

(19) **DANMARK**

(10) **DK/EP 3354748 T3**



(12) **Oversættelse af  
europæisk patentskrift**

Patent- og  
Varemærkestyrelsen

- 
- (51) Int.Cl.: **C 12 Q 1/6869 (2018.01)** **C 12 Q 1/6809 (2018.01)** **C 12 Q 1/6827 (2018.01)**
- (45) Oversættelsen bekendtgjort den: **2020-01-02**
- (80) Dato for Den Europæiske Patentmyndigheds bekendtgørelse om meddelelse af patentet: **2019-10-16**
- (86) Europæisk ansøgning nr.: **17209781.8**
- (86) Europæisk indleveringsdag: **2013-03-08**
- (87) Den europæiske ansøgnings publiceringsdag: **2018-08-01**
- (30) Prioritet: **2012-03-08 US 201261608623 P** **2012-04-06 US 201261621451 P**
- (62) Stamansøgningsnr: **13757943.9**
- (84) Designerede stater: **AL AT BE BG CH CY CZ DE DK EE ES FI FR GB GR HR HU IE IS IT LI LT LU LV MC MK MT NL NO PL PT RO RS SE SI SK SM TR**
- (73) Patenthaver: **The Chinese University Of Hong Kong, Knowledge Transfer Office , Sha Tin , New Territories, Hong Kong, Kina**
- (72) Opfinder: **LO, Yuk Ming Dennis, 4th Floor 7 King Tak Street, Homantin, Kowloon, Hong Kong, Kina**  
**CHAN, Kwan Chee, Flat A 13/F Block 34 Broadway Street, Mei Foo Sun Chuen, Kowloon, Hong Kong, Kina**  
**ZHENG, Wenli, 147 Red Globe Street, North Augusta, SC 29860, USA**  
**JIANG, Peiyong, Flat 7 1st Floor of Block B, Kwong Lam Court, Nos 62-66 Siu Lok Yuen Road, Shatin, New Territories, Hong Kong, Kina**  
**LIAO, Jiawei, Flat 16 10/F Yat Yan House, Yat Nga court, Tai Po Market, New Territories, Hong Kong, Kina**  
**CHIU, Wai Kwun Rossa, House 31, Double Haven, 52 Ma Lok Path, Sha Tin, New Territories, Hong Kong, Kina**
- (74) Fuldmægtig i Danmark: **Zacco Denmark A/S, Arne Jacobsens Allé 15, 2300 København S, Danmark**
- (54) Benævnelse: **Størrelsesbaseret DNA-analyse til klassificering af cancer**
- (56) Fremdragne publikationer:  
**FLORENT MOULIERE ET AL: "High fragmentation characterizes tumour-derived circulating DNA", PLOS ONE, PUBLIC LIBRARY OF SCIENCE, US, vol. 6, no. 9, 6 September 2011 (2011-09-06), pages e23418-1, XP002730500, ISSN: 1932-6203, DOI: 10.1371/JOURNAL.PONE.0023418**  
**H. C. FAN ET AL: "Analysis of the Size Distributions of Fetal and Maternal Cell-Free DNA by Paired-End Sequencing", CLINICAL CHEMISTRY, vol. 56, no. 8, 1 August 2010 (2010-08-01) , pages 1279-1286, XP055026439, ISSN: 0009-9147, DOI: 10.1373/clinchem.2010.144188**



## DESCRIPTION

**[0001]** The discovery of cell-free fetal DNA in maternal plasma has opened up new possibilities for noninvasive prenatal diagnosis (Lo YMD et al. *Lancet*1997;350:485-487). The mean/median fractional fetal DNA concentration has been reported to be approximately 3% to 10% (Lo YMD et al. *Am J Hum Genet* 1998;62:768-775; Lun FMF et al. *Clin Chem* 2008;54:1664-1672). The fractional fetal DNA concentration is an important parameter which affects the performance of noninvasive prenatal diagnostic tests using maternal plasma DNA. For example, for the noninvasive prenatal diagnosis of fetal chromosomal aneuploidies (e.g. trisomy 21, trisomy 18 or trisomy 13), the higher the fractional fetal DNA concentration is, the higher will be the overrepresentation of DNA sequences derived from the aneuploid chromosome in maternal plasma. Indeed, it has been demonstrated that for every two times reduction in the fractional fetal DNA concentration in maternal plasma, the number of molecules that one would need to count to achieve aneuploidy detection would be four times (Lo YMD et al. *Proc Natl Acad Sci USA* 2007;104:13116-13121).

**[0002]** For the noninvasive prenatal detection of fetal trisomy by random massively parallel sequencing, the fractional fetal DNA concentration of a sample would affect the amount of sequencing that one would need to perform to achieve a robust detection (Fan HC and Quake SR. *PLoS One* 2010;5:e10439). Indeed, a number of groups have included a quality control step in which the fractional fetal DNA concentration is first measured and only samples that contain more than a minimum fractional fetal DNA concentration would be eligible to generate a diagnostic result (Palomaki GE et al. *Genet Med* 2011;13:913-920). Other groups have included the fractional fetal DNA concentration in their diagnostic algorithm for estimating the risk that a particular maternal plasma sample is obtained from an aneuploid pregnancy (Sparks AB et al. *Am J Obstet Gynecol* 2012; 206: 319.e1-9). Fan et al., "Analysis of Size Distributions of Fetal and Maternal Cell-Free DNA by Paired-End Sequencing", *Clinical Chemistry*, vol. 56, no. 8, 1 August 2010, pp. 1279-1286 discloses measurement of the length distribution of cell-free maternal plasma with single-base resolution using paired end sequencing for increasing the sensitivity in fetal aneuploidy detection.

**[0003]** In addition to aneuploidy detection, the fractional fetal DNA concentration also similarly affects noninvasive prenatal diagnostic tests conducted using maternal plasma DNA for detecting monogenic diseases, e.g. the hemoglobinopathies (Lun FMF et al. *Proc Natl Acad Sci USA*2008;105:19920-19925) and hemophilia (Tsui NBY et al. *Blood* 2011;117:3684-3691). The fractional fetal DNA concentration also affects the depth of sequencing that one would need to perform for constructing a fetal genomewide genetic and mutational map, as well as fetal whole genome sequencing (Lo YMD et al. *Sci Transl Med* 2010;2:61ra91 and US-2011/0105353).

**[0004]** A number of methods have been described for measuring the fractional fetal DNA concentration. One approach is to measure the concentration of a fetal-specific, paternally-inherited sequence that is absent from the maternal genome. Examples of such sequences

include the sequences on the Y chromosome that are present in male fetuses and sequences from the *RHD* gene in a Rhesus D positive fetus carried by a Rhesus D negative pregnant woman. One could also measure the total maternal plasma DNA using sequences that are present in both the mother and the fetus. To arrive at a fractional fetal DNA concentration, one could then calculate the ratio of the concentration of the fetal-specific, paternally-inherited sequence over the concentration of the total maternal plasma DNA.

**[0005]** Another example of sequences that one could use includes the use of single nucleotide polymorphisms (Lo YMD et al. *Sci Transl Med* 2010;2:61ra91). A disadvantage of using genetic markers for the measurement of the fractional fetal DNA concentration is that no single set of genetic markers would be informative for all fetus-mother pair. Yet another method that one could employ is the use of DNA sequences that exhibit fetal or placental-specific DNA methylation patterns in maternal plasma (Nygren AO et al. *Clin Chem* 2010;56:1627-1635). The potential disadvantage of the use of DNA methylation markers is that there may be inter-individual variation in the level of DNA methylation. Furthermore, methods that are used for the detection of DNA methylation markers are typically complex, including the use of methylation-sensitive restriction enzyme digestion (Chan KCA et al. *Clin Chem* 2008;52:2211-2218) or bisulfite conversion (Chim SSC et al. *Proc Natl Acad Sci USA* 2005;102:14753-14758) or methylated DNA immunoprecipitation (MeDIP) (Papageorgiou EA et al. *Nat Med* 2011; 17: 510-513).

**[0006]** Mouliere et al., "High fragmentation characterizes tumour-derived circulating DNA", *PLOS One*, vol. 6, no. 9, 6 September 2011 discloses a study in which the size distribution of circulating DNA (ctDNA) measured using Q-PCR analysis of plasma samples from cancer patients showed that ctDNA exhibits a specific amount profile based on ctDNA size and significant higher ctDNA fragmentation.

**[0007]** According to the present invention, there is provided a method of analyzing a biological sample of a subject, the biological sample including DNA originating from normal cells and potentially from cells associated with cancer, wherein at least some of the DNA is cell-free in the biological sample, the method comprising: (a) sequencing a plurality of DNA fragments to obtain a plurality of sequence reads comprising the outermost nucleotides at each end of a plurality of DNA fragments; (b) aligning the sequence reads to a reference genome, thereby obtaining a set of genomic coordinates including genomic coordinates of the outermost nucleotides defining a size of a DNA fragment for each of the plurality of DNA fragments; (c) calculating a value of a parameter based on amounts of sequence reads of DNA fragments aligning to the set of genomic coordinates of (b) at multiple sizes; (d) comparing the value to a reference value; and (e) determining a classification of a level of cancer in the subject based on comparing of the value to the reference value.

**[0008]** Other embodiments are directed to systems, portable consumer devices, and computer readable media associated with methods described herein.

**[0009]** A better understanding of the nature and advantages of the present invention may be

gained with reference to the following detailed description and the accompanying drawings.

FIG. 1 shows a plot 100 of a size distribution of circulating cell-free DNA in maternal plasma according to embodiments of the present invention.

FIG. 2A shows a plot 200 of size distributions of fetal DNA in two maternal plasma samples (1<sup>st</sup> trimester pregnancies) with different fractional fetal DNA concentrations according to embodiments of the present invention.

Figure 2B shows a plot 250 of size distributions of DNA fragments in two maternal plasma samples (2<sup>nd</sup> trimester pregnancies) with different fractional fetal DNA concentrations according to embodiments of the present invention.

FIG. 3 is a flowchart of a method 300 illustrating a method of estimating a fractional concentration of clinically-relevant DNA in a biological sample according to embodiments of the present invention.

FIG. 4 is a plot 400 showing a size distribution (electropherogram) of maternal plasma DNA obtained using electrophoresis according to embodiments of the present invention.

FIG. 5A is a plot 500 showing a proportion of DNA fragments that are 150 bp or below for samples having various fetal DNA percentage in maternal plasma according to embodiments of the present invention.

FIG. 5B is a plot 550 showing a size ratio of the amounts of DNA fragments of  $\leq 150$  bp and DNA from 163 bp to 169 bp, which is labeled as  $(CF(\text{size} \leq 150) / \text{size}(163-169))$ .

FIG. 6A is a plot 600 showing a size ratio of the amounts of DNA fragments from 140 bp to 146 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-146) / \text{size}(163-169))$ .

FIG. 6B is a plot 650 showing a size ratio of the amounts of DNA fragments from 140 bp to 154 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-154) / \text{size}(163-169))$ .

FIG. 7 is a plot 700 showing a size ratio of the amounts of DNA fragments from 100 bp to 150 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(100-150) / \text{size}(163-169))$ .

FIG. 8 is a plot 800 showing a proportion of DNA fragments of 150 bp or below for samples having various fetal DNA percentages in maternal plasma according to embodiments of the present invention.

FIG. 9A is a plot 900 showing a size ratio of the amounts of DNA fragments of  $\leq 150$  bp and DNA from 163 bp to 169 bp, which is labeled as  $(CF(\text{size} \leq 150) / \text{size}(163-169))$ .

FIG. 9B is a plot 950 showing a size ratio of the amounts of DNA fragments from 140 bp to 146 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-146) / \text{size}(163-169))$ .

FIG. 10A is a plot 1000 showing a size ratio of the amounts of DNA fragments from 140 bp to 154 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-154) / \text{size}(163-169))$ .

FIG. 10B is a plot 1005 showing a size ratio of the amounts of DNA fragments from 100 bp to 150 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(100-150)/\text{size}(163-169))$ .

FIG. 11 is a plot showing a size ratio plotted vs. fetal DNA percentage for the size of repeat elements according to embodiments of the present invention.

FIG. 12A is an electropherogram 1200 that may be used to determine a size ratio according to embodiments of the present invention.

FIG. 12B a plot 1250 showing a size ratio of the amounts of DNA fragments from 200 bp to 267 bp and DNA from 290 bp to 294 bp for samples having various fetal DNA percentage in maternal plasma according to embodiments of the present invention.

FIG. 13 is a flowchart of a method 1300 for determining calibration data points from measurements made from calibration samples according to embodiments of the present invention.

FIG. 14A is a plot 1400 of a size ratio against the fractional concentration of fetal DNA for the training set according to embodiments of the present invention.

FIG. 14B is a plot 1450 of fractional concentrations deduced (estimated) from linear function 1410 of FIG. 14A against the fractional concentrations measured using fetal-specific sequences according to embodiments of the present invention.

FIG. 15A is a plot 1500 showing a proportion of DNA fragments of 150 bp or below for samples having various tumor DNA percentages in plasma of two hepatocellular carcinoma (HCC) patients before and after tumor resection according to embodiments of the present invention.

FIG. 15B is a plot 1550 showing a size ratio of the amounts of DNA fragments of  $\leq 150$  bp and DNA from 163 bp to 169 bp, which is labeled as  $(CF(\text{size} \leq 150)/\text{size}(163-169))$ , for two HCC patients before and after tumor resection.

FIG. 16A is a plot 1600 showing a size ratio of the amounts of DNA fragments from 140 bp to 146 bp and DNA from 163 bp to 169 bp, which is labeled  $(\text{size}(140-146)/\text{size}(163-169))$ , for two HCC patients before and after tumor resection.

FIG. 16B is a plot 1650 showing a size ratio of the amounts of DNA fragments from 140 bp to 154 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-154)/\text{size}(163-169))$ , for two HCC patients before and after tumor resection.

FIG. 17 is a plot 1700 showing a size ratio of the amounts of DNA fragments from 100 bp to 150 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(100-150)/\text{size}(163-169))$ , for two HCC patients before and after tumor resection.

FIG. 18A is a plot 1800 showing a proportion of DNA fragments of 150 bp or below for HCC patients before and after tumor resection.

FIG. 18B is a plot 1850 showing a size ratio of the amounts of DNA fragments of  $\leq 150$  bp and

DNA from 163 bp to 169 bp, which is labeled as  $(CF(\text{size} \leq 150)/\text{size}(163-169))$ , for HCC patients before and after tumor resection.

FIG. 19A is a plot 1900 showing a size ratio of the amounts of DNA fragments from 140 bp to 146 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-146)/\text{size}(163-169))$ , for HCC patients before and after tumor resection.

FIG. 19B is a plot 1950 showing a size ratio of the amounts of DNA fragments from 140 bp to 154 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-154)/\text{size}(163-169))$ , for HCC patients before and after tumor resection.

FIG. 20 is a plot 2000 showing a size ratio of the amounts of DNA fragments from 100 bp to 150 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(100-150)/\text{size}(163-169))$ , for HCC patients before and after tumor resection.

FIG. 21 is a flowchart illustrating a method 2100 for analyzing a biological sample of an organism to determine a classification of a level of cancer according to embodiments of the present invention.

FIG. 22 is a table 2200 showing some common chromosomal aberrations seen in various types of cancers.

FIG. 23 shows a block diagram of an example computer system 2300 usable with system and methods according to embodiments of the present invention.

**[0010]** The term "*biological sample*" as used herein refers to any sample that is taken from a subject (e.g., a human, such as a pregnant woman) and contains one or more nucleic acid molecule(s) of interest. Examples include plasma, saliva, pleural fluid, sweat, ascitic fluid, bile, urine, serum, pancreatic juice, stool and cervical smear samples. The biological sample may be obtained from a human, an animal, or other suitable organism. A "*calibration sample*" corresponds to a biological sample whose clinically-relevant DNA fraction is known or determined via a calibration method, e.g., using an allele specific to the clinically relevant DNA. Examples of clinically-relevant DNA are fetal DNA in maternal plasma or tumor DNA in a patient's plasma.

**[0011]** As used herein, the term "*locus*" or its plural form "*loci*" is a location or address of any length of nucleotides (or base pairs) which has a variation across genomes. The term "*sequence read*" refers to a sequence obtained from all or part of a nucleic acid molecule, e.g., a DNA fragment. In one embodiment, just one end of the fragment is sequenced. Alternatively, both ends (e.g., about 30 bp from each end) of the fragment can be sequenced to generate two sequence reads. The paired sequence reads can then be aligned to a reference genome, which can provide a length of the fragment. In yet another embodiment, a linear DNA fragment can be circularized, e.g., by ligation, and the part spanning the ligation site can be sequenced.

**[0012]** The term "*universal sequencing*" refers to sequencing where adapters are added to the end of a fragment, and the primers for sequencing attached to the adapters. Thus, any fragment can be sequenced with the same primer, and thus the sequencing can be random.

**[0013]** The term fractional fetal DNA concentration is used interchangeably with the terms fetal DNA proportion and fetal DNA fraction, and refers to the proportion of fetal DNA molecules that are present in a biological sample (e.g., maternal plasma or serum sample) that is derived from the fetus (Lo YMD et al. *Am J Hum Genet* 1998;62:768-775; Lun FMF et al. *Clin Chem* 2008;54:1664-1672). Similarly, the terms fractional tumor DNA concentration may be used interchangeably with the terms tumor DNA proportion and tumor DNA fraction, and refers to the proportion of tumor DNA molecules that are present in a biological sample.

**[0014]** The term "*size profile*" generally relates to the sizes of DNA fragments in a biological sample. A size profile may be a histogram that provides a distribution of an amount of DNA fragments at a variety of sizes. Various statistical parameters (also referred to as size parameters or just parameter) can be used to distinguish one size profile to another. One parameter is the percentage of DNA fragment of a particular size or range of sizes relative to all DNA fragments or relative to DNA fragments of another size or range.

**[0015]** Examples of "clinically-relevant" DNA include fetal DNA in maternal plasma and tumor DNA in the patient's plasma. Another example include the measurement of the amount of graft-associated DNA in the plasma of a transplant patient. A further example include the measurement of the relative amounts of hematopoietic and nonhematopoietic DNA in the plasma of a subject. This latter embodiment can be used for detecting or monitoring or prognosticating pathological processes or injuries involving hematopoietic and/or nonhematopoietic tissues.

**[0016]** A "*calibration data point*" includes a "*calibration value*" and a measured or known fractional concentration of the DNA of interest (i.e., the clinically-relevant DNA). The calibration value is a value of a size parameter as determined for a calibration sample, for which the fractional concentration of the clinically-relevant DNA is known. The calibration data points may be defined in a variety of ways, e.g., as discrete points or as a calibration function (also called a calibration curve or calibrations surface).

**[0017]** The term "*level of cancer*" can refer to whether cancer exists, a stage of a cancer, a size of tumor, how many deletions or amplifications of a chromosomal region are involved (e.g. duplicated or tripled), and/or other measure of a severity of a cancer. The level of cancer could be a number or other characters. The level could be zero. The level of cancer also includes premalignant or precancerous conditions associated with deletions or amplifications.

**[0018]** It is known that cell-free fetal DNA molecules in maternal plasma are generally shorter than the maternally-derived ones (Chan KCA et al. *Clin Chem* 2004;50:88-92; Lo YMD et al. *Sci Transl Med* 2010;2:61ra91). The presence of fetal DNA results in a shift in the overall size distribution of maternal plasma DNA and the degree of shifting is associated with the fractional

concentration of fetal DNA. By measuring particular values of the size profile of maternal plasma DNA, embodiments can obtain the fractional fetal DNA concentration in maternal plasma.

