π).

29. The method of claim 28, wherein the prior joint probability is calculated using marginal distributions $P(g_{1i} | \pi)$ and $P(g_{2i} | \pi)$ that satisfy the Hardy-Weinberg equilibrium.

30. The method of claim 29, wherein the prior joint probability is calculated using genetic relationship between the two contributors.

31. The method of claim 26, wherein the probabilistic mixture model accounts for nucleic acid molecule copy number errors resulting from extracting the nucleic acid molecules performed in (a), as well as the read count errors resulting from the sequencing operation in (c).

32. The method of claim 31, wherein the probabilistic mixture model uses a second binomial distribution to model allele counts of the extracted nucleic acid molecules for alleles at the one or more polymorphism loci.

33. The method of claim 32, wherein the second binomial distribution is expressed as follows:

$$n_{li}'' \sim BN(n_i'', p_{li})$$

wherein

n_{li}'' is an allele count of extracted nucleic acid molecules for allele l at locus i ;

n_i'' is a total nucleic acid molecule count at locus i ; and

p_{li} is a probability parameter indicating the probability of allele l at locus i .

34. The method of claim 33, wherein the first binomial distribution is conditioned on an allele fraction n_{li}''/n_i'' .

35. The method of claim 34, wherein the first binomial distribution is re-parameterized as follows:

$$n_{li} \sim BN(n_i, n_{li}''/n_i'')$$

wherein

n_{li} is an allele count of nucleic acid sequence reads for allele l at locus i ;

n_i'' is a total number of nucleic acid molecules at locus i , which equals to a total genome copy numbers n'' ;

n_i is a total read count at locus i ; and

n_{li}'' is a number of extracted nucleic acid molecules for allele l at locus i .

36. The method of claim 35, wherein the probabilistic mixture model uses a first beta distribution to approximate a distribution of n_{li}''/n'' .

37. The method of claim 36, wherein the first beta distribution has a mean and a variance that match a mean and a variance of the second binomial distribution.

38. The method of claim 36, wherein locus i is modeled as biallelic and the first beta distribution is expressed as follows:

$$n_{i1}/n'' \sim \text{Beta}((n''-1)p_{1i}, (n''-1)p_{2i})$$

wherein

p_{1i} is a probability parameter indicating the probability of a first allele at locus i ; and

p_{2i} is a probability parameter indicating the probability of a second allele at locus i .

39. The method of claim 36, wherein (f) comprises combining the first binomial distribution, modeling sequencing read counts, and the first beta distribution, modeling extracted nucleic acid molecule number, to obtain the single-locus likelihood function of n_{1i} that follows a first beta-binomial distribution.

40. The method of claim 39, wherein the first beta-binomial distribution has the form:

$$n_{1i} \sim \text{BB}(n_i, (n''-1) \cdot p_{1i}, (n''-1) \cdot p_{2i}),$$

or an alternative approximation:

$$n_{1i} \sim \text{BB}(n_i, n'' \cdot p_{1i}, n'' \cdot p_{2i}).$$

41. The method of claim 40, wherein the multiple-loci likelihood function comprises:

$$L(\beta, n'', \lambda, \pi; n_1, n_2) = \prod_i [\sum g_i \text{BB}(n_{1i} | n_i, (n''-1) \cdot p_{1i}, (n''-1) \cdot p_{2i}) \cdot P(g_i | \pi)]$$

wherein $L(\beta, n'', \lambda, \pi; n_1, n_2)$ is the likelihood of observing allele count vectors n_1 and n_2 for alleles 1 and 2 at all loci, and $p_{1i} = p(g_i, \lambda, \beta)$, $p_{2i} = 1 - p_{1i}$.

42. The method of claim 41, wherein the contributors comprise two contributors, and the multiple-loci likelihood function comprises:

$$L(\beta, n'', \lambda, \pi; n_1, n_2) =$$

$$\prod_i \sum_{g_{1i} g_{2i}} \text{BB}(n_{1i}, n_{2i} | n_i, (n''-1) \cdot p_{1i}(g_{1i}, g_{2i}, \lambda, \beta), (n''-1) \cdot p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)$$

wherein $L(\beta, n'', \lambda, \pi; n_1, n_2)$ is the likelihood of observing an allele count vector for the first allele of all loci (n_1) and an allele count vector for the second allele of all loci (n_2) given parameters β , n'' , λ , and π ;

$p_{1i}(g_{1i}, g_{2i}, \lambda, \beta)$ is a probability parameter, taken as p_{1i}' from Table 3, indicating a probability of allele 1 at locus i based on the two contributors' genotypes (g_{1i}, g_{2i});

$p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)$ is a probability parameter, taken as p_2' from Table 3, indicating a probability of allele 2 at locus i based on the two contributors' genotypes (g_{1i}, g_{2i}); and

$P(g_{1i}, g_{2i} | \pi)$ is a prior joint probability of observing the first contributor's genotype for the first allele (g_{1i}) and the second contributor's genotype for the first allele (g_{2i}) at locus i given a population allele frequency (π).

43. The method of claim 35, wherein (f) comprises estimating the total extracted genome copy number n'' from a mass of the extracted nucleic acid molecules.

44. The method of claim 43, wherein the estimated total extracted genome copy number n'' is adjusted according to fragment size of the extracted nucleic acid molecules.

45. The method of claim 26, wherein the probabilistic mixture model accounts for nucleic acid molecule number errors resulting from amplifying the nucleic acid molecules performed in (b), as well as the read count errors resulting from the sequencing operation in (c).

46. The method of claim 45, the amplification process of (b) is modeled as follows:

$$x_{t+1} = x_t + y_{t+1}$$

wherein

x_{t+1} is the nucleic acid copies of a given allele after cycle $t+1$ of amplification;

x_t is the nucleic acid copies of a given allele after cycle t of amplification;

y_{t+1} is the new copies generated at cycle $t+1$, and it follows a binomial distribution $y_{t+1} \sim BN(x_t, r_{t+1})$; and

r_{t+1} is the amplification rate for cycle $t+1$.

47. The method of claim 45, wherein the probabilistic mixture model uses a second beta distribution to model allele fractions of the amplified nucleic acid molecules for alleles at the one or more polymorphism loci.

48. The method of claim 47, wherein locus i is biallelic and the second beta distribution is expressed as follows:

$$n_{1i}' / (n_{1i}' + n_{2i}') \sim \text{Beta}(n'' \cdot \rho_i \cdot p_{1i}, n'' \cdot \rho_i \cdot p_{2i})$$

wherein

n_{1i}' is an allele count of amplified nucleic acid molecules for a first allele at locus i ;

n_{2i}' is an allele count of amplified nucleic acid molecules for a second allele at locus i ;

n'' is a total nucleic acid molecule count at any locus;

ρ_i is a constant related to an average amplification rate r ;

p_{1i} is the probability of the first allele at locus i ; and

p_{2i} is the probability of the second allele at locus i .

49. The method of claim 48, wherein ρ_i is $(1+r)/(1-r) / [1-(1+r)^{-t}]$, and r is the average amplification rate per cycle.

50. The method of claim 48, wherein ρ_i is approximated as $(1+r)/(1-r)$.

51. The method of claim 48, wherein (f) comprises combining the first binomial distribution and the second beta distribution to obtain the single-locus likelihood function for n_{1i} that follows a second beta-binomial distribution.

52. The method of claim 51, wherein the second beta-binomial distribution has the form:

$$n_{1i} \sim \text{BB}(n_i, n'' \cdot \rho_i \cdot p_{1i}, n'' \cdot \rho_i \cdot p_{2i})$$

wherein

n_{1i} is an allele count of nucleic acid sequence reads for the first allele at locus i ;

p_{1i} is a probability parameter indicating the probability of a first allele at locus i ; and

p_{2i} is a probability parameter indicating the probability of a second allele at locus i .

