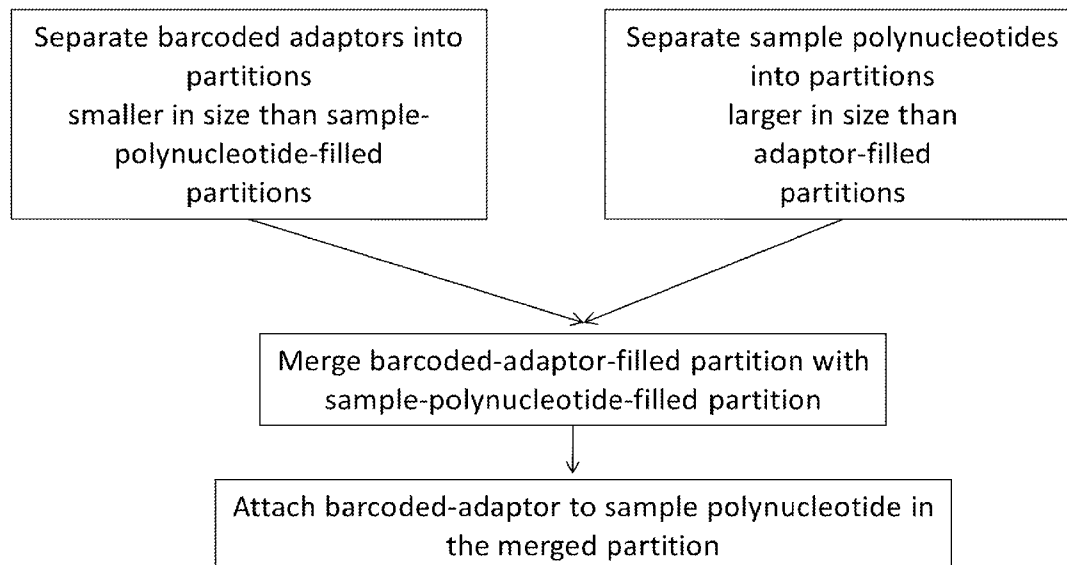




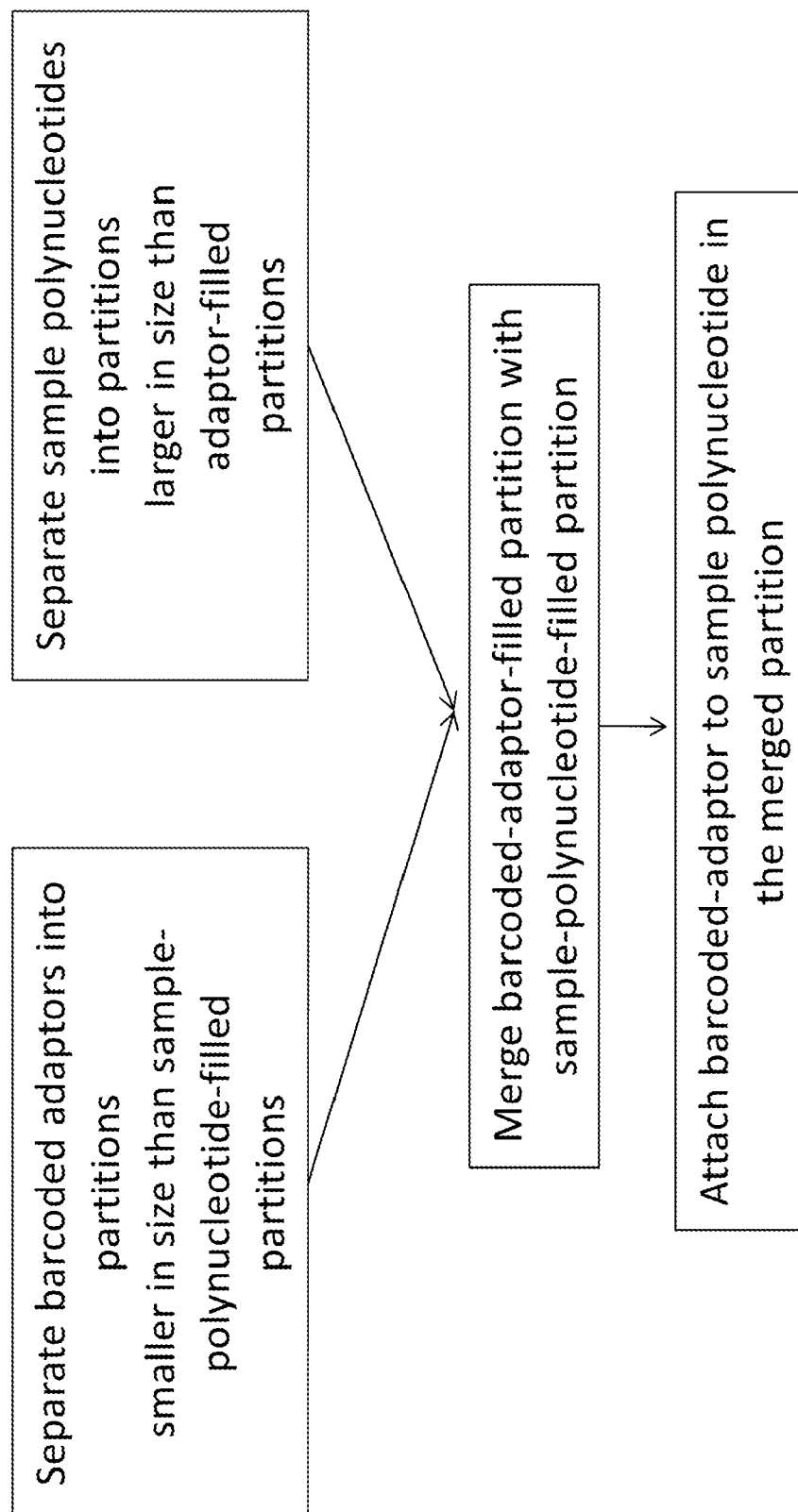
US 20120316074A1

(19) **United States**(12) **Patent Application Publication**  
**Saxonov**(10) **Pub. No.: US 2012/0316074 A1**(43) **Pub. Date: Dec. 13, 2012**(54) **METHODS AND COMPOSITIONS FOR  
NUCLEIC ACID ANALYSIS**(75) Inventor: **Serge Saxonov**, Oakland, CA (US)(73) Assignee: **Bio-Rad Laboratories, Inc.**,  
Hercules, CA (US)(21) Appl. No.: **13/456,121**(22) Filed: **Apr. 25, 2012****Related U.S. Application Data**(60) Provisional application No. 61/478,777, filed on Apr.  
25, 2011.**Publication Classification**(51) **Int. Cl.****C40B 50/06** (2006.01)**C40B 20/00** (2006.01)(52) **U.S. Cl.** ..... **506/2; 506/26**(57) **ABSTRACT**

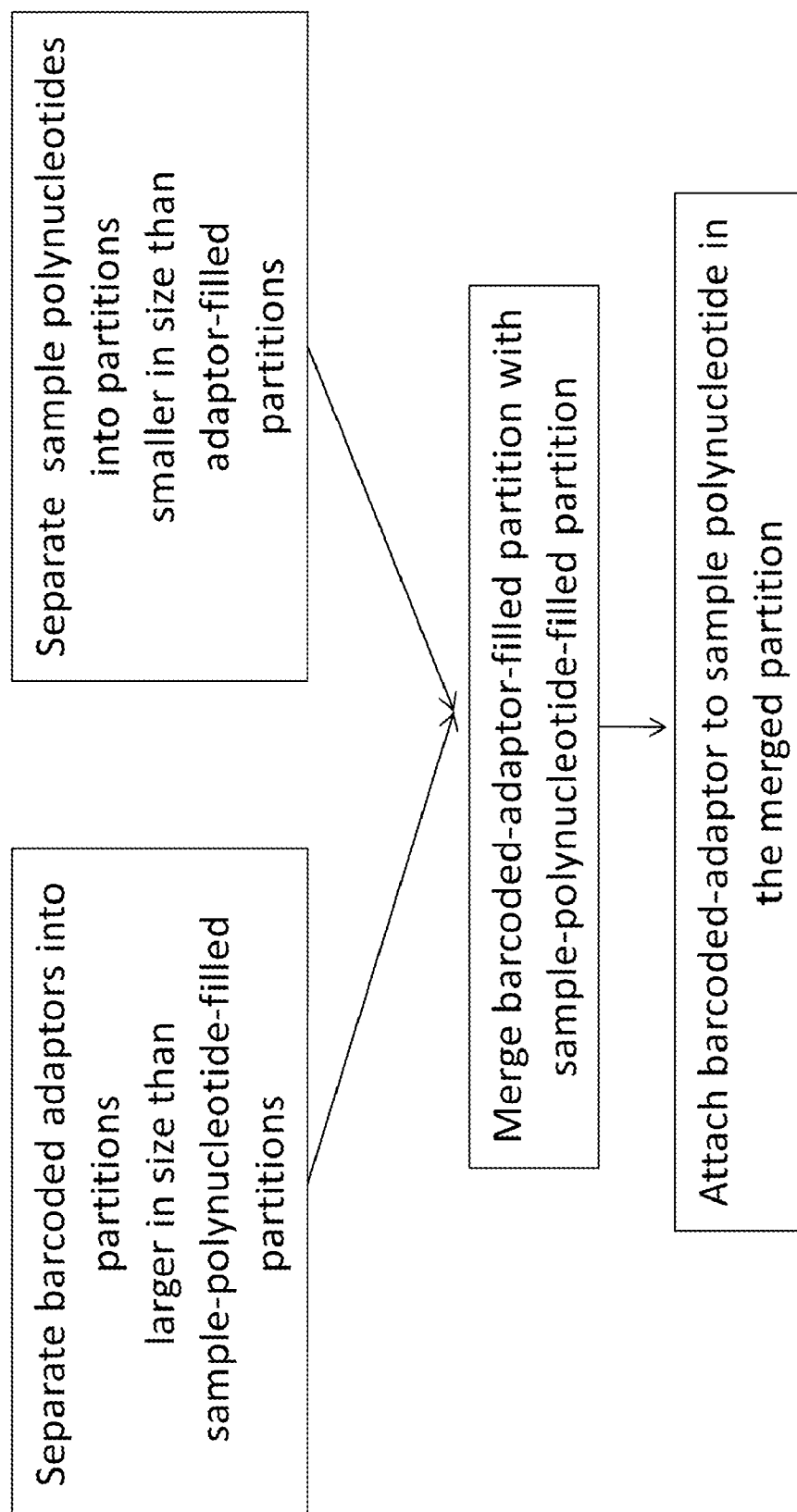
Provided herein are methods, compositions, and kits for assays, many of which involve amplification reactions such as digital PCR or droplet digital PCR. The assays may be used for such applications as sequencing, copy number variation analysis, and others. In some cases, the assays involve subdividing a sample into multiple partitions (e.g., droplets) and merging the partitions with other partitions that comprise adaptors with barcodes.

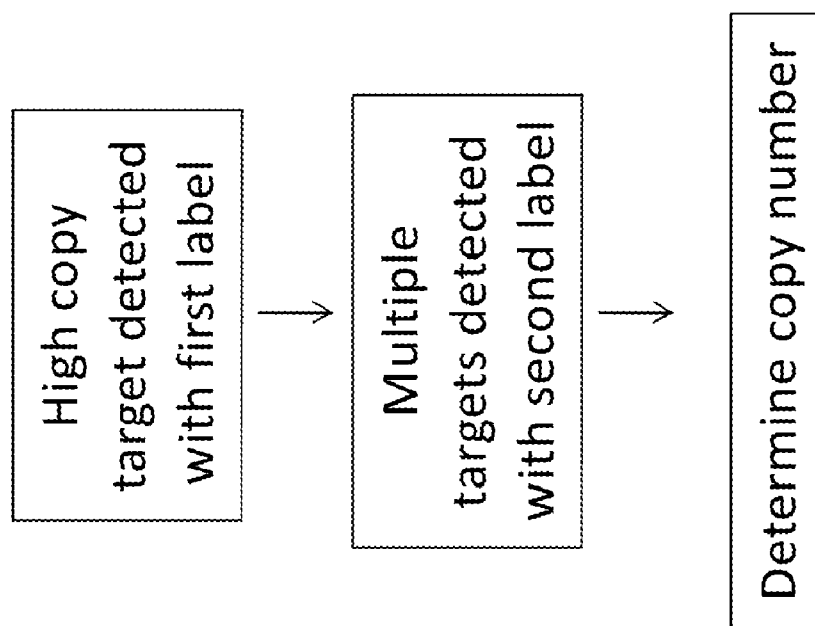


**Figure 1A**



**Figure 1B**



**Figure 2**

## METHODS AND COMPOSITIONS FOR NUCLEIC ACID ANALYSIS

### CROSS REFERENCE

**[0001]** This application claims the benefit of U.S. Provisional Patent Application No. 61/478,777, filed Apr. 25, 2011, which application is incorporated herein by reference in its entirety.

### BACKGROUND

**[0002]** Next generation sequencing has many useful applications and can be used to analyze multiple samples. There is a need for improved methods of multiplexing samples for applications of next generation sequencing. There is also a need for improved methods of barcode tagging partitioned polynucleotides and analyzing the barcode tagged polynucleotides.

**[0003]** Determining the copy number of a target sequence can have many useful applications. There is a need for improved methods of determining the copy number of a target sequence.

### SUMMARY

**[0004]** This disclosure provides methods that can be used in sequencing and other applications. In some instances, this disclosure provides a method comprising: a. subdividing a plurality of adaptors into a plurality of first partitions, wherein each of the first partitions has on average a first volume and wherein the adaptors comprise unique barcodes; b. subdividing a sample comprising multiple polynucleotides into a plurality of second partitions, wherein each of the second partitions has on average a second volume, wherein the second volume is greater than the first volume; c. merging at least one of the first partitions with at least one of the second partitions to form a merged partition; and d. tagging one of the multiple polynucleotides, or fragment thereof, with at least one of the adaptors.

**[0005]** The method may comprise: a. subdividing a plurality of adaptors into a plurality of first partitions, wherein each of the first partitions has on average a first volume and wherein the adaptors comprise unique barcodes; b. subdividing a sample comprising multiple polynucleotides into a plurality of second partitions, wherein each of the second partitions has on average a second volume, wherein said second volume is less than said first volume; c. merging at least one of said first partitions with at least one of said second partitions to form a merged partition; and d. tagging one of said multiple polynucleotides, or fragment thereof, with at least one of said adaptors.

**[0006]** Often, in a method disclosed herein, the first partitions are droplets. In some instances, said second partitions are droplets. In some cases, said droplets are within an immiscible fluid.

**[0007]** In some cases, the polynucleotides are genomic DNA. For example, the genomic DNA may be high molecular weight DNA. In some cases, the sample of genomic DNA is partitioned so that it is unlikely that a given partition comprises two or more polynucleotides, or fragments thereof, from the same locus but from different chromosomes.

**[0008]** In some cases, the first partitions are first droplets and the second partitions are second droplets; and prior to the merging, the at least one second droplet comprises the at least one first droplet. In other cases, the first partitions are first

droplets and the second partitions are second droplets; and prior to the merging, the at least one second droplet does not comprise the at least one first droplet. In some cases, the first partitions are first droplets and the second partitions are second droplets; and prior to the merging, the at least one first droplet comprises the at least one second droplet. In some cases, the first partitions are first droplets and the second partitions are second droplets; and prior to the merging, the at least one first droplet does not comprise the at least one second droplet.

**[0009]** The volumes of the partitions containing the sample may be different than the volumes of the partitions containing the adaptors. For example, the second volume is at least two times the volume of the first volume. In other cases, the first volume is at least two times the volume of the second volume. The methods disclosed herein may further comprise modifying the temperature of the droplets.

**[0010]** In some cases, the method further comprises merging droplets by a method comprising use of a controller such that each of the first droplets merges with each of the second droplets. In some cases, merging comprises randomly merging droplets comprising polynucleotides with droplets comprising adaptors. The methods may further comprise pooling the adaptor-tagged polynucleotides, or fragments thereof.

**[0011]** Often, the method further comprises analyzing the adaptor-tagged polynucleotides, or fragments thereof. The analyzing may involve sequencing the adaptor-tagged polynucleotides, or fragments thereof. The analyzing may comprise determining whether the adaptor-tagged polynucleotides, or fragments thereof, were located in the same partition; or, in some cases, estimating the likelihood that any two sequence reads generated by the sequencing came from the same or different partitions.

**[0012]** In some cases, the method further comprises fragmenting the polynucleotides within the second partitions to form polynucleotide fragments. The polynucleotides fragments may be generated by fragmenting the polynucleotides with an endonuclease.

**[0013]** In some cases, the polynucleotides are tagged by ligating the adaptors to the polynucleotides within a plurality of the merged partitions. The tagging may be accomplished by multiple means; for example, tagging can be accomplished using transposons.

**[0014]** Often, the methods herein include an amplification reaction. Often, the amplification comprises a polymerase chain reaction; or, the amplification can be a different type of reaction such as multiple-displacement amplification. Often, tagged polynucleotides are amplified; and, in some cases, they are amplified before tagging.

**[0015]** In some cases, each of the first partitions comprises, on average, less than five adaptors. Often, each of said second partitions comprises, on average, less than five of the multiple polynucleotides. In some cases, the subdividing of the sample comprises emulsifying or mixing the sample with the second partitions. Often, the subdividing of the plurality of adaptors comprises emulsifying or mixing the plurality of adaptors with the second partitions.

**[0016]** In some aspects, this disclosure provides a method comprising: a. partitioning organelles into a plurality of partitions, wherein each partition comprises on average less than five organelles per partition; b. lysing the extracellular organelles in the plurality of partitions, wherein the lysing releases RNA from the organelles; c. generating tagged cDNA from the released RNA in the plurality of partitions

with adaptors comprising a barcode, wherein each partition in the plurality of partitions comprises adaptors with a unique barcode. In some cases, the organelles are extracellular organelles such as exosomes. In some cases, the generating tagged cDNA comprises reverse transcription of the released RNA with partition-specific barcoded primers. The method may further comprise sequencing the tagged cDNA and/or determining if the tagged cDNA is from the same organelle.

**[0017]** This disclosure also provides a method comprising: a. partitioning microorganisms into a plurality of partitions; b. obtaining polynucleotides from the microorganisms in the plurality of partitions; and c. tagging the polynucleotides in the plurality of partitions with adaptors comprising a barcode, wherein each partition in the plurality of partitions comprises adaptors with a unique barcode. In some cases, each of said partitions comprises, on average, less than five microorganisms. The method may further comprise sequencing the tagged polynucleotides and/or determining if the tagged polynucleotide fragments are from the same partition.

#### INCORPORATION BY REFERENCE

**[0018]** All publications, patents, and patent applications mentioned in this specification are herein incorporated by reference to the same extent as if each individual publication, patent, or patent application was specifically and individually indicated to be incorporated by reference.

#### BRIEF DESCRIPTION OF THE DRAWINGS

**[0019]** Novel features are set forth with particularity in the appended claims. A better understanding of the features and advantages will be obtained by reference to the following detailed description that sets forth illustrative embodiments, in which principles are utilized, and the accompanying drawings of which:

**[0020]** FIGS. 1A and 1B illustrate methods of merging droplets comprising a sample with droplets comprising adaptors with barcodes.