**[0019]** Apart from applications in noninvasive prenatal diagnosis, embodiments can also be used for measuring the fractional concentration of clinically useful nucleic acid species of different sizes in biological fluids, which can be useful for cancer detection, transplantation, and medical monitoring. It has previously been shown that tumor-derived DNA is shorter than the non-cancer-derived DNA in a cancer patient's plasma (Diehl F et al. Proc Natl Acad Sci USA 2005;102:16368-16373). In the transplantation context, it has been shown hematopoietic-derived DNA is shorter than nonhematopoietic DNA (Zheng YW et al. Clin Chem 2012;58:549-558). For example, if a patient receives a liver from a donor, then the DNA derived from the liver (a nonhematopoietic organ in the adult) will be shorter than hematopoietic-derived DNA in the plasma (Zheng YW et al. Clin Chem 2012;58:549-558). Similarly, in a patient with myocardial infarction or stroke, the DNA released by the damaged nonhematopoietic organs (i.e. the heart and brain, respectively) would be expected to result in a shift in the size profile of plasma DNA towards the shorter spectrum.

## I. SIZE DISTRIBUTION

**[0020]** To demonstrate embodiments, we show in the following examples that one can measure the size profile, for example, by paired-end massively parallel sequencing or by electrophoresis (e.g. using a Bioanalyzer). The latter example is particularly useful because electrophoresis using a Bioanalyzer is a quick and relatively cheap procedure. This would allow one to rapidly perform this analysis as a quality control measure before one would subject a plasma DNA sample to the relatively expensive sequencing process.

**[0021]** FIG. 1 shows a plot 100 of a size distribution of circulating cell-free DNA in maternal plasma according to embodiments of the present invention. A size distribution can be obtained by measuring a size of DNA fragments and then counting the number of DNA fragments at various sizes, e.g., within the range of 50 bases to about 220 bases. Plot 100 shows two distributions. Distribution 110 is for all of the DNA fragments in the maternal plasma sample, and distribution 120 is only for DNA that is from the fetus. The horizontal axis is the size in base pairs (bp) of the DNA fragments. The vertical axis is the percentage of measured DNA fragments

**[0022]** In FIG. 1, the size distribution of fetal-derived DNA in maternal plasma has been shown to be shorter than that of the maternally derived ones (Chan KC et al. ClinChem 2004; 50:88-92.) Recently, we have used paired-end massively parallel sequencing analysis to determine the high-resolution size distribution of the fetal-specific DNA and total DNA (mainly derived from the mother) in a pregnant woman. We showed that a main difference between the two species of DNA is that there is a reduction in the fraction of 166 bp DNA fragments and an increase proportion of shorter DNA of below 150 bp for the fetal-derived DNA (Lo YM et al. Sci

Transl Med 2010 2:61ra91).

**[0023]** Herein, we outline how an analysis of a size distribution of total DNA fragments in a maternal plasma sample (an example of a biological sample) would be useful for determining the fractional concentration of fetal DNA in maternal plasma. The increased fractional concentration of fetal DNA in maternal plasma would result in the shortening of the overall size distribution of the total DNA. In one embodiment, the relative abundance (an example of a parameter) of DNA fragments of approximately 144 bp and DNA fragments of approximately 166 bp could be used to reflect the fractional concentration of fetal DNA. In another embodiment, other parameters or combination of parameters regarding a size profile can be used to reflect the size distribution of plasma DNA.

**[0024]** FIG. 2A shows a plot 200 of size distributions of fetal DNA in two maternal plasma samples (1<sup>st</sup> trimester pregnancies) with different fractional fetal DNA concentrations according to embodiments of the present invention. Both of these two pregnant women were carrying male fetuses. The fractional fetal DNA concentrations were determined from the proportion of sequences from the Y chromosome among the total sequenced DNA fragments. Both samples were taken from pregnant women during the first trimester of their pregnancies. Case 338 (solid line, fractional fetal DNA concentration 10%) had a lower fractional fetal DNA concentration than Case 263 (dotted line, fractional fetal DNA concentration 20%). When compared with Case 263, Case 338 had a higher peak at 166 bp but the peaks for size below 150 bp were lower. In other words, DNA fragments shorter than 150 bp were more abundant in Case 263 whereas the fragments of approximately 166 bp were more abundant in Case 338. These observations are consistent with the hypothesis that the relative amounts of short and long DNA may be correlated to the fractional fetal DNA concentration.

**[0025]** Figure 2B shows a plot 250 of size distributions of DNA fragments in two maternal plasma samples (2<sup>nd</sup> trimester pregnancies) with different fractional fetal DNA concentrations according to embodiments of the present invention. Both samples were taken from pregnant women during the second trimester. Both of these two pregnant women were carrying male fetuses. The fractional fetal DNA concentrations were determined from the proportion of sequences from the Y chromosome among the total sequenced DNA fragments. Similar to the previous example, case 5415 (dotted line, with higher fractional fetal DNA concentration 19%) had higher peaks for sizes below 150 bp whereas case 5166 (solid line, with lower fractional fetal DNA concentration 12%) had a higher peak at 166 bp.

**[0026]** The correlation of different values of a size parameter to values of fractional fetal DNA concentration is shown in data plots below. Additionally, the size of fragments of tumor DNA is correlated to the percentage of tumor DNA fragments in a sample with tumor DNA fragments and DNA fragments from normal cells. Thus, the size of tumor fragments can also be used to determine the percentage of tumor fragments in the sample.

## II. METHOD

**[0027]** Since the size of DNA fragments is correlated to a fractional concentration (also referred to as a percentage), embodiments can use this correlation to determine a fractional concentration of a particular type of DNA (e.g., fetal DNA or DNA from a tumor) in a sample. The particular type of DNA is clinically-relevant as that is the fractional concentration being estimated. Accordingly, a method can estimate a fractional concentration of clinically-relevant DNA in a biological sample based on a measured size of the DNA fragments.

**[0028]** FIG. 3 is a flowchart of a method 300 illustrating a method of estimating a fractional concentration of clinically-relevant DNA in a biological sample according to embodiments of the present invention. The biological sample includes the clinically-relevant DNA and other DNA. The biological sample may be obtained from a patient, e.g., a female subject pregnant with a fetus. In another embodiment, the patient may have or be suspected of having a tumor. In one implementation, the biological sample may be received at a machine, e.g., a sequencing machine, which outputs measurement data (e.g., sequence reads) that can be used to determine sizes of the DNA fragments. Method 300 may be performed wholly or partially with a computer system, as can other methods described herein.

**[0029]** At block 310, amounts of DNA fragments corresponding to various sizes are measured. For each size of a plurality of sizes, an amount of a plurality of DNA fragments from the biological sample corresponding to the size can be measured. For instance, the number of DNA fragments having a length of 140 bases may be measured. The amounts may be saved as a histogram. In one embodiment, a size of each of the plurality of nucleic acids from the biological sample is measured, which may be done on an individual basis (e.g., by single molecule sequencing) or on a group basis (e.g., via electrophoresis). The sizes may correspond to a range. Thus, an amount can be for DNA fragments that have a size within a particular range.

**[0030]** The plurality of DNA fragments may be chosen at random or preferentially selected from one or more predetermined regions of a genome. For example, targeted enrichment may be performed, as described above. In another embodiment, DNA fragments may be randomly sequenced (e.g., using universal sequencing), and the resulting sequence reads can be aligned to a genome corresponding to the subject (e.g., a reference human genome). Then, only DNA fragments whose sequence reads align to the one or more predetermined regions may be used to determine the size.

**[0031]** In various embodiments, the size can be mass, length, or other suitable size measures. The measurement can be performed in various ways, as described herein. For example, paired-end sequencing and alignment of DNA fragments may be performed, or electrophoresis may be used. A statistically significant number of DNA fragments can be measured to provide an accurate size profile of the biological sample. Examples of a statistically significant number of DNA fragments include greater than 100,000; 1,000,000; 2,000,000, or other suitable values, which may depend on the precision required.

**[0032]** In one embodiment, the data obtained from a physical measurement, such as paired-end sequencing or electrophoresis, can be received at a computer and analyzed to accomplish the measurement of the sizes of the DNA fragments. For instance, the sequence reads from the paired-end sequencing can be analyzed (e.g., by alignment) to determine the sizes. As another example, the electropherogram resulting from electrophoresis can be analyzed to determine the sizes. In one implementation, the analyzing of the DNA fragments does include the actual process of sequencing or subjecting DNA fragments to electrophoresis, while other implementations can just perform an analysis of the resulting data.

**[0033]** At block 320, a first value of a first parameter is calculated based on the amounts of DNA fragments at multiple sizes. In one aspect, the first parameter provides a statistical measure of a size profile (e.g., a histogram) of DNA fragments in the biological sample. The parameter may be referred to as a size parameter since it is determined from the sizes of the plurality of DNA fragments.

**[0034]** The first parameter can be of various forms. Such a parameter is a number of DNA fragments at a particular size divided by the total number of fragments, which may be obtained from a histogram (any data structure providing absolute or relative counts of fragments at particular sizes). As another example, a parameter could be a number of fragments at a particular size or within a particular range divided by a number of fragments of another size or range. The division can act as a normalization to account for a different number of DNA fragments being analyzed for different samples. A normalization can be accomplished by analyzing a same number of DNA fragments for each sample, which effectively provides a same result as dividing by a total number fragments analyzed. Other examples of parameters are described herein.

**[0035]** At block 330, one or more first calibration data points are obtained. Each first calibration data point can specify a fractional concentration of clinically-relevant DNA corresponding to a particular value (a calibration value) of the first parameter. The fractional concentration can be specified as a particular concentration or a range of concentrations. A calibration value may correspond to a value of the first parameter (i.e., a particular size parameter) as determined from a plurality of calibration samples. The calibration data points can be determined from calibration samples with known fractional concentrations, which may be measured via various techniques described herein. At least some of the calibration samples would have a different fractional concentration, but some calibration samples may have a same fractional concentration

**[0036]** In various embodiments, one or more calibration points may be defined as one discrete point, a set of discrete points, as a function, as one discrete point and a function, or any other combination of discrete or continuous sets of values. As an example, a calibration data point could be determined from one calibration value of a size parameter (e.g., number of fragments in a particular size or size range) for a sample with a particular fractional concentration. A plurality of histograms can be used, with a different histogram for each calibration sample, where some of the calibration samples may have the same fractional concentration.

**[0037]** In one embodiment, measured values of a same size parameter from multiple samples at the same fractional concentration could be combined to determine a calibration data point for a particular fractional concentration. For example, an average of the values of the size parameter may be obtained from the size data of samples at the same fractional concentration to determine a particular calibration data point (or provide a range that corresponds to the calibration data point). In another embodiment, multiple data points with the same calibration value can be used to determine an average fractional concentration.

**[0038]** In one implementation, the sizes of DNA fragments are measured for many calibration samples. A calibration value of the same size parameter is determined for each calibration sample, where the size parameter may be plotted against the known fractional concentration of the sample. A function may then be fit to the data points of the plot, where the functional fit defines the calibration data points to be used in determining the fractional concentration for a new sample.

**[0039]** At block 340, the first value is compared to a calibration value of at least one calibration data point. The comparison can be performed in a variety of ways. For example, the comparison can be whether the first value is higher or lower than the calibration value. The comparison can involve comparing to a calibration curve (composed of the calibration data points), and thus the comparison can identify the point on the curve having the first value of the first parameter. For example, a calculated value  $X$  of the first parameter (as determined from the measured sizes of DNA in the new sample) can be used as input into a function  $F(X)$ , where  $F$  is the calibration function (curve). The output of  $F(X)$  is the fractional concentration. An error range can be provided, which may be different for each  $X$  value, thereby providing a range of values as an output of  $F(X)$ .

**[0040]** In step 350, the fractional concentration of the clinically-relevant DNA in the biological sample is estimated based on the comparison. In one embodiment, one can determine if the first value of the first parameter is above or below a threshold calibration value, and thereby determine if the estimated fractional concentration of the instant sample is above or below the fractional concentration corresponding to the threshold calibration value. For example, if the calculated first value  $X_1$  for the biological is above a calibration value  $X_c$  then the fractional concentration  $FC_1$  of the biological sample can be determined as being above the fractional concentration  $FC_c$  corresponding to  $X_c$ . This comparison can be used to determine if a sufficient fractional concentration exists in the biological sample to perform other tests, e.g., testing for a fetal aneuploidy. This relationship of above and below can depend on how the parameter is defined. In such an embodiment, only one calibration data point may be needed.

**[0041]** In another embodiment, the comparison is accomplished by inputting the first value into a calibration function. The calibration function can effectively compare the first value to calibration values by identifying the point on a curve corresponding to the first value. The estimated fractional concentration is then provided as the output value of the calibration function.

**[0042]** In one embodiment, the value of more than one parameter can be determined for the biological sample. For example, a second value can be determined for a second parameter, which corresponds to a different statistical measure of the size profile of DNA fragments in the biological sample. The second value can be determined using the same size measurements of the DNA fragments, or different size measurements. Each parameter can correspond to a different calibration curve. In one implementation, the different values can be compared independently to different calibration curves to obtain a plurality of estimated fractional concentrations, which may then be averaged or used to provide a range as an output.

**[0043]** In another implementation, a multidimensional calibration curve can be used, where the different values of the parameters can effectively be input to a single calibration function that outputs the fractional concentration. The single calibration function can result from a functional fit of all of the data points obtained from the calibration samples. Thus, in one embodiment, the first calibration data points and the second calibration data points can be points on a multidimensional curve, where the comparison includes identifying the multidimensional point having coordinates corresponding to the first value and the one or more second values.

### **III. DETERMINING SIZE**

**[0044]** The size distribution of plasma DNA can be determined, for example, but not limited to, using real-time PCR, electrophoresis and mass spectrometry analysis. In various embodiments, the measured size is a length, a molecular mass, or a measured parameter that is proportional to the length or mass, such as the mobility in a electrophoretogram and the time required to travel a fixed distance in electrophoresis or mass spectrometer. In another example, one can stain the DNA with an intercalating fluorescence dye, e.g. ethidium bromide or SYBR Green, where the amount of dye bound will be proportional to the length of the DNA molecule. One can determine the amount of dye bound by the intensity of the emitted fluorescence when UV light is shone on the sample. Some examples for measuring size and resulting data are described below.

#### ***A. First Fetal Sample Set Using Sequencing***

**[0045]** Table 1 shows sample information and sequencing analyses for an example involving a fetal DNA fraction. Plasma samples were collected from 80 pregnant women, each carrying a single male fetus. Among these 80 pregnant women, 39 were carrying a euploid fetus, 18 were carrying a trisomy 21 (T21) fetus, 10 were carrying a trisomy 18 (T18) fetus, and 13 were carrying a trisomy 13 (T13) fetus. A size distribution of plasma DNA was determined using paired-end massively parallel sequencing. Sequencing libraries of maternal plasma DNA were constructed as previously described (Lo YM et al. *Sci Transl Med* 2010; 2:61ra91), except that a 6-base barcode was introduced to the DNA molecules of each plasma sample through a

triple-primer PCR amplification.

**[0046]** Two samples were introduced into one sequencing lane (i.e. 2-plex sequencing). In other embodiments, more than two samples can be introduced into one sequencing lane, e.g. 6, or 12, or 20, or more than 20. All libraries were sequenced by a Genome Analyzer IIx (Illumina) using the 36-bp × 2 PE format. An additional 7 cycles of sequencing were performed to decode the index sequence on each sequenced plasma DNA molecule. The 36-bp sequence reads were aligned to the non-repeat-masked human reference genome (Hg18) (genome.ucsc.edu), using the Short Oligonucleotide Alignment Program 2 (SOAP2) (soap.genomics.org.cn). Paired end (PE) reads with individual members sequenced on the same cluster position on the flow cell and uniquely aligned to a single location in the human genome with the correct orientation and without any nucleotide mismatch were identified. In other embodiments, alignment may not be unique and mismatches may be allowed.

**[0047]** Only the PE reads that demonstrated an insert size  $\leq 600$  bp were retrieved for analysis. With these criteria, the size of the analyzed plasma DNA fragments in these experiments ranged from 36 bp to 600 bp. The size of each sequenced DNA fragment was inferred from the coordinates of the outermost nucleotides at each end of the sequenced fragments.

Table 1 shows data for samples of various aneuploidy status. The data includes the number of cases, gestational age median and range, along with number of paired-end reads median and range and the fetal DNA fraction.

Case type	No. of cases	Gestational age (weeks) median (range)	No. of PE reads (millions) median (range)	Fetal DNA fraction (%) median (range)
Euploid	39	13.2 (11.3 - 15.1)	4.7 (1.8 - 12.0)	15.7 (5.9 - 25.7)
T21	18	13.0 (12.1 - 17.9)	5.2 (2.5 - 8.9)	13.8 (7.4 - 27.2)
T18	10	13.3 (12.1 - 14.2)	4.9 (3.6 - 6.2)	7.2 (4.8 - 16.7)
T13	13	12.4 (11.5 - 16.4)	5.3 (2.7 - 7.7)	7.5 (3.2 - 14.1)
All	80	13.1 (11.3 - 17.9)	4.9 (1.8 - 12.0)	13.7 (3.2 - 27.2)

**[0048]** The fractional concentrations of fetal DNA in the maternal plasma samples were deduced from the amount of sequences aligning to chromosome Y as previously described (Chiu RW et al. BMJ 2011;342:c7401). This technique is an example of a calibration method. Thus, the measured fetal DNA fraction in Table 1 can be used in calibration data points to estimate a fetal DNA fraction in a new sample. The samples used to collect the data in Table 1 may be considered calibration samples.

### ***B. Second Fetal Sample Set Using Targeted Sequencing***

**[0049]** Table 2 shows sample information and targeted enrichment of maternal plasma DNA according to embodiments of the present invention. Plasma samples were collected from 48 pregnant women, each carrying a single fetus. Among these 48 pregnant women, 21 were carrying a euploid fetus, 17 were carrying a trisomy 21 (T21) fetus, 9 were carrying a trisomy 18 (T18) fetus, and 1 was carrying a trisomy 13 (T13) fetus. These data, along with examples below, illustrate that embodiments can use targeted techniques. The size distribution of plasma DNA can be determined using paired-end massively parallel sequencing. In other embodiments, the size distribution of plasma DNA can be determined for example by not limited to using real-time PCR, electrophoresis and mass spectrometry analysis.

**[0050]** To obtain high-fold sequencing coverage of the target regions, the Agilent SureSelect Target Enrichment System was employed in one embodiment to design probes to capture DNA molecules from chr7 (0.9 Mb region), chr13 (1.1 Mb region), chr18 (1.2 Mb region) and chr21 (1.3 Mb region). In the probe design, exons on chr7, chr13, chr18, and the Down syndrome critical region on chr21 (21q22.1-q22.3) were first selected as target regions. Because chr13, chr18 and chr21 have less exonic regions than chr7, additional non-exonic regions on chr13, chr18, and the Down syndrome critical region on chr21 were introduced to balance the total length of the targeted region among the above four chromosomes. The selected non-exonic regions were 120 bp in length, uniquely mappable, with GC content close to 0.5 and evenly distributed over the targeted chromosomes

**[0051]** Coordinates of all of the above exonic and non-exonic regions were submitted to the Agilent eArray platform for probe design. 500 ng of each maternal plasma DNA library was incubated with the capture probes for 24 h at 65°C. After hybridization, the targeted DNA molecules were eluted and amplified by a 12-cycle PCR according to manufacturer's instructions. Libraries with target enrichment were indexed and sequenced on a GA IIx (Illumina) using the 50-bp × 2 PE format. An additional 7 cycles of sequencing were performed to decode the index sequence on each sequenced plasma DNA molecule. The 50-bp sequence reads were aligned to the non-repeat-masked human reference genome (Hg18) (genome.ucsc.edu), using the Short Oligonucleotide Alignment Program 2 (SOAP2) (soap.genomics.org.cn). PE reads with individual members were sequenced on the same cluster position on the flow cell and uniquely aligned to a single location in the human genome with the correct orientation. Two mismatches were allowed; complexity of the sequencing library was significantly reduced after target enrichment.