53. The method of claim 52, wherein (f) comprises, by assuming the one or more polymorphism loci have a same amplification rate, re-parameterizing the second beta-binomial distribution as:

$$n_{1i} \sim \text{BB}(n_i, n'' \cdot (1+r)/(1-r) \cdot p_{1i}, n'' \cdot (1+r)/(1-r) \cdot p_{2i})$$

wherein r is an amplification rate.

54. The method of claim 53, wherein the multiple-loci likelihood function comprises:

$$L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i [\sum g_i BB(n_{1i} | n_i, n'' \cdot (1+r)/(1-r) \cdot p_{1i}, n'' \cdot (1+r)/(1-r) \cdot p_{2i}) \cdot P(g_i | \pi)]$$

55. The method of claim 53, wherein the contributors comprise two contributors and the multiple-loci likelihood function comprises:

$$L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i \sum_{g_{1i} g_{2i}} [BB(n_{1i} | n_i, n'' \cdot (1+r)/(1-r) \cdot p_{1i}(g_{1i}, g_{2i}, \lambda, \beta), n'' \cdot (1+r)/(1-r) \cdot p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)]$$

wherein $L(\beta, n'', r, \lambda, \pi; n_1, n_2)$ is the likelihood of observing an allele count vector for the first allele of all loci (n_1) and an allele count vector for the second allele of all loci (n_2) given parameters β, n'', r, λ , and π .

56. The method of claim 52, wherein (f) comprises, by defining a relative amplification rate of each polymorphism locus to be proportional to a total reads of the locus, re-parameterizing the second beta-binomial distribution as:

$$n_{1i} \sim BB(n_i, c' \cdot n_i \cdot p_{1i}, c' \cdot n_i \cdot p_{2i})$$

wherein

c' is a parameter to be optimized; and

n_i is the total reads at locus i .

57. The method of claim 56, wherein the multiple-loci likelihood function comprises:

$$L(\beta, n'', c', \lambda, \pi; n_1, n_2) = \prod_i [\sum g_i BB(n_{1i} | n_i, c' \cdot n_i \cdot p_{1i}, c' \cdot n_i \cdot p_{2i}) \cdot P(g_i | \pi)]$$

58. The method of claim 26, wherein the probabilistic mixture model accounts for nucleic acid molecule number errors resulting from extracting the nucleic acid molecules performed in (a) and amplifying the nucleic acid molecules performed in (b), as well as the read count errors resulting from the sequencing operation in (c).

59. The method of claim 58, wherein the probabilistic mixture model uses a third beta distribution to model allele fractions of the amplified nucleic acid molecules for alleles at the one or more polymorphism loci, accounting for the sampling errors resulting from extracting the

nucleic acid molecules performed in (a) and amplifying the nucleic acid molecules performed in (b).

60. The method of claim 59, wherein locus i is biallelic and the third beta distribution has the form of:

$$n_{1i}' / (n_{1i}' + n_{2i}') \sim \text{Beta}(n'' \cdot (1 + r_i)/2 \cdot p_{1i}, n'' \cdot (1 + r_i)/2 \cdot p_{2i})$$

wherein

n_{1i}' is an allele count of amplified nucleic acid molecules for a first allele at locus i ;

n_{2i}' is an allele count of amplified nucleic acid molecules for a second allele at locus i ;

n'' is a total nucleic acid molecule count;

r_i is the average amplification rate for locus i ;

p_{1i} is the probability of the first allele at locus i ; and

p_{2i} is a probability of the second allele at locus i .

61. The method of claim 60, wherein (f) comprises combining the first binomial distribution and the third beta distribution to obtain the single-locus likelihood function of n_{1i} that follows a third beta-binomial distribution.

62. The method of claim 61, wherein the third beta-binomial distribution has the form:

$$n_{1i} \sim \text{BB}(n_i, n'' \cdot (1 + r_i)/2 \cdot p_{1i}, n'' \cdot (1 + r_i)/2 \cdot p_{2i})$$

wherein r_i is an amplification rate.

63. The method of claim 62, wherein the multiple-loci likelihood function comprises:

$$L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i [\sum_{g_i} \text{BB}(n_{1i} | n_i, n'' \cdot (1 + r)/2 \cdot p_{1i}, n'' \cdot (1 + r)/2 \cdot p_{2i}) \cdot P(g_i | \pi)],$$

wherein r is an amplification rate assumed to be equal for all loci.

64. The method of claim 62, wherein the contributors comprise two contributors, and wherein the multiple-loci likelihood function comprises:

$$L(\beta, n'', r, \lambda, \pi; n_1, n_2) = \prod_i \sum_{g_{1i}, g_{2i}} \text{BB}(n_{1i} | n_i, n'' \cdot (1 + r)/2 \cdot p_{1i}(g_{1i}, g_{2i}, \lambda, \beta), n'' \cdot (1 + r)/2 \cdot p_{2i}(g_{1i}, g_{2i}, \lambda, \beta)) \cdot P(g_{1i}, g_{2i} | \pi)$$

wherein $L(n_1, n_2 | \beta, n'', r, \lambda, \pi)$ is the likelihood of observing allele counts for the first allele vector n_1 and an allele count for the second allele vector n_2 given parameters β, n'', r, λ , and π .

65. The method of claim 1, further comprising: (g) estimating one or more confidence intervals of the one or more fractions of nucleic acid of the one or more contributors using the hessian matrix of the log-likelihood using numerical differentiation.

66. The method of claim 1, wherein the mapping of (d) comprises identifying, by the one or more processors using computer hashing and computer dynamic programming, reads among the nucleic acid sequence reads matching any sequence of a plurality of unbiased target sequences, wherein the plurality of unbiased target sequences comprises sub-sequences of the reference sequence and sequences that differ from the subsequences by a single nucleotide.

67. The method of claim 66, wherein the plurality of unbiased target sequences comprises five categories of sequences encompassing each polymorphic site of a plurality of polymorphic sites:

(i) a reference target sequence that is a sub-sequence of the reference sequence, the reference target sequence having a reference allele with a reference nucleotide at the polymorphic site;

(ii) alternative target sequences each having an alternative allele with an alternative nucleotide at the polymorphic site, the alternative nucleotide being different from the reference nucleotide;

(iii) mutated reference target sequences comprising all possible sequences that each differ from the reference target sequence by only one nucleotide at a site that is not the polymorphic site;

(iv) mutated alternative target sequences comprising all possible sequences that each differ from an alternative target sequence by only one nucleotide at a site that is not the polymorphic site; and

(v) unexpected allele target sequences each having an unexpected allele different from the reference allele and the alternative allele, and each having a sequence different from the previous four categories of sequences.

68. The method of claim 67, further comprising estimating a sequencing error rate λ at the variant site base on a frequency of observing the unexpected allele target sequences of (v).

69. The method of claim 67, wherein (e) comprises using the identified reads and their matching unbiased target sequences to determine allele counts of the nucleic acid sequence reads for the alleles at the one or more polymorphism loci.

70. The method of claim 67, wherein the plurality of unbiased target sequences comprises sequences that are truncated to have the same length as the nucleic acid sequence reads.

71. The method of claim 67, wherein the plurality of unbiased target sequences comprises sequences stored in one or more hash tables, and the reads are identified using the hash tables.

72. A system quantifying a nucleic acid sample comprising nucleic acid of one or more contributors, the system comprising:

(a) a sequencer configured to (i) receive nucleic acid molecules extracted from the nucleic acid sample, (ii) amplify the extracted nucleic acid molecules, and (iii) sequence the amplified nucleic acid molecules under conditions that produce nucleic acid sequence reads; and

(b) a computer comprising one or more processors configured to:

map the nucleic acid sequence reads to one or more polymorphism loci on a reference sequence;

determine, using the mapped nucleic acid sequence reads, allele counts of nucleic acid sequence reads for one or more alleles at the one or more polymorphism loci; and

quantify, using a probabilistic mixture model, one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample,

wherein

using the probabilistic mixture model comprises applying a probabilistic mixture model to the allele counts of nucleic acid sequence reads, and

the probabilistic mixture model uses probability distributions to model the allele counts of nucleic acid sequence reads at the one or more polymorphism loci, the probability distributions accounting for errors in the nucleic acid sequence reads.