**[0021]** FIG. 2 illustrates a method of determining copy number of a high copy number target using references detected with a common label.

#### DETAILED DESCRIPTION

**[0022]** In general, described herein are methods, compositions, and kits for library preparation for sequencing polynucleotides. The methods, compositions, and kits can be used to separate a sample of polynucleotides into a plurality of partitions, and each of the plurality of partitions can be provided with a unique set of adaptors comprising a barcode. Library preparation can be performed in each of the plurality of partitions (e.g., droplets). The contents of the partitions can be pooled and sequenced to generate sequence reads, and the barcodes can be used to identify which sequence reads came from the same partition. A number of embodiments of methods, compositions, systems, and kits are provided herein.

**[0023]** Overview

**[0024]** In general, barcoding (or “tagging”) can enable one to pool samples of nucleic acids in order to reduce the cost of sequencing per sample, yet retain the ability to determine from which sample a sequence read is derived. Separate library preparations can be prepared for each sample, and each sample can have its own unique barcode. The separately prepared libraries with unique barcodes can then be pooled

and sequenced. Each sequence read of the resulting dataset can be traced back to an original sample via the barcode in the sequence read.

**[0025]** In methods provided herein, polynucleotides in a sample can be separated into a plurality of partitions, e.g., droplets. Adaptors with a unique barcode (or “tag”) can be supplied to each of a plurality of partitions comprising polynucleotides. Polynucleotides with barcode adaptors can be sequenced, and the barcodes can be used to determine if two or more sequence reads were generated from one or more polynucleotides in the same partition.

**[0026]** Barcode adaptors can be bundled within a partition, e.g., an aqueous phase of an emulsion, e.g., a droplet. Barcode tagging may be accomplished by merging adaptor-filled partitions (e.g., droplets) with sample-polynucleotide-containing partitions (e.g., droplets). In some cases, adaptor-filled partitions are smaller than sample-polynucleotide-containing partitions (see e.g., FIG. 1A). Barcoded-adaptors can be separated into a plurality of partitions smaller in size than sample-polynucleotide-containing partitions. Larger sample-polynucleotide-containing partitions can be formed. A barcoded-adaptor-filled partition can be merged with a sample-polynucleotide-containing partition, and an adaptor can be attached to a polynucleotide. For example, the partitions containing sample polynucleotide may be, on average, greater than 1.5-fold, 2-fold, 2.5-fold, 3-fold, 3.5-fold, 4-fold, 4.5 fold, 5-fold, 6-fold, 7-fold, 8-fold, 9-fold, 10-fold, 20-fold, 50-fold, 75-fold, 100-fold, 500-fold, 1000-fold, 10,000-fold, 50,000-fold, or 100,000-fold the average size of the partitions containing the adaptors. The partitions containing sample polynucleotide may be, on average, greater than 1.5-fold, 2-fold, 2.5-fold, 3-fold, 3.5-fold, 4-fold, 4.5 fold, 5-fold, 6-fold, 7-fold, 8-fold, 9-fold, 10-fold, 20-fold, 50-fold, 75-fold, 100-fold, 500-fold, 1000-fold, 10,000-fold, 50,000-fold, or 100,000-fold the average volume of the partitions containing the adaptors. In some cases, sample-polynucleotide-containing partitions are formed so that they contain adaptor-filled partitions. For example, adaptor-filled partitions (e.g., droplets) can be emulsified with a polynucleotide sample so that sample-polynucleotide-containing partitions (e.g., droplets) end up containing adaptor-filled partitions. The adaptor-filled droplets can be burst (e.g., through a temperature adjustment) to release reaction components (e.g., PCR or ligation components) that can be used for library preparation. In some embodiments, the temperature adjustment comprises raising the temperature to about, more than about, or at least about 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, or 100° C. for about, more than about, or at least about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, or 60 min. In some embodiments, the temperature adjustment can last for about, more than about, or at least about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, or 24 hrs. In some cases, the adaptor-filled droplets are not contained within the sample-containing droplets. In such cases, separate droplets may be merged together.

**[0027]** In some cases, an adaptor-filled partition is larger than a sample-polynucleotide-containing partition (see e.g., FIG. 1B). Barcoded-adaptors can be separated into a plurality of partitions larger in size than sample-polynucleotide-con-

taining partitions. Smaller sample-polynucleotide-containing partitions can be formed. A barcoded-adaptor-filled partition can be merged with a sample-polynucleotide-containing partition, and an adaptor can be attached to a polynucleotide. For example, partitions containing adaptors may be, on average, greater than 1.5-fold, 2-fold, 2.5-fold, 3-fold, 3.5-fold, 4-fold, 4.5 fold, 5-fold, 6-fold, 7-fold, 8-fold, 9-fold, 10-fold, 20-fold, 50-fold, 75-fold, 100-fold, 500-fold, 1000-fold, 10,000-fold, 50,000-fold, or 100,000-fold the average size of the partitions containing the samples. The partitions containing adaptors can be, on average, greater than 1.5-fold, 2-fold, 2.5-fold, 3-fold, 3.5-fold, 4-fold, 4.5 fold, 5-fold, 6-fold, 7-fold, 8-fold, 9-fold, 10-fold, 20-fold, 50-fold, 75-fold, 100-fold, 500-fold, 1000-fold, 10,000-fold, 50,000-fold, or 100,000-fold the average volume of the partitions containing the samples. In some cases, adaptor-containing partitions are formed so that they contain sample-containing partitions. For example, in some embodiments, sample-polynucleotide-filled partitions (e.g., droplets) can be emulsified with adaptors so that adaptor-containing partitions (e.g., droplets) end up enveloping sample-containing partitions. In such cases, the sample-containing droplets can be burst (e.g., through a temperature adjustment) so that the contents of the different types of droplets can merge. In some embodiments, the temperature adjustment comprises raising the temperature to about, more than about, or at least about 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, or 100° C. for about, more than about, or at least about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, or 60 min. In some embodiments, the temperature adjustment can last for about, more than about, or at least about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, or 24 hrs. In some cases, the sample-polynucleotide-filled droplets are not contained within the adaptor-containing droplets. In such cases, separate droplets may be merged together.

**[0028]** A microfluidic device can be used to merge pre-made adaptor reagents with a plurality of sample-polynucleotide partitions such that every sample-polynucleotide partition comprises adaptor reagents. For example, a square-shaped device can be used with 1000×1000 (one million) partitions, and each polynucleotide can be tagged with two barcodes. One million unique identifiers can be constructed with 2,000 different barcodes. Reagents with 1,000 different barcodes can be loaded in horizontal channels of the device and reagents for another set of 1,000 different barcodes can be loaded in vertical channels of the device. Every one of the million partitions can have its own unique combination of barcodes.

**[0029]** In some cases, sample-polynucleotide-containing partitions (e.g., droplets) and adaptor-filled partitions (e.g., droplets) are merged in a controlled manner, e.g., one droplet of sample polynucleotides with one droplet of unique adaptors. In some cases, adaptor-filled partitions are randomly merged with sample-polynucleotide-containing partitions.

**[0030]** The following example illustrates an embodiment of a method. A large set of droplets of a number (N) types can be made. Each type of droplet can be loaded with its own barcode. The value N can be determined in part by the length of the barcode (L). For example, N can be as large as  $4^L$ .

Thus, if  $L=10$ , around 1 million different droplet types can be generated. Standard sequencing library preparation within each partition can be performed. Once the libraries are prepared, the contents of all the partitions can be merged (e.g., by breaking droplets) and loaded onto a nucleic acid sequencer. The sequencer can generate sequence reads for many of the library polynucleotides. Polynucleotides prepared within the same droplet can contain the same barcode. If the number of barcodes is sufficiently large, it can be surmised that molecules containing the same barcode came from the same partition. If N is sufficiently large (e.g., larger than the number of adaptor-filled partitions actually used in the experiment), it can be expected that any two sample-polynucleotide-containing partitions can be tagged by different adaptor-filled partitions. However, if N is not very large, sample-polynucleotide partitions can be tagged with the same adaptors. In that case, the likelihood that any two reads came from the same or different sample containing partitions can be estimated probabilistically. For many applications, a probabilistic assessment can be sufficient.

**[0031]** In some embodiments, the number of different samples that can be multiplexed, e.g., in a sequencing reaction, can be about, more than about, less than about, or at least about 1000, 5000, 10,000, 20,000, 30,000, 40,000, 50,000, 60,000, 70,000, 80,000, 90,000, 100,000, 200,000, 300,000, 400,000, 500,000, 600,000, 700,000, 800,000, 900,000, 1,000,000, 2,000,000, 3,000,000, 4,000,000, 5,000,000, 6,000,000, 7,000,000, 8,000,000, 9,000,000, or 10,000,000 samples. In some embodiments, about 1000 to about 10,000, about 10,000 to about 100,000, about 10,000 to about 500,000, about 100,000 to about 500,000, about 100,000 to about 1,000,000, about 500,000 to about 1,000,000, about 1,000,000 to about 5,000,000, or about 1,000,000 to about 10,000,000 samples are multiplexed in the methods described herein.

**[0032]** Some methods of barcode tagging are described, for example, in U.S. Patent Application Publication No. 20110033854.

**[0033]** Fusing Smaller Droplets (E.g., Inner Droplets) with Larger Droplets (E.g., Outer Droplets)

**[0034]** Some methods for fusing smaller droplets (e.g., inner droplets) with larger droplets (e.g., outer droplets) are disclosed, e.g., in U.S. Patent Application Publication No. 20110053798. In some cases, an inner droplet (or partition) can be fused with an outer droplet (or partition) by heating/cooling to change temperature, applying pressure, altering composition (e.g., via a chemical additive), applying acoustic energy (e.g., via sonication), exposure to light (e.g., to stimulate a photochemical reaction), applying an electric field, or any combination thereof. In some cases, the inner droplet may fuse to the outer droplet spontaneously. The treatment may be continuous or may vary temporally (e.g., pulsatile, shock, and/or repetitive treatment). The treatment may provide a gradual or rapid change in an emulsion parameter, to effect steady state or transient initiation of droplet fusion. The stability of the partitions, and their responsiveness to a treatment to induce droplet fusion, may be determined during their formation by selection of an appropriate surfactant type, surfactant concentration, critical micelle concentration, ionic strength, etc., for one or more phases of the inner/outer partition.

**[0035]** The fusing can result in a fused emulsion. Fusion may occur spontaneously, such that no treatment, other than a sufficient time delay (or no delay), is necessary before pro-

cessing fused droplets. Alternatively, the inner/outer droplet may be treated to controllably induce fusion of droplets to form assay mixtures.

**[0036]** The fused emulsion may be processed. Processing may include subjecting the fused emulsion to any condition or set of conditions under which at least one reaction of interest can occur (and/or is stopped), and for any suitable time period. Accordingly, processing may include maintaining the temperature of the fused emulsion near a predefined set point, varying the temperature of the fused emulsion between two or more predefined set points (such as thermally cycling the fused emulsion), exposing the fused emulsion to light, changing a pressure exerted on the fused emulsion, adding at least one chemical substance to the fused emulsion, applying an electric field to the fused emulsion, or any combination thereof, among others.

**[0037]** Signals may be detected from the fused emulsion after and/or during processing. Detection is described further in other sections herein. The signals may be detected optically, electrically, chemically, or a combination thereof, among others. The detected signals may include test signals that correspond to at least one reaction of interest performed in the fused emulsion. Alternatively, or in addition, the detected signals may include code signals that correspond to codes present in the fused emulsion. Test signals and code signals generally are distinguishable and may be detected using the same or distinct detectors. For example, the test signals and code signals each may be detected as fluorescence signals, which may be distinguishable based on excitation wavelength (or spectrum), emission wavelength (or spectrum), and/or distinct positions in a fused droplet (e.g., code signals may be detectable as more localized than test signals with respect to fused droplets), among others. As another example, the test signals and code signals may be detected as distinct optical characteristics, such as test signals detected as fluorescence and code signals detected as optical reflectance. As a further example, the test signals may be detected optically and the code signals electrically, or vice versa.

**[0038]** Adaptors

**[0039]** Barcodes can be present on adaptors, and an adaptor with a barcode can be attached to a polynucleotide by ligation. A variety of types of adaptors can be used in the methods, compositions, systems, and kits described herein. For example, an adaptor can comprise double stranded sequence. An adaptor with double stranded sequence can comprise one blunt end. In some cases, an adaptor with double stranded sequence comprises two blunt ends. An adaptor with double stranded sequence can comprise one 3' overhang. An adaptor with double stranded sequence can comprise two 3' overhangs. An adaptor with double stranded sequence can comprise one 5' overhang. In some cases, an adaptor with double stranded sequence can comprise two 5' overhangs. An adaptor with double stranded sequence can comprise a 5' overhang and a 3' overhang. In some cases, an adaptor comprises only single stranded nucleic acid.

**[0040]** When an adaptor has one or more overhangs, the overhang can be about, more than about, less than about, or at least about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, or 20 bases. For example, a 3' overhang can be about, more than about, less than about, or at least about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, or 20 bases. A 5' overhang can be about, more than about, at least about, or less than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16,

17, 18, 19, or 20 bases. If an adaptor comprises two overhangs, the overhangs can comprise the same or different number of bases.

**[0041]** The longest strand of an adaptor can be about, more than about, less than about, or at least about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, or 100 bases. If an adaptor comprises a double-stranded portion, the double stranded portion can be about, more than about, at least about, or less than about 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, or 100 base-pairs.

**[0042]** An adaptor can comprise DNA and/or RNA. In some cases, an adaptor comprises DNA. In some cases, an adaptor comprises RNA. In some cases, an adaptor comprises DNA and RNA.

**[0043]** An adaptor can comprise double stranded nucleic acid. In some cases, an adaptor comprises double stranded DNA. In some cases, an adaptor comprises double stranded RNA. In some cases, an adaptor comprises a DNA/RNA hybrid duplex.

**[0044]** An adaptor can comprise single stranded nucleic acid. In some embodiments, an adaptor comprises single stranded RNA. In some cases, an adaptor comprises single stranded DNA. In some cases, an adaptor comprises single stranded RNA and DNA.

**[0045]** When an adaptor comprises double stranded sequence, one strand of the adaptor can comprise only DNA and one strand of the adaptor can comprise only RNA. A first strand can comprise DNA and RNA and a second strand can comprise DNA only. A first strand can comprise DNA and RNA, and a second strand can comprise RNA only. If a strand of an adaptor comprises both DNA and RNA, the DNA can be 5' of the RNA or the DNA can be 3' of the RNA. In some embodiments, an adaptor is single stranded and comprises DNA and RNA, and the DNA is 5' of the RNA or 3' of the RNA.

**[0046]** An adaptor can comprise a hairpin (or hairpin loop). A hairpin can comprise DNA and/or RNA. The number of nonbase-paired bases in a loop of a hairpin can be about, more than about, or at least about 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, or 30 bases. The number of nonbase-paired bases in a loop of a hairpin can be about 4 to about 8, about 4 to about 10, about 4 to about 14, about 4 to about 16, about 4 to about 20, about 4 to about 24, about 4 to about 26, or about 4 to about 30 bases. The length of the stem (base-paired portion) of the adaptor can be about, more than about, or at least about 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, or 30 base-pairs.

**[0047]** In some cases, a hairpin adaptor is ligated to only one end of a polynucleotide. In some cases, a first hairpin adaptor is ligated to one end of a polynucleotide and a second hairpin adaptor is ligated to the other end of the polynucleotide. The hairpin adaptors that are ligated to each end of a polynucleotide can comprise the same nucleic acid sequence or different nucleic acid sequence. The hairpin adaptors that



[0049] Barcode

**[0051]** A barcode can be double stranded in an adaptor. In some cases, a barcode is single stranded in an adaptor. A barcode can comprise double stranded and single stranded sequence in an adaptor. An adaptor can comprise about, more than about, at least about, or less than about 1, 2, 3, 4, 5, 6, 7, 8, 9, or different barcodes. If an adaptor comprises more than one barcode, the barcodes can be separated from each other by about, more than about, or at least about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, or 20 bases or base pairs on the adaptor.

**[0052]** Commercially available kits comprising adaptors with barcodes can be used in the methods described herein. For example, a kit comprising adaptors with barcodes can include the ENCORE™ 384 Multiplex System (NUGEN®) which can comprise 384 molecularly barcoded library adaptors. The ENCORE™ NGS Multiplex Library Systems for ION TORRENT™ can comprise adaptors with barcodes that can be ligated to fragments. The ENCORE™ Complete RNA-Seq IL Multiplex System 1-8 (NUGEN®) and ENCORE™ Complete RNA-Seq IL Multiplex System 9-16 (NUGEN®) can provide barcoded adaptors for multiplex sequencing. The ENCORE™ Complete RNA-Seq DR Multiplex system 1-8 (NUGEN®) and ENCORE™ Complete RNA-Seq DR Multiplex system 9-16 (NUGEN®) can provide a dedicated read (DR) barcode design. Examples of kits with adaptors with barcodes from LIFE TECHNOLOGIES™ include 5500 SOLiD™ Fragment Library Barcode Adaptors 1-16, 5500 SOLiD™ Fragment Library Barcode Adaptors 1-96, 5500 SOLiD™ Fragment Library Barcode Adaptors 17-32, 5500 SOLiD™ Fragment Library Barcode Adaptors 33-48, 5500 SOLiD™ Fragment Library Barcode Adaptors 49-64, 5500 SOLiD™ Fragment Library Barcode Adaptors 65-80, 5500 SOLiD™ Fragment Library Barcode Adaptors 81-96, 5500 SOLiD™ Fragment Library Core Kit, 5500 SOLiD™ Fragment Library Standard Adaptors.

**[0053]** Other commercially available kits with adaptors with barcodes include SureSelect AB Barcode Adaptor Kit (AGILENT TECHNOLOGIES), Bioo Scientific's AIR™ Barcoded Adapters, NEXTFLEX™ DNA Barcodes, ILLUMINA® TRUSEQ™ RNA and DNA Sample Preparation Kits, RAINDANCE® Technologies DEEPSEQ™ FFPE solution, NEBNext® Multiplex Oligos for ILLUMINA® (Index Primers 1-12), or NEBNext® Multiplex Small RNA Library Prep set for ILLUMINA® (Index Primers 1-12).

**[0054]** A polynucleotide can receive a barcode by being ligated to an adaptor comprising a barcode. The ligation can involve use of one or more ligases. A barcode can be attached to a polynucleotide by amplification with a primer comprising a barcode.

**[0055]** A barcode can be adjacent to a primer binding site. A barcode can be 0 or about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, or 20 bases 3' of a primer binding (annealing, hybridization) site.

[0056] Primer/Probe Binding Site

**[0057]** An adaptor can comprise one or more primer, probe, or oligonucleotide hybridization sites. The one or more primer, probe, or oligonucleotide hybridization sites can be about, more than about, less than about, or at least about 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, or 50 bases. A primer, probe, or oligonucleotide hybridization site can be used to anneal an oligonucleotide primer to the adaptor for amplification or to anneal a primer to the adaptor for a sequencing reaction. An adaptor can comprise sequence for annealing of more than one oligonucleotide primer or probes, e.g., 2, 3, 4, 5, 6, 7, 8, 9, 10 or more oligonucleotide primers or probes. An adaptor can have a site for annealing a sequencing primer and an amplification primer. A primer or probe that anneals to an adaptor can be about, more than about, less than about, or at least about 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, or 30 bases in length.

