



US012387819B2

(12) **United States Patent**
Mouliere et al.

(10) **Patent No.:** US 12,387,819 B2
(45) **Date of Patent:** Aug. 12, 2025

(54) **ENHANCED DETECTION OF TARGET DNA BY FRAGMENT SIZE ANALYSIS**

(71) Applicant: **Cancer Research Technology Limited**, London (GB)

(72) Inventors: **Florent Mouliere**, London (GB); **Dineika Chandrananda**, London (GB); **Anna Piskorz**, London (GB); **James Brenton**, London (GB); **Nitzan Rosenfeld**, London (GB)

(73) Assignee: **Cancer Research Technology Limited**, London (GB)

(*) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 1104 days.

(21) Appl. No.: 17/290,924

(22) PCT Filed: Nov. 7, 2019

(86) PCT No.: PCT/EP2019/080506

§ 371 (c)(1),
(2) Date: May 3, 2021

(87) PCT Pub. No.: WO2020/094775

PCT Pub. Date: May 14, 2020

(65) **Prior Publication Data**

US 2023/0014674 A1 Jan. 19, 2023

(30) **Foreign Application Priority Data**

Nov. 7, 2018 (GB) 1818159

(51) **Int. Cl.**

G16B 40/20 (2019.01)
C12N 15/10 (2006.01)

(Continued)

(52) **U.S. Cl.**

CPC G16B 40/20 (2019.02); C12N 15/1093 (2013.01); G16B 30/00 (2019.02);

(Continued)

(58) **Field of Classification Search**

CPC G16B 40/20; G16B 30/00; G16B 40/00; C12N 15/1093; G16H 10/40; G16H 50/20; G16H 50/70; G16H 70/60; C12Q 1/6886

See application file for complete search history.

(56) **References Cited**

U.S. PATENT DOCUMENTS

2018/0307796 A1 10/2018 Jiang et al.

2019/0287645 A1 * 9/2019 Abdueva G16B 30/10
(Continued)

FOREIGN PATENT DOCUMENTS

WO WO 2018/009723 A1 1/2018

WO WO-2019222657 A1 * 11/2019 A61P 35/00

OTHER PUBLICATIONS

Hovelson, Daniel. Precision Oncology Opportunities and Disease Insights From Next-Generation-Sequencing Profiling of Routine Clinical Biospecimens. Diss. 2017. (Year: 2017).*

(Continued)

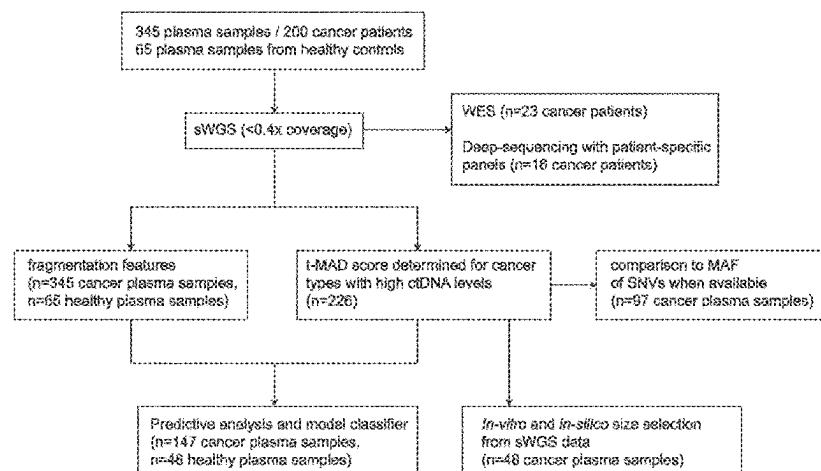
Primary Examiner — Jesse P Frumkin

(74) *Attorney, Agent, or Firm* — Wolf, Greenfield & Sacks, P.C.

(57) **ABSTRACT**

The present invention provides a computer-implemented method for detecting variant nucleic acid from a cell-free nucleic acid-containing sample. The method comprises (a) providing data representing fragment sizes of nucleic acid fragments obtained from said sample and/or representing a measure of deviation from copy number neutrality of the nucleic acid fragments obtained from said sample; b) processing the data from step a) according to a classification algorithm, wherein said classification algorithm operates to classify sample data into one of at least a first class containing the variant nucleic acid and a second class not containing the variant nucleic acid, based on a plurality of cell-free nucleic acid fragment size features and/or a devia-

(Continued)



tion from copy number neutrality feature; and c) outputting the classification of the sample from step b, thereby determining whether the sample contains the variant nucleic acid or not, or a probability that the sample contains the variant nucleic acid. Related methods are also provided.

18 Claims, 32 Drawing Sheets

(51) Int. Cl.

G16B 30/00	(2019.01)
G16H 10/40	(2018.01)
G16H 50/20	(2018.01)
G16H 50/70	(2018.01)
G16H 70/60	(2018.01)

(52) U.S. Cl.

CPC	G16H 10/40 (2018.01); G16H 50/20 (2018.01); G16H 50/70 (2018.01); G16H 70/60 (2018.01)
-----------	--

(56)

References Cited

U.S. PATENT DOCUMENTS

- 2021/0002728 A1* 1/2021 Landau G16B 20/00
2021/0174958 A1* 6/2021 Drake G16H 20/10

OTHER PUBLICATIONS

- Kurtz, David M. Personalized Risk Assessment and Disease Monitoring in Non-Hodgkin Lymphoma from Circulating Tumor DNA. Diss. Stanford University, 2017. (Year: 2017).*
- Lapin, Morten, et al. "Fragment size and level of cell-free DNA provide prognostic information in patients with advanced pancreatic cancer." *Journal of translational medicine* 16 (2018): 1-10. (Year: 2018).*
- PCT/EP2019/080506, Feb. 26, 2020, International Search Report and Written Opinion.
- PCT/EP2019/080506, May 20, 2021, International Preliminary Report on Patentability.
- International Preliminary Report on Patentability for International Application No. PCT/EP2019/080506 mailed May 20, 2021.
- Abbosh et al., Phylogenetic ctDNA analysis depicts early-stage lung cancer evolution. *Nature*. May 25, 2017;545(7655):446-51.
- Adalsteinsson et al., Scalable whole-exome sequencing of cell-free DNA reveals high concordance with metastatic tumors. *Nature communications*. Nov. 6, 2017;8(1):1324. 13 pages.
- Belic et al., Rapid identification of plasma DNA samples with increased ctDNA levels by a modified FAST-SeqS approach. *Clinical chemistry*. Jun. 1, 2015;61(6):838-49.
- Best et al., RNA-Seq of tumor-educated platelets enables blood-based pan-cancer, multiclass, and molecular pathway cancer diagnostics. *Cancer cell*. Nov. 9, 2015;28(5):666-76.
- Best et al., Swarm intelligence-enhanced detection of non-small-cell lung cancer using tumor-educated platelets. *Cancer cell*. Aug. 14, 2017;32(2):238-52.
- Bettegowda et al., Detection of circulating tumor DNA in early-and late-stage human malignancies. *Science translational medicine*. Feb. 19, 2014;6(224):224ra24-. 12 pages.
- Bronkhorst et al., Characterization of the cell-free DNA released by cultured cancer cells. *Biochimica et Biophysica Acta (BBA)-Molecular Cell Research*. Jan. 1, 2016;1863(1):157-65.
- Burnham et al., Single-stranded DNA library preparation uncovers the origin and diversity of ultrashort cell-free DNA in plasma. *Scientific reports*. Jun. 14, 2016;6(1):27859. 9 pages.
- Chandrananda et al., High-resolution characterization of sequence signatures due to non-random cleavage of cell-free DNA. *BMC medical genomics*. Dec. 2015;8:1-9.
- Chaudhuri et al., Early detection of molecular residual disease in localized lung cancer by circulating tumor DNA profiling. *Cancer discovery*. Dec. 1, 2017;7(12):1394-403.
- Cohen et al., Detection and localization of surgically resectable cancers with a multi-analyte blood test. *Science*. Feb. 23, 2018;359(6378):926-30. 5 pages.
- Dawson et al., Analysis of circulating tumor DNA to monitor metastatic breast cancer. *New England Journal of Medicine*. Mar. 28, 2013;368(13):1199-209.
- Diehl et al., Circulating mutant DNA to assess tumor dynamics. *Nature medicine*. Sep. 2008;14(9):985-90.
- Diehl et al., Detection and quantification of mutations in the plasma of patients with colorectal tumors. *Proceedings of the National Academy of Sciences*. Nov. 8, 2005;102(45):16368-73.
- Forshew et al., Noninvasive identification and monitoring of cancer mutations by targeted deep sequencing of plasma DNA. *Science translational medicine*. May 30, 2012;4(136):136ra68-. 13 pages.
- Genovese et al., Clonal hematopoiesis and blood-cancer risk inferred from blood DNA sequence. *New England Journal of Medicine*. Dec. 25, 2014;371(26):2477-87.
- Giaccona et al., Cell-free DNA in human blood plasma: length measurements in patients with pancreatic cancer and healthy controls. *Pancreas*. Jul. 1, 1998;17(1):89-97.
- Hanahan et al., Hallmarks of cancer: the next generation. *cell*. Mar. 4, 2011;144(5):646-74.
- Haque et al., Challenges in using ctDNA to achieve early detection of cancer. *BioRxiv*. Dec. 21, 2017:237578. 20 pages.
- Heitzer et al., Tumor-associated copy number changes in the circulation of patients with prostate cancer identified through whole-genome sequencing. *Genome medicine*. Dec. 2013;5(4):1-6.
- Hu et al., False-positive plasma genotyping due to clonal hematopoiesis. *Clinical Cancer Research*. Sep. 15, 2018;24(18):4437-43.
- Jahr et al., DNA fragments in the blood plasma of cancer patients: quantitations and evidence for their origin from apoptotic and necrotic cells. *Cancer research*. Feb. 2, 2001;61(4):1659-65.
- Jiang et al., The long and short of circulating cell-free DNA and the ins and outs of molecular diagnostics. *Trends in Genetics*. Jun. 1, 2016;32(6):360-71.
- Jiang et al., Lengthening and shortening of plasma DNA in hepatocellular carcinoma patients. *Proceedings of the National Academy of Sciences*. Mar. 17, 2015;112(11):E1317-25.
- Li et al., Fast and accurate short read alignment with Burrows-Wheeler transform. *bioinformatics*. Jul. 15, 2009;25(14):1754-60.
- Lo et al., Maternal plasma DNA sequencing reveals the genome-wide genetic and mutational profile of the fetus. *Science translational medicine*. Dec. 8, 2010;2(61):61ra91-. 14 pages.
- Lun et al., Noninvasive prenatal diagnosis of monogenic diseases by digital size selection and relative mutation dosage on DNA in maternal plasma. *Proceedings of the National Academy of Sciences*. Dec. 16, 2008;105(50):19920-5.
- Macintyre et al., Copy number signatures and mutational processes in ovarian carcinoma. *Nature genetics*. Sep. 2018;50(9):1262-70.
- Minarik et al., Utilization of benchtop next generation sequencing platforms ion torrent PGM and MiSeq in noninvasive prenatal testing for chromosome 21 trisomy and testing of impact of in silico and physical size selection on its analytical performance. *PLoS one*. Dec. 15, 2015;10(12):e0144811. 12 pages.
- Mouliere et al., Multi-marker analysis of circulating cell-free DNA toward personalized medicine for colorectal cancer. *Molecular oncology*. Jul. 1, 2014;8(5):927-41.
- Mouliere et al., High fragmentation characterizes tumour-derived circulating DNA. *PLoS one*. Sep. 6, 2011;6(9):e23418. 10 pages.
- Murtaza et al., Non-invasive analysis of acquired resistance to cancer therapy by sequencing of plasma DNA. *Nature*. May 2, 2013;497(7447):108-12.
- Newman et al., Integrated digital error suppression for improved detection of circulating tumor DNA. *Nature biotechnology*. May 2016;34(5):547-55.
- Parkinson et al., Exploratory analysis of TP53 mutations in circulating tumour DNA as biomarkers of treatment response for patients with relapsed high-grade serous ovarian carcinoma: a retrospective study. *PLoS medicine*. Dec. 20, 2016;13(12):e1002198. 25 pages.

(56)

References Cited**OTHER PUBLICATIONS**

- Patel et al., Association of plasma and urinary mutant DNA with clinical outcomes in muscle invasive bladder cancer. *Scientific reports.* Jul. 17, 2017;7(1):5554. 12 pages.
- Riediger et al., Mutation analysis of circulating plasma DNA to determine response to EGFR tyrosine kinase inhibitor therapy of lung adenocarcinoma patients. *Scientific reports.* Sep. 19, 2016;6(1):33505. 8 pages.
- Routy et al., Gut microbiome influences efficacy of PD-1-based immunotherapy against epithelial tumors. *Science.* Jan. 5, 2018;359(6371):91-7. 8 pages.
- Scheinin et al., DNA copy number analysis of fresh and formalin-fixed specimens by shallow whole-genome sequencing with identification and exclusion of problematic regions in the genome assembly. *Genome research.* Dec. 1, 2014;24(12):2022-32.
- Siravegna et al., Integrating liquid biopsies into the management of cancer. *Nature reviews Clinical oncology.* Sep. 2017;14(9):531-48.
- Snyder et al., Cell-free DNA comprises an *in vivo* nucleosome footprint that informs its tissues-of-origin. *Cell.* Jan. 14, 2016;164(1):57-68.
- Stover et al., Association of cell-free DNA tumor fraction and somatic copy number alterations with survival in metastatic triple-negative breast cancer. *Journal of Clinical Oncology.* Feb. 2, 2018;36(6):543. 21 pages.
- Thierry et al., Origins, structures, and functions of circulating DNA in oncology. *Cancer and metastasis reviews.* Sep. 2016;35:347-76.
- Tie et al., Circulating tumor DNA analysis detects minimal residual disease and predicts recurrence in patients with stage II colon cancer. *Science translational medicine.* Jul. 6, 2016;8(346):346ra92-. 11 pages.
- Ulz et al., Inferring expressed genes by whole-genome sequencing of plasma DNA. *Nature genetics.* Oct. 2016;48(10):1273-8.
- Umetani et al., Prediction of breast tumor progression by integrity of free circulating DNA in serum. *Journal of clinical oncology.* Sep. 10, 2006;24(26):4270-6.
- Underhill et al., Fragment length of circulating tumor DNA. *PLoS genetics.* Jul. 18, 2016;12(7):e1006162. 24 pages.
- Wan et al., Liquid biopsies come of age: towards implementation of circulating tumour DNA. *Nature Reviews Cancer.* Apr. 2017;17(4):223-38.
- Yu et al., Size-based molecular diagnostics using plasma DNA for noninvasive prenatal testing. *Proceedings of the National Academy of Sciences.* Jun. 10, 2014;111(23):8583-8.
- Zill et al., The landscape of actionable genomic alterations in cell-free circulating tumor DNA from 21,807 advanced cancer patients. *Clinical Cancer Research.* Aug. 1, 2018;24(15):3528-38. International Search Report and Written Opinion for International Application No. PCT/EP2019/080506 mailed Feb. 26, 2020.
- Adalsteinsson et al., Scalable whole-exome sequencing of cell-free DNA reveals high concordance with metastatic tumors. *Nature communications.* Nov. 6, 2017;8(1):1-26. Together with Supplemental Materials available at <https://www.nature.com/articles/s41467-017-00965-y>:521 pages.
- Mouliere et al., Selecting short DNA fragments in plasma improves detection of circulating tumour DNA. *BioRxiv.* Jan. 1, 2017:134437. Together with Supplemental Materials available at <https://www.biorxiv.org/content/10.1101/134437v1>:364 pages.
- Underhill et al., Fragment length of circulating tumor DNA. *PLoS genetics.* Jul. 18, 2016;12(7):e1006162. 29 pages.

* cited by examiner

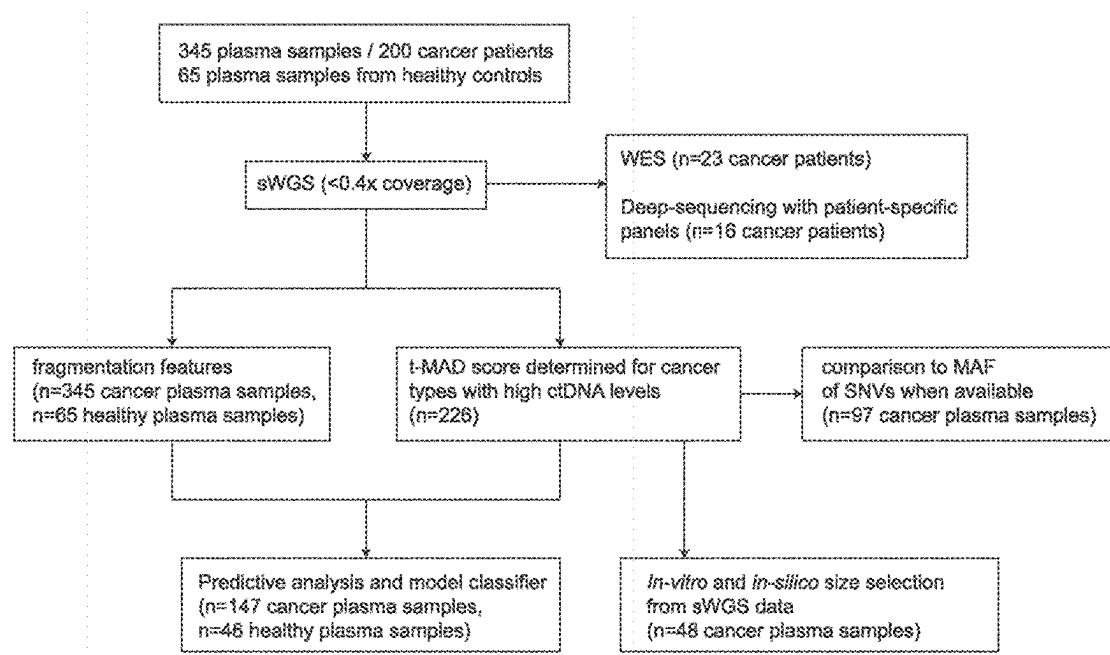


FIG. 1

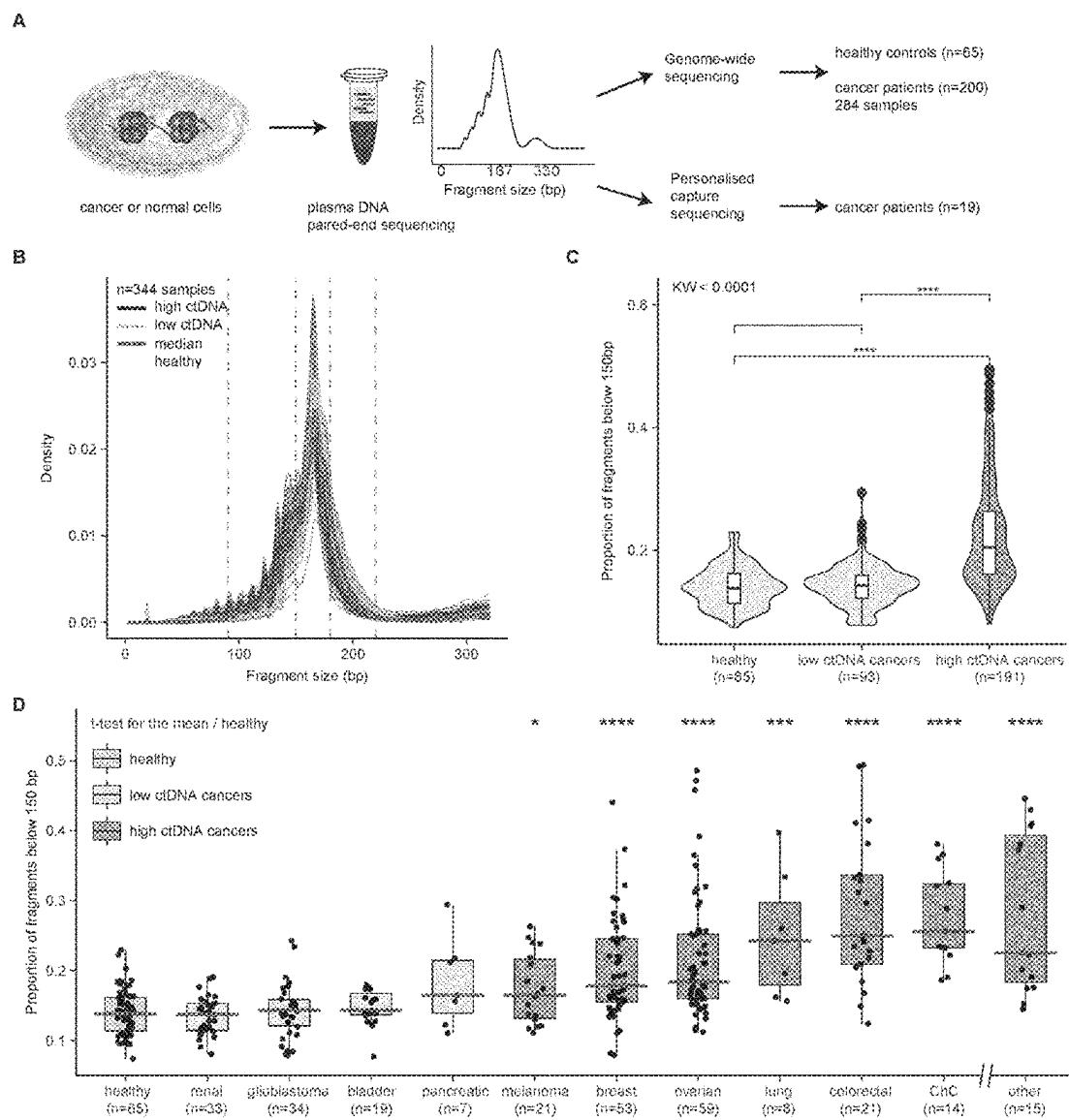


FIG. 2

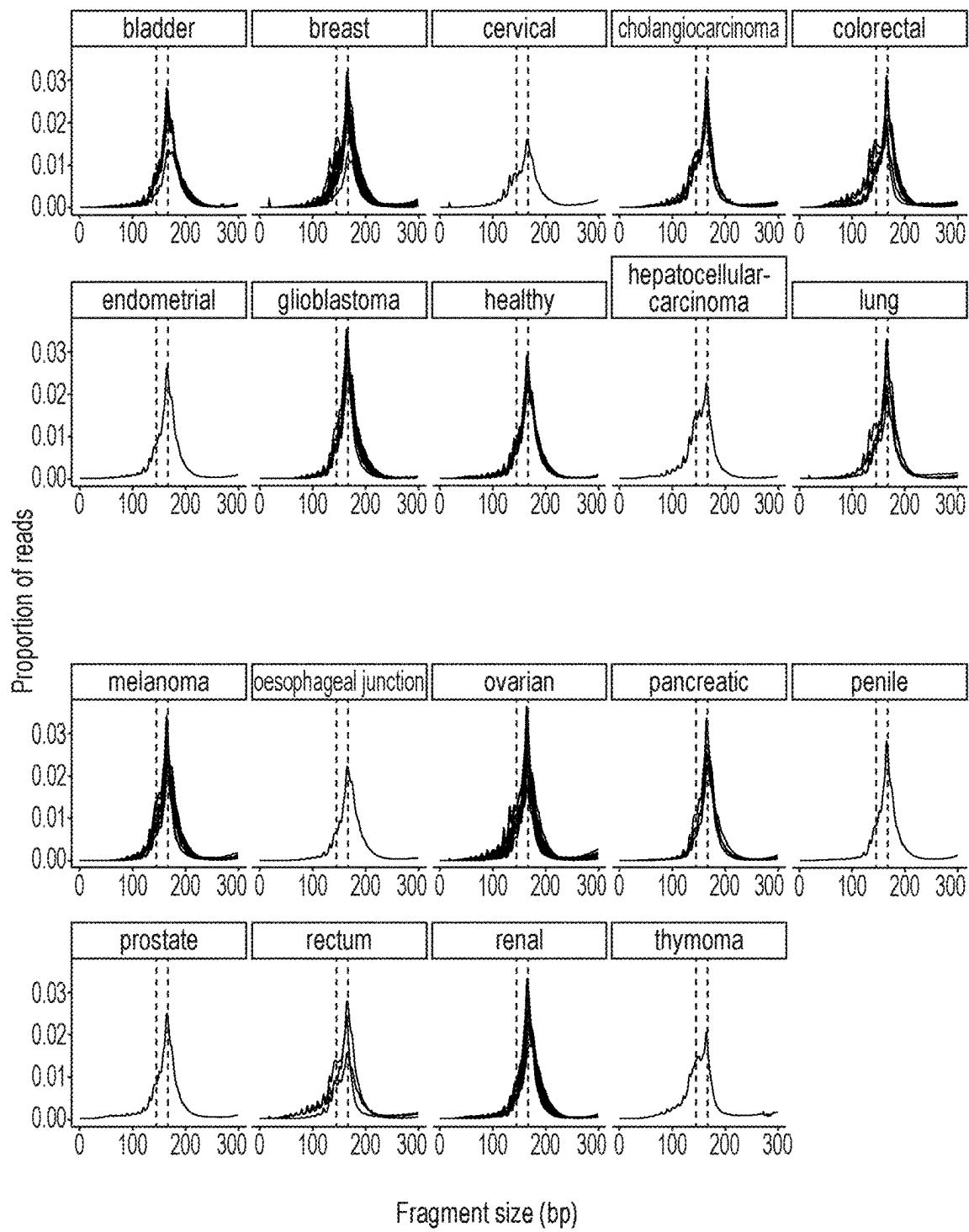


FIG. 3

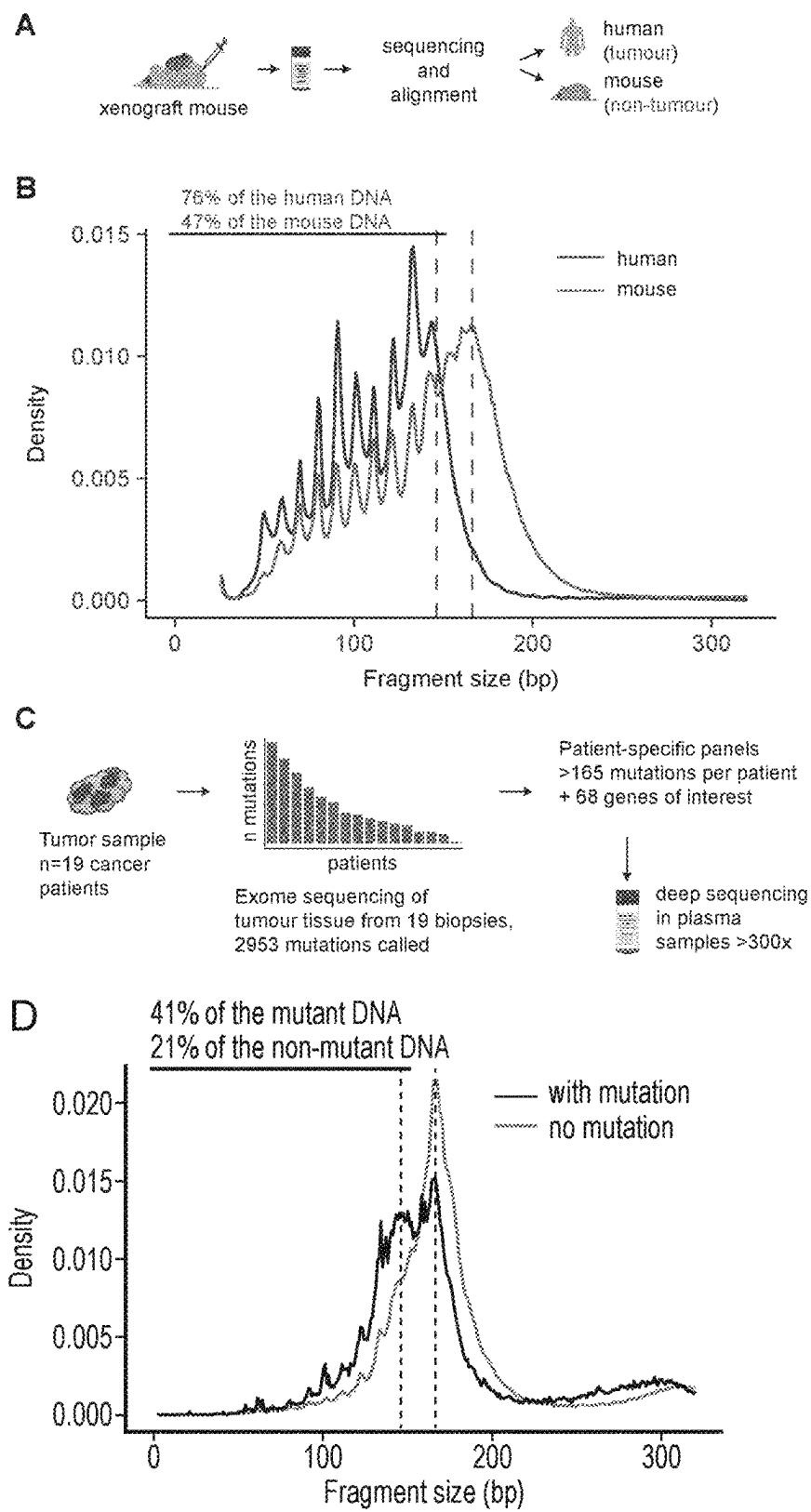


FIG. 4

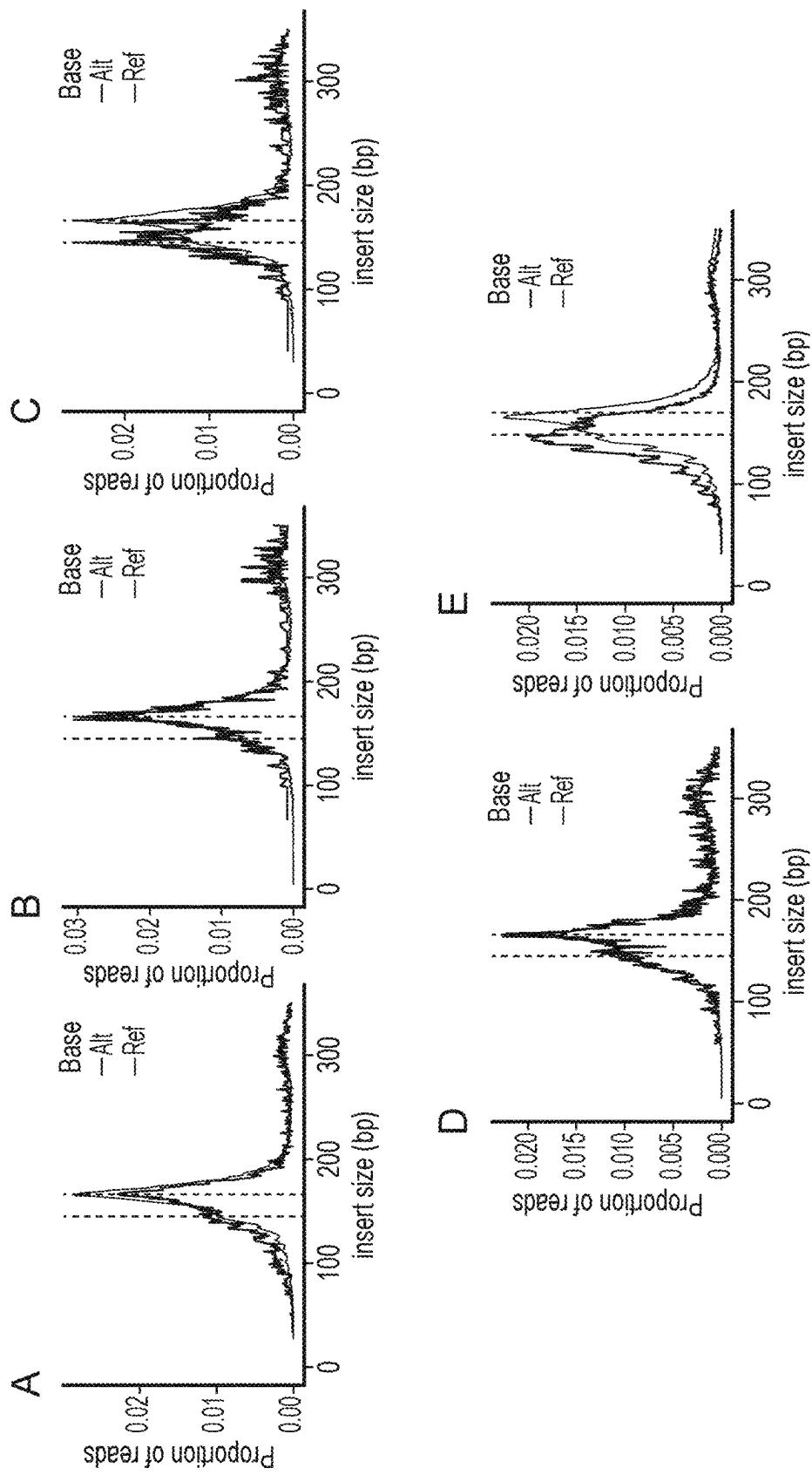


FIG. 5

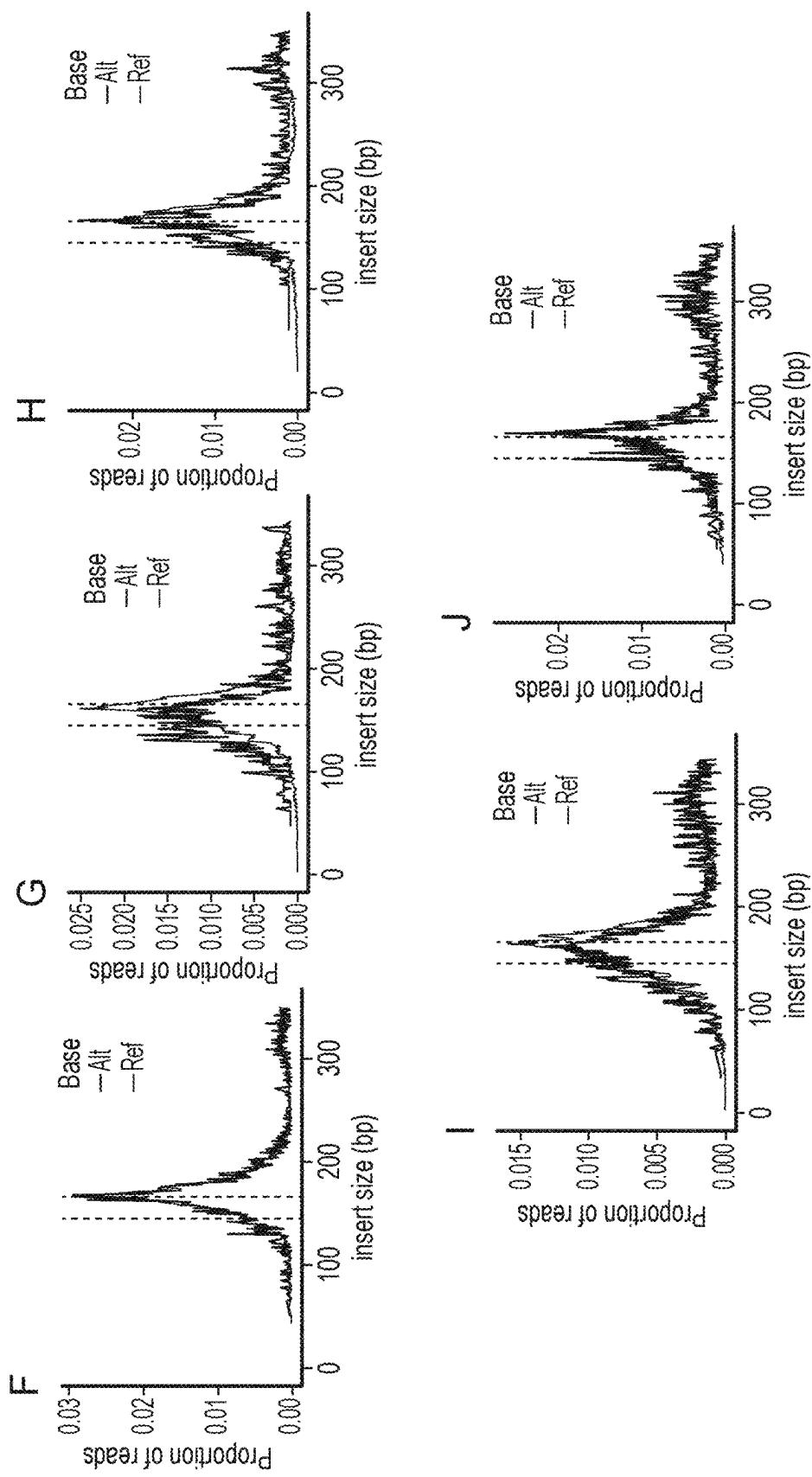


FIG. 5 (Continued)

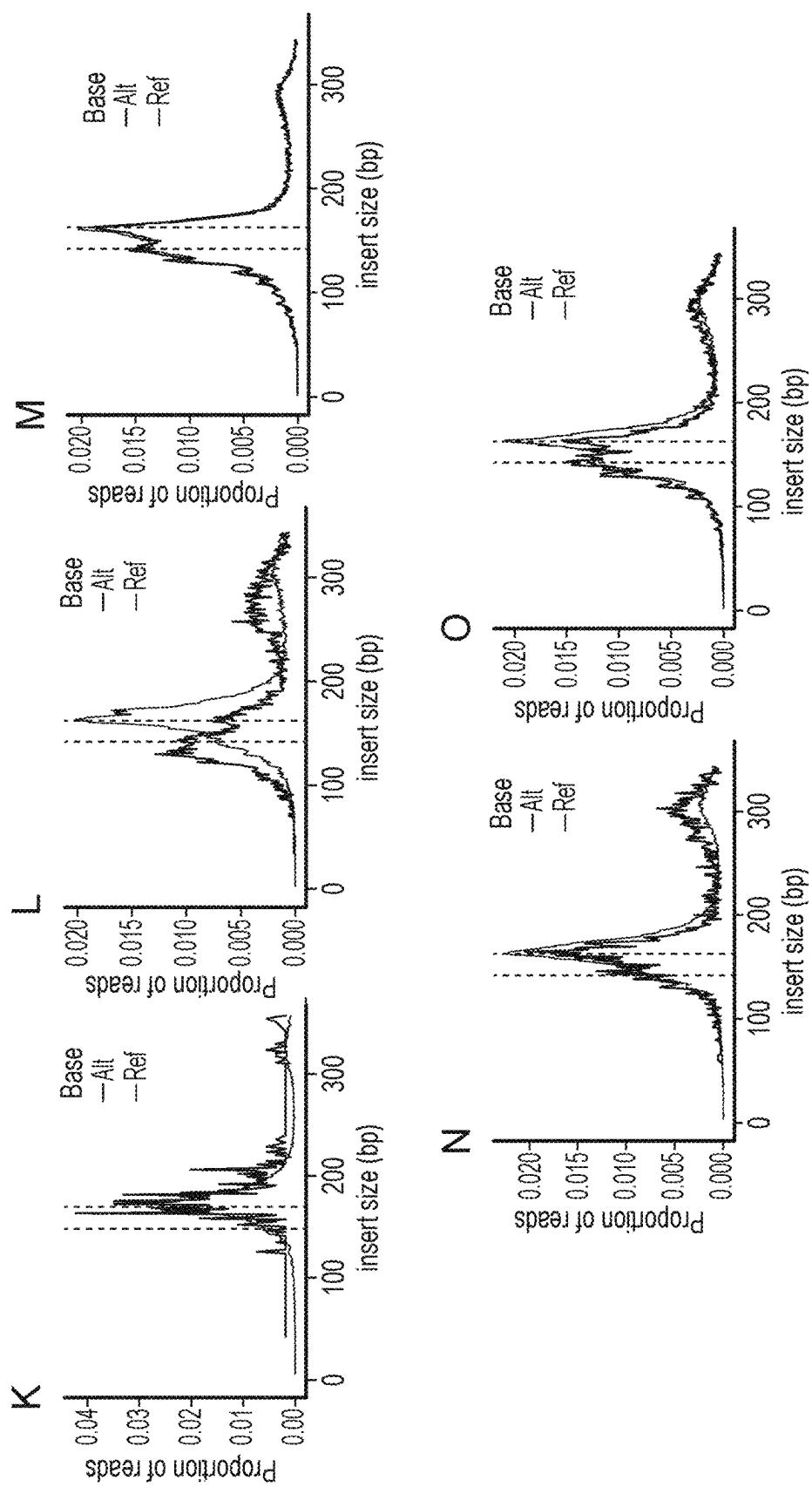


FIG. 5 (Continued)