73. The system of claim 72, further comprising a tool for extracting nucleic acid molecules from the nucleic acid sample.

74. The system of claim 72, wherein the probability distributions comprise a first binomial distribution as follows:

$$n_{1i} \sim BN(n_i, p_{1i})$$

wherein

n_{1i} is an allele count of nucleic acid sequence reads for allele 1 at locus i ;

n_i is a total read count at locus i , which equals to a total genome copy numbers n "; and

p_{1i} is a probability parameter indicating the probability of allele 1 at locus i .

75. A computer program product comprising a non-transitory machine readable medium storing program code that, when executed by one or more processors of a computer system, causes the computer system to implement a method of quantifying a nucleic acid sample comprising nucleic acid of one or more contributors, said program code comprising:

code for mapping the nucleic acid sequence reads to one or more polymorphism loci on a reference sequence;

code for determining, using the mapped nucleic acid sequence reads, allele counts of nucleic acid sequence reads for one or more alleles at the one or more polymorphism loci; and

code for quantifying, using a probabilistic mixture model, one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample,

wherein

using the probabilistic mixture model comprises applying a probabilistic mixture model to the allele counts of nucleic acid sequence reads, and

the probabilistic mixture model uses probability distributions to model the allele counts of nucleic acid sequence reads at the one or more polymorphism loci, the probability distributions accounting for errors in the nucleic acid sequence reads.

76. A method, implemented at a computer system that includes one or more processors and system memory, of quantifying a nucleic acid sample comprising nucleic acid of one or more contributors, the method comprising:

(a) receiving, by the one or more processors, nucleic acid sequence reads obtained from the nucleic acid sample;

(b) mapping, by the one or more processors, using computer hashing and computer dynamic programming, the nucleic acid sequence reads to one or more polymorphism loci on a reference sequence;

(c) determining, using the mapped nucleic acid sequence reads and by the one or more processors, allele counts of nucleic acid sequence reads for one or more alleles at the one or more polymorphism loci; and

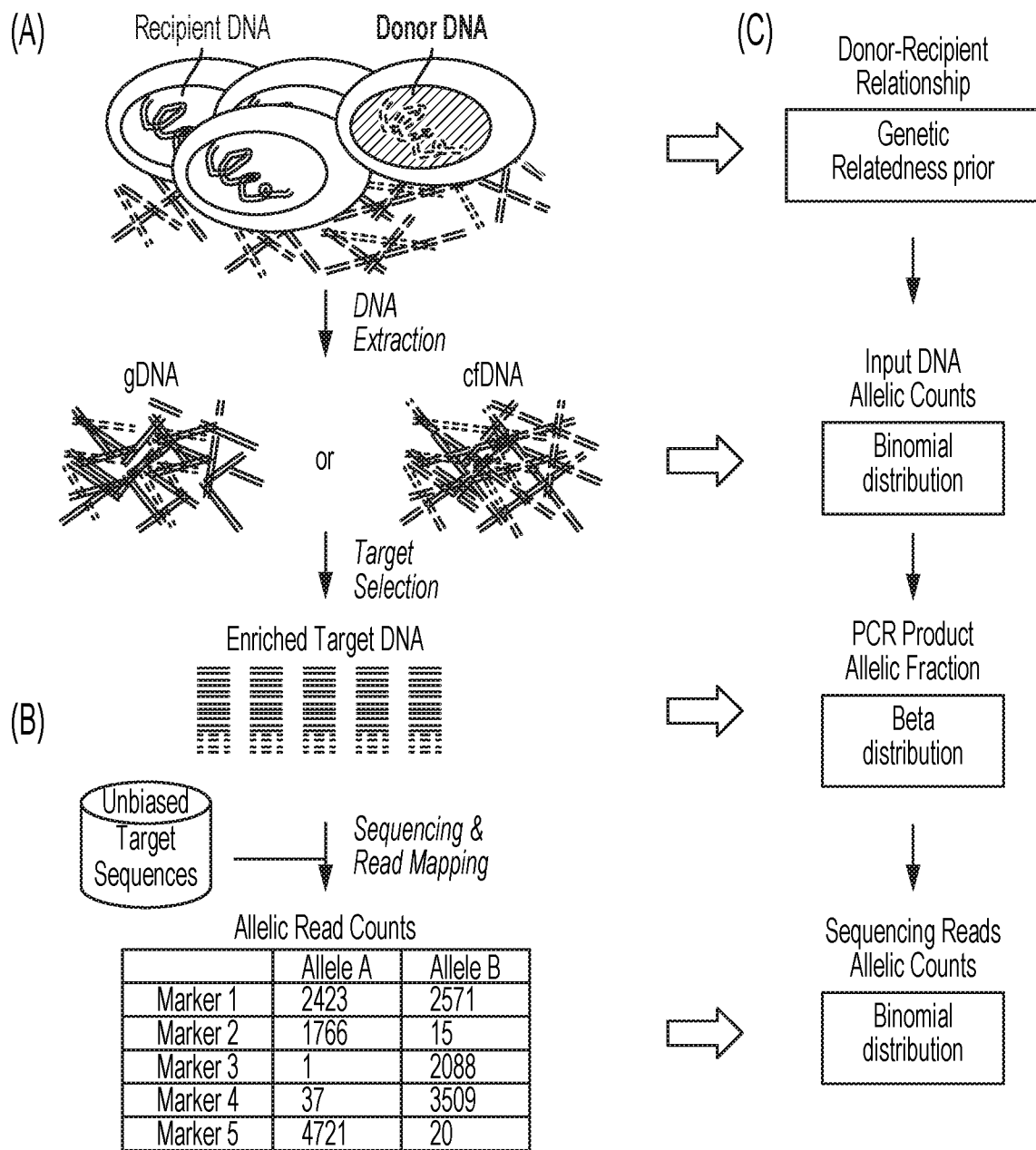
(d) quantifying, using a probabilistic mixture model and by the one or more processors, one or more fractions of nucleic acid of the one or more contributors in the nucleic acid sample and confidence of the fractions,

wherein using the probabilistic mixture model comprises applying a probabilistic mixture model to the allele counts of nucleic acid sequence reads,