**[0052]** Only the PE reads that demonstrated an insert size  $\leq 600$  bp were retrieved for analysis. With these criteria, the size of the analyzed plasma DNA fragments in the current study ranged from 36 bp to 600 bp. The size of each sequenced DNA fragment was inferred from the coordinates of the outermost nucleotides at each end of the sequenced fragments. The fractional concentrations of fetal DNA in the maternal plasma samples were estimated from the ratio of fragments carrying the fetal-specific alleles and the alleles shared with the respective mothers.

Table 2 shows data from targeted sequencing for samples of various aneuploidy status.

---

Case type	No. of cases	Gestational age (weeks) median (range)	No. of PE reads (millions) median (range)	Fetal DNA fraction (%) median (range)
Euploid	21	13.0 (12.0 -13.3)	2.2 (1.7-3.0)	13.5 (8.4-22.0)
T21	17	13.6 (12.6 - 20.9)	2.1 (1.5-2.7)	15.4 (8.7-22.7)
T18	9	12.7 (11.9 -13.7)	1.9 (1.7-3.1)	10.5 (7.2-16.3)
T13	1	13	1.6	9.2
All	48	13.1 (11.9 - 20.9)	2.1 (1.5-3.1)	13.4 (7.2-22.7)

### C. Electrophoresis For Fetal Sample

[0053] In addition to using massively parallel sequencing, the analysis of the size distribution of plasma DNA can be achieved by electrophoresis. Electrophoresis measures a time for a fragment to move through a medium. Particles of different sizes take different times to move through the medium. Thus, in one embodiment, microfluidic electrophoresis of sequencing library of maternal plasma DNA can be performed to determine the size distribution of the maternal plasma DNA.

[0054] FIG. 4 is a plot 400 showing a size distribution (electropherogram) of maternal plasma DNA obtained using electrophoresis according to embodiments of the present invention. The micro fluidic electrophoresis was performed using the Agilent 2100 Bioanalyzer. The electropherograms of the sequencing libraries of two samples are shown in plot 400. The X-axis represents the time duration the DNA taken to reach the sensor and corresponds to the size of the DNA fragments. The Y-axis represents the fluorescence units (FU) of the DNA fragments passing through the sensor at a particular time.

[0055] The time duration a DNA fragment takes to reach the sensor is positively correlated with the size of the DNA fragment. The Bioanalyzer can automatically convert the time duration to fragment size by comparing the running time of the test sample to those of a mixture of DNA fragments with known lengths (i.e., a DNA ladder). The DNA sequencing libraries were subsequently sequenced using massively parallel sequencing and the fraction of chromosome Y sequences were used to determine the fractional fetal DNA concentrations of these samples.

[0056] In plot 400, the solid line 410 represents the sample UK92797 which had a fractional fetal DNA concentration of 8.3% and the dashed line 420 represents the sample UK94884 which had a fractional fetal DNA concentration of 20.3%. Comparing with sample UK92797, sample UK94884 (the sample with the higher fractional fetal DNA) had a relatively higher amount of DNA at electrophoretic time interval from 63 seconds to 73 seconds (region A) which corresponds to DNA size from 200 bp to 267 bp and a relatively lower amount of DNA at electrophoretic time of 76s (region B), corresponding to a DNA size of ~292 bp

**[0057]** According to the manufacturer's protocol, DNA adaptors and primer sets which had a total size of 122 bp were introduced to the plasma DNA for sequencing library construction. Therefore, the region A corresponds to plasma DNA fragments approximately from 78 bp to 145 bp, and region B corresponds to plasma DNA fragments of approximately 170 bp. Such deduction can be adapted to different protocols for DNA library construction. For example, during the Illumina single-read sequencing library preparation, a total size of 92 bp from adapter/primer sets would be introduced, while this size would be 119 bp for the standard paired-end sequencing library preparation.

**[0058]** In another embodiment, the plasma DNA can be amplified by a whole genome amplification system known to those skilled in the art, e.g. the Rubicon Genomics PlasmaPlex WGA kit ([www.rubicongenomics.com/products](http://www.rubicongenomics.com/products)). The amplified products can then be analyzed by the Bioanalyzer. In yet other embodiments, the amplified products can be analyzed by an electrophoretic system from e.g. Caliper ([www.caliperls.com/products/labchip-systems](http://www.caliperls.com/products/labchip-systems)). In yet other embodiments, the size distribution of plasma DNA can be analyzed directly, without amplification, using for example, a nanopore-based sequencer (e.g. from Oxford Nanopore Technologies ([www.nanoporetech.com](http://www.nanoporetech.com))), or a Helico DNA sequencer ([www.helicosbio.com](http://www.helicosbio.com)).

#### IV. SIZE PARAMETERS

**[0059]** As mentioned above, various parameters can provide a statistical measure of a size profile of DNA fragments in the biological sample. A parameter can be defined using the sizes of all of the DNA fragments analyzed, or just a portion. In one embodiment, a parameter provides a relative abundance of short and long DNA fragments, where the short and long DNA may correspond to specific sizes or ranges of sizes.

**[0060]** To investigate if the overall size distribution of maternal plasma DNA can be used for reflecting the fractional fetal DNA concentration, we have used different parameters to quantify the relative abundance of short and long DNA, and determined the correlation between these parameters and fractional fetal DNA concentrations. The results of these investigations are provided in sections below. Parameters that we used, for illustration purposes, for reflecting the relative abundance of short DNA include:

1. i. Proportion of DNA fragments of 150 bp or below, which is labeled CF (size  $\leq$  150). CF refers to cumulative frequency. Thus, CF (size  $\leq$  150) refers to the cumulative frequency of fragments less than or equal to 150 bp;
2. ii. Ratio of the amounts of DNA fragments of  $\leq$ 150 bp and DNA from 163 bp to 169 bp, which is labeled  $(CF(\text{size} \leq 150) / \text{size}(163-169))$ ;
3. iii. Ratio of the amounts of DNA fragments from 140 bp to 146 bp and DNA from 163 bp to 169 bp, which is labeled  $(\text{size}(140-146) / \text{size}(163-169))$ ;
4. iv. Ratio of the amounts of DNA fragments from 140 bp to 154 bp and DNA from 163 bp to 169 bp, which is labeled  $(\text{size}(140-154) / \text{size}(163-169))$ ; and
5. v. Ratio of the amounts of DNA fragments from 100 bp to 150 bp and DNA from 163 bp

to 169 bp, which is labeled  $(\text{size}(100-150)/\text{size}(163-169))$ .

**[0061]** Other examples of parameters are the frequency counters of a histogram. In one embodiment, multiple parameters may be used. For example, the value of each parameter may give a difference percentage and then an average percentage can be determined. In another embodiment, each parameter corresponds to a different dimension of a multidimensional calibration function, where the values of the parameters for a new sample corresponds to a coordinate on the corresponding multidimensional surface.

## V. CORRELATION OF SIZE TO FRACTIONAL CONCENTRATION

**[0062]** The two samples sets using sequencing are used to illustrate the correlation of various size parameters to fractional concentration. An analysis of the size of repeat elements is also provided. The electrophoresis data also shows a correlation between size parameters and fractional concentration.

### A. First Sample Set

**[0063]** FIG. 5A is a plot 500 showing a proportion of DNA fragments that are 150 bp or below for samples having various fetal DNA percentage in maternal plasma according to embodiments of the present invention. The proportion of DNA  $\leq 150$  bp is plotted against the fractional fetal DNA concentration for the 80 maternal plasma samples. The euploid samples are represented by filled circles. The trisomy 13 (T13) samples are represented by unfilled triangles. The trisomy 18 (T18) samples are represented by unfilled rhombus and the trisomy 21 (T21) samples are represented by inverted unfilled triangles.

**[0064]** There is a positive relationship between the fractional fetal DNA concentration and the proportion of DNA  $\leq 150$  bp for all samples (Pearson correlation coefficient = 0.787). The positive correlation between the size parameter and the fractional fetal DNA concentration appears to be consistent across samples with different fetal chromosomal status. These results suggest that the analysis of the size parameter is useful for estimating the fractional fetal DNA concentration in a maternal plasma sample. Accordingly, the data points in FIG. 5 can be used as the calibration data points of method 300. Then, if the parameter  $CF(\text{size} \leq 150)$  is determined to be 30 for a new sample, the fetal DNA percentage can be estimated as being between about 7% and 16%. The data points in FIG. 5 can also be used to determine a calibration function that fits the raw data points shown.

**[0065]** FIG. 5B is a plot 550 showing a size ratio of the amounts of DNA fragments of  $\leq 150$  bp and DNA from 163 bp to 169 bp, which is labeled as  $(CF(\text{size} \leq 150)/\text{size}(163-169))$ . The  $CF(\text{size} \leq 150)/\text{size}(163-169)$  ratio is plotted against the fractional fetal DNA concentration for

the 80 maternal plasma samples. There is a positive relationship between the fractional fetal DNA concentration and the  $CF(\text{size} \leq 150)/\text{size}(163-169)$  ratio for all samples (Pearson correlation coefficient = 0.815). The positive correlation between the size parameter and the fractional fetal DNA concentration is consistent across samples with different fetal chromosomal ploidy status.

**[0066]** FIG. 6A is a plot 600 showing a size ratio of the amounts of DNA fragments from 140 bp to 146 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-146))/\text{size}(163-169)$ ). The  $\text{size}(140-146)/\text{size}(163-169)$  ratio is plotted against the fractional fetal DNA concentration for the 80 maternal plasma samples. There is a positive relationship between the fractional fetal DNA concentration and the  $\text{size}(140-146)/\text{size}(163-169)$  ratio for all samples (Pearson correlation coefficient = 0.808). The positive correlation between the size parameter and the fractional fetal DNA concentration is consistent across samples with different fetal chromosomal ploidy status.

**[0067]** FIG. 6B is a plot 650 showing a size ratio of the amounts of DNA fragments from 140 bp to 154 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-154))/\text{size}(163-169)$ ). The  $\text{size}(140-154)/\text{size}(163-169)$  ratio is plotted against the fractional fetal DNA concentration for the 80 maternal plasma samples. There is a positive relationship between the fractional fetal DNA concentration and the  $\text{size}(140-154)/\text{size}(163-169)$  ratio for all samples (Pearson correlation coefficient = 0.802). The positive correlation between the size parameter and the fractional fetal DNA concentration appears to be consistent across samples with different fetal chromosomal ploidy status.

**[0068]** FIG. 7 is a plot 700 showing a size ratio of the amounts of DNA fragments from 100 bp to 150 bp and DNA from 163 bp to 169 bp, which is labeled  $(\text{size}(100-150))/\text{size}(163-169)$ . The  $\text{size}(100-150)/\text{size}(163-169)$  ratio is plotted against the fractional fetal DNA concentration for the 80 maternal plasma samples. There is a positive relationship between the fractional fetal DNA concentration and the  $\text{size}(100-150)/\text{size}(163-169)$  ratio for all samples (Pearson correlation coefficient = 0.831). The positive correlation between the size parameter and the fractional fetal DNA concentration is consistent across samples with different fetal chromosomal ploidy status.

### ***B. Second Sample Set***

**[0069]** FIG. 8 is a plot 800 showing a proportion of DNA fragments of 150 bp or below for samples having various fetal DNA percentage in maternal plasma according to embodiments of the present invention. The proportion of  $\text{DNA} \leq 150$  bp is plotted against the fractional fetal DNA concentration for the 48 maternal plasma samples which were massively parallel paired-end sequenced after target enrichment. The euploid samples are represented by filled circles. The trisomy 13 (T13) samples are represented by unfilled triangles. The trisomy 18 (T18) samples are represented by unfilled rhombus and the trisomy 21 (T21) samples are represented by inverted unfilled triangles. There is a positive relationship between the

fractional fetal DNA concentration and the proportion of DNA  $\leq 150$  bp for all samples (Pearson correlation coefficient = 0.816). The positive correlation between the size parameter and the fractional fetal DNA concentration is consistent across samples with different fetal chromosomal status. These results suggest that the analysis of the size parameter is useful for estimating the fractional fetal DNA concentration in a maternal plasma sample.

**[0070]** FIG. 9A is a plot 900 showing a size ratio of the amounts of DNA fragments of  $\leq 150$  bp and DNA from 163 bp to 169 bp, which is labeled as  $(CF(\text{size} \leq 150)/\text{size}(163-169))$ . The  $CF(\text{size} \leq 150)/\text{size}(163-169)$  ratio is plotted against the fractional fetal DNA concentration for the 48 maternal plasma samples. There is a positive relationship between the fractional fetal DNA concentration and the  $CF(\text{size} \leq 150)/\text{size}(163-169)$  ratio for all samples (Pearson correlation coefficient = 0.776). The positive correlation between the size parameter and the fractional fetal DNA concentration is consistent across samples with different fetal chromosomal ploidy status.

**[0071]** FIG. 9B is a plot 950 showing a size ratio of the amounts of DNA fragments from 140 bp to 146 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-146)/\text{size}(163-169))$ . The  $\text{size}(140-146)/\text{size}(163-169)$  ratio is plotted against the fractional fetal DNA concentration for the 48 maternal plasma samples. There is a positive relationship between the fractional fetal DNA concentration and the  $\text{size}(140-146)/\text{size}(163-169)$  ratio for all samples (Pearson correlation coefficient = 0.790). The positive correlation between the size parameter and the fractional fetal DNA concentration is consistent across samples with different fetal chromosomal ploidy status.

**[0072]** FIG. 10A is a plot 1000 showing a size ratio of the amounts of DNA fragments from 140 bp to 154 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-154)/\text{size}(163-169))$ . The  $\text{size}(140-154)/\text{size}(163-169)$  ratio is plotted against the fractional fetal DNA concentration for the 48 maternal plasma samples. There is a positive relationship between the fractional fetal DNA concentration and the  $\text{size}(140-154)/\text{size}(163-169)$  ratio for all samples (Pearson correlation coefficient = 0.793). The positive correlation between the size parameter and the fractional fetal DNA concentration is consistent across samples with different fetal chromosomal ploidy status.

**[0073]** FIG. 10B is a plot 1005 showing a size ratio of the amounts of DNA fragments from 100 bp to 150 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(100-150)/\text{size}(163-169))$ . The  $\text{size}(100-150)/\text{size}(163-169)$  ratio is plotted against the fractional fetal DNA concentration for the 48 maternal plasma samples. There is a positive relationship between the fractional fetal DNA concentration and the  $\text{size}(100-150)/\text{size}(163-169)$  ratio for all samples (Pearson correlation coefficient = 0.798). The positive correlation between the size parameter and the fractional fetal DNA concentration is consistent across samples with different fetal chromosomal ploidy status

### ***C. Repeats***

**[0074]** Above, we have demonstrated that the size of all mappable DNA fragments in the maternal plasma is correlated with the fractional fetal DNA concentration. In this section, we investigate if the analysis of the size of the repeat elements in the genome can also be used for the estimation of fractional fetal DNA concentration in plasma. In the current example, we analyzed the size distribution of the DNA fragments mapping to the *Alu* repeats of the genome.

**[0075]** FIG. 11 is a plot showing a size ratio plotted vs. fetal DNA percentage for the size of repeat elements according to embodiments of the present invention. This example uses the ratio of the amounts of DNA fragments from 100 bp to 150 bp and DNA from 163 bp to 169 bp ( $\text{size}(100-150)/\text{size}(163-169)$ ) to reflect the alteration in the size distribution vs. fetal DNA percentage. There is a positive correlation between the size ratio and the fractional fetal DNA concentration (Pearson correlation coefficient = 0.829). This result suggests that the size analysis of the repeat elements can also be used to determine the fractional fetal DNA concentration in a maternal sample.

**[0076]** In addition to using massively parallel sequencing, other methods, e.g. PCR, real-time PCR and mass spectrometry analysis can also be used to determine the size distribution of the repeat elements (e.g., *Alu* repeats) in maternal plasma. In one embodiment, the DNA in a maternal plasma sample can be ligated to a linker. Then, PCR can be performed using one primer specific to the *Alu* sequences and the other primer specific to the linker. Following PCR, the PCR products could be analyzed for their sizes, e.g. by electrophoresis, mass spectrometry, or massively parallel sequencing. This would allow a readout of the sizes of sequences derived from the *Alu* repeats in maternal plasma. This strategy can be used for other target sequence or sequence family. Furthermore, the PCR can be followed by a nested PCR involving another *Alu*-specific primer, in combination with either the same linker-specific primer or a nested primer within the linker. Such nested PCR would have the advantage of increasing the specificity of the amplification towards the sequence of interest (in this case being the *Alu* sequences).

**[0077]** One advantage of using repeat elements is that they have a relatively high copy number and so they may be easier to analyze. For example, one may be able to use fewer cycles of amplification. Also, with a higher copy number, the analytical precision is potentially higher. A potential disadvantage is that certain classes of repeat elements may have copy numbers that vary from individual to individual.

#### ***D. Electrophoresis***

**[0078]** FIG. 12A is an electropherogram 1200 that may be used to determine a size ratio according to embodiments of the present invention. For all of the analyzed DNA libraries, there was a sharp peak at approximately 292 bp, followed by a secondary peak ranging from 300 bp to 400 bp. As the area under curve for a size range can represent the relative amount of DNA

fragments from that region, we used the ratio of the area of regions A (from 200 bp to 267 bp) and B (from 290 bp to 294 bp) to quantify the relative abundance of short and long DNA fragments. We first manually adjusted the baseline of fluorescence units (FU) to 0 and then generated the area for the selected region.

[0079] FIG. 12B is a plot 1250 showing a size ratio of the amounts of DNA fragments from 200 bp to 267 bp and DNA from 290 bp to 294 bp (i.e., the ratio of the areas of region A and B shown on the electropherogram) for samples having various fetal DNA percentages in maternal plasma according to embodiments of the present invention. There was one T13 case showing a low 292-bp peak with the FU value of 6.1, whereas all other cases showed a FU value  $\geq 20$  FUs. As the low FU value would make the area measurement imprecise, this case was ignored from the analysis. The ratio of the areas of region A and B is plotted against the fractional fetal DNA concentration for the all other 79 maternal plasma samples. There is a positive relationship between the fractional fetal DNA concentration and the area A and B ratio for these samples (Pearson correlation coefficient = 0.723).

## VI. DETERMINING CALIBRATION DATA POINTS

[0080] As mentioned above, the calibration data points may be defined in a variety of ways. Additionally, the calibration data points may be obtained in a variety of ways. For example, the calibration data points may simply be read from memory as a series of calibration values of a parameter along with the corresponding fractional concentration. Also, a calibration function can be read from memory (e.g., a linear or non-linear function with a predetermined functional form), where the function defines the calibration data points. In some embodiments, the calibration data points can be calculated from data measured from calibration samples.