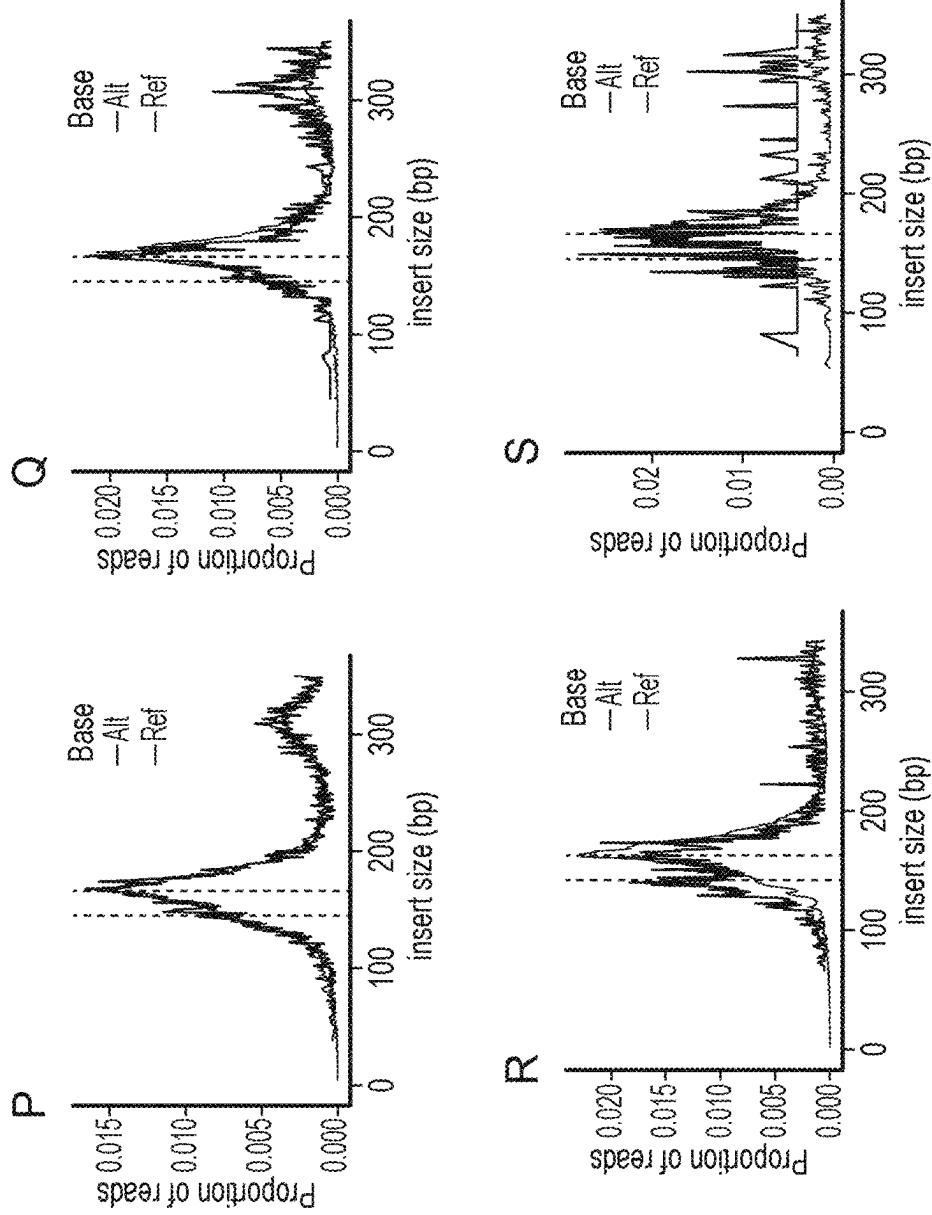


FIG. 5 (Continued)

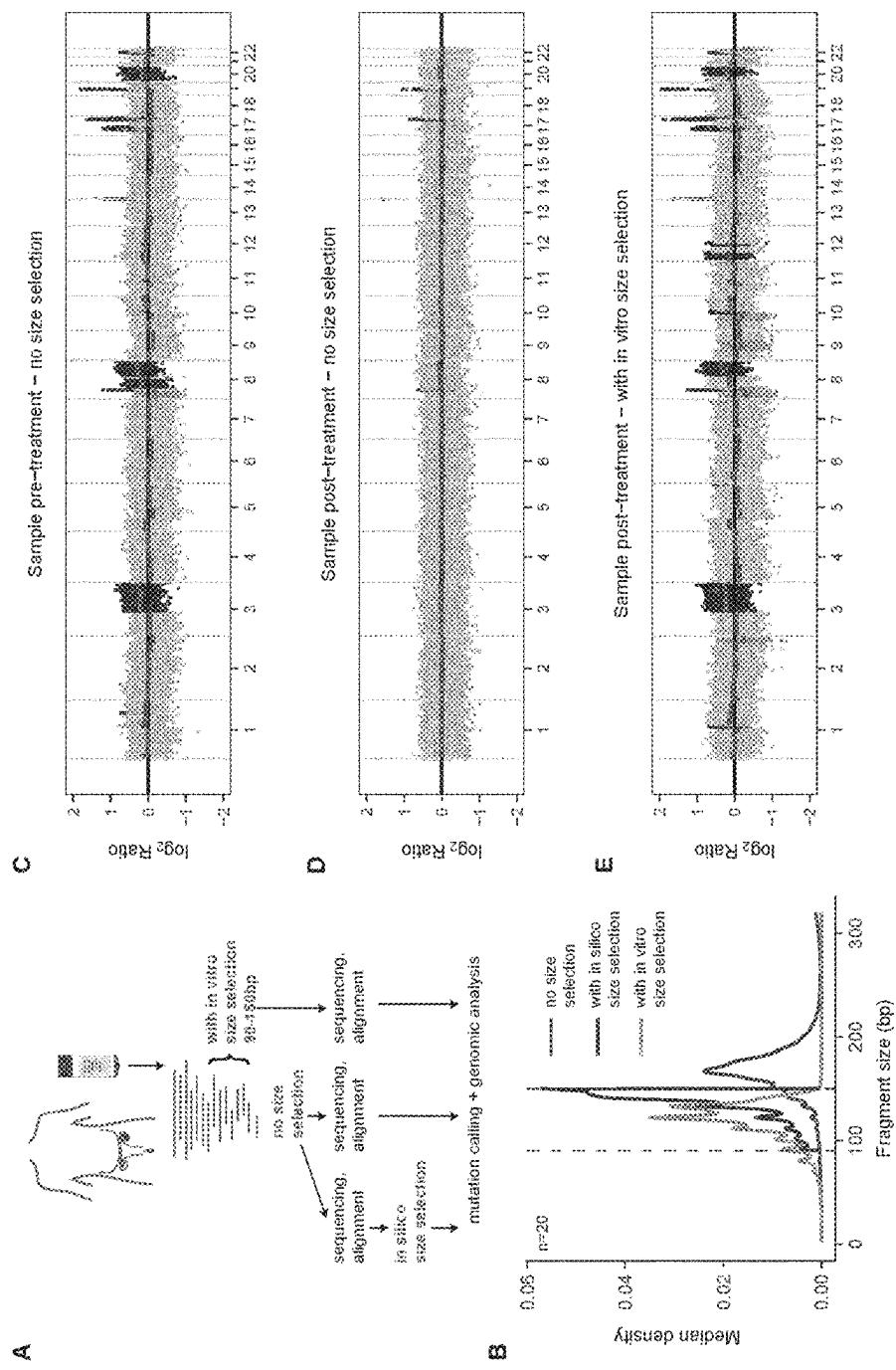


FIG. 6

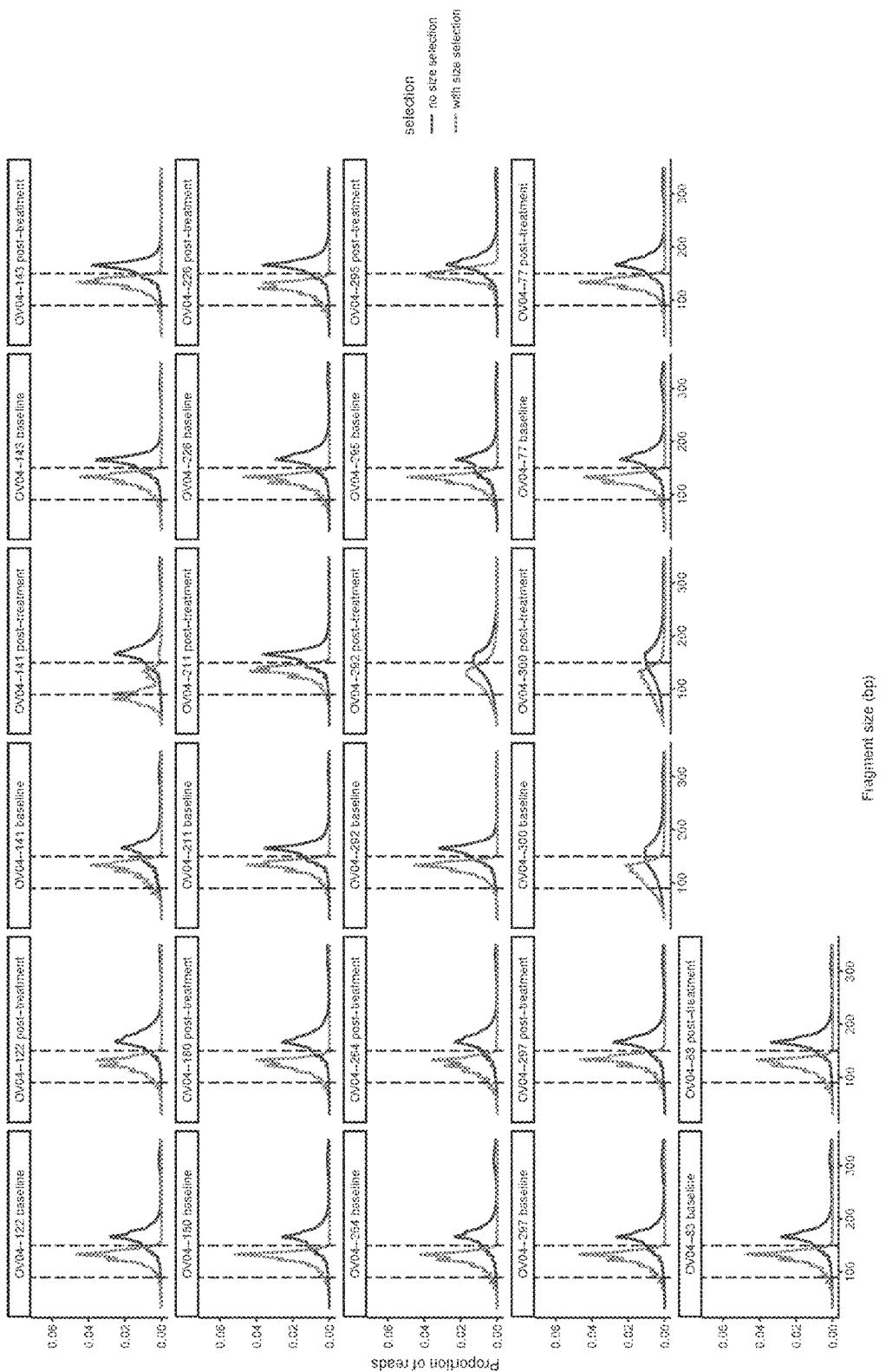


FIG. 7

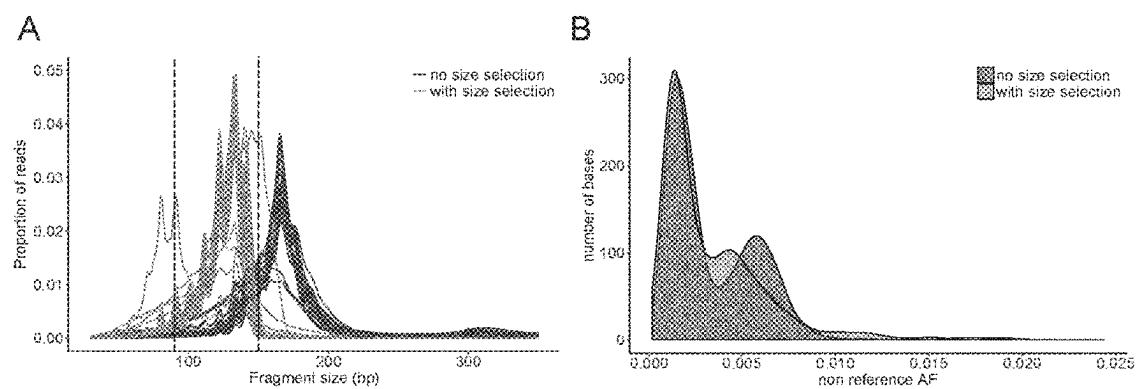


FIG. 8

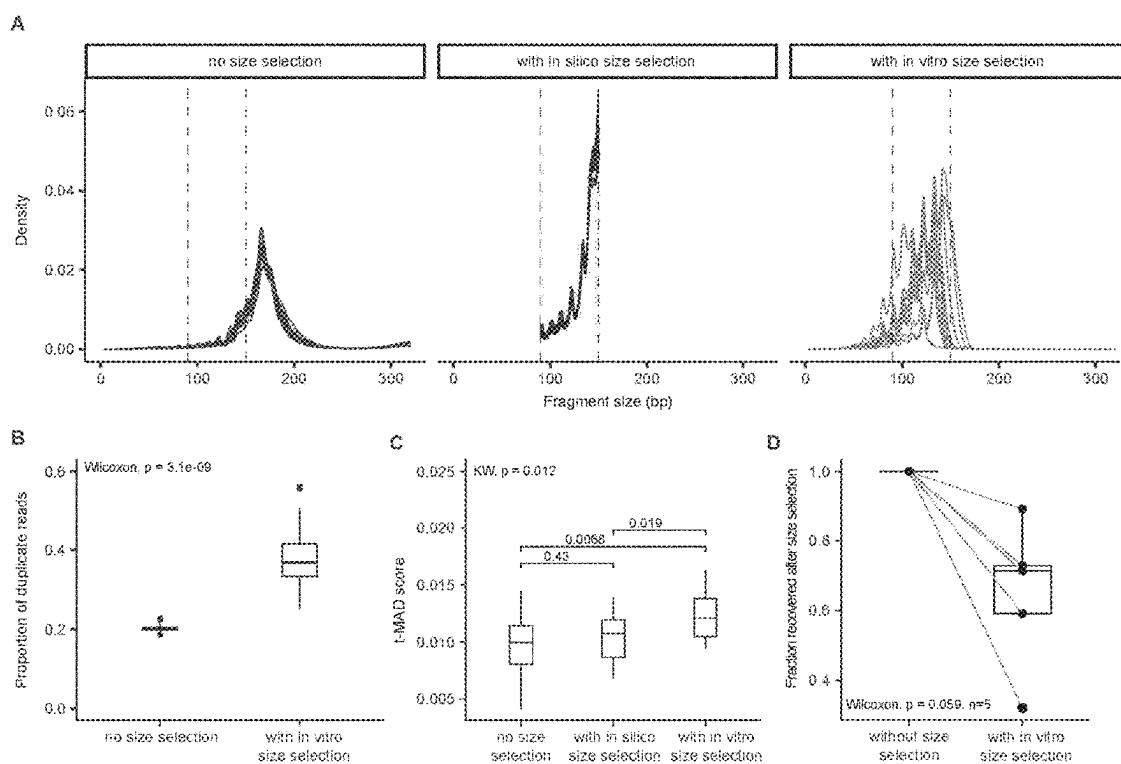


FIG. 9

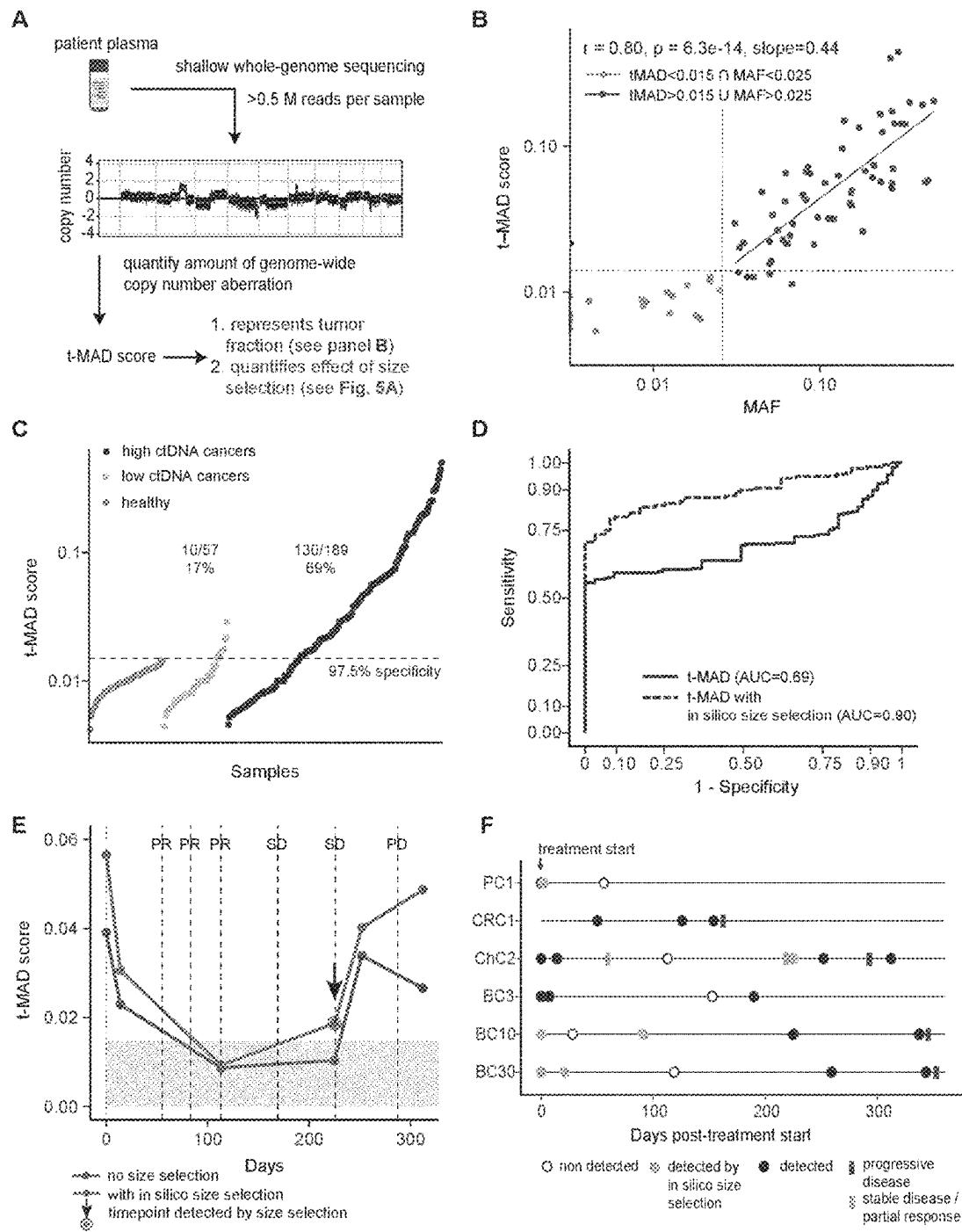


FIG. 10

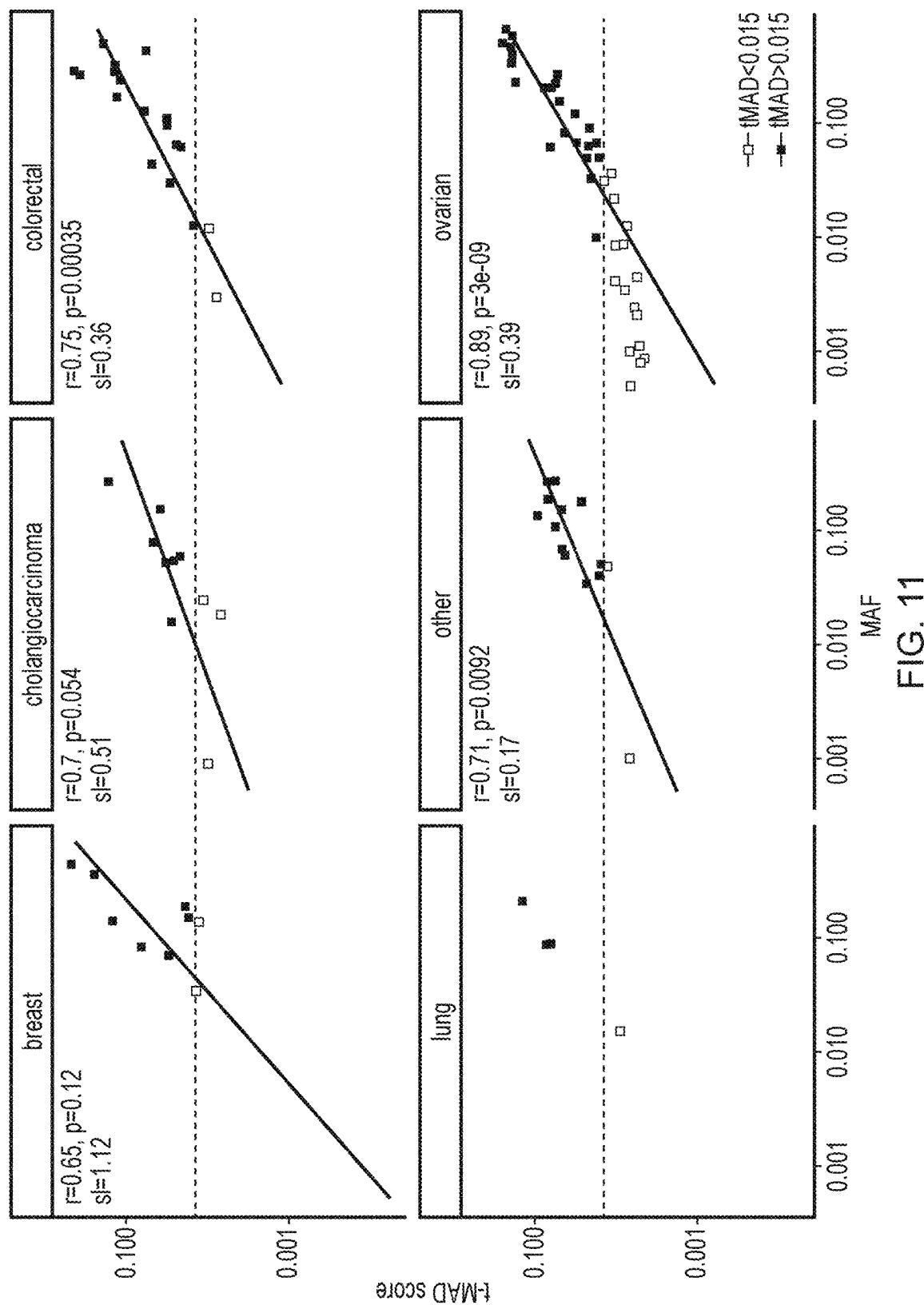


FIG. 11

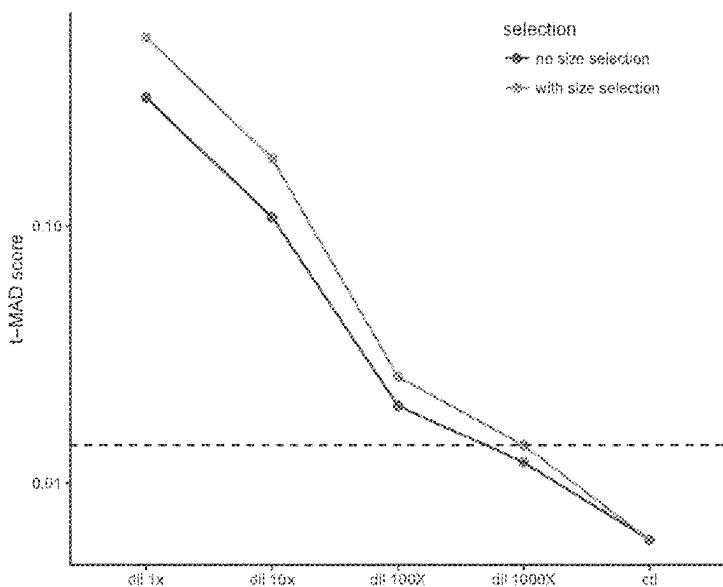


FIG. 12

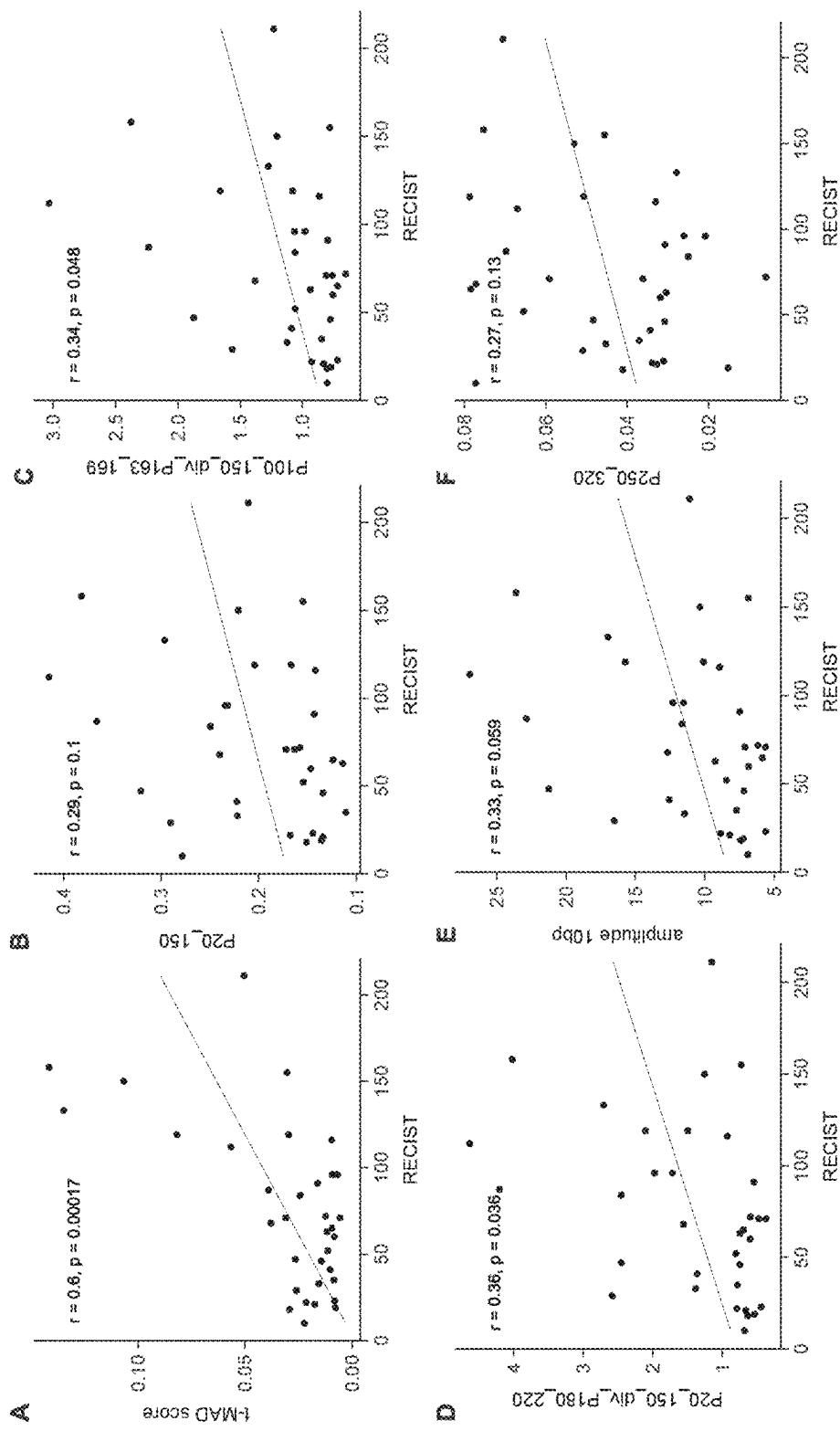
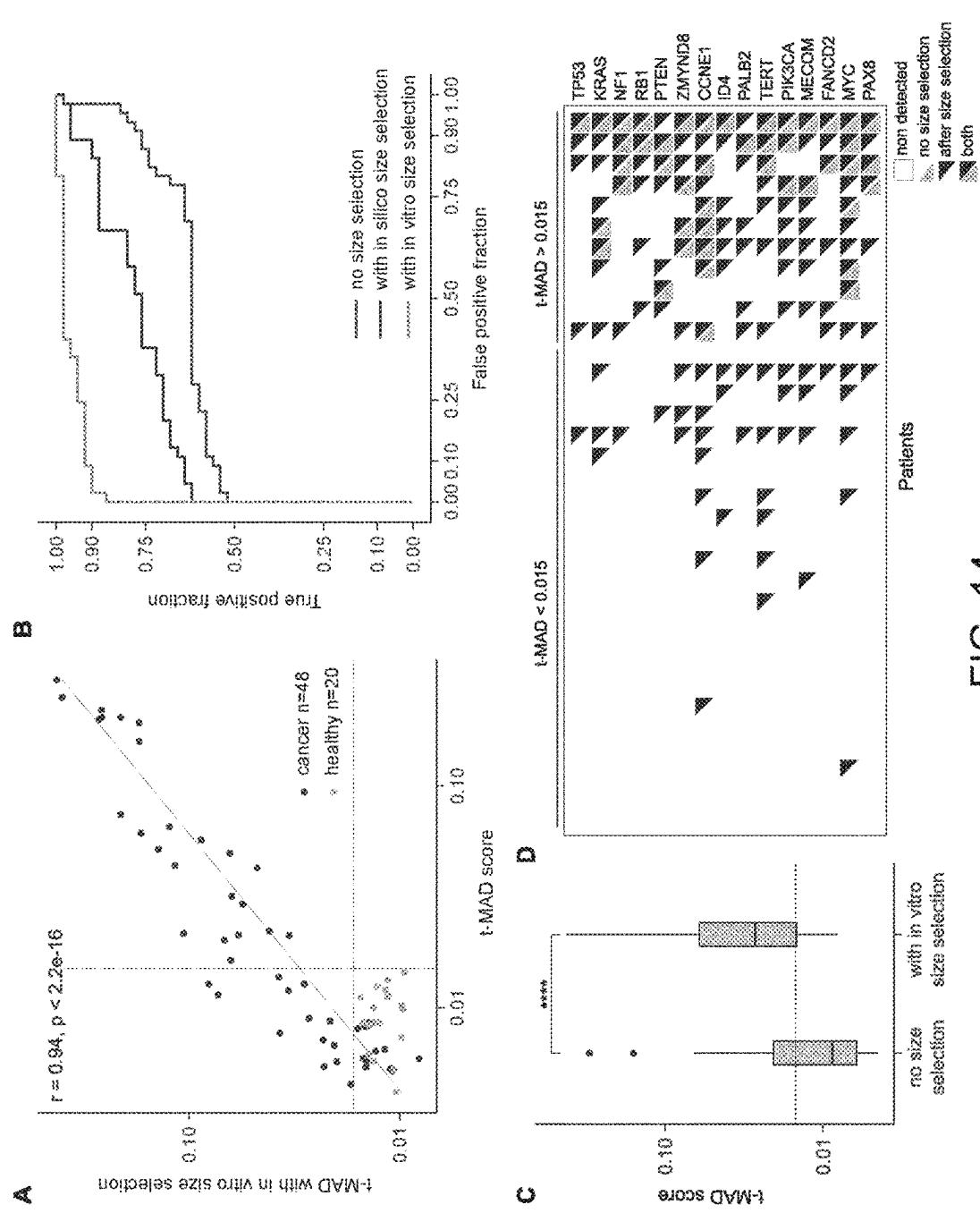