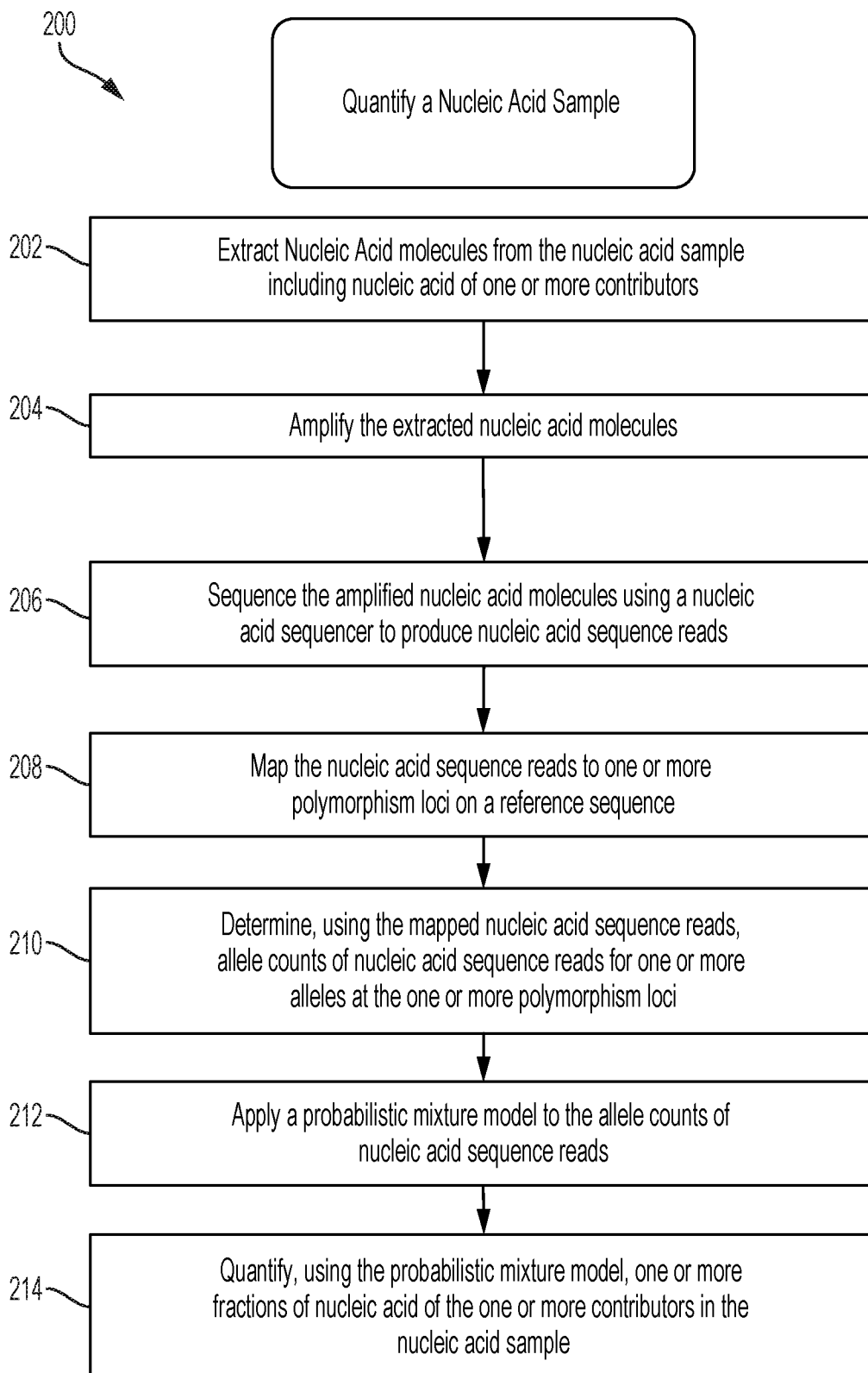
wherein the probabilistic mixture model uses probability distributions to model the allele counts of nucleic acid sequence reads at the one or more polymorphism loci, the probability distributions accounting for errors in the mapped nucleic acid sequence reads,

and wherein the quantifying employs (i) a computer optimization method combining multi-iteration grid searching and a BFGS - quasi-Newton method, or an iterative weighted linear regression, and (ii) a numerical differentiation method.