**[0058]** Restriction Enzyme Site

**[0059]** An adaptor can comprise one or more restriction enzyme binding sites and/or cleavage sites. A restriction enzyme that can bind or cleave an adaptor can be, e.g., AatII, Acc65I, AccI, Acil, AcII, AcuI, Afel, AfIII, AfIII, AfgI, Ahdl, Alel, Alul, AlwI, AlwNI, ApaI, ApaLI, ApeKI, AgeI, AscI, Asel, AsiSI, Aval, AvaII, AvrII, BaeGI, BaeI, BamHI, BanI, BanII, BbsI, BbvCI, BbvI, BccI, BceAI, BcgI, BciVI, BclI, Bfal, BfuAI, BfuCI, BglI, BgIII, BlpI, BmgBI, Bmrl, BmtI, BpmI, Bpu10I, BpuEI, BsaAI, BsaBI, BsaHI, BsaI, BsaJI, BsaWI, BsaXI, BseRI, BseYI, BsgI, BsiEI, BsiHKAI, BsiWI, BslI, BsmAI, BsmBI, BsmFI, BsmI, BSOBI,

Bsp1286I, BspCNI, BspDI, BspEI, BspHI, BspMI, BspQI, BsrBI, BsrDI, BsrFI, BsrGI, BsrI, BssHII, BssKI, BssSI, BstAPI, BstBI, BstEII, BstNI, BstUI, BstXI, BstYI, BstZ17I, Bsu36I, BtgI, BtgZI, BtsCI, BtsI, Cac8I, ClaI, CspCI, CviAI, CviKI-1, CviQI, DdeI, DpnI, DpnII, DraI, DraIII, DrdI, EaeI, EagI, EarI, EciI, Eco53kI, EcoNI, EcoO109I, EcoP15I, EcoRI, EcoRV, FatI, FauI, Fnu4HI, FokI, FseI, FspI, HaeII, HaeIII, HgaI, HhaI, HincII, HindIII, HinfI, HinPII, HpaI, HpaII, HphI, Hpy166I, Hpy188I, Hpy188 III, Hpy99I, HpyAV, HpyCH4III, HpyCH4IV, HpyCH4V, KasI, KpnI, MboI, MboII, MfeI, MluI, MlyI, MmeU, MnlI, MscI, MseI, MslI, MspAI, MspI, MwoI, NaeI, NarI, Nb.BbvCI, Nb.BsmI, Nb.BsrDI, Nb.BtsI, NciI, NcoI, NdeI, NgoMIV, NheI, NlaII, NlaIV, NmeAIII, NotI, NruI, NsiI, NspI, Nt.AlwI, Nt.BbvCI, Nt.BsmAI, Nt.BspQI, Nt.BstNBI, Nt.CviPII, PacI, PaeR7I, PciI, PfiFI, PfiMI, PhoI, PleI, PmeI, PmlI, PpuMI, PshAI, PsiI, PspGI, PspOMI, PspXI, PstI, PvuI, PvuII, RsaI, RsrII, SacI, SacII, SalI, SapI, Sau3AI, Sau96I, SbfI, ScaI, ScrFI, SexAI, SfaNI, SfcI, SfiI, SfoI, SgrAI, SmaI, SmlI, SnaBI, SpeI, SphI, SspI, StuI, StyD4I, StyI, Swal, T, TaqI, TfiI, TliI, TseI, Tsp45I, Tsp509I, TspMI, TspRI, Tth11I, XbaI, XcmI, XhoI, XmaI, XmnI, or ZraI.

**[0060]** An adaptor can comprise a Type IIS restriction enzyme binding site. A Type IIS restriction enzyme can cleave DNA at a defined distance from a non-palindromic asymmetric recognition site. Examples of Type IIS restriction enzymes include AarI, Acc36I, AccBSI, AciI, AclWI, AcuI, AloI, Alw26I, AlwI, AsuHPI, BaeI, BbsI, BbvCI, BbvI, BccI, BceAI, BcgI, BciVI, BfiI, BfuAI, BfuI, BmgBI, Bmrl, BpiI, BpmI, Bpu10I, Bpu10I, BpuAI, BpuEI, BsaI, BsaMI, BsaXI, BseII, Bse3DI, BseGI, BseMI, BseMII, BseNI, BseRI, BseXI, BseYI, BsgI, BsmAI, BsmBI, BsmFI, BsmI, Bso31I, BspCNI, BspMI, BspQI, BspTNI, BsrBI, BsrDI, BsrI, BsrSI, BssSI, Bst2BI, Bst6I, BstF5I, BstMAI, BstV1I, BstV2I, BtgZI, BtrI, BtsCI, BtsI, CspCI, Eam1104I, EarI, EciI, Eco31I, Eco57I, Eco57MI, Esp3I, FauI, FauI, FokI, GsuI, HgaI, Hin4I, HphI, HpyAV, Ksp632I, LweI, MbiI, MboII, MlyI, Mmel, MnlI, MvaI269I, NmeAIII, PctI, PleI, PpiI, PpsI, PsrI, SapI, SchI, SfaNI, SmuI, TspDTI, TspGWI, or Taq II. A restriction enzyme can bind recognition sequence within an adaptor and cleave sequence outside the adaptor (e.g., in a polynucleotide).

**[0061]** The restriction enzyme can be a methylation sensitive restriction enzyme. The methylation sensitive restriction enzyme can specifically cleave methylated DNA. The methylation sensitive restriction enzyme can specifically cleave unmethylated DNA. A methylation sensitive enzyme can include, e.g., DpnI, Acc65I, KpnI, ApaI, Bsp120I, Bsp143I, MboI, BspOI, NheI, Cfr9I, SmaI, Csp6I, RsaI, Ecl136II, SacI, EcoRII, MvaI, HpaII, MSpJI, LpnPI, FsnEI, DpnII, McrBc, or MspI.

**[0062]** An adaptor can comprise one or more recognition sites for one or more nicking endonucleases, Type I endonucleases, or Type III endonucleases. A nicking endonuclease can hydrolyze only one strand of a duplex to produce DNA molecules that are "nicked" rather than cleaved. The nicking can result in a 3'-hydroxyl and a 5'-phosphate. Examples of nicking enzymes include Nt.CviPII, Nb.BsmI, Nb.BbvCI, Nb.BsrDI, Nb.BtsI, Nt.BsmAI, Nt.BspQI, Nt.AlwI, Nt.BbvCI, or Nt.BstNBI. A Type I endonuclease can cleave at a site that differs and is at a random distance away from the recognition site. A Type III endonuclease can recognize two

separate non-palindromic sequences that are inversely oriented. Examples of Type III restriction enzymes include EcoP15 and EcoP1.

**[0063]** One or more restriction enzymes used in the methods, compositions and/or kits described herein can be a component of a hybrid or chimeric protein. For example, a domain of a restriction enzyme comprising an enzymatic activity (e.g., endonuclease activity) can be fused to another protein, e.g., a DNA binding protein. The DNA binding protein can target the hybrid to a specific sequence on a DNA. The nucleic acid cleavage activity of the domain with enzymatic activity can be sequence specific or sequence non-specific. For example, the non-specific cleavage domain from the Type IIS restriction endonuclease FokI can be used as the enzymatic (cleavage) domain of the hybrid nuclease. The sequence the domain with the enzymatic activity can cleave can be limited by the physical tethering of the hybrid to DNA by the DNA binding domain. The DNA binding domain can be from a eukaryotic or prokaryotic transcription factor. The DNA binding domain can recognize about, or at least about 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, or 25 bases or base pairs of continuous nucleic acid sequence. The DNA binding domain can recognize about 9 to about 18 bases or base pairs of sequence. The DNA binding domain can be, e.g., a zinc finger DNA binding domain. The DNA binding domain can be from a naturally occurring protein. The DNA binding domain can be engineered to specifically bind any desired nucleotide sequence. The hybrid can be a zinc finger nuclease (e.g., zinc finger nuclease). The hybrid protein can function as a multimer (e.g., dimer, trimer, tetramer, pentamer, hexamer, etc.).

**[0064]** Modifications

**[0065]** An adaptor can comprise one or more end modifications. An adaptor can comprise one 5' phosphate. An adaptor can comprise two 5' phosphates. An adaptor can comprise one 3' hydroxyl. An adaptor can comprise two 3' hydroxyls. An adaptor can lack a 3' hydroxyl.

**[0066]** An adaptor can comprise one or more 3' end modifications. The 3' end modification can be, e.g., 3'-amino, 3'-black hole quencher (e.g., BHQ-0, BHQ-1, BHQ-2), 3'-biotin, 3'-cholesterol, 3'-dabcyl CPG, 3' dabsyl CPG, 3'-dye (e.g., fluorescein-CPG, Tamra-CPG, Rox-CPG, Cal Fluor 560-CPG, Quasar 570 (Cy3 substitute)-CPG, Quasar 670 (Cy5 substitute)-CPG, Quasar 705 (Cy5.5 substitute)-CPG, Pulsar 650-CPG, Epoch Richmond Red-CPG, Epoch Yakima Yellow-CPG, Acridine-CPG, 3'-inverted linkage (with 5'-OH attached to support and 3'-OH available for chain extension), 3'-phosphate. An adaptor can comprise any fluorescent dye described herein.

**[0067]** An adaptor can comprise one or more 5' end modifications. The 5' end modification can be, e.g., a 5'-amino group, 5'-biotin, 5'-digoxigenin (DIG), 5' phosphate group, or 5'-thiol. An adaptor can comprise a 5' aldehyde, 5' alkaline phosphatase, 5' amine, 5' horse radish peroxidase (HRP), 5' fluorescein, 5' HEX, 5' ROX, 5' TET, or 5' TAMRA. The 5' modification can be a molecular probe dye, e.g., Alexa Fluor 488, Alexa Fluor 532, Alexa Fluor 546, Alexa Fluor 555, Alexa Fluor 594, Alexa Fluor 647, Alexa Fluor 660, Alexa Fluor 750, BODIPY® FL, BODIPY® 530/550, BODIPY® 493/503, BODIPY® 558/569, BODIPY® 564/570, BODIPY® 576/589, BODIPY® 581/591, BODIPY® FL-X, BODIPY® TR-X, BODIPY® TMR, BODIPY® R6G, BODIPY® R6G-X, BODIPY® 630/650, BODIPY® 650/665, CASCADE BLUE™ Dye, MARINA BLUE™,

OREGON GREEN® 514, OREGON GREEN® 488, OREGON GREEN® 488-X, PACIFIC BLUE™ Dye, RHODAMINE GREEN™ Dye, RHODOL GREEN™ Dye, RHODAMINE GREEN™-X, RHODAMINE RED™-X, TEXAS RED®, or TEXAS RED™-X.

**[0068]** A modification can be attached to a nucleic acid strand through a linker, e.g., C1, C2, C3, C4, C5, C6, C7, C8, C9, C10, C11, C12, C13, C14, C15, C16, C17, C18, C19, or C20. A linker can be, e.g., PC (photo cleavable) spacer, hexanediol, spacer 9 (a triethylene glycol spacer), spacer 18 (an 18-atom hexa-ethyleneglycol spacer), 1', 2' dideoxyribose (dspacer), or I-linker (from Exiqon).

**[0069]** An adaptor can comprise one or more methyl groups.

**[0070]** An adaptor can be synthesized with canonical nucleotides (dATP, dCTP, dGTP, and dTTP). An adaptor can be made one or more noncanonical nucleotides. The noncanonical nucleotide can be dUTP. An adaptor can comprise deoxyuracil or deoxyinosine.

**[0071]** An adaptor can comprise one or more RNA-like nucleosides, e.g., ANA (arabino), LNA (locked), 2'-O methyl RNA, FANA (2'-fluoroarabino), or 2'-fluoro RNA. An adaptor can comprise a DNA-like nucleoside, e.g.,  $\beta$ -D-DNA,  $\beta$ -L-DNA, or  $\alpha$ -D-DNA. An adaptor can comprise one or more 5'-3' phosphorothioate linkages or inverted linkages (5'-5' or 3'-3'). An adaptor can comprise A-phosphorothioate, C-phosphorothioate, G-phosphorothioate, or T-phosphorothioate.

**[0072]** Modified bases in an adaptor can include, e.g., LNA (locked nucleic acid), 2-aminopurine, trimer-20, fluoro bases, 2,6-diaminopurine (2-amino-dA), 5-bromo dU, deoxyuridine, inverted dT, dideoxy C, 5-methyl dC, deoxyinosine, 5-nitroindole, ribo A, ribo C, ribo G, ribo U, or -2'-O-methyl RNA bases. An adaptor can have any type of nucleic acid modification described herein.

**[0073]** In some embodiments, an adaptor is chemically synthesized. In some embodiments, an adaptor is not chemically synthesized.

**[0074]** The modifications described herein can be found on sample polynucleotides.

**[0075]** Partitions

**[0076]** A partition can be formed by any mode of separating that can be used for digital PCR. A partition can be a microfluidic channel, a well on a nano- or microfluidic device or on a microtiter plate, or a reaction chamber in a microfluidic device. A partition can be an area on an array surface. A partition can be an aqueous phase of an emulsion (e.g., a droplet). Methods of generating droplets are described herein.

**[0077]** Droplet Generation

**[0078]** The present disclosure includes compositions, methods, and kits for manipulation of genetic material in droplets, e.g., using droplet digital PCR. The droplets described herein can include emulsion compositions (or mixtures of two or more immiscible fluids) described in U.S. Pat. No. 7,622,280, and droplets generated by devices described in International Application Publication No. WO/2010/036352, first inventor: Colston, each of which is hereby incorporated by reference in its entirety. The term emulsion, as used herein, can refer to a mixture of immiscible liquids (such as oil and water). Oil-phase and/or water-in-oil emulsions can allow for the compartmentalization of reaction mixtures within aqueous droplets. In some embodiments, the emulsions can comprise aqueous droplets within a continuous oil

phase. In other embodiments, the emulsions provided herein are oil-in-water emulsions, wherein the droplets are oil droplets within a continuous aqueous phase. The droplets provided herein can be used to prevent mixing between compartments, and each compartment can protect its contents from evaporation and coalescing with the contents of other compartments. One or more enzymatic reactions can occur in a droplet.

**[0079]** The mixtures or emulsions described herein can be stable or unstable. The emulsions can be relatively stable and have minimal coalescence. Coalescence can occur when small droplets combine to form progressively larger droplets. Less than about 0.00001%, 0.00005%, 0.00010%, 0.00050%, 0.001%, 0.005%, 0.01%, 0.05%, 0.1%, 0.5%, 1%, 2%, 2.5%, 3%, 3.5%, 4%, 4.5%, 5%, 6%, 7%, 8%, 9%, or 10% of droplets generated from a droplet generator can coalesce with other droplets. The emulsions can also have limited flocculation, a process by which the dispersed phase comes out of suspension in flakes.

**[0080]** Splitting a sample into small reaction volumes as described herein can enable the use of reduced amounts of reagents, thereby lowering the material cost of the analysis. Reducing sample complexity by partitioning can also improve the dynamic range of detection, since higher-abundance molecules can be separated from low-abundance molecules in different compartments, thereby allowing lower-abundance molecules greater proportional access to reaction reagents, which in turn can enhance the detection of lower-abundance molecules.

**[0081]** Droplets can be generated having an average diameter of about, more than about, less than about, or at least about 0.001, 0.01, 0.05, 0.1, 1, 5, 10, 20, 30, 40, 50, 60, 70, 80, 100, 120, 130, 140, 150, 160, 180, 200, 300, 400, or 500 microns. The average diameter of the droplets can be about 0.001 microns to about 0.01 microns, about 0.001 microns to about 0.005 microns, about 0.001 microns to about 0.1 microns, about 0.001 microns to about 1 micron, about 0.001 microns to about 10 microns, about 0.001 microns to about 100 microns, about 0.001 microns to about 500 microns, about 0.01 microns to about 0.1 microns, about 0.01 microns to about 1 micron, about 0.01 microns to about 10 microns, about 0.01 microns to about 100 microns, about 0.01 microns to about 500 microns, about 0.1 microns to about 1 micron, about 0.1 microns to about 10 microns, about 0.1 microns to about 100 microns, about 0.1 microns to about 500 microns, about 1 micron to about 10 microns, about 1 micron to about 100 microns, 1 micron to about 500 microns, about 10 microns to about 100 microns, about 10 microns to about 500 microns, or about 100 microns to about 500 microns.

**[0082]** Droplet volume can be about, more than about, less than about, or at least about 0.001 nL, 0.01 nL, 0.1 nL, 1 nL (100  $\mu\text{m}^3$ ), 10 nL, or 100 nL. Droplets can be generated using, e.g., RAINSTORM™ (RAINANCE™), microfluidics from ADVANCED LIQUID LOGIC, or ddPCR™ (BIO-RAD).

**[0083]** Microfluidic methods of producing emulsion droplets using microchannel cross-flow focusing or physical agitation can produce either monodisperse or polydisperse emulsions. The droplets can be monodisperse droplets. The droplets can be generated such that the size of said droplets does not vary by more than plus or minus 5% of the average size of said droplets. The droplets can be generated such that the size of said droplets does not vary by more than plus or minus 2% of the average size of said droplets. A droplet generator can generate a population of droplets from a single

sample, wherein none of the droplets can vary in size by more than plus or minus 0.1%, 0.5%, 1%, 1.5%, 2%, 2.5%, 3%, 3.5%, 4%, 4.5%, 5%, 5.5%, 6%, 6.5%, 7%, 7.5%, 8%, 8.5%, 9%, 9.5%, or 10% of the average size of the total population of droplets.

**[0084]** Both the flow rate in a droplet generator and the length of nucleic acids in a sample can have an impact on droplet generation. One way to decrease extension is to decrease flow rate; however, this can have the undesirable side effect of lower throughput and also increased droplet size. Long nucleic acids can disrupt droplet formation in extreme cases, resulting in a steady flow rather than discrete droplets. Reducing nucleic acid size in a sample can improve droplet formation when nucleic acid loads are high. Samples with high nucleic acid loads (e.g., high DNA loads, high RNA loads, etc.) can be used. Reducing the length of nucleic acids in a sample (e.g., by digestion, sonication, heat treatment, or shearing) can improve droplet formation.

**[0085]** Higher mechanical stability can be useful for microfluidic manipulations and higher-shear fluidic processing (e.g., in microfluidic capillaries or through 90 degree turns, such as valves, in a fluidic path). Pre- and post-thermally treated droplets or capsules can be mechanically stable to standard pipette manipulations and centrifugation.

**[0086]** A droplet can be formed by flowing an oil phase through an aqueous sample. A partition, e.g., an aqueous phase of an emulsion, can comprise a buffered solution and reagents for performing an amplification reaction, e.g., a PCR reaction, including nucleotides, primers, probe(s) for fluorescent detection, template nucleic acids, DNA polymerase enzyme, and/or reverse transcriptase enzyme.