FIG. 13

**FIG. 14**

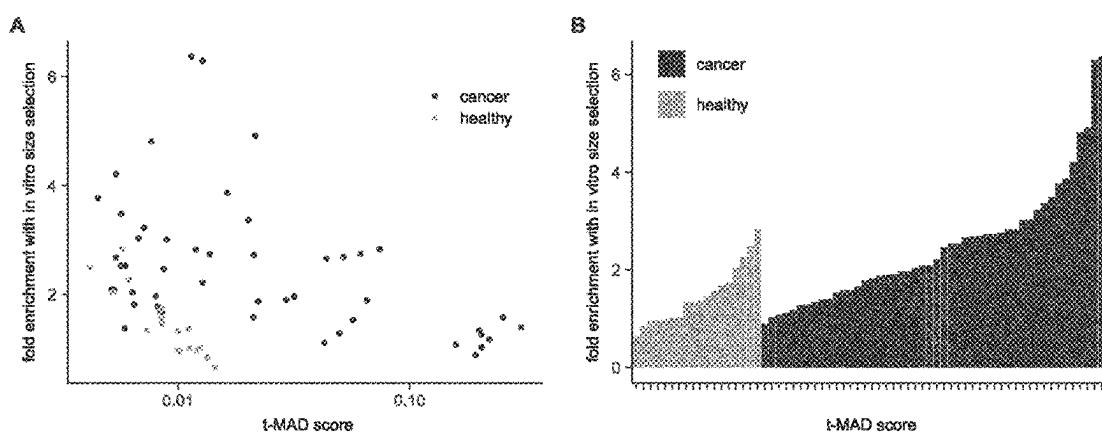


FIG. 15

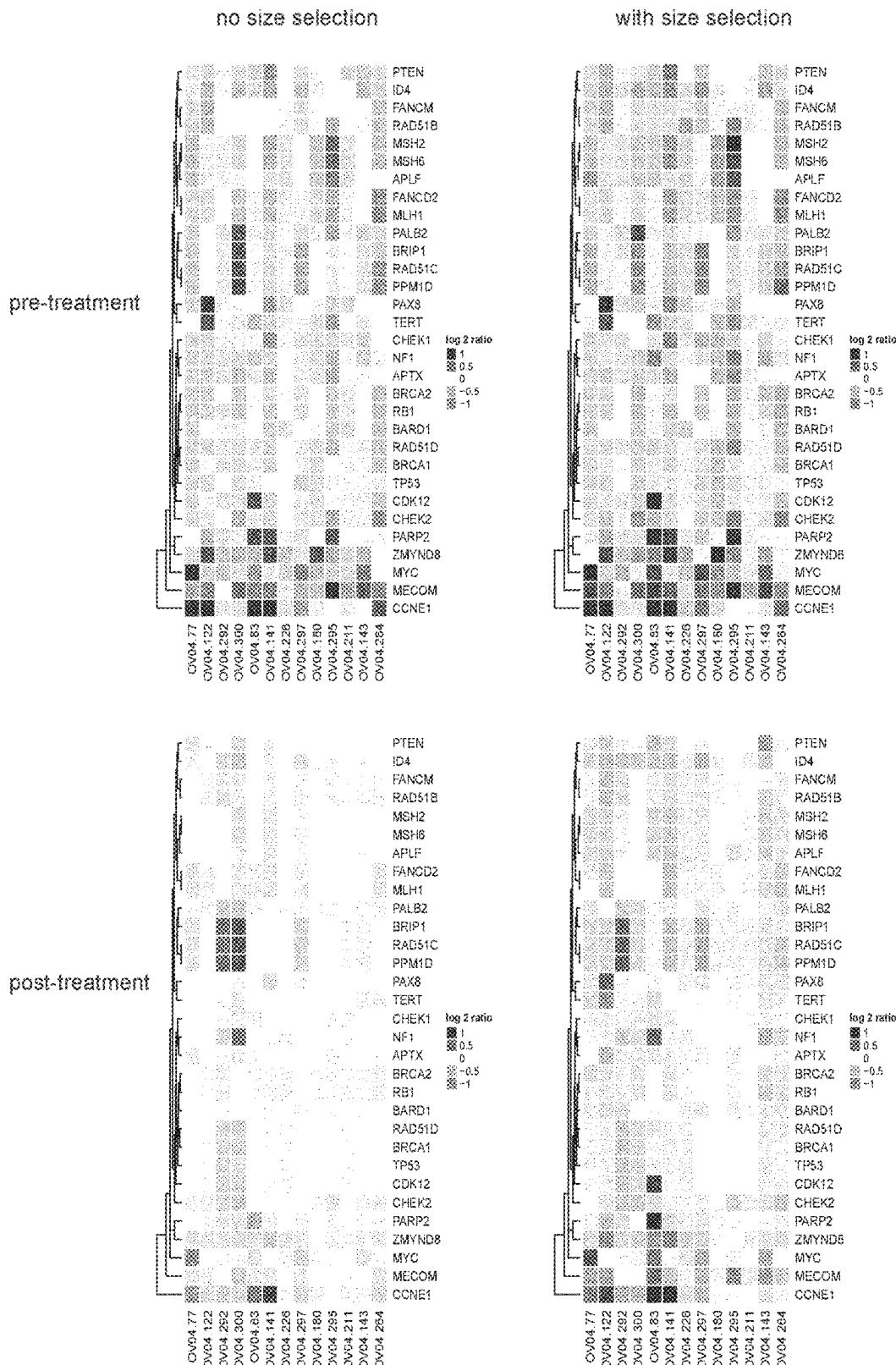


FIG. 16

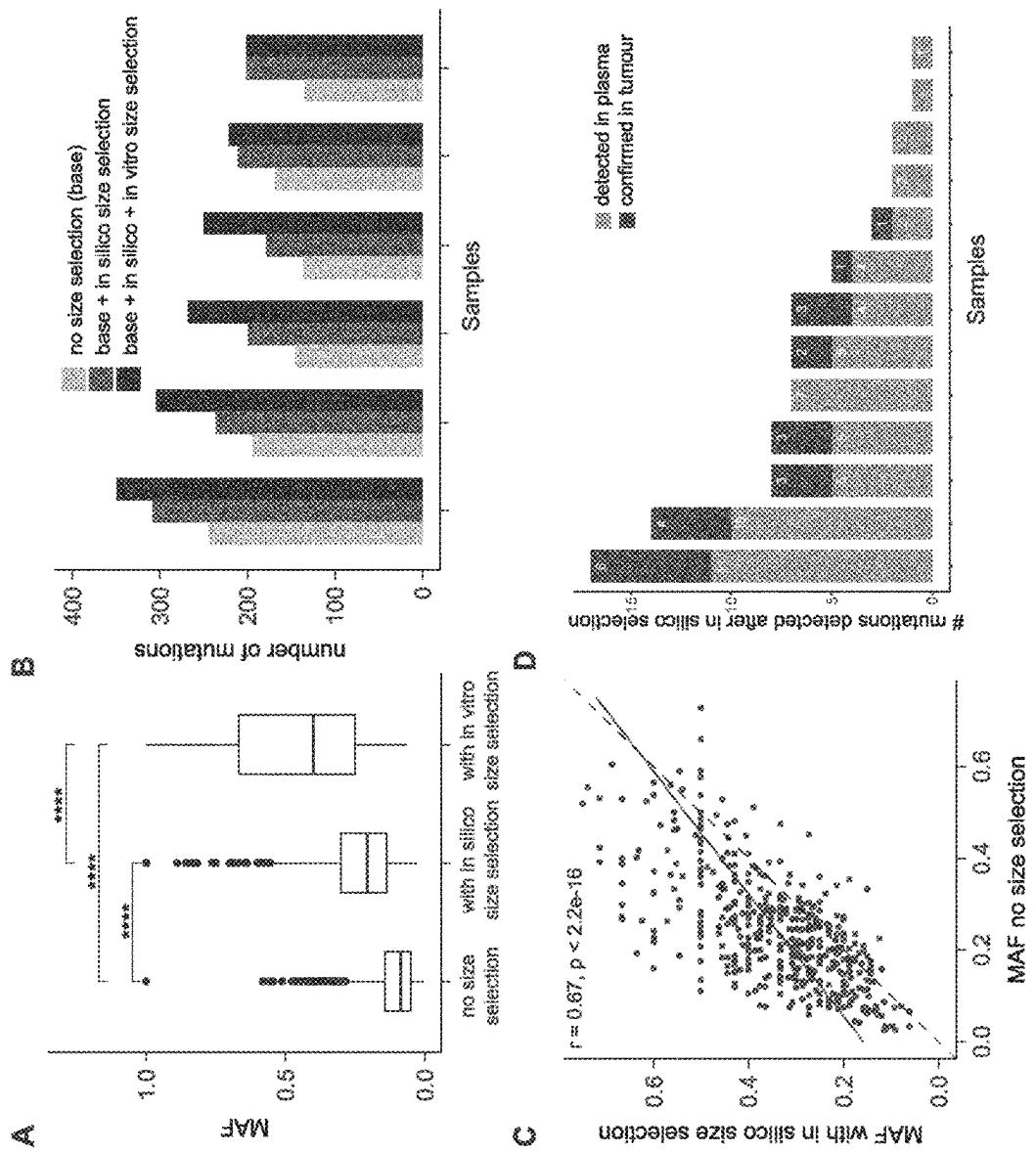


FIG. 17

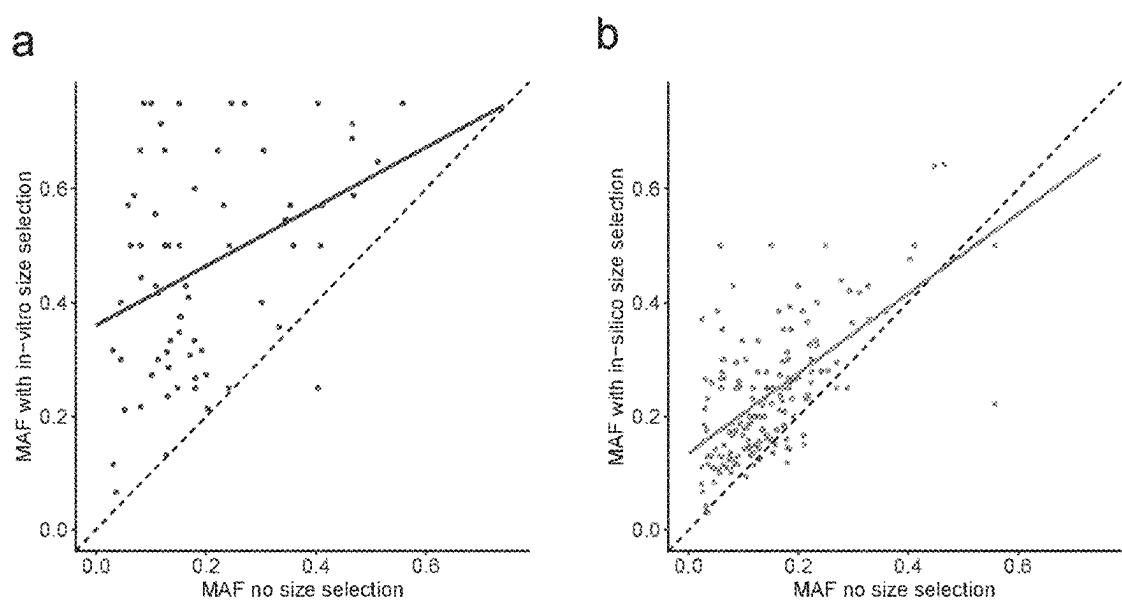


FIG. 18

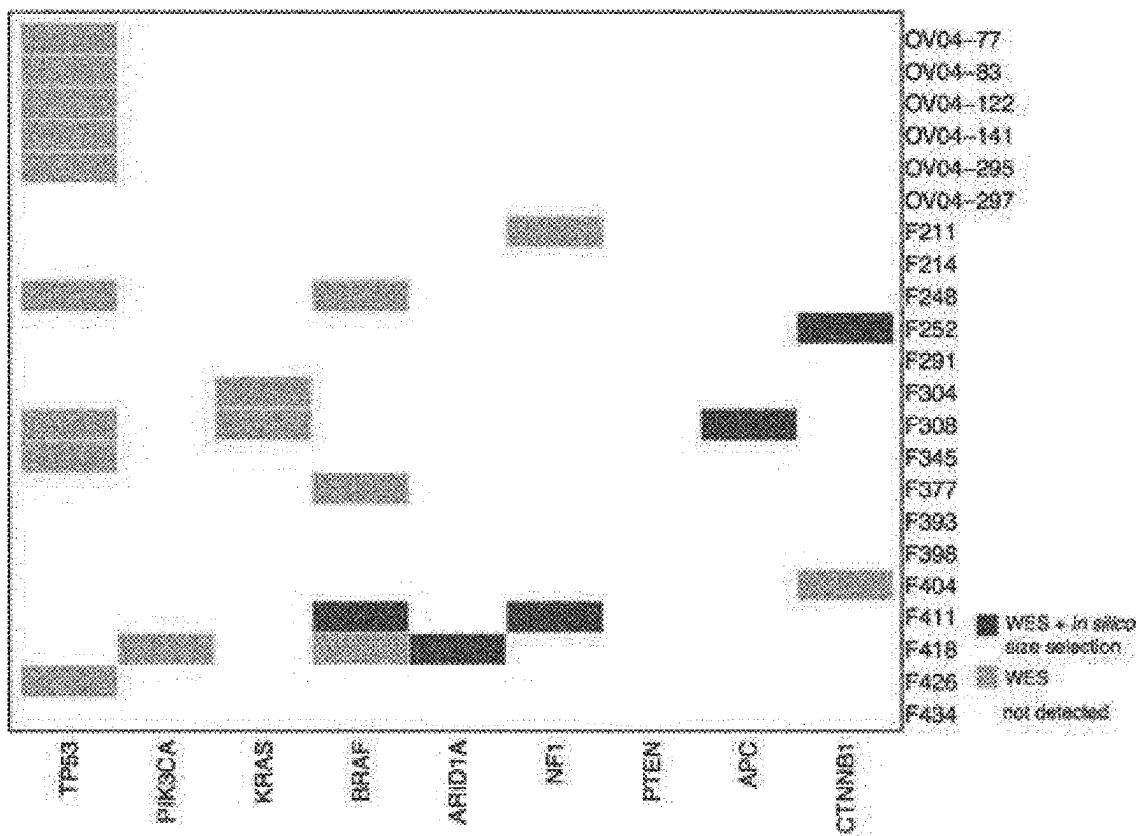


FIG. 19

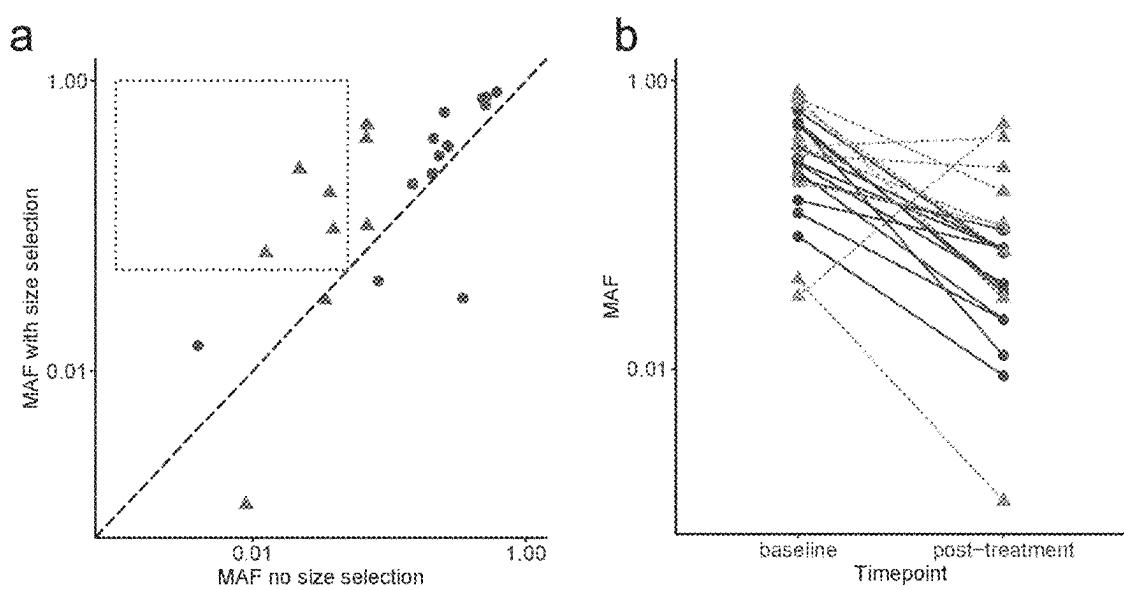


FIG. 20

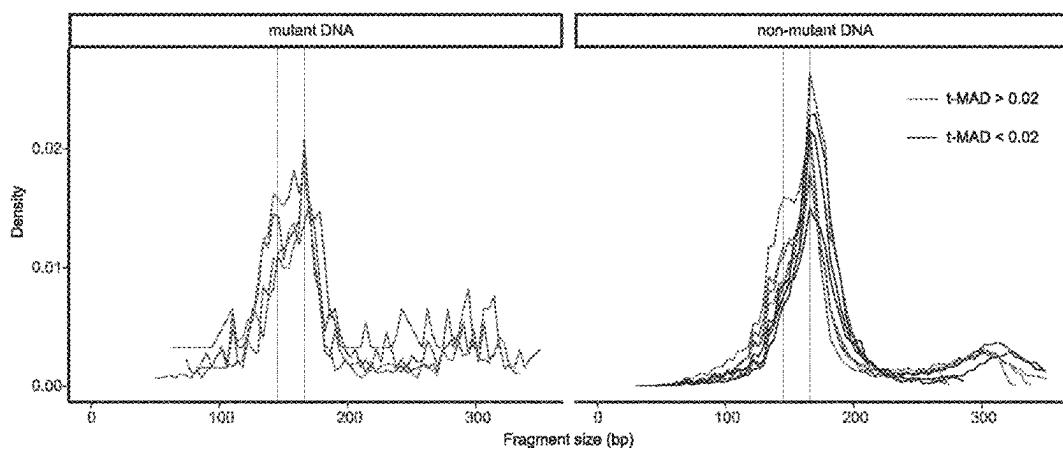


FIG. 21

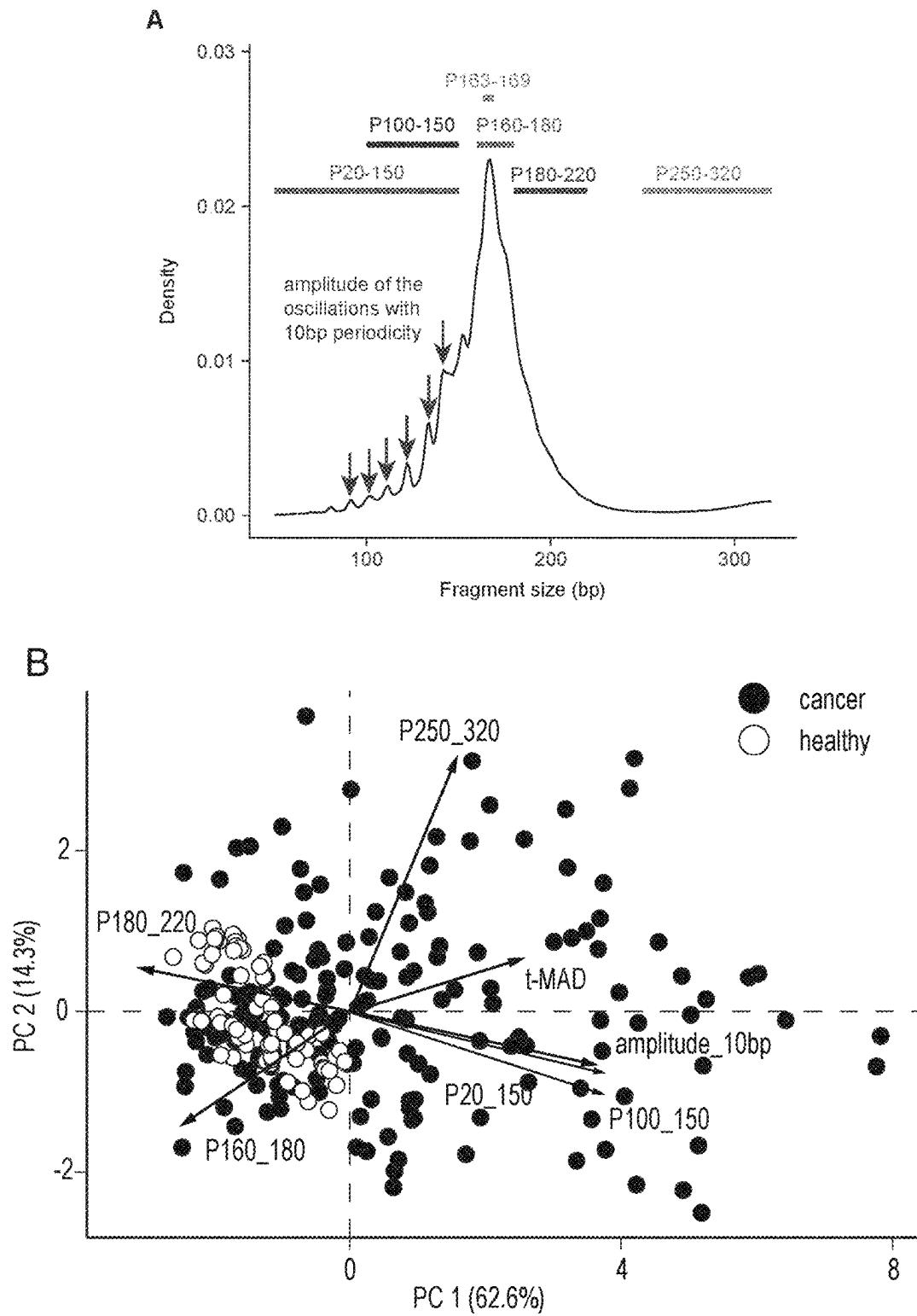


FIG. 22

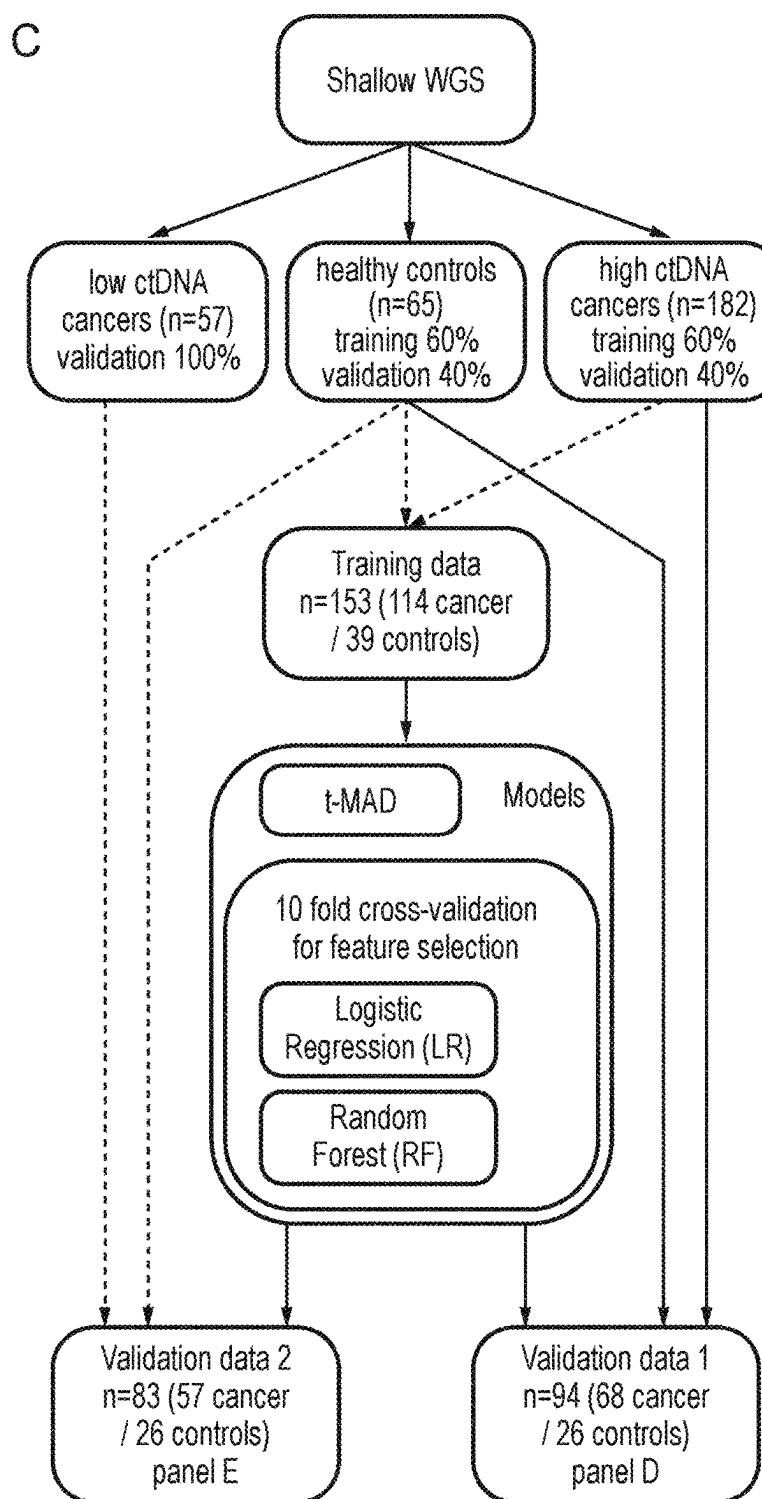


FIG. 22 (Continued)

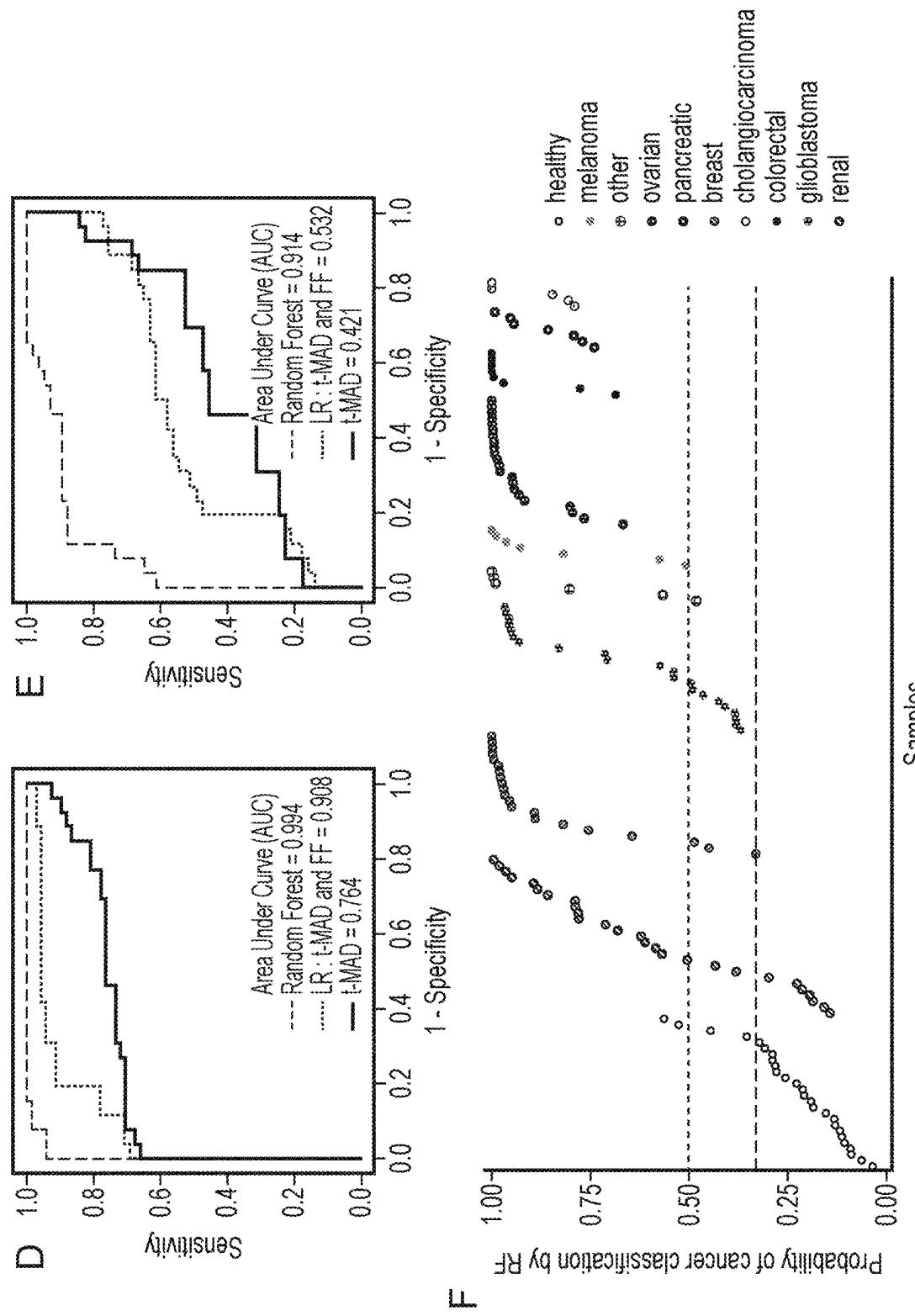


FIG. 22 (Continued)

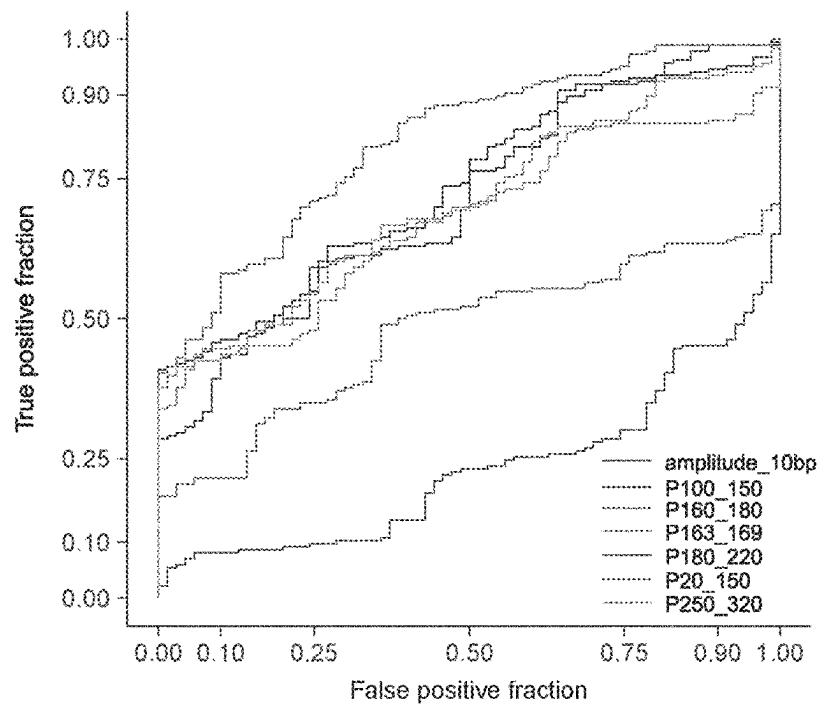


FIG. 23

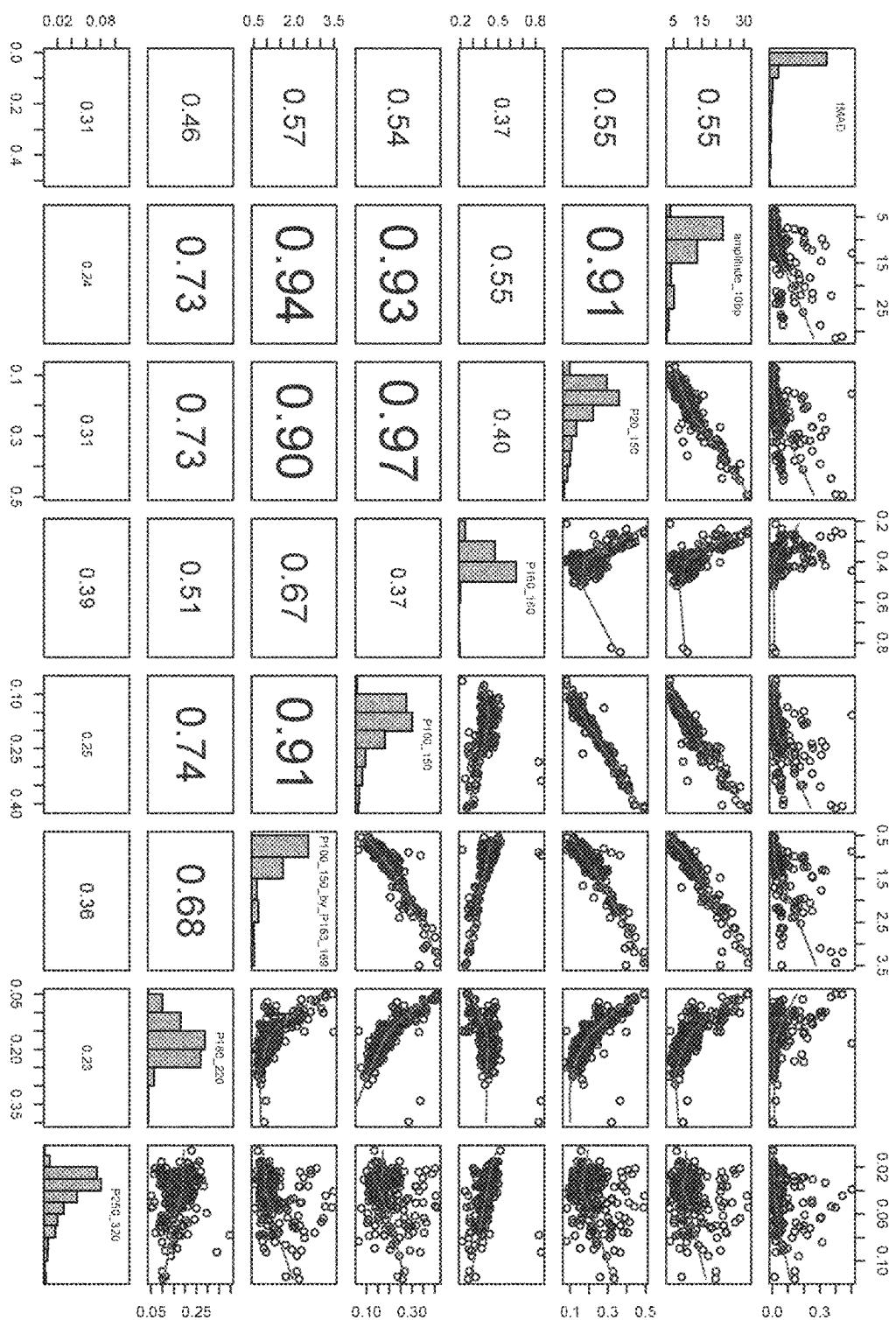


FIG. 24

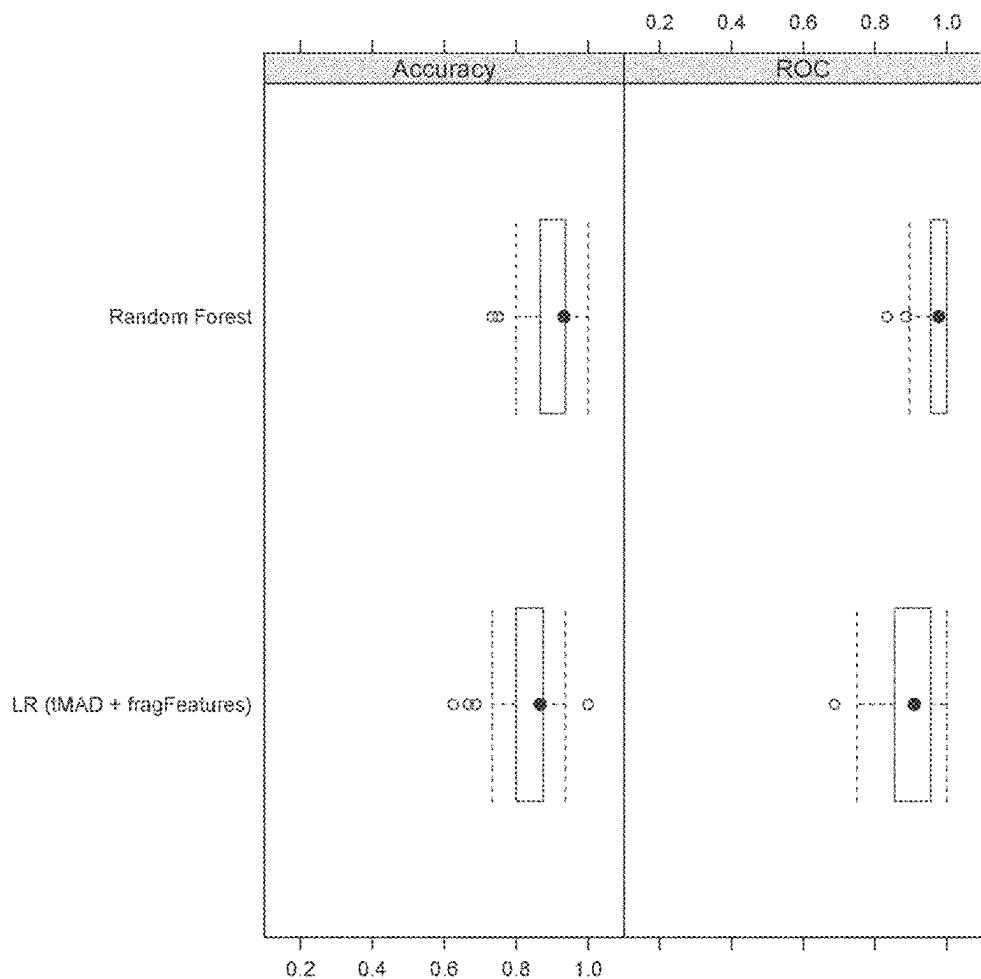


FIG. 25

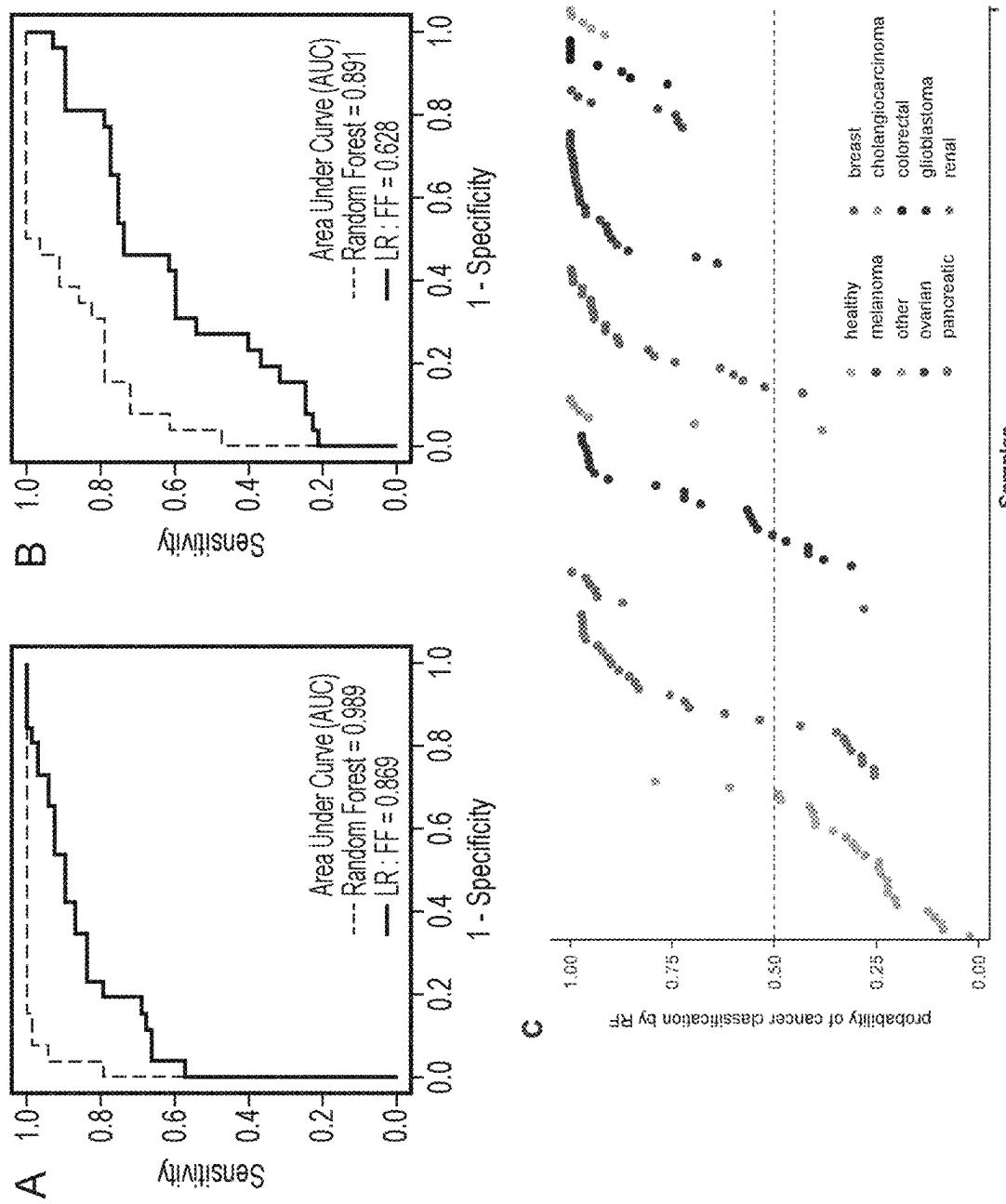


FIG. 26

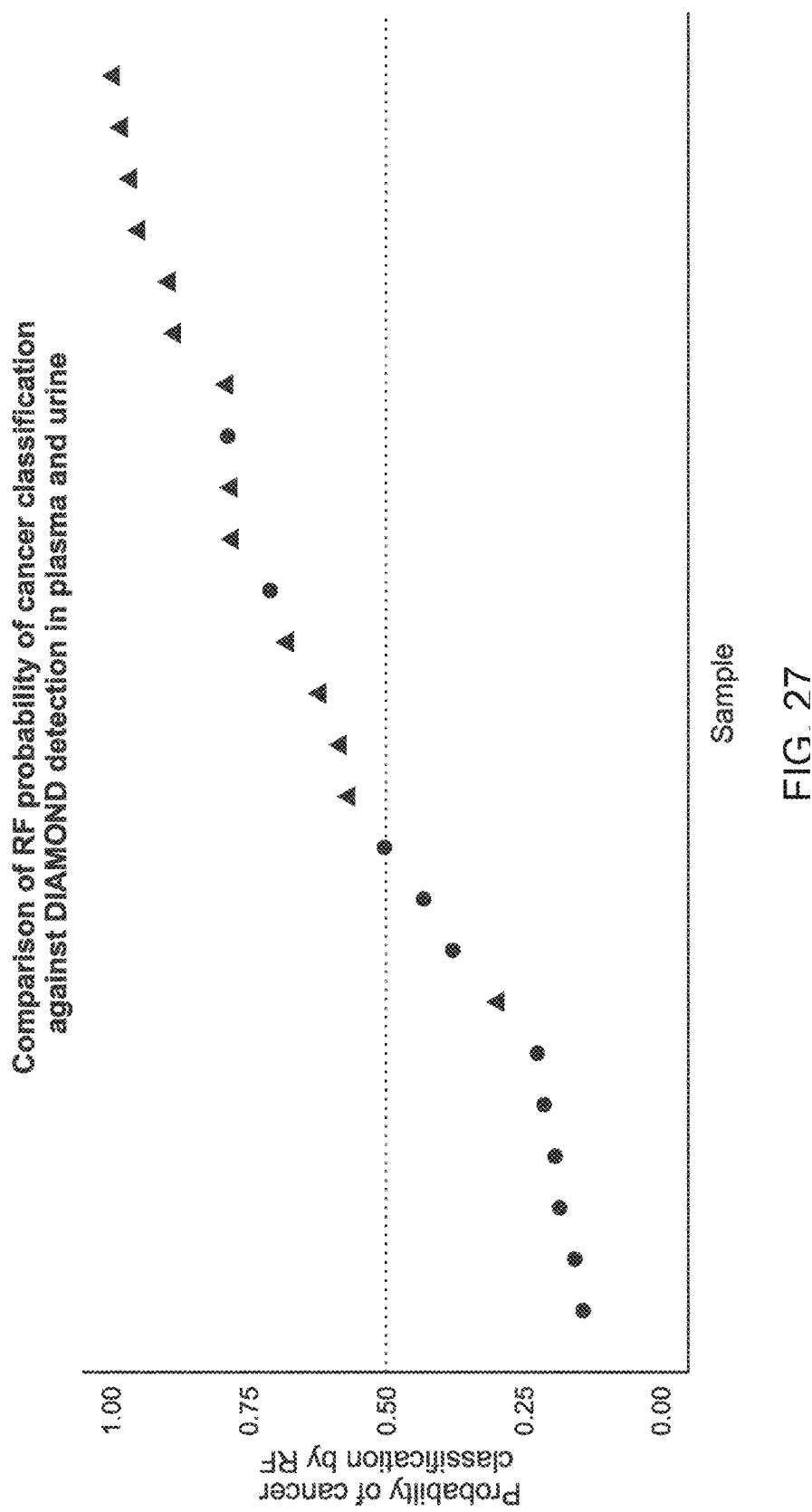


FIG. 27

1**ENHANCED DETECTION OF TARGET DNA
BY FRAGMENT SIZE ANALYSIS****RELATED APPLICATIONS**

This Application is a National Stage filing under 35 U.S.C. § 371 of International Patent Application Serial No. PCT/EP2019/080506, filed Nov. 7, 2019, which claims priority from British Application No. GB1818159.4, filed Nov. 7, 2018, each of which is incorporated by reference herein in its entirety.

FIELD OF THE INVENTION

The present invention relates in part to methods for detecting the presence of target DNA, such as circulating tumour DNA (ctDNA) from, e.g., a cell-free DNA (cfDNA) source, such as blood plasma or other biological fluid. In particular, the methods of the invention find use in the diagnosis, treatment and especially monitoring of cancer. 15

BACKGROUND TO THE INVENTION

Blood plasma of cancer patients contains circulating tumor DNA (ctDNA), but this valuable source of information is diluted by much larger quantities of DNA of non-cancerous origins: ctDNA therefore represents only a small fraction of the total cell-free DNA (cfDNA) (1, 2). High-depth targeted sequencing of selected genomic regions can be used to detect low levels of ctDNA, but broader analysis with methods such as whole exome sequencing (WES) and shallow whole genome sequencing (sWGS) are only generally informative when ctDNA levels are ~10% or greater (3-5). The concentration of ctDNA can exceed 10% of the total cfDNA in patients with advanced-stage cancers (6-8), but is much lower in patients with low tumor burden (9-12) and in patients with some cancer types such as gliomas and renal cancers (6). Current strategies to improve ctDNA detection rely on increasing depth of sequencing coupled with various error-correction methods (2, 13, 14). However, approaches that focus only on mutation analysis do not take advantage of the potential differences in chromatin organization or fragment size in ctDNA (15-17). Results of ever-deeper sequencing are also confounded by the likelihood of false positive results from detection of mutations from non-cancerous cells or clonal expansions in normal epithelia, or clonal hematopoiesis of indeterminate potential (CHIP) (13, 18, 19).

The cell of origin and the mechanism of cfDNA release into blood can mark cfDNA with specific fragmentation signatures, potentially providing precise information about cell type, gene expression, oncogenic potential or action of treatment (15, 16, 20). cfDNA fragments commonly show a prominent mode at 167 bp, suggesting release from apoptotic caspase-dependent cleavage (21-24). Circulating fetal DNA has been shown to be shorter than maternal DNA in plasma, and these size differences have been used to improve sensitivity of non-invasive prenatal diagnosis (22, 25-27). The size distribution of tumor-derived cfDNA has only been investigated in a few studies, encompassing a small number of cancer types and patients, and shows conflicting results (28-33). A limitation of previous studies is that determining the specific sizes of tumor-derived DNA fragments requires detailed characterization of matched tumor-derived alterations (30, 33), and the broader understanding and implications of potential biological differences have not previously been explored. Mouliere, Pikorz, Chan-

2

drananda, Moore et al., 2017, *BioRxiv Preprint*, doi: dx.doi.org/10.1101/134437 reports that selecting short fragments in plasma improves detection of circulating tumour DNA (ctDNA) in patients having recurrent high-grade 5 serous ovarian cancer.

While detection of ctDNA shows promise in the field of cancer care, there remains an unmet need for methods and systems that maximise signal-to-noise ratio in the context of ctDNA detection. A related problem is the need to distinguish somatic cancer mutations from mutations present in non-cancerous cells, clonal expansions of normal epithelia or CHIP. The present invention seeks to provide solutions to these needs and provides further related advantages.

BRIEF DESCRIPTION OF THE INVENTION

The present inventors hypothesised that differences in fragment lengths of circulating DNA could be exploited to enhance sensitivity for detecting the presence of ctDNA and 15 for non-invasive genomic analysis of cancer. As described in 20 detail herein, analysis of size-selected cfDNA identified clinically actionable mutations and copy number alterations 25 that were otherwise not detected. Identification of patients with advanced cancer was improved by predictive models integrating fragment length and copy number analysis of cfDNA with $AUC > 0.99$ compared to $AUC < 0.80$ without fragmentation features. Increased detection of ctDNA from 30 patients with glioma, renal and pancreatic cancer patients was achieved with $AUC > 0.91$, compared to $AUC < 0.5$ without fragmentation features. Detection of ctDNA from glioma, which does not metastasize beyond the central nervous system (CNS) has previously been reported to be very challenging (6). Fragment-size analysis and selective sequencing of specific fragment sizes can boost ctDNA 35 detection, and could be an alternative to deeper mutation sequencing for clinical applications, earlier diagnosis and to study tumor biology.

Accordingly, in a first aspect the present invention provides a computer-implemented method for detecting variant 40 nucleic acid (e.g. DNA or RNA) from a cell-free nucleic acid (e.g. DNA or RNA)-containing sample, comprising:

- a) providing data representing fragment sizes of nucleic acid fragments obtained from said sample and/or representing a measure of deviation from copy number neutrality of the nucleic acid fragments obtained from said sample;
- b) causing a processor of the computer to process the data from step a) according to a classification algorithm that has been trained on a training set comprising a plurality of samples of cell-free nucleic acid containing the variant nucleic acid and a plurality of samples not containing the variant nucleic acid, wherein said classification algorithm operates to classify sample data into one of at least two classes, the at least two classes comprising a first class containing the variant nucleic acid and a second class not containing the variant nucleic acid, based on a plurality of cell-free nucleic acid fragment size features and/or a deviation from copy number neutrality feature; and
- c) outputting the classification of the sample from step b) and thereby determining whether the sample contains the variant nucleic acid or not, or determining a probability that the sample contains the variant nucleic acid.

In some embodiments the cell-free nucleic acid-containing sample is a cell-free DNA (cfDNA)-containing sample, and wherein the variant nucleic acid is variant DNA. In particular, the variant DNA may be selected from the group

consisting of: circulating tumour DNA (ctDNA), circulating bacterial DNA, circulating pathogen DNA, circulating mitochondrial DNA, circulating foetal DNA, circulating DNA derived from a donor organ or donor tissue, circulating DNA release by a cell or tissue with an altered physiology, circulating extra chromosomal DNA, and a double minute of circular DNA. In a particularly preferred embodiment the variant DNA is ctDNA.

In some embodiments the data representing fragment sizes of the nucleic acid fragments (e.g. DNA or RNA fragments) comprise fragment sizes inferred from sequence reads, fragment sizes determined by fluorimetry, or fragment sizes determined by densitometry.

In some embodiments the present invention provides a computer-implemented method for detecting variant DNA from a cell-free DNA (cfDNA)-containing sample, comprising:

- a) providing sequence data representing fragment sizes of cfDNA fragments obtained from said sample and/or representing a measure of deviation from copy number neutrality of the cfDNA fragments obtained from said sample;
- b) causing a processor of the computer to process the sequence data from step a) according to a classification algorithm that has been trained on a training set comprising a plurality of samples of cfDNA containing the variant DNA and a plurality of samples not containing the variant DNA, wherein said classification algorithm operates to classify sample data into one of at least two classes, the at least two classes comprising a first class containing the variant DNA and a second class not containing the variant DNA, based on a plurality of cfDNA fragment size features and/or a deviation from copy number neutrality feature; and
- c) outputting the classification of the sample from step b) and thereby determining whether the sample contains the variant DNA or not, or determining a probability that the sample contains the variant DNA. As described in the Examples herein, classification algorithms can learn from cfDNA fragmentation features and somatic copy number alterations (SCNAs) analysis and improve the detection of ctDNA with a relatively low-cost and shallow sequencing approach. Moreover, the cfDNA fragmentation features and/or SCNA analysis can be leveraged to classify cancer and healthy samples with high accuracy.

In some embodiments the classification algorithm operates to classify sample data into one of said at least two, three, four, or at least five classes based on at least a plurality of cfDNA fragment size features selected from the group consisting of:

- (i) the proportion of fragments in the size range 20-150 bp (P20-150);
- (ii) the proportion of fragments in the size range 100-150 bp (P100-150);
- (iii) the proportion of fragments in the size range 160-180 bp (P160-180);
- (iv) the proportion of fragments in the size range 180-220 bp (P180-220);
- (v) the proportion of fragments in the size range 250-320 bp (P250-320);
- (vi) the ratio of the proportions P(20-150)/P(160-180);
- (vii) the ratio of the proportion P(100-150) divided by the proportion of fragment in the size range 163-169 bp;
- (viii) the ratio of the proportions P(20-150)/P180-220); and

(ix) the amplitude oscillations in fragment size density with 10 bp periodicity. It will be appreciated that the sequence data representing fragment sizes of cfDNA fragments in step a) includes the cfDNA fragment size features used by the classification algorithm.