### *A. Method*

[0081] FIG. 13 is a flowchart of a method 1300 for determining calibration data points from measurements made from calibration samples according to embodiments of the present invention. The calibration samples include the clinically-relevant DNA and other DNA.

[0082] At block 1310, a plurality of calibration samples are received. The calibration samples may be obtained as described herein. Each sample can be analyzed separately via separate experiments or via some identification means (e.g., tagging a DNA fragment with a bar code) to identify which sample a molecule was from. For example, a calibration sample may be received at a machine, e.g., a sequencing machine, which outputs measurement data (e.g., sequence reads) that can be used to determine sizes of the DNA fragments, or is received at an electrophoresis machine.

[0083] At block 1320, the fractional concentration of clinically-relevant DNA is measured in

each of the plurality of calibration samples. In various embodiments measuring a fetal DNA concentration, a paternally-inherited sequence or a fetal-specific epigenetic markers may be used. For example, a paternally-inherited allele would be absent from a genome of the pregnant female and can be detected in maternal plasma at a percentage that is proportional to the fractional fetal DNA concentration. Fetal-specific epigenetic markers can include DNA sequences that exhibit fetal or placental-specific DNA methylation patterns in maternal plasma.

**[0084]** At block 1330, amounts of DNA fragments from each calibration sample are measured for various sizes. The sizes may be measured as described herein. The sizes may be counted, plotted, used to create a histogram, or other sorting procedure to obtain data regarding a size profile of the calibration sample.

**[0085]** At block 1340, a calibration value is calculated for a parameter based on the amounts of DNA fragments at multiple sizes. A calibration value can be calculated for each calibration sample. In one embodiment, the same parameter is used for each calibration value. However, embodiments may use multiple parameters as described herein. For example, the cumulative fraction of DNA fragments less than 150 bases may be used as the parameter, and samples with different fractional concentration would likely have different calibration values. A calibration data point may be determined for each sample, where the calibration data point includes the calibration value and the measured fractional concentration for the sample. These calibration data points can be used in method 300, or can be used to determine the final calibration data points (e.g., as defined via a functional fit).

**[0086]** At block 1350, a function that approximates the calibration values across a plurality of fractional concentrations is determined. For example, a linear function could be fit to the calibration values as a function of fractional concentration. The linear function can define the calibration data points to be used in method 300.

**[0087]** In some embodiments, calibration values for multiple parameters can be calculated for each sample. The calibration values for a sample can define a multidimensional coordinate (where each dimension is for each parameter) that along with the fractional concentration can provide a data point. Thus, in one implementation, a multidimensional function can be fit to all of the multidimensional data points. Accordingly, a multidimensional calibration curve can be used, where the different values of the parameters can effectively be input to a single calibration function that outputs the fractional concentration. And, the single calibration function can result from a functional fit of all of the data points obtained from the calibration samples.

### ***B. Measuring Tumoral DNA concentration***

**[0088]** As mentioned, embodiments can also be applied to concentration of tumor DNA in a biological sample. An example involving determining the fractional concentration of tumoral DNA follows.

**[0089]** We collected the plasma samples from two patients suffering from hepatocellular carcinoma (HCC) before and after surgical resection of the tumors. The size analysis was performed using paired-end (PE) massively parallel sequencing. Sequencing libraries of maternal plasma DNA were constructed as previously described (Lo YM et al. *Sci Transl Med* 2010; 2:61ra91). All libraries were sequenced by a HiSeq 2000 (Illumina) using the 50-bp × 2 PE format. The 50-bp sequence reads were aligned to the non-repeat-masked human reference genome (Hg18) (<http://genome.ucsc.edu>), using the Short Oligonucleotide Alignment Program 2 (SOAP2) ([soap.genomics.org.cn](http://soap.genomics.org.cn)). The size of each sequenced fragments was inferred from the coordinates of the outermost nucleotides at each end of the aligned fragments.

**[0090]** We genotyped the DNA extracted from the blood cells and the tumor sample of the HCC patients using the Affymetrix SNP6.0 microarray system. For each case, the regions demonstrating loss of heterozygosity (LOH) in the tumor tissue were identified using the Affymetrix Genotyping Console v4.0 based on the intensities of the different alleles of the SNP loci. The fractional concentrations of tumor-derived DNA (F) were estimated from the difference in amounts of sequences carrying the deleted and non-deleted alleles at the LOH regions using the following formula:  $F = (A-B) / A \times 100\%$ , where A is the number of sequence reads carrying the non-deleted alleles at the heterozygous SNPs in the LOH regions, and B is the number of sequence reads carrying the deleted alleles for the heterozygous SNPs in the LOH regions. Table 3 shows the results.

Table 3 shows sequencing information and measured fractional concentration of tumor DNA in the plasma samples.

Case No.	Sampling time	No. of sequenced reads	Fractional concentration of tumor DNA in plasma (%)
1	before tumor resection	448M	51.60
	after tumor resection	486M	0.90
2	before tumor resection	479M	5.60
	after tumor resection	542M	0.90

**[0091]** In another embodiment, a locus that exhibits duplication can be used. For example, a tumor can exhibit a gain of one copy of one of the two homologous chromosomes such that an allele is duplicated. Then, one can determine a first amount A of sequence reads having a non-duplicated allele at the one or more heterozygous loci (e.g., SNPs) and a second amount B of sequence reads having a duplicated allele at the heterozygous loci. The fractional concentration F of clinically-relevant DNA can be calculated as a ratio of the first amount and the second amount using a ratio  $(B - A) / A$ .

**[0092]** In another embodiment, one or more homozygous loci may be used. For example, one can identify one or more loci where the patient is homozygous and where a single nucleotide mutation is present in the tumor tissue. Then, a first amount A of sequence reads having a wildtype allele at the one or more homozygous loci can be determined. And, a second amount B of sequence reads having a mutant allele at one or more homozygous loci can be determined. The fractional concentration F of clinically-relevant DNA can be calculated as a ratio of the first amount and the second amount using a ratio  $2B/(A+B)$ .

### ***C. Example of Functional Fit to Data Points***

**[0093]** An example of performing a functional fit to the parameter values determined from calibration samples is now described. Plasma samples from 80 pregnant women each carrying a singleton male fetus were analyzed. Among these 80 pregnant women, 39 were carrying euploid fetuses, 13 were carrying trisomy 13 (T13) fetuses, 10 were carrying trisomy 18 (T18 fetuses) and 18 were carrying trisomy 21 (T21) fetuses. The median gestational age of the pregnant women was 13 weeks and 1 day. DNA was extracted from the plasma samples and sequenced using the Illumina HiSeq2000 platform as described (Zheng YW et al. Clin Chem. 2012;58:549-58.) except that the sequencing was performed in an 8-plex format. For each DNA molecule, 50 nucleotides were sequenced from each of the two ends and aligned to a reference genome (hg18).

**[0094]** The size of each sequenced molecule was then deduced from the coordinates of the outermost nucleotides at both ends. For each sample, a median of 11.1 million fragments were sequenced and aligned uniquely to the reference genome. A ratio was calculated by dividing the proportion of DNA molecules with size 100 bp to 150 bp by the proportion of DNA molecules with size 163 bp to 169 bp and this ratio is termed the size ratio. As all the 80 pregnancies were carrying a male fetus, the proportion of sequence reads that were uniquely aligned to the chromosome Y was used to determine the fractional concentration of fetal DNA in each plasma DNA sample.

**[0095]** The samples were randomly divided into two sets, namely the training set and validation set. The relationship between the fractional fetal DNA concentration and the size ratio was established based the samples in the training set using linear regression. Then, the size ratio was used to deduce the fractional fetal DNA concentration for the validation set using the linear regression formula. The validation is discussed in the next section.

**[0096]** FIG. 14A is a plot 1400 of a size ratio against the fractional concentration of fetal DNA for the training set according to embodiments of the present invention. As mentioned above, the size ratio is calculated by dividing the proportion of DNA molecules with size 100 bp to 150 bp by the proportion of DNA molecules with size 163 bp to 169 bp. The size ratio is plotted against the fractional concentration of fetal DNA, as shown by data points 1405. The unfilled circles represent the euploid cases. The filled symbols represent the aneuploidy cases (square for T13, circle for T18 and triangle for T21). The linear regression line 1410 results from the

functional fit to the data points. The functional fit can be performed via any suitable techniques, e.g., least squares. The line 1410 can be used to estimate values of parameters measured for other samples, not in the training set. Each part of line 1410 can be considered a calibration data point.

## VII. COMPARISON TO CALIBRATION DATA POINTS

**[0097]** As mentioned above, the calibration data points can be used to determine the fractional concentration of the clinically relevant DNA. For example, the raw data points 1405 in FIG. 14A may be used to provide a range of fractional DNA concentration for a particular calibration value (labeled size ratio in FIG. 14A), where the range can be used to determine if the fractional concentration is above a threshold amount. Instead of a range, an average of the fractional concentrations at a particular size ratio can be used. For example, the fractional concentration corresponding to a measurement of 1.3 as the size ratio in a new sample can be determined as the average concentration calculated from the two data points at 1.3. In one embodiment, a functional fit (e.g., line 1410) may be used.

**[0098]** FIG. 14B is a plot 1450 of fractional concentrations deduced (estimated) from linear function 1410 of FIG. 14A against the fractional concentrations measured using fetal-specific sequences according to embodiments of the present invention. Using the regression equation (i.e., line 1410) determined based on the data of the training set, the size ratio determined for a validation sample was used to deduce the fractional concentration of fetal DNA for the samples of the validation set. The measured fractional concentrations correspond to the proportion of chromosome Y sequences in the plasma DNA sample (i.e., proportion of sequence reads aligning to the chromosome Y).

**[0099]** The line 1460 represents the perfect correlation between the two sets of values. The deviation of a data point 1455 indicates how accurate the estimate was, with points on line 1460 being perfectly accurate. As noted herein, the estimate does not have to be perfectly accurate, as the desired test may simply be to determine whether a sufficient percentage of clinically-relevant DNA is in the biological sample. The unfilled circles represent the euploid cases. The filled symbols represent the aneuploidy cases (square for T13, circle for T18 and triangle for T21). The median difference between the fractional fetal DNA concentration deduced from the size ratio and that measured from the proportion of chromosome Y sequences was 2.1%. The difference was less than 4.9% in 90% of the samples.

**[0100]** Samples with different ploidy status were used in both the calibration set and the validation set. As shown in FIG. 14A, the relationship between the size ratio and the fractional fetal DNA concentration were consistent across samples with different ploidy status. As a result, the fractional fetal DNA concentration can be deduced from the size ratio of the sample without a prior knowledge of the ploidy status of the sample as illustrated in FIG. 14B. One calibration curve was used for samples with different ploidy status and, hence, we do not need

to know the ploidy status of the sample before using embodiments to determine the fractional fetal DNA concentration.

## VIII. CANCER

[0101] As described herein, embodiments can be used to estimate the fractional concentration of tumor DNA in a biological sample. As with the fetal examples, calibration samples can be used to determine correlation data points, e.g., by fitting a function (e.g., a linear function) to data points showing a correlation between a value of a size parameter and a measured fractional concentration.

### *A. Correlation of Size to Tumoral DNA concentration*

[0102] FIG. 15A is a plot 1500 showing a proportion of DNA fragments of 150 bp or below for samples having various tumor DNA percentages in plasma of two HCC patients before and after tumor resection according to embodiments of the present invention. The proportion of DNA  $\leq 150$  bp is plotted against the fractional tumoral DNA concentrations for the two HCC patients before (filled circles) and after (unfilled circles) tumor resection. The two unfilled circles are very close in location to one another (effectively on top of each other). These results suggest that the analysis of the size parameter is useful for estimating the fractional tumoral DNA concentration in the plasma sample of HCC patients. There is a reduction in both the fractional tumor DNA concentration and the proportion of DNA fragments of  $\leq 150$  bp after tumor resection. The filled circle 1505 corresponds to a sample with much lower tumor DNA percentage, which is related to a smaller size of the tumor. In other words, the patient with a larger tumor has a higher proportion of short DNA which is reflected by a higher ratio of CF( $\leq 150$  bp) compared with the patient with a smaller tumor.

[0103] FIG. 15B is a plot 1550 showing a size ratio of the amounts of DNA fragments of  $\leq 150$  bp and DNA from 163 bp to 169 bp, which is labeled as  $(CF(\text{size} \leq 150) / \text{size}(163-169))$ , for two HCC patients before and after tumor resection. The  $CF(\text{size} \leq 150) / \text{size}(163-169)$  ratio is plotted against the fractional tumoral DNA concentrations for the two HCC patients before (filled circles) and after (unfilled circles) tumor resection. The two unfilled circles are very close in location to one another. There is a reduction in both the fractional tumor DNA concentration and the size ratio after tumor resection.

[0104] FIG. 16A is a plot 1600 showing a size ratio of the amounts of DNA fragments from 140 bp to 146 bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-146) / \text{size}(163-169))$ , for two HCC patients before and after tumor resection. The  $\text{size}(140-146) / \text{size}(163-169)$  ratio is plotted against the fractional tumoral DNA concentrations for the two HCC patients before (filled circles) and after (unfilled circles) tumor resection. There is a reduction in both the fractional tumor DNA concentration and the size ratio after tumor resection.

**[0105]** FIG. 16B is a plot 1650 showing a size ratio of the amounts of DNA fragments from 140 bp to 154bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-154)/\text{size}(163-169))$ , for two HCC patients before and after tumor resection. The  $\text{size}(140-154)/\text{size}(163-169)$  ratio is plotted against the fractional tumoral DNA concentrations for the two HCC patients before (filled circles) and after (unfilled circles) tumor resection. There is a reduction in both the fractional tumor DNA concentration and the size ratio after tumor resection.

**[0106]** FIG. 17 is a plot 1700 showing a size ratio of the amounts of DNA fragments from 100 bp to 150bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(100-150)/\text{size}(163-169))$ , for two HCC patients before and after tumor resection. The  $\text{size}(100-150)/\text{size}(163-169)$  ratio is plotted against the fractional tumoral DNA concentrations for the two HCC patients before (filled circles) and after (unfilled circles) tumor resection. There is a reduction in both the fractional tumor DNA concentration and the size ratio after tumor resection.

### ***B. Size Decrease due to Treatment***

**[0107]** FIG. 18A is a plot 1800 showing a proportion of DNA fragments of 150 bp or below for HCC patients before and after tumor resection. The pair of samples from the same cancer patient is depicted by identical symbols connected by a dashed line. There is a general decrease in the proportion of  $\text{DNA} \leq 150$  bp for the plasma DNA in cancer patients after tumor resection.

**[0108]** The separation in the values of the proportion for pre-treatment and post-treatment illustrate a correlation between the existence of a tumor and the value of the size parameter. The separation in the values for pre-treatment and post-treatment can be used to determine how successful the treatment was, e.g., by comparing the proportion to a threshold, where a proportion below the threshold can indicate success. In another example, a difference between the pre-treatment and post-treatment can be compared to a threshold.

**[0109]** The proportion (or any other value of a size parameter) can also be used to detect an occurrence of a tumor. For example, a baseline value for a size parameter can be determined. Then, at a later time, a value for the size parameter can be measured again. If the value of the size parameter shows a significant change, then the patient may be at a higher risk of having a tumor. If the value of the size parameter does not vary much among individuals, which FIG. 18A indicates that the proportion does not (i.e., since post-treatment values are the same), then the same baseline value can be used for other patients. Thus, a baseline value does not need to be taken for each patient.

**[0110]** FIG. 18B is a plot 1850 showing a size ratio of the amounts of DNA fragments of  $\leq 150$  bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{CF}(\text{size} \leq 150)/\text{size}(163-169))$ , for HCC patients before and after tumor resection. The pair of samples from the same cancer patient is depicted by identical symbols connected by a dashed line. There is a decrease in this size ratio

for the two cases after tumor resection.

[0111] FIG. 19A is a plot 1900 showing a size ratio of the amounts of DNA fragments from 140bp to 146bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-146)/\text{size}(163-169))$ , for HCC patients before and after tumor resection. The pair of samples from the same cancer patient is depicted by identical symbols connected by a dashed line. There is decrease in this size ratio for the two cases after tumor resection.

[0112] FIG. 19B is a plot 1950 showing a size ratio of the amounts of DNA fragments from 140bp to 154bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(140-154)/\text{size}(163-169))$ , for HCC patients before and after tumor resection. The pair of samples from the same cancer patient is depicted by identical symbols connected by a dashed line. There is a decrease in this size ratio for the two cases after tumor resection.

[0113] FIG. 20 is a plot 2000 showing a size ratio of the amounts of DNA fragments from 100 bp to 150bp and DNA from 163 bp to 169 bp, which is labeled as  $(\text{size}(100-150)/\text{size}(163-169))$ , for HCC patients before and after tumor resection. The pair of samples from the same cancer patient is depicted by identical symbols connected by a dashed line. There is a decrease in this size ratio for the two cases after tumor resection.

### **C. Method**

[0114] FIG. 21 is a flowchart illustrating a method 2100 for analyzing a biological sample of an organism to determine a classification of a level of cancer according to embodiments of the present invention. Method 2100 can analyze the biological sample of the organism (e.g., a human). The biological sample includes DNA originating from normal cells and potentially from cells associated with cancer. At least some of the DNA is cell-free in the biological sample. Aspects of methods 300 and 1300 can be used with embodiments of method 2100.

[0115] At block 2110, amounts of DNA fragments corresponding to various sizes are measured. For each size of a plurality of sizes, an amount of a plurality of DNA fragments from the biological sample corresponding to the size can be measured, as described for method 300. The plurality of DNA fragments may be chosen at random or preferentially selected from one or more predetermined regions of a genome. For example, targeted enrichment may be performed or selection of sequence reads that are from particular regions of the genome may be used, e.g., as described above.

[0116] At block 2120, a first value of a first parameter is calculated based on the amounts of DNA fragments at multiple sizes. In one aspect, the first parameter provides a statistical measure of a size profile (e.g., a histogram) of DNA fragments in the biological sample. The parameter may be referred to as a size parameter since it is determined from the sizes of the plurality of DNA fragments. Examples of parameter are provided herein. Multiple parameters may be used, as is also described herein.

[0117] At block 2130, the first value is compared to a reference value. Examples of a reference value include a normal value and a cutoff value that is a specified distance from a normal value (e.g., in units of standard deviation). The reference value may be determined from a different sample from the same organism (e.g., when the organism was known to be healthy). Thus, the reference value may correspond to a value of the first parameter determined from a sample when the organism is presumed to have no cancer. In one embodiment, the biological sample is obtained from the organism after treatment and the reference value corresponds to a value of the first parameter determined from a sample taken before treatment (e.g., illustrated above). The reference value may also be determined from samples of other healthy organisms.

[0118] At block 2140, a classification of a level of cancer in the organism is determined based on the comparison. In various embodiments, the classification may be numerical, textual, or any other indicator. The classification can provide a binary result of yes or no as to cancer, a probability or other score, which may be absolute or a relative value, e.g., relative to a previous classification of the organism at an earlier time. In one implementation, the classification is that the organism does not have cancer or that the level of cancer has decreased. In another implementation, the classification is that the organism does have cancer or that a level of cancer has increased.