**[0087]** A partition, e.g., an aqueous phase of an emulsion, can comprise a buffered solution and reagents for performing an enzymatic reaction (e.g., a PCR) without solid-state beads, such as magnetic-beads. The buffered solution can comprise about, more than about, at least about, or less than about 1, 5, 10, 15, 20, 30, 50, 100, or 200 mM Tris. A partition, e.g., an aqueous phase of an emulsion, can comprise one or more buffers including, e.g., TAPS, bicine, Tris, Tricine, TAPSO, HEPES, TES, MOPS, PIPES, cacodylate, SSC, ADA, ACES, choline chloride, acetamidoglycine, glycine, maleate, phosphate, CABS, piperidine, glycine, citrate, glycylglycine, malate, formate, succinate, acetate, propionate, pyridine, piperazine, histidine, bis-tris, ethanolamine, carbonate, MOPSO, imidazole, BIS-TRIS propane, BES, MOBS, triethanolamine (TEA), HEPPSO, POPSO, hydrazine, Trizma (tris), EPPS, HEPPS, bicine, HEPBS, AMPPO, taurine (AES), borate, CHES, 2-amino-2-methyl-1-propanol (AMP), ammonium hydroxide, methylamine, or MES. The pH of the partition, e.g., an aqueous phase of an emulsion, can be about 2, 2.5, 3, 3.5, 4, 4.5, 5, 5.5, 6, 6.5, 6.6, 6.7, 6.8, 6.9, 7, 7.1, 7.2, 7.3, 7.4, 7.5, 7.6, 7.7, 7.8, 7.9, 8, 8.1, 8.2, 8.3, 8.4, 8.5, 8.6, 8.7, 8.8, 8.9, 9, 9.5, 10, 10.5, 11, 11.5, 12, or 12.5. The pH of the partition, e.g., an aqueous phase of an emulsion, can be about 5 to about 9, about 5 to about 8, about 5 to about 7, about 6.5 to about 8, about 6.5 to about 7.5, about 6 to about 7, about 6 to about 9, or about 6 to about 8.

**[0088]** A partition, e.g., an aqueous phase of an emulsion, can comprise a salt, e.g., potassium acetate, potassium chloride, magnesium acetate, magnesium chloride, sodium acetate, or sodium chloride. The concentration of potassium chloride can be about, more than about, at least about, or less

than about 10, 20, 30, 40, 50, 60, 80, 100, 200 mM. The buffered solution can comprise about 15 mM Tris and about 50 mM KCl.

**[0089]** A partition, e.g., an aqueous phase of an emulsion, can comprise nucleotides. The nucleotides can comprise deoxyribonucleotide triphosphate molecules, including dATP, dCTP, dGTP, dTTP, in concentrations of about, more than about, less than about, or at least about 50, 100, 200, 300, 400, 500, 600, or 700  $\mu$ M each. dUTP can be added within a partition, e.g., an aqueous phase of an emulsion, to a concentration of about, less than about, more than about, or at least about 50, 100, 200, 300, 400, 500, 600, 700, 800, 900, or 1000  $\mu$ M. The ratio of dUTP to dTTP in a partition, e.g., an aqueous phase of an emulsion, can be about 1:1000, 1:500, 1:250, 1:100, 1:75, 1:50, 1:40, 1:30, 1:20, 1:10, 1:5, 1:4, 1:3, 1:2, or 1:1.

**[0090]** A partition, e.g., an aqueous phase of an emulsion, can comprise one or more divalent cations. The one or more divalent cations can be, e.g.,  $Mg^{2+}$ ,  $Mn^{2+}$ ,  $Cu^{2+}$ ,  $Co^{2+}$ , or  $Zn^{2+}$ . Magnesium chloride ( $MgCl_2$ ) can be added to a partition, e.g., an aqueous phase of an emulsion, at a concentration of about, more than about, at least about, or less than about 1.0, 2.0, 3.0, 4.0, or 5.0 mM. The concentration of  $MgCl_2$  can be about 3.2 mM. Magnesium sulfate can be substituted for magnesium chloride, at similar concentrations. A partition, e.g., an aqueous phase of an emulsion, can comprise both magnesium chloride and magnesium sulfate. A wide range of common, commercial PCR buffers from varied vendors can be substituted for the buffered solution.

**[0091]** A non-ionic Ethylene Oxide/Propylene Oxide block copolymer can be added to a partition, e.g., an aqueous phase of an emulsion, in a concentration of about, more than about, less than about, or at least about 0.1%, 0.2%, 0.3%, 0.4%, 0.5%, 0.6%, 0.7%, 0.8%, 0.9%, or 1.0%. A partition, or aqueous phase, can comprise a biosurfactant. Common biosurfactants include non-ionic surfactants such as Pluronic F-68, Tetronics, Zonyl FSN. Pluronic F-68 can be present at a concentration of about 0.5% w/v.

**[0092]** Additives

**[0093]** A partition, e.g., an aqueous phase of an emulsion, can comprise one or more additives including, but not limited to, non-specific background/blocking nucleic acids (e.g., salmon sperm DNA), biopreservatives (e.g., sodium azide), PCR enhancers (e.g., betaine (N,N,N-trimethylglycine; [carboxymethyl]trimethylammonium), trehalose, etc.), and/or inhibitors (e.g., RNase inhibitors). A GC-rich additive comprising, e.g., betaine and DMSO, can be added to samples assayed in the methods provided herein.

**[0094]** The one or more additives can include a non-specific blocking agent such as BSA or gelatin from bovine skin. The gelatin or BSA can be present in a concentration range of approximately 0.1 to about 0.9% w/v. Other blocking agents can include betalactoglobulin, casein, dry milk, or other common blocking agents. In some cases, the concentration of BSA and gelatin are about 0.1% w/v.

**[0095]** The one or more additives can include 2-pyrrolidone, acetamide, N-methylpyrrolidone (NMP), B-hydroxyethylpyrrolidone (HEP), propionamide, NN-dimethylacetamide (DMA), N-methylformamide (MMP), NN-dimethylformamide (DMF), formamide, N-methylacetamide (MMA), polyethylene glycol, tetramethylammonium chloride (TMAC), 7-deaza-2'-deoxyguanosine, T4 gene 32 protein, glycerol, or nonionic detergent (Triton X-100, Tween 20, Nonidet P-40 (NP-40), Tween 40, SDS

(e.g., about 0.1% SDS)), salmon sperm DNA, sodium azide, formamide, dithiothreitol (DTT), betamercaptoethanol (BME), 2-mercaptoethylamine-HCl, tris(2-carboxyethyl) phosphine (TCEP), cysteine-HCl, or a plant polysaccharide. The one or more additives can be ethanol, ethylene glycol, dimethylacetamide, dimethylformamide, or suphalane.

#### [0096] Primers

[0097] A partition, e.g., an aqueous phase of an emulsion, can comprise oligonucleotide primers. The oligonucleotide primers can be used for amplification. Primers for amplification within a partition, e.g., an aqueous phase of an emulsion, can have a concentration of about, more than about, less than about, or at least about 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, or 1.0  $\mu\text{M}$ . The concentration of each primer can be about 0.5  $\mu\text{M}$ . Primers can be designed according to known parameters for avoiding secondary structures and self-hybridization. Different primer pairs can anneal and melt at about the same temperatures, for example, within about 1, 2, 3, 4, 5, 6, 7, 8, 9 or 10° C. of another primer pair. In some cases, greater than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 15, 20, 25, 30, 35, 40, 45, 50, 100, 200, 500, 1000, 5000, 10,000 or more primers are initially used. Such primers may be able to hybridize to the genetic targets described herein. About 2 to about 10,000, about 2 to about 5,000, about 2 to about 2,500, about 2 to about 1,000, about 2 to about 500, about 2 to about 100, about 2 to about 50, about 2 to about 20, about 2 to about 10, or about 2 to about 6 primers can be used.

[0098] Primers can be prepared by a variety of methods including but not limited to cloning of appropriate sequences and direct chemical synthesis using methods well known in the art (Narang et al., *Methods Enzymol.* 68:90 (1979); Brown et al., *Methods Enzymol.* 68:109 (1979)). Primers can also be obtained from commercial sources such as Integrated DNA Technologies, Operon Technologies, Amersham Pharmacia Biotech, Sigma, or Life Technologies. The primers can have an identical melting temperature. The melting temperature of a primer can be about 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 81, 82, 83, 84, or 85° C. The melting temperature of a primer can be about 30 to about 85° C., about 30 to about 80° C., about 30 to about 75° C., about 30 to about 70° C., about 30 to about 65° C., about 30 to about 60° C., about 30 to about 55° C., about 30 to about 50° C., about 40 to about 85° C., about 40 to about 80° C., about 40 to about 75° C., about 40 to about 70° C., about 40 to about 65° C., about 40 to about 60° C., about 40 to about 55° C., about 40 to about 50° C., about 50 to about 85° C., about 50 to about 80° C., about 50 to about 75° C., about 50 to about 70° C., about 50 to about 65° C., about 50 to about 60° C., about 50 to about 55° C., about 52 to about 60° C., about 52 to about 58° C., about 52 to about 56° C., or about 52 to about 54° C.

[0099] The lengths of the primers can be extended or shortened at the 5' end or the 3' end to produce primers with desired melting temperatures. One of the primers of a primer pair can be longer than the other primer. The 3' annealing lengths of the primers, within a primer pair, can differ. Also, the annealing position of each primer pair can be designed such that the sequence and length of the primer pairs yield the desired melting temperature. An equation for determining the melting temperature of primers smaller than 25 base pairs is the Wallace Rule ( $T_d = 2(A+T) + 4(G+C)$ ). Computer programs can also be used to design primers, including but not limited to Array Designer Software (Arrayit Inc.), Oligonucleotide

Probe Sequence Design Software for Genetic Analysis (Olympus Optical Co.), NetPrimer, and DNAsis from Hitachi Software Engineering. The TM (melting or annealing temperature) of each primer can be calculated using software programs such as Net Primer (free web based program at <http://www.premierbiosoft.com/netprimer/index.html>). The annealing temperature of the primers can be recalculated and increased after any cycle of amplification, including but not limited to about cycle 1, 2, 3, 4, 5, about cycle 6 to about cycle 10, about cycle 10 to about cycle 15, about cycle 15 to about cycle 20, about cycle 20 to about cycle 25, about cycle 25 to about cycle 30, about cycle 30 to about cycle 35, or about cycle 35 to about cycle 40. After the initial cycles of amplification, the 5' half of the primers can be incorporated into the products from each loci of interest; thus the TM can be recalculated based on both the sequences of the 5' half and the 3' half of each primer.

[0100] The annealing temperature of the primers can be recalculated and increased after any cycle of amplification, including but not limited to about cycle 1, 2, 3, 4, 5, about cycle 6 to about cycle 10, about cycle 10 to about cycle 15, about cycle 15 to about cycle 20, about cycle 20 to about cycle 25, about cycle 25 to about cycle 30, about cycle 30 to about cycle 35, or about cycle 35 to about cycle 40. After the initial cycles of amplification, the 5' half of the primers can be incorporated into the products from each loci of interest, thus the TM can be recalculated based on both the sequences of the 5' half and the 3' half of each primer.

#### [0101] Probes

[0102] A partition, e.g., an aqueous phase of an emulsion, can comprise one or more probes for fluorescent detection. The concentration of each of the one or more probes can be about, more than about, at least about, or less than about 0.1, 0.2, 0.3, 0.4, or 0.5  $\mu\text{M}$ . The concentration of the one or more probes for fluorescent detection can be about 0.25  $\mu\text{M}$ . Amenable ranges for target nucleic acid concentrations in PCR can be between about 1 pg and about 500 ng. A probe can be about, more than about, less than about, or at least about, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, or 40 bases long. A probe can be about 8 to about 40, about 10 to about 40, about 10 to about 35, about 10 to about 30, about 10 to about 25, about 10 to about 20, about 15 to about 40, about 15 to about 35, about 15 to about 30, about 15 to about 25, about 15 to about 20, about 18 to about 40, about 18 to about 35, about 18 to about 30, about 18 to about 25, or about 18 to 22 bases. A label (fluorophore, dye) used on a probe (e.g., a Taqman probe) can be, e.g., 6-carboxyfluorescein (FAM), tetrachlorofluorescein (TET), 4,7,2'-trichloro-7'-phenyl-6-carboxyfluorescein (VIC), HEX, Cy3, Cy 3.5, Cy 5, Cy 5.5, Cy 7, tetramethylrhodamine, ROX, and JOE, Alexa Fluor dye, e.g., Alexa Fluor 350, 405, 430, 488, 532, 546, 555, 568, 594, 633, 647, 660, 680, 700, and 750; Cascade Blue, Marina Blue, Oregon Green 500, Oregon Green 514, Oregon Green 488, Oregon Green 488-X, Pacific Blue, Rhodamine Green, Rhodol Green, Rhodamine Green-X, Rhodamine Red-X, and Texas Red-X. The label can be at the 5' end of a probe, 3' end of the probe, at both the 5' and 3' end of a probe, or internal to the probe. A unique label can be used to detect each different locus in an experiment. A probe, e.g., a Taqman probe, can comprise a quencher, e.g., a 3' quencher. The 3' quencher can be, e.g., TAMARA, DABCYL, BHQ-1, BHQ-2, or BHQ-3. In some cases, a quencher used in the methods provided herein is a black hole quencher (BHQ). In some cases, the

quencher is a minor groove binder (MGB). In some cases, the quencher is a fluorescent quencher. In other cases, the quencher is a non-fluorescent quencher (NFQ).

**[0103] Polymerases**

**[0104]** A partition, e.g., an aqueous phase of an emulsion, can comprise a polymerase. The polymerase can be a DNA polymerase. The DNA polymerase can be, e.g., T4 DNA polymerase, DEEP VENT™ DNA polymerase, LONGAMP® Tag, PHUSION® High Fidelity DNA polymerase, LONGAMP® Hot Start Taq, Crimson LONGAMP® Taq, Taq DNA polymerase, Crimson Taq DNA polymerase, ONE-TAQ® DNA polymerase, QUICK-LOAD® DNA polymerase, VENTR® DNA polymerase, Hemo KLENTAQ®, Bsu DNA polymerase, DNA polymerase I, DNA Polymerase I, Large (Klenow), Klenow Fragment, Phi29 DNA polymerase, Pfu DNA polymerase, Pfx DNA polymerase, Tth DNA polymerase, Vent DNA polymerase, bacteriophage 29, REDTaq™, or T7 DNA polymerase. The DNA polymerase can comprise 3' to 5' exonuclease activity. The DNA polymerase can comprise 5' to 3' exonuclease activity. The DNA polymerase can comprise both 3' to 5' exonuclease activity and 5' to 3' exonuclease activity. The DNA polymerase can comprise neither 3' to 5' exonuclease activity nor 5' to 3' exonuclease activity. The DNA polymerase can comprise strand displacement activity. In some cases, the DNA polymerase does not comprise strand displacement activity. The error rate of the DNA polymerase can be less than  $1 \times 10^{-6}$  bases.

**[0105]** A partition, e.g., an aqueous phase of an emulsion, can comprise a reverse transcriptase. The reverse transcriptase can be AMV reverse transcriptase or M-MuLV reverse transcriptase. The RNA polymerase can comprise 5' to 3' exonuclease activity. The reverse transcriptase can comprise both 3' to 5' exonuclease activity and 5' to 3' exonuclease activity. The reverse transcriptase can comprise neither 3' to 5' exonuclease activity nor 5' to 3' exonuclease activity. The reverse transcriptase can comprise strand displacement activity. In some embodiments, the reverse transcriptase does not comprise strand displacement activity.

**[0106]** A partition, e.g., an aqueous phase of an emulsion, can comprise an RNA polymerase. The RNA polymerase can be, e.g., phi6 RNA polymerase, SP6 RNA polymerase, or T7 RNA polymerase.

**[0107]** In some embodiments, a partition, e.g., an aqueous phase of an emulsion, comprises Poly(U) polymerase or Poly(A) polymerase.

**[0108] Oil Phase**

**[0109]** The oil phase can comprise a fluorinated base oil which can be additionally stabilized by combination with a fluorinated surfactant such as a perfluorinated polyether. The base oil can be one or more of HFE 7500, FC-40, FC-43, FC-70, or another common fluorinated oil. The anionic surfactant can be Ammonium Krytox (Krytox-AM), the ammonium salt of Krytox FSH, or morpholino derivative of Krytox-FSH. Krytox-AS can be present at a concentration of about, more than about, less than about, or at least about 0.1%, 0.2%, 0.3%, 0.4%, 0.5%, 0.6%, 0.7%, 0.8%, 0.9%, 1.0%, 2.0%, 3.0%, or 4.0% w/w. The concentration of Krytox-AS can be about 1.8%. The concentration of Krytox-AS can be about 1.62%. Morpholino derivative of Krytox-FSH can be present at a concentration of about, more than about, less than about, or at least about 0.1%, 0.2%, 0.3%, 0.4%, 0.5%, 0.6%, 0.7%, 0.8%, 0.9%, 1.0%, 2.0%, 3.0%, or 4.0% w/w. The concentration of morpholino derivative of Krytox-FSH can be about

1.8%. In some embodiments, the concentration of morpholino derivative of Krytox-FSH is about 1.62%.

**[0110]** The oil phase can comprise an additive for tuning the oil properties, such as vapor pressure or viscosity or surface tension. Nonlimiting examples include perfluoro-ocetanol and 1H,1H,2H,2H-Perfluorodecanol. In some embodiments, 1H,1H,2H,2H-Perfluorodecanol is added to a concentration of about 0.05%, 0.06%, 0.07%, 0.08%, 0.09%, 1.00%, 1.25%, 1.50%, 1.75%, 2.00%, 2.25%, 2.50%, 2.75%, or 3.00% w/w. In some embodiments, 1H,1H,2H,2H-Perfluorodecanol is added to a concentration of 0.18% w/w.

**[0111] Microcapsules**

**[0112]** In some embodiments, the emulsion is formulated to produce highly monodisperse droplets having a liquid-like interfacial film that can be converted by heating into microcapsules having a solid-like interfacial film; such microcapsules can behave as bioreactors able to retain their contents through a reaction process such as PCR amplification. The conversion to microcapsule form can occur upon heating. For example, such conversion can occur at a temperature of greater than about, more than about, or at least about 50, 60, 70, 80, 90, or 95 degrees Celsius. In some embodiments this heating occurs using a thermocycler. During the heating process, a fluid or mineral oil overlay can be used to prevent evaporation. Excess continuous phase oil may or may not be removed prior to heating. The biocompatible capsules can be resistant to coalescence and/or flocculation across a wide range of thermal and mechanical processing.

**[0113]** Following conversion, the capsules can be stored at about, more than about, less than about, or at least about 3, 4, 5, 6, 7, 8, 9, 10, 15, 20, 25, 30, 35, or 40 degrees Celsius, with one embodiment comprising storage of capsules at less than about 25 degrees Celsius. In some embodiments, these capsules are useful in biomedical applications, such as stable, digitized encapsulation of macromolecules, particularly aqueous biological fluids containing a mix of nucleic acids or protein, or both together; drug and vaccine delivery; biomolecular libraries; clinical imaging applications, and others.