In some embodiments the plurality of cfDNA fragment size features comprise: P(160-180), P(180-220), P(250-320) and the amplitude oscillations in fragment size density with 10 bp periodicity. As described in the Examples herein, both a linear and a non-linear machine learning algorithm independently identified the same four fragment size features P(160-180), P(180-220), P(250-320) and the amplitude oscillations in fragment size density with 10 bp periodicity, along with the SCNA feature (i.e. trimmed Median Absolute Deviation from copy number neutrality (t-MAD) score), albeit with some differences in the rank order of the features. Classification with high accuracy was obtained using only the four fragmentation features (see FIG. 26).

In some embodiments the classification algorithm operates to classify sample data into one of said at least two classes based on at least a deviation from copy number neutrality feature which is a trimmed Median Absolute Deviation from copy number neutrality (t-MAD) score or an ichorCNA feature.

ichorCNA is a tool for estimating the fraction of tumor in cell-free DNA from ultra-low-pass whole genome sequencing (ULP-WGS, 0.1x coverage). The code for ichorCNA is available at the following URL: github.com/broadinstitute/ichorCNA. ichorCNA uses a probabilistic model, implemented as a hidden Markov model (HMM), to simultaneously segment the genome, predict large-scale copy number alterations, and estimate the tumor fraction of a ultra-low-pass whole genome sequencing sample (ULP-WGS). The methodology and probabilistic model are described in: Adalsteinsson, Ha, Freeman, et al. Scalable whole-exome sequencing of cell-free DNA reveals high concordance with metastatic tumors. (2017) Nature Communications November 6; 8(1):1324. doi: 10.1038/s41467-017-00965-y (the contents of which are incorporated herein by reference). The analysis workflow consists of 2 tasks:

GC-Content Bias Correction (Using HMMcopy)
 a. Computing read coverage from ULP-WGS
 b. Data correction and normalization

CNA Prediction and Estimation of Tumor Fraction of cfDNA.

In particular, when the deviation from copy number neutrality feature comprise a t-MAD score, the score may be determined by trimming regions of genome that exhibit high copy number variability in whole genome datasets derived from healthy subjects and then calculating the median absolute deviation from $\log_2 R=0$ of the non-trimmed regions of the genome.

In some embodiments in accordance with the present invention the classification algorithm performs random forests (RF) analysis, logistic regression (LR) analysis, or support vector machine (SVM) analysis. The classification algorithm may provide an output that is a probability of correct classification, e.g., a probability that the sample in question has been classified correctly to the healthy class or cancerous class per the training set on which the classification algorithm has been trained.

In some embodiments the performance of the classification algorithm when trained on the training set is assessed by the area under the curve (AUC) value from a receiver operating characteristic (ROC) analysis. Generally the classification algorithm model showing the highest AUC value is selected as having the best performance.

In some embodiments the classification algorithm has been trained on a training set comprising at least 10, 20, 30, 40 or at least 50 samples from healthy subjects and at least 10, 20, 30, 40 or at least 50 samples from subjects known to have a cancer. In particular, the samples employed in the training set may be those shown in Table 2.

In some embodiments the sequence data provided in step a) represent whole-genome sequence (WGS) reads, Tailored Panel Sequencing (TAPAS) sequence reads, Integration of Variant Reads (INVAR) TAPAS (see co-pending patent application GB1803596.4 filed 6 Mar. 2018, incorporated herein by reference), hybrid-capture sequence reads, Tagged-Amplicon Deep Sequencing (TAm-Seq) reads, focused-exome sequence reads or whole-exome sequence reads. In particular, the sequence data provided in step a) may represent shallow whole-genome sequence (sWGS) reads, optionally 0.4× depth WGS reads.

In some embodiments the data provided in step a) represent fragment sizes of multiple nucleic acid fragments (e.g. DNA fragments) from a substantially cell-free liquid sample from a subject having or suspected as having a cancer.

In some embodiments the sequence data provided in step a) represent sequence reads of multiple DNA fragments from a substantially cell-free liquid sample from a subject having or suspected as having a cancer.

In some embodiments, the cancer may be selected from melanoma, lung cancer, cholangiocarcinoma, bladder cancer, oesophageal cancer, colorectal cancer, ovarian cancer, glioma, pancreatic cancer, renal cancer and breast cancer.

In some embodiments the sample is a plasma sample, a urine sample, a saliva sample, a cerebrospinal fluid sample, a serum sample or other nucleic acid containing (e.g. DNA-containing) biological liquid sample.

In some embodiments, wherein the variant DNA is ctDNA, the method is for detecting the presence of, growth of, prognosis of, regression of, treatment response of, or recurrence of a cancer in a subject from which the sample has been obtained.

In some embodiments the presence of ctDNA in the sample is distinguished from cfDNA containing somatic mutations of non-cancerous origin. It is specifically contemplated herein that including fragment size information on each read may enhance mutation calling algorithms from high depth sequencing so as to distinguish tumour-derived mutations from other sources of somatic variants (including clonal expansions of non-cancerous cells) or background sequencing noise. In certain embodiments the method may distinguish variant sequence reads representing clonal expansions of normal epithelia or clonal haematopoiesis of indeterminate potential (CHIP) from variant sequence reads representing ctDNA.

In certain embodiments the fragment size data provided in step a) represent sequence reads of multiple DNA fragments from a substantially cell-free liquid sample from a subject and wherein the method is for determining whether the sample contains ctDNA or contains cfDNA from CHIP. In particular, the classification algorithm may have been trained on a training set further comprising a plurality of samples of cfDNA obtained from subjects having CHIP, and wherein said at least two classes further comprise a third class containing CHIP-derived cfDNA based on a plurality of cfDNA fragment size features and/or a deviation from copy number neutrality feature.

In a second aspect the present invention provides a method for detecting variant nucleic acid from a cell-free nucleic acid-containing sample, comprising:

analysing a cell-free nucleic acid-containing sample, or a library derived from a cell-free nucleic acid-containing sample, wherein the sample has been obtained from a subject, to determine fragment sizes of nucleic acid fragments in said sample or said library; and carrying out the method of the first aspect of the invention using the fragment sizes.

In some embodiments said analysing comprises: sequencing nucleic acids from the nucleic acid-containing sample or the library and inferring fragment sizes from the sequence reads;

measuring fragment sizes of nucleic acids from the nucleic acid-containing sample or the library by fluorimetry; and/or

measuring fragment sizes of nucleic acids from the nucleic acid-containing sample or the library by densitometry.

In some embodiments the present invention provides a method for detecting variant DNA from a cell-free DNA (cfDNA)-containing sample, comprising:

sequencing a cfDNA-containing sample, or a library derived from a cfDNA-containing sample, that has been obtained from a subject to obtain a plurality of sequence reads;

processing the sequence reads to determine sequence data representing fragment sizes of cfDNA fragments obtained from said sample and/or representing a measure of deviation from copy number neutrality of the cfDNA fragments obtained from said sample; and carrying out the method of the first aspect of the invention using the sequence data.

In some embodiments the sequencing comprises generating a sequencing library from the sample and performing whole-genome sequencing, Tailored Panel Sequencing (TAPAS) sequencing, hybrid-capture sequencing, TAm-Seq sequencing, focused-exome sequencing or whole-exome sequencing, optionally generating an indexed sequencing library and performing shallow whole genome sequencing (e.g. to a depth of 0.4×).

In some embodiments processing the sequence reads comprises one or more of the following steps:

removal of contaminating adapter sequences;

removal of PCR and optical duplicates;

removal of sequence reads of low mapping quality; and if multiplex sequencing, de-multiplexing by excluding mismatches in sequencing barcodes.

In some embodiments the variant DNA is selected from the group consisting of: circulating tumour DNA (ctDNA), circulating bacterial DNA, circulating pathogen DNA, circulating mitochondrial DNA, circulating foetal DNA, and circulating DNA derived from a donor organ or donor tissue, circulating DNA release by a cell or tissue with an altered physiology, circulating extra chromosomal DNA, and a double minute of circular DNA.

In some embodiments processing the sequence reads to determine sequence data representing fragment sizes of cfDNA fragments obtained from said sample and/or representing a measure of deviation from copy number neutrality of the cfDNA fragments obtained from said sample comprises determining one or more (e.g. 2, 3, 4, 5 or more) features selected from the group consisting of:

(i) the proportion of fragments in the size range 20-150 bp (P20-150);

- (ii) the proportion of fragments in the size range 100-150 bp (P100-150);
- (iii) the proportion of fragments in the size range 160-180 bp (P160-180);
- (iv) the proportion of fragments in the size range 180-220 bp (P180-220);
- (v) the proportion of fragments in the size range 250-320 bp (P250-320);
- (vi) the ratio of the proportions P(20-150)/P(160-180);
- (vii) the ratio of the proportion P(100-150) divided by the proportion of fragment in the size range 163-169 bp;
- (viii) the ratio of the proportions P(20-150)/P(180-220); and
- (ix) the amplitude oscillations in fragment size density with 10 bp periodicity.

In some embodiments the plurality of cfDNA fragment size features comprise: P(160-180), P(180-220), P(250-320) and the amplitude oscillations in fragment size density with 10 bp periodicity.

In some embodiments the fragment sizes of cfDNA fragments are inferred from sequence reads using the mapping locations of the read ends in the genome following alignment of the sequence reads with the reference genome of the species from which the sample was obtained.

In some embodiments processing the sequence reads to determine sequence data representing a measure of deviation from copy number neutrality of the cfDNA fragments obtained from said sample comprises determining a trimmed Median Absolute Deviation from copy number neutrality (t-MAD) score or an ichorCNA score. In particular, the t-MAD score may be determined by trimming regions of genome that exhibit high copy number variability in whole genome datasets derived from healthy subjects and then calculating the median absolute deviation from $\log_2 R=0$ of the non-trimmed regions of the genome.

In some embodiments the sample contains multiple DNA fragments from a substantially cell-free liquid from a subject having or suspected as having a cancer. In particular cases, the cancer may be selected from melanoma, lung cancer, cholangiocarcinoma, bladder cancer, oesophageal cancer, colorectal cancer, ovarian cancer, glioma, pancreatic cancer, renal cancer and breast cancer.

In some embodiments the sample is a plasma sample, a urine sample, a saliva sample, a cerebrospinal fluid sample, a serum sample or other DNA-containing biological liquid sample.

In accordance with any aspect of the present invention the sample may be or may have been subjected to one or more processing steps to remove whole cells, for example by centrifugation.

In certain embodiments, wherein the variant DNA is ctDNA, the method may be for detecting the presence of, growth of, prognosis of, regression of, treatment response of, or recurrence of a cancer in a subject from which the sample has been obtained.

In some embodiments the presence of ctDNA is distinguished from the presence of cfDNA containing somatic mutations of non-cancerous origin, optionally from CHIP origin.

In some embodiments a somatic mutation containing cfDNA fragment is classified as being of tumour origin or being of CHIP origin based on a plurality of fragment size features determined from the sequence reads.

In some embodiments the variant DNA is ctDNA and the classification of the sample as containing ctDNA or not, or the determined probability that the sample contains ctDNA

is used to predict whether said sample or a further sample from the same subject will be susceptible to further ctDNA analysis.

In some cases the further ctDNA analysis comprises sequencing to a greater sequencing depth and/or targeted sequencing of ctDNA in said sample.

In some embodiments, when the probability that the sample contains ctDNA as determined by the classification algorithm is at least 0.5 (e.g. at least 0.6 or at least 0.75), the sample is subjected to said further ctDNA analysis.

In some embodiments:

said sample is a plasma sample and the probability that the sample contains ctDNA as determined by the classification algorithm is used to determine whether ctDNA will be detectable in a urine sample; or

said sample is a urine sample and wherein the probability that the sample contains ctDNA as determined by the classification algorithm is used to determine whether ctDNA will be detectable in a plasma sample. As shown in Example 8, a relatively high probability shown by the classification algorithm that a plasma sample contains ctDNA was associated with an increased probability that useful detection of ctDNA was possible with a urine sample (see also FIG. 27).

In a third aspect the present invention provides a method for improving the detection of circulating tumour DNA (ctDNA) in a cell-free DNA (cfDNA) containing sample, comprising performing an in vitro and/or in silico size selection to enrich for DNA fragments of less than 167 bp in length and/or to enrich for DNA fragments in the size range 250 to 320 bp. In some embodiments the size selection is to enrich for DNA fragments in the range 90 to 150 bp in length. In some cases the size selection may comprise excluding high molecular weight DNA such as that derived from white blood cells when the sample comprises a serum sample.

In some embodiments the sample may have been obtained from a subject having or suspected as having a cancer selected from the group consisting of melanoma, cholangiocarcinoma, colorectal cancer, glioma, pancreatic cancer, renal cancer and breast cancer.

In some embodiments the size selection comprises an in vitro size selection that is performed on DNA extracted from a cfDNA containing sample and/or is performed on a library created from DNA extracted from a cfDNA containing sample. In particular, the in vitro size selection may comprise agarose gel electrophoresis.

In some embodiments the size selection comprises an in silico size selection that is performed on sequence reads.

In particular cases the sequence reads may comprise paired-end reads generated by sequencing DNA from both ends of the fragments present in a library generated from the cfDNA containing sample. The original length of the DNA fragments in the cfDNA containing sample may be inferred using the mapping locations of the read ends in the genome following alignment of the sequence reads with the reference genome of the species from which the sample was obtained (e.g. the human reference genome GRCh37 for a human subject).

In some embodiments DNA fragments outside the range 90 to 150 bp in length are substantially excluded (see, e.g., FIG. 6B).

In some embodiments the size selection is performed on a genome wide basis or an exome wide basis. As described herein, the present inventors identified size differences between mutant and non-mutant cfDNA on a genome-wide

and pan-cancer scale in contrast to previous studies that were limited to specific genomic loci, cancer types or cases (30, 32, 33).

In certain embodiments the in vitro size selection is performed prior to shallow whole genome sequencing (sWGS) or the in silico size selection is performed on sWGS sequencing reads.

In certain embodiments the method further comprises performing somatic copy number aberration analysis and/or mutation calling on the sequence reads subsequent to the size selection. In particular cases somatic copy number aberration analysis may comprise processing the sequence reads to determine a trimmed Median Absolute Deviation from copy number neutrality (t-MAD) score or an iChorCNA score. For example, the t-MAD score may be determined by trimming regions of genome that exhibit high copy number variability in whole genome datasets derived from healthy subjects and then calculating the median absolute deviation from $\log_2 R=0$ of the non-trimmed regions of the genome.

In certain embodiments somatic copy number aberration analysis may comprise detecting amplifications in one or more genes selected from NF1, TERT, and MYC. As described in the Examples herein, analysis of plasma cfDNA after size selection revealed a large number of SCNA that were not observed in the same samples without size selection.

In certain embodiments mutation calling comprises detecting mutations in one or more genes selected from BRAF, ARID1A, and NF1. As described in the Examples herein, size selection enriched the mutant allele fraction (MAF) for nearly all mutations.

In some embodiments the cancer is a high ctDNA cancer selected from the group consisting of: colorectal, cholangiocarcinoma, breast and melanoma.

In some embodiments the cancer is a low ctDNA cancer selected from the group consisting of: pancreatic cancer, renal cancer and glioma.

In certain embodiments the sample may be a plasma sample, a urine sample, a saliva sample, a cerebrospinal fluid sample, a serum sample or other DNA-containing biological liquid sample.

In some embodiments the method further comprises detecting the presence of, growth of, prognosis of, regression of, treatment response of, or recurrence of a cancer in a subject from which the sample has been obtained. Improving the detection of ctDNA, mutation calling and/or SCNA detection in accordance with the methods of this aspect of the invention may assist with the early detection of cancer and with ongoing cancer monitoring, and may inform treatment strategies.

In some embodiments the method may be carried out on a sample obtained prior to a cancer treatment of the subject and on a sample obtained following the cancer treatment of the subject. As described herein, size selected samples indicated tumour progression 69 and 87 days before detection by imaging or non-size selected t-MAD analysis (see FIGS. 10E and F).

In accordance with any aspect of the present invention, the subject may be a human, a companion animal (e.g. a dog or cat), a laboratory animal (e.g. a mouse, rat, rabbit, pig or non-human primate), a domestic or farm animal (e.g. a pig, cow, horse or sheep).

Preferably, the subject is a human patient. In some cases, the subject is a human patient who has been diagnosed with, is suspected of having or has been classified as at risk of developing, a cancer.

Embodiments of the present invention will now be described by way of example and not limitation with reference to the accompanying figures. However various further aspects and embodiments of the present invention will be apparent to those skilled in the art in view of the present disclosure.

The present invention includes the combination of the aspects and preferred features described except where such a combination is clearly impermissible or is stated to be expressly avoided. These and further aspects and embodiments of the invention are described in further detail below and with reference to the accompanying examples and figures.

15 BRIEF DESCRIPTION OF THE FIGURES

FIG. 1 shows a flowchart summarizing the different experiments done in this study and the corresponding samples numbers used at each step.

FIG. 2 shows a survey of plasma DNA fragmentation with genome-wide sequencing on a pan-cancer scale. A, The size profile of cfDNA can be determined from paired-end sequencing of plasma samples and reflects its organization around the nucleosome. cfDNA is released in the blood circulation by various means, each of which leaves a signature on the fragment sizes. The size profile of cfDNA was inferred by analyzing with sWGS ($n=344$ plasma samples from 65 healthy controls and 200 cancer patients), and the size profile of mutant ctDNA by personalized capture sequencing ($n=18$ plasma samples). B, Fragment size distributions of 344 plasma samples from 200 cancer patients. Patients are split into two groups based on previous literature (3), cancer samples previously observed to have low levels of ctDNA (renal, bladder, pancreatic, and glioma) and 30 cancer samples observed to have higher ctDNA levels (breast, melanoma, ovarian, lung, colorectal, cholangiocarcinoma, and others, see Table 1). C, Proportion of cfDNA fragments below 150 bp by cancer grouping defined in B. The Kruskal-Wallis test for difference in size distributions 35 indicated a significant difference between the group of cancer types releasing high amounts of ctDNA, and the group releasing low amounts as well as the group of healthy individuals ($p<0.001$). D, Proportion of cfDNA fragments below 150 bp by cancer type (all samples). Cancer types represented by fewer than 4 individuals are grouped in the “other” category. The line indicates the median proportion per cancer type.

FIG. 3 shows the size distribution of cfDNA for all plasma samples of healthy individuals and cancer patients included in this study depending on their cancer type, determined by sWGS. The plasma samples showed here were collected from renal cancer ($n=33$), glioblastoma ($n=11$), bladder cancer ($n=19$), breast cancer ($n=34$), melanoma ($n=21$), pancreatic ($n=7$), ovarian ($n=59$), lung ($n=8$), colorectal ($n=21$), cholangiocarcinoma ($n=14$), cervical ($n=1$), penile ($n=1$), endometrial ($n=1$), thymoma ($n=1$), hepatocellular carcinoma ($n=1$). The size profile of cfDNA from healthy individuals ($n=46$) is also shown.

FIG. 4 depicts the determination of the size profile of mutant ctDNA with animal models and personalized capture sequencing. A, A mouse model with xenografted human tumor cells enabled the discrimination of DNA fragments released by cancer cells (reads aligning to the human genome) from the DNA released by healthy cells (reads aligning to the mouse genome), with the use of sWGS. B, Fragment size distribution, from the plasma extracted from a mouse xenografted with a human ovarian tumor, showing

ctDNA originating from tumor cells and cfDNA from non-cancerous cells. Two vertical lines indicate 145 bp and 167 bp. The fraction of reads shorter than 150 bp is indicated. C, Design of personalized hybrid-capture sequencing panels developed to specifically determine the size profiles of mutant DNA and non-mutant DNA in plasma from 19 patients with late stage cancers. Capture panels included somatic mutations identified in tumor tissue by WES. A mean of 165 mutations per patient were then analyzed from matched plasma samples. Reads were aligned and separated into fragments that carry either the reference or the mutant sequence. Fragment sizes for paired-end reads were calculated. D, Size profiles of mutant DNA and non-mutant DNA in plasma from 19 patients with late stage cancers were determined by tumor-guided capture sequencing. The fraction of reads shorter than 150 bp is indicated.

FIG. 5 shows the insert size distribution determined with hybrid-capture sequencing for 19 patients included in the mutant DNA size distribution analysis (A-S). The size distribution of mutant DNA fragments is shown in red and the distribution of non-tumour reference cfDNA from the same sample is shown in grey. The vertical dashed lines represent 145 bp and 167 bp. The insert sizes were determined by aggregating the insert sizes observed from mutant DNA and reference DNA of all samples for each patient.

FIG. 6 shows the enhancement of the tumor fraction from plasma sequencing with size selection. A, Plasma samples collected from ovarian cancer patients were analyzed in parallel without size selection, or using either in silico and in vitro size selection. B, Accuracy of the in vitro and in silico size selection determined on a cohort of 20 healthy controls. Shows the size distribution before size selection, after in silico size selection (with sharp cutoff at 90 and 150 bp) and after in vitro size selection. C, SCNA analysis with sWGS from plasma DNA of an ovarian cancer patient collected before initiation of treatment, when ctDNA MAF was 0.271 for a TP53 mutation as determined by TAm-Seq. Shows inferred amplifications and deletions. Copy number neutral regions are in grey. D, SCNA analysis of a plasma sample from the same patient as panel C collected three weeks after treatment start. The MAF for the TP53 mutation was 0.068, and ctDNA was not detected at this time-point by sWGS (before size selection). E, Analysis of the same plasma sample as D after in vitro size selection of fragments between 90 bp and 150 bp in length. The MAF for the TP53 mutation increased to 0.402 after in vitro size selection, and SCNAs were clearly apparent by sWGS. More SCNAs are detected in comparison to C and D (e.g. in chr2, chr9, chr10).

FIG. 7 shows the distribution of insert sizes determined with sWGS for each plasma sample from the 13 ovarian patients of the OV04 cohort, collected before and after treatment. The distribution of cell-free DNA (cfDNA) without size selection is shown and the distribution of the same cfDNA samples after size selection is shown. The vertical lines represent the range of fragments selected with the PippinHT cassettes, between 90 and 150 bp. To note that patient OV04-292 and OV04-300 exhibit an altered fragmentation profile indicating a possible issue with the preparation or pre-analytical preservation of the samples.

FIG. 8 shows the quality control assessment of the in vitro size selection, estimated with sWGS and targeted sequencing. A, Size distribution of DNA fragments from the plasma samples included in the size selection study, assessed by sWGS, before size-selection and after in vitro size-selection. The two dotted vertical lines indicate the size selection range between 90 bp and 150 bp. B, Proportion of non-reference

allele fractions corresponding to the sequencing background noise as determined during targeted sequencing (TAm-Seq) of plasma DNA sample from ovarian cancer patients, with and without in vitro size selection.

FIG. 9 shows the second quality control assessment of the in vitro and in silico size selection. 20 plasmas were selected from healthy controls, extracted DNA and performed sWGS without size selection, with in vitro and in silico size selection on these samples. A, The size profile determined for each samples and condition. B, There was an increase in the fraction of duplicated reads, and therefore these were removed for any downstream size selection analysis. In order to determine if the size selection could introduce more sequencing noise during the analysis, a QC metric called the median absolute pairwise difference (MAPD) algorithm was used to find the sequencing noise. MAPD measured the absolute difference between the log 2 CN ratios of every pair of neighboring bins and found the median across all bins. Higher MAPD scores reflected greater noise, typically associated with poor-quality samples. All samples exhibited a MAPD score of 0.01 (+/-0.01), irrespective of the size selection condition. C, In addition to the noise estimation the ctDNA fraction between the 20 controls samples as estimated by the t-MAD score were compared. The t-MAD score from the samples without size selected was not significant different with the t-MAD determined after in silico size selection (t-test, p=0.43), but a significant difference with the samples after in vitro size selection (t-test, p=0.0068) was observed. Even if the t-MAD value was increased after in vitro size selection, the mean (0.011) and the maxima (0.016) detected were still constrained in the threshold limit determined empirically from the whole cohort of controls (n=65). D, The yield of DNA recovered after in vitro size selection was determined (as in silico size selection is not affected by this technical bias).

FIG. 10 shows the quantification of the ctDNA enrichment by sWGS with in silico size selection and t-MAD. A, Workflow to quantify tumor fraction from SCNA as a genome-wide score named t-MAD. B, Correlation between the MAF of SNVs determined by digital PCR or hybrid-capture sequencing and t-MAD score determined by sWGS. Data included 97 samples from cancer patients of multiples cancer types with matched MAF measurements and t-MAD scores. Pearson correlation (coefficient r) between MAF and t-MAD scores was calculated for all cases with MAF>0.025 and t-MAD>0.015. Linear regression indicated a fit with a slope of 0.44 (solid line). C, Comparison of t-MAD scores determined from sWGS between healthy samples, samples collected from patients with cancer types that exhibited low amounts of ctDNA in circulation and from patients with cancer types that exhibited high amounts of ctDNA in circulation. All samples for which t-MAD could be calculated have been included. D, ROC analysis comparing the classification of these plasma samples from high ctDNA cancer samples (n=189) and plasma samples from healthy controls (n=65) using t-MAD had an area under curve (AUC) of 0.69 without size selection (black solid curve). After applying in silico size selection to the samples from the cancer patients, we observed an AUC of 0.90 (black dashed curve). E, Determination of t-MAD from longitudinal plasma samples of a colorectal cancer patient. t-MAD was analyzed before and after in silico size selection of the DNA fragments 90-150 bp, and then compared to the RECIST status for this patient. F, Application of in silico size selection to 6 patients with long follow-up. t-MAD score was determined before and after in silico size selection of the short DNA fragments. Dark circles indicate samples

in which ctDNA was detected both with and without in silico size selection. Light circles indicate samples where ctDNA was detected only after in silico size selection. Empty circles indicate samples where ctDNA was not detected by either analysis. Times when RECIST status was assessed are indicated by a bar for progression, or a bar for regression or stable disease.

FIG. 11 shows a comparison of the MAF and t-MAD score depending on the cancer type for available matched data. Data from ovarian, breast, cholangiocarcinoma, colorectal and lung are detailed. Other cancer types are grouped in the category “other”. Samples are labelled depending on their t-MAD score, with t-MAD<0.015, and t-MAD>0.015. Pearson correlations, p values and slopes are indicated when n>5 and t-MAD>0.015.

FIG. 12 shows plasma DNA from a breast cancer patient, which was spiked into pooled plasma DNA derived from healthy individual. This was serially diluted in steps of 10-, 100- and 1000-fold. A total of 10ng of DNA was used for the initial DNA library preparation. The allele fraction for a TP53 mutation of the neat sample was estimated by both WES and TAm-Seq to be ~45.6%, and was used as the reference for the dilution. In the dilution series data, the t-MAD score appears to detect SCNA with very low coverage and mutant AF (down to ~0.4% AF, or 100x diluted sample). In addition the sequencing data has been in silico size selected for the short fragments (90-150 bp), improving the t-MAD score for the lower AF.

FIG. 13 shows a comparison of the available RECIST volume (in mm) determined by CT-scan to the tMAD score and fragmentation features. The RECIST volume was compared to the tMAD score (A), the proportion of fragments between 20 and 150 bp (B), the ratio of the proportion of fragments between 100-150 bp and the proportion of fragments between 163-169 bp (C), the ratio of the proportion of fragments between 20-150 bp and the proportion of fragments between 180-220 bp (D), the statistic amplitude of the 10 bp peaks and valleys (E), and the proportion of fragments between 250-350 bp (F). Correlation and p values are calculated for each comparison.

FIG. 14 shows the quantification of the ctDNA enrichment by sWGS with in vitro size selection. A, The effect of in vitro size selection on the t-MAD score. For each of 48 plasma samples collected from 35 patients, the t-MAD score was determined from the sWGS after in vitro size selection (y axis) and without size selection (x axis). In vitro size selection increased the t-MAD score for nearly all samples, with a median increase of 2.1-fold (range from 1.1 to 6.4 fold). t-MAD scores determined from sWGS for 46 samples from healthy individuals were all <0.015 both before and after in vitro size selection. B, ROC analysis comparing the classification of these plasma samples from cancer samples (n=48) and plasma samples from healthy controls (n=46) using t-MAD had an area under curve (AUC) of 0.64 without size selection. After applying in silico size selection to the samples from the cancerous and healthy patients, an AUC of 0.78 was observed, and after in vitro size selection, an AUC of 0.97. C, Comparison of t-MAD scores determined from sWGS between matched ovarian cancer samples with and without in vitro size selection. The t-test for the difference in means indicate a significant increase in tumor fraction (measured by t-MAD) with in vitro size selection ($p<0.0001$). D, Detection of SCNA across 15 genes frequently mutated in recurrent ovarian cancer, measured in plasma samples collected during treatment for 35 patients. Patients were ranked from left to right by increasing tumor fraction as quantified by tMAD (before in vitro size selec-

tion). SCNA are labelled as detected for a gene if the relative copy number in that region was greater than 0.05. Empty squares represent copy number neutral regions, bottom left triangles indicate that SCNA were detected without size selection and top right triangles in represent SCNA detected after in vitro size selection.

FIG. 15 shows the analysis of each of the 48 plasma samples collected from 35 ovarian patients with and without size selection. A, There is a negative correlation between the ctDNA fraction represented by the t-MAD score, and the level of enrichment (Pearson, -0.49 , $p<0.001$. B, The t-MAD score determined from the sWGS with size selection was higher than without size selection for nearly all samples, with a median increase of 2.1-fold. The enrichment factor with size selection, determined by t-MAD, varied per sample but was higher for samples with low initial t-MAD score. Values from healthy individuals are added for comparison purposes.

FIG. 16 shows the SCNA analysis of the segmental log 2 ratio determined after sWGS. This was performed using a list of 29 genes frequently mutated in recurrent ovarian cancer from the plasma samples collected at baseline and after treatment for 13 patients. The log 2 ratio are represented for the samples without size selection and with in vitro size selection of the shorter DNA.

FIG. 17 shows the improvement in the detection of somatic alterations by WES in multiple cancer types with size selection. A, Analysis of the MAF of mutations detected by WES in 6 patients with HGSOC without size selection and with in vitro and in silico size selection. B, Comparison of size-selected WES data with non-selected WES data to assess the number of mutations detected in plasma samples from 6 patients with HGSOC. For each patient, the first bar shows the number of mutations called without size selection, the second bar quantifies the number of mutations called after the addition of those identified with in silico size selection, and the third bar shows the number of mutations called after addition of mutations called after in vitro size selection. C, Patients (n=16) were retrospectively selected from a cohort with different cancer types (colorectal, cholangiocarcinoma, pancreatic, prostate) enrolled in early phase clinical trials. Matched tumor tissue DNA was available for each plasma sample, and 2 patients also had a biopsy collected at relapse. WES was performed on tumor tissue DNA and plasma DNA samples, and in silico size selection was applied to the data. 2061/2133, 97% of the shared mutations detected by WES showed higher MAF after in silico size selection. D, Mutations detected only after in silico selection of WES data from 16 patients (as in C) compared to mutations called by WES of the matched tumor tissue. Three of 16 patients had no additional mutations identified after in silico size selection. Of the 82 mutations detected in plasma after in silico size selection, 23 (28%) had low signal levels in tumor WES data and were not initially identified in those samples.

FIG. 18 shows the Mutant allelic fraction (MAF) for each single nucleotide variants (SNVs) called by WES on the OV04 samples without and with size-selection. A, The MAF determined by WES with in vitro size selection (vertical) was higher than without in vitro size selection (horizontal) for most of the mutations detected from the plasma samples of 6 HGSOC patients. B, Enrichment is also observed in the same samples after in silico size selection from WES data.

FIG. 19 depicts the mutations detected for 9 genes of clinical importance by WES with and without size selection of the short DNA fragments. All the plasma samples submitted to WES (6 ovarian cancer cases from OV04 study,

and 16 cancers from the CoPPO study) were analysed. Mutations called by without size selection were integrated, and also the new mutations called by WES after in vitro and in-silico size selection.

FIG. 20 shows A, The MAF for TP53 mutations determined by TAm-Seq with in vitro size selection was higher than without size selection for most samples, including samples collected at baseline (circles) and after initiation of treatment (triangles). Only the 26 samples collected from 13 patients with a sample collected before and after treatment are shown. The dotted area highlights samples which had initially low MAF (<5%), where methods such as whole-exome sequencing (at sequencing depth of ~100x) would not be effective, and where in vitro size selection enriched the MAF to >5% and therefore accessible for wide-scale analysis. B, Comparison of the MAF detected by TAm-Seq before treatment and after initiation of treatment with in vitro size selection (triangles) and without size selection (circles).

FIG. 21 shows the size distribution of mutant and non-mutant DNA obtained from the personalised sequencing. A fraction of 10 patients from this figure were sub-selected. The loci selected corresponded to clinically validated variants (based on the WES of the tumor tissue DNA). The left panel exhibit the size distribution of mutant DNA, and the right panel the size distribution of the corresponding non-mutant DNA. The mutant ctDNA confirm enrichment in the size range 90-150 bp (as previously described in the manuscript). The non-mutant exhibited a lower enrichment in the size range 90-150 bp, but with variations depending on the patient. The patient with the highest concentration of ctDNA as determined by t-MAD, had an enrichment in shorter non-mutant DNA, whereas the patients with a lower value of t-MAD, have less short fragments. This suggests that even in the non-mutant DNA, tumor signal (=non-mutant ctDNA) can be detected by analysing the size of the cfDNA fragments.

FIG. 22 depicts enhancing the potential for ctDNA detection by combining SCNAs and fragment-size features. A, Schematic illustrating the selection of different size ranges and features in the distribution of fragment sizes. For each sample, fragmentation features included the proportion (P) of fragments in specific size ranges, the ratio between certain ranges and a quantification of the amplitude of the 10 bp oscillations in the 90-145 size bp range calculated from the periodic “peaks” and “valleys”. B, Principal Component Analysis (PCA) comparing cancer and healthy samples using data from t-MAD scores and the fragmentation features. Fragmentation features shown in grey are not included in the following steps. C, Workflow for the predictive analysis combining SCNAs and fragment size features. Plasma DNA sWGS data from healthy controls was split into a training set (60% of samples) and a validation set (used in both Validation data 1 and Validation set 2). sWGS data from plasma samples from a pan-cancer cohort of 182 samples from patients with cancer types with high levels of ctDNA (colorectal, cholangiocarcinoma, lung, ovarian, breast) was split into a training set (60% of samples) and a validation set (Validation data 1, together with the healthy individual validation set). A further dataset of sWGS from 57 samples from cancer types exhibiting low levels of ctDNA (glioma, renal, pancreatic) was used as Validation data 2, together with the healthy individual validation set. D, ROC curves for Validation data 1 (samples from cancer patients with high ctDNA levels=68, healthy=26) for 3 predictive models built on the pan-cancer training cohort (cancer=114, healthy=39). The curve represents the ROC

curve for classification with t-MAD only, the long dashed line represents the logistic regression model combining the top 5 features based on recursive feature elimination (t-MAD score, 10 bp amplitude, P(160-180), P(180-220) and P(250-320)), and the dashed line shows the result for a random forest classifier trained on the combination of the same 5 features, independently chosen for the best RF predictive model. E, ROC curves for Validation data 2 (samples from cancer patients with low ctDNA levels=57, healthy=26) for the same 3 classifiers as D. The curve represents the model using t-MAD only, the long-dashed represents the logistic regression model combining the top 5 features (t-MAD score, 10 bp amplitude, P(160-180), P(180-220), and P(250-320)), and the dashed shows the result for a random forest classifier trained on the combination of same 5 predictive features. F, Plot representing the probability of classification as cancer with the RF model for all samples in both validation datasets. Samples are separated by cancer type and sorted within each by the RF probability of classification as cancer. The dashed horizontal line indicates 50% probability and the light long-dashed line indicates 33% probability.

FIG. 23 shows the ROC analysis of the cfDNA fragmentation features between healthy samples and samples from patients with high ctDNA cancers.

FIG. 24 shows a comparison of t-MAD score to the 9 fragmentation features determined by sWGS from the 147 plasma samples from cancer patients included in the training and validation dataset of the classifier models. The correlation score was estimated for each cross-comparison, and the value displayed on the bottom left side of the figure.

FIG. 25 shows the performance metrics for the different algorithms: logistic regression (on t-MAD score and the fragmentation features), and random forest (RF) on training set data from sWGS ($n=153$; 114 cancer samples, and 39 healthy controls). The median ROC score and accuracy values are displayed for each models, as well as the 0.95 confidence level.

FIG. 26 shows LR and RF models, which detect cancer from healthy samples with the fragmentation features alone. A, ROC curves from the first validation sample set (cancer=68, healthy=26) for 2 classifiers built on the pan-cancer training cohort (cancer=114, healthy=39). The curve represents the ROC for a logistic regression model trained only with the fragmentation features without t-MAD and the dashed curve shows the result for a random forest classifier trained on the combination of the best 3 predictive fragmentation features (amplitude_10 bp, P(160-180), and P(250-320)). B, ROC curves from the second validation sample set (cancer=57, healthy=26) for 2 classifiers built on the same training set as A. The curve represents the logistic regression model trained only with the fragmentation features and the dashed curve shows the result for a random forest classifier trained on the combination of 3 predictive features (amplitude_10 bp, P(160-180), and P(250-320)). C, plot representing the probability of classification as cancer with the RF model for the second validation dataset (described in B). Samples are ranked by cancer-type and by probability of classification as cancer. The dashed horizontal line represents the 50% probability.

FIG. 27 shows the probability of cancer classification by the random forest (RF) model, for a given renal cell carcinoma (RCC) patient plasma sample, as indicated on the y-axis. Patient plasma samples are indicated on the x-axis. For each patient, this same plasma sample (and in some cases matched urine supernatant) were assessed for ctDNA content by INVAR-TAPAS and t-MAD analysis. Circles

indicate patients in which ctDNA was not detected in either fluid by either approach. Triangles indicate patients in which ctDNA was detected in either fluid by either method.

DETAILED DESCRIPTION OF THE INVENTION

Aspects and embodiments of the present invention will now be discussed with reference to the accompanying figures. Further aspects and embodiments will be apparent to those skilled in the art. All documents mentioned in this text are incorporated herein by reference.

In describing the present invention, the following terms will be employed, and are intended to be defined as indicated below.

“Computer-implemented method” where used herein is to be taken as meaning a method whose implementation involves the use of a computer, computer network or other programmable apparatus, wherein one or more features of the method are realised wholly or partly by means of a computer program.

A “sample” as used herein may be a biological sample, such as a cell-free DNA sample, a cell (including a circulating tumour cell) or tissue sample (e.g. a biopsy), a biological fluid, an extract (e.g. a protein or DNA extract obtained from the subject). In particular, the sample may be a tumour sample, a biological fluid sample containing DNA, a blood sample (including plasma or serum sample), a urine sample, a cervical smear, a cerebrospinal fluid sample, or a non-tumour tissue sample. It has been found that urine and cervical smears contain cells, and so may provide a suitable sample for use in accordance with the present invention. Other sample types suitable for use in accordance with the present invention include fine needle aspirates, lymph nodes, surgical margins, bone marrow or other tissue from a tumour microenvironment, where traces of tumour DNA may be found or expected to be found. The sample may be one which has been freshly obtained from the subject (e.g. a blood draw) or may be one which has been processed and/or stored prior to making a determination (e.g. frozen, fixed or subjected to one or more purification, enrichment or extractions steps, including centrifugation). The sample may be derived from one or more of the above biological samples via a process of enrichment or amplification. For example, the sample may comprise a DNA library generated from the biological sample and may optionally be a barcoded or otherwise tagged DNA library. A plurality of samples may be taken from a single patient, e.g. serially during a course of treatment. Moreover, a plurality of samples may be taken from a plurality of patients. Sample preparation may be as described in the Materials and Methods section herein.

“and/or” where used herein is to be taken as specific disclosure of each of the two specified features or components with or without the other. For example “A and/or B” is to be taken as specific disclosure of each of (i) A, (ii) B and (iii) A and B, just as if each is set out individually herein. Providing Sequence Reads

The sequence reads data may be provided or obtained directly, e.g., by sequencing the cfDNA sample or library or by obtaining or being provided with sequencing data that has already been generated, for example by retrieving sequence read data from a non-volatile or volatile computer memory, data store or network location. Where the sequence reads are obtained by sequencing a sample, the median mass of input DNA may in some cases be in the range 1-100 ng, e.g., 2-50 ng or 3-10 ng. The DNA may be amplified to obtain a library

having, e.g. 100-1000 ng of DNA. The sequence reads may be in a suitable data format, such as FASTQ.

Sequence Data Processing and Error Suppression

The sequence read data, e.g., FASTQ files, may be subjected to one or more processing or clean-up steps prior to or as part of the step of reads collapsing into read families. For example, the sequence data files may be processed using one or more tools selected from as FastQC v0.11.5, a tool to remove adaptor sequences (e.g. cutadapt v1.9.1). The sequence reads (e.g. trimmed sequence reads) may be aligned to an appropriate reference genome, for example, the human reference genome GRCh37 for a human subject.

As used herein “read” or “sequencing read” may be taken to mean the sequence that has been read from one molecule and read once. Each molecule can be read any number of times, depending on the sequencing performed.

“Classifier” or “classification algorithm” may be a model or algorithm that maps input data, such as a cfDNA fragment size features, to a category, such as cancerous or non-cancerous origin. In some embodiments, the present invention provides methods for detecting, classifying, prognosticating, or monitoring cancer in subjects. In particular, data obtained from sequence analysis, such as fragment length and/or copy number (e.g. trimmed median absolute deviation from copy-number neutrality “t-MAD”) of may be evaluated using one or more pattern recognition algorithms. Such analysis methods may be used to form a predictive model, which can be used to classify test data. For example, one convenient and particularly effective method of classification employs multivariate statistical analysis modelling, first to form a model (a “predictive mathematical model”) using data (“modelling data”) from samples of known category (e.g., from subjects known to have a particular cancer), and second to classify an unknown sample (e.g., “test sample”) according to category.

Pattern recognition is the use of multivariate statistics, both parametric and non-parametric, to analyse data, and hence to classify samples and to predict the value of some dependent variable based on a range of observed measurements. There are two main approaches. One set of methods is termed “unsupervised” and these simply reduce data complexity in a rational way and also produce display plots which can be interpreted by the human eye. However, this type of approach may not be suitable for developing a clinical assay that can be used to classify samples derived from subjects independent of the initial sample population used to train the prediction algorithm.