1/11

**Figure 1**

2/11

**Figure 2A**

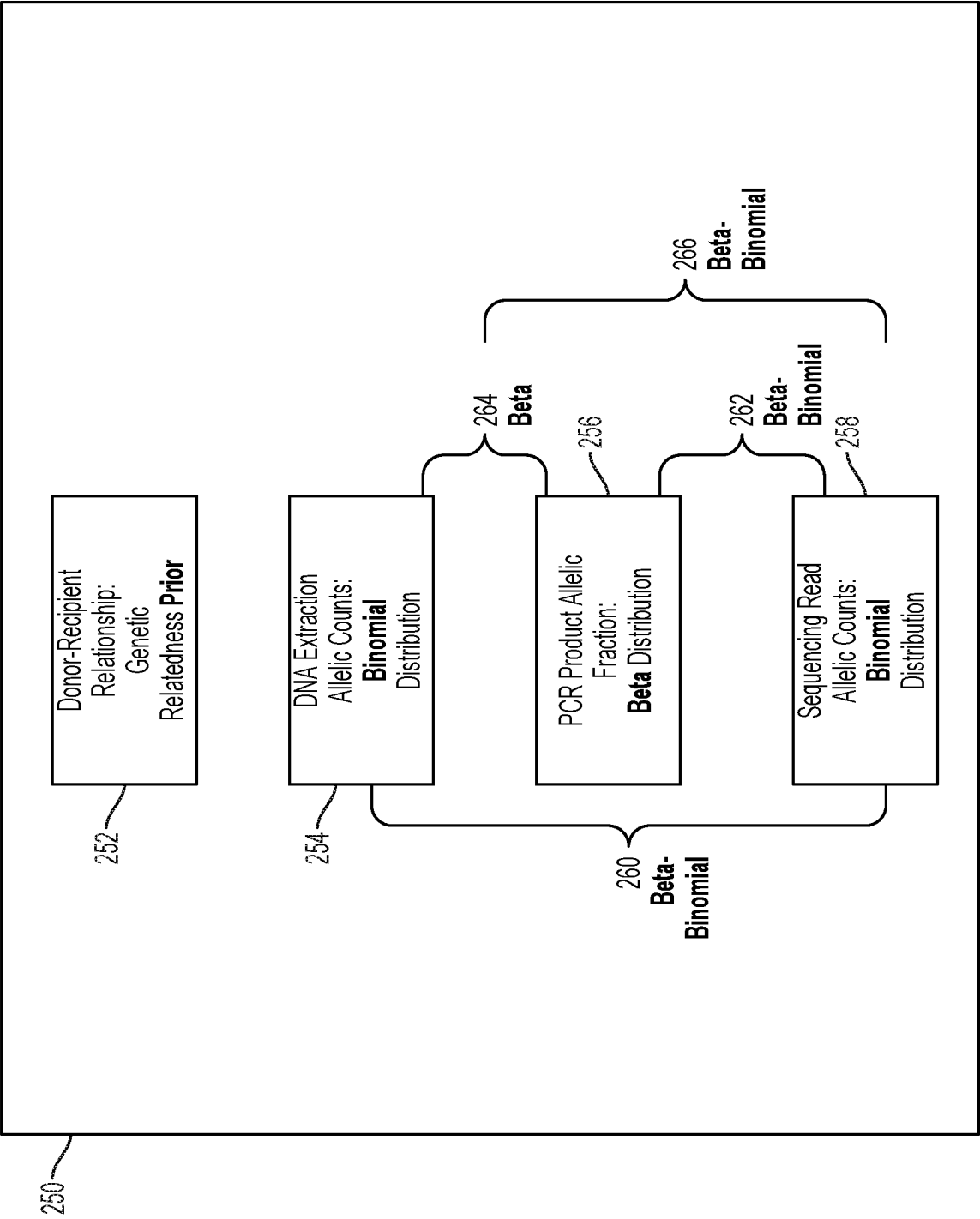


Figure 2B

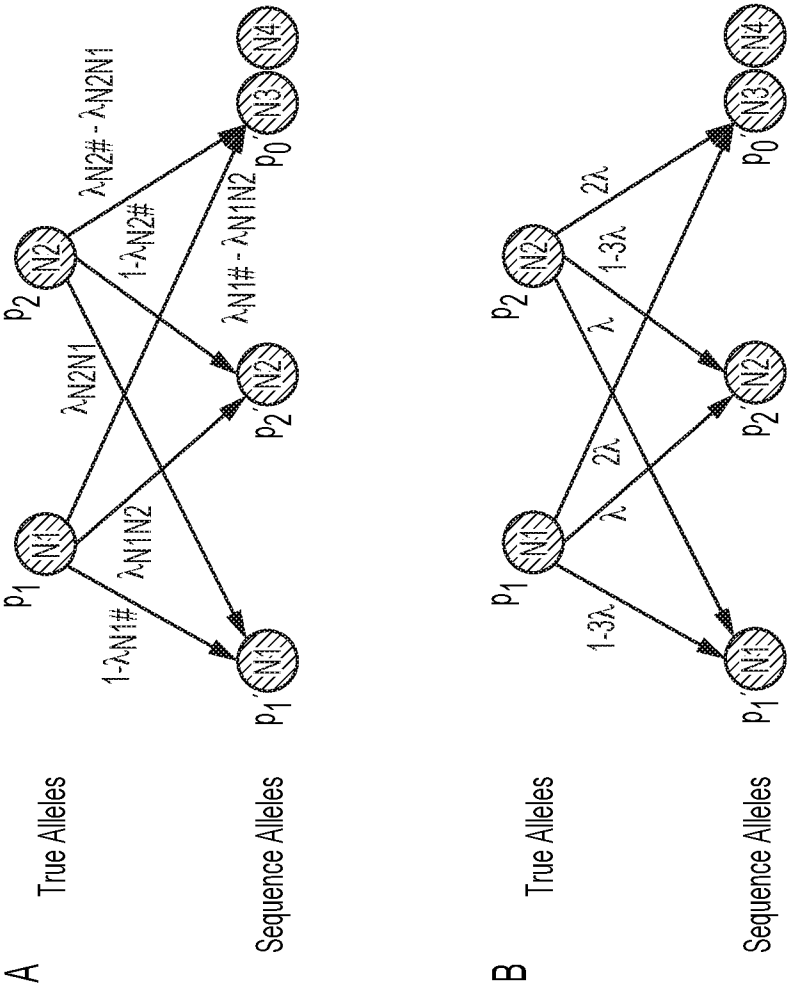
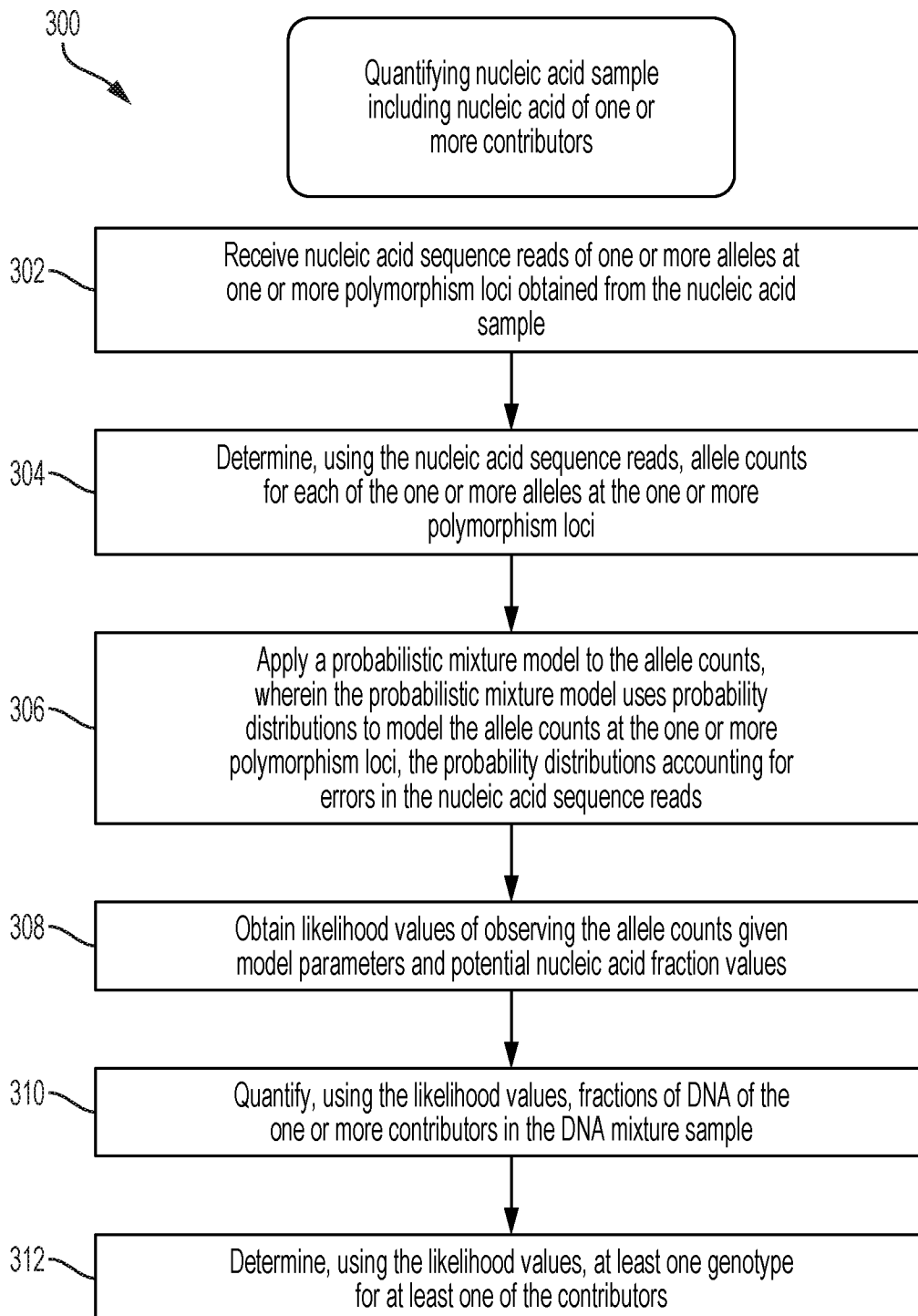
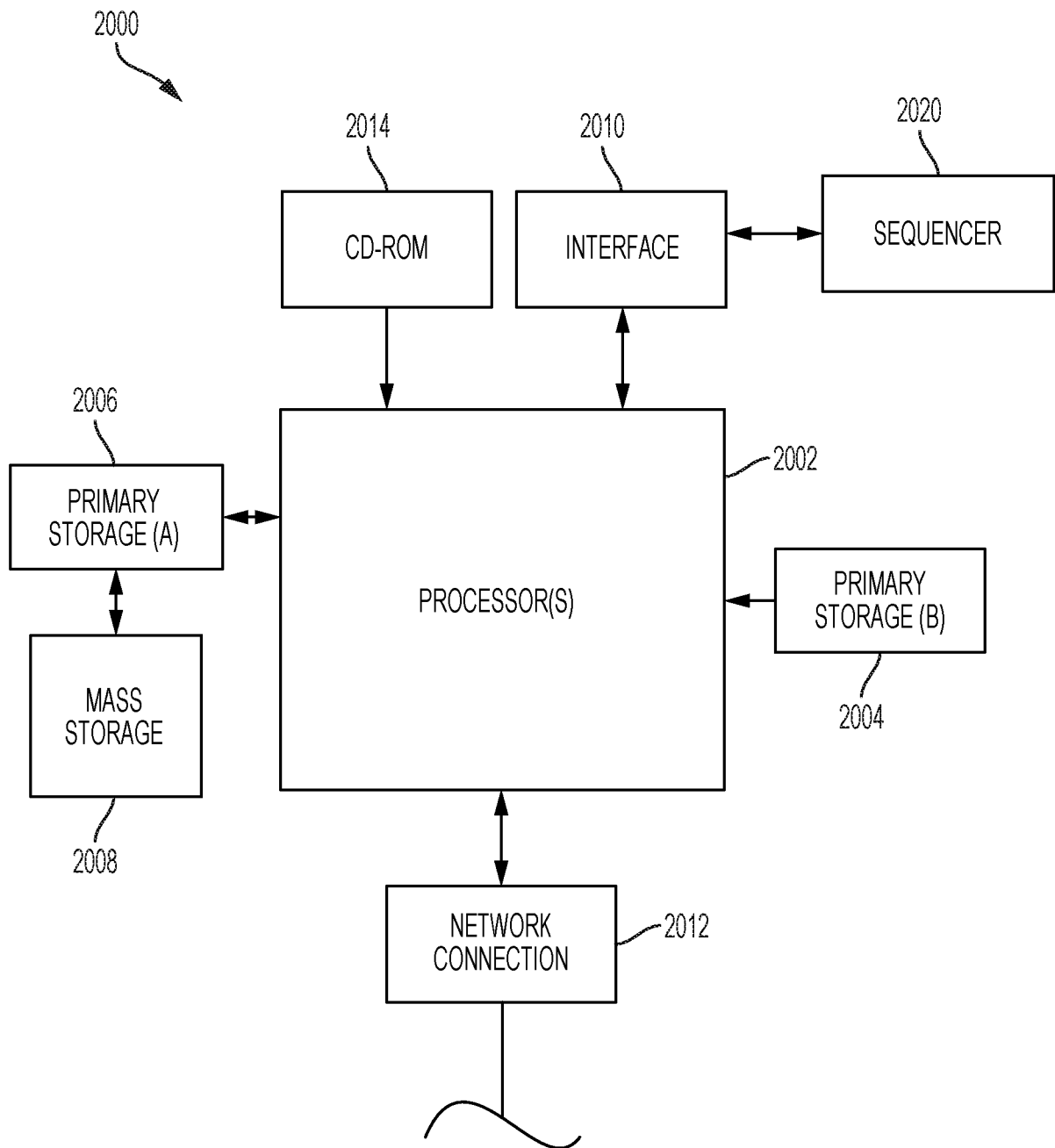