[0119] As described herein, the level of cancer can include an existence of cancer, a stage of the cancer, or a size of a tumor. For example, whether the first value exceeds (e.g., greater than or less than, depending on how the first parameter is define) can be used to determine if cancer exists, or at least a likelihood (e.g., a percentage likelihood). The extent above the threshold can provide an increasing likelihood, which can lead to the use of multiple thresholds. Additionally, the extent above can correspond to a different level of cancer, e.g., more tumors or larger tumors. Thus, embodiments can diagnose, stage, prognosticate, or monitor progress of a level of cancer in the organism.

#### ***D. Determining size distribution for particular regions***

[0120] As with other embodiments, the first set of DNA fragments can correspond to one or more predetermined regions of a genome of the organism. Thus, the size analysis can also be performed for select regions, e.g., specific chromosomes, arms of chromosomes, or multiple regions (bins) of the same length, e.g., 1 Mb. For example, one can focus on regions that are commonly altered in a cancer type of interest. Table 2200 of FIG. 22 shows some common chromosomal aberrations seen in various types of cancers. The gain refers to an amplification of a chromosome with one or more additional copies within a particular segment and loss refers to deletions of one or both homologous chromosome within a particular segment.

[0121] In one embodiment, additional sets of DNA fragments can be identified from the biological sample. Each set of DNA fragments can correspond to different predetermined

regions, such as the regions specified in table 2200. Regions that are not associated with cancer could also be used, e.g., to determine a reference value. The amount of DNA fragments corresponding to various sizes can be determined and size value of a parameter can be determined for each additional set of DNA fragments, as described herein. Thus, a different size value can be determined for each genomic region, where there is a one-one correspondence between a set of DNA fragments and a genomic region.

**[0122]** Each of the size values can be compared to a respective reference value. Predetermined regions where the corresponding size value is statistically different than the respective reference value can be identified. When a reference value is a normal value, the determination of statistical difference can be made by comparing a size value to a cutoff (e.g., where the cutoff value is a specific number of standard deviations from the normal value, based on an assumed or measured statistical distribution). The respective reference values may be the same or different for different regions. For example, different regions may have different normal values for size.

**[0123]** In one embodiment, the number of regions statistically different than the reference value may be used to determine the classification. Thus, one can determine the number of identifying predetermined regions where the corresponding size value is statistically different than the respective reference value. The number can be compared to a threshold number of regions to determine the classification of the level of cancer in the organism. The threshold number can be determined based on a variance within normal samples and within cancer samples.

**[0124]** As highlighted in table 2200, different cancers are associated with different parts of the genome. Thus, which regions that statistically different can be used to determine one or more possible types of cancer when the possible types of cancer are associated with the identified regions. For example, if a size value for DNA fragments from chromosomal segment 7p is found to be significantly lower than a normal value (e.g., as determined by a cutoff value), then colorectal cancer can be identified as a likely cancer when the classification indicates that cancer exists. Note that the size value for chromosomal segment 7p may be used as a sole indicator to determine the classification, or multiple regions may be used. In one embodiment, only if an overall classification indicates cancer would the size value for chromosomal segment 7p be used to identify colorectal cancer as a likely cancer.

## **IX. COMPUTER SYSTEM**

**[0125]** Any of the computer systems mentioned herein may utilize any suitable number of subsystems. Examples of such subsystems are shown in FIG. 23 in computer apparatus 2300. In some embodiments, a computer system includes a single computer apparatus, where the subsystems can be the components of the computer apparatus. In other embodiments, a computer system can include multiple computer apparatuses, each being a subsystem, with internal components.

**[0126]** The subsystems shown in FIG. 23 are interconnected via a system bus 2375. Additional subsystems such as a printer 2374, keyboard 2378, fixed disk 2379, monitor 2376, which is coupled to display adapter 2382, and others are shown. Peripherals and input/output (I/O) devices, which couple to I/O controller 2371, can be connected to the computer system by any number of means known in the art, such as serial port 2377. For example, serial port 2377 or external interface 2381 (e.g. Ethernet, Wi-Fi, etc.) can be used to connect computer system 2300 to a wide area network such as the Internet, a mouse input device, or a scanner. The interconnection via system bus 2375 allows the central processor 2373 to communicate with each subsystem and to control the execution of instructions from system memory 2372 or the fixed disk 2379, as well as the exchange of information between subsystems. The system memory 2372 and/or the fixed disk 2379 may embody a computer readable medium. Any of the values mentioned herein can be output from one component to another component and can be output to the user.

**[0127]** A computer system can include a plurality of the same components or subsystems, e.g., connected together by external interface 2381 or by an internal interface. In some embodiments, computer systems, subsystem, or apparatuses can communicate over a network. In such instances, one computer can be considered a client and another computer a server, where each can be part of a same computer system. A client and a server can each include multiple systems, subsystems, or components.

**[0128]** It should be understood that any of the embodiments of the present invention can be implemented in the form of control logic using hardware (e.g. an application specific integrated circuit or field programmable gate array) and/or using computer software with a generally programmable processor in a modular or integrated manner. Based on the disclosure and teachings provided herein, a person of ordinary skill in the art will know and appreciate other ways and/or methods to implement embodiments of the present invention using hardware and a combination of hardware and software.

**[0129]** Any of the software components or functions described in this application may be implemented as software code to be executed by a processor using any suitable computer language such as, for example, Java, C++ or Perl using, for example, conventional or object-oriented techniques. The software code may be stored as a series of instructions or commands on a computer readable medium for storage and/or transmission, suitable media include random access memory (RAM), a read only memory (ROM), a magnetic medium such as a hard-drive or a floppy disk, or an optical medium such as a compact disk (CD) or DVD (digital versatile disk), flash memory, and the like. The computer readable medium may be any combination of such storage or transmission devices.

**[0130]** Such programs may also be encoded and transmitted using carrier signals adapted for transmission via wired, optical, and/or wireless networks conforming to a variety of protocols, including the Internet. As such, a computer readable medium according to an embodiment of the present invention may be created using a data signal encoded with such programs.

Computer readable media encoded with the program code may be packaged with a compatible device or provided separately from other devices (e.g., via Internet download). Any such computer readable medium may reside on or within a single computer program product (e.g. a hard drive, a CD, or an entire computer system), and may be present on or within different computer program products within a system or network. A computer system may include a monitor, printer, or other suitable display for providing any of the results mentioned herein to a user.

**[0131]** Any of the methods described herein may be totally or partially performed with a computer system including one or more processors, which can be configured to perform the steps. Thus, embodiments can be directed to computer systems configured to perform the steps of any of the methods described herein, potentially with different components performing a respective steps or a respective group of steps. Although presented as numbered steps, steps of methods herein can be performed at a same time or in a different order. Additionally, portions of these steps may be used with portions of other steps from other methods. Also, all or portions of a step may be optional. Additionally, any of the steps of any of the methods can be performed with modules, circuits, or other means for performing these steps.

**[0132]** The specific details of particular embodiments may be combined in any suitable manner without departing from the scope of the invention. However, other embodiments of the invention may be directed to specific embodiments relating to each individual aspect, or specific combinations of these individual aspects.

**[0133]** The above description of exemplary embodiments of the invention has been presented for the purposes of illustration and description. It is not intended to be exhaustive or to limit the invention to the precise form described, and many modifications and variations are possible in light of the teaching above. The embodiments were chosen and described in order to best explain the principles of the invention and its practical applications to thereby enable others skilled in the art to best utilize the invention in various embodiments and with various modifications as are suited to the particular use contemplated.

**[0134]** A recitation of "a", "an" or "the" is intended to mean "one or more" unless specifically indicated to the contrary.

## **REFERENCES CITED IN THE DESCRIPTION**

This list of references cited by the applicant is for the reader's convenience only. It does not form part of the European patent document. Even though great care has been taken in compiling the references, errors or omissions cannot be excluded and the EPO disclaims all liability in this regard.

### **Patent documents cited in the description**

- [US20110105353A](#) [0003]

#### Non-patent literature cited in the description

- LO YMD et al. *Lancet*, 1997, vol. 350, 485-487 [0001]
- LO YMD et al. *Am J Hum Genet*, 1998, vol. 62, 768-775 [0001] [0013]
- LUN FMF et al. *Clin Chem*, 2008, vol. 54, 1664-1672 [0001] [0013]
- LO YMD et al. *Proc Natl Acad Sci USA*, 2007, vol. 104, 13116-13121 [0001]
- FAN HCQUAKE SR. *PLoS One*, 2010, vol. 5, e10439- [0002]
- PALOMAKI GE et al. *Genet Med*, 2011, vol. 13, 913-920 [0002]
- SPARKS AB et al. *Am J Obstet Gynecol*, 2012, vol. 206, 319- [0002]
- FAN et al. Analysis of Size Distributions of Fetal and Maternal Cell-Free DNA by Paired-End Sequencing *Clinical Chemistry*, 2010, vol. 56, 81279-1286 [0002]
- LUN FMF et al. *Proc Natl Acad Sci USA*, 2008, vol. 105, 19920-19925 [0003]
- TSUI NBY et al. *Blood*, 2011, vol. 117, 3684-3691 [0003]
- LO YMD et al. *Sci Transl Med*, 2010, vol. 2, 61ra91- [0003] [0005] [0018]
- NYGREN AO et al. *Clin Chem*, 2010, vol. 56, 1627-1635 [0005]
- CHAN KCA et al. *Clin Chem*, 2008, vol. 52, 2211-2218 [0005]
- CHIM SSC et al. *Proc Natl Acad Sci USA*, 2005, vol. 102, 14753-14758 [0005]
- PAPAGEORGIUO EA et al. *Nat Med*, 2011, vol. 17, 510-513 [0005]
- MOULIERE et al. High fragmentation characterizes tumour-derived circulating DNAPLOS *One*, 2011, vol. 6, 9 [0006]
- CHAN KCA et al. *Clin Chem*, 2004, vol. 50, 88-92 [0018]
- DIEHL F et al. *Proc Natl Acad Sci USA*, 2005, vol. 102, 16368-16373 [0019]
- ZHENG YW et al. *Clin Chem*, 2012, vol. 58, 549-558 [0019] [0019]
- CHAN KC et al. *Clin Chem*, 2004, vol. 50, 88-92 [0022]
- LO YM et al. *Sci Transl Med*, 2010, vol. 2, 61ra91- [0022] [0045] [0089]
- CHIU RW et al. *BMJ*, 2011, vol. 342, c7401- [0048]
- ZHENG YW et al. *Clin Chem.*, 2012, vol. 58, 549-58 [0093]

PATENTKRAV

1. Fremgangsmåde til analyse af en biologisk prøve fra et individ, hvilken biologisk prøve indbefatter DNA, der stammer fra normale celler og eventuelt fra celler forbundet med cancer, hvor mindst noget af DNA'et er cellefrit i den biologiske prøve, hvilken fremgangsmåde omfatter:
- 5
- (a) sekvensering af en flerhed af DNA-fragmenter for at opnå en flerhed af sekvenslæsninger, der omfatter de yderste nukleotider ved hver ende af en flerhed af DNA-fragmenter;
  - 10 (b) aligning af sekvenslæsninger med et referencegenom, hvorved der opnås et sæt af genomiske koordinater, der indbefatter genomiske koordinater for de yderste nukleotider, der definerer en størrelse af et DNA-fragment for hvert af flerheden af DNA-fragmenter;
  - (c) beregning af en værdi af et parameter baseret på mængder af sekvenslæsninger af DNA-fragmenter, der alignes med sættet af genomiske koordinater fra (b) ved multiple størrelser;
  - 15 (d) sammenligning af værdien med en referenceværdi og
  - (e) bestemmelse af en klassificering af et niveau af cancer hos individet baseret på sammenligning af værdien med referenceværdien.
- 20       **2.** Fremgangsmåde ifølge krav 1, hvor den biologiske prøve er udvalgt fra gruppen bestående af blod, plasma, serum, urin og saliva.
- 3.** Fremgangsmåde ifølge krav 1 eller 2, hvor sættet af genomiske koordinater svarer til et område af referencegenomet.
- 4.** Fremgangsmåde ifølge krav 3, hvor området af referencegenomet er
- 25 forudbestemt.
- 5.** Fremgangsmåde ifølge et hvilket som helst af de foregående krav, hvor sekvenseringen omfatter parret-ende-sekvensering.

**6.** Fremgangsmåde ifølge et hvilket som helst af de foregående krav, hvor hver sekvenslæsning af sekvenslæsningerne omfatter mindst ca. 30 nukleotider fra mindst én ende af DNA-fragmentet.

5 **7.** Fremgangsmåde ifølge et hvilket som helst af de foregående krav, hvor værdien bestemmes ved anvendelse af et antal sekvenslæsninger med ender, der alignes med genomiske koordinater af sættet af genomiske koordinater for størrelse af DNA-fragmenterne.

10 **8.** Fremgangsmåde ifølge et hvilket som helst af krav 1 til 6, hvor værdien bestemmes ved anvendelse af et første antal sekvenslæsninger med ender, der alignes med genomiske koordinater af sætte af genomiske koordinater, der er normaliseret ved hjælp af et andet antal sekvenslæsninger med forskellige genomiske koordinater for størrelser af DNA-fragmenterne.

15 **9.** Fremgangsmåde ifølge et hvilket som helst af de foregående krav, hvor klassificeringen svarer til en størrelse eller eksistens af cancer, et stadie for canceren eller en størrelse på en tumor.

**10.** Fremgangsmåde ifølge et hvilket som helst af de foregående krav, hvor mindst en del af flerheden af DNA-fragmenter er afledt af en tumor.

20 **11.** Fremgangsmåde ifølge et hvilket som helst af de foregående krav, hvor den biologiske prøve opnås fra individet efter en behandling, og hvor referenceværdien svarer til en værdi bestemt ud fra en prøve taget fra individet før behandlingen.

**12.** Fremgangsmåde ifølge et hvilket som helst af krav 1 til 10, hvor referenceværdien svarer til en værdi bestemt ud fra en prøve, når individet ikke formodes at have cancer.

25 **13.** Fremgangsmåde ifølge et hvilket som helst af krav 1 til 10, hvor referenceværdien etableres ud fra én eller flere biologiske prøver opnået fra ét eller flere sunde individer.