The other approach is termed “supervised” whereby a training set of samples with known class or outcome is used to produce a mathematical model which is then evaluated with independent validation data sets. Here, a “training set” of sequence information, e.g. fragmentation features and/or copy number features, is used to construct a statistical model that predicts correctly the class of each sample. This training set is then tested with independent data (referred to as a test or validation set) to determine the robustness of the computer-based model. These models are sometimes termed “expert systems,” but may be based on a range of different mathematical procedures such as support vector machine (SVM), decision trees, k-nearest neighbour and naïve Bayes, each of which are contemplated herein for use in accordance with the present invention. As detailed in the Examples herein, logistic regression (LR) and Random Forests (RF) were used for variable selection and the classification of samples as “healthy” or “cancer”. Supervised methods can use a data set with reduced dimensionality (for example, the first few principal components), but typically use unreduced

19

data, with all dimensionality. The robustness of the predictive models can also be checked using cross-validation, by leaving out selected samples from the analysis.

Tailored Panel Sequencing (TAPAS)

As used herein tailored panel sequencing refers to sequencing of targeted regions and/or genes. This may employ selected or custom capture panels that target genes of interest, such as genes commonly mutated in cancer and/or genes found to carry mutations in a tumour of the subject of interest (e.g. identified by sequencing matched tumor tissue DNA and plasma DNA samples). In some cases the capture panels may range in size from 0.5-5 Mb, e.g. 1-3 Mb.

The following is presented by way of example and is not to be construed as a limitation to the scope of the claims.

EXAMPLES

Materials and Methods

Study Design

344 plasma samples from 200 patients with multiple cancer types, and 65 plasma samples from 65 healthy controls, were collected. Among the patients, 172 individuals were recruited through prospective clinical studies at Addenbrooke's Hospital, Cambridge, UK, approved by the local research ethics committee (REC reference numbers: 07/Q0106/63; and NRES Committee East of England—Cambridge Central 03/018). Written informed consent was obtained from all patients and blood samples were collected before and after initiation of treatment with surgery or chemotherapeutic agents. DNA was extracted from 2 mL of plasma using the QIAamp circulating nucleic acid kit (Qiagen) or QIASymphony (Qiagen) according to the manufacturer's instructions. In addition, 28 patients were recruited as part of the Copenhagen Prospective Personalized Oncology (CoPPO) program (Ref: PMID: 25046202) at Rigshospitalet, Copenhagen, Denmark, approved by the local research ethics committee. Baseline tumor tissue biopsies were available from all 28 patients, together with re-biopsies collected at relapse from two patients, including matched plasma samples. Brain tumor patients were recruited at the Addenbrooke's Hospital, Cambridge, UK, as part of the BLING study (REC-15/EE/0094). Bladder cancer patients were recruited at the Netherlands Cancer Institute, Amsterdam, The Netherlands, and approval was in accordance with national guidelines (N13KCM/CFMPB250) (47). 65 plasma samples were obtained from healthy control individuals using a similar protocol (Seralab). Plasma samples were freeze-thawed no more than 2 times to reduce artifactual fragmentation of cfDNA. FIG. 1 describes the study as a flowchart.

In Vitro Size Selection

Between 8-20 ng of DNA were loaded into a 3% agarose cassette (HTC3010, Sage Bioscience) and size selection was performed on a PippinHT (Sage Bioscience) according to the manufacturer's protocol. Quality controls of in vitro size selection were performed on 20 healthy controls samples. Duplicate reads observed with in vitro selection were removed for any downstream size selection analysis. A QC metric called the median absolute pairwise difference (MAPD) algorithm was used to determine the sequencing noise. MAPD measured the absolute difference between the log₂ CN ratios of every pair of neighboring bins and determined the median across all bins. Higher MAPD scores reflected greater noise, typically associated with poor-quality

20

samples. All samples exhibited a MAPD score of 0.01 (+−0.01), irrespective of the size selection condition.

TAm-Seq

Tagged-Amplicon Deep Sequencing libraries were prepared as previously described (34), using primers designed to assess single nucleotide variants (SNV) and small indels across selected hotspots and the entire coding regions of TP53. Libraries were sequenced using MiSeq or HiSeq 4000 (Illumina).

10 Shallow Whole Genome Sequencing (sWGS)

Indexed sequencing libraries were prepared using commercially available kits (ThruPLEX-Plasma Seq and/or Tag-Seq, Rubicon Genomics). Libraries were pooled in equimolar amounts and sequenced to <0.4x depth of coverage on a 15 HiSeq 4000 (Illumina) generating 150-bp paired-end reads. Sequence data were analyzed using an in-house pipeline. Paired end sequence reads were aligned to the human reference genome (GRCh37) using BWA-mem following the removal of contaminating adapter sequences (48). PCR 20 and optical duplicates were marked using MarkDuplicates (Picard Tools) feature and these were excluded from downstream analysis along with reads of low mapping quality and supplementary alignments. When necessary, reads were down-sampled to 10 million in all samples for comparison 25 purposes.

Somatic Copy Number Aberration Analysis

The analysis was performed in R using a software suite for shallow Whole Genome Sequencing copy number analysis named CNAclinic (github.com/sdchandra/CNAclinic) as 30 well as the QDNaseq pipeline (49). Sequencing reads were randomly sampled to 10 million reads per dataset and allocated into equally sized (30 Kbp) non-overlapping bins throughout the length of the genome. Read counts in each bin were corrected to account for sequence GC content and 35 mappability. Bins overlapping 'blacklisted' regions (derived from the ENCODE Project and the 1000 Genomes Project database) prone to artefacts were excluded from downstream analysis. Read counts in test samples were normalized by the 40 counts from an identically processed healthy individual and log₂ transformed to obtain copy number ratio values per genomic bin. Read counts in healthy controls were normalized by their median genome-wide count. Bins were then 45 segmented using both Circular Binary Segmentation and Hidden Markov Model algorithms. An averaged log₂ R value per bin was calculated.

An in-house empirical blacklist of aberrant read count 50 regions was constructed. Firstly, 65 sWGS datasets from healthy plasma were used to calculate median read counts per 30 Kbp genomic bin as a function of GC content and mappability. A 2D LOESS surface was then applied and the difference between the actual count and the LOESS fitted 55 values were calculated. The median of these residual values across the 65 controls were calculated per genomic bin and regions with median residuals greater than 4 standard deviations were blacklisted. The averaged segmental log₂ R values in each test sample that overlap this cfDNA blacklist were trimmed and the median absolute value was calculated. This score was defined as the trimmed median absolute deviation (t-MAD) from log₂ R=0. The R code to reproduce 60 this analysis is provided in github.com/sdchandra/tMAD (incorporated herein by reference in its entirety).

Whole Exome Sequencing (WES)

Indexed sequencing libraries were prepared as described above (see Methods, sWGS). Plasma DNA libraries from 65 each sample were made and pooled together for exome capture (TruSeq Exome Enrichment Kit, Illumina). Pools were concentrated using a SpeedVac vacuum concentrator

(Eppendorf). Exome enrichment was performed following the manufacturer's protocol. Enriched libraries were quantified using quantitative PCR (KAPA library quantification, KAPA Biosystems), and DNA fragments sizes observed by Bioanalyzer (2100 Bioanalyzer, Agilent Genomics) and pooled in equimolar ratios for paired-end next generation sequencing on a HiSeq4000 (Illumina). Sequencing reads were de-multiplexed allowing zero mismatches in barcodes. Paired-end alignment to the GRCh37 reference genome was performed using BWA-mem for all exome sequencing data (germline/plasma/tumor tissue DNA). PCR duplicates were marked using Picard. Base quality score recalibration and local realignment were performed using Genome Analysis Tool Kit (GATK).

Mutation Calling

Mutation allele fractions (MAFs) for each single-base locus were calculated with MuTect2 for all bases with PHRED quality ≥ 30 . Filtering parameters were then applied so that a mutation was called if no mutant reads for an allele were observed in germline DNA at a locus that was covered at least 10 \times , and if at least 4 reads supporting the mutant were found in the plasma data with at least 1 read on each strand (forward and reverse). At loci with <10 \times coverage in normal DNA and no mutant reads, mutations were called in plasma if a prior plasma sample showed no evidence of a mutation and was covered adequately (10 \times or more). A method called Integrated Signal Amplification for Non-invasive Interrogation of Tumors was used to aggregate mutations called before and after size selection. This method combined different subsets of mutations called from the same plasma DNA sample using different processing approaches. The mutation aggregation as used in this study was formalized as follows: aggregated mutations=mutations detected without size selection U (mutations detected with in vitro size selection U mutations detected with in silico size selection).

In Silico Size Selection

Paired-end reads are generated by sequencing DNA from both ends of the fragments present in the library. The original length of the DNA can be inferred using the mapping locations of the read ends in the genome. Once alignment is complete, Samtools software is used to select paired reads that correspond to fragment lengths in a specific range. Mutect2 is used to call mutations from this in silico size selected data as described in the previous section.

Tumor-Guided Capture Sequencing

Matched tumor tissue DNA and plasma DNA samples of 19 patients collected from the RigsHospitalet (Copenhagen, Denmark) with advanced cancer were sequenced by WES. Variants were called from these samples by mutation calling (see above). Hybrid-based capture for longitudinal plasma samples analysis were designed to cover these variants for each patient using SureDesign (Agilent). A median of 160 variants were included per patient, and in addition, 41 common genes of interest for pan-cancer analysis were included in the tumor-guided sequencing panel. Indexed sequencing libraries were prepared as per sWGS (see above). Plasma DNA libraries from each sample were made and pooled together for tumor-guided capture sequencing (SureSelect, Agilent). Pools were concentrated using a SpeedVac vacuum concentrator (Eppendorf). Capture enrichment was performed following the manufacturer's protocol. Enriched libraries were quantified using quantitative PCR (KAPA library quantification, KAPA Biosystems), and DNA fragments sizes controlled by Bioanalyzer (2100 Bioanalyzer, Agilent Genomics) and pooled in equimolar ratio for paired-end next generation sequencing on a

HiSeq4000 (Illumina). Sequencing reads were de-multiplexed allowing zero mismatches in barcodes. Paired-end alignment to the GRCh37 reference genome was performed using BWA-mem for all exome sequencing data including germline, plasma and tumor tissue DNA where generated. PCR duplicates were marked using Picard. Base quality score recalibration and local realignment were performed using Genome Analysis Tool Kit (GATK).

Classification Analysis

The preliminary analysis was carried out on 304 samples (182 high ctDNA cancer samples, 57 low ctDNA cancer samples and 65 healthy controls). For each sample the following features were calculated from sWGS data: t-MAD, amplitude_10 bp, P(20-150), P(160-180), P(20-150)/P(160-180), P(100-150), P(100-150)/P(163-169), P(180-220), P(250-320), P(20-150)/P(180-220) (see Table 2). The data was arranged in a matrix where the rows represent each sample and the columns held the aforementioned features with an extra "class" column with the binary labels of "cancer"/"healthy". The following analysis was carried out in R utilising RandomForest, caret, and pROC packages. The caret package is available and is described at the following URL: topepo.github.io/caret/index.html. Exemplary source code for the classification algorithms described in the Examples herein is shown below in the section headed "Code". The pairwise correlations between the features were calculated to assess multi-collinearity in the dataset. A single variable was selected for removal from pairs with Pearson correlation >0.75 . Highly correlated fragmentation features that were composite of individual variables already in the dataset such as P(20-150)/P(180-220), were prioritized for removal. The features were also assessed for zero variance and linear dependencies but none were flagged. After this pre-processing the following 5 variables were selected for further analysis: t-MAD, amplitude_10 bp, P(160-180), P(180-220) and P(250-320) (see Table 2). All 57 low ctDNA samples were set aside for validation of the models. The data matrix for the remaining high ctDNA cancer samples and healthy controls ($n=247$) were randomly partitioned in a 60:40 split into 1 training and 1 validation dataset with the different cancer types and healthy samples represented in similar proportions. Hence, the training data contained 153 samples (cancer=114, healthy=39) while the first validation set of high ctDNA cancers contained 94 samples (cancer=68, healthy=26). This validation dataset was only utilized for final assessment of the classifiers.

Classification of samples as healthy or cancer was performed using one linear and one non-linear machine learning algorithm, namely logistic regression (LR), and random forest (RF). Each algorithm was paired with recursive feature selection in order to identify the best predictor variables. This analysis was carried out with caret within the framework of 5 repeats of 10-fold cross-validation on the training set. The algorithm was configured to explore all possible subsets of the features. The optimal model for each classifier was selected using ROC metric. Separately, a logistic regression model was trained only using t-MAD as a predictor in order to assess the difference in performance without the addition of fragmentation features. Finally, the 68 high ctDNA cancer samples, 57 low ctDNA cancer samples and 26 healthy controls set aside for validation were used to test the classifiers, utilizing area under the curve in a ROC analysis to quantify their performance.

A secondary analysis was carried out on the same training and validation cohorts with the only difference being the features used in the model. Here, we tested predictive ability

23

of fragmentation features without the addition of information from SCNAs (i.e. t-MAD). Hence the features utilized were: amplitude_10 bp, P(160-180), P(180-220) and P(250-320).

Quantification of the 10 bp Periodic Oscillation

The amplitude of the 10 bp periodic oscillation observed in the size distribution of cfDNA samples was determined from the sWGS data as follows. Local maxima and minima in the range 75 bp to 150 bp were identified. The average of their positions across the samples was calculated (for minima: 84, 96, 106, 116, 126, 137, 148, and maxima: 81, 92, 102, 112, 122, 134, 144). To compute the amplitude of the oscillations with 10 bp periodicity observed below 150 bp, the sum of the minima were subtracted from the sum of the heights of the maxima. The larger this difference, the more distinct the peaks. The height of the x bp peak is defined as the number of fragments with length x divided by the total number of fragments. To define local maxima, y positions were selected such that y was the largest value in the interval [y-2, y+2]. The same rationale was used to pick minima.

Example 1: Surveying the Fragmentation Features of Tumour cfDNA

A catalogue of cfDNA fragmentation features was generated using 344 plasma samples from 200 patients with 18 different cancer types, and an additional 65 plasma samples from healthy controls (FIG. 1 and FIG. 2A). The size distribution of cfDNA fragments in cancer patients differed in the size ranges of 90-150 bp, 180-220 bp and 250-320 bp compared to healthy individuals (FIG. 2B and FIG. 3). cfDNA fragment sizes in plasma of healthy individuals, and in plasma of patients with late stage glioma, renal, pancreatic and bladder cancers, were significantly longer than in other late stage cancer types including breast, ovarian, lung, melanoma, colorectal and cholangiocarcinoma ($p<0.001$, Kruskal-Wallis; FIG. 2C). Sorting the 18 cancer types according to the proportion of cfDNA fragments in the size range 20-150 bp was very similar to ordering by Bettegowda et al. based on the concentrations of ctDNA measured by individual mutation assays (FIG. 2D) (6). In contrast to previous reports (6, 34), this sorting analysis was performed without any prior knowledge of the presence of mutations or somatic copy number alterations (SCNAs), yet allowed the investigation of ctDNA content in different cancers.

Example 2: Sizing Up Mutant ctDNA

The size profile of mutant ctDNA in plasma was determined using two high specificity approaches. First, the specific size profile of ctDNA and non-tumor cfDNA was inferred with sWGS from the plasma of mice bearing human ovarian cancer xenografts (FIG. 4A). There was a shift in ctDNA fragment sizes to less than 167 bp (FIG. 4B). Second, the size profile of mutant ctDNA was determined in plasma from 19 cancer patients, using deep sequencing with patient-specific hybrid-capture panels developed from whole-exome profiling of matched tumor samples (FIG. 4C). By sequencing hundreds of mutations at a depth $>300\times$ in cfDNA, allele-specific reads from mutant and normal DNA were obtained. Enrichment of DNA fragments carrying tumor-mutated alleles was observed in fragments \sim 20-40 bp shorter than nucleosomal DNA sizes (multiples of 167 bp) (FIG. 4D). Mutant ctDNA was generally more fragmented than non-mutant cfDNA, with a maximum enrichment of ctDNA in fragments between 90 and 150 bp (FIG.

24

5), as well as enrichment in the size range 250-320 bp. These data also indicated that mutant DNA in plasma of patients with advanced cancer (pre-treatment) is consistently shorter than predicted mono-, and di-nucleosomal DNA fragment lengths (FIG. 4D).

Example 3: Selecting Tumour-Derived DNA Fragments

These data indicated that ctDNA is shorter than non-tumor cfDNA and suggested that biological differences in fragment lengths could be harnessed to improve ctDNA detection. The feasibility of selective sequencing of shorter fragments was determined using in vitro size selection with a bench-top microfluidic device followed by sWGS, in 48 plasma samples from 35 patients with high-grade serous ovarian cancer (HGSOC) (FIG. 6A, FIG. 7 and FIG. 8). The accuracy and quality of the size selection was assessed using the plasma from 20 healthy individuals (FIG. 6B and FIG. 9). The utility of in silico size selection of fragmented DNA was also explored using read-pair positioning from unprocessed sWGS data (FIG. 6A). In silico size selection was performed once reads were aligned to the genome reference, by selecting the paired-end reads that corresponded to the fragments lengths in a 90-150 bp size range. FIG. 6C, FIG. 6D and FIG. 6E illustrate the effect of in vitro size selection for one HGSOC case. SCNAs in plasma cfDNA before treatment were identified, when the concentration of ctDNA was high (FIG. 6C). Only a small number of focal SCNAs were observed in the subsequent plasma sample collected 3 weeks after initiation of chemotherapy (without size selection, FIG. 6D). In vitro size selection of the same post-treatment plasma sample showed a median increase of 6.4 times in the amplitude of detectable SCNAs without size selection. Selective sequencing of shorter fragments in this sample resulted in the detection of multiple other SCNAs that were not observed without size selection (FIG. 6E), and a genome-wide copy-number profile that was similar to that obtained before treatment when ctDNA levels were 4 times higher (FIG. 6C). It was concluded that selecting short DNA fragments in plasma can enrich tumor content on a genome-wide scale.

Example 4: Quantifying the Impact of Size Selection

To quantitatively assess the enrichment after size selection on a genome-wide scale, a metric from sWGS data ($<0.4\times$ coverage) called t-MAD (trimmed Median Absolute Deviation from copy-number neutrality, see FIG. 10A) was developed. All sWGS data were downsampled to 10 million sequencing reads for comparison. To define the detection threshold, the t-MAD score for sWGS data from 65 plasma samples from 46 healthy individuals was measured and the maximal value found (median=0.01, range 0.004-0.015). On comparison of the t-MAD to the mutant allele fraction (MAF) in the high ctDNA cancer types assessed by digital PCR (dPCR) or WES in 97 samples, there was a high correlation (Pearson correlation, $r=0.80$) (FIG. 10B) between t-MAD and MAF, for samples with t-MAD greater than the detection threshold (0.015), or with $MAF>0.025$. FIG. 11 shows that the slope of t-MAD versus MAF fit lines differed between cancer types (range 0.17-1.12) reflecting likely differences in the extent of SCNAs. The sensitivity of t-MAD for detecting low ctDNA levels was estimated using a spike-in dilution of DNA from a patient with a TP53 mutation into DNA from a pool of 7 healthy individuals

(FIG. 12) which confirmed that the t-MAD score was linear with ctDNA levels down to MAF of ~0.01. In addition, t-MAD scores greater than the detection threshold (0.015) for samples were present even in samples with a MAF as low as 0.004. t-MAD was also strongly correlated with tumor volume determined by RECIST1.1 (Pearson correlation, $r=0.6$, $p<0.0001$, $n=35$) (FIG. 13).

Using t-MAD ctDNA was detected from 69% (130/189) of the samples from cancer types where ctDNA levels have been shown to be high (FIG. 10C). From cancer types for which ctDNA levels are suspected to be low (glioma, renal, bladder, pancreatic), ctDNA was detected in 17% (10/57) of the cases (FIG. 10C). To improve the sensitivity for detecting t-MAD in silico size selection of the DNA fragments between 90-150 bp from the high ctDNA cancers ($n=189$) and healthy controls ($n=65$) was used (FIG. 10D). Receiver operating characteristic (ROC) analysis comparing the t-MAD score for the samples revealed an area under the curve (AUC) of 0.90 after in silico size selection, against an AUC of 0.69 without size selection (FIG. 10D).

To explore whether size selected sequencing could improve the detection of response or disease progression, sWGS of longitudinal plasma samples from six cancer patients (FIGS. 10E and F) and in silico size selection of the cfDNA fragments between 90-150 bp was used. In two patients, size selected samples indicated tumor progression 60 and 87 days before detection by imaging or unselected t-MAD analysis (FIGS. 10E and F). Other longitudinal samples exhibited improvements in the detection of ctDNA with t-MAD and size selection (FIG. 10F). Confirmation in large clinical studies will be necessary to determine the potential of selective sequencing of ctDNA for clinical applications.

Example 6: Identifying More Clinically Relevant Mutations with Size Selection

The ability of size selection to increase the sensitivity for detecting new mutations in cfDNA was examined. To test effects on copy number aberrations, 35 patients with HGSOC were studied as this is the archetypal copy-number driven cancer (35). t-MAD was used to quantify the enrichment of ctDNA with in vitro size selection in 48 plasma samples, including samples collected before and after initiation of chemotherapy treatment. In vitro size selection resulted in an increase in the calculated t-MAD score from the sWGS data for 47/48 of the plasma samples (98%, t-test, $p=0.06$) with a mean 2.5 and median 2.1-fold increase (FIG. 14A). The t-MAD scores were then compared against those obtained by sWGS for the plasma samples from healthy individuals. 44 of the 48 size-selected HGSOC plasma samples (92%) had a t-MAD score greater than the highest t-MAD value determined in the in vitro size selected healthy plasma samples (FIG. 14A and FIG. 15), compared to only 24 out of 48 without size selection (50%). ROC analysis comparing the t-MAD score for the samples from the cancer patients (pre- and post-treatment initiation, $n=48$) and healthy controls ($n=46$) revealed an AUC of 0.97 after in vitro size selection, with maximal sensitivity and specificity of 90% and 98%, respectively. This was significantly superior to detection by sWGS without size selection (AUC=0.64) (FIG. 14B).

This was then investigated to determine if improved sensitivity resulted in the detection of SCNAs with potential clinical value. Across the genome, t-MAD scores evaluating SCNAs were higher after size selection in 33/35 (94%) HGSOC patients, and the absolute level of the copy number

(log 2 ratio) values significantly increased after in vitro size selection (t-test for the means, $p=0.003$) (FIG. 14C). The relative copy number values were then compared for 15 genes frequently altered in HGSOC (Table 3). Analysis of plasma cfDNA after size selection revealed a large number of SCNA that were not observed in the same samples without size selection (FIG. 14D), including amplifications in key genes such as NF1, TERT, and MYC (FIG. 16).

To exclude the possibility that size selection might only increase the sensitivity for sWGS analysis, it was examined if enrichment was seen for substitutions. Whole exome sequencing of plasma cfDNA from 23 patients with 7 cancer types was performed (FIG. 2). A comparison of the size distributions of fragments carrying mutant or non-mutant alleles (FIG. 17A) could be made using the WES data, and indicated whether size selection could identify additional mutations. 6 patients with HGSOC were selected and WES of plasma DNA with and without in vitro size selection in the 90-150 bp range was performed, analysing time-points before and after initiation of treatment (36). In addition, in silico size selection for the same range of fragment sizes was performed (FIG. 17A). Analysis of the MAF of SNVs revealed statistically significant enrichment of the tumor fraction with both in vitro size selection (mean 4.19-fold, median 4.27-fold increase, t-test, $p<0.001$) and in silico size selection (mean 2.20-fold, median 2.25-fold increase, t-test, $p<0.001$) (FIG. 17A and FIG. 18). Three weeks after initiation of treatment, ctDNA levels are often lower (36), and therefore post-treatment plasma samples were further analyzed using Tagged-Amplicon Deep Sequencing (TAm-Seq) (37). Enrichment of MAFs by in vitro size selection was observed to be between 0.9 and 118 times (mean 2.1 times, median 1.5 times) compared to the same samples without size selection (FIG. 19).

Size selection with both in vitro and in silico methods increased the number of mutations detected by WES by an average of 53% compared to no size selection (FIG. 17B). A total of 1023 mutations in the non-size-selected samples were identified. An additional 260 mutations were detected by in vitro size selection, and an additional 310 mutations were called after in silico size selection (FIG. 17B and Table 4). New mutations were also detectable in tumor specimens, which excludes the possibility that the improved sensitivity for mutation detection was a result of sequencing artefacts. In silico size selection was then used in an independent cohort of 16 patients, where matched tumor tissue DNA was available. In silico size selection enriched the MAF for nearly all mutations (2061/2133, 97%), with an average increase of MAF of $\times 1.7$ (FIG. 17C). For 13 of 16 patients (81%) additional mutations in plasma after in silico size selection were identified. Of these 82 additional mutations, 23 (28%) were confirmed to be present in the matched tumor tissue DNA (FIG. 17D). Notably, this included mutations in key cancer genes including BRAF, ARID1A, and NF1 (FIG. 20).

Example 7: Detecting Cancer by Supervised Machine Learning Combining cfDNA Fragmentation and Somatic Alteration Analysis

It is important to note that although in vitro and in silico size selection increase the sensitivity of detection, they also result in a loss of cfDNA for analysis. Regions of the cancer genome which are not altered by mutation also excluded and cannot contribute to the analysis (FIG. 21). It was hypoth-

esized that leveraging other biological properties of the cfDNA fragmentation profile could enhance the detection of ctDNA.

The sWGS data defined other cfDNA fragmentation features including (1) the proportion of fragments in multiple size ranges, (2) the ratios of proportions of fragments in different sizes and (3) the amplitude of oscillations in fragment-size density with 10 bp periodicity (FIG. 22A). These fragmentation features were compared between cancer patients and healthy individuals (FIG. 23) and the feature representing the proportion (P) of fragments between 20-150 bp exhibited the highest AUC (0.819). Principal component analysis (PCA) of the samples represented by t-MAD and fragmentation features showed a separation between healthy and cancerous samples and that fragment features clustered with t-MAD scores (FIG. 22B).

Furthermore, the potential of fragmentation features to enhance the detection of tumor DNA in plasma samples was explored. A predictive analysis was performed using the t-MAD score and 9 fragmentation features across 304 samples (239 from cancers patients and 65 from healthy controls) (FIG. 22C and FIG. 24 and Table 2). The 9 fragmentation features determined from sWGS included five features based on the proportion (P) of fragments in defined size ranges: P(20-150), P(100-150), P(160-180), P(180-220), P(250-320); three features based on ratios of those proportions: P(20-150)/P(160-180), P(100-150)/P(163-169), P(20-150)/P(180-220); and a further feature based on the amplitude of the oscillations having 10 bp periodicity observed below 150 bp.

Variable selection and the classification of samples as “healthy” or “cancer” were performed using logistic regression (LR) and random forests (RF) trained on 153 samples, and validated on two datasets of 94 and 83 independent samples (FIG. 22C). The best feature set for the LR model included t-MAD, 10 bp amplitude, P(160-180), P(180-220) and P(250-320). The same five variables were independently identified using the RF model (with some differences in their ranking). FIG. 25 shows performance metrics for the different algorithms on training set data using cross-validation. The source code for the classification algorithms is shown below in the section headed “Code”. Using t-MAD alone in the validation pan-cancer dataset (FIG. 22D and FIG. 24), cancer samples could be distinguished from healthy individuals with AUC=0.764. Using the LR model improved the classification of the samples to AUC=0.908. The RF model (trained on the 153-sample training set) could distinguish cancer from healthy individuals even more accurately in the validation data set (n=94) with AUC=0.994. On the second validation dataset containing low-ctDNA cancer samples (n=83) (FIG. 22E), t-MAD alone or the LR performed less well, with AUC values of 0.421 and 0.532 respectively. However, the RF model was still able to distinguish samples from low-ctDNA cancer samples from healthy controls with AUC=0.914. At a specificity of 95%, the RF model correctly classified as cancer 64/68 (94%) of the samples from high-ctDNA cancers (colorectal, cholangiocarcinoma, ovarian, breast, melanoma), and 37/57 (65%) of the samples from

low-ctDNA cancers (pancreatic, renal, glioma) (FIG. 22F). In a second iteration of model training, t-MAD was omitted, using only the 4 fragmentation features (FIG. 26). The RF model could still distinguish cancer from healthy controls albeit with slightly reduced AUCs (0.989 for cancer types with high levels of ctDNA and 0.891 for cancer types with low levels of ctDNA), suggesting that the cfDNA fragmentation pattern is most important predictive component.

Example 8: Use of Random Forest (RF) Model to Predict Detection of ctDNA in Cancer Patient Fluid

A random forest (RF) model in accordance with the present invention and as described in Example 7 was based on the density or proportion of plasma cell-free DNA fragments with length 20-150, 100-150, 160-180, 163-169, 180-220 and 250-320 bp, as well as the amplitude of the oscillations with 10 bp periodicity and can predict the probability that a given plasma sample has been collected from an individual with cancer.

In addition, our data indicates that the output of this same RF classification model might allow for the triage of cancer patient fluid samples into those with sufficiently high levels of ctDNA for detection by other methods (including those with greater sensitivity and/or that allow targeted analysis of specific somatic mutations), and those without.

After applying the RF model to plasma samples from patients with renal cell carcinoma (RCC), of those with >50% probability of cancer by the RF model:

- ~62% had detectable ctDNA in plasma by our INTEGRATION of VARIANT READS of TAIlor PAnel Sequencing (INVAR TAPAS) method (see co-pending patent application GB1803596.4 filed 6 Mar. 2018, the contents of which are incorporated herein by reference);
- ~63% had detectable ctDNA in plasma by INVAP and/or t-MAD (the latter of which is as described above);
- ~81% had detectable ctDNA in plasma and/or urine by INVAP and/or t-MAD. Conversely, only 11% of plasma samples with <50% probability of cancer by RF model, had detectable ctDNA. This is summarised in FIG. 27.

In summary, this analysis has the potential to highlight those cancer patients in which ctDNA analysis (by more sensitive or targeted methods such as INVAP-TAPAS) is more likely to yield informative output. In-turn these samples are more likely to prove clinically useful, potentially allowing, for example, prediction of response to therapy through identification of resistance mutations, disease prognostication, and assessment of clonal evolution through application of targeted methods. This may prove particularly relevant in those cancer types in which ctDNA detection is unreliable (such as renal cancer and glioblastoma), even at later stages of disease at which ctDNA detection would be expected to be reliable (based on equivalent data from other cancer types). Moreover, preliminary results (not shown) suggest that the above findings for RCC are corroborated in a glioblastoma cohort.

Tables

TABLE 1

summary table of the samples and patients included in the study								
index	patient	sample	SLX	barcode	cancer	cancer_type	timepoint	RECIST_volume
1	GB2	GB2_1	SLX-11868	D710-D505	glioblastoma	low_ctDNA_cancer	baseline	NA
2	GB3	GB3_1	SLX-11868	D710-D506	glioblastoma	low_ctDNA_cancer	baseline	NA
3	GB4	GB4_1	SLX-11868	D710-D507	glioblastoma	low_ctDNA_cancer	baseline	NA
4	GB5	GB5_1	SLX-11868	D710-D508	glioblastoma	low_ctDNA_cancer	baseline	NA

TABLE 1-continued

summary table of the samples and patients included in the study								
index	patient	sample	SLX	barcode	cancer	cancer_type	timepoint	RECIST_volume
5	GB6	GB6_1	SLX-11868	D711-D505	glioblastoma	low_ctDNA_cancer	baseline	NA
6	GB7	GB7_1	SLX-11868	D711-D506	glioblastoma	low_ctDNA_cancer	baseline	NA
7	GB8	GB8_1	SLX-11868	D711-D507	glioblastoma	low_ctDNA_cancer	baseline	NA
8	GB9	GB9_1	SLX-11868	D711-D508	glioblastoma	low_ctDNA_cancer	baseline	NA
9	GB10	GB10_1	SLX-11868	D712-D505	glioblastoma	low_ctDNA_cancer	baseline	NA
10	GB11	GB11_1	SLX-11868	D712-D506	glioblastoma	low_ctDNA_cancer	baseline	NA
11	GB12	GB12_1	SLX-11868	D712-D507	glioblastoma	low_ctDNA_cancer	baseline	NA
12	GB13	GB13_1	SLX-11868	D712-D508	glioblastoma	low_ctDNA_cancer	baseline	NA
13	Os1	Os1_1	SLX-11870	D707-D505	esophageal junction	low_ctDNA_cancer	baseline	NA
14	B1	B1_1	SLX-11034	A019	breast	high_ctDNA_cancer	baseline	NA
15	L1	L1_1	SLX-11870	D711-D504	lung	high_ctDNA_cancer	baseline	NA
16	Ov1	Ov1_1	SLX-11870	D712-D502	ovarian	high_ctDNA_cancer	baseline	NA
17	Ov2	Ov2_1	SLX-11870	D708-D505	ovarian	high_ctDNA_cancer	baseline	NA
18	Ren1	Ren1_1	SLX-11870	D708-D507	renal	low_ctDNA_cancer	baseline	NA
19	B2	B2_1	SLX-11870	D710-D501	breast	high_ctDNA_cancer	baseline	NA
20	L2	L2_1	SLX-11870	D712-D504	lung	high_ctDNA_cancer	baseline	NA
21	L3	L3_1	SLX-11870	D712-D503	lung	high_ctDNA_cancer	baseline	NA
22	T1	T1_1	SLX-11870	D709-D506	thymoma	high_ctDNA_cancer	baseline	NA
23	R1	R1_1	SLX-11870	D710-D504	rectum	high_ctDNA_cancer	baseline	NA
24	B3	B3_1	SLX-11870	D711-D502	breast	high_ctDNA_cancer	baseline	NA
25	L4	L4_1	SLX-13710	D708-D508	lung	high_ctDNA_cancer	baseline	NA
26	R2	R2_1	SLX-13710	D707-D502	rectum	high_ctDNA_cancer	baseline	NA
27	B4	B4_1	SLX-13710	D706-D503	breast	high_ctDNA_cancer	baseline	NA
28	P1	P1_1	SLX-13710	D705-D504	pancreatic	low_ctDNA_cancer	baseline	NA
29	Ov3	Ov3_1	SLX-13710	D704-D505	ovarian	high_ctDNA_cancer	baseline	NA
30	B5	B5_1	SLX-13710	D702-D507	breast	high_ctDNA_cancer	baseline	NA
31	B6	B6_1	SLX-13710	D701-D508	breast	high_ctDNA_cancer	baseline	NA
32	L5	L5_1	SLX-12841	D701-D501	lung	high_ctDNA_cancer	baseline	NA
33	ChC1	ChC1_1	SLX-12841	D701-D502	cholangiocarcinoma	high_ctDNA_cancer	baseline	96
34	B7	B7_1	SLX-12841	D701-D503	breast	high_ctDNA_cancer	baseline	NA
35	C1	C1_1	SLX-12841	D701-D504	colorectal	high_ctDNA_cancer	baseline	NA
36	ChC2	ChC2_1	SLX-12841	D702-D501	cholangiocarcinoma	high_ctDNA_cancer	baseline	87
37	HCC1	HCC1_1	SLX-12841	D702-D502	hepatocellular	high_ctDNA_cancer	baseline	NA
38	C2	C2_1	SLX-12841	D702-D503	colorectal	high_ctDNA_cancer	baseline	NA
39	P2	P2_1	SLX-12841	D702-D504	pancreatic	low_ctDNA_cancer	baseline	NA
40	ChC3	ChC3_1	SLX-12841	D703-D505	cholangiocarcinoma	high_ctDNA_cancer	baseline	NA
41	P3	P3_1	SLX-12841	D703-D506	pancreatic	low_ctDNA_cancer	baseline	NA
42	R3	R3_1	SLX-12841	D703-D507	rectum	high_ctDNA_cancer	baseline	NA
43	ChC4	ChC4_1	SLX-12841	D703-D508	cholangiocarcinoma	high_ctDNA_cancer	baseline	NA
44	ChC5	ChC5_1	SLX-12841	D704-D505	cholangiocarcinoma	high_ctDNA_cancer	baseline	NA
45	P4	P4_1	SLX-12841	D704-D506	pancreatic	low_ctDNA_cancer	baseline	NA
46	C3	C3_1	SLX-12841	D704-D507	colorectal	high_ctDNA_cancer	baseline	158
47	Ov4	Ov4_1	SLX-12841	D704-D508	ovarian	high_ctDNA_cancer	baseline	NA
48	Ov5	Ov5_1	SLX-12841	D705-D501	ovarian	high_ctDNA_cancer	baseline	NA
49	B8	B8_1	SLX-12841	D705-D502	breast	high_ctDNA_cancer	baseline	NA
50	L6	L6_1	SLX-12841	D705-D503	lung	high_ctDNA_cancer	baseline	NA
51	C4	C4_1	SLX-12841	D705-D504	colorectal	high_ctDNA_cancer	baseline	NA
52	Pe1	Pe1_1	SLX-12841	D706-D501	penile	high_ctDNA_cancer	baseline	NA
53	Pr1	Pr1_1	SLX-12841	D706-D502	prostate	high_ctDNA_cancer	baseline	33
54	Ce1	Ce1_1	SLX-12841	D706-D503	cervical	high_ctDNA_cancer	baseline	NA
55	C5	C5_1	SLX-12841	D706-D504	colorectal	high_ctDNA_cancer	baseline	112
56	Ov6	Ov6_1	SLX-12841	D707-D505	ovarian	high_ctDNA_cancer	baseline	NA
57	En1	En1_1	SLX-12841	D707-D506	endometrial	high_ctDNA_cancer	baseline	NA
58	C6	C6_1	SLX-12841	D707-D507	colorectal	high_ctDNA_cancer	baseline	22
59	C7	C7_1	SLX-12841	D707-D508	colorectal	high_ctDNA_cancer	baseline	NA
60	OV04-77	JBLAB_5688	SLX-13223	D701-D501	ovarian	high_ctDNA_cancer	baseline	NA
61	OV04-77	JBLAB_5689	SLX-13223	D701-D502	ovarian	high_ctDNA_cancer	post-treatment	NA
62	OV04-83	JBLAB_5203	SLX-13223	D703-D501	ovarian	high_ctDNA_cancer	baseline	NA
63	OV04-83	JBLAB_5205	SLX-13223	D703-D502	ovarian	high_ctDNA_cancer	post-treatment	NA
64	OV04-122	JBLAB_5712	SLX-13223	D701-D503	ovarian	high_ctDNA_cancer	baseline	NA
65	OV04-122	JBLAB_5713	SLX-13223	D701-D504	ovarian	high_ctDNA_cancer	post-treatment	NA
66	OV04-141	JBLAB_5392	SLX-13223	D703-D503	ovarian	high_ctDNA_cancer	baseline	NA
67	OV04-141	JBLAB_5393	SLX-13223	D703-D504	ovarian	high_ctDNA_cancer	post-treatment	NA
68	OV04-143	JBLAB_5587	SLX-11873	D707-D501	ovarian	high_ctDNA_cancer	baseline	NA
69	OV04-143	JBLAB_5588	SLX-11873	D707-D502	ovarian	high_ctDNA_cancer	post-treatment	NA