Figure 2C

5/11

**Figure 3**

6/11

**Figure 4**

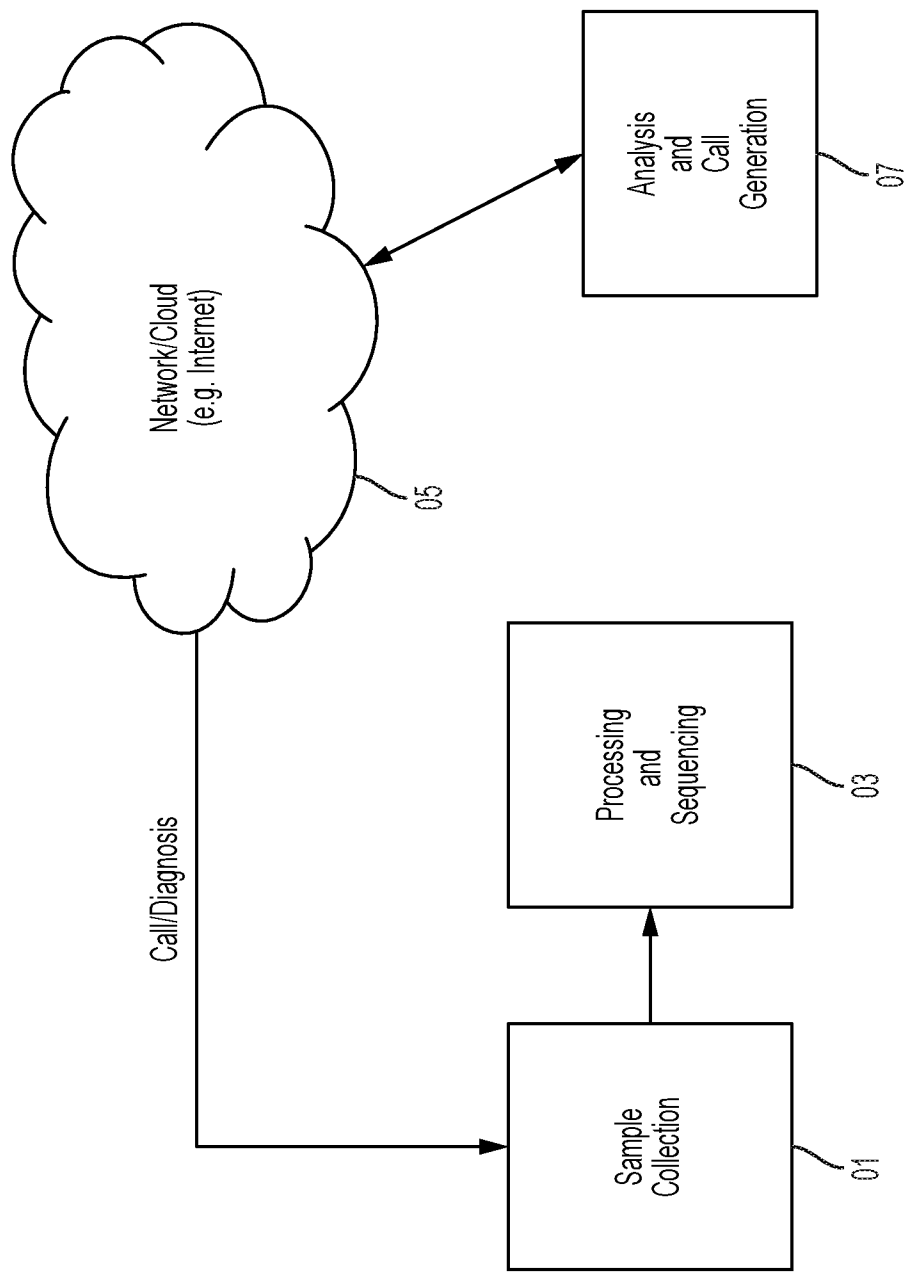


Figure 5