**[0114]** The microcapsules can contain one or more nucleic acid probes (e.g., molecular inversion probe, ligation probe, etc.) and can resist coalescence, particularly at high temperatures. Accordingly, PCR amplification reactions can occur at a very high density (e.g., number of reactions per unit volume). In some embodiments, greater than 100,000, 500,000, 1,000,000, 1,500,000, 2,000,000, 2,500,000, 5,000,000, or 10,000,000 separate reactions can occur per ml. In some embodiments, the reactions occur in a single well, e.g., a well of a microtiter plate, without inter-mixing between reaction volumes. The microcapsules can also contain other components necessary to enable an enzymatic reaction (e.g., a PCR reaction) to occur, e.g., nucleotides, primers, probes, dNTPs, DNA or RNA polymerases, reverse transcriptases, restriction enzymes, etc. These capsules exhibit resistance to coalescence and flocculation across a wide range of thermal and mechanical processing.

**[0115]** The compositions described herein can include compositions comprising mixtures of two or more immiscible fluids such as oil and water that contain a type of nucleic acid probe (e.g., TaqMan probe, molecular inversion probe, ligation probe, etc.). In some cases, the composition comprises a restriction enzyme described herein, e.g., a droplet comprising a restriction enzyme (e.g., methylation-sensitive enzyme). In other embodiments, the compositions described herein comprise microcapsules that contain a type of nucleic

[0117] Library preparation within partitions (e.g., droplets) can entail fragmentation of polynucleotides in a sample and ligation of adaptors. Generally, the fragmentation occurs within a partition (e.g., droplet); but, in some applications, the fragmentation may occur prior to the partitioning. Fragmentation can be accomplished enzymatically, e.g., using an endonuclease. The endonuclease can be, e.g., AatII, Acc65I, AccI, AciI, AclI, AclI, AclI, AfeI, AflII, AflIII, AgeI, AhdI, AleI, AluI, AlwI, AlwNI, ApaI, ApaLI, ApeKI, ApoI, AscI, AseI, AsiSI, Aval, AvaII, AvrII, BaeGI, BaeI, BamHI, BanI, BanII, BbsI, BbvCI, BbvI, Bed, BceAI, BcgI, BciVI, BclI, BfaI, BfuAI, BfuCI, BglI, BglII, BglIII, BlnI, BmgBI, BmrI, BmtI, BpmI, Bpu10I, BpuEI, BsaAI, BsaBI, BsaHI, BsaI, BsaJI, BsaWI, BsaXI, BseRI, BseYI, BsgI, BsiEI, BsiHKAI, BsiWI, BslI, BsmAI, BsmBI, BsmFI, BsmI, BsoBI, Bsp1286I, BspCNI, BspDI, BspEI, BspHI, BspMI, BspQI, BsrBI, BsrDI, BsrFI, BsrGI, BsrI, BssHII, BssKI, BssSI, BstAPI, BstBI, BstEII, BstNI, BstUI, BstXI, BstYI, BstZI 7I, Bsu36I, BtgI, BtgZI, BtsCI, BtsI, Cac8I, ClaI, CspCI, CviAI, CviKI-1, CviQI, DdeI, DpnI, DpnII, DraI, DraII, DrdI, EaeI, EagI, EarI, EciI, Eco53kI, EcoNU, EcoO109I, EcoP15I, EcoRI, EcoRV, FatI, FaeI, Fnu4HI, FokI, FseI, FspI, HaeII, HaeIII, HgaI, HhaI, HincII, HincIII, HinfI, HinfPI, HpaI, HpaII, HphI, Hpy166I, Hpy188I, Hpy188III, Hpy99I, HpyAV, HpyCH4III, HpyCH4IV, HpyCH4V, KasI, KpnI, MboI, MboII, MfeI, MluI, MlyI, MmeI, MnlI, MscI, MseI, MslI, MspAI, MspI, MwoI, NaeI, NatI, Nb.BbvCI, Nb.BsmI, Nb.BsrDI, Nb.BtsI, NciI, NcoI, NdeI, NgoMIV, NheI, NlaIII, NlaIV, NmeAIII, NotI, NruI, NsiI, NspI, Nt.AiwI, Nt.BbvCI, Nt.BsmAI, Nt.BspQI, Nt.BstNBI, Nt.CviPII, PacI, PaeR7I, PciI, PfiFI, PfiMI, PfoI, PleI, PmeI, PmlI, PpuMI, PshAI, PsiI, PspGI, PspOMI, PspXI, PstI, PvuI, PvuII, RsaI, RsrII, SacI, SacII, SalI, SapI, Sau3AI, Sau96I, SbfI, Scal, ScrFI, SexAI, SfaNI, SfeI, SfcI, SfoI, SgrAI, SmaI, SmlI, SnaBI, SphI, SphI, SspI, StuI, StyD4I, StyO, SwaI, T, TaqoI, TfiI, TfiI, TseI, Tsp45I, Tsp509I, TspMI, TspRI, Tth111I, XbaI, XcmI, XhoI, XmaI, XmnI, or ZraI

**[0118]** In some embodiments, the fragmentation is mechanical fragmentation. In some embodiments, shear forces created during lysis or extraction can mechanically fragment polynucleotides. Fragmentation can be accomplished by, e.g., sonication, heat treatment, or shearing. In some embodiments, mechanical fragmentation is by nebulization.

**[0119]** In some embodiments, the endonuclease is a methylation sensitive restriction enzyme. In some embodiments, the methylation sensitive restriction enzyme specifically cleaves methylated polynucleotides. In some embodiments, the methylation sensitive restriction enzyme specifically cleaves unmethylated polynucleotides. A methylation sensitive enzyme can include, e.g., DpnI, Acc65I, KpnI, ApaI, Bsp120I, Bsp143I, MboI, BspOI, NheI, Cfr9I, SmaI, Csp6I, RsaI, Ecl136II, SacI, EcoRII, MvaI, HpaII, MspJI, LpnPI, FsnEI, DpnII, McrBc, or MspI.

**[0120]** In some embodiments, fragmentation of a polynucleotide is accomplished by introducing one or more non-canonical nucleotides (e.g., dUTP) into a polynucleotide, generating one or more abasic sites by cleaving the base of the

non-canonical nucleotide (e.g., using, e.g., Uracil N-Glycosylase (UNG) or Uracil DNA glycosylase (UDG)), and fragmenting the polynucleotide at the one or more abasic sites. The fragmenting can be by an enzymatic agent or a chemical agent. The chemical agent can be, e.g., a polyamine, e.g., N,N'-dimethylethylenediamine (DMED). The enzymatic agent can be, e.g., apurinic/apyrimidinic endonuclease (APE 1). In some embodiments, fragmentation can be accomplished as described in U.S. Patent Application Publication Nos. 20110033854 or 20100022403.

[0122] Fragmentation can be followed by a step of ligating adaptors to polynucleotides. In some embodiments, a ligation step does not following a fragmentation step. A partition, e.g., an aqueous phase of an emulsion, can comprise a ligase. The ligase can be, e.g., T4 DNA ligase, *E. coli* DNA ligase, Taq DNA ligase, 9° N<sup>TM</sup> DNA ligase, T4 RNA ligase 1 (ssRNA ligase), T4 RNA ligase 2 (dsRNA ligase), or T4 RNA Ligase 2, truncated (NEB).

**[0123]** A partition, e.g., an aqueous phase of an emulsion, can comprise reagents for a ligation reaction, e.g., buffer, salt, and/or reducing agent. Ligase and other reagents can be supplied in a partition, e.g., an aqueous phase of an emulsion, separate from a partition, e.g., an aqueous phase of an emulsion, comprising polynucleotides. A partition, e.g., an aqueous phase of an emulsion, comprising ligase can be merged with a partition, e.g., an aqueous phase of an emulsion, comprising polynucleotides.

**[0124]** The ligation reaction can occur at a temperature of about, more than about, less than about, or at least about 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, or 99° C. The ligation can occur at about 4° C. to about 16° C., about 16° C. to about 25° C., about 25° C. to about 30° C., about 25° C. to about 37° C., about 37° C. to about 45° C., about 37° C. to about 50° C., or about 50° C. to about 65° C.

**[0125]** The ligation reaction can occur for a time period of about, more than about, less than about, or at least about 5 min, 15 min, 30 min, 45 min, or 60 min. The ligation reaction can occur for a time period of about, more than about, less than about, or at least about 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, or 48 hr. The ligation reaction can last for about 5 min to about 15 min, about 5 min to about 30 min, about 5 min to about 45 min, about 5 min to about 60 min, about 30 min to about 60 min, about 30 min to about 90 min, about 30 min to about 120 min, about 1 hr to about 2 hr, about 1 hr to about 6 hr, about 1 hr to about 12 hr, about 12 hr to about 24 hr, or about 12 hr to about 48 hr.

**[0126]** A transposon-based approach such as that provided by NEXTERA™ can be used in which DNA is fragmented and an adaptor can be ligated in a single step reaction (see e.g., [http://www.epibio.com/newsletter/16-3\\_4-6.pdf](http://www.epibio.com/newsletter/16-3_4-6.pdf)). A TRANSPOSOME™ complex can comprise free transposon ends and a transposase. A TRANSPOSOME™ complex can be incubated with target double strand DNA, and the target can be fragmented. The transferred strand of a transposon end oligonucleotide can be covalently attached to the 5' end of a target fragment. Transposon integration and strand transfer



can occur via a staggered, dsDNA break within a target polynucleotide. The resulting fragments can have single-stranded gaps. The concentration of TRANSPOSOME™ complexes can be varied to control the size distribution of the fragmented and tagged DNA library. The transposon ends can comprise barcodes. Adaptor ligation can be followed with PCR amplification of ligated products to increase their concentrations.

**[0127]** The NEXTERA™ technology can be used to generate di-tagged libraries. The libraries can be optionally bar-coded. The libraries can be compatible, e.g., with Roche/454 or ILLUMINA®/SOLEXA® sequencing platforms. To generate platform-specific libraries, free transposon ends or appended transposon ends can be used. Platform specific tags, and optional barcoding, can be introduced by, e.g., PCR. Amplification can occur by, e.g., emulsion PCR (emPCR) or bridge PCR (bPCR).

**[0128]** In some embodiments, the methods of ligating adaptors to polynucleotides are those described in U.S. Pat. Nos. 5,789,206 or Arneson et al. (2008) Whole-Genome Amplification by Adaptor-Ligation PCR of Randomly Sheared Genomic DNA (PRSG) *Cold Spring Harbor Protocols*.

**[0129]** Sizes of fragments of polynucleotides that can be generated can be about, more than about, at least about, or less than about 10, 25, 50, 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10,000, 20,000, 30,000, 40,000, 50,000, 60,000, 70,000, 80,000, 90,000, 100,000, 200,000, 300,000, 400,000, 500,000, 600,000, 700,000, 800,000, 900,000, 1,000,000, 2,000,000, 3,000,000, 4,000,000, 5,000,000, 6,000,000, 7,000,000, 8,000,000, 9,000,000, or 10,000,000 bases or base pairs. In some embodiments, the size of fragmented polynucleotides is about 50 to about 100, about 50 to about 150, about 50 to about 200, about 100 to about 150, about 100 to about 200, about 100 to about 300, about 150 to about 200, about 150 to about 250, about 200 to about 300, about 200 to about 400, about 300 to about 400, about 300 to about 500, about 400 to about 500, about 400 to about 600, about 500 to about 600, about 500 to about 700, about 600 to about 700, about 600 to about 800, about 700 to about 800, about 700 to about 900, about 800 to about 1000, about 50 to about 500, about 100 to about 500, about 100 to about 1000, about 50 to about 1500, about 50 to about 2000, about 1000 to about 2000, or about 1500 to about 2000 bases or base pairs. In some embodiments, the size of fragmented polynucleotides is about 1000 to about 5000, about 1000 to about 10,000, about 10,000 to about 20,000, about 10,000 to about 50,000, about 10,000 to about 100,000, about 50,000 to about 100,000, about 100,000 to about 200,000, about 100,000 to about 500,000, about 100,000 to about 1,000,000, about 200,000 to about 1,000,000, about 300,000 to about 1,000,000, about 400,000 to about 1,000,000, about 500,000 to about 1,000,000, or about 750,000 to about 1,000,000 bases or base pairs.

**[0130]** Amplification

**[0131]** Polynucleotides may be amplified before they are partitioned. In some embodiments, polynucleotides are amplified while in a partition (e.g., aqueous phase of an emulsion, e.g., droplet). In some embodiments, polynucleotides are amplified before fragmentation in a partition. In some embodiments, polynucleotides are amplified after fragmentation in a partition. In some embodiments, polynucleotides are amplified both before and after fragmentation in a partition. In some embodiments, polynucleotides are amplified in a partition before ligating an adaptor to a polynucle-

otide in a partition. In some embodiments, polynucleotides are amplified in a partition after ligating an adaptor to the polynucleotide in the partition. In some embodiments, polynucleotides are amplified after ligating an adaptor to the polynucleotides and pooling polynucleotides from different partitions.

**[0132]** In some embodiments, the amplification comprises polymerase chain reaction (PCR), digital PCR, reverse-transcription PCR, quantitative PCR, real-time PCR, isothermal amplification, linear amplification, or isothermal linear amplification, quantitative fluorescent PCR (QF-PCR), multiplex fluorescent PCR (MF-PCR), single cell PCR, restriction fragment length polymorphism PCR (PCR-RFLP), PCR-RFLP/RT-PCR-RFLP, hot start PCR, nested PCR, in situ polony PCR, in situ rolling circle amplification (RCA), bridge PCR (bPCR), picotiter PCR, digital PCR, droplet digital PCR, or emulsion PCR (emPCR). Other suitable amplification methods include ligase chain reaction (LCR (oligonucleotide ligase amplification (OLA))), transcription amplification, cycling probe technology (CPT), molecular inversion probe (MIP)PCR, self-sustained sequence replication, selective amplification of target polynucleotide sequences, consensus sequence primed polymerase chain reaction (CP-PCR), arbitrarily primed polymerase chain reaction (AP-PCR), transcription mediated amplification (TMA), degenerate oligonucleotide-primed PCR (DOP-PCR), multiple-displacement amplification (MDA), strand displacement amplification (SDA), and nucleic acid based sequence amplification (NABSA). Other amplification methods that can be used herein include those described in U.S. Pat. Nos. 5,242,794; 5,494,810; 4,988,617; and 6,582,938.

**[0133]** In some embodiments, a multiple-displacement amplification (MDA) step can be performed within a partition (e.g., droplet) prior to fragmentation of polynucleotides and adaptor ligation to amplify the amount of DNA in each droplet in order to cover more of the captured polynucleotides. MDA can be a non-PCR based amplification technique that can involve annealing multiple primers (e.g., hexamer primers) to a polynucleotide template, and initiating DNA synthesis (e.g., using Phi 29 polymerase). When DNA synthesis proceeds to the next synthesis starting site, the polymerase can displace the newly produced DNA strand and continues its strand elongation. Strand displacement can generate newly synthesized single stranded DNA template to which other primers can anneal. Further primer annealing and strand displacement on the newly synthesized template can result in a hyper-branched DNA network. The sequence debranching during amplification can result in a high yield of products. To separate the DNA branching network, one or more S1 nucleases can be used to cleave the fragments at displacement sites. The nicks on the resulting DNA fragments can be repaired by DNA polymerase I. The generated DNA fragments can be directly used for analysis or be ligated to generate genomic libraries for further sequencing analysis. MDA is described, e.g., in U.S. Pat. No. 7,074,600.

**[0134]** Amplification of polynucleotides can occur on a bead. In other embodiments, amplification does not occur on a bead. A hot start PCR can be performed wherein the reaction is heated to 95° C. for two minutes prior to addition of the polymerase or the polymerase can be kept inactive until the first heating step in cycle 1. Hot start PCR can be used to minimize nonspecific amplification. Other strategies for and aspects of amplification suitable for use in the methods described herein are described in U.S. Patent Application



Publication No. 2010/0173394 A1, published Jul. 8, 2010, which is incorporated herein by reference.

**[0135]** Any number of PCR cycles can be used to amplify the DNA, e.g., about, more than about, at least about, or less than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44 or 45 cycles. The number of amplification cycles can be about 1 to about 45, about 10 to about 45, about 20 to about 45, about 30 to about 45, about 35 to about 45, about 10 to about 40, about 10 to about 30, about 10 to about 25, about 10 to about 20, about 10 to about 15, about 20 to about 35, about 25 to about 35, about 30 to about 35, or about 35 to about 40.

**[0136]** Thermocycling reactions can be performed on samples contained in droplets. The droplets can remain intact during thermocycling. Droplets can remain intact during thermocycling at densities of greater than about 10,000 droplets/mL, 100,000 droplets/mL, 200,000 droplets/mL, 300,000 droplets/mL, 400,000 droplets/mL, 500,000 droplets/mL, 600,000 droplets/mL, 700,000 droplets/mL, 800,000 droplets/mL, 900,000 droplets/mL or 1,000,000 droplets/mL. Droplets can remain intact during thermocycling at densities of greater than about 10,000 droplets/mL to about 100,000 droplets/mL, 10,000 droplets/mL to about 1,000,000 droplets/mL, or about 100,000 droplets/mL to about 1,000,000 droplets/mL. In other cases, two or more droplets may coalesce during thermocycling. In other cases, greater than about 100 or greater than about 1,000 droplets may coalesce during thermocycling.