TABLE 1-continued

summary table of the samples and patients included in the study								
index	patient	sample	SLX	barcode	cancer	cancer_type	timepoint	RECIST_volume
70	OV04-180	JBLAB_5432	SLX-13223	D705-D505	ovarian	high_ctDNA_cancer	baseline	NA
71	OV04-180	JBLAB_5433	SLX-13223	D705-D506	ovarian	high_ctDNA_cancer	post-treatment	NA
72	OV04-211	JBLAB_5471	SLX-13223	D706-D505	ovarian	high_ctDNA_cancer	baseline	NA
73	OV04-211	JBLAB_5472	SLX-13223	D706-D506	ovarian	high_ctDNA_cancer	post-treatment	NA
74	OV04-226	JBLAB_5507	SLX-13223	D704-D505	ovarian	high_ctDNA_cancer	baseline	NA
75	OV04-226	JBLAB_5508	SLX-13223	D704-D506	ovarian	high_ctDNA_cancer	post-treatment	NA
76	OV04-264	JBLAB_5622	SLX-11873	D707-D503	ovarian	high_ctDNA_cancer	baseline	NA
77	OV04-264	JBLAB_5623	SLX-11873	D707-D504	ovarian	high_ctDNA_cancer	post-treatment	NA
78	OV04-292	JBLAB_5742	SLX-13223	D702-D501	ovarian	high_ctDNA_cancer	baseline	NA
79	OV04-292	JBLAB_5743	SLX-13223	D702-D502	ovarian	high_ctDNA_cancer	post-treatment	NA
80	OV04-295	JBLAB_5420	SLX-13223	D705-D507	ovarian	high_ctDNA_cancer	baseline	NA
81	OV04-295	JBLAB_5422	SLX-13223	D705-D508	ovarian	high_ctDNA_cancer	post-treatment	NA
82	OV04-297	JBLAB_5288	SLX-13223	D704-D507	ovarian	high_ctDNA_cancer	baseline	NA
83	OV04-297	JBLAB_5289	SLX-13223	D704-D508	ovarian	high_ctDNA_cancer	post-treatment	NA
84	OV04-300	JBLAB_5754	SLX-13223	D702-D503	ovarian	high_ctDNA_cancer	baseline	NA
85	OV04-300	JBLAB_5755	SLX-13223	D702-D504	ovarian	high_ctDNA_cancer	post-treatment	NA
86	X76	X76_T1_pre	SLX-13621	D701-D501	ovarian	high_ctDNA_cancer	baseline	NA
87	X75_2	X75_T13_pre	SLX-13621	D702-D501	ovarian	high_ctDNA_cancer	baseline	NA
88	X52	X52_T1_pre	SLX-13621	D703-D501	ovarian	high_ctDNA_cancer	baseline	NA
89	X150	X150_T1_pre	SLX-13621	D704-D501	ovarian	high_ctDNA_cancer	baseline	NA
90	X129	X129_T8_pre	SLX-13621	D705-D501	ovarian	high_ctDNA_cancer	baseline	NA
91	X57	X57_T1_pre	SLX-13621	D706-D501	ovarian	high_ctDNA_cancer	baseline	NA
92	X73	X73_T3B_pre	SLX-13621	D707-D501	ovarian	high_ctDNA_cancer	baseline	NA
93	JG090	JG090_T6_12_pre	SLX-13621	D708-D501	ovarian	high_ctDNA_cancer	baseline	NA
94	X145	X145_T8_pre	SLX-13621	D709-D501	ovarian	high_ctDNA_cancer	baseline	NA
95	X112	X112_T1_pre	SLX-13621	D710-D501	ovarian	high_ctDNA_cancer	baseline	NA
96	X75_1	X75_T1_pre	SLX-13621	D711-D501	ovarian	high_ctDNA_cancer	baseline	NA
97	X72	X72_T1_pre	SLX-13621	D712-D501	ovarian	high_ctDNA_cancer	baseline	NA
98	X74	X74_T1_pre	SLX-13621	D701-D502	ovarian	high_ctDNA_cancer	baseline	NA
99	X127	X127_T1_pre	SLX-13621	D702-D502	ovarian	high_ctDNA_cancer	baseline	NA
100	X30	X30_T1_pre	SLX-13621	D703-D502	ovarian	high_ctDNA_cancer	baseline	NA
101	JBLAB_5180	JBLAB.5180_pre	SLX-13621	D704-D502	ovarian	high_ctDNA_cancer	baseline	NA
102	JBLAB_5027	JBLAB.5027_pre	SLX-13621	D705-D502	ovarian	high_ctDNA_cancer	baseline	NA
103	JBLAB_5595	JBLAB.5595_pre	SLX-13621	D706-D502	ovarian	high_ctDNA_cancer	baseline	NA
104	JBLAB_5599	JBLAB.5599_pre	SLX-13621	D707-D502	ovarian	high_ctDNA_cancer	baseline	NA
105	JBLAB_5611	JBLAB.5611_pre	SLX-13621	D708-D502	ovarian	high_ctDNA_cancer	baseline	NA
106	JBLAB_5477	JBLAB.5477_pre	SLX-13621	D709-D502	ovarian	high_ctDNA_cancer	baseline	NA
107	JBLAB_5632	JBLAB.5632_pre	SLX-13621	D710-D502	ovarian	high_ctDNA_cancer	baseline	NA
108	B9	B9_1	SLX-11043	D705-D506	breast	high_ctDNA_cancer	baseline	119
109	B10	B10_1	SLX-11043	D702-D501	breast	high_ctDNA_cancer	baseline	46
110	B11	B11_1	SLX-11043	D701-D501	breast	high_ctDNA_cancer	baseline	52
111	B12	B12_1	SLX-11043	D705-D508	breast	high_ctDNA_cancer	baseline	23
112	B13	B13_1	SLX-11043	D704-D508	breast	high_ctDNA_cancer	baseline	35
113	B14	B14_1	SLX-11043	D704-D505	breast	high_ctDNA_cancer	baseline	60
114	B15	B15_1	SLX-11043	D703-D503	breast	high_ctDNA_cancer	baseline	116
115	B16	B16_1	SLX-11042	D703-D508	breast	high_ctDNA_cancer	baseline	10
116	B17	B17_1	SLX-11042	D704-D504	breast	high_ctDNA_cancer	baseline	71
117	B18	B18_1	SLX-11042	D704-D502	breast	high_ctDNA_cancer	baseline	19
118	B19	B19_1	SLX-11042	D705-D502	breast	high_ctDNA_cancer	baseline	63
119	B20	B20_1	SLX-11042	D705-D504	breast	high_ctDNA_cancer	baseline	72
120	B21	B21_1	SLX-11042	D701-D505	breast	high_ctDNA_cancer	baseline	21
121	B22	B22_1	SLX-11042	D701-D507	breast	high_ctDNA_cancer	baseline	71
122	B23	B23_1	SLX-11042	D702-D506	breast	high_ctDNA_cancer	baseline	68
123	B24	B24_1	SLX-11042	D702-D508	breast	high_ctDNA_cancer	baseline	18
124	B25	B25_1	SLX-11042	D703-D506	breast	high_ctDNA_cancer	baseline	150
125	B26	B26_1	SLX-11042	D706-D502	breast	high_ctDNA_cancer	baseline	211
126	B27	B27_1	SLX-11042	D706-D503	breast	high_ctDNA_cancer	baseline	91
127	B28	B28_1	SLX-11042	D706-D504	breast	high_ctDNA_cancer	baseline	155
128	B29	B29_1	SLX-11043	D703-D502	breast	high_ctDNA_cancer	baseline	NA
129	B30	B30_1	SLX-11043	D701-D504	breast	high_ctDNA_cancer	post-treatment	NA
130	B31	B31_1	SLX-11043	D704-D507	breast	high_ctDNA_cancer	post-treatment	NA
131	B32	B32_1	SLX-11042	D703-D507	breast	high_ctDNA_cancer	post-treatment	NA
132	B11	B11_1	SLX-10991		bladder	low_ctDNA_cancer	baseline	NA
133	B12	B12_1	SLX-10991		bladder	low_ctDNA_cancer	baseline	NA
134	B13	B13_1	SLX-11094	D708-D501	bladder	low_ctDNA_cancer	baseline	NA

TABLE 1-continued

summary table of the samples and patients included in the study								
index	patient	sample	SLX	barcode	cancer	cancer_type	timepoint	RECIST_volume
135	B14	B14_1	SLX-10575	iPCRtagT014	bladder	low_ctDNA_cancer	baseline	NA
136	B15	B15_1	SLX-11904	D709-D507	bladder	low_ctDNA_cancer	baseline	NA
137	B16	B16_1	SLX-10572	D704-D505	bladder	low_ctDNA_cancer	baseline	NA
138	B17	B17_1	SLX-10572	D708-D507	bladder	low_ctDNA_cancer	baseline	NA
139	B18	B18_1	SLX-11896	D708-D504	bladder	low_ctDNA_cancer	baseline	NA
140	B19	B19_1	SLX-11896	D707-D507	bladder	low_ctDNA_cancer	baseline	NA
141	B110	B110_1	SLX-11896	D707-D508	bladder	low_ctDNA_cancer	baseline	NA
142	B111	B111_1	SLX-11896	D709-D506	bladder	low_ctDNA_cancer	baseline	NA
143	B112	B112_1	SLX-11904	D708-D504	bladder	low_ctDNA_cancer	baseline	NA
144	B113	B113_1	SLX-11904	D709-D501	bladder	low_ctDNA_cancer	baseline	NA
145	B114	B114_1	SLX-11986	D709-D504	bladder	low_ctDNA_cancer	baseline	NA
146	B115	B115_1	SLX-10572	D708-D508	bladder	low_ctDNA_cancer	baseline	NA
147	B116	B116_1	SLX-11896	D707-D502	bladder	low_ctDNA_cancer	baseline	NA
148	B117	B117_1	SLX-10572	D708-D505	bladder	low_ctDNA_cancer	baseline	NA
149	B118	B118_1	SLX-11896	D709-D503	bladder	low_ctDNA_cancer	baseline	NA
150	B119	B119_1	SLX-11896	D708-D503	bladder	low_ctDNA_cancer	baseline	NA
151	Ren2	Ren2_1	SLX-13900	D707-D501	renal	low_ctDNA_cancer	baseline	NA
152	Ren3	Ren3_1	SLX-13900	D707-D502	renal	low_ctDNA_cancer	baseline	NA
153	Ren4	Ren4_1	SLX-13900	D707-D503	renal	low_ctDNA_cancer	baseline	NA
154	Ren5	Ren5_1	SLX-13900	D707-D504	renal	low_ctDNA_cancer	baseline	NA
155	Ren6	Ren6_1	SLX-13900	D708-D501	renal	low_ctDNA_cancer	baseline	NA
156	Ren7	Ren7_1	SLX-13900	D708-D502	renal	low_ctDNA_cancer	baseline	NA
157	Ren8	Ren8_1	SLX-13900	D708-D503	renal	low_ctDNA_cancer	baseline	NA
158	Ren9	Ren9_1	SLX-13900	D708-D504	renal	low_ctDNA_cancer	baseline	NA
159	Ren10	Ren10_1	SLX-13900	D708-D505	renal	low_ctDNA_cancer	baseline	NA
160	Ren11	Ren11_1	SLX-13900	D708-D506	renal	low_ctDNA_cancer	baseline	NA
161	Ren12	Ren12_1	SLX-13900	D708-D507	renal	low_ctDNA_cancer	baseline	NA
162	Ren13	Ren13_1	SLX-13900	D708-D508	renal	low_ctDNA_cancer	baseline	NA
163	Ren14	Ren14_1	SLX-13900	D709-D501	renal	low_ctDNA_cancer	baseline	NA
164	Ren15	Ren15_1	SLX-13900	D709-D502	renal	low_ctDNA_cancer	baseline	NA
165	Ren16	Ren16_1	SLX-13900	D709-D503	renal	low_ctDNA_cancer	baseline	NA
166	Ren17	Ren17_1	SLX-13900	D709-D504	renal	low_ctDNA_cancer	baseline	NA
167	Ren18	Ren18_1	SLX-13900	D709-D505	renal	low_ctDNA_cancer	baseline	NA
168	Ren19	Ren19_1	SLX-13900	D709-D506	renal	low_ctDNA_cancer	baseline	NA
169	Ren20	Ren20_1	SLX-13900	D710-D501	renal	low_ctDNA_cancer	baseline	NA
170	Ren21	Ren21_1	SLX-13900	D710-D502	renal	low_ctDNA_cancer	baseline	NA
171	Ren22	Ren22_1	SLX-13900	D710-D503	renal	low_ctDNA_cancer	baseline	NA
172	Ren23	Ren23_1	SLX-13900	D710-D504	renal	low_ctDNA_cancer	baseline	NA
173	Ren24	Ren24_1	SLX-13900	D710-D505	renal	low_ctDNA_cancer	baseline	NA
174	Ren25	Ren25_1	SLX-13900	D710-D506	renal	low_ctDNA_cancer	baseline	NA
175	Ren26	Ren26_1	SLX-13900	D710-D507	renal	low_ctDNA_cancer	baseline	NA
176	Ren27	Ren27_1	SLX-13900	D710-D508	renal	low_ctDNA_cancer	baseline	NA
177	Ren28	Ren28_1	SLX-13900	D711-D501	renal	low_ctDNA_cancer	baseline	NA
178	Ren29	Ren29_1	SLX-13900	D711-D502	renal	low_ctDNA_cancer	baseline	NA
179	Ren30	Ren30_1	SLX-13900	D711-D503	renal	low_ctDNA_cancer	baseline	NA
180	Ren31	Ren31_1	SLX-13900	D711-D504	renal	low_ctDNA_cancer	baseline	NA
181	Ren32	Ren32_1	SLX-13900	D711-D505	renal	low_ctDNA_cancer	baseline	NA
182	Ren33	Ren33_1	SLX-13900	D711-D506	renal	low_ctDNA_cancer	baseline	NA
183	HIP_1	HIP_1	SLX-12531	D703-D501	healthy	healthy	baseline	NA
184	HIP_10	HIP_10	SLX-12531	D705-D506	healthy	healthy	baseline	NA
185	HIP_11	HIP_11	SLX-12531	D705-D507	healthy	healthy	baseline	NA
186	HIP_12	HIP_12	SLX-12531	D705-D508	healthy	healthy	baseline	NA
187	HIP_13	HIP_13	SLX-12531	D706-D505	healthy	healthy	baseline	NA
188	HIP_14	HIP_14	SLX-12531	D706-D506	healthy	healthy	baseline	NA
189	HIP_15	HIP_15	SLX-12531	D706-D507	healthy	healthy	baseline	NA
190	HIP_16	HIP_16	SLX-12531	D706-D508	healthy	healthy	baseline	NA
191	HIP_17	HIP_17	SLX-12531	D707-D501	healthy	healthy	baseline	NA
192	HIP_18	HIP_18	SLX-12531	D707-D502	healthy	healthy	baseline	NA
193	HIP_19	HIP_19	SLX-12531	D707-D503	healthy	healthy	baseline	NA
194	HIP_2	HIP_2	SLX-12531	D703-D502	healthy	healthy	baseline	NA
195	HIP_20	HIP_20	SLX-12531	D707-D504	healthy	healthy	baseline	NA
196	HIP_21	HIP_21	SLX-12531	D708-D501	healthy	healthy	baseline	NA
197	HIP_22	HIP_22	SLX-12531	D708-D502	healthy	healthy	baseline	NA
198	HIP_23	HIP_23	SLX-12531	D708-D503	healthy	healthy	baseline	NA
199	HIP_24	HIP_24	SLX-12531	D708-D504	healthy	healthy	baseline	NA
200	HIP_27	HIP_27	SLX-12534	D707-D502	healthy	healthy	baseline	NA
201	HIP_28	HIP_28	SLX-12534	D707-D503	healthy	healthy	baseline	NA
202	HIP_29	HIP_29	SLX-12534	D707-D504	healthy	healthy	baseline	NA
203	HIP_3	HIP_3	SLX-12531	D703-D503	healthy	healthy	baseline	NA
204	HIP_30	HIP_30	SLX-12534	D708-D501	healthy	healthy	baseline	NA
205	HIP_31	HIP_31	SLX-12534	D708-D502	healthy	healthy	baseline	NA
206	HIP_32	HIP_32	SLX-12534	D708-D503	healthy	healthy	baseline	NA
207	HIP_33	HIP_33	SLX-12534	D708-D504	healthy	healthy	baseline	NA
208	HIP_34	HIP_34	SLX-12534	D709-D501	healthy	healthy	baseline	NA
209	HIP_35	HIP_35	SLX-12534	D709-D503	healthy	healthy	baseline	NA
210	HIP_36	HIP_36	SLX-12534	D709-D504	healthy	healthy	baseline	NA

TABLE 1-continued

summary table of the samples and patients included in the study								
index	patient	sample	SLX	barcode	cancer	cancer_type	timepoint	RECIST_volume
211	HIP_37	HIP_37	SLX-12534	D710-D501	healthy	healthy	baseline	NA
212	HIP_38	HIP_38	SLX-12534	D710-D502	healthy	healthy	baseline	NA
213	HIP_39	HIP_39	SLX-12534	D710-D503	healthy	healthy	baseline	NA
214	HIP_4	HIP_4	SLX-12531	D703-D504	healthy	healthy	baseline	NA
215	HIP_40	HIP_40	SLX-12534	D710-D504	healthy	healthy	baseline	NA
216	HIP_41	HIP_41	SLX-12534	D711-D505	healthy	healthy	baseline	NA
217	HIP_42	HIP_42	SLX-12534	D711-D506	healthy	healthy	baseline	NA
218	HIP_43	HIP_43	SLX-12534	D711-D507	healthy	healthy	baseline	NA
219	HIP_44	HIP_44	SLX-12534	D711-D508	healthy	healthy	baseline	NA
220	HIP_45	HIP_45	SLX-12534	D712-D505	healthy	healthy	baseline	NA
221	HIP_46	HIP_46	SLX-12534	D712-D506	healthy	healthy	baseline	NA
222	HIP_47	HIP_47	SLX-12534	D712-D507	healthy	healthy	baseline	NA
223	HIP_48	HIP_48	SLX-12534	D712-D508	healthy	healthy	baseline	NA
224	HIP_5	HIP_5	SLX-12531	D704-D501	healthy	healthy	baseline	NA
225	HIP_6	HIP_6	SLX-12531	D704-D502	healthy	healthy	baseline	NA
226	HIP_7	HIP_7	SLX-12531	D704-D503	healthy	healthy	baseline	NA
227	HIP_8	HIP_8	SLX-12531	D704-D504	healthy	healthy	baseline	NA
228	HIP_9	HIP_9	SLX-12531	D705-D505	healthy	healthy	baseline	NA
229	M1	M1_1	SLX-11379	D701-D502	melanoma	high_ctDNA_cancer	baseline	23.8895
230	M1	M1_2	SLX-11379	D701-D501	melanoma	high_ctDNA_cancer	post-treatment	11.3665
231	M4	M4_1	SLX-11379	D702-D501	melanoma	high_ctDNA_cancer	baseline	4.61105
232	M4	M4_2	SLX-12758	D704-D501	melanoma	high_ctDNA_cancer	post-treatment	1.02111
233	M4	M4_3	SLX-12759	D708-D501	melanoma	high_ctDNA_cancer	post-treatment	1.29681
234	M4	M4_4	SLX-12758	D709-D502	melanoma	high_ctDNA_cancer	post-treatment	5.49329
235	M4	M4_5	SLX-12758	D702-D501	melanoma	high_ctDNA_cancer	post-treatment	28.2798
236	M4	M4_6	SLX-11383	D701-D506	melanoma	high_ctDNA_cancer	post-treatment	157.486
237	M4	M4_7	SLX-11379	D701-D503	melanoma	high_ctDNA_cancer	post-treatment	307.577
238	M12	M12_1	SLX-11379	D703-D502	melanoma	high_ctDNA_cancer	baseline	991.038
239	M12	M12_2	SLX-11847	D704-D502	melanoma	high_ctDNA_cancer	post-treatment	135.874
240	M12	M12_3	SLX-11847	D704-D503	melanoma	high_ctDNA_cancer	post-treatment	186.259
241	M12	M12_4	SLX-11847	D707-D507	melanoma	high_ctDNA_cancer	post-treatment	499.186
242	M14	M14_1	SLX-11383	D708-D503	melanoma	high_ctDNA_cancer	baseline	0.95626
243	M14	M14_2	SLX-12758	D706-D506	melanoma	high_ctDNA_cancer	post-treatment	0.46476
244	M22	M22_1	SLX-11379	D704-D507	melanoma	high_ctDNA_cancer	baseline	34.9164
245	M22	M22_2	SLX-12758	D706-D507	melanoma	high_ctDNA_cancer	post-treatment	19.8097
246	M22	M22_3	SLX-11379	D704-D508	melanoma	high_ctDNA_cancer	post-treatment	21.37
247	M22	M22_4	SLX-12758	D704-D508	melanoma	high_ctDNA_cancer	post-treatment	46.8143
248	M32	M32_1	SLX-11379	D705-D506	melanoma	high_ctDNA_cancer	baseline	70.2068
249	M32	M32_2	SLX-11847	D705-D503	melanoma	high_ctDNA_cancer	baseline	123.343
250	C8	C8_T1	SLX-12832	D709-D501	colorectal	high_ctDNA_cancer	post-treatment	133
251	C8	C8_T2	SLX-12832	D709-D502	colorectal	high_ctDNA_cancer	post-treatment	84
252	L5	L5_T2	SLX-12832	D709-D503	lung	high_ctDNA_cancer	post-treatment	NA
253	ChC1	ChC1_3	SLX-12832	D709-D504	cholangiocarcinoma	high_ctDNA_cancer	post-treatment	96
254	ChC1	ChC1_4	SLX-12832	D710-D501	cholangiocarcinoma	high_ctDNA_cancer	post-treatment	NA
255	ChC2	ChC2_2	SLX-12832	D710-D502	cholangiocarcinoma	high_ctDNA_cancer	post-treatment	NA
256	ChC2	ChC2_3	SLX-12832	D710-D503	cholangiocarcinoma	high_ctDNA_cancer	post-treatment	NA
257	HCC1	HCC1_2	SLX-12832	D710-D504	hepatocellular	high_ctDNA_cancer	post-treatment	NA
258	HCC1	HCC1_3	SLX-12832	D711-D505	hepatocellular	high_ctDNA_cancer	post-treatment	NA
259	HCC1	HCC1_4	SLX-12832	D711-D506	hepatocellular	high_ctDNA_cancer	post-treatment	NA
260	HCC1	HCC1_5	SLX-12832	D711-D507	hepatocellular	high_ctDNA_cancer	post-treatment	NA

TABLE 1-continued

summary table of the samples and patients included in the study								
index	patient	sample	SLX	barcode	cancer	cancer_type	timepoint	RECIST_volume
261	P2	P2_2	SLX-12832	D711-D508	pancreatic	low_ctDNA_cancer	post-treatment	NA
262	P4	P4_2	SLX-12832	D712-D505	pancreatic	low_ctDNA_cancer	post-treatment	NA
263	C4	C4_2	SLX-12832	D712-D506	colorectal	high_ctDNA_cancer	post-treatment	NA
264	Pr1	Pr1_4	SLX-12832	D712-D507	prostate	high_ctDNA_cancer	post-treatment	29
265	Ov6	Ov6_2	SLX-12832	D712-D508	ovarian	high_ctDNA_cancer	post-treatment	NA
266	ChC2	ChC2_6	SLX-12838	D701-D505	cholangio-carcinoma	high_ctDNA_cancer	post-treatment	47
267	ChC3	ChC3_2	SLX-12838	D701-D506	cholangio-carcinoma	high_ctDNA_cancer	post-treatment	NA
268	C3	C3_5	SLX-12838	D701-D507	colorectal	high_ctDNA_cancer	post-treatment	NA
269	L6	L6_2	SLX-12838	D701-D508	lung	high_ctDNA_cancer	post-treatment	NA
270	Pr1	Pr1_3	SLX-12838	D702-D505	prostate	high_ctDNA_cancer	post-treatment	NA
271	B7	B7_2	SLX-12838	D702-D506	breast	high_ctDNA_cancer	post-treatment	NA
272	C1	C1_2	SLX-12838	D702-D507	colorectal	high_ctDNA_cancer	post-treatment	NA
273	ChC2	ChC2_4	SLX-12838	D702-D508	cholangio-carcinoma	high_ctDNA_cancer	post-treatment	41
274	ChC2	ChC2_5	SLX-12838	D703-D501	cholangio-carcinoma	high_ctDNA_cancer	post-treatment	NA
275	P4	P4_3	SLX-12838	D703-D502	pancreatic	low_ctDNA_cancer	post-treatment	NA
276	C3	C3_4	SLX-12838	D703-D503	colorectal	high_ctDNA_cancer	post-treatment	119
277	Ov4	Ov4_2	SLX-12838	D703-D504	ovarian	high_ctDNA_cancer	post-treatment	NA
278	Ov5	Ov5_2	SLX-12838	D704-D501	ovarian	high_ctDNA_cancer	post-treatment	NA
279	B8	B8_2	SLX-12838	D704-D502	breast	high_ctDNA_cancer	post-treatment	NA
280	C5	C5_3	SLX-12838	D704-D503	colorectal	high_ctDNA_cancer	post-treatment	65
281	En1	En1_2	SLX-12838	D704-D504	endometrial	high_ctDNA_cancer	post-treatment	NA
282	C6	C6_2	SLX-12838	D705-D505	colorectal	high_ctDNA_cancer	post-treatment	NA
283	ChC1	ChC1_2	SLX-12838	D705-D506	cholangio-carcinoma	high_ctDNA_cancer	post-treatment	NA
284	C3	C3_2	SLX-12838	D705-D507	colorectal	high_ctDNA_cancer	post-treatment	NA
285	C3	C3_3	SLX-12838	D705-D508	colorectal	high_ctDNA_cancer	post-treatment	NA
286	Ov4	Ov4_3	SLX-12838	D706-D505	ovarian	high_ctDNA_cancer	post-treatment	NA
287	Ov5	Ov5_3	SLX-12838	D706-D506	ovarian	high_ctDNA_cancer	post-treatment	NA
288	Pr1	Pr1_2	SLX-12838	D706-D507	prostate	high_ctDNA_cancer	post-treatment	NA
289	C5	C5_2	SLX-12838	D706-D508	colorectal	high_ctDNA_cancer	post-treatment	NA
290	B33	B33_1	SLX-15332	D707-D505	breast	high_ctDNA_cancer	baseline	NA
291	B34	B34_1	SLX-15332	D707-D506	breast	high_ctDNA_cancer	baseline	NA
292	B35	B35_1	SLX-15332	D707-D508	breast	high_ctDNA_cancer	baseline	NA
293	B36	B36_1	SLX-15332	D708-D505	breast	high_ctDNA_cancer	baseline	NA
294	B37	B37_1	SLX-15332	D708-D506	breast	high_ctDNA_cancer	baseline	NA
295	B38	B38_1	SLX-15332	D708-D507	breast	high_ctDNA_cancer	baseline	NA
296	B39	B39_1	SLX-15332	D709-D502	breast	high_ctDNA_cancer	baseline	NA
297	B40	B40_1	SLX-15332	D708-D508	breast	high_ctDNA_cancer	baseline	NA
298	B41	B41_1	SLX-15332	D709-D501	breast	high_ctDNA_cancer	baseline	NA
299	B42	B42_1	SLX-15332	D709-D503	breast	high_ctDNA_cancer	baseline	NA
300	B43	B43_1	SLX-15332	D709-D504	breast	high_ctDNA_cancer	baseline	NA
301	B44	B44_1	SLX-13227	D704-D506	breast	high_ctDNA_cancer	baseline	NA
302	B45	B45_1	SLX-13227	D704-D508	breast	high_ctDNA_cancer	baseline	NA
303	B46	B46_1	SLX-13227	D705-D506	breast	high_ctDNA_cancer	baseline	NA
304	B47	B47_1	SLX-13227	D701-D502	breast	high_ctDNA_cancer	baseline	NA
305	B48	B48_1	SLX-13227	D701-D504	breast	high_ctDNA_cancer	baseline	NA
306	B49	B49_1	SLX-13227	D702-D502	breast	high_ctDNA_cancer	baseline	NA
307	B50	B50_1	SLX-13227	D702-D504	breast	high_ctDNA_cancer	baseline	NA

TABLE 1-continued

summary table of the samples and patients included in the study								
index	patient	sample	SLX	barcode	cancer	cancer_type	timepoint	RECIST_volume
308	B51	B51_1	SLX-13227	D703-D502	breast	high_ctDNA_cancer	baseline	NA
309	GB14	GB14_1	SLX-12839	D701-D501	glioblastoma	low_ctDNA_cancer	baseline	NA
310	GB15	GB15_1	SLX-12839	D701-D502	glioblastoma	low_ctDNA_cancer	baseline	NA
311	GB16	GB16_1	SLX-12839	D701-D503	glioblastoma	low_ctDNA_cancer	baseline	NA
312	GB17	GB17_1	SLX-12839	D701-D504	glioblastoma	low_ctDNA_cancer	baseline	NA
313	GB18	GB18_1	SLX-12839	D702-D501	glioblastoma	low_ctDNA_cancer	baseline	NA
314	GB19	GB19_1	SLX-12839	D702-D502	glioblastoma	low_ctDNA_cancer	baseline	NA
315	GB20	GB20_1	SLX-12839	D702-D503	glioblastoma	low_ctDNA_cancer	baseline	NA
316	GB21	GB21_1	SLX-12839	D702-D504	glioblastoma	low_ctDNA_cancer	baseline	NA
317	GB22	GB22_1	SLX-12839	D703-D505	glioblastoma	low_ctDNA_cancer	baseline	NA
318	GB23	GB23_1	SLX-12839	D703-D506	glioblastoma	low_ctDNA_cancer	baseline	NA
319	GB24	GB24_1	SLX-12839	D704-D505	glioblastoma	low_ctDNA_cancer	baseline	NA
320	GB25	GB25_1	SLX-12839	D704-D506	glioblastoma	low_ctDNA_cancer	baseline	NA
321	GB26	GB26_1	SLX-12839	D705-D507	glioblastoma	low_ctDNA_cancer	baseline	NA
322	GB27	GB27_1	SLX-12839	D705-D508	glioblastoma	low_ctDNA_cancer	baseline	NA
323	GB28	GB28_1	SLX-12839	D704-D507	glioblastoma	low_ctDNA_cancer	baseline	NA
324	GB29	GB29_1	SLX-12839	D704-D508	glioblastoma	low_ctDNA_cancer	baseline	NA
325	GB30	GB30_1	SLX-12839	D705-D501	glioblastoma	low_ctDNA_cancer	baseline	NA
326	GB31	GB31_1	SLX-12839	D705-D502	glioblastoma	low_ctDNA_cancer	baseline	NA
327	GB32	GB32_1	SLX-12839	D705-D503	glioblastoma	low_ctDNA_cancer	baseline	NA
328	GB33	GB33_1	SLX-12839	D706-D501	glioblastoma	low_ctDNA_cancer	baseline	NA
329	GB34	GB34_1	SLX-12839	D706-D502	glioblastoma	low_ctDNA_cancer	baseline	NA
330	GB35	GB35_1	SLX-12839	D706-D503	glioblastoma	low_ctDNA_cancer	baseline	NA
331	batch2_ctl1	batch2_ctl1	SLX-13222	D701-D501	healthy	healthy	baseline	NA
332	batch2_ctl2	batch2_ctl2	SLX-13222	D701-D502	healthy	healthy	baseline	NA
333	batch2_ctl3	batch2_ctl3	SLX-13222	D701-D503	healthy	healthy	baseline	NA
334	batch2_ctl4	batch2_ctl4	SLX-13222	D701-D504	healthy	healthy	baseline	NA
335	batch2_ctl5	batch2_ctl5	SLX-13222	D702-D501	healthy	healthy	baseline	NA
336	batch2_ctl6	batch2_ctl6	SLX-13222	D702-D502	healthy	healthy	baseline	NA
337	batch2_ctl7	batch2_ctl7	SLX-13222	D702-D503	healthy	healthy	baseline	NA
338	batch2_ctl8	batch2_ctl8	SLX-13222	D702-D504	healthy	healthy	baseline	NA
339	batch2_ctl9	batch2_ctl9	SLX-13222	D703-D501	healthy	healthy	baseline	NA
340	batch2_ctl10	batch2_ctl10	SLX-13222	D703-D502	healthy	healthy	baseline	NA
341	batch2_ctl11	batch2_ctl11	SLX-13222	D703-D503	healthy	healthy	baseline	NA
342	batch2_ctl12	batch2_ctl12	SLX-13222	D703-D504	healthy	healthy	baseline	NA
343	batch2_ctl13	batch2_ctl13	SLX-13222	D704-D505	healthy	healthy	baseline	NA
344	batch2_ctl14	batch2_ctl14	SLX-13222	D704-D506	healthy	healthy	baseline	NA
345	batch2_ctl15	batch2_ctl15	SLX-13222	D704-D507	healthy	healthy	baseline	NA
346	batch2_ctl16	batch2_ctl16	SLX-13222	D704-D508	healthy	healthy	baseline	NA
347	batch2_ctl17	batch2_ctl17	SLX-13222	D705-D505	healthy	healthy	baseline	NA
348	batch2_ctl18	batch2_ctl18	SLX-13222	D705-D506	healthy	healthy	baseline	NA
349	batch2_ctl19	batch2_ctl19	SLX-13222	D705-D507	healthy	healthy	baseline	NA
350	batch2_ctl20	batch2_ctl20	SLX-13222	D705-D508	healthy	healthy	baseline	NA
351	batch2_ctl21	batch2_ctl21	SLX-13222	D706-D505	healthy	healthy	baseline	NA
352	batch2_ctl22	batch2_ctl22	SLX-13222	D706-D506	healthy	healthy	baseline	NA
353	batch2_ctl23	batch2_ctl23	SLX-13222	D706-D507	healthy	healthy	baseline	NA
354	batch2_ctl24	batch2_ctl24	SLX-13222	D706-D508	healthy	healthy	baseline	NA

TABLE 2

values for 9 fragmentation features determined from shallow Whole Genome Sequencing (sWGS) data for the samples included in the study.

index	patient	sample	SLX	barcode	cancer	tMAD	MAF	amplitude_10_bp
1	GB2	GB2_1	SLX-11868	D710-D505	glioblastoma	NA	NA	8,288894
2	GB3	GB3_1	SLX-11868	D710-D506	glioblastoma	NA	NA	7,066083
3	GB4	GB4_1	SLX-11868	D710-D507	glioblastoma	NA	NA	11,734284
4	GB5	GB5_1	SLX-11868	D710-D508	glioblastoma	NA	NA	7,039499
5	GB6	GB6_1	SLX-11868	D711-D505	glioblastoma	NA	NA	11,29576
6	GB7	GB7_1	SLX-11868	D711-D506	glioblastoma	NA	NA	8,584404
7	GB8	GB8_1	SLX-11868	D711-D507	glioblastoma	NA	NA	6,550569
8	GB9	GB9_1	SLX-11868	D711-D508	glioblastoma	NA	NA	6,966088
9	GB10	GB10_1	SLX-11868	D712-D505	glioblastoma	NA	NA	8,034286
10	GB11	GB11_1	SLX-11868	D712-D506	glioblastoma	NA	NA	6,35459
11	GB12	GB12_1	SLX-11868	D712-D507	glioblastoma	NA	NA	9,182074
12	GB13	GB13_1	SLX-11868	D712-D508	glioblastoma	NA	NA	5,20761
13	Other1	Os1_1	SLX-11870	D707-D505	esophageal junction	0.00662352	0.001	7,951253
14	B1	B1_1	SLX-11034	A019	breast	0.25477547	0.355	21,5673
15	L1	L1_1	SLX-11870	D711-D504	lung	0.14086039	0.21	22,320015
16	Ov1	Ov1_1	SLX-11870	D712-D502	ovarian	0.01414883	0	8,014098

TABLE 2-continued

values for 9 fragmentation features determined from shallow Whole Genome Sequencing (sWGS) data for the samples included in the study.						
17	Ov2	Ov2_1	SLX-11870	D708-D505	ovarian	0.0069475
18	Ren1	Ren1_1	SLX-11870	D708-D507	renal	0.01326047
19	B2	B2_1	SLX-11870	D710-D501	breast	0.00749228
20	L2	L2_1	SLX-11870	D712-D504	lung	0.00857841
21	L3	L3_1	SLX-11870	D712-D503	lung	0.10416469
22	T1	T1_1	SLX-11870	D709-D506	thymoma	0.04634427
23	R1	R1_1	SLX-11870	D710-D504	rectum	0.19414737
24	B3	B3_1	SLX-11870	D711-D502	breast	0.50279607
25	L4	L4_1	SLX-13710	D708-D508	lung	0.009
26	R2	R2_1	SLX-13710	D707-D502	rectum	0.00763274
27	B4	B4_1	SLX-13710	D706-D503	breast	0.18705825
28	P1	P1_1	SLX-13710	D705-D504	pancreatic	0.00595467
29	Ov3	Ov3_1	SLX-13710	D704-D505	ovarian	0.01732876
30	B5	B5_1	SLX-13710	D702-D507	breast	0.17913012
31	B6	B6_1	SLX-13710	D701-D508	breast	0.08931304
32	L5	L5_1	SLX-12841	D701-D501	lung	0.06389893
33	ChC1	ChC1_1	SLX-12841	D701-D502	cholangio-carcinoma	0.00692924
34	B7	B7_1	SLX-12841	D701-D503	breast	0.06720376
35	C1	C1_1	SLX-12841	D701-D504	colorectal	0.04858582
36	ChC2	ChC2_1	SLX-12841	D702-D501	cholangio-carcinoma	0.03907079
37	HCC1	HCC1_1	SLX-12841	D702-D502	hepatocellular	0.04818769
38	C2	C2_1	SLX-12841	D702-D503	colorectal	0.00692044
39	P2	P2_1	SLX-12841	D702-D504	pancreatic	0.0070876
40	ChC3	ChC3_1	SLX-12841	D703-D505	cholangio-carcinoma	0.04646124
41	P3	P3_1	SLX-12841	D703-D506	pancreatic	0.02184309
42	R3	R3_1	SLX-12841	D703-D507	rectum	0.12517655
43	ChC4	ChC4_1	SLX-12841	D703-D508	cholangio-carcinoma	NA
44	ChC5	ChC5_1	SLX-12841	D704-D505	cholangio-carcinoma	0.17356419
45	P4	P4_1	SLX-12841	D704-D506	pancreatic	0.01773972
46	C3	C3_1	SLX-12841	D704-D507	colorectal	0.14143417
47	Ov4	Ov4_1	SLX-12841	D704-D508	ovarian	0.017
48	Ov5	Ov5_1	SLX-12841	D705-D501	ovarian	0.03797909
49	B8	B8_1	SLX-12841	D705-D502	breast	0.0223823
50	L6	L6_1	SLX-12841	D705-D503	lung	0.06512785
51	C4	C4_1	SLX-12841	D705-D504	colorectal	0.40146873
52	Pe1	Pe1_1	SLX-12841	D706-D501	penile	0.0242622
53	Pr1	Pr1_1	SLX-12841	D706-D502	prostate	0.01561834
54	Ce1	Ce1_1	SLX-12841	D706-D503	cervical	0.07434257
55	C5	C5_1	SLX-12841	D706-D504	colorectal	0.05664277
56	Ov6	Ov6_1	SLX-12841	D707-D505	ovarian	0.16596734
57	En1	En1_1	SLX-12841	D707-D506	endometrial	0.0418592
58	C6	C6_1	SLX-12841	D707-D507	colorectal	0.02161484
59	C7	C7_1	SLX-12841	D707-D508	colorectal	0.03247175
60	OV04-77	JBLAB_5688	SLX-13223	D701-D501	ovarian	0.19930844
61	OV04-77	JBLAB_5689	SLX-13223	D701-D502	ovarian	0.02929487
62	OV04-83	JBLAB_5203	SLX-13223	D703-D501	ovarian	0.05179566
63	OV04-83	JBLAB_5205	SLX-13223	D703-D502	ovarian	0.017
64	OV04-122	JBLAB_5712	SLX-13223	D701-D503	ovarian	0.20397411
65	OV04-122	JBLAB_5713	SLX-13223	D701-D504	ovarian	0.011
66	OV04-141	JBLAB_5392	SLX-13223	D703-D503	ovarian	0.2039022
67	OV04-141	JBLAB_5393	SLX-13223	D703-D504	ovarian	0.02154792
68	OV04-143	JBLAB_5587	SLX-11873	D707-D501	ovarian	0.05706915
69	OV04-143	JBLAB_5588	SLX-11873	D707-D502	ovarian	0.01
70	OV04-180	JBLAB_5432	SLX-13223	D705-D505	ovarian	0.07421503
71	OV04-180	JBLAB_5433	SLX-13223	D705-D506	ovarian	0.00647481
72	OV04-211	JBLAB_5471	SLX-13223	D706-D505	ovarian	0.04274618
73	OV04-211	JBLAB_5472	SLX-13223	D706-D506	ovarian	0.00853438
74	OV04-226	JBLAB_5507	SLX-13223	D704-D505	ovarian	0.03174241
75	OV04-226	JBLAB_5508	SLX-13223	D704-D506	ovarian	0.011
76	OV04-264	JBLAB_5622	SLX-11873	D707-D503	ovarian	0.22037788
77	OV04-264	JBLAB_5623	SLX-11873	D707-D504	ovarian	0.02013793
78	OV04-292	JBLAB_5742	SLX-13223	D702-D501	ovarian	0.04971341
79	OV04-292	JBLAB_5743	SLX-13223	D702-D502	ovarian	0.06534916
80	OV04-295	JBLAB_5420	SLX-13223	D705-D507	ovarian	0.25240821
81	OV04-295	JBLAB_5422	SLX-13223	D705-D508	ovarian	0.00713784
82	OV04-297	JBLAB_5288	SLX-13223	D704-D507	ovarian	0.06130302
83	OV04-297	JBLAB_5289	SLX-13223	D704-D508	ovarian	0.0212589
84	OV04-300	JBLAB_5754	SLX-13223	D702-D503	ovarian	0.19251179
85	OV04-300	JBLAB_5755	SLX-13223	D702-D504	ovarian	0.15867713
86	X76	X76_T1_pre	SLX-13621	D701-D501	ovarian	0.02212855
87	X75_2	X75_T13_pre	SLX-13621	D702-D501	ovarian	0.00516137
88	X52	X52_T1_pre	SLX-13621	D703-D501	ovarian	0.00569295