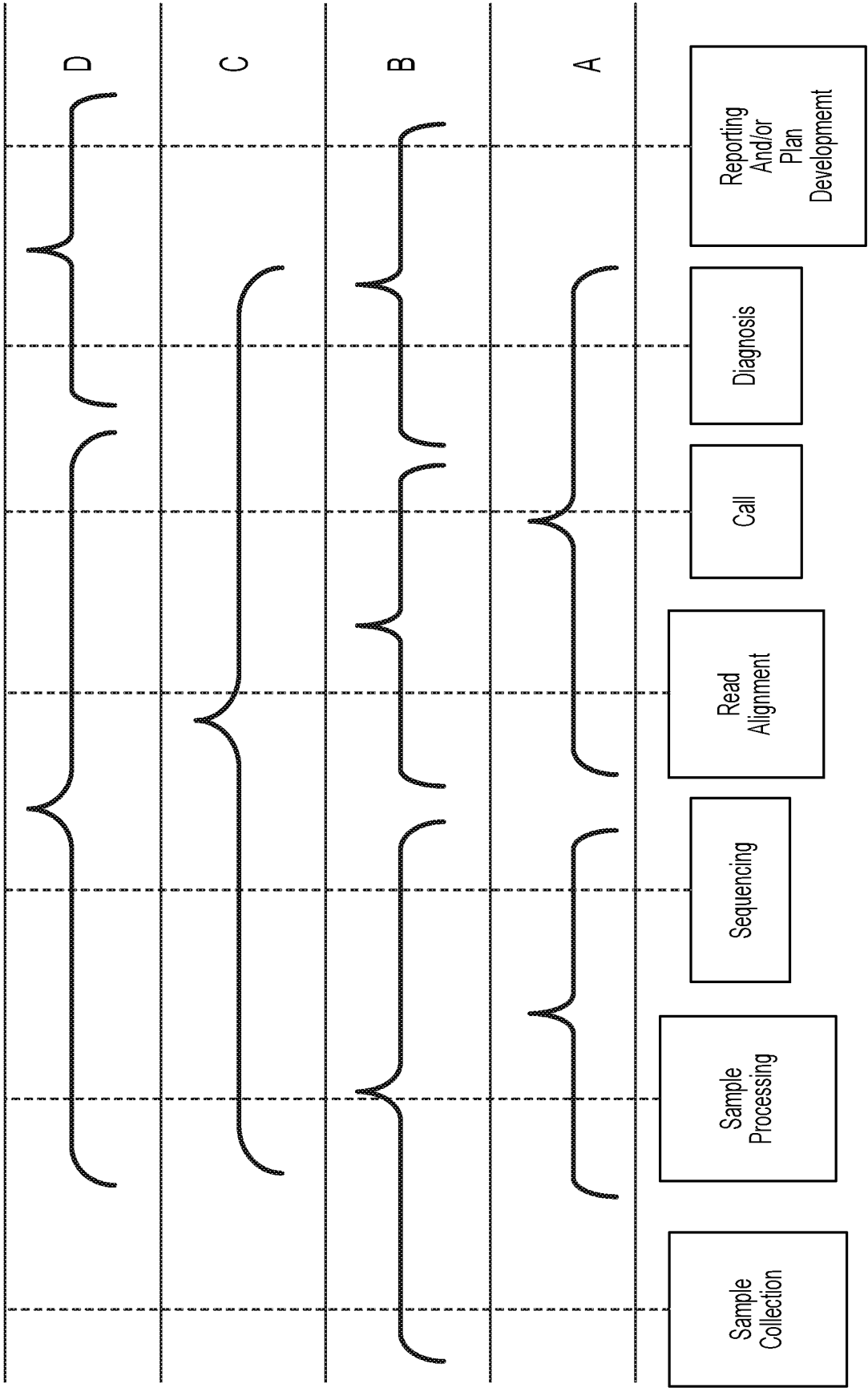


Figure 6

9/11

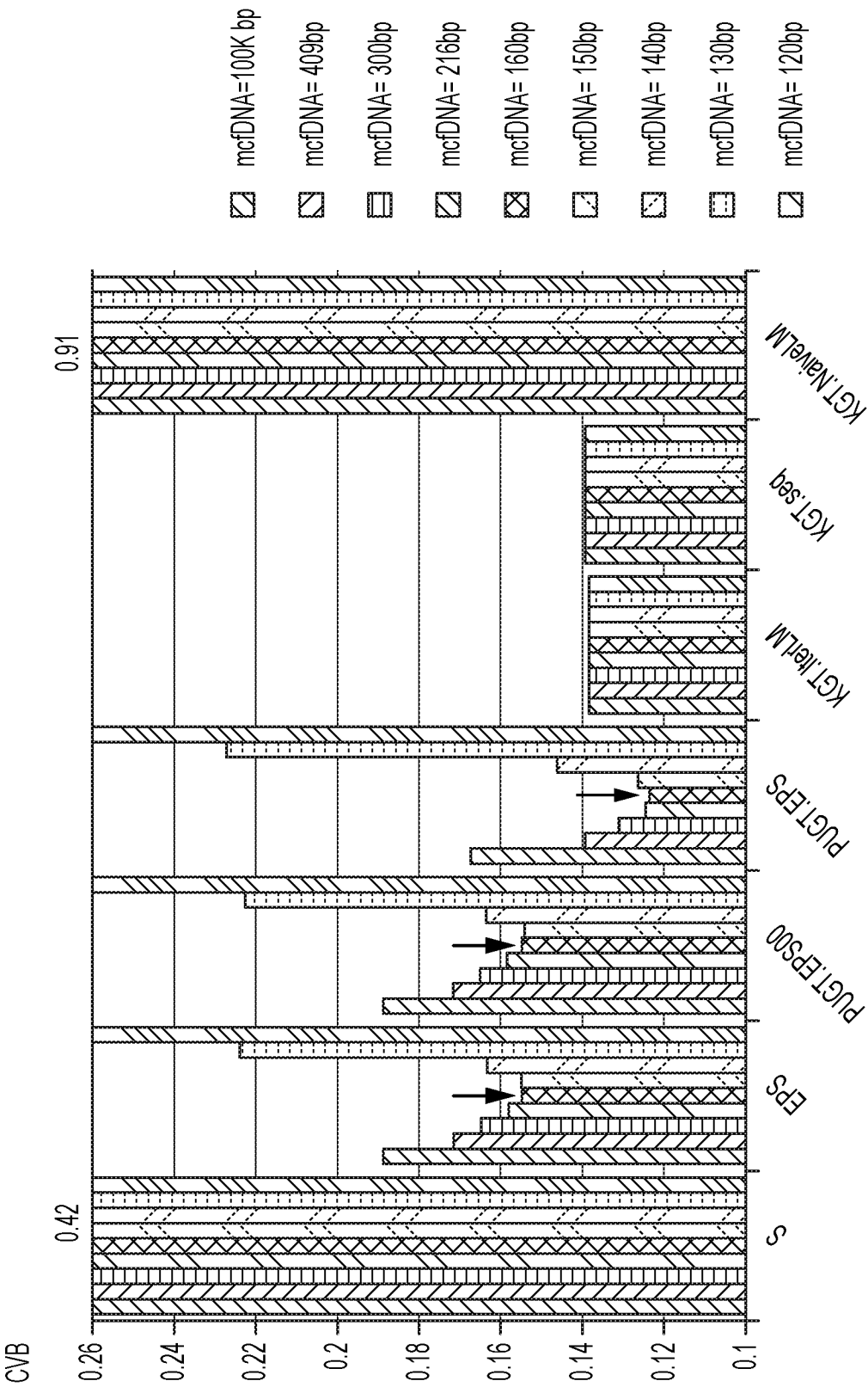
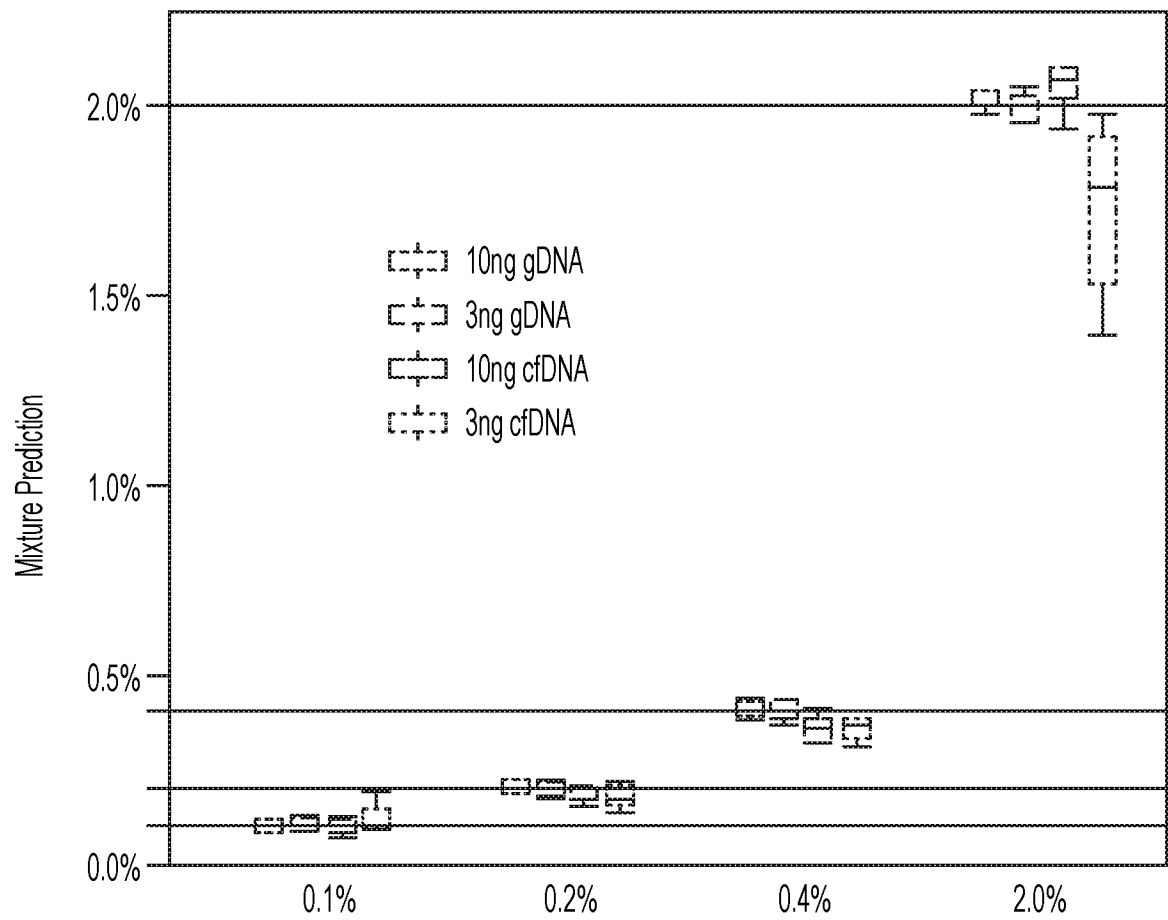