# DRAWINGS

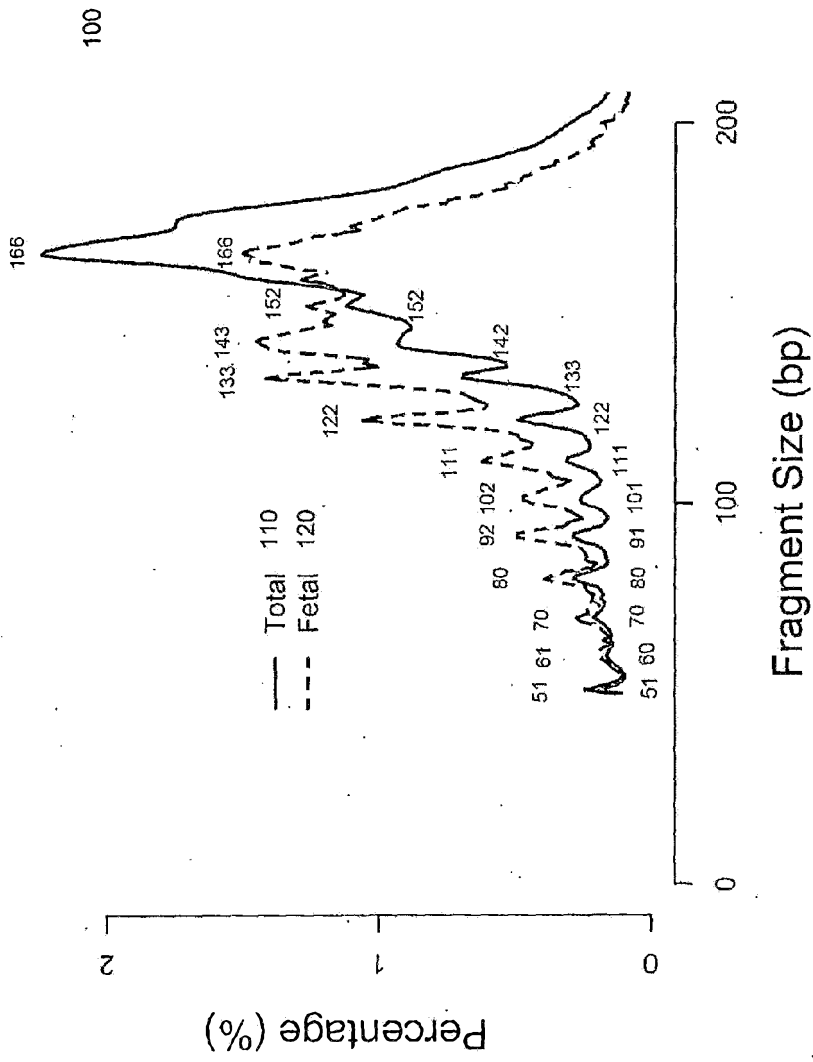


FIG. 1

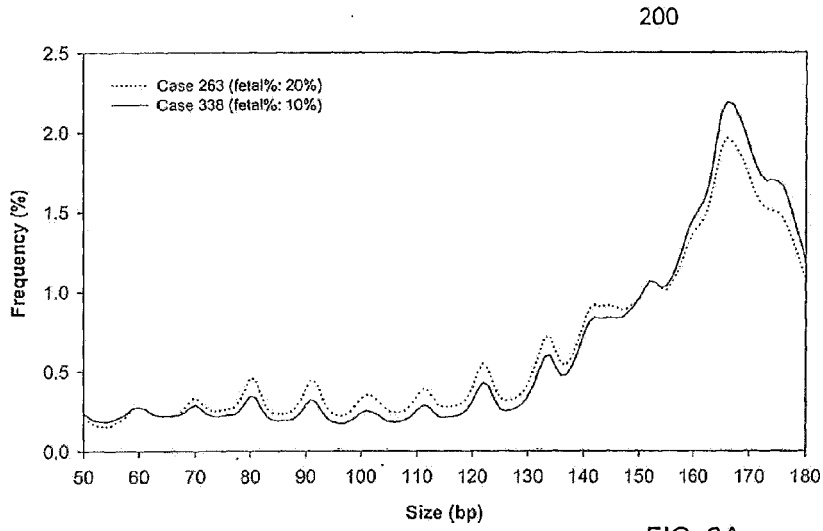


FIG. 2A

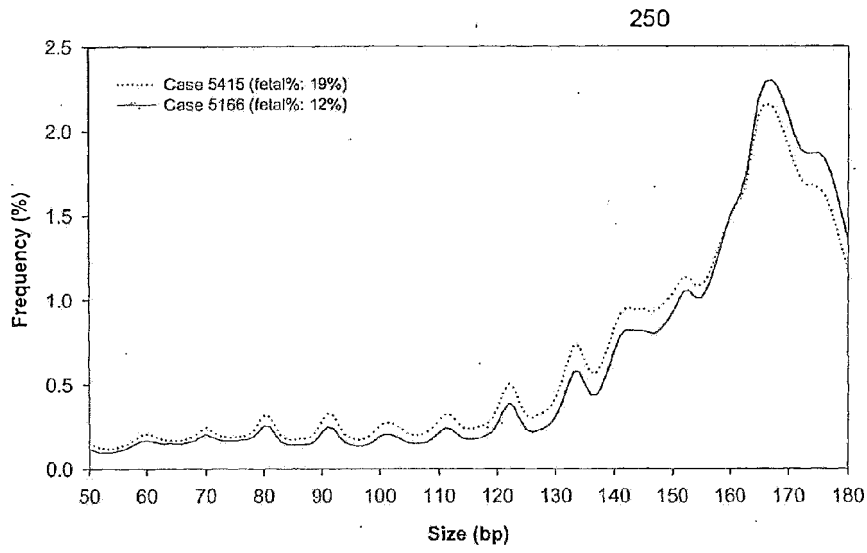


FIG. 2B

300

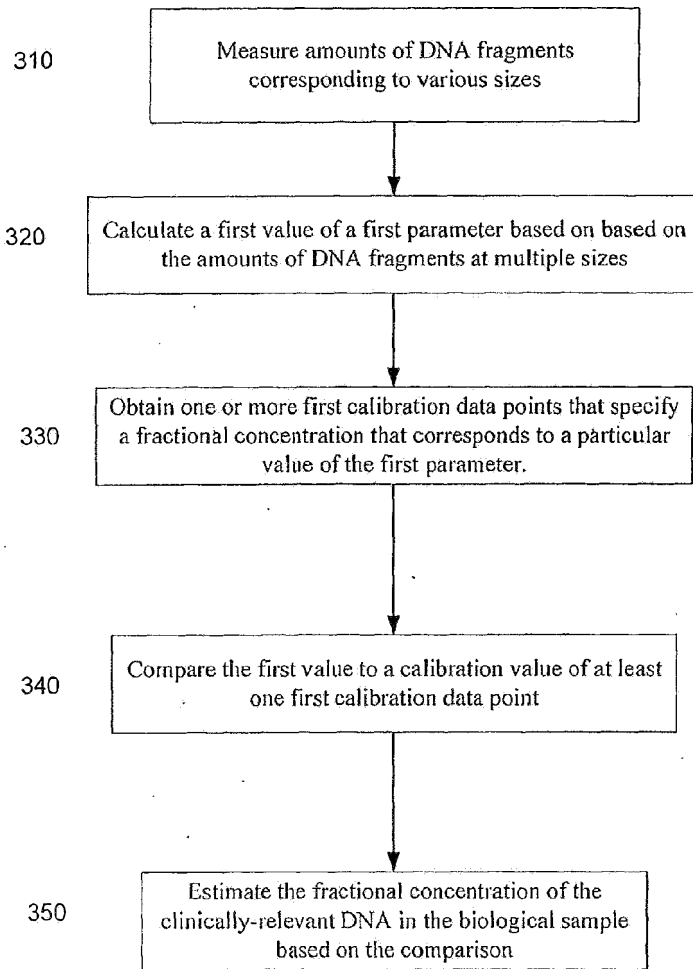


FIG. 3

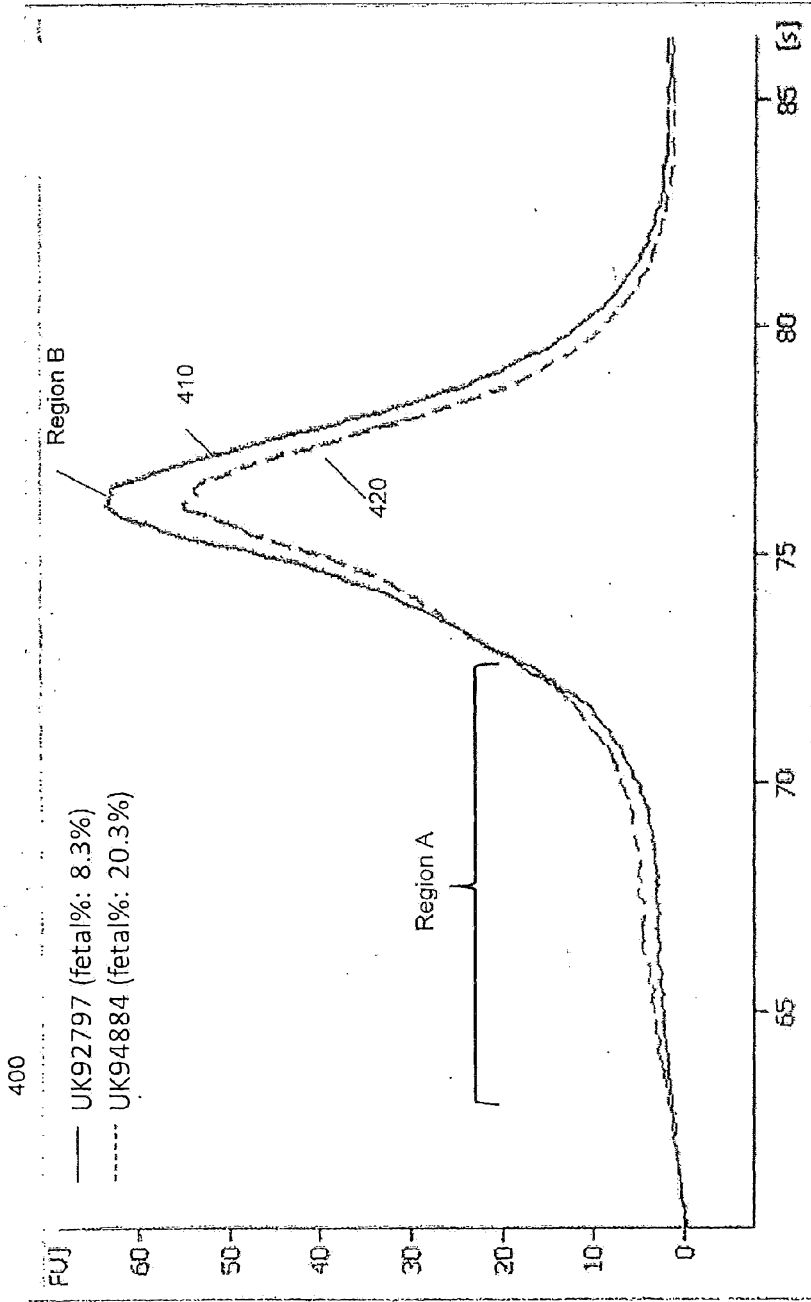
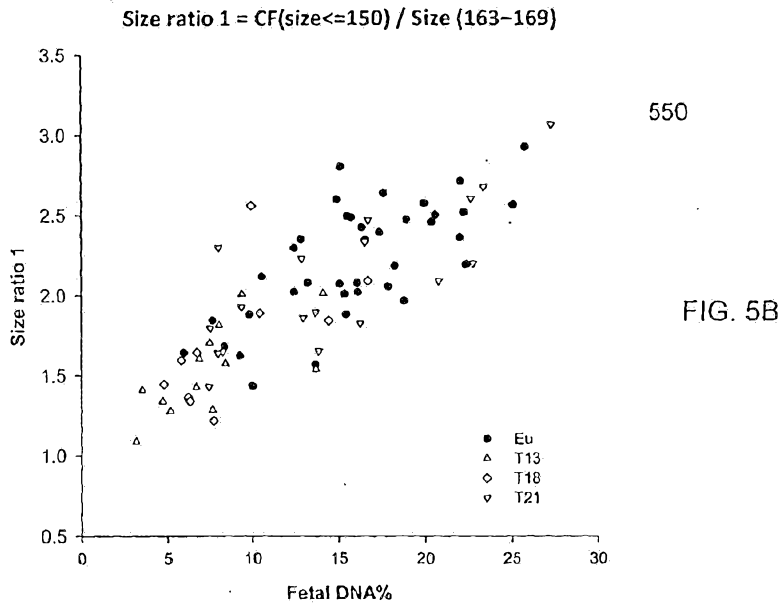
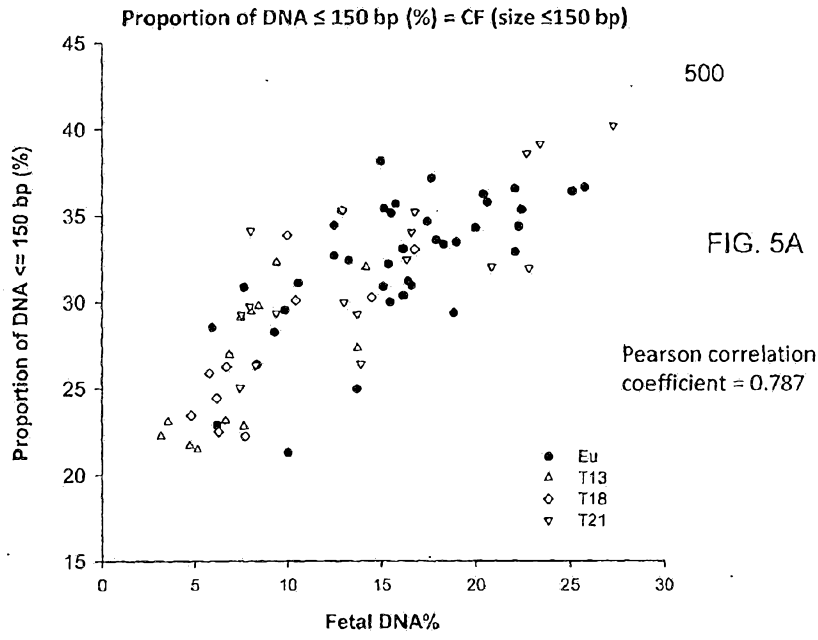
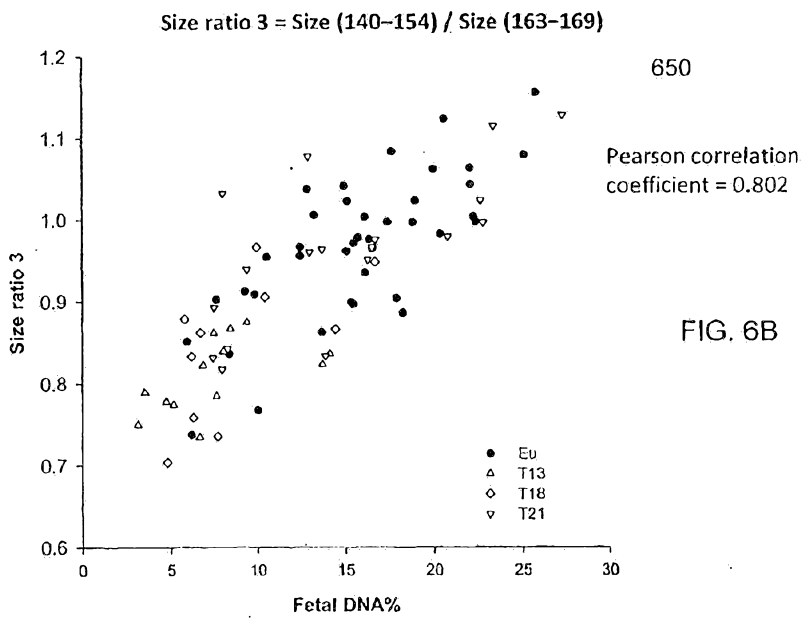
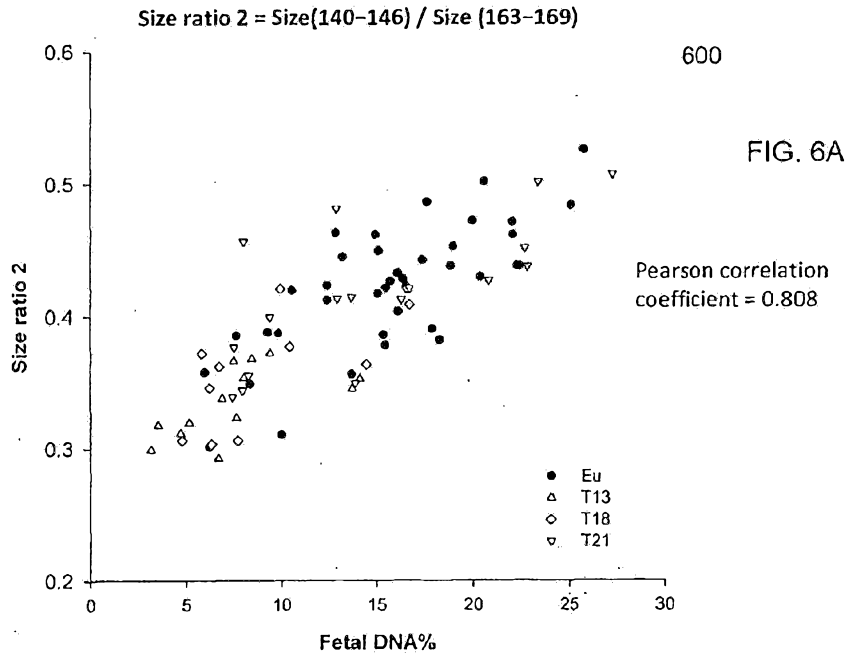


FIG. 4





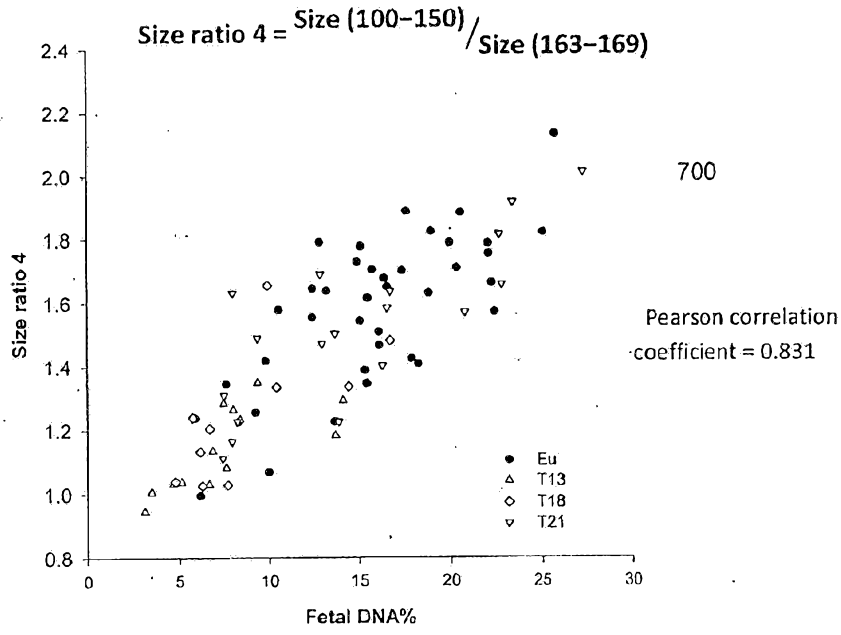


FIG. 7

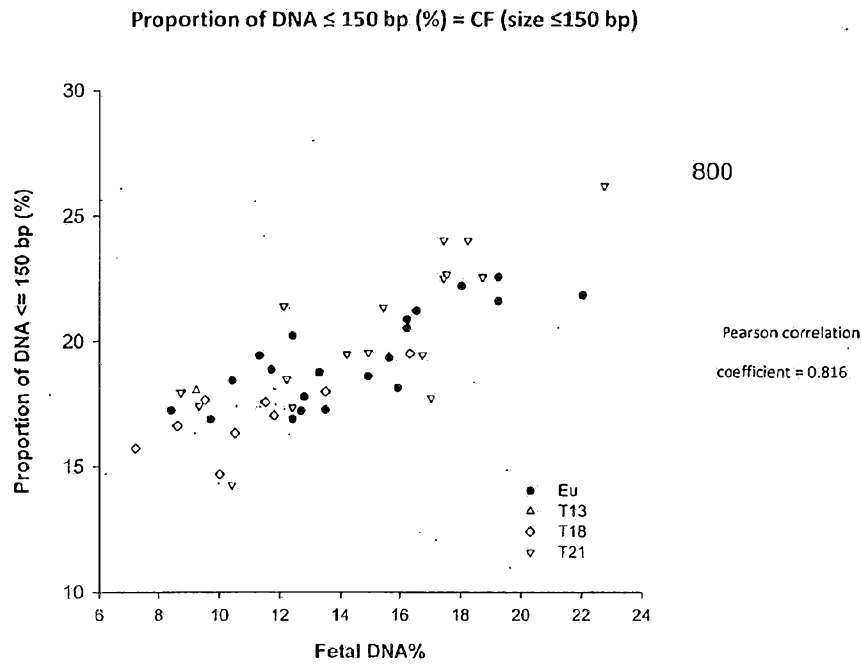
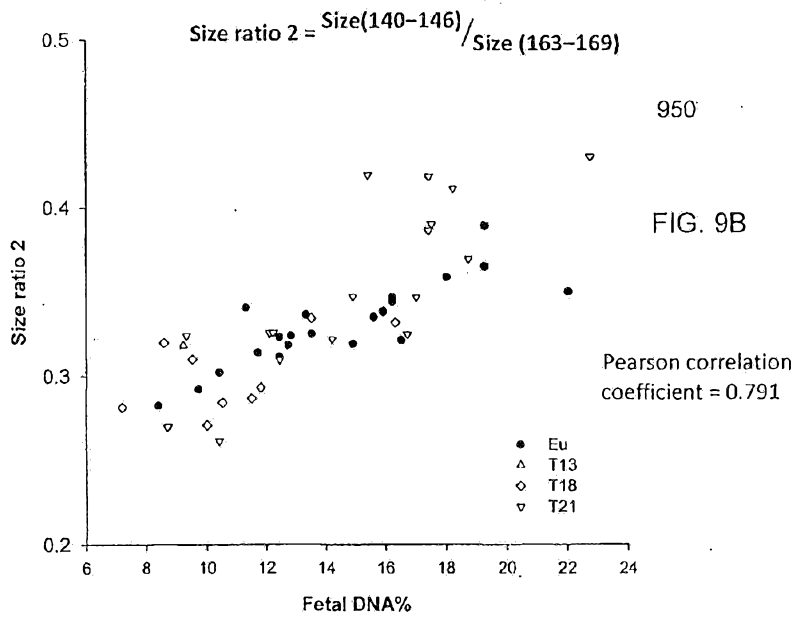
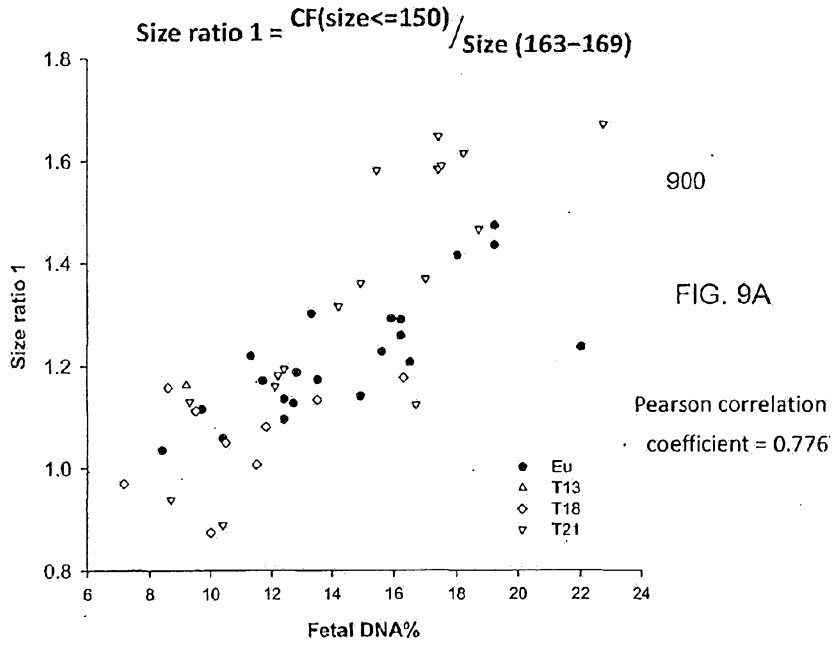
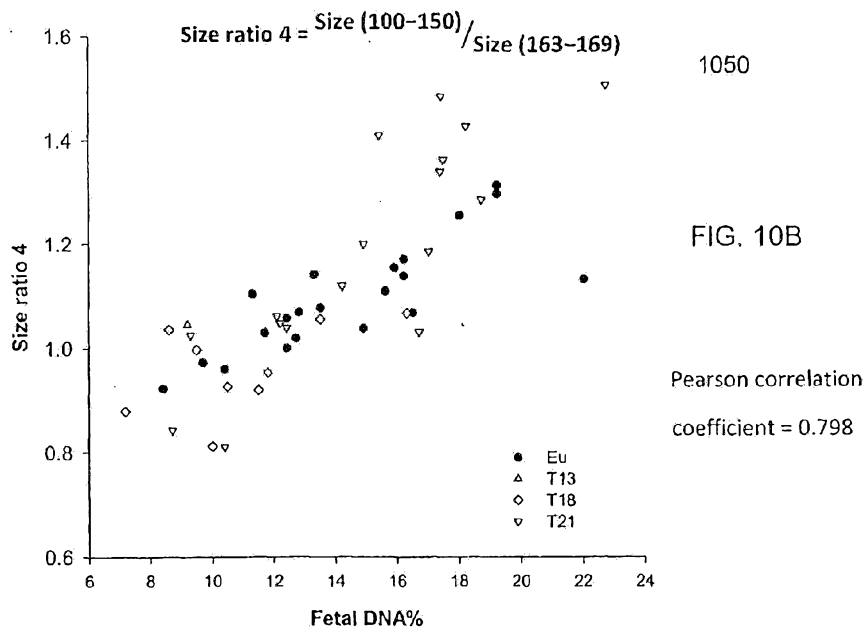
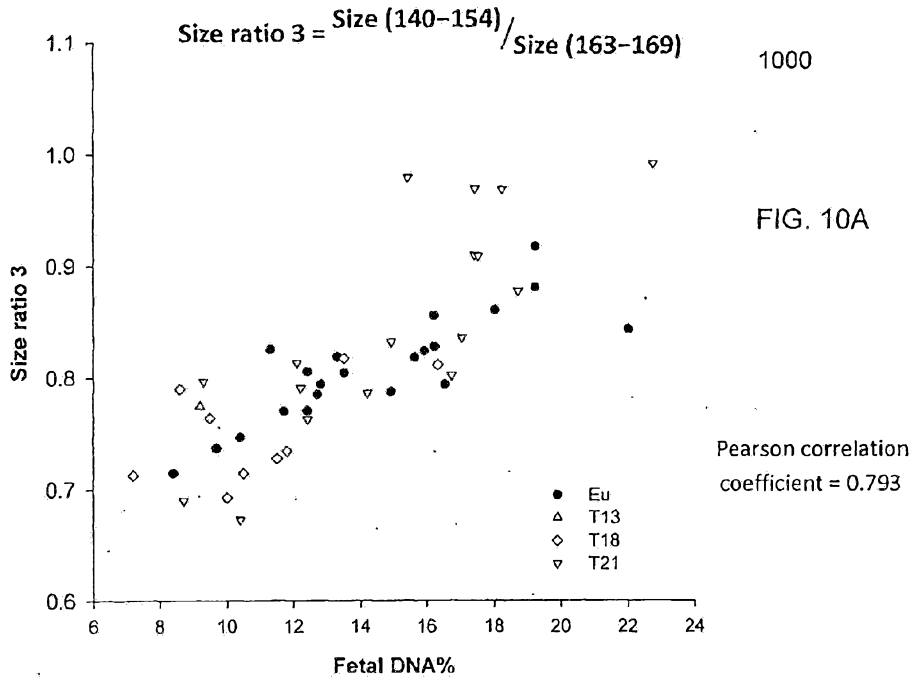