**[0137]** Polynucleotides

**[0138]** The methods described herein can be used for manipulating and analyzing polynucleotides. The term polynucleotide, or grammatical equivalents, can refer to at least two nucleotides covalently linked together. A nucleic acid described herein can contain phosphodiester bonds, although in some cases, as outlined herein (for example in the construction of primers and probes such as label probes), nucleic acid analogs are included that can have alternate backbones, comprising, for example, phosphoramidate (see e.g., Beaucage et al., *Tetrahedron* 49(10):1925 (1993) and references therein; Letsinger, *J. Org. Chem.* 35:3800 (1970); Sprinzl et al., *Eur. J. Biochem.* 81:579 (1977); Letsinger et al., *Nucl. Acids Res.* 14:3487 (1986); Sawai et al., *Chem. Lett.* 805 (1984); Letsinger et al., *J. Am. Chem. Soc.* 110:4470 (1988); and Pauwels et al., *Chemica Scripta* 26:141 (1986)), phosphorothioate (Mag et al., *Nucleic Acids Res.* 19:1437 (1991); and U.S. Pat. No. 5,644,048), phosphorodithioate (Briu et al., *J. Chem. Soc.* 111:2321 (1989), O-methyl phosphoramidate linkages (see e.g., Eckstein, *Oligonucleotides and Analogues: A Practical Approach*, Oxford University Press), and peptide nucleic acid (also referred to herein as "PNA") backbones and linkages (see e.g., Egholm, *J. Am. Chem. Soc.* 114:1895 (1992); Meier et al., *Chem. Int. Ed. Engl.* 31:1008 (1992); Nielsen, *Nature*, 365:566 (1993); Carlsson et al., *Nature* 380:207 (1996), all of which are incorporated by reference). Other analog nucleic acids include those with bicyclic structures including locked nucleic acids (LNAs are a class of nucleic acid analogues in which the ribose ring is "locked" by a methylene bridge connecting the 2'-O atom with the 4'-C atom), also referred to herein as "LNA" (see e.g., Koshkin et al., *J. Am. Chem. Soc.* 120:13252 (1998)); positive backbones (Denpoy et al., *Proc. Natl. Acad. Sci. USA* 92:6097 (1995); non-ionic backbones (see e.g., U.S. Pat. Nos. 5,386,023, 5,637,684, 5,602,240, 5,216,141 and 4,469,863;

Kiedrowski et al., *Angew. Chem. Intl. Ed. English* 30:423 (1991); Letsinger et al., *J. Am. Chem. Soc.* 110:4470 (1988); Letsinger et al., *Nucleoside & Nucleotide* 13:1597 (1994); Chapters 2 and 3, ASC Symposium Series 580, "Carbohydrate Modifications in Antisense Research", Ed. Y. S. Sanghui and P. Dan Cook; Mesmaeker et al., *Bioorganic & Medicinal Chem. Lett.* 4:395 (1994); Jeffs et al., *J. Biomolecular NMR* 34:17 (1994); *Tetrahedron Lett.* 37:743 (1996)), and non-ribose backbones, including those described in U.S. Pat. Nos. 5,235,033 and 5,034,506, and Chapters 6 and 7, ASC Symposium Series 580, "Carbohydrate Modifications in Antisense Research", Ed. Y. S. Sanghui and P. Dan Cook. Nucleic acids containing one or more carbocyclic sugars are also included within the definition of nucleic acids (see e.g., Jenkins et al., *Chem. Soc. Rev.* (1995) pp 169-176). Several nucleic acid analogs are described in Rawls, *C & E News* Jun. 2, 1997 page 35. All of these references are hereby expressly incorporated by reference. These modifications of the ribose-phosphate backbone can be done to increase the stability and half-life of such molecules in physiological environments. For example, PNA:DNA and LNA-DNA hybrids can exhibit higher stability and thus can be used in some embodiments. The target nucleic acids can be single stranded or double stranded, as specified, or contain portions of both double stranded or single stranded sequence. Depending on the application, the nucleic acids can be DNA (including, e.g., genomic DNA, mitochondrial DNA, and cDNA), RNA (including, e.g., mRNA and rRNA) or a hybrid, where the nucleic acid can contain any combination of deoxyribo- and ribo-nucleotides, and any combination of bases, including uracil, adenine, thymine, cytosine, guanine, inosine, xanthine hypoxanthine, isocytosine, isoguanine, etc.

**[0139]** The methods, compositions, and kits provided herein can be used to analyze polynucleotides (e.g., DNA, RNA, mitochondrial DNA, genomic DNA, mRNA, siRNA, miRNA, crRNA, single-stranded DNA, double-stranded DNA, single-stranded RNA, double-stranded RNA, tRNA, rRNA, cDNA, etc.). The methods, compositions and kits can be used to evaluate a quantity of a first polynucleotide compared to the quantity of a second polynucleotide. The methods can be used to analyze the quantity of synthetic plasmids in a solution; to detect a pathogenic organism (e.g., microbe, bacteria, virus, parasite, retrovirus, lentivirus, HIV-1, HIV-2, influenza virus, etc.) within a sample obtained from a subject or obtained from an environment. The methods also can be used in other applications wherein a rare population of polynucleotides exists within a larger population of polynucleotides. Polynucleotides can be obtained through cloning, e.g., cloning into plasmids, yeast, or bacterial artificial chromosomes. A polynucleotide can be obtained by reverse transcription of isolated mRNA.

**[0140]** In some embodiments, genomic DNA is analyzed. In some embodiments, the genomic DNA is from a mammal, e.g., a human. The genomic DNA can be obtained from normal somatic tissue, germinal tissue, or diseased tissue (e.g., tumor tissue). In some embodiments, about, more than about, at least about, or less than about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, or 100 genome equivalents are used. A genome equivalent can be the amount of DNA in a single copy of a

genome (e.g., a single diploid cell has 2 genome equivalents of DNA). In some embodiments, about 1 to about 10, about 1 to about 15, about 1 to about 20, about 1 to about 25, about 1 to about 30, about 1 to about 35, about 1 to about 40, about 1 to about 45, about 1 to about 50, about 1 to about 55, about 1 to about 60, about 5 to about 10, about 5 to about 15, about 5 to about 20, about 5 to about 25, about 5 to about 30, about 5 to about 35, about 5 to about 40, about 5 to about 45, about 5 to about 50, about 5 to about 55, about 5 to about 60, about 10 to about 20, about 10 to about 30, about 10 to about 40, about 10 to about 50, or about 10 to about 50 genome equivalents are used. In some embodiments, about, more than about, at least about, or less than about 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10,000, 20,000, 30,000, 40,000, 50,000, 60,000, 70,000, 80,000, 90,000, 100,000, 200,000, 300,000, 400,000, 500,000, 600,000, 700,000, 800,000, 900,000, 1,000,000, 2,000,000, 3,000,000, 4,000,000, 5,000,000, 6,000,000, 7,000,000, 8,000,000, 9,000,000, or 10,000,000 genome equivalents are used. In some embodiments, about 100 to about 1000, about 100 to about 10,000, about 100 to about 100,000, about 100 to about 1,000,000, about 100 to about 10,000,000, about 1000 to about 10,000, about 1000 to about 100,000, about 1000 to about 1,000,000, about 1000 to about 10,000,000, about 10,000 to about 100,000, about 10,000 to about 1,000,000, about 10,000 to about 10,000,000, about 100,000 to about 1,000,000, about 100,000 to about 10,000,000, or about 1,000,000 to about 10,000,000 genome equivalents are used.

**[0141]** In some embodiments, polynucleotides are protected from shearing. Additives that can protect polynucleotides from shearing include, e.g., spermidine, spermine, poly(N-vinylpyrrolidone) 40 (PVP40), or  $\text{Co}(\text{NH}_3)_6\text{Cl}_3$ . In some embodiments, wide pore pipettes are used to avoid shearing of polynucleotides, e.g., when polynucleotides are transferred from one receptacle to another. Methods and compositions for protecting polynucleotides from shearing are described, e.g., in Kovacic et al. (1995) *Nucleic Acids Research* 23: 3999-4000 and Gurrieri S and Bustamante C. (1997) *Biochem J.* 326: 131-138.

**[0142]** The length of polynucleotides, or fragments of polynucleotides, that can be partitioned (e.g., in droplets) as described herein can be about, more than about, at least about, or less than about 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10,000, 20,000, 30,000, 40,000, 50,000, 60,000, 70,000, 80,000, 90,000, 100,000, 200,000, 300,000, 400,000, 500,000, 600,000, 700,000, 800,000, 900,000, 1,000,000, 2,000,000, 3,000,000, 4,000,000, 5,000,000, 6,000,000, 7,000,000, 8,000,000, 9,000,000, 10,000,000, 20,000,000, 30,000,000, 40,000,000, 50,000,000, 60,000,000, 70,000,000, 80,000,000, 90,000,000, 100,000,000, 110,000,000, 120,000,000, 130,000,000, 140,000,000, 150,000,000, 160,000,000, 170,000,000, 180,000,000, 190,000,000, 200,000,000, 210,000,000, 220,000,000, 230,000,000, 240,000,000, or 250,000,000 nucleotides or base pairs in length.

**[0143]** Individual chromosomes can be separated into individual partitions. Human chromosomes that can be partitioned as described herein can include chromosomes 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, X or Y.

**[0144]** In some embodiments, gentle processing steps are used to obtain large polynucleotides from a sample. The gentle processing steps can include, e.g., low speed centrifur-

gation, release of genomic DNA using proteinase K and/or RNase digestion, or dialysis. In some embodiments, steps such as vortexing, high speed centrifugation, or ethanol precipitation are not performed.

**[0145]** Next Generation Sequencing

**[0146]** The methods, compositions, and kits described herein can be used with next generation sequence platforms. For example, adaptors with barcodes can be ligated to polynucleotides, different samples of polynucleotides with different barcodes can be pooled, the pooled polynucleotides can be sequenced using next generation sequencing, and barcodes can be used to determine which sequence reads are generated from polynucleotides in the same partition (e.g., droplet).

**[0147]** In some embodiments, the next generation sequencing technique is 454 sequencing (Roche) (see e.g., Margulies, M et al. (2005) *Nature* 437: 376-380). 454 sequencing can involve two steps. In the first step, DNA can be sheared into fragments of approximately 300-800 base pairs, and the fragments can be blunt ended. Oligonucleotide adaptors can then be ligated to the ends of the fragments. The adaptors can serve as sites for hybridizing primers for amplification and sequencing of the fragments. The fragments can be attached to DNA capture beads, e.g., streptavidin-coated beads using, e.g., Adaptor B, which can contain 5'-biotin tag. The fragments can be attached to DNA capture beads through hybridization. A single fragment can be captured per bead. The fragments attached to the beads can be PCR amplified within droplets of an oil-water emulsion. The result can be multiple copies of clonally amplified DNA fragments on each bead. The emulsion can be broken while the amplified fragments remain bound to their specific beads. In a second step, the beads can be captured in wells (pico-liter sized; PicoTiterPlate (PTP) device). The surface can be designed so that only one bead fits per well. The PTP device can be loaded into an instrument for sequencing. Pyrosequencing can be performed on each DNA fragment in parallel. Addition of one or more nucleotides can generate a light signal that can be recorded by a CCD camera in a sequencing instrument. The signal strength can be proportional to the number of nucleotides incorporated. Pyrosequencing can make use of pyrophosphate (PPi) which can be released upon nucleotide addition. PPi can be converted to ATP by ATP sulfurylase in the presence of adenosine 5' phosphosulfate. Luciferase can use ATP to convert luciferin to oxyluciferin, and this reaction can generate light that is detected and analyzed.

**[0148]** In some embodiments, the next generation sequencing technique is SOLiD technology (Applied Biosystems; Life Technologies). In SOLiD sequencing, genomic DNA can be sheared into fragments, and adaptors can be attached to the 5' and 3' ends of the fragments to generate a fragment library. Alternatively, internal adaptors can be introduced by ligating adaptors to the 5' and 3' ends of the fragments, circularizing the fragments, digesting the circularized fragment to generate an internal adaptor, and attaching adaptors to the 5' and 3' ends of the resulting fragments to generate a mate-paired library. Next, clonal bead populations can be prepared in microreactors containing beads, primers, template, and PCR components. Following PCR, the templates can be denatured and beads can be enriched to separate the beads with extended templates. Templates on the selected beads can be subjected to a 3' modification that permits bonding to a glass slide. A sequencing primer can bind to adaptor sequence. A set of four fluorescently labeled di-base probes can compete for ligation to the sequencing primer. Specificity of the di-

base probe can be achieved by interrogating every first and second base in each ligation reaction. The sequence of a template can be determined by sequential hybridization and ligation of partially random oligonucleotides with a determined base (or pair of bases) that can be identified by a specific fluorophore. After a color is recorded, the ligated oligonucleotide can be cleaved and removed and the process can be then repeated. Following a series of ligation cycles, the extension product can be removed and the template can be reset with a primer complementary to the n-1 position for a second round of ligation cycles. Five rounds of primer reset can be completed for each sequence tag. Through the primer reset process, most of the bases can be interrogated in two independent ligation reactions by two different primers. Up to 99.99% accuracy can be achieved by sequencing with an additional primer using a multi-base encoding scheme.

**[0149]** In some embodiments, the next generation sequencing technique is SOLEXA sequencing (Illumina sequencing). SOLEXA sequencing can be based on the amplification of DNA on a solid surface using fold-back PCR and anchored primers. SOLEXA sequencing can involve a library preparation step. Genomic DNA can be fragmented, and sheared ends can be repaired and adenylated. Adaptors can be added to the 5' and 3' ends of the fragments. The fragments can be size selected and purified. SOLEXA sequence can comprise a cluster generation step. DNA fragments can be attached to the surface of flow cell channels by hybridizing to a lawn of oligonucleotides attached to the surface of the flow cell channel. The fragments can be extended and clonally amplified through bridge amplification to generate unique clusters. The fragments become double stranded, and the double stranded molecules can be denatured. Multiple cycles of the solid-phase amplification followed by denaturation can create several million clusters of approximately 1,000 copies of single-stranded DNA molecules of the same template in each channel of the flow cell. Reverse strands can be cleaved and washed away. Ends can be blocked, and primers can be hybridized to DNA templates. SOLEXA sequencing can comprise a sequencing step. Hundreds of millions of clusters can be sequenced simultaneously. Primers, DNA polymerase and four fluorophore-labeled, reversibly terminating nucleotides can be used to perform sequential sequencing. All four bases can compete with each other for the template. After nucleotide incorporation, a laser is used to excite the fluorophores, and an image is captured and the identity of the first base is recorded. The 3' terminators and fluorophores from each incorporated base are removed and the incorporation, detection and identification steps are repeated. A single base can be read each cycle.

**[0150]** In some embodiments, the next generation sequencing technique comprises real-time (SMRT™) technology by Pacific Biosciences. In SMRT, each of four DNA bases can be attached to one of four different fluorescent dyes. These dyes can be phospholinked. A single DNA polymerase can be immobilized with a single molecule of template single stranded DNA at the bottom of a zero-mode waveguide (ZMW). A ZMW can be a confinement structure which enables observation of incorporation of a single nucleotide by DNA polymerase against the background of fluorescent nucleotides that can rapidly diffuse in an out of the ZMW (in microseconds). It can take several milliseconds to incorporate a nucleotide into a growing strand. During this time, the fluorescent label can be excited and produce a fluorescent signal, and the fluorescent tag can be cleaved off. The ZMW

can be illuminated from below. Attenuated light from an excitation beam can penetrate the lower 20-30 nm of each ZMW. A microscope with a detection limit of 20 zeptoliters ( $10^{-21}$  liters) can be created. The tiny detection volume can provide 1000-fold improvement in the reduction of background noise. Detection of the corresponding fluorescence of the dye can indicate which base was incorporated. The process can be repeated.

**[0151]** In some embodiments, the next generation sequencing is nanopore sequencing (See e.g., Soni GV and Meller A. (2007) *Clin Chem* 53: 1996-2001). A nanopore can be a small hole, of the order of about one nanometer in diameter. Immersion of a nanopore in a conducting fluid and application of a potential across it can result in a slight electrical current due to conduction of ions through the nanopore. The amount of current which flows can be sensitive to the size of the nanopore. As a DNA molecule passes through a nanopore, each nucleotide on the DNA molecule can obstruct the nanopore to a different degree. Thus, the change in the current passing through the nanopore as the DNA molecule passes through the nanopore can represent a reading of the DNA sequence. The nanopore sequencing technology can be from Oxford Nanopore Technologies; e.g., a GridION system. A single nanopore can be inserted in a polymer membrane across the top of a microwell. Each microwell can have an electrode for individual sensing. The microwells can be fabricated into an array chip, with 100,000 or more microwells per chip. An instrument (or node) can be used to analyze the chip. Data can be analyzed in real-time. One or more instruments can be operated at a time. The nanopore can be a protein nanopore, e.g., the protein alpha-hemolysin, a heptameric protein pore. The nanopore can be a solid-state nanopore made, e.g., a nanometer sized hole formed in a synthetic membrane (e.g.,  $\text{SiN}_x$ , or  $\text{SiO}_2$ ). The nanopore can be a hybrid pore (e.g., an integration of a protein pore into a solid-state membrane). The nanopore can be a nanopore with an integrated sensors (e.g., tunneling electrode detectors, capacitive detectors, or graphene based nano-gap or edge state detectors (see e.g., Garaj et al. (2010) *Nature* vol. 67, doi:10.1038/nature09379)). A nanopore can be functionalized for analyzing a specific type of molecule (e.g., DNA, RNA, or protein). Nanopore sequencing can comprise "strand sequencing" in which intact DNA polymers can be passed through a protein nanopore with sequencing in real time as the DNA translocates the pore. An enzyme can separate strands of a double stranded DNA and feed a strand through a nanopore. The DNA can have a hairpin at one end, and the system can read both strands. In some embodiments, nanopore sequencing is "exonuclease sequencing" in which individual nucleotides can be cleaved from a DNA strand by a processive exonuclease, and the nucleotides can be passed through a protein nanopore. The nucleotides can transiently bind to a molecule in the pore (e.g., cyclodextran). A characteristic disruption in current can be used to identify bases. Nanopore sequencing technology from GENIA or NABsys can be used. In GENIA's technology, an engineered protein pore can be embedded in a lipid bilayer membrane, and "Active Control" technology can enable efficient nanopore-membrane assembly and control of DNA movement through the channel. In some embodiments, the next generation sequencing comprises ion semiconductor sequencing (e.g., using technology from Life Technologies (Ion Torrent)). Ion semiconductor sequencing can take advantage of the fact that when a nucleotide is incorporated into a strand of DNA, an ion can be released. To perform ion semi-

conductor sequencing, a high density array of micromachined wells can be formed. Each well can hold a single DNA template. Beneath the well can be an ion sensitive layer, and beneath the ion sensitive layer can be an ion sensor. When a nucleotide is added to a DNA,  $H^+$  is released, which can be measured as a change in pH. The  $H^+$  ion can be converted to voltage and recorded by the semiconductor sensor. An array chip can be sequentially flooded with one nucleotide after another. No scanning, light, or cameras can be required.

**[0152]** In some embodiments, the next generation sequencing is DNA nanoball sequencing (as performed, e.g., by Complete Genomics; see e.g., Drmanac et al. (2010) *Science* 327: 78-81). DNA can be isolated, fragmented, and size selected. For example, DNA can be fragmented (e.g., by sonication) to a mean length of about 500 bp. Adaptors (Ad1) can be attached to the ends of the fragments. The adaptors can be used to hybridize to anchors for sequencing reactions. DNA with adaptors bound to each end can be PCR amplified. The adaptor sequences can be modified so that complementary single strand ends bind to each other forming circular DNA. The DNA can be methylated to protect it from cleavage by a Type IIS restriction enzyme used in a subsequent step. An adaptor (e.g., the right adaptor) can have a restriction recognition site, and the restriction recognition site can remain non-methylated. The non-methylated restriction recognition site in the adaptor can be recognized by a restriction enzyme (e.g., *AclI*), and the DNA can be cleaved by *AclI* 13 bp to the right of the right adaptor to form linear double stranded DNA. A second round of right and left adaptors (Ad2) can be ligated onto either end of the linear DNA, and all DNA with both adaptors bound can be PCR amplified (e.g., by PCR). Ad2 sequences can be modified to allow them to bind each other and form circular DNA. The DNA can be methylated, but a restriction enzyme recognition site can remain non-methylated on the left Ad1 adaptor. A restriction enzyme (e.g., *AclI*) can be applied, and the DNA can be cleaved 13 bp to the left of the Ad1 to form a linear DNA fragment. A third round of right and left adaptor (Ad3) can be ligated to the right and left flank of the linear DNA, and the resulting fragment can be PCR amplified. The adaptors can be modified so that they can bind to each other and form circular DNA. A type III restriction enzyme (e.g., *EcoP15*) can be added; *EcoP15* can cleave the DNA 26 bp to the left of Ad3 and 26 bp to the right of Ad2. This cleavage can remove a large segment of DNA and linearize the DNA once again. A fourth round of right and left adaptors (Ad4) can be ligated to the DNA, the DNA can be amplified (e.g., by PCR), and modified so that they bind each other and form the completed circular DNA template. Rolling circle replication (e.g., using *Phi 29* DNA polymerase) can be used to amplify small fragments of DNA. The four adaptor sequences can contain palindromic sequences that can hybridize and a single strand can fold onto itself to form a DNA nanoball (DNB™) which can be approximately 200-300 nanometers in diameter on average. A DNA nanoball can be attached (e.g., by adsorption) to a microarray (sequencing flowcell). The flow cell can be a silicon wafer coated with silicon dioxide, titanium and hexamethyldisilazane (HMDS) and a photoresist material. Sequencing can be performed by unchained sequencing by ligating fluorescent probes to the DNA. The color of the fluorescence of an interrogated position can be visualized by a high resolution camera. The identity of nucleotide sequences between adaptor sequences can be determined.