Figure 7

10/11

**Figure 8**

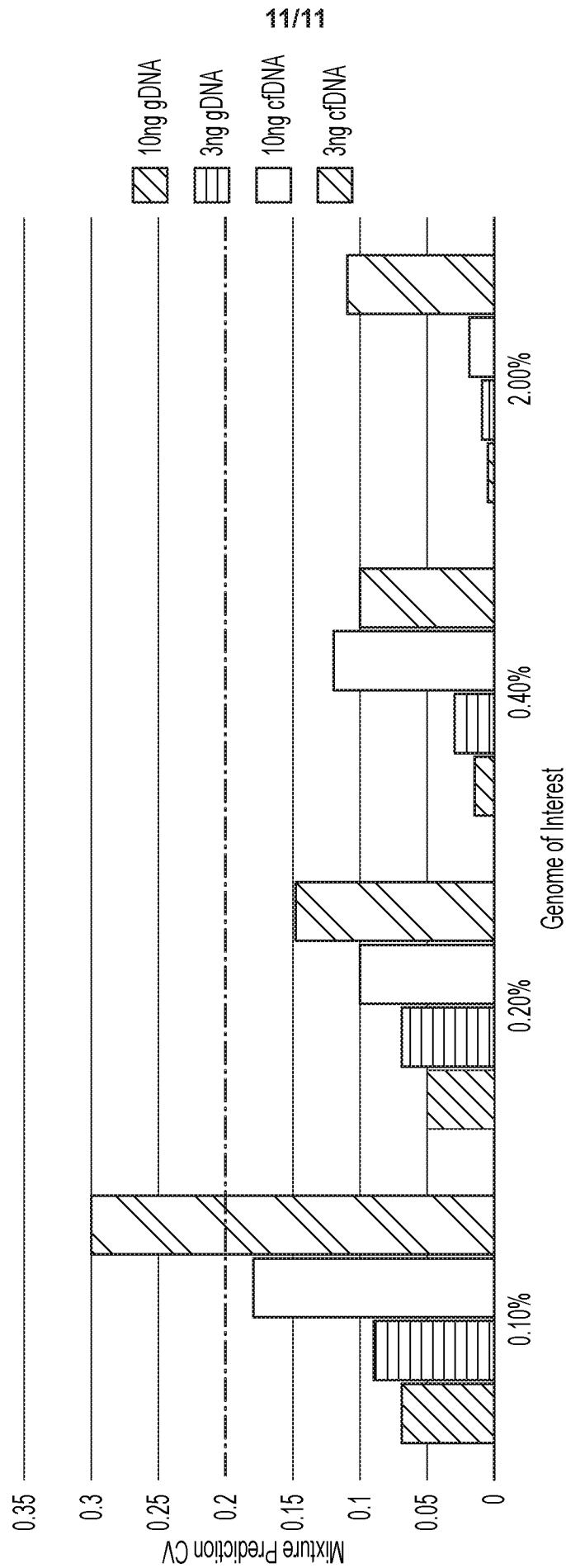


Figure 9

INTERNATIONAL SEARCH REPORT

International application No
PCT/US2018/038342

A. CLASSIFICATION OF SUBJECT MATTER
INV. G06F19/22
ADD.

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)
G06F

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

EPO-Internal, WPI Data

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	WO 2013/130848 A1 (NATERA INC [US]; RABINOWITZ MATTHEW [US]; BANER JOHAN [US]; GEMELOS GE) 6 September 2013 (2013-09-06)	1-76
Y	p.5 l.29-31, p.6 l.4-5, l.16-17, p.8 l.21-23, p.21 l.3-21, p.23 par.2, p.28 l.17-p.29 l.31, p.71 l.5-20, p.78 l.10-p.79 l.7, p.99 l.5-11, p.105 l.24-31, p.123 l.28 et seq., "Ploidy Calling Informatics Methods", p.133 l.16-p.135 l.15; claims 3, 21, 23 ----- -/--	65,66,76



Further documents are listed in the continuation of Box C.



See patent family annex.

* Special categories of cited documents :

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier application or patent but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art

"&" document member of the same patent family

Date of the actual completion of the international search

19 September 2018

Date of mailing of the international search report

26/09/2018

Name and mailing address of the ISA/

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040,
Fax: (+31-70) 340-3016

Authorized officer

Werner, Andreas

INTERNATIONAL SEARCH REPORT

International application No
PCT/US2018/038342

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	JIANG P ET AL: "FetalQuant: deducing fractional fetal DNA concentration from massively parallel sequencing of DNA in maternal plasma", BIOINFORMATICS, vol. 28, no. 22, 8 September 2012 (2012-09-08), pages 2883-2890, XP055127069, ISSN: 1367-4803, DOI: 10.1093/bioinformatics/bts549	1,2,4-6, 10-76
Y	the whole document	65,66,76
A	NIELSEN R ET AL: "Genotype and SNP calling from next-generation sequencing data", NATURE REVIEWS GENETICS, vol. 12, no. 6, 1 June 2011 (2011-06-01), pages 443-451, XP055046801, ISSN: 1471-0056, DOI: 10.1038/nrg2986 abstract p.445, col.1 par.3 and middle paragraph in the left-hand margin	66,76
Y	KIM S Y ET AL: "Estimation of allele frequency and association mapping using next-generation sequencing data", BMC BIOINFORMATICS, BIOMED CENTRAL, LONDON, GB, vol. 12, no. 1, 11 June 2011 (2011-06-11), page 231, XP021102534, ISSN: 1471-2105, DOI: 10.1186/1471-2105-12-231 abstract p.3 col.1 par.4	65,76
A	WO 2016/167408 A1 (EONE DIAGNOMICS GENOME CENT CO LTD) 20 October 2016 (2016-10-20) [35], [41]; claims 1, 5, 15; figure 1	3,7-9
A,P	-& EP 3 285 193 A1 (EONE DIAGNOMICS GENOME CENTER CO LTD [KR]) 21 February 2018 (2018-02-21) [0014], [0015]; claims 1, 5, 15; figure 1	3,7-9
A	WO 2014/014498 A1 (VERINATA HEALTH INC [US]; RAVA RICHARD P [US]; CHINNAPPA MANJULA [US];) 23 January 2014 (2014-01-23) p.1 par.5-p.2 par.2, p.193 l.19-p.200 l.7; claim 44; figures 1, 14-17	1-76

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No

PCT/US2018/038342

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
WO 2013130848	A1	06-09-2013	NONE
WO 2016167408	A1	20-10-2016	CN 107660234 A 02-02-2018
		EP 3285193 A1 21-02-2018	
		JP 2018514205 A 07-06-2018	
		KR 20160122563 A 24-10-2016	
		WO 2016167408 A1 20-10-2016	
EP 3285193	A1	21-02-2018	CN 107660234 A 02-02-2018
			EP 3285193 A1 21-02-2018
			JP 2018514205 A 07-06-2018
			KR 20160122563 A 24-10-2016
			WO 2016167408 A1 20-10-2016
WO 2014014498	A1	23-01-2014	AU 2013204615 A1 06-02-2014
			CA 2878246 A1 23-01-2014
			EP 2875149 A1 27-05-2015
			EP 2877594 A1 03-06-2015
			WO 2014014497 A1 23-01-2014
			WO 2014014498 A1 23-01-2014