TABLE 2-continued

values for 9 fragmentation features determined from shallow Whole Genome Sequencing (sWGS) data for the samples included in the study.						
89 X150	X150_T1_pre	SLX-13621	D704-D501	ovarian	0.00567981	0
90 X129	X129_T8_pre	SLX-13621	D705-D501	ovarian	0.00801224	0.0087
91 X57	X57_T1_pre	SLX-13621	D706-D501	ovarian	0.00538757	0.0045
92 X73	X73_T3B_pre	SLX-13621	D707-D501	ovarian	0.00590527	0.0026
93 JG090	JG090_T6_12_pre	SLX-13621	D708-D501	ovarian	0.30281177	0.0035
94 X145	X145_T8_pre	SLX-13621	D709-D501	ovarian	0.04365296	0.0815
95 X112	X112_T1_pre	SLX-13621	D710-D501	ovarian	0.00530119	0.0011
96 X75_1	X75_T1_pre	SLX-13621	D711-D501	ovarian	0.01	0.0041
97 X72	X72_T1_pre	SLX-13621	D712-D501	ovarian	0.00541364	0.0021
98 X74	X74_T1_pre	SLX-13621	D701-D502	ovarian	0.01631991	0.051
99 X127	X127_T1_pre	SLX-13621	D702-D502	ovarian	0.01	0.0085
100 X30	X30_T1_pre	SLX-13621	D703-D502	ovarian	0.01369393	0.0325
101 JBLAB_5180	JBLAB_5180_pre	SLX-13621	D704-D502	ovarian	0.00451049	0.000868
102 JBLAB_5027	JBLAB_5027_pre	SLX-13621	D705-D502	ovarian	0.00636608	0
103 JBLAB_5595	JBLAB_5595_pre	SLX-13621	D706-D502	ovarian	0.00674627	0.001
104 JBLAB_5599	JBLAB_5599_pre	SLX-13621	D707-D502	ovarian	0.00587396	0.00015
105 JBLAB_5611	JBLAB_5611_pre	SLX-13621	D708-D502	ovarian	0.02116335	NA
106 JBLAB_5477	JBLAB_5477_pre	SLX-13621	D709-D502	ovarian	0.00767838	0.0035
107 JBLAB_5632	JBLAB_5632_pre	SLX-13621	D710-D502	ovarian	0.00817832	NA
108 B9	B9_1	SLX-11043	D705-D506	breast	0.08182814	0
109 B10	B10_1	SLX-11043	D702-D501	breast	0.0144354	0.0336
110 B11	B11_1	SLX-11043	D701-D501	breast	0.013	0.14
111 B12	B12_1	SLX-11043	D705-D508	breast	0.00826536	NA
112 B13	B13_1	SLX-11043	D704-D508	breast	0.00851616	NA
113 B14	B14_1	SLX-11043	D704-D505	breast	0.0083561	NA
114 B15	B15_1	SLX-11043	D703-D503	breast	0.016	NA
115 B16	B16_1	SLX-11042	D703-D508	breast	0.02232398	NA
116 B17	B17_1	SLX-11042	D704-D504	breast	0.03101881	NA
117 B18	B18_1	SLX-11042	D704-D502	breast	0.00787396	NA
118 B19	B19_1	SLX-11042	D705-D502	breast	0.011	NA
119 B20	B20_1	SLX-11042	D705-D504	breast	0.008	NA
120 B21	B21_1	SLX-11042	D701-D505	breast	0.01747348	NA
121 B22	B22_1	SLX-11042	D701-D507	breast	0.00567912	0
122 B23	B23_1	SLX-11042	D702-D506	breast	0.03790757	NA
123 B24	B24_1	SLX-11042	D702-D508	breast	0.02927472	NA
124 B25	B25_1	SLX-11042	D703-D506	breast	0.10663707	NA
125 B26	B26_1	SLX-11042	D706-D502	breast	0.05045255	NA
126 B27	B27_1	SLX-11042	D706-D503	breast	0.01616385	NA
127 B28	B28_1	SLX-11042	D706-D504	breast	0.03047302	NA
128 B29	B29_1	SLX-11043	D703-D502	breast	0.01713247	0.15
129 B30	B30_1	SLX-11043	D701-D504	breast	0.01909028	0.187
130 B31	B31_1	SLX-11043	D704-D507	breast	0.021	NA
131 B32	B32_1	SLX-11042	D703-D507	breast	0.03009715	0.069
132 B11	B11_1	SLX-10991		bladder	NA	NA
133 B12	B12_1	SLX-10991		bladder	NA	NA
134 B13	B13_1	SLX-11094	D708-D501	bladder	NA	NA
135 B14	B14_1	SLX-10575	iPCRtagT014	bladder	NA	NA
136 B15	B15_1	SLX-11904	D709-D507	bladder	NA	NA
137 B16	B16_1	SLX-10572	D704-D505	bladder	NA	NA
138 B17	B17_1	SLX-10572	D708-D507	bladder	NA	NA
139 B18	B18_1	SLX-11896	D708-D504	bladder	NA	NA
140 B19	B19_1	SLX-11896	D707-D507	bladder	NA	NA
141 B110	B110_1	SLX-11896	D707-D508	bladder	NA	NA
142 B111	B111_1	SLX-11896	D709-D506	bladder	NA	NA
143 B112	B112_1	SLX-11904	D708-D504	bladder	NA	NA
144 B113	B113_1	SLX-11904	D709-D501	bladder	NA	NA
145 B114	B114_1	SLX-11986	D709-D504	bladder	NA	NA
146 B115	B115_1	SLX-10572	D708-D508	bladder	NA	NA
147 B116	B116_1	SLX-11896	D707-D502	bladder	NA	NA
148 B117	B117_1	SLX-10572	D708-D505	bladder	NA	NA
149 B118	B118_1	SLX-11896	D709-D503	bladder	NA	NA
150 B119	B119_1	SLX-11896	D708-D503	bladder	NA	NA
151 Ren2	Ren2_1	SLX-13900	D707-D501	renal	0.009	NA
152 Ren3	Ren3_1	SLX-13900	D707-D502	renal	0.01	NA
153 Ren4	Ren4_1	SLX-13900	D707-D503	renal	NA	NA
154 Ren5	Ren5_1	SLX-13900	D707-D504	renal	0.016	NA
155 Ren6	Ren6_1	SLX-13900	D708-D501	renal	0.011	NA
156 Ren7	Ren7_1	SLX-13900	D708-D502	renal	0.013	NA
157 Ren8	Ren8_1	SLX-13900	D708-D503	renal	0.011	NA
158 Ren9	Ren9_1	SLX-13900	D708-D504	renal	0.016	NA
159 Ren10	Ren10_1	SLX-13900	D708-D505	renal	0.021	NA
160 Ren11	Ren11_1	SLX-13900	D708-D506	renal	0.008	NA
161 Ren12	Ren12_1	SLX-13900	D708-D507	renal	0.015	NA
162 Ren13	Ren13_1	SLX-13900	D708-D508	renal	0.01	NA
163 Ren14	Ren14_1	SLX-13900	D709-D501	renal	0.017	NA
164 Ren15	Ren15_1	SLX-13900	D709-D502	renal	NA	NA
165 Ren16	Ren16_1	SLX-13900	D709-D503	renal	NA	NA

TABLE 2-continued

values for 9 fragmentation features determined from shallow Whole Genome Sequencing (sWGS) data for the samples included in the study.						
166	Ren17	Ren17_1	SLX-13900	D709-D504	renal	0.01
167	Ren18	Ren18_1	SLX-13900	D709-D505	renal	NA
168	Ren19	Ren19_1	SLX-13900	D709-D506	renal	NA
169	Ren20	Ren20_1	SLX-13900	D710-D501	renal	0.01
170	Ren21	Ren21_1	SLX-13900	D710-D502	renal	0.013
171	Ren22	Ren22_1	SLX-13900	D710-D503	renal	0.011
172	Ren23	Ren23_1	SLX-13900	D710-D504	renal	0.011
173	Ren24	Ren24_1	SLX-13900	D710-D505	renal	0.009
174	Ren25	Ren25_1	SLX-13900	D710-D506	renal	0.009
175	Ren26	Ren26_1	SLX-13900	D710-D507	renal	0.01
176	Ren27	Ren27_1	SLX-13900	D710-D508	renal	0.017
177	Ren28	Ren28_1	SLX-13900	D711-D501	renal	0.012
178	Ren29	Ren29_1	SLX-13900	D711-D502	renal	NA
179	Ren30	Ren30_1	SLX-13900	D711-D503	renal	0.008
180	Ren31	Ren31_1	SLX-13900	D711-D504	renal	0.01
181	Ren32	Ren32_1	SLX-13900	D711-D505	renal	0.029
182	Ren33	Ren33_1	SLX-13900	D711-D506	renal	0.01
183	HIP_1	HIP_1	SLX-12531	D703-D501	healthy	0.01365609
184	HIP_10	HIP_10	SLX-12531	D705-D506	healthy	0.00999028
185	HIP_11	HIP_11	SLX-12531	D705-D507	healthy	0.01083427
186	HIP_12	HIP_12	SLX-12531	D705-D508	healthy	0.01109017
187	HIP_13	HIP_13	SLX-12531	D706-D505	healthy	0.01131455
188	HIP_14	HIP_14	SLX-12531	D706-D506	healthy	0.00870144
189	HIP_15	HIP_15	SLX-12531	D706-D507	healthy	0.00967468
190	HIP_16	HIP_16	SLX-12531	D706-D508	healthy	0.00967468
191	HIP_17	HIP_17	SLX-12531	D707-D501	healthy	0.01094406
192	HIP_18	HIP_18	SLX-12531	D707-D502	healthy	0.00912639
193	HIP_19	HIP_19	SLX-12531	D707-D503	healthy	0.01262082
194	HIP_2	HIP_2	SLX-12531	D703-D502	healthy	0.00692027
195	HIP_20	HIP_20	SLX-12531	D707-D504	healthy	0.01190763
196	HIP_21	HIP_21	SLX-12531	D708-D501	healthy	0.01254617
197	HIP_22	HIP_22	SLX-12531	D708-D502	healthy	0.01158689
198	HIP_23	HIP_23	SLX-12531	D708-D503	healthy	0.0100046
199	HIP_24	HIP_24	SLX-12531	D708-D504	healthy	0.00925125
200	HIP_27	HIP_27	SLX-12534	D707-D502	healthy	0.01217069
201	HIP_28	HIP_28	SLX-12534	D707-D503	healthy	0.00878362
202	HIP_29	HIP_29	SLX-12534	D707-D504	healthy	0.01030374
203	HIP_3	HIP_3	SLX-12531	D703-D503	healthy	0.01246399
204	HIP_30	HIP_30	SLX-12534	D708-D501	healthy	0.00751474
205	HIP_31	HIP_31	SLX-12534	D708-D502	healthy	0.0105142
206	HIP_32	HIP_32	SLX-12534	D708-D503	healthy	0.00923109
207	HIP_33	HIP_33	SLX-12534	D708-D504	healthy	0.00824142
208	HIP_34	HIP_34	SLX-12534	D709-D501	healthy	0.00603306
209	HIP_35	HIP_35	SLX-12534	D709-D503	healthy	0.00704468
210	HIP_36	HIP_36	SLX-12534	D709-D504	healthy	0.01441797
211	HIP_37	HIP_37	SLX-12534	D710-D501	healthy	0.00760246
212	HIP_38	HIP_38	SLX-12534	D710-D502	healthy	0.00764811
213	HIP_39	HIP_39	SLX-12534	D710-D503	healthy	0.01278262
214	HIP_4	HIP_4	SLX-12531	D703-D504	healthy	0.00885683
215	HIP_40	HIP_40	SLX-12534	D710-D504	healthy	0.0126438
216	HIP_41	HIP_41	SLX-12534	D711-D505	healthy	0.00779714
217	HIP_42	HIP_42	SLX-12534	D711-D506	healthy	0.01226728
218	HIP_43	HIP_43	SLX-12534	D711-D507	healthy	0.00886215
219	HIP_44	HIP_44	SLX-12534	D711-D508	healthy	0.01102103
220	HIP_45	HIP_45	SLX-12534	D712-D505	healthy	0.01151546
221	HIP_46	HIP_46	SLX-12534	D712-D506	healthy	0.01069675
222	HIP_47	HIP_47	SLX-12534	D712-D507	healthy	0.01326822
223	HIP_48	HIP_48	SLX-12534	D712-D508	healthy	0.01307578
224	HIP_5	HIP_5	SLX-12531	D704-D501	healthy	0.00640521
225	HIP_6	HIP_6	SLX-12531	D704-D502	healthy	0.00943859
226	HIP_7	HIP_7	SLX-12531	D704-D503	healthy	0.01017749
227	HIP_8	HIP_8	SLX-12531	D704-D504	healthy	0.0097156
228	HIP_9	HIP_9	SLX-12531	D705-D505	healthy	0.00951729
229	M1	M1_1	SLX-11379	D701-D502	melanoma	0.31468668
230	M1	M1_2	SLX-11379	D701-D501	melanoma	0.086146
231	M4	M4_1	SLX-11379	D702-D501	melanoma	0.009
232	M4	M4_2	SLX-12758	D704-D501	melanoma	0.00607225
233	M4	M4_3	SLX-12759	D708-D501	melanoma	0.01
234	M4	M4_4	SLX-12758	D709-D502	melanoma	0.0059634
235	M4	M4_5	SLX-12758	D702-D501	melanoma	0.009
236	M4	M4_6	SLX-11383	D701-D506	melanoma	0.00622659
237	M4	M4_7	SLX-11379	D701-D503	melanoma	0.008
238	M12	M12_1	SLX-11379	D703-D502	melanoma	0.06257905
239	M12	M12_2	SLX-11847	D704-D502	melanoma	0.00825359
240	M12	M12_3	SLX-11847	D704-D503	melanoma	0.02188627
241	M12	M12_4	SLX-11847	D707-D507	melanoma	0.02521355
242	M14	M14_1	SLX-11383	D708-D503	melanoma	0.01

TABLE 2-continued

values for 9 fragmentation features determined from shallow Whole Genome Sequencing (sWGS) data for the samples included in the study.						
243	M14	M14_2	SLX-12758	D706-D506	melanoma	0.03887853
244	M22	M22_1	SLX-11379	D704-D507	melanoma	0.05850595
245	M22	M22_2	SLX-12758	D706-D507	melanoma	0.00659093
246	M22	M22_3	SLX-11379	D704-D508	melanoma	0.1123879
247	M22	M22_4	SLX-12758	D704-D508	melanoma	0.11091958
248	M32	M32_1	SLX-11379	D705-D506	melanoma	0.01892249
249	M32	M32_2	SLX-11847	D705-D503	melanoma	0.013
250	C8	C8_T1	SLX-12832	D709-D501	colorectal	0.13461166
251	C8	C8_T2	SLX-12832	D709-D502	colorectal	0.02433155
252	L5	L5_T2	SLX-12832	D709-D503	lung	0.05910309
253	ChC1	ChC1_3	SLX-12832	D709-D504	cholangio- carcinoma	0.01
254	ChC1	ChC1_4	SLX-12832	D710-D501	cholangio- carcinoma	0.029
255	ChC2	ChC2_2	SLX-12832	D710-D502	cholangio- carcinoma	0.04069151
256	ChC2	ChC2_3	SLX-12832	D710-D503	cholangio- carcinoma	0.02290481
257	HCC1	HCC1_2	SLX-12832	D710-D504	hepatocellular	0.05593432
258	HCC1	HCC1_3	SLX-12832	D711-D505	hepatocellular	0.05623691
259	HCC1	HCC1_4	SLX-12832	D711-D506	hepatocellular	0.07020201
260	HCC1	HCC1_5	SLX-12832	D711-D507	hepatocellular	0.06769479
261	P2	P2_2	SLX-12832	D711-D508	pancreatic	0.00737544
262	P4	P4_2	SLX-12832	D712-D505	pancreatic	0.00845528
263	C4	C4_2	SLX-12832	D712-D506	colorectal	0.44317612
264	Pr1	Pr1_4	SLX-12832	D712-D507	prostate	0.02602964
265	Ov6	Ov6_2	SLX-12832	D712-D508	ovarian	0.23784565
266	ChC2	ChC2_6	SLX-12838	D701-D505	cholangio- carcinoma	0.02660187
267	ChC3	ChC3_2	SLX-12838	D701-D506	cholangio- carcinoma	0.01405692
268	C3	C3_5	SLX-12838	D701-D507	colorectal	0.03204027
269	L6	L6_2	SLX-12838	D701-D508	lung	0.07217697
270	Pr1	Pr1_3	SLX-12838	D702-D505	prostate	0.01337188
271	B7	B7_2	SLX-12838	D702-D506	breast	0.14971349
272	C1	C1_2	SLX-12838	D702-D507	colorectal	0.06302754
273	ChC2	ChC2_4	SLX-12838	D702-D508	cholangio- carcinoma	0.012
274	ChC2	ChC2_5	SLX-12838	D703-D501	cholangio- carcinoma	0.03388701
275	P4	P4_3	SLX-12838	D703-D502	pancreatic	0.01492043
276	C3	C3_4	SLX-12838	D703-D503	colorectal	0.02969907
277	Ov4	Ov4_2	SLX-12838	D703-D504	ovarian	0.01768853
278	Ov5	Ov5_2	SLX-12838	D704-D501	ovarian	0.03000071
279	B8	B8_2	SLX-12838	D704-D502	breast	0.01711789
280	C5	C5_3	SLX-12838	D704-D503	colorectal	0.015
281	En1	En1_2	SLX-12838	D704-D504	endometrial	0.09648123
282	C6	C6_2	SLX-12838	D705-D505	colorectal	0.01
283	ChC1	ChC1_2	SLX-12838	D705-D506	cholangio- carcinoma	0.00657679
284	C3	C3_2	SLX-12838	D705-D507	colorectal	0.14260432
285	C3	C3_3	SLX-12838	D705-D508	colorectal	0.14314493
286	Ov4	Ov4_3	SLX-12838	D706-D505	ovarian	0.00620281
287	Ov5	Ov5_3	SLX-12838	D706-D506	ovarian	0.02161473
288	Pr1	Pr1_2	SLX-12838	D706-D507	prostate	0.016
289	C5	C5_2	SLX-12838	D706-D508	colorectal	0.05837149
290	B33	B33_1	SLX-15332	D707-D505	breast	0.00834566
291	B34	B34_1	SLX-15332	D707-D506	breast	0.01937858
292	B35	B35_1	SLX-15332	D707-D508	breast	0.3099655
293	B36	B36_1	SLX-15332	D708-D505	breast	0.2510418
294	B37	B37_1	SLX-15332	D708-D506	breast	0.37214783
295	B38	B38_1	SLX-15332	D708-D507	breast	0.0073204
296	B39	B39_1	SLX-15332	D709-D502	breast	0.01750562
297	B40	B40_1	SLX-15332	D708-D508	breast	0.04741394
298	B41	B41_1	SLX-15332	D709-D501	breast	0.02476021
299	B42	B42_1	SLX-15332	D709-D503	breast	0.33542756
300	B43	B43_1	SLX-15332	D709-D504	breast	0.09644121
301	B44	B44_1	SLX-13227	D704-D506	breast	0.14065498
302	B45	B45_1	SLX-13227	D704-D508	breast	0.00602283
303	B46	B46_1	SLX-13227	D705-D506	breast	0.06773296
304	B47	B47_1	SLX-13227	D701-D502	breast	0.06050266
305	B48	B48_1	SLX-13227	D701-D504	breast	0.01216387
306	B49	B49_1	SLX-13227	D702-D502	breast	0.0714198
307	B50	B50_1	SLX-13227	D702-D504	breast	0.19923403
308	B51	B51_1	SLX-13227	D703-D502	breast	0.01111396
309	GB14	GB14_1	SLX-12839	D701-D501	glioblastoma	0.00722063
310	GB15	GB15_1	SLX-12839	D701-D502	glioblastoma	0.00999163

TABLE 2-continued

values for 9 fragmentation features determined from shallow Whole Genome Sequencing (sWGS) data for the samples included in the study.

			P(20_150)/	P(100_150)/		P(20_150)/		
index	P(20_150)	P(160_180)	P(160_180)	P(100_150)	P(163_169)	P(180_220)	P(250_320)	P(180_220)
1	0.15593628	0.474759905	0.328452926	0.150467716	0.797230669	0.242257259	0.01344566	0.643680527
2	0.153305045	0.517651152	0.296155132	0.151170814	0.716017076	0.210272306	0.003406292	0.729078633
3	0.190293559	0.42569701	0.447016433	0.17607343	0.1097693843	0.242204427	0.031265598	0.785673333
4	0.153458877	0.532513429	0.288178417	0.151180028	0.676444354	0.191137675	0.0031301	0.802870897
5	0.234162421	0.481611843	0.486205695	0.228334925	0.1022916892	0.123317396	0.003111174	1.898859598
6	0.182383923	0.500662425	0.364285222	0.178274903	0.797659738	0.170358622	0.004612955	1.070588159
7	0.125970767	0.435908671	0.288984311	0.119404664	0.735520087	0.300845212	0.026021089	0.418722857
8	0.150216458	0.506601991	0.296517702	0.146539351	0.685757089	0.211747161	0.004645669	0.710330082
9	0.150859409	0.445272059	0.338802775	0.143907267	0.829470351	0.268316043	0.021243303	0.562245205
10	0.134771126	0.507443882	0.265588237	0.132198502	0.669215178	0.239142669	0.007797143	0.563559513
11	0.168015932	0.470466497	0.357126242	0.16075196	0.871359965	0.253543337	0.012234413	0.713645461
12	0.119421664	0.516409351	0.231253876	0.117728563	0.590052228	0.255304555	0.004916546	0.467761588
13	0.144461769	0.384670633	0.375546654	0.13414291	0.931861214	0.227571414	0.039721361	0.63479752
14	0.270943962	0.403405095	0.671642389	0.24852795	1.472004778	0.132617318	0.030309159	2.043051131
15	0.333745777	0.341029821	0.978641035	0.316934675	2.213760282	0.124345996	0.036590701	2.684009046
16	0.258242277	0.321164069	0.804082093	0.237227487	1.700196737	0.112764896	0.113684512	2.290094577
17	0.161376514	0.472136335	0.341800667	0.157467767	0.80410171	0.180049905	0.024531581	0.896287691
18	0.155759138	0.432486714	0.360147798	0.150598877	0.853785668	0.186174631	0.040511621	0.836629232
19	0.159149606	0.457356112	0.347977433	0.155320273	0.83295436	0.18830049	0.028989014	0.845189548
20	0.161875577	0.441582658	0.366580467	0.156550382	0.880567715	0.203818787	0.027965978	0.79421323
21	NA	NA	NA	NA	NA	NA	NA	NA
22	0.406794901	0.271498664	1.498331135	0.353708315	2.778961716	0.07056407	0.076066888	5.764901312
23	0.410998565	0.31613605	1.300068642	0.348798067	2.281498855	0.064663009	0.050209365	6.356007447
24	0.161643441	0.443226021	0.364697543	0.157566047	0.868614807	0.184260953	0.038850541	0.877252821
25	0.156543642	0.484951752	0.322802508	0.149583553	0.778941478	0.215896693	0.016460595	0.725085873
26	0.183928705	0.453968867	0.405157089	0.176759494	0.97390446	0.195448459	0.022055332	0.941059886
27	0.178035293	0.432936842	0.411226941	0.171816011	1.030458873	0.214017304	0.029019478	0.831873358
28	0.211253249	0.485403371	0.435211747	0.203936205	0.952144517	0.137147125	0.017063151	1.540340336
29	0.183987884	0.444977085	0.413477211	0.17537584	0.990412234	0.208852123	0.022703094	0.880948117

TABLE 2-continued

values for 9 fragmentation features determined from shallow Whole Genome Sequencing (sWGS) data for the samples included in the study.

338	0.179043752	0.40279764	0.4445005	0.156516401	0.995156689	0.195795219	0.032222904	0.914443945
339	0.153789797	0.384579996	0.399890481	0.139614476	0.939119253	0.197631967	0.043190347	0.778162973
340	0.151687512	0.390862158	0.388084415	0.138701569	0.917641434	0.198408563	0.042495211	0.764520995
341	0.202828454	0.38526313	0.526467337	0.183556194	1.223146834	0.19508449	0.031409195	1.039695435
342	0.214794203	0.376699468	0.570200447	0.181221462	1.238168985	0.185805333	0.031992862	1.156017439
343	0.137307321	0.429936632	0.319366415	0.119809959	0.738055807	0.231063427	0.02681307	0.594240823
344	0.133695268	0.428157382	0.312257299	0.117586519	0.727508555	0.231158515	0.027550362	0.578370509
345	0.132021502	0.432975105	0.304917073	0.1117675672	0.72313081	0.235034549	0.026771733	0.561711042
346	0.163394787	0.402317221	0.406134209	0.151179887	0.978981206	0.210723204	0.034161649	0.775400067
347	0.163034832	0.396959389	0.410709098	0.15057059	0.986908486	0.207057515	0.035907789	0.78738911
348	0.215344642	0.446933098	0.481827467	0.186152356	0.992390381	0.152964257	0.023425551	1.407810211
349	0.134232226	0.4290997862	0.312824271	0.119169742	0.737495251	0.232784327	0.027005169	0.576637731
350	0.209426922	0.446487965	0.469053902	0.181170839	0.967777005	0.155600599	0.025020308	1.345926195
351	0.110719663	0.374297401	0.295806657	0.088920222	0.672869906	0.274581447	0.030802323	0.403230676
352	0.229414294	0.424955266	0.539855161	0.20464533	1.16131228	0.160407165	0.028328176	1.430199789
353	0.222260906	0.41089728	0.540915982	0.196935416	1.179464903	0.170880962	0.028556532	1.300676818
354	0.224566246	0.454347959	0.494260492	0.202018713	1.027491671	0.139032167	0.025108809	1.615210712

TABLE 3

t-MAD score for the 48 plasma samples of the OV04 cohort before and after in vitro size selection.

index	SLXID	binSize	control	Sample Names	median TP53 MAF	median_tMAD_no_size_selection	selection	treatment	patient	median_tMAD_with_fold_size_selection	enrichment
1	SLX-11873	30	K5042	R146	0.232	0.057069147	no	before	OV04-143	0.087364547	1.530854264
2	SLX-11873	30	K5042	R147	0.022	0.012773248	no	post	OV04-143	0.028316869	2.216888688
3	SLX-11873	30	K5042	R148	0.514	0.220377876	no	before	OV04-264	0.258905932	1.174827241
4	SLX-11873	30	K5042	R149	0.034	0.020137929	no	post	OV04-264	0.067751424	3.364368997
7	SLX-13223	30	K5042	JBLAB_5688	0.346385	0.199308443	no	before	OV04-77	0.266627416	1.337762776
8	SLX-13223	30	K5042	JBLAB_5689	0.068603	0.029294865	no	post	OV04-77	0.055629976	1.898966798
9	SLX-13223	30	K5042	JBLAB_5712	0.483385	0.203974112	no	before	OV04-122	0.210309045	1.031057534
10	SLX-13223	30	K5042	JBLAB_5713	0.036652	0.012782907	no	post	OV04-122	0.080429849	6.29198421
11	SLX-13223	30	K5042	JBLAB_5742	0.14797	0.049713406	no	before	OV04-292	0.063867761	1.284719076
12	SLX-13223	30	K5042	JBLAB_5743	0.069141	0.065349155	no	post	OV04-292	0.123748162	1.893645939
13	SLX-13223	30	K5042	JBLAB_5754	0.266115	0.192511793	no	before	OV04-300	0.171876244	0.89280891
14	SLX-13223	30	K5042	JBLAB_5755	0.03915	0.15867713	no	post	OV04-300	0.171629671	1.081628279
15	SLX-13223	30	K5042	JBLAB_5203	0.2712105	0.05179566	no	before	OV04-83	0.139343378	2.690252002
16	SLX-13223	30	K5042	JBLAB_5205	0.0687565	0.011382743	no	post	OV04-83	0.072524334	6.371428574
17	SLX-13223	30	K5042	JBLAB_5342	0.610217	0.203902197	no	before	OV04-141	0.259249767	1.271441754
18	SLX-13223	30	K5042	JBLAB_5343	0.064836	0.021547924	no	post	OV04-141	0.105868625	4.913170522
19	SLX-13223	30	K5042	JBLAB_5507	0.123199135	0.031742405	no	before	OV04-226	0.062392469	1.965587327
20	SLX-13223	30	K5042	JBLAB_5508	0.022327219	0.011923695	no	post	OV04-226	0.033677313	2.824402419
21	SLX-13223	30	K5042	JBLAB_5288	0.20705	0.061303019	no	before	OV04-297	0.168597772	2.750236036
22	SLX-13223	30	K5042	JBLAB_5289	0.092029	0.0212589	no	post	OV04-297	0.05805594	2.73090047
23	SLX-13223	30	K5042	JBLAB_5432	0.212771398	0.074215033	no	before	OV04-180	0.210353293	2.834375793
24	SLX-13223	30	K5042	JBLAB_5433	0.001046472	0.006474814	no	post	OV04-180	0.011753831	1.815315621
25	SLX-13223	30	K5042	JBLAB_5420	0.5065815	0.252408213	no	before	OV04-295	0.399111409	1.581214035
26	SLX-13223	30	K5042	JBLAB_5422	0.0124825	0.007137838	no	post	OV04-295	0.023034569	3.227107284

TABLE 3-continued

t-MAD score for the 48 plasma samples of the OV04 cohort before and after in vitro size selection.											
index	SLXID	binSize	control	Sample Names	median TP53 MAF	median_tMAD_no_size_selection	selection	treatment	patient	median_tMAD_with_fold_size_selection	enrichment
27	SLX-13223	30	K5042 310_1	JBLAB_5471	0.082816831	0.04274618	no	before	OV04-211	0.047433825	1.109662314
28	SLX-13223	30	K5042 310_1	JBLAB_5472	0.008998983	0.008534381	no	post	OV04-211	0.014143088	1.657189666
29	SLX-13621	30	K5042 310_1	X76_T1_pre	0	0.022128547	no		OV04-76	0.041468333	1.873974509
30	SLX-13621	30	K5042 310_1	X75_T13_pre	0.0007705	0.005161371	no		OV04-75	0.01079341	2.0911905
31	SLX-13621	30	K5042 310_1	X52_T1_pre	0.0024735	0.005692945	no		OV04-52	0.019834069	3.483973409
32	SLX-13621	30	K5042 310_1	X150_T1_pre	0	0.005679811	no		OV04-150	0.014364408	2.529029223
33	SLX-13621	30	K5042 310_1	X129_T8pre	0.00119	0.008012243	no		OV04-129	0.015789503	1.970672008
34	SLX-13621	30	K5042 310_1	X57_T1_pre	0.00119	0.005387574	no		OV04-57	0.014437579	2.67979224
35	SLX-13621	30	K5042 310_1	X73_T3B_pre	0.0021	0.005905265	no		OV04-73	0.014933244	2.528801671
36	SLX-13621	30	K5042 310_1	JG090_T612_pre	0.003092	0.302811769	no		JG090	0.423426811	1.39831689
37	SLX-13621	30	K5042 310_1	X145_T8_pre	0	0.043652958	no		OV04-145	0.116005436	2.657447314
38	SLX-13621	30	K5042 310_1	X112_T1_pre	0	0.005301188	no		OV04-112	0.011067067	2.087657899
39	SLX-13621	30	K5042 310_1	X75_T1_pre	0.0041885	0.008682287	no		OV04-75	0.021401469	2.464957562
40	SLX-13621	30	K5042 310_1	X72_T1_pre	0	0.005413644	no		OV04-72	0.022785962	4.208987883
41	SLX-13621	30	K5042 310_1	X74_T1_pre	0.001392	0.016319911	no		OV04-74	0.063135101	3.868593462
42	SLX-13621	30	K5042 310_1	X127_T1_pre	0.0022355	0.008930611	no		OV04-127	0.026903941	3.012553228
43	SLX-13621	30	K5042 310_1	X30_T1_pre	0.032437	0.013693931	no		OV04-30	0.037435405	2.733722333
44	SLX-13621	30	K5042 310_1	JBLAB.5180_pre	0	0.004510492	no		JBLAB.5180	0.017007543	3.770662491
45	SLX-13621	30	K5042 310_1	JBLAB.5027_pre	0	0.006366084	no		JBLAB.5027	0.012995165	2.04131221
46	SLX-13621	30	K5042 310_1	JBLAB.5595_pre	0	0.006746273	no		JBLAB.5595	0.020444819	3.030535379
47	SLX-13621	30	K5042 310_1	JBLAB.5599_pre	0	0.005873961	no		JBLAB.5599	0.00810866	1.380441579
48	SLX-13621	30	K5042 310_1	JBLAB.5611_pre	0.045	0.021163354	no		JBLAB.5611	0.033449519	1.580539597
49	SLX-13621	30	K5042 310_1	JBLAB.5477_pre	0	0.007678384	no		JBLAB.5477	0.036978881	4.815971824
50	SLX-13621	30	K5042 310_1	JBLAB.5632_pre	0	0.008178321	no		JBLAB.5632	0.014573466	1.78196307

Our results indicate that exploiting fundamental properties of cfDNA with fragment specific analyses can provide more sensitive analysis of ctDNA. We based the selection criteria on a biological observation that ctDNA fragment size distribution is shifted from normal cfDNA. Our work builds on a comprehensive survey of plasma cfDNA fragmentation patterns across 200 patients with multiple cancer types and 65 healthy individuals. We identified features that could determine the presence and amount of ctDNA in plasma samples, without a priori knowledge of somatic aberrations. Although this catalogue is the first of its kind, we note that it employed double-stranded DNA from plasma samples, and is subject to potential biases incurred by the DNA extraction and sequencing methods we used. Additional biological effects could contribute to further selective analysis of cfDNA. Other bodily fluids (urine, cerebrospinal fluid, saliva), different nucleic acids and structures, altered mechanisms of release into circulation, or sample processing methods could exhibit varying fragment size signatures and could offer additional exploitable biological patterns for selective sequencing.

Previous work has reported the size distributions of mutant ctDNA, but only considered limited genomic loci, cancer types, or cases (30, 32, 33). We identified the size differences between mutant and non-mutant DNA on a genome-wide and pan-cancer scale. We developed a method to size mutant ctDNA without using high-depth WGS. By sequencing >150 mutations per patient at high depth we obtained large numbers of reads that could be unequivocally identified as tumor-derived, and thus determined the size distribution of mutant ctDNA and non-mutant cfDNA in cancer patients. A potential limitation of our approach is that capture-based sequencing is biased by probe capture efficiency and therefore our data may not accurately reflect ctDNA fragments <100 bp or >300 bp.

Our work provides strong evidence that the modal size of ctDNA for many cancer types is less than 167 bp, which is the length of DNA wrapped around the chromatosome. In addition, our work also shows that there is a high level of enrichment of mutant DNA fragments at sizes greater than 167 bp, notably in the range 250-320 bp. These longer fragments may explain previous observations that longer ctDNA can be detected in the plasma of cancer patients (29, 32). The origin of these long fragments is still unknown, and their observation could be linked to technical factors. However, it is likely that mechanisms of compaction and release of cfDNA into circulation, which may differ depending on its origin, will be reflected by different fragment sizes (38). Improving the characterization of these fragments will be important, especially for future work combining ctDNA analysis with other entities in blood such as microvesicles and tumor-educated platelets (39, 40). Fragment specific analyses not only increase the sensitivity for detection of rare mutations, but could be used to track modifications in the size distribution of ctDNA. Future work should address whether this approach could be used to elucidate mechanistic effects of treatment on tumor cells, for example by distinguishing between necrosis and apoptosis based on fragment size (41).

Genome-wide and exome sequencing of plasma DNA at multiple time-points during cancer treatment have been proposed as non-invasive means to study cancer evolution and for the identification of possible resistance mechanisms to treatment (3). However, WGS and WES approaches are costly and have thus far been applicable only in samples for

which the tumor DNA fraction was >5-10% (3-5, 42). We demonstrated that we could exploit the differences in fragment lengths using in vitro and in silico size selection to enrich for tumor content in plasma samples which improved mutation and SCNA detection in sWGS and WES data. We demonstrated that size selection improved the detection of mutations that are present in plasma at low allelic fractions, while maintaining low sequencing depth by sWGS and WES. Size selection can be achieved with simple means and at low cost, and is compatible with a wide range of downstream genome-wide and targeted genomic analyses, greatly increasing the potential value and utility of liquid biopsies.

Size selection can be applied in silico, which incurs no added costs, or in vitro, which adds a simple and low-cost intermediate step that can be applied to either the extracted DNA or the libraries created from it. This approach, applied prospectively to new studies, could boost the clinical utility of ctDNA detection and analysis, and creates an opportunity for re-analysis of large volumes of existing data (4, 34, 43). The limitation of this technique is a potential loss of material and information, since some of the informative fragments may be found in size ranges that are filtered out or deprioritized in the analysis. This may be particularly problematic if only a few copies of the fragments of interest are present in plasma. Despite potential loss of material, we demonstrated that classification algorithms can learn from cfDNA fragmentation features and SCNAs analysis and improve the detection of ctDNA with a cheap sequencing approach (FIG. 22). Moreover, the cfDNA fragmentation features alone can be leveraged to classify cancer and healthy samples with a high accuracy (AUC=0.989 for high ctDNA cancers, and AUC=0.891 for low ctDNA cancers) (FIG. 26).

Analysis of fragment sizes could provide improvements in other applications. Introducing fragment size information on each read could enhance mutation-calling algorithms from high depth sequencing, to identify tumor-derived mutations from other sources such as somatic variants or background sequencing noise. In addition, cfDNA analysis in patients with CHIP is likely to be structurally different from ctDNA released during tumor cell proliferation (18, 19). Thus, fragmentation analysis or selective sequencing strategies could be applied to distinguish clinically relevant tumor mutations from those present in clonal expansions of normal cells. This will be critical for the development of cfDNA-based methods for identification of patients with early stage cancer.

Size selection could also have an impact on the detection of other types of DNA in body fluids or to enrich signals for circulating bacterial or pathogen DNA and mitochondrial DNA. These DNA fragments are not associated with nucleosomes and are often highly fragmented below 100 bp. Filtering such fragments may prove to be important in light of the recently established link between the microbiome and treatment efficiency (17, 44). Moreover, recent work highlights a stronger correlation between ctDNA detection and cellular proliferation, rather than cell-death (45). We hypothesize that the mode of the distribution of ctDNA fragment sizes at 145 bp could reflect cfDNA released during cell proliferation, and the fragments at 167 bp may reflect cfDNA released by apoptosis or maturation/turnover of blood cells. The effect of other cancer hallmarks (46) on ctDNA biology, structure, concentration and release is yet unknown.

In summary, ctDNA fragment size analysis, via size selection and machine learning approaches, boosts non-invasive genomic analysis of tumor DNA. Size selection of

shorter plasma DNA fragments enriches ctDNA, and leads to the identification of a greater number of genomic alterations with both targeted and untargeted sequencing at a minimal additional cost. Combining cfDNA fragment size analysis and the detection of SCNAs with a non-linear classification algorithm improved the discrimination between samples from cancer patients and healthy individuals. As the analysis of fragment sizes is based on the structural property of ctDNA, size selection could be used with any downstream sequencing applications. Our work could help overcome current limitations of sensitivity for liquid biopsy, supporting expanded clinical and research

applications. Our results indicate that exploiting the endogenous biological properties of cfDNA provides an alternative paradigm to deeper sequencing of ctDNA.