**[0153]** In some embodiments, the next generation sequencing technique is Helicos True Single Molecule Sequencing (tSMS) (see e.g., Harris T. D. et al. (2008) *Science* 320:106-109). In the tSMS technique, a DNA sample can be cleaved into strands of approximately 100 to 200 nucleotides, and a polyA sequence can be added to the 3' end of each DNA strand. Each strand can be labeled by the addition of a fluorescently labeled adenosine nucleotide. The DNA strands can then be hybridized to a flow cell, which can contain millions of oligo-T capture sites immobilized to the flow cell surface. The templates can be at a density of about 100 million templates/cm<sup>2</sup>. The flow cell can then be loaded into an instrument, e.g., HELISCOPE™ sequencer, and a laser can illuminate the surface of the flow cell, revealing the position of each template. A CCD camera can map the position of the templates on the flow cell surface. The template fluorescent label can then be cleaved and washed away. The sequencing reaction can begin by introducing a DNA polymerase and a fluorescently labeled nucleotide. The oligo-T nucleic acid can serve as a primer. The DNA polymerase can incorporate the labeled nucleotides to the primer in a template directed manner. The DNA polymerase and unincorporated nucleotides can be removed. The templates that have directed incorporation of the fluorescently labeled nucleotide can be detected by imaging the flow cell surface. After imaging, a cleavage step can remove the fluorescent label, and the process can be repeated with other fluorescently labeled nucleotides until a desired read length is achieved. Sequence information can be collected with each nucleotide addition step. The sequencing can be asynchronous. The sequencing can comprise at least 1 billion bases per day or per hour.

**[0154]** In some embodiments, the sequencing technique can comprise paired-end sequencing in which both the forward and reverse template strand can be sequenced. In some embodiments, the sequencing technique can comprise mate pair library sequencing. In mate pair library sequencing, DNA can be fragments, and 2-5 kb fragments can be end-repaired (e.g., with biotin labeled dNTPs). The DNA fragments can be circularized, and non-circularized DNA can be removed by digestion. Circular DNA can be fragmented and purified (e.g., using the biotin labels). Purified fragments can be end-repaired and ligated to sequencing adaptors.

**[0155]** In some embodiments, a sequence read is about, more than about, less than about, or at least about 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278,

279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 525, 550, 575, 600, 625, 650, 675, 700, 725, 750, 775, 800, 825, 850, 875, 900, 925, 950, 975, 1000, 1100, 1200, 1300, 1400, 1500, 1600, 1700, 1800, 1900, 2000, 2100, 2200, 2300, 2400, 2500, 2600, 2700, 2800, 2900, or 3000 bases. In some embodiments, a sequence read is about 10 to about 50 bases, about 10 to about 100 bases, about 10 to about 200 bases, about 10 to about 300 bases, about 10 to about 400 bases, about 10 to about 500 bases, about 10 to about 600 bases, about 10 to about 700 bases, about 10 to about 800 bases, about 10 to about 900 bases, about 10 to about 1000 bases, about 10 to about 1500 bases, about 10 to about 2000 bases, about 50 to about 100 bases, about 50 to about 150 bases, about 50 to about 200 bases, about 50 to about 500 bases, about 50 to about 1000 bases, about 100 to about 200 bases, about 100 to about 300 bases, about 100 to about 400 bases, about 100 to about 500 bases, about 100 to about 600 bases, about 100 to about 700 bases, about 100 to about 800 bases, about 100 to about 900 bases, or about 100 to about 1000 bases.

**[0156]** In some embodiments, the sequencing depth is about, more than about, at least about, or less than about 1x, 2x, 3x, 4x, 5x, 6x, 7x, 8x, 9x, 10x, 11x, 12x, 13x, 14x, 15x, 16x, 17x, 18x, 19x, 20x, 21x, 22x, 23x, 24x, 25x, 26x, 27x, 28x, 29x, 30x, 31x, 32x, 33x, 34x, 35x, 36x, 37x, 38x, 39x, 40x, 41x, 42x, 43x, 44x, 45x, 46x, 47x, 48x, 49x, 50x, 51x, 52x, 53x, 54x, 55x, 56x, 57x, 58x, 59x, 60x, 61x, 62x, 63x, 64x, 65x, 66x, 67x, 68x, 69x, 70x, 71x, 72x, 73x, 74x, 75x, 76x, 77x, 78x, 79x, 80x, 81x, 82x, 83x, 84x, 85x, 86x, 87x, 88x, 89x, 90x, 91x, 92x, 93x, 94x, 95x, 96x, 97x, 98x, 99x, 100x, 110x, 120x, 130x, 140x, 150x, 160x, 170x, 180x, 190x, 200x, 210x, 220x, 230x, 240x, 250x, 260x, 270x, 280x, 290x, 300x, 350x, 400x, 450x, 500x, 550x, 600x, 650x, 700x, 750x, 800x, 850x, 900x, 950x, 1000x, 2000x, 3000x, 4000x, 5000x, 6000x, 7000x, 8000x, 9000x, 10,000x. In some embodiments, the sequencing depth is about 1x to about 4x, about 1x to about 5x, about 1x to about 8x, about 1x to about 10x, about 2x to about 4x, about 2x to about 8x, about 2x to about 10x, about 5x to about 10x, about 3x to about 6x, about 10x to about 15x, about 10x to about 20x, about 15x to about 20x, about 15x to about 25x, about 15x to about 30x, about 20x to about 30x, about 25x to about 30x, about 25x to about 50x, about 25x to about 75x, about 25x to about 100x, about 50x to about 100x, about 100x to about 200x, about 100x to about 500x, about 100x to about 1000x, about 200x to about 500x, about 500x to about 750x, about 500x to about 1000x, about 750x to about 1000x, about

1000x to about 2000x, about 1000x to about 5000x, about 1000x to about 10,000x, about 2000x to about 5000x, or about 5000x to about 10,000x. Depth of sequencing can be the number of times a sequence (e.g., a genome) is sequenced. In some embodiments, the Lander/Waterman equation is used for computing coverage. The general equation can be:  $C=LN/G$ , where C=coverage; G=haploid genome length; L=read length; and N=number of reads.

**[0157]** Applications

**[0158]** Long Reads, Phasing and De Novo Sequencing

**[0159]** In some embodiments, the methods, compositions, and kits described herein can be used for haplotype phasing. In some embodiments, short read sequencers, such as those made by Illumina and ABI, can be unable to provide phasing information. These sequencers can produce reads of 100-200 bases and as short as 30 bases. 454 sequencing (Roche) can produce sequence reads of about 400 bases. In some embodiments, 400 bases can be too short to yield sufficient phasing information. Sequencing using technology from Pacific Biosciences can produce sequence reads of about 1000 bases. In some embodiments, 1000 bases is too short to provide phasing information.

**[0160]** Short sequence reads can make it challenging to sequence a large genome de-novo. Short sequence reads can make it difficult to determine phasing information for all but a very small fraction of polymorphisms. The partitioning and barcoding schemes described herein can be used to re-construct longer reads using long range assembly and supply phasing information while making use of existing sequencing approaches.

**[0161]** Next generation sequencing platforms can entail a library preparation step. Genomic DNA can be fragmented, optionally sized, and ligated to nucleic acid sequence (e.g., an adaptor) that can provide hybridization sites for a common set of primers. A common set of primers can be used for massive clonal amplification, e.g., in solution or on a solid support. In some embodiments, these clones can then be sequenced because the presence of a massive amount of identical sequence in a tightly confined space can allow for the amplification of a fluorescent (or other) signal emitted by the sequencing reaction.

**[0162]** Tag sequences can be appended to regions that serve as binding sites for primers so that a common barcode can be ligated to every sequence from a particular sample. Libraries from different samples can be mixed and sequenced in a single run. Because every read can contain a barcode, it can be determined which sample produced any given sequence read. This process can be known as sample multiplexing and can allow for much more cost effective pricing per sample for many sequencing applications. In some embodiments, part of every sequence read includes barcode sequence.

**[0163]** In some embodiments, a high molecular weight DNA sample can be partitioned so that a given partition is unlikely to contain two fragments from the same locus but different chromosomes. In some embodiments, high molecular weight DNA can comprise polynucleotides of greater than about 10,000, 100,000, 1,000,000, 10,000,000, 100,000,000, or 200,000,000 bases or base pairs. In some embodiments, polynucleotides are separated such that it is a rare event to have any given region of a genome of both a maternal and paternal polynucleotide in the same partition. In some embodiments, less than 10%, 9%, 8%, 7%, 6%, 5%, 4%, 3%, 2%, 1%, 0.1%, 0.01%, or 0.001% of partitions have two fragments from the same locus but from different chromo-

somes. In some embodiments, a sample is partitioned such that about, or less than about 20%, 19%, 18%, 17%, 16%, 15%, 14%, 13%, 12%, 11%, 10%, 9%, 8%, 7%, 6%, 5%, 4%, 3%, 2%, 1%, 0.1%, 0.01%, or 0.001% of a haploid genome is found per partition (e.g., droplet). In some embodiments, a sample is partitioned such that about 0.1% to about 1%, about 0.5% to about 1%, about 0.25% to about 0.75%, about 1% to about 5%, about 1% to about 2%, about 1% to about 10%, or about 5% to about 10% of a haploid genome is found per partition (e.g., droplet).

**[0164]** Library preparation can be performed within partitions (e.g., droplets) as described herein. Sequence reads that map somewhat close to each other in a genome and are determined to be from the same partition (e.g., in the same droplet) are likely linked to each other and thus reside on the same chromosome. In this fashion individual, short reads can be strung together into longer sequence fragments. See e.g., Example 1.

**[0165] Single Cell Analysis**

**[0166]** In some embodiments, the methods and compositions described herein can be used to analyze cells, e.g., individual cells. For example, individual cells can be separated into unique partitions, uniquely barcoded adaptors can be added to each partition, polynucleotides, or fragments of polynucleotides, within each partition can be barcoded by ligating adaptors to the polynucleotides or fragments of polynucleotides, barcoded polynucleotides from each partition can be pooled, the pooled polynucleotides can be sequenced, and barcodes can be used to determine if sequence reads were generated in the same or different partitions, and thus, in the same or different cells. In some embodiments, the methods and compositions described herein are used for single cell transcriptome sequencing, single cell genomic sequencing, or single cell methylome sequencing.

**[0167]** There are approximately 210 different types of cells in the human body. The individual cells that are partitioned can be any type of cell in the human body. A cell can be, e.g., a hormone secreting cell, an exocrine secretory epithelial cell, a keratinizing epithelial cell, a wet stratified barrier epithelial cell, a sensory transducer cell, an autonomic neuron cell, a sense organ or peripheral neuron supporting cell, a central nervous system neuron or glial cell, a lens cell, a metabolism or storage cell, a kidney cell, an extracellular matrix cell, a contractile cell, a blood or immune system cell, a pigment cell, a germ cell, a nurse cell, or an interstitial cell. The blood or immune system cell can be, e.g., erythrocyte (red blood cell), megakaryocyte (platelet precursor), monocyte, connective tissue macrophage, epidermal Langerhans cell, osteoclast, dendritic cell, microglial cell, neutrophil granulocyte, eosinophil granulocyte, basophil granulocyte, mast cell, Helper T cell, Suppressor T cell, Cytotoxic T cell, Natural Killer T cell, B cell, Natural killer cell, reticulocyte, or stem cell.

**[0168]** Individual cells can be from other types of samples described herein.

**[0169]** In some embodiments, the individual cell is from an environmental sample. In some embodiments, an environmental sample is separated into a plurality of partitions. The environmental sample can be, e.g., air, water, agricultural, or soil. The environmental sample can be, e.g., from a creek, river, pond, lake, lagoon, run, delta, marsh, salt marsh, swamp, mangrove swamp, mill pond, moat, sea, barachois, basin, bayou, beck, boil, canal, cove, estuary, gulf, harbor, inlet, ocean, bay, sewage treatment facility, slough, sound,

spring, stream, tide pool, wash, wetland, Superfund site, coal mine, farm, field, desert, glacier, mountain, or mere. In some embodiments, a sample is from a pool (e.g., swimming pool), gymnasium, school, workplace, office, lobby, elevator, restroom, hospital, medical office, ventilation shaft, or restaurant. In some embodiments, an environmental sample can be from a surface, e.g., floor, table, skin, keyboard, computer, laptop, crime scene evidence (e.g., a weapon, e.g., gun or knife), or doorknob. In some embodiments, the sample is from a bioterrorist attack. In some embodiments, the sample comprises bacteria and/or viruses. In some embodiments, the sample comprises about, at least about, more than about, or less than about 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10,000, 20,000, 30,000, 40,000, 50,000, 60,000, 70,000, 80,000, 90,000, or 100,000 different species and/or types of viruses. In some embodiments, and environmental sample comprises about 10 to about 100, about 10 to about 1000, about 100 to about 1000, about 100 to about 10,000, about 1000 to 10,000, about 10,000 to about 50,000, or about 10,000 to about 100,000 different species and/or types of viruses.

**[0170] Single Cell Transcriptome Sequencing**

**[0171]** In one aspect, single cells can be captured within separate partitions (e.g., droplets), and the single cells can be lysed. Messenger RNA from the individual cells in each partition (e.g., droplet) can be reverse transcribed with partition-specific barcoded primers. In some embodiments, the appropriate reagents (e.g., reverse transcriptase, nucleotides) can be sequestered in a partition (e.g., droplet) that is inside a larger droplet. The inner droplet can be burst (e.g., by heating) when desired to allow the reverse transcriptase to contact the messenger RNA. A reverse transcription (RT) reaction can be followed by library prep, which can incorporate unique barcodes.

**[0172]** Calculations for the number of droplets and barcodes to be used in single cell transcriptome sequencing can be similar to those described in Example 1 for analyzing phasing. For example, for analyzing 2,000 cells, sufficient partitioning can be performed to capture each cell in a separate droplet. For example, the 2,000 cells can be partitioned among 20,000 partitions (e.g., droplets). Steps can be taken to ensure that each of the partitions (e.g., droplets) with cells receives a unique barcode (e.g., on adaptors). This goal can be accomplished, e.g., by using 10,000 different barcodes.

**[0173]** After partitioning, lysing, barcoding, and sequencing, sequence read data can be analyzed to determine which transcripts came from the same cell. In this way, the massive capacity of next generation sequencing can be applied to large collections of cells while preserving single cell resolution.

**[0174] Single Cell Genomic Sequencing**

**[0175]** In some instances, individual cells can be captured in separate partitions (e.g., droplets), and genomic DNA from a partition with a single cell can be uniquely barcoded (e.g., using adaptors). Barcoded genomic DNA from different partitions can be pooled and sequenced, and the barcodes can be used to determine which sequence reads are from the same cell. In some embodiments, genomic DNA is fragmented in a partition. In some embodiments, genomic DNA is amplified

before and/or after adding adaptors with barcodes. In some embodiments, genomic DNA is not amplified before or after adaptors are ligated to the genomic DNA.

**[0176]** In some embodiments, the sequence coverage per cell can be shallow (e.g., few reads per locus). In some embodiments, single cell genomic DNA sequencing can be used to determine copy number variation (CNV).

**[0177]** In some embodiments, MDA is performed within a droplet on a cell's genome prior to fragmentation and adaptor ligation. In some embodiments, performing MDA can provide more genetic material from a cell to sequence. In some embodiments, MDA may introduce bias. In some embodiments, amplification may result in loss of some copy number variation (CNV) information. In some embodiments, MDA is not performed within a droplet on a cell's genome prior to fragmentation and adaptor ligation.

#### **[0178] Single Cell Methylome Sequencing**

**[0179]** In some embodiments, methods and compositions described herein can be used for analyzing genomic methylation. For example, methods described herein can be used for single cell methylome sequencing. In some embodiments, individual cells are partitioned, e.g., into droplets. The partitions can be comprised of methyl-sensitive enzymes (e.g., endonucleases). In some embodiments, the methyl-sensitive enzymes digest methylated sites. Each of the partitions can comprise uniquely barcoded adaptors. For example, partitions (e.g., droplets) comprising methyl-sensitive enzymes are merged with partitions comprising sample polynucleotides. The adaptors can be ligated to polynucleotides in the partition before or after digestion with the methyl-sensitive enzyme. Barcode tagged polynucleotides can be pooled, and the polynucleotides can be sequenced. Sequence reads from the same partition can be determined. Absence of sequence reads can indicate digestion of polynucleotides in a partition. In some embodiments, the methyl-sensitive enzymes do not digest methylated DNA, but digest unmethylated DNA.

#### **[0180] Genomic Methylation**

**[0181]** In some embodiments, methods and compositions described herein can be used for genomic methylation analysis. In some embodiments, polynucleotides can be treated with bisulfite. Bisulfite can convert unmethylated cytosines to uracil. Bisulfite does not convert methylated cytosines to uracil. Treated and untreated polynucleotides can be partitioned into a plurality of partitions (e.g., droplets). Polynucleotides can be fragmented in the partitions. Uniquely barcoded adaptors can be provided to each partition and ligated to bisulfite treated polynucleotides. The tagged polynucleotides can be pooled and sequenced to determine the methylation status of nucleic acids from the same and different partitions.

#### **[0182] Exosome Sequencing**

**[0183]** Exosomes are generally organelles such as small extracellular vesicles that can contain RNA. Exosomes can contain mRNA and/or miRNA. In some embodiments, individual exosomes are partitioned into separate partitions (e.g., droplets). Exosomes can be partitioned such that on average, each partition comprises less than about five, four, three, two, or 1 exosomes. Reverse transcription can be used to convert RNA in the exosome into cDNA. Uniquely barcoded adaptors can be added to polynucleotides from a partitioned exosome. Polynucleotides from the partitions can be pooled, the pooled polynucleotides can be sequenced, and the barcodes can be used to determine which sequence reads were derived from

the same exosome. Other types of organelles that can be analyzed can include mitochondria (e.g., mitochondrial DNA can be analyzed).

#### **[0184] Metagenomics Sequencing**

**[0185]** In another aspect, the methods and compositions described herein can be used for metagenomic analysis. Metagenomics can be the study of genetic material in an environmental sample. In some embodiments, individual viruses and/or bacteria in a sample, e.g., an environmental sample, can be partitioned into a plurality of partitions, adaptors with unique barcodes can be added to each partition, and individual organisms or viruses can have their genomes and/or transcriptomes sequenced. Sequence reads with the same barcode can be assembled to determine the sequence of genomes or transcriptomes of the organisms and/or viruses.