FIG. 8





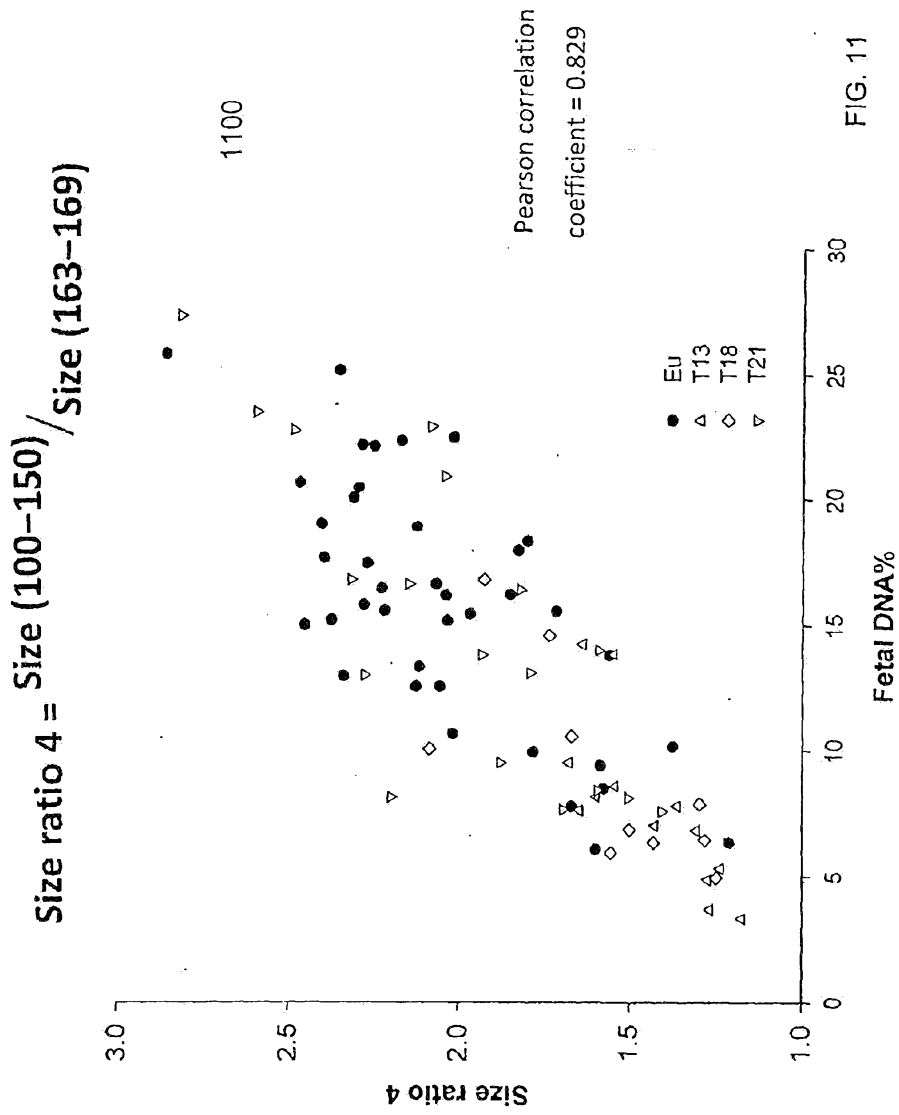


FIG. 11

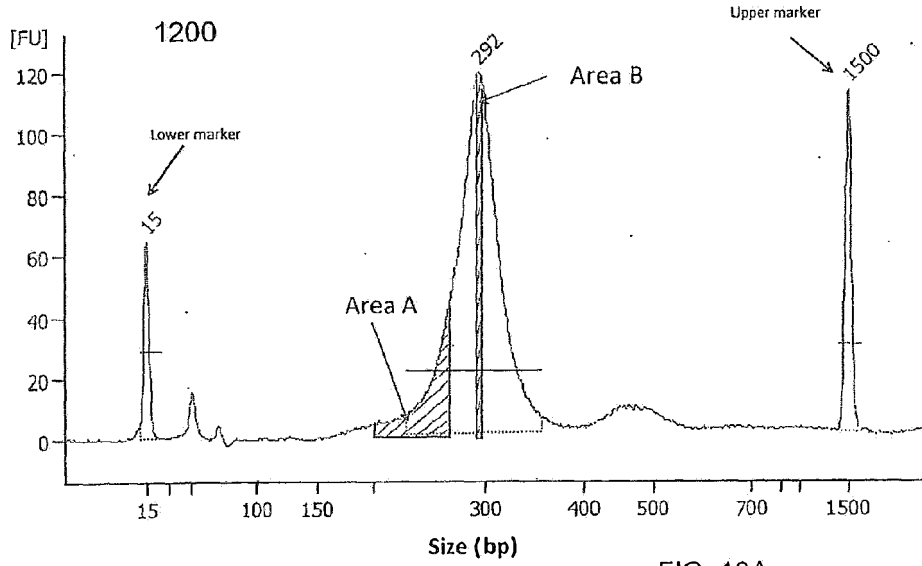


FIG. 12A

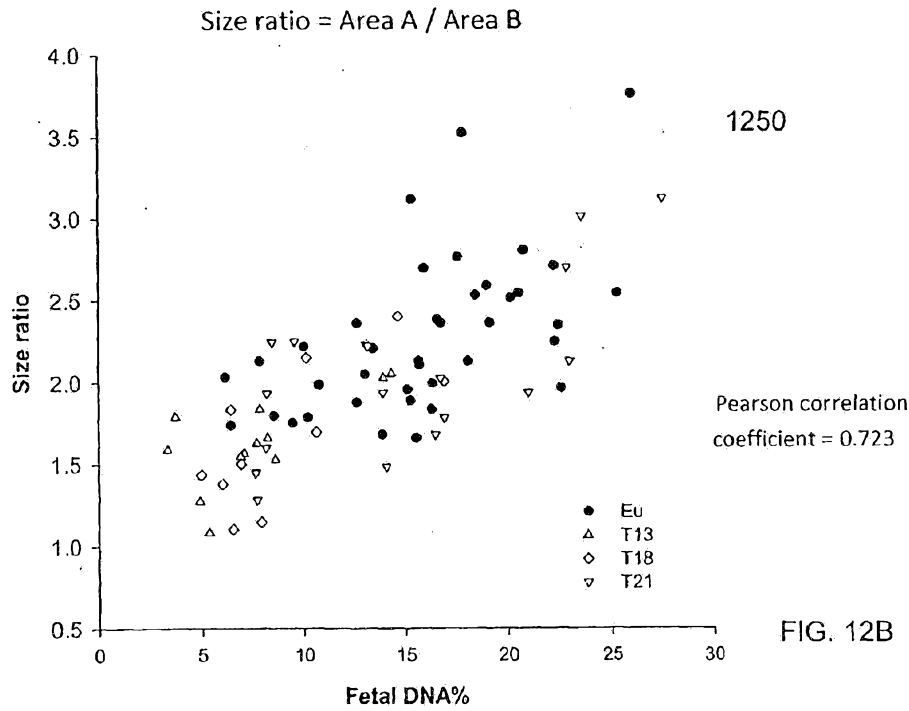


FIG. 12B

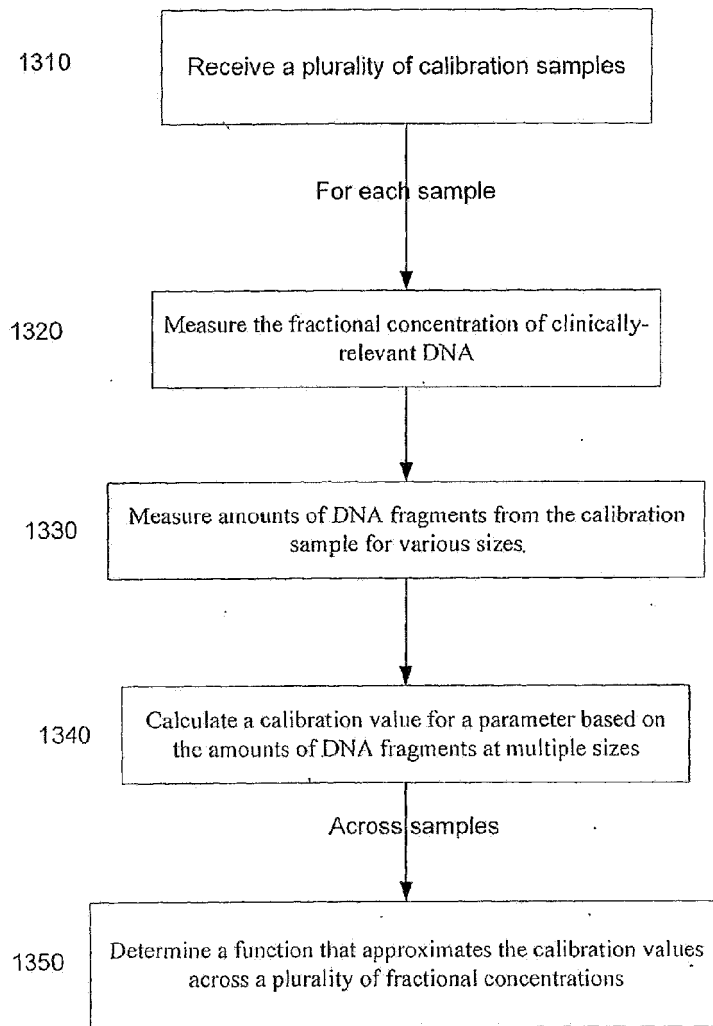


FIG. 13

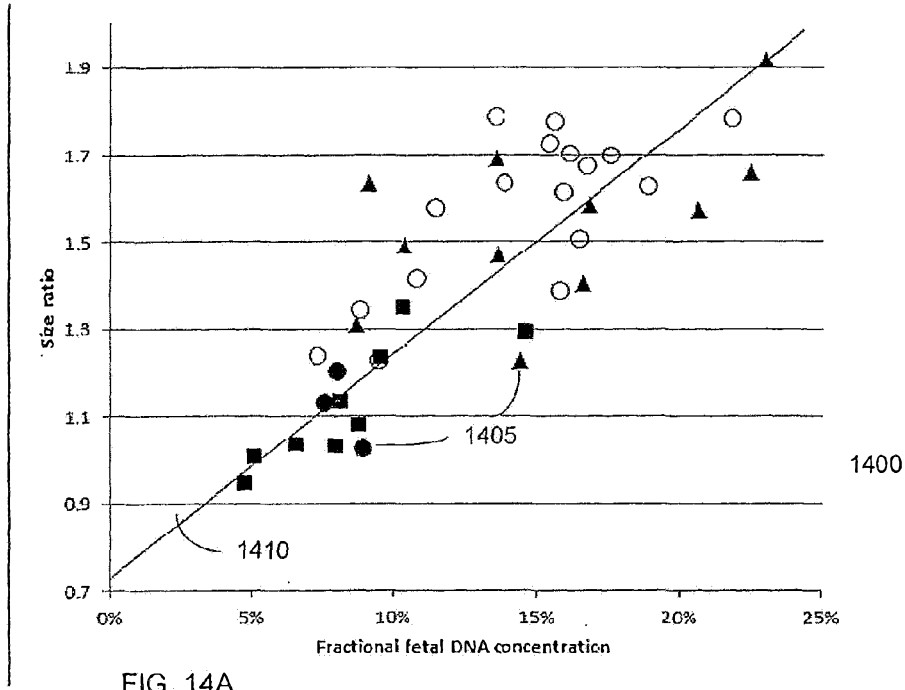


FIG. 14A

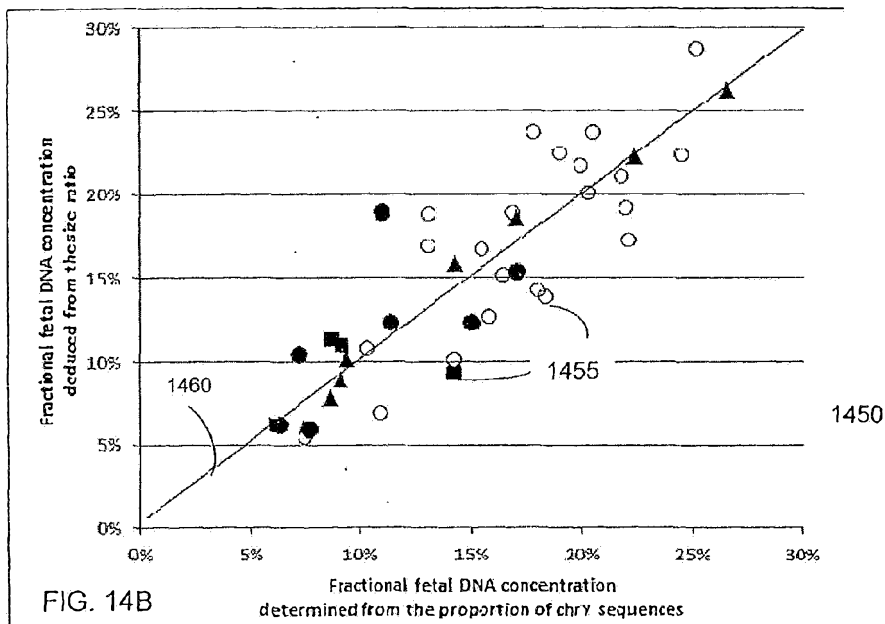


FIG. 14B

FIG. 15A

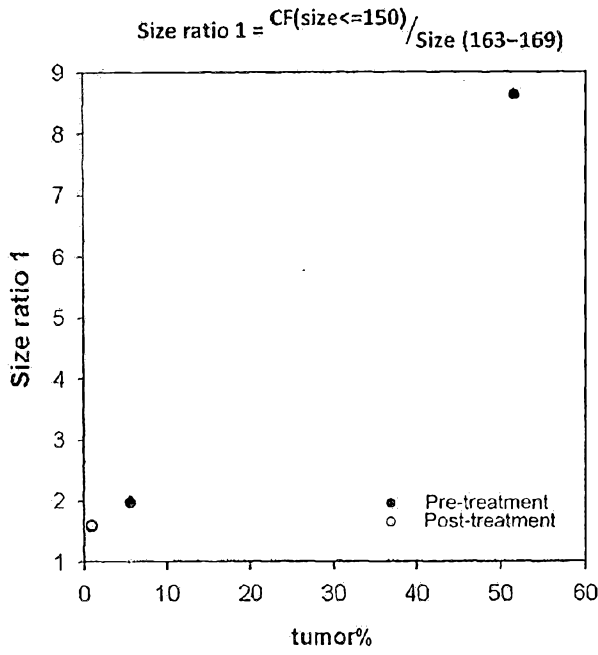
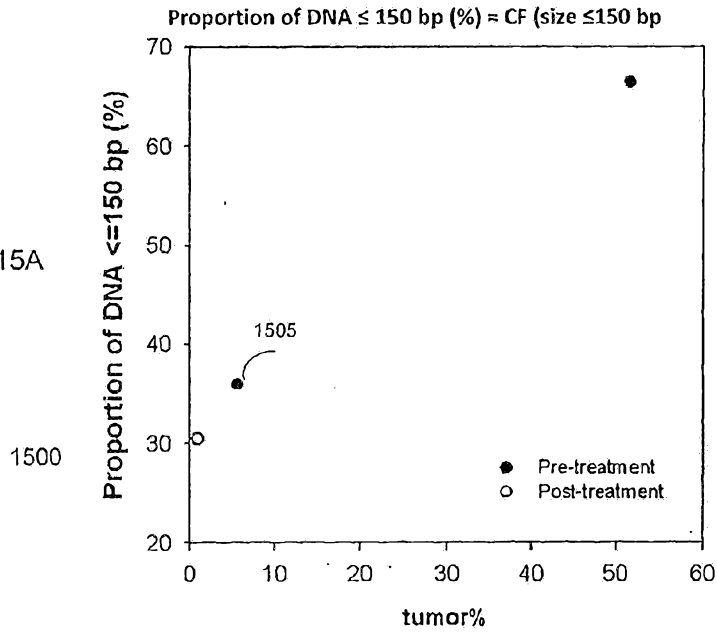