Code

The following exemplary analysis code for the classification algorithms described in the Examples above is in the R programming environment (see www.r-project.org/about.html). The features may be taken from Table 2, wherein the samples are separated into group A cancers (“high ctDNA cancers”) and group B (“low ctDNA cancer”), and wherein healthy controls are used in each (i.e. a copy in each of the files).

```

---  

title: "PAN-CANCER classifier"  

author: "Dineika Chandrananda"  

date: "20 November 2017"  

output: html_document  

---  

# Data pre-processing  

* Separating out cancer types into Group A  

* containing "healthy", "breast", "melanoma", "ovarian", "lung",  

"colorectal", "cholangiocarcinoma"  

* and Group B the low ctDNA cancers  

* Only plasma  

* No size selection  

* Timepoints mixed (baseline and post-treatment)  

* Remove degraded DNA  

# Run feature selection and model the training data  

` ` ` {r feature selection}  

library(caret)  

library(pROC)  

MY_SEED <- 666  

filename_NO_SZ <- "./2018_Group_A_cancers_noSZ.csv"  

full_data_NO_SZ<- read.csv(filename_NO_SZ, header=TRUE,  

stringsAsFactors=FALSE)  

stopifnot(!anyNA(full_data_NO_SZ))  

# breast          cervical      cholangiocarcinoma    colorectal  

# 53              1            13  

# endometrial     healthy       hepatocellular  

lung  

# 2                65           5               7  

# melanoma        ovarian      penile           prostate  

# 18              56           1               4  

# rectum          thymoma  

# 3                1  

# partition data so that the cancerTypes + healthy are evenly  

separated  

# Use a 60:40 split in all cancer + healthy categories  

full_data_NO_SZ$cancer <- factor(full_data_NO_SZ$cancer)  

set.seed(MY_SEED)  

intrain <- createDataPartition(y=full_data_NO_SZ$cancer, p=0.6,  

list = FALSE)  

#####  

# Convert multiple cancer classes into cancer/healthy  

#####  

full_data_NO_SZ$cancer <- as.character(full_data_NO_SZ$cancer)  

full_data_NO_SZ$cancer[full_data_NO_SZ$cancer != "healthy"] <-  

"cancer"  

full_data_NO_SZ$cancer <- factor(full_data_NO_SZ$cancer,  

levels=c("healthy", "cancer"))  

#####  

names (full_data_NO_SZ) [names(full_data_NO_SZ) == "cancer"] <-  

"Class"  

# Split the test/train data sets  

neat_train <- full_data_NO_SZ[intrain,]  

neat_test <- full_data_NO_SZ[-intrain,]  

table (neat_train$Class)  

#  

# healthy          cancer  

# 39              114  

table (neat_test$Class)  

#  

# healthy          cancer  

# 26              68  

## The baseline set of predictors,  

b1 <- c("tMAD",  

"amplitude_10bp",

```

```

“P160_180”,
“P180_220”,
“P250_320”)
training <- neat_train[, c(“sample”, “Class”, b1)]
testing <- neat_test [, c(“sample”, “Class”, b1)]
saveRDS(training, “training”)
saveRDS(testing, “testing”)
predVars <- names(training) [!(names(training) %in%
  c(“sample”, “Class”))]
saveRDS(predVars, “predVars”)
## This summary function is used to evaluate the models.
fiveStats <- function(. . .) c(twoClassSummary(. . .),
  defaultSummary(. . .))
## We create the cross-validation data as a list to use with
different
## functions
index <- createMultiFolds(training$Class, times = 5)
## The candidate set of the number of predictors to evaluate
varSeq <- seq(1, length(predVars) -1)
## We can also use parallel processing to run each resampled RFE
## iteration
library(doMC)
registerDoMC(20)
set.seed(MY_SEED)
ctrl <- rfeControl(method = “repeatedcv”, repeats = 5,
  saveDetails = TRUE,
  index = index,
  returnResamp = “final”)
set.seed(MY_SEED)
fullCtrl <- trainControl (method = “repeatedcv”,
  repeats = 5,
  summaryFunction = fiveStats,
  classProbs = TRUE,
  index = index)
#####
## Fit the RFE models
#####
ctrl$functions <- rfFuncs
ctrl$functions$summary<- fiveStats
set.seed(MY_SEED)
rfRFE <- rfe(training[, predVars],
  training$Class,
  sizes = varSeq,
  metric = “ROC”,
  ntree = 1000,
  rfeControl = ctrl
) # keep.forest=TRUE
rfRFE
saveRDS(rfRFE, file=“rfRFE”)
ctrl$functions <- lrFuncs
ctrl$functions$summary <- fiveStats
set.seed(MY_SEED)
lrRFE <- rfe(training[, predVars],
  training$Class,
  sizes = varSeq,
  metric = “ROC”,
  rfeControl = ctrl)
lrRFE
saveRDS(lrRFE, file=“lrRFE”)
#####
# Plotting ROC curves for test set (high ctDNA)
library(caret)
library(pROC)
library(ggplot2)
library(randomForest)
MY_SEED <- 666
testing <- training <- lrRFE <- rfRFE <- NULL
testing <- readRDS(“testing”)
training <- readRDS(“training”)
lrRFE <- readRDS(“lrRFE”)
rfRFE <- readRDS(“rfRFE”)
predVars <- c( “tMAD”,
  “amplitude_10bp”,
  “P160_180”,
  “P180_220”,
  “P250_320”)
# Get ROC curves for the different models
#1) Only t-MAD
training_binary <- training
testing_binary <- testing

```

```

training_binary$Class <- as.character(training_binary$Class)
testing_binary$Class <- as.character(testing_binary$Class)
training_binary$Class[training_binary$Class == "healthy"] <- 0
training_binary$Class[training_binary$Class != "0"] <- 1
training_binary$Class <- factor(as.numeric(training_binary$Class))
testing_binary$Class[testing_binary$Class == "healthy"] <- 0
testing_binary$Class[testing_binary$Class != "0"] <- 1
testing_binary$Class <- factor(as.numeric(testing_binary$Class))
lr_tMAD <- glm(Class ~ tMAD,
                 data = training_binary,
                 family = binomial)
saveRDS(lr_tMAD, file="lr_tMAD")
prob <- predict(lr_tMAD, newdata=testing_binary, type="response")
pred <- ROCR::prediction(prob, testing_binary$Class)
perf <- ROCR::performance(pred, measure = "tpr", x.measure = "fpr")
tMAD_AUC <- ROCR::performance(pred, measure = "auc")@y.values[[1]]
df_tMAD <- data.frame(Sensitivity=perf@y.values[[1]],
                        Specificity=perf@y.values[[1]])
# Logistic regression, recursive feature elimination
ROC_lrRFE <- roc(testing$Class,
                   predict(lrRFE, testing[,predVars])$cancer)
df_lrRFE <- data.frame(Sensitivity=ROC_lrRFE$sensitivities,
                        Specificity=1-ROC_lrRFE$specificities)
# Random Forest RFE
library(randomForest)
ROC_rfRFE <- roc(testing$Class,
                   predict(rfRFE, testing[,predVars])$cancer,
                   levels=c("healthy", "cancer"))
ROC_rfRFE
df_rfRFE <- data.frame(Sensitivity=ROC_rfRFE$sensitivities,
                        Specificity=1-ROC_rfRFE$specificities)
# Plotting ROC curves
pdf("Model_Comparison_on_TestData_high_ctDNA.pdf")
plot(x=df_rfRFE$Specificity,
      y=df_rfRFE$Sensitivity,
      xlab="1 - Specificity",
      ylab="Sensitivity", type="l",
      col="blue")
points(x=df_lrRFE$Specificity,
       y=df_lrRFE$Sensitivity,
       type="l",
       col="red")
points(x=df_tMAD$Specificity,
       y=df_tMAD$Sensitivity,
       type="l",
       col="black")
AUC_values <- c(
  paste0("RF ", paste(rfRFE$optVariables, collapse=","), ") = ",
         round(ROC_rfRFE$auc, 3)),
  paste0("cancer ~ ",
         paste(lrRFE$optVariables, collapse="+"), " = ",
         round(ROC_lrRFE$auc, 3)),
  paste0("cancer ~ tMAD = ", round(tMAD_AUC, 3)))
legend(0.08, 0.3, title=" Area Under Curve (AUC) ", title.adj=0.1,
       legend = AUC_values,
       col=c("blue", "red", "black"),
       text.col=c("blue", "red", "black"),
       title.col="black",
       cex=0.8, bty="n")
dev.off()
##### Get the resampling results for all the models in the training
data
rfeResamples <- resamples(list("Random Forest" = rfRFE,
                                 "LR (tMAD + fragFeatures)" = lrRFE))
saveRDS(rfeResamples, "rfeResamples")
pdf("Supplementary_Model_Comparison_on_trainingData_crossValidation.
pdf")
print(bwplot(rfeResamples, metric=c("ROC", "Accuracy"),
            xlim=c(0.1, 1.1)))
dev.off()
summary(rfeResamples)
```
Predict low-ctDNA cancers with test control cohort (n = 26)
`{r}
#####
Plotting for training & test
library(ggplot2)
library(dplyr)
library(caret)

```

```

library(pROC)
library(ggplot2)
library(randomForest)
MY_SEED <- 666
groupB <- read.csv(file=".~/2018_Group_B_cancers_noSZ.csv",
 header=T,
 stringsAsFactors = F)
Convert multiple cancer classes into cancer/healthy
groupB$cancer <- as.character(groupB$cancer)
groupB$cancer[groupB$cancer != "healthy"] <- "cancer"
groupB$cancer <- factor(groupB$cancer,
 levels=c("healthy", "cancer"))
names(groupB)[names(groupB) == "cancer"] <- "Class"
testing <- training <- lrRFE <- rfRFE <- NULL
testing <- readRDS("testing")
training <- readRDS("training")
lrRFE <- readRDS("lrRFE")
rfRFE <- readRDS("rfRFE")
predVars <- c("tMAD", "amplitude_10bp",
 "P160_180",
 "P180_220",
 "P250_320")
lowctDNA cancer data combined with healthy samples from test
cohort
testing <- rbind(testing[testing$Class == "healthy",],
 groupB[groupB$Class == "cancer", c("sample",
 "Class", predVars)])
testing$Class <- factor(testing$Class, levels = c("healthy",
"cancer"))
Get ROC curves for the different models
#1) Only t-MAD
training_binary <- training
testing_binary <- testing
training_binary$Class <- as.character(training_binary$Class)
testing_binary$Class <- as.character(testing_binary$Class)
training_binary$Class[training_binary$Class == "healthy"] <- 0
training_binary$Class[training_binary$Class != "0"] <- 1
training_binary$Class <- factor(as.numeric(training_binary$Class))
testing_binary$Class[testing_binary$Class == "healthy"] <- 0
testing_binary$Class[testing_binary$Class != "0"] <- 1
testing_binary$Class <- factor(as.numeric(testing_binary$Class))
lr_tMAD <- glm(Class ~ tMAD,
 data = training_binary,
 family = binomial)
saveRDS(lr_tMAD , file="lr_tMAD_groupB_26Controls")
prob <- predict(lr_tMAD, newdata=testing_binary, type="response")
pred <- ROCR::prediction(prob, testing_binary$Class)
perf <- ROCR::performance(pred, measure = "tpr", x.measure = "fpr")
tMAD_AUC <- ROCR::performance(pred, measure = "auc")@y.values[[1]]
df_tMAD <- data.frame(Sensitivity=perf@x.values[[1]],
 Specificity=perf@y.values[[1]])
Logistic regression, recursive feature elimination
ROC_lrRFE <- roc(testing$Class,
 predict(lrRFE, testing[,predVars])$cancer)
ROC_lrRFE
df_lrRFE <- data.frame(Sensitivity=ROC_lrRFE$sensitivities,
 Specificity=1-ROC_lrRFE$specificities)
Random Forest RFE
library(randomForest)
ROC_rfRFE <- roc(testing$Class,
 predict(rfRFE, testing[,predVars])$cancer,
 levels=c("healthy", "cancer"))
ROC_rfRFE
df_rfRFE <- data.frame(Sensitivity=ROC_rfRFE$sensitivities,
 Specificity=1-ROC_rfRFE$specificities)
Plotting ROC curves
pdf("Model_Comparison_on_GroupB_26Controls.pdf")
plot(x=df_rfRFE$Specificity,
 y=df_rfRFE$Sensitivity,
 xlab="1 - Specificity",
 ylab="Sensitivity", type="l",
 col="red4")
points(x=df_lrRFE$Specificity,
 y=df_lrRFE$Sensitivity,
 type="l",
 col="orange3")
points(x=df_tMAD$Specificity,
 y=df_tMAD$Sensitivity,
 type="l",
 col="blue4")

```

```

col="black")
AUC_values <- c(
 paste0("RF ", paste(rfRFE$optVariables, collapse=","), " = ",
 round(ROC_rfRFE$auc, 3)),
 paste0("cancer ~ ",
 paste(lrRFE$optVariables, collapse="+"), " = ",
 round(ROC_lrRFE$auc, 3)),
 paste0("cancer ~ tMAD = ", round(tMAD_AUC, 3)))
legend(0.08, 0.3, title=" Area Under Curve (AUC) ", title.adj=0.1,
 legend = AUC_values,
 col=c("red4", "orange3", "black"),
 text.col=c("red4", "orange3", "black"),
 title.col="black",
 cex=0.8, bty="n")
dev.off()
#####

```

All references cited herein are incorporated herein by reference in their entirety and for all purposes to the same extent as if each individual publication or patent or patent application was specifically and individually indicated to be incorporated by reference in its entirety.

The specific embodiments described herein are offered by way of example, not by way of limitation. Any sub-titles herein are included for convenience only, and are not to be construed as limiting the disclosure in any way.

#### REFERENCES

1. G. Siravegna, S. Marsoni, S. Siena, A. Bardelli, Integrating liquid biopsies into the management of cancer, *Nat. Rev. Clin. Oncol.* (2017), doi:10.1038/nrclinonc.2017.14.
2. J. C. M. Wan, C. Massie, J. Garcia-Corbacho, F. Mouliere, J. D. Brenton, C. Caldas, S. Pacey, R. Baird, N. Rosenfeld, Liquid biopsies come of age: towards implementation of circulating tumour DNA, *Nat. Rev. Cancer* 17, 223-238 (2017).
3. M. Murtaza, S.-J. Dawson, D. W. Y. Tsui, D. Gale, T. Forshaw, A. M. Piskorz, C. Parkinson, S.-F. Chin, Z. Kingsbury, A. S. C. Wong, F. Marass, S. Humphray, J. Hadfield, D. Bentley, T. M. Chin, J. D. Brenton, C. Caldas, N. Rosenfeld, Non-invasive analysis of acquired resistance to cancer therapy by sequencing of plasma DNA, *Nature* 497, 108-112 (2013).
4. V. A. Adalsteinsson, G. Ha, S. S. Freeman, A. D. Choudhury, D. G. Stover, H. A. Parsons, G. Gydush, S. C. Reed, D. Rotem, J. Rhoades, D. Loginov, D. Livitz, D. Rosebrock, I. Leshchiner, J. Kim, C. Stewart, M. Rosenberg, J. M. Francis, C.-Z. Zhang, O. Cohen, C. Oh, H. Ding, P. Polak, M. Lloyd, S. Mahmud, K. Helvie, M. S. Merrill, R. A. Santiago, E. P. O'Connor, S. H. Jeong, R. Leeson, R. M. Barry, J. F. Kramkowski, Z. Zhang, L. Polacek, J. G. Lohr, M. Schleicher, E. Lipscomb, A. Saltzman, N. M. Oliver, L. Marini, A. G. Waks, L. C. Harshman, S. M. Tolaney, E. M. Van Allen, E. P. Winer, N. U. Lin, M. Nakabayashi, M.-E. Taplin, C. M. Johannessen, L. A. Garraway, T. R. Golub, J. S. Boehm, N. Wagle, G. Getz, J. C. Love, M. Meyerson, Scalable whole-exome sequencing of cell-free DNA reveals high concordance with metastatic tumors, *Nat. Commun.* 8, 1324 (2017).
5. E. Heitzer, P. Ulz, J. Belic, S. Gutschi, F. Quehenberger, K. Fischereder, T. Benezeder, M. Auer, C. Pischler, S. Mannweiler, M. Pichler, F. Eisner, M. Haeusler, S. Riethdorf, K. Pantel, H. Samonigg, G. Hoefler, H. Augustin, J. B. Geigl, M. R. Speicher, Tumor-associated copy number changes in the circulation of patients with prostate cancer identified through whole-genome sequencing, *Genome Med.* 5, 30 (2013).
6. C. Bettegowda, M. Sausen, R. J. Leary, I. Kinde, Y. Wang, N. Agrawal, B. R. Bartlett, H. Wang, B. Luber, R. M. Alani, E. S. Antonarakis, N. S. Azad, A. Bardelli, H. Brem, J. L. Cameron, C. C. Lee, L. A. Fecher, G. L. Gallia, P. Gibbs, D. Le, R. L. Giuntoli, M. Goggins, M. D. Hogarty, M. Holdhoff, S.-M. Hong, Y. Jiao, H. H. Juhl, J. J. Kim, G. Siravegna, D. A. Laheru, C. Lauricella, M. Lim, E. J. Lipson, S. K. N. Marie, G. J. Netto, K. S. Oliner, A. Olivi, L. Olsson, G. J. Riggins, A. Sartore-Bianchi, K. Schmidt, I.-M. Shih, S. M. Oba-Shinjo, S. Siena, D. Theodorescu, J. Tie, T. T. Harkins, S. Veronese, T.-L. Wang, J. D. Weingart, C. L. Wolfgang, L. D. Wood, D. Xing, R. H. Hruban, J. Wu, P. J. Allen, C. M. Schmidt, M. A. Choti, V. E. Velculescu, K. W. Kinzler, B. Vogelstein, N. Papadopoulos, L. A. Diaz, Detection of Circulating Tumor DNA in Early- and Late-Stage Human Malignancies, *Sci. Transl. Med.* 6, 224ra24-224ra24 (2014).
7. F. Diehl, M. Li, D. Dressman, Y. He, D. Shen, S. Szabo, L. A. Diaz, S. N. Goodman, K. A. David, H. Juhl, K. W. Kinzler, B. Vogelstein, Detection and quantification of mutations in the plasma of patients with colorectal tumors, *Proc. Natl. Acad. Sci.* 102, 16368-16373 (2005).
8. S.-J. Dawson, D. W. Y. Tsui, M. Murtaza, H. Biggs, O. M. Rueda, S.-F. Chin, M. J. Dunning, D. Gale, T. Forshaw, B. Mahler-Araujo, S. Rajan, S. Humphray, J. Becq, D. Hall-sall, M. Wallis, D. Bentley, C. Caldas, N. Rosenfeld, Analysis of Circulating Tumor DNA to Monitor Metastatic Breast Cancer, *N. Engl. J. Med.* 368, 1199-1209 (2013).
9. F. Diehl, K. Schmidt, M. A. Choti, K. Romans, S. Goodman, M. Li, K. Thornton, N. Agrawal, L. Sokoll, S. A. Szabo, K. W. Kinzler, B. Vogelstein, L. A. Diaz, Circulating mutant DNA to assess tumor dynamics., *Nat. Med.* 14, 985-90 (2008).
10. J. Tie, Y. Wang, C. Tomasetti, L. Li, S. Springer, I. Kinde, N. Silliman, M. Tacey, H.-L. Wong, M. Christie, S. Kosmider, I. Skinner, R. Wong, M. Steel, B. Tran, J. Desai, I. Jones, A. Haydon, T. Hayes, T. J. Price, R. L. Strausberg, L. A. Diaz, N. Papadopoulos, K. W. Kinzler, B. Vogelstein, P. Gibbs, Circulating tumor DNA analysis detects minimal residual disease and predicts recurrence in patients with stage II colon cancer., *Sci. Transl. Med.* 8, 346ra92 (2016).
11. A. A. Chaudhuri, J. J. Chabon, A. F. Lovejoy, A. M. Newman, H. Stehr, T. D. Azad, M. S. Khodadoust, M. S.

- Esfahani, C. L. Liu, L. Zhou, F. Scherer, D. M. Kurtz, C. Say, J. N. Carter, D. J. Merriott, J. C. Dudley, M. S. Binkley, L. Modlin, S. K. Padda, M. F. Gensheimer, R. B. West, J. B. Shrager, J. W. Neal, H. A. Wakelee, B. W. Loo, A. A. Alizadeh, M. Diehn, Early Detection of Molecular Residual Disease in Localized Lung Cancer by Circulating Tumor DNA Profiling., *Cancer Discov.* 7, 1394-1403 (2017).
12. J. D. Cohen, L. Li, Y. Wang, C. Thoburn, B. Afsari, L. Danilova, C. Douville, A. A. Javed, F. Wong, A. Mattox, R. H. Hruban, C. L. Wolfgang, M. G. Goggins, M. Dal Molin, T.-L. Wang, R. Roden, A. P. Klein, J. Ptak, L. Dobbyn, J. Schaefer, N. Silliman, M. Popoli, J. T. Vogelstein, J. D. Browne, R. E. Schoen, R. E. Brand, J. Tie, P. Gibbs, H.-L. Wong, A. S. Mansfield, J. Jen, S. M. Hanash, M. Falconi, P. J. Allen, S. Zhou, C. Bettegowda, L. A. Diaz, C. Tomasetti, K. W. Kinzler, B. Vogelstein, A. M. Lennon, N. Papadopoulos, Detection and localization of surgically resectable cancers with a multi-analyte blood test., *Science* 359, 926-930 (2018).
13. I. S. Haque, O. Elemento, Challenges in Using ctDNA to Achieve Early Detection of Cancer, *bioRxiv*, 237578 (2017).
14. A. M. Newman, A. F. Lovejoy, D. M. Klass, D. M. Kurtz, J. J. Chabon, F. Scherer, H. Stehr, C. L. Liu, S. V Bratman, C. Say, L. Zhou, J. N. Carter, R. B. West, G. W. Sledge Jr, J. B. Shrager, B. W. Loo, J. W. Neal, H. A. Wakelee, M. Diehn, A. A. Alizadeh, Integrated digital error suppression for improved detection of circulating tumor DNA, *Nat. Biotechnol.* 34, 547-555 (2016).
15. P. Ulz, G. G. Thallinger, M. Auer, R. Graf, K. Kashofer, S. W. Jahn, L. Abete, G. Pristauz, E. Petru, J. B. Geigl, E. Heitzer, M. R. Speicher, Inferring expressed genes by whole-genome sequencing of plasma DNA, *Nat. Genet.* 48, 1273-1278 (2016).
16. M. W. Snyder, M. Kircher, A. J. Hill, R. M. Daza, J. Shendure, Cell-free DNA Comprises an In Vivo Nucleosome Footprint that Informs Its Tissues-Of-Origin., *Cell* 164, 57-68 (2016).
17. P. Burnham, M. S. Kim, S. Agbor-Enoh, H. Luikart, H. A. Valantine, K. K. Khush, I. De Vlaminck, Single-stranded DNA library preparation uncovers the origin and diversity of ultrashort cell-free DNA in plasma, *Sci. Rep.* 6, 27859 (2016).
18. G. Genovese, A. K. Kahler, R. E. Handsaker, J. Lindberg, S. A. Rose, S. F. Bakhour, K. Chambert, E. Mick, B. M. Neale, M. Fromer, S. M. Purcell, O. Svantesson, M. Landén, M. Höglund, S. Lehmann, S. B. Gabriel, J. L. Moran, E. S. Lander, P. F. Sullivan, P. Sklar, H. Grönberg, C. M. Hultman, S. A. McCarroll, Clonal Hematopoiesis and Blood-Cancer Risk Inferred from Blood DNA Sequence, *N. Engl. J. Med.* 371, 2477-2487 (2014).
19. Y. Hu, B. Ulrich, J. Supplee, Y. Kuang, P. H. Lizotte, N. Feeney, N. Guibert, M. M. Awad, K.-K. Wong, P. A. Janne, C. P. Paweletz, G. R. Oxnard, False positive plasma genotyping due to clonal hematopoiesis., *Clin. Cancer Res.*, clincanres.0143.2018 (2018).
20. A. J. Bronkhorst, J. F. Wentzel, J. Aucamp, E. van Dyk, L. du Plessis, P. J. Pretorius, Characterization of the cell-free DNA released by cultured cancer cells, *Biochim. Biophys. Acta—Mol. Cell Res.* 1863, 157-165 (2016).
21. S. Jahr, H. Hentze, S. Englisch, D. Hardt, F. O. Fackelmayer, R. D. Hesch, R. Knippers, DNA fragments in the blood plasma of cancer patients: quantitations and evidence for their origin from apoptotic and necrotic cells., *Cancer Res.* 61, 1659-65 (2001).

22. Y. M. D. Lo, K. C. A. Chan, H. Sun, E. Z. Chen, P. Jiang, F. M. F. Lun, Y. W. Zheng, T. Y. Leung, T. K. Lau, C. R. Cantor, R. W. K. Chiu, Maternal plasma DNA sequencing reveals the genome-wide genetic and mutational profile of the fetus., *Sci. Transl. Med.* 2, 61ra91 (2010).
23. D. Chandrananda, N. P. Thorne, M. Bahlo, L.-S. Tam, G. Liao, E. Li, High-resolution characterization of sequence signatures due to non-random cleavage of cell-free DNA, *BMC Med. Genomics* 8, 29 (2015).
10. 24. P. Jiang, Y. M. D. Lo, The Long and Short of Circulating Cell-Free DNA and the Ins and Outs of Molecular Diagnostics, *Trends Genet.* 32, 360-371 (2016).
15. 25. S. C. Y. Yu, K. C. A. Chan, Y. W. L. Zheng, P. Jiang, G. J. W. Liao, H. Sun, R. Akolekar, T. Y. Leung, A. T. J. I. Go, J. M. G. van Vugt, R. Minekawa, C. B. M. Oudejans, K. H. Nicolaides, R. W. K. Chiu, Y. M. D. Lo, Size-based molecular diagnostics using plasma DNA for noninvasive prenatal testing., *Proc. Natl. Acad. Sci. U.S.A* 111, 8583-8 (2014).
20. 26. F. M. F. Lun, N. B. Y. Tsui, K. C. A. Chan, T. Y. Leung, T. K. Lau, P. Charoenkwan, K. C. K. Chow, W. Y. W. Lo, C. Wanapirak, T. Sanguansermsri, C. R. Cantor, R. W. K. Chiu, Y. M. D. Lo, Noninvasive prenatal diagnosis of monogenic diseases by digital size selection and relative mutation dosage on DNA in maternal plasma., *Proc. Natl. Acad. Sci. U.S.A* 105, 19920-5 (2008).
27. G. Minarik, G. Repiska, M. Hyblova, E. Nagyova, K. Solty, J. Budis, F. Duris, R. Sysak, M. Gerykova Bujalkova, B. Vlkova-Izrael, O. Biro, B. Nagy, T. Szemes, Utilization of Benchtop Next Generation Sequencing Platforms Ion Torrent PGM and MiSeq in Noninvasive Prenatal Testing for Chromosome 21 Trisomy and Testing of Impact of In Silico and Physical Size Selection on Its Analytical Performance., *PLoS One* 10, e0144811 (2015).
30. 28. M. B. Giacoma, G. C. Ruben, K. A. Iczkowski, T. B. Roos, D. M. Porter, G. D. Sorenson, Cell-Free DNA in Human Blood Plasma, *Pancreas* 17, 89-97 (1998).
35. 29. N. Umetani, A. E. Giuliano, S. H. Hiramatsu, F. Amersi, T. Nakagawa, S. Martino, D. S. B. Hoon, Prediction of breast tumor progression by integrity of free circulating DNA in serum., *J. Clin. Oncol.* 24, 4270-6 (2006).
40. 30. F. Mouliere, B. Robert, E. Arnau Peyrotte, M. Del Rio, M. Ychou, F. Molina, C. Gongora, A. R. Thierry, T. Lee, Ed. High Fragmentation Characterizes Tumour-Derived Circulating DNA, *PLoS One* 6, e23418 (2011).
45. 31. F. Mouliere, S. El Messaoudi, D. Pang, A. Dritschilo, A. R. Thierry, Multi-marker analysis of circulating cell-free DNA toward personalized medicine for colorectal cancer, *Mol. Oncol.* 8, 927-941 (2014).
50. 32. P. Jiang, C. W. M. Chan, K. C. A. Chan, S. H. Cheng, J. Wong, V. W.-S. Wong, G. L. H. Wong, S. L. Chan, T. S. K. Mok, H. L. Y. Chan, P. B. S. Lai, R. W. K. Chiu, Y. M. D. Lo, Lengthening and shortening of plasma DNA in hepatocellular carcinoma patients., *Proc. Natl. Acad. Sci. U.S.A* 112, E1317-25 (2015).
55. 33. H. R. Underhill, J. O. Kitzman, S. Hellwig, N. C. Welker, R. Daza, D. N. Baker, K. M. Gligorich, R. C. Rostomily, M. P. Bronner, J. Shendure, D. J. Kwiatkowski, Ed. Fragment Length of Circulating Tumor DNA, *PLOS Genet.* 12, e1006162 (2016).
60. 34. O. A. Zill, K. C. Banks, S. R. Fairclough, S. A. Mortimer, J. V Vowles, R. Mokhtari, D. R. Gandara, P. C. Mack, J. I. Odegaard, R. J. Nagy, A. M. Baca, H. Eltoukhy, D. I. Chudova, R. B. Lanman, A. Talasaz, The Landscape of Actionable Genomic Alterations in Cell-Free Circulating Tumor DNA from 21,807 Advanced Cancer Patients., *Clin. Cancer Res.*, clincanres.3837.2017 (2018).

35. G. Macintyre, T. E. Goranova, D. De Silva, D. Ennis, A. M. Piskorz, M. Eldridge, D. Sie, L.-A. Lewisley, A. Hanif, C. Wilson, S. Dowson, R. M. Glasspool, M. Lockley, E. Brockbank, A. Montes, A. Walther, S. Sundar, R. Edmondson, G. D. Hall, A. Clamp, C. Gourley, M. Hall, C. Fotopoulou, H. Gabra, J. Paul, A. Supernat, D. Millan, A. Hoyle, G. Bryson, C. Nourse, L. Mincarelli, L. N. Sanchez, B. Ylstra, M. Jimenez-Linan, L. Moore, O. Hofmann, F. Markowitz, I. A. McNeish, J. D. Brenton, Copy number signatures and mutational processes in ovarian carcinoma, *Nat. Genet.*, 1 (2018).
36. C. A. Parkinson, D. Gale, A. M. Piskorz, H. Biggs, C. Hodgkin, H. Addley, S. Freeman, P. Moyle, E. Sala, K. Sayal, K. Hosking, I. Gounaris, M. Jimenez-Linan, H. M. Earl, W. Qian, N. Rosenfeld, J. D. Brenton, E. R. Mardis, Ed. Exploratory Analysis of TP53 Mutations in Circulating Tumour DNA as Biomarkers of Treatment Response for Patients with Relapsed High-Grade Serous Ovarian Carcinoma: A Retrospective Study, *PLOS Med.* 13, e1002198 (2016).
37. T. Forshaw, M. Murtaza, C. Parkinson, D. Gale, D. W. Y. Tsui, F. Kaper, S.-J. Dawson, A. M. Piskorz, M. Jimenez-Linan, D. Bentley, J. Hadfield, A. P. May, C. Caldas, J. D. Brenton, N. Rosenfeld, Noninvasive identification and monitoring of cancer mutations by targeted deep sequencing of plasma DNA., *Sci. Transl. Med.* 4, 136ra68 (2012).
38. A. R. Thierry, S. El Messaoudi, P. B. Gahan, P. Anker, M. Stroun, Origins, structures, and functions of circulating DNA in oncology, *Cancer Metastasis Rev.* 35, 347-376 (2016).
39. M. G. Best, N. Sol, B. A. Tannous, P. Wesseling, T. Wurdinger, RNA-Seq of Tumor-Educated Platelets Enables Blood-Based Pan-Cancer, Multiclass, and Molecular Pathway Cancer Diagnostics, *Cancer Cell* 28, 666-676 (2015).
40. M. G. Best, N. Sol, S. G. J. G. In't Veld, A. Vancura, M. Muller, A.-L. N. Niemeijer, A. V Fejes, L.-A. Tjon Kon Fat, A. E. Huis In't Veld, C. Leurs, T. Y. Le Large, L. L. Meijer, I. E. Kooi, F. Rustenburg, P. Schellen, H. Verschueren, E. Post, L. E. Wedekind, J. Bracht, M. Esenkbrink, L. Wils, F. Favaro, J. D. Schoonhoven, J. Tannous, H. Meijers-Heijboer, G. Kazemier, E. Giovannetti, J. C. Reijneveld, S. Idema, J. Killestein, M. Heger, S. C. de Jager, R. T. Urbanus, I. E. Hoefer, G. Pasterkamp, C. Mannhalter, J. Gomez-Arroyo, H.-J. Bogaard, D. P. Noske, W. P. Vandertop, D. van den Broek, B. Ylstra, R. J. A. Nilsson, P. Wesseling, N. Karachaliou, R. Rosell, E. Lee-Lewandrowski, K. B. Lewandrowski, B. A. Tannous, A. J. de Langen, E. F. Smit, M. M. van den Heuvel, T. Wurdinger, Swarm Intelligence-Enhanced Detection of Non-Small-Cell Lung Cancer Using Tumor-Educated Platelets., *Cancer Cell* 32, 238-252.e9 (2017).
41. A. L. Riediger, S. Dietz, U. Schirmer, M. Meister, I. Heinzmamn-Groth, M. Schneider, T. Muley, M. Thomas, H. Siltmann, Mutation analysis of circulating plasma DNA to determine response to EGFR tyrosine kinase inhibitor therapy of lung adenocarcinoma patients, *Sci. Rep.* 6, 33505 (2016).
42. J. Belic, M. Koch, P. Ulz, M. Auer, T. Gerhalter, S. Mohan, K. Fischereder, E. Petru, T. Bauernhofer, J. B. Geigl, M. R. Speicher, E. Heitzer, Rapid Identification of Plasma DNA Samples with Increased ctDNA Levels by a Modified FAST-SeqS Approach, *Clin. Chem.* 61, 838-849 (2015).
43. D. G. Stover, H. A. Parsons, G. Ha, S. S. Freeman, W. T. Barry, H. Guo, A. D. Choudhury, G. Gyudush, S. C.

- Reed, J. Rhoades, D. Rotem, M. E. Hughes, D. A. Dillon, A. H. Partridge, N. Wagle, I. E. Krop, G. Getz, T. R. Golub, J. C. Love, E. P. Winer, S. M. Tolaney, N. U. Lin, V. A. Adalsteinsson, Association of Cell-Free DNA Tumor Fraction and Somatic Copy Number Alterations With Survival in Metastatic Triple-Negative Breast Cancer, *J. Clin. Oncol.* 36, 543-553 (2018).
44. B. Routy, E. Le Chatelier, L. Derosa, C. P. M. Duong, M. T. Alou, R. Daillere, A. Fluckiger, M. Messaoudene, C. Rauber, M. P. Roberti, M. Fidelle, C. Flament, V. Poirier-Colame, P. Opolon, C. Klein, K. Iribarren, L. Mondragon, N. Jacquemet, B. Qu, G. Ferrere, C. Clémenson, L. Mezquita, J. R. Masip, C. Naltet, S. Brosseau, C. Kadherhai, C. Richard, H. Rizvi, F. Levenez, N. Galleron, B. Quinquis, N. Pons, B. Ryffel, V. Minard-Colin, P. Gonin, J.-C. Soria, E. Deutsch, Y. Loriot, F. Ghiringhelli, G. Zalcman, F. Goldwasser, B. Escudier, M. D. Hellmann, A. Eggermont, D. Raoult, L. Albiges, G. Kroemer, L. Zitvogel, Gut microbiome influences efficacy of PD-1-based immunotherapy against epithelial tumors., *Science* 359, 91-97 (2018).
45. C. Abbosh, N. J. Birkbak, G. A. Wilson, M. Jamal-Hanjani, T. Constantin, R. Salari, J. Le Quesne, D. A. Moore, S. Veeriah, R. Rosenthal, T. Marafioti, E. Kirkizlar, T. B. K. Watkins, N. McGranahan, S. Ward, L. Martinson, J. Riley, F. Fraioli, M. Al Bakir, E. Grönroos, F. Zambrana, R. Endozo, W. L. Bi, F. M. Fennelly, N. Spomer, D. Johnson, J. Laycock, S. Shafi, J. Czyzewska-Khan, A. Rowan, T. Chambers, N. Matthews, S. Turajlic, C. Hiley, S. M. Lee, M. D. Forster, T. Ahmad, M. Falzon, E. Borg, D. Lawrence, M. Hayward, S. Kolvekar, N. Panagiotopoulos, S. M. Janes, R. Thakrar, A. Ahmed, F. Blackhall, Y. Summers, D. Hafez, A. Naik, A. Ganguly, S. Kareht, R. Shah, L. Joseph, A. Marie Quinn, P. A. Crosbie, B. Naidu, G. Middleton, G. Langman, S. Trotter, M. Nicolson, H. Remmen, K. Kerr, M. Chetty, L. Gomersall, D. A. Fennell, A. Nakas, S. Rathinam, G. Anand, S. Khan, P. Russell, V. Ezhil, B. Ismail, M. Irvin-Sellers, V. Prakash, J. F. Lester, M. Kornaszewska, R. Attanoos, H. Adams, H. Davies, D. Oukrif, A. U. Akarca, J. A. Hartley, H. L. Lowe, S. Lock, N. Iles, H. Bell, Y. Ngai, G. Elgar, Z. Szallas, R. F. Schwarz, J. Herrero, A. Stewart, S. A. Quezada, K. S. Peggs, P. Van Loo, C. Dive, C. J. Lin, M. Rabinowitz, H. J. W. L. Aerts, A. Hackshaw, J. A. Shaw, B. G. Zimmermann, TRACERx consortium, PEACE consortium, C. Swanton, Phylogenetic ctDNA analysis depicts early-stage lung cancer evolution., *Nature* 545, 446-451 (2017).
46. D. Hanahan, R. A. Weinberg, Hallmarks of cancer: the next generation., *Cell* 144, 646-74 (2011).
47. K. M. Patel, K. E. van der Vos, C. G. Smith, F. Mouliere, D. Tsui, J. Morris, D. Chandrananda, F. Marass, D. van den Broek, D. E. Neal, V. J. Gnanapragasam, T. Forshaw, B. W. van Rhijn, C. E. Massie, N. Rosenfeld, M. S. van der Heijden, Association Of Plasma And Urinary Mutant DNA With Clinical Outcomes In Muscle Invasive Bladder Cancer, *Sci. Rep.* 7, 5554 (2017).
48. H. Li, R. Durbin, Fast and accurate short read alignment with Burrows-Wheeler transform, *Bioinformatics* 25, 1754-1760 (2009).
49. I. Scheinin, D. Sie, H. Bengtsson, M. A. van de Wiel, A. B. Olshen, H. F. van Thuijl, H. F. van Essen, P. P. Eijk, F. Rustenburg, G. A. Meijer, J. C. Reijneveld, P. Wesseling, D. Pinkel, D. G. Albertson, B. Ylstra, DNA copy number analysis of fresh and formalin-fixed specimens by shallow whole-genome sequencing with identification and exclu-

sion of problematic regions in the genome assembly, *Genome Res.* 24, 2022-2032 (2014).

The invention claimed is:

1. A computer-implemented method for detecting variant nucleic acid from a cell-free nucleic acid-containing sample from a subject, wherein the variant nucleic acid is circulating tumor DNA (ctDNA), the method comprising:

- a) providing data representing fragment sizes of nucleic acid fragments obtained from said sample;
- b) causing a processor to process the data from step a) according to a classification algorithm that has been trained on a training set comprising a plurality of samples of cell-free nucleic acid containing the variant nucleic acid and a plurality of samples not containing the variant nucleic acid, wherein said classification algorithm operates to classify sample data into one of at least two classes, the at least two classes comprising a first class containing the variant nucleic acid and a second class not containing the variant nucleic acid, wherein said classification algorithm operates to classify sample data into one of said at least two classes based on at least a plurality of cell-free DNA (cfDNA) fragment size features selected from the group consisting of:
  - (i) a proportion of fragments in a 20-150 bp size range (P(20-150));
  - (ii) a proportion of fragments in a 100-150 bp size range (P(100-150));
  - (iii) a proportion of fragments in a 160-180 bp size range (P(160-180));
  - (iv) a proportion of fragments in a 180-220 bp size range (P(180-220));
  - (v) a proportion of fragments in a 250-320 bp size range (P(250-320));
  - (vi) a ratio of the proportions P(20-150)/P(160-180);
  - (vii) a ratio of the proportion P(100-150) divided by the proportion of fragments in a 163-169 bp size range;
  - (viii) a ratio of the proportions P(20-150)/P(180-220); and
  - (ix) amplitude oscillations in fragment size density with 10 bp periodicity,

wherein the data representing fragment sizes of nucleic acid fragments in step a) includes the plurality of cfDNA fragment size features used by the classification algorithm; and

c) outputting the classification of the sample from step b) and thereby determining whether the sample contains ctDNA,

wherein the classification of the sample as containing ctDNA or not is used to predict whether said sample or a further sample from the subject will be susceptible to further ctDNA analysis,

wherein said sample is classified as containing ctDNA, said further ctDNA analysis comprises sequencing to a greater sequencing depth and/or targeted sequencing of ctDNA in said sample or said further sample, and the sample or further sample is subjected to said further ctDNA analysis.

2. The method of claim 1, wherein the data representing fragment sizes of the nucleic acid fragments comprise fragment sizes inferred from sequence reads, fragment sizes determined by fluorimetry, or fragment sizes determined by densitometry, or wherein the fragment sizes of cfDNA fragments are inferred from sequence reads using mapping locations of read ends in a reference genome of a species from which the sample was obtained following alignment of the sequence reads with the reference genome.

3. The method of claim 1, wherein the plurality of cfDNA fragment size features comprise: P(160-180), P(180-220), P(250-320) and the amplitude oscillations in fragment size density with 10 bp periodicity.

4. The method of claim 1, wherein said classification algorithm operates to classify sample data into one of said at least two classes based on said plurality of cell-free DNA (cfDNA) fragment size features and a deviation from copy number neutrality feature which is a trimmed Median Absolute Deviation from copy number neutrality (t-MAD) score or an ichorCNA score.

5. The method of claim 4, wherein the t-MAD score is determined by trimming regions of genome that exhibit high copy number variability in whole genome datasets derived from healthy subjects and then calculating a median absolute deviation from  $\log_2 R=0$  of non-trimmed regions of the genome.

6. The method of claim 1, wherein the classification algorithm performs Random Forests (RF) analysis, logistic regression (LR) analysis, or support vector machine (SVM) analysis, wherein the performance of the classification algorithm when trained on the training set is assessed by an area under the curve (AUC) value from a receiver operating characteristic (ROC) analysis, or wherein the classification algorithm that has been trained on a training set comprising at least 10 samples from healthy subjects and at 10 samples from subjects known to have a cancer.

7. The method according to claim 1, wherein the data provided in step a) represent whole-genome sequence (WGS) reads, Tailored Panel Sequencing (TAPAS) sequence reads, Tagged-Amplicon Deep Sequencing (Tam-Seq) reads, hybrid-capture sequence reads, focused-exome sequence reads or whole-exome sequence reads.

8. The method according to claim 7, wherein the data provided in step a) represent shallow whole-genome sequence (sWGS) reads, optionally 0.4x depth WGS reads.

9. The method according to claim 1, wherein the data provided in step a) represent fragment sizes of multiple DNA fragments from a substantially cell-free liquid sample from a subject having or suspected as having a cancer selected from melanoma, lung cancer, cholangiocarcinoma, bladder cancer, oesophageal cancer, colorectal cancer, ovarian cancer, glioma, pancreatic cancer, renal cancer and breast cancer, or

wherein the sample is a plasma sample, a urine sample, a saliva sample, a cerebrospinal fluid sample, a serum sample, or other DNA-containing biological liquid sample.

10. The method of claim 1, wherein the method is for detecting a presence of, growth of, prognosis of, regression of, treatment response of, or recurrence of a cancer in a subject from which the sample has been obtained.

11. The method of claim 10, wherein the presence of ctDNA in the sample is distinguished from cfDNA containing somatic mutations of non-cancerous origin, wherein the non-cancerous origin comprises clonal expansions of normal epithelia or clonal hematopoiesis of indeterminate potential (CHIP).

12. The method of claim 10, wherein the fragment size data provided in step a) represent sequence reads of multiple DNA fragments from a substantially cell-free liquid sample from a subject and wherein the method is for determining whether the sample contains ctDNA or contains cfDNA from CHIP, wherein the classification algorithm has been trained on a training set further comprising a plurality of samples of cfDNA obtained from subjects having CHIP, and

**85**

wherein said at least two classes further comprise a third class containing CHIP-derived cfDNA.

- 13.** The method of claim 1, further comprising:  
 analysing the cell-free nucleic acid-containing sample, or  
 a library derived from the cell-free nucleic acid-con-  
 taining sample, wherein the sample has been obtained  
 from the subject, to determine fragment sizes of nucleic  
 acid fragments in said sample or said library;  
 wherein said analysing comprises:  
 sequencing nucleic acids from the nucleic acid-con-  
 taining sample or the library to obtain sequence reads  
 and inferring the fragment sizes from the sequence  
 reads;  
 measuring the fragment sizes of nucleic acids from the  
 nucleic acid-containing sample or the library by  
 fluorimetry; or  
 measuring the fragment sizes of nucleic acids from the  
 nucleic acid-containing sample or the library by  
 densitometry.

- 14.** The method of claim 1, further comprising:  
 sequencing the cell-free nucleic acid-containing sample,  
 or a library derived from the cell-free nucleic acid-  
 containing to obtain a plurality of sequence reads; and  
 processing the plurality of sequence reads to determine  
 sequence data representing fragment sizes of cfDNA  
 fragments obtained from said sample and/or represent-  
 ing a measure of deviation from copy number neutrality  
 of the cfDNA fragments obtained from said sample.

- 15.** The method of claim 14, wherein the sequencing  
 comprises generating a sequencing library from the sample  
 and performing whole-genome sequencing, Tailored Panel  
 Sequencing (TAPAS) sequencing, hybrid-capture sequenc-  
 ing, TAm-Seq sequencing, focused-exome sequencing,  
 whole-exome sequencing, or wherein the sequencing com-  
 prises generating an indexed sequencing library and per-  
 forming shallow whole genome sequencing (sWGS),

**86**

optionally sWGS to a depth of 0.4x, or wherein processing  
 the sequence reads comprises one or more of the following  
 steps:

- 5
- aligning sequence reads to a reference genome of a  
 species of the subject;
  - removal of contaminating adapter sequences;
  - removal of PCR and optical duplicates;
  - removal of sequence reads of low mapping quality; and
  - if multiplex sequencing, de-multiplexing by excluding  
 mismatches in sequencing barcodes.

- 16.** The method of claim 1, wherein the sample is a  
 plasma sample, a urine sample, a saliva sample, a cerebro-  
 spinal fluid sample, a serum sample, or other DNA-contain-  
 ing biological liquid sample.

- 17.** The method of claim 14,  
 wherein the method is for detecting a presence of, growth  
 of, prognosis of, regression of, treatment response of,  
 or recurrence of a cancer in a subject from which the  
 sample has been obtained, wherein the presence of  
 ctDNA is distinguished from the presence of cfDNA  
 containing somatic mutations of non-cancerous origin,  
 wherein a somatic mutation containing cfDNA frag-  
 ment is classified as being of tumour origin or being of  
 CHIP origin based on a plurality of fragment size  
 features determined from the plurality of sequence  
 reads.

- 18.** The method of claim 1, wherein:  
 said sample is a plasma sample and wherein a probability  
 that the sample contains ctDNA as determined by the  
 classification algorithm is used to determine whether  
 ctDNA will be detectable in a urine sample; or  
 said sample is a urine sample and wherein a probability  
 that the sample contains ctDNA as determined by the  
 classification algorithm is used to determine whether  
 ctDNA will be detectable in a plasma sample.

\* \* \* \* \*