#### **[0186] Microfluidics**

**[0187]** In another aspect, a microfluidics device can be devised that can partition a sample comprising cells so that every cell ends up in a unique partition (e.g., chamber) with its own set of barcodes. The contents of each chamber can then be processed separately to dilute and further partition (e.g., through an emulsion) in order to enable whole genome or transcriptome amplification separately for each cell. Whole genome amplification or other amplification schemes can benefit from partitioning because of a reduction in competition between different parts of the genome or transcriptome.

#### **[0188] Slugs**

**[0189]** In another aspect, slugs can be made to capture individual cells and supply them with their own barcodes (e.g., by ligating adaptors with unique barcodes). Slugs can be serial slugs of reagent that completely fill the diameter of a flow tube. Those slugs can be broken into many (e.g., thousands or more) smaller droplets in order to perform unbiased whole genome/transcriptome amplification in droplets. In some embodiments, a slug can be broken down into about, at least about, more than about, or less than about 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, or 10,000 droplets. In some embodiments, a slug can be broken down into about 100 to about 500, about 100 to about 1000, about 500 to about 1000, about 1000 to about 1500, about 1000 to about 2000, about 1000 to about 5000, about 1000 to about 10,000, or about 5,000 to about 10,000 droplets. In some embodiments, whole genome amplification can work better in droplets than in bulk. The droplets from all the slugs can be mixed together because they are already furnished with adaptors with unique barcodes. Sequencing information can be used to determine which reads came from which slug.

#### **[0190] Protein Expression and Nucleic Acid Information**

**[0191]** In another embodiment, methods described herein can be used to capture cells with specific cell surface markers and analyze polynucleotides (e.g., DNA or RNA) from the cells. In some embodiments, antibodies can be linked to beads coated with short DNA fragments with a unique barcode. Each antibody can be associated with its own unique sequence. The antibodies could also be linked to droplets containing DNA fragments—which can be burst as appropriate. Cells can be pre-coated with these antibodies, then captured in larger droplets along with droplet/cell-specific barcode adaptors. Library prep can ensue as described herein, contents of the droplets can be sequenced, and it can be inferred which reads came from which cell by the barcodes. Thus, this technique allows one to, in addition to sequencing a cell's genome or transcriptome, obtain information about



their proteins. In some embodiments, some of the same information can be captured via FACS.

**[0192]** An antibody can include polyclonal and monoclonal antibodies, as well as antigen binding fragments of such antibodies. An antibody, or antigen binding fragment of such an antibody, can be characterized by having specific binding activity for a polypeptide or a peptide portion thereof of at least about  $1 \times 10^5 \text{ M}^{-1}$ . Fab, F(ab').sub.2, Fd, Fv, single chain Fv (scFv) fragments of an antibody and the like, which retain specific binding activity for a polypeptide, can be used. Specific binding activity of an antibody for a polypeptide can be readily determined by one skilled in the art, for example, by comparing the binding activity of an antibody to a particular polypeptide versus a control polypeptide that is not the particular polypeptide. Methods of preparing polyclonal or monoclonal antibodies are well known to those skilled in the art (see, for example, Harlow and Lane, *Antibodies: A Laboratory Manual*, Cold Spring Harbor Laboratory Press (1988)).

**[0193]** An antibody can include naturally occurring antibodies as well as non-naturally occurring antibodies, including, for example, single chain antibodies, chimeric, bifunctional and humanized antibodies, as well as antigen-binding fragments thereof. Such non-naturally occurring antibodies can be constructed using solid phase peptide synthesis, can be produced recombinantly or can be obtained, for example, by screening combinatorial libraries consisting of variable heavy chains and variable light chains as described by Huse et al. (*Science* 246:1275-1281 (1989)). These and other methods of making functional antibodies are well known to those skilled in the art (Winter and Harris, *Immunol. Today* 14:243-246 (1993); Ward et al., *Nature* 341:544-546 (1989); Harlow and Lane, *supra*, (1988); Hilyard et al., *Protein Engineering: A practical approach* (IRL Press 1992); Borrabeck, *Antibody Engineering*, 2d ed. (Oxford University Press 1995)).

**[0194]** Samples

**[0195]** Samples to be analyzed using the methods, compositions, and kits provided herein can be derived from a non-cellular entity comprising nucleic acid (e.g., a virus) or from a cell-based organism (e.g., member of archaea, bacteria, or eukarya domains). A sample can be obtained in some cases from a hospital, laboratory, clinical or medical laboratory. The sample can comprise nucleic acid, e.g., RNA or DNA. The sample can comprise cell-free nucleic acid. In some cases, the sample is obtained from a swab of a surface, such as a door or bench top.

**[0196]** The sample can from a subject, e.g., a plant, fungi, eubacteria, archaeobacteria, protest, or animal. The subject may be an organism, either a single-celled or multi-cellular organism. The subject may be cultured cells, which may be primary cells or cells from an established cell line, among others. The sample may be isolated initially from a multi-cellular organism in any suitable form. The animal can be a fish, e.g., a zebrafish. The animal can be a bird, e.g., a chicken. The animal can be a mammal. The mammal can be, e.g., a dog, cat, horse, cow, mouse, rat, or pig. The mammal can be a primate, e.g., a human, chimpanzee, orangutan, or gorilla. The human can be a male or female. The sample can be from a human embryo or human fetus. In some embodiments, the human can be an infant, child, teenager, adult, or elderly person. The female can be pregnant, can be suspected of being pregnant, or planning to become pregnant. In some embodiments, the sample is from a plant. In some embodiments, the samples comprises one or more viruses.

**[0197]** The sample can be from a subject (e.g., human subject) who is healthy. In some embodiments, the sample is taken from a subject (e.g., an expectant mother) at least 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, or 26 weeks of gestation. In some embodiments, the subject is affected by a genetic disease, a carrier for a genetic disease or at risk for developing or passing down a genetic disease, where a genetic disease is any disease that can be linked to a genetic variation such as mutations, insertions, additions, deletions, translocation, point mutation, trinucleotide repeat disorders and/or single nucleotide polymorphisms (SNPs).

**[0198]** The sample can be from a subject who has a specific disease, disorder, or condition, or is suspected of having (or at risk of having) a specific disease, disorder or condition. For example, the sample can be from a cancer patient, a patient suspected of having cancer, or a patient at risk of having cancer. The cancer can be, e.g., acute lymphoblastic leukemia (ALL), acute myeloid leukemia (AML), adrenocortical carcinoma, Kaposi Sarcoma, anal cancer, basal cell carcinoma, bile duct cancer, bladder cancer, bone cancer, osteosarcoma, malignant fibrous histiocytoma, brain stem glioma, brain cancer, craniopharyngioma, ependymoblastoma, ependymoma, medulloblastoma, medulloepithelioma, pineal parenchymal tumor, breast cancer, bronchial tumor, Burkitt lymphoma, Non-Hodgkin lymphoma, carcinoid tumor, cervical cancer, chordoma, chronic lymphocytic leukemia (CLL), chronic myelogenous leukemia (CML), colon cancer, colorectal cancer, cutaneous T-cell lymphoma, ductal carcinoma in situ, endometrial cancer, esophageal cancer, Ewing Sarcoma, eye cancer, intraocular melanoma, retinoblastoma, fibrous histiocytoma, gallbladder cancer, gastric cancer, glioma, hairy cell leukemia, head and neck cancer, heart cancer, hepatocellular (liver) cancer, Hodgkin lymphoma, hypopharyngeal cancer, kidney cancer, laryngeal cancer, lip cancer, oral cavity cancer, lung cancer, non-small cell carcinoma, small cell carcinoma, melanoma, mouth cancer, myelodysplastic syndromes, multiple myeloma, medulloblastoma, nasal cavity cancer, paranasal sinus cancer, neuroblastoma, nasopharyngeal cancer, oral cancer, oropharyngeal cancer, osteosarcoma, ovarian cancer, pancreatic cancer, papillomatosis, paraganglioma, parathyroid cancer, penile cancer, pharyngeal cancer, pituitary tumor, plasma cell neoplasm, prostate cancer, rectal cancer, renal cell cancer, rhabdomyosarcoma, salivary gland cancer, Sezary syndrome, skin cancer, nonmelanoma, small intestine cancer, soft tissue sarcoma, squamous cell carcinoma, testicular cancer, throat cancer, thymoma, thyroid cancer, urethral cancer, uterine cancer, uterine sarcoma, vaginal cancer, vulvar cancer, Waldenstrom Macroglobulinemia, or Wilms Tumor. The sample can be from the cancer and/or normal tissue from the cancer patient.

**[0199]** The sample can be from a subject who is known to have a genetic disease, disorder or condition. In some cases, the subject is known to be wild-type or mutant for a gene, or portion of a gene, e.g., CFTR, Factor VIII (F8 gene), beta globin, hemachromatosis, G6PD, neurofibromatosis, GAPDH, beta amyloid, or pyruvate kinase gene. In some cases, the status of the subject is either known or not known, and the subject is tested for the presence of a mutation or genetic variation of a gene, e.g., CFTR, Factor VIII (F8 gene), beta globin, hemachromatosis, G6PD, neurofibromatosis, GAPDH, beta amyloid, or pyruvate kinase gene.

**[0200]** In other embodiments, the sample is taken from a female patient of child-bearing age and, in some cases, the



female patient is not pregnant or of unknown pregnancy status. In still other cases, the subject is a male patient, a male expectant father, or a male patient at risk of, diagnosed with, or having a specific genetic abnormality. In some cases, the female patient is known to be affected by, or is a carrier of, a genetic disease or genetic variation, or is at risk of, diagnosed with, or has a specific genetic abnormality. In some cases, the status of the female patient with respect to a genetic disease or genetic variation may not be known. In further embodiments, the sample is taken from any child or adult patient of known or unknown status with respect to copy number variation of a genetic sequence. In some cases, the child or adult patient is known to be affected by, or is a carrier of, a genetic disease or genetic variation.

**[0201]** The sample can be aqueous humour, vitreous humour, bile, whole blood, blood serum, blood plasma, breast milk, cerebrospinal fluid, cerumen, enolymph, perilymph, gastric juice, mucus, peritoneal fluid, saliva, sebum, semen, sweat, perspiration, tears, vaginal secretion, vomit, feces, or urine. The sample can be obtained from a hospital, laboratory, clinical or medical laboratory. The sample can be taken from a subject. The sample can comprise nucleic acid. The nucleic acid can be, e.g., mitochondrial DNA, genomic DNA, mRNA, siRNA, miRNA, cRNA, single-stranded DNA, double-stranded DNA, single-stranded RNA, double-stranded RNA, tRNA, rRNA, or cDNA. The sample can comprise cell-free nucleic acid. The sample can be a cell line, genomic DNA, cell-free plasma, formalin fixed paraffin embedded (FFPE) sample, or flash frozen sample. A formalin fixed paraffin embedded sample can be deparaffinized before nucleic acid is extracted. The sample can be from an organ, e.g., heart, skin, liver, lung, breast, stomach, pancreas, bladder, colon, gall bladder, brain, etc.

**[0202]** In some embodiments, the sample is an environmental sample, e.g., air, water, agricultural, or soil.

**[0203]** When the nucleic acid is RNA, the source of the RNA can be any source described herein. For example, the RNA can be a cell-free mRNA, can be from a tissue biopsy, core biopsy, fine needle aspirate, flash frozen, or formalin-fixed paraffin embedded (FFPE) sample. The FFPE sample can be deparaffinized before the RNA is extracted. The extracted RNA can be heated to about, more than about, less than about, or at least about 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, or 99° C. before analysis. The extracted RNA can be heated to any of these temperatures for about, or at least about, 15 min, 30 min, 45 min, 60 min, 1.5 hr, 2 hr, 2.5 hr, 3 hr, 3.5 hr, 4 hr, 4.5 hr, 5 hr, 5.5 hr, 6 hr, 6.5 hr, 7 hr, 7.5 hr, 8 hr, 8.5 hr, 9 hr, 9.5 hr, or 10 hr.

**[0204]** RNA can be used for a variety of downstream applications. For example, the RNA can be converted to cDNA with a reverse transcriptase and the cDNA can optionally be subject to PCR, e.g., real-time PCR. The RNA or cDNA can be used in an isothermal amplification reaction, e.g., an isothermal linear amplification reaction. The RNA, resulting cDNA, or molecules amplified therefrom can be used in a microarray experiment, gene expression experiment, Northern analysis, Southern analysis, sequencing reaction, next generation sequencing reaction, etc. Specific RNA sequences can be analyzed, or RNA sequences can be globally analyzed.

**[0205]** Nucleic acids can be extracted from a sample by means available to one of ordinary skill in the art.

**[0206]** The sample may be processed to render it competent for amplification. Exemplary sample processing can include lysing cells of the sample to release nucleic acid, purifying the sample (e.g., to isolate nucleic acid from other sample components, which may inhibit amplification), diluting/concentrating the sample, and/or combining the sample with reagents for amplification, such as a DNA/RNA polymerase (e.g., a heat-stable DNA polymerase for PCR amplification), dNTPs (e.g., dATP, dCTP, dGTP, and dTTP (and/or dUTP)), a primer set for each allele sequence or polymorphic locus to be amplified, probes (e.g., fluorescent probes, such as TAQMAN probes or molecular beacon probes, among others) capable of hybridizing specifically to each allele sequence to be amplified,  $Mg^{2+}$ , DMSO, BSA, a buffer, or any combination thereof, among others. In some examples, the sample may be combined with a restriction enzyme, uracil-DNA glycosylase (UNG), reverse transcriptase, or any other enzyme of nucleic acid processing.

**[0207]** Computers

**[0208]** A computer can be used to store and process the data. A computer-executable logic can be employed to perform such functions grouping sequence reads by barcode sequence. A computer can be useful for displaying, storing, retrieving, or calculating diagnostic results from the molecular profiling; displaying, storing, retrieving, or calculating raw data from genomic or nucleic acid expression analysis; or displaying, storing, retrieving, or calculating any sample or patient information useful in the methods described herein. Provided herein are systems comprising computer readable instructions for performing methods described herein. Provided herein are computer readable medium comprising instructions which, when executed by a computer, cause the computer to perform methods described herein.

**[0209]** Kits

**[0210]** Provided herein are kits for performing methods described herein. The kits can comprise one or more restriction enzymes, endonucleases, exonucleases, ligases, polymerases, RNA polymerases, DNA polymerases, reverse transcriptases, topoisomerases, kinases, phosphatases, buffers, salts, metal ions, reducing agents, BSA, spermine, spermidine, glycerol, oligonucleotides, primers, probes, or labels (e.g., fluorescent labels). The kits can comprise one or more sets of instructions.

**[0211]** Multiplexing to Align the Dynamic Range of Targets Whose Concentrations are Different and to Smooth Out Biological Variation of Reference Genes

**[0212]** Also provided herein are methods for estimating copy number variation (CNV). Copy number variation of one or more target sequences can play a role in a number of diseases and disorders. One method to analyze copy number variation of a target sequence is through a digital analysis, such as digital PCR, or droplet digital PCR. However, digital analysis of copy number of a target sequence can underestimate the number of copies of a target nucleic acid sequence in a sample if multiple copies of the target nucleic acid sequence are on the same polynucleotide in a sample. For example, in a digital PCR assay that has multiple compartments (e.g., partitions, spatially isolated regions), nucleic acids in a sample can be partitioned such that each compartment receives on average about 0, 1, 2, or several target polynucleotides. Each partition can have, on average, less than 5, 4, 3, 2, or 1 copies of a target nucleic acid per partition (e.g., droplet). In some cases, at least 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 125, 150, 175, or 200

partitions (e.g., droplets) have zero copies of a target nucleic acid. The number of compartments that contain a polynucleotide can be enumerated. However, if two copies of a target nucleic acid sequence are on a single polynucleotide, a compartment containing that polynucleotide can be counted as having only one target sequence.

**[0213]** Methods of analyzing CNVs are disclosed, e.g., in U.S. patent application Ser. No. 13/385,277, filed Feb. 9, 2012. For example, methods can be used for physically separating target nucleic acids sequences. Often, the methods can avoid underestimating copy numbers of a target sequence due to the presence of multiple copies of the target sequence on a single polynucleotide. In some embodiments, a first sample of polynucleotides is obtained; the first sample can be, e.g., a genomic DNA sample. The target nucleic acid sequences in the first sample can be physically separated (e.g., by contacting the first sample with one or more restriction enzymes). The first sample can be separated into a plurality of partitions. The number of partitions with the target sequence can be enumerated. The copy number of the target can then be estimated.

**[0214]** The target nucleic acids can be identical; or, in other cases, the target nucleic acids can be different. In some cases, the target nucleic acids are located within the same gene. In some cases, the target nucleic acids are each located in a different copy (identical or near identical copy) of a gene. In still other cases, the target sequences are located within introns, or in a region between genes. Sometimes, one target sequence is located in a gene; and the second target sequence is located outside of the gene. In some cases, a target sequence is located within an exon.

**[0215]** Different targets within a sample may often be present at different copy numbers. In such cases, the target that is present at a lower copy number level, may be probed with multiple probes, each recognizing a different locus or region with the target sequence. For example, target A may be present at copy number 3, while target B is present at copy number 1. In such cases, target B may be probed with 3 primer/probe pairs to increase the number of B-positive droplets, or to increase the signal from droplets that comprise target B. The probes may be directed to different regions within target B. Often, the probes that target target B are labeled with the same label; but, in some cases, different labels may be used. Thus, such methods enable alignment of the dynamic range of targets with different copy numbers. Target B can be a different target or may be a reference sample, as described further herein. Therefore, such methods can also enable alignment of the dynamic range of targets with reference samples.

**[0216]** In some cases, a genome comprises one target sequence. In some cases, a genome comprises two or more target sequences. When a genome comprises two or more target sequences, the target sequences can be about, or more than about 50%, 55%, 60%, 65%, 70%, 75%, 80%, 85%, 90%, 95%, or 100% identical.

**[0217]** Separating two target sequences can comprise separating the target sequences by cleaving a specific site on the nucleic acid sequence. In some cases, the separating target nucleic acid sequences can comprise contacting the first sample with one or more restriction enzymes. Separating the target nucleic acid sequences can comprise digesting a polynucleotide at a site located between the target nucleic acid sequences. In some cases, the target nucleic acid sequences are each located within a gene. In some cases, the site that is

targeted for digestion is located between the two genes. In some cases, the site selected for digestion is located in a gene; and, in some cases, the gene is the same gene as the gene which contains the target sequences. In other cases, the site selected for digestion is located in a different gene from that of the target sequence. In some cases, a target sequence and the site targeted for digestion are located in the same gene; and the target sequence is located upstream of the site targeted for digestion. In other cases, a target sequence and the site targeted for digestion are located in the same gene; but the target sequence is located downstream of the site targeted for digestion. In some cases, target nucleic acids can be separated by treatment of a nucleic acid sample with one or more restriction enzymes. In some cases, target nucleic acids can be separated by shearing. In some cases, target nucleic acids