FIG. 15B

1550

FIG. 16A

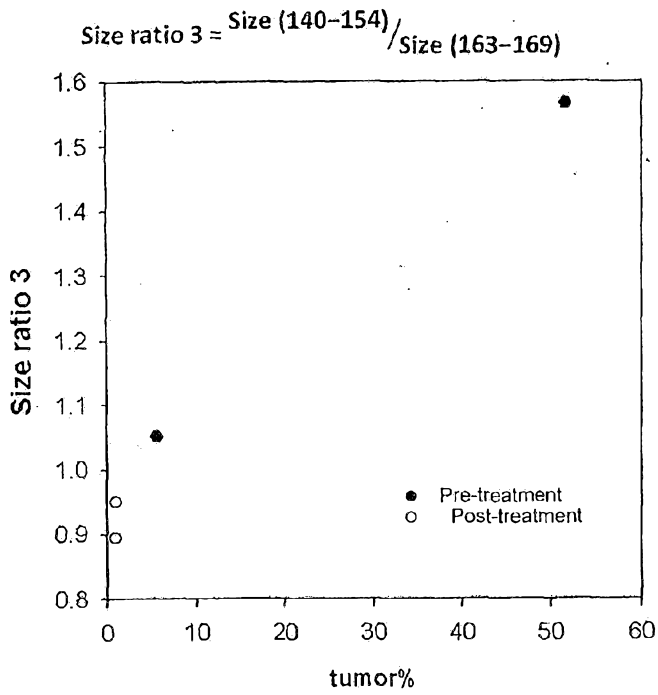
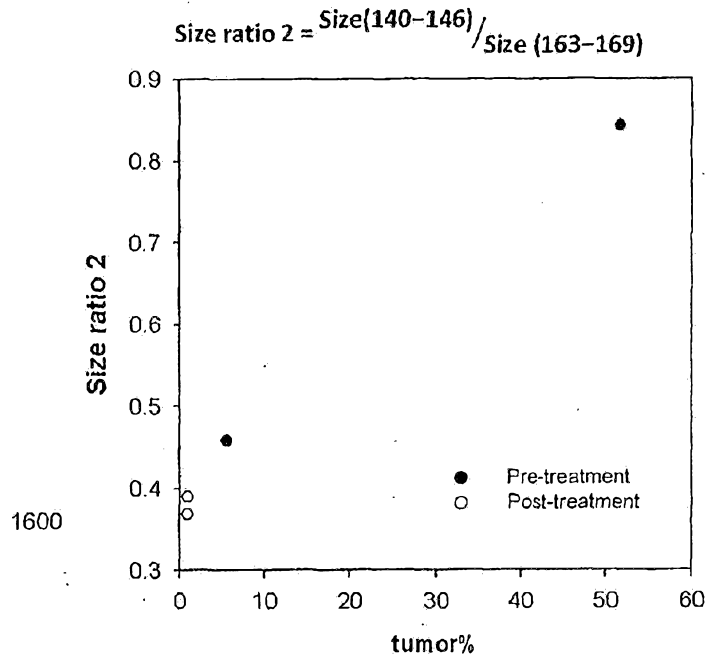


FIG. 16B

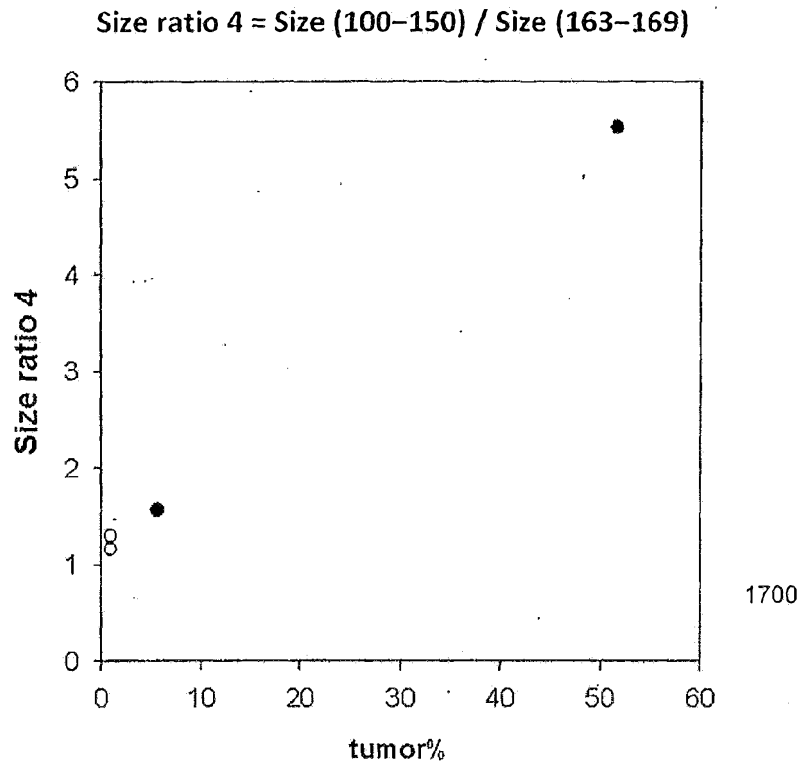
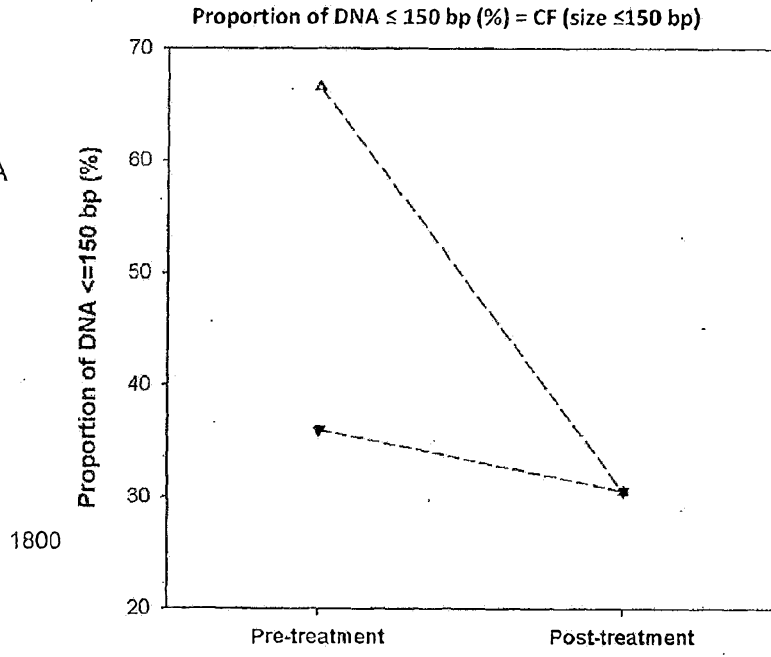


FIG. 17

FIG. 18A



Size ratio 1 = CF(size  $\leq 150$ ) / Size (163-169)

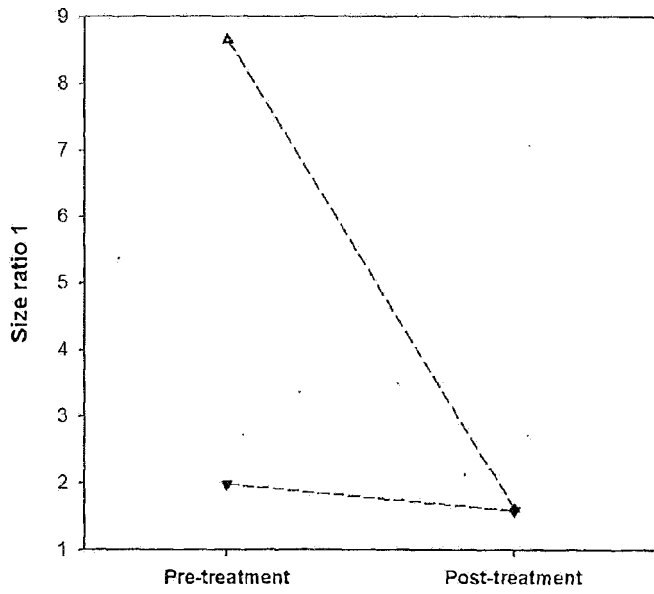
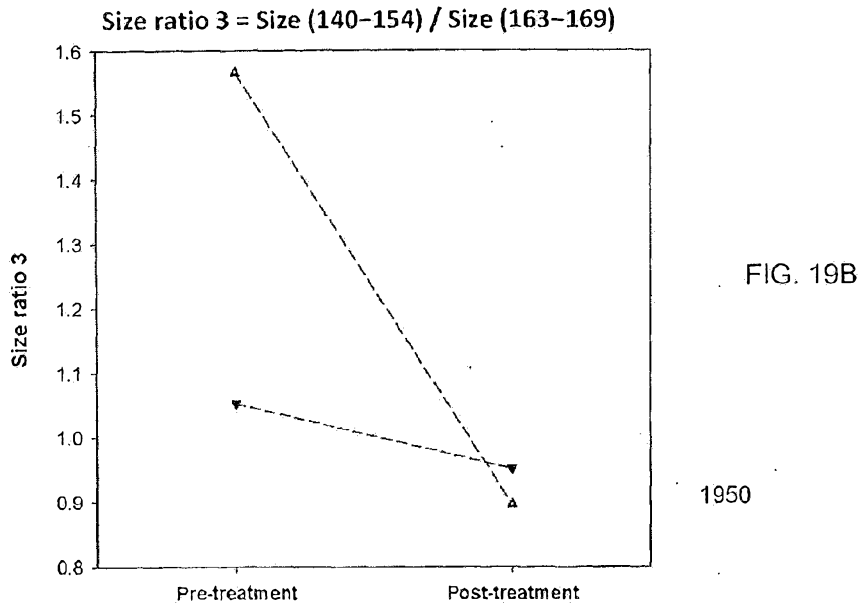
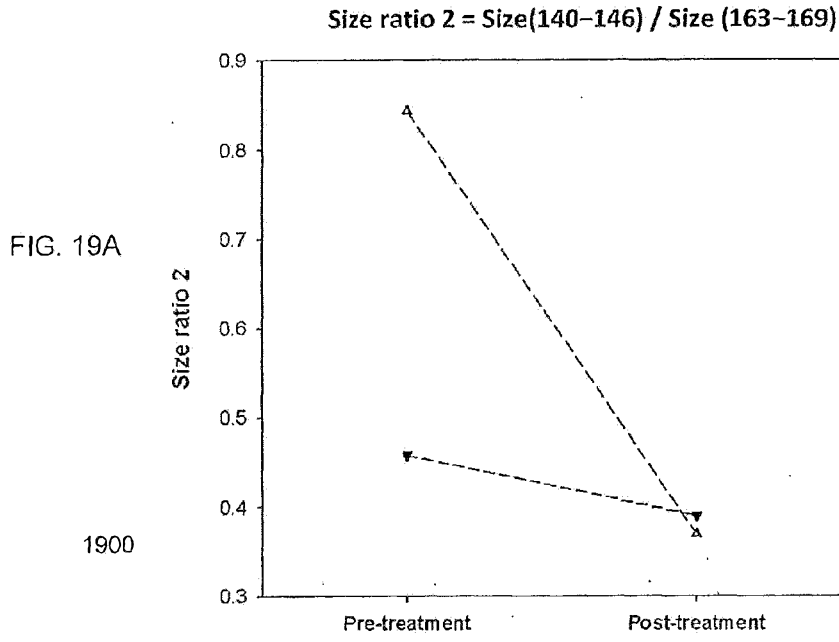


FIG. 18B

1850



$$\text{Size ratio 4} = \frac{\text{Size (100-150)}}{\text{Size (163-169)}}$$

2000

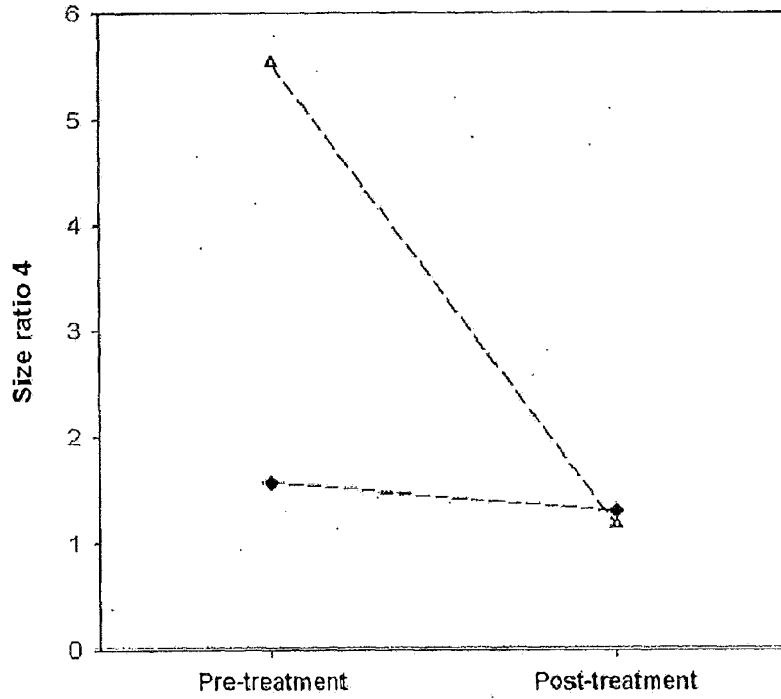


FIG. 20

2100

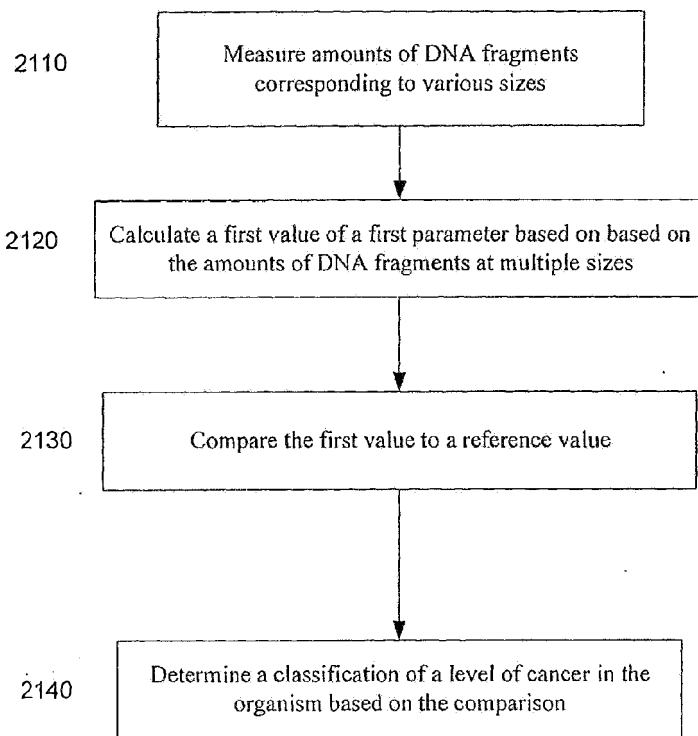


FIG. 21

	Gain	Loss	Reference
Colorectal	7p, 7q, 8q, 11q, 13q, and 20q,	5q, 8p, 17p, 18p, 18q and 20p,	(Nakao et al. Carcinogenesis 2004;25: 1345-1357.) (Tsafrir et al. Cancer Res 2006; 66: 2129-2137)
Breast	1q, 6p, 8q, 11q, 16p, 17q, 19, and 20q	6q, 13q, 16q, 17p, and 22q	(Tirkkonen et al. Gene Chromosome Canc 1998; 21: 177-184) (Richard et al. Int J Cancer 2000; 89: 305-310) (Pinkel et al. Nat Genet 1998; 20: 207-211) (Persson et al. Gene Chromosome Canc 1999; 25: 115-122) (Nishizaki et al. Int J Cancer 1997; 74: 513 - 517)
Lung	1q, 3q, 5p, and 8q	3p, 6q, 8p, 9p, 13q, and 17p	(Berrieman et al. Brit J Cancer 2004; 90: 900-905) (Luk et al. Cancer Genet Cytogen 2001; 125: 87 - 99) (Petersen et al. Cancer Res 1997; 57: 2331-2335) (Pei et al. Gene Chromosome Canc 2001; 31: 282-287)
HCC	1q, 8q, 17q and 20q	4q, 6q, 8p, 13q, 16q and 17p	(Kusano et al. Cancer 2002; 94: 746-751) (Laurent-Puig et al. Gastroenterology 2001; 120: 1763-1773) (Moinzadeh et al. Brit J Cancer 2005; 92: 935-941)
Ovarian	20q, 3q, 1q, 8q, 12p, 11q, and 17q	Xp, 18q, 4q, 9p, and 13q	(Tactile et al. Gene Chromosome Canc 1999; 25: 290-300) (Schraml et al. Am J Pathol 2003; 163: 985 - 992) (Sonoda et al. Gene Chromosome Canc 1997; 20: 320-328)

FIG. 22

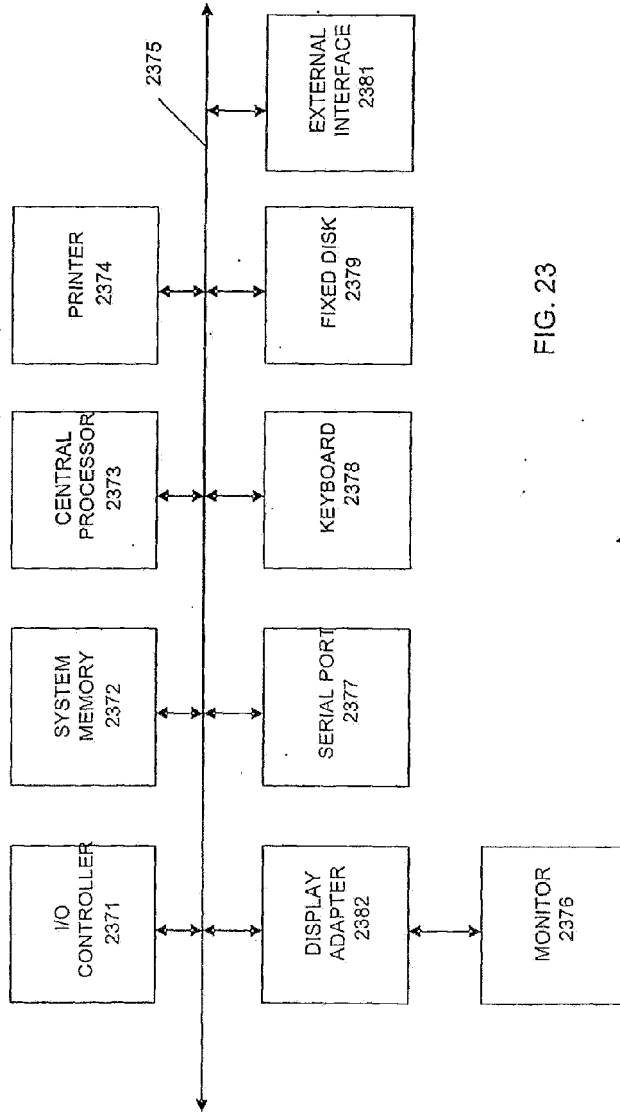


FIG. 23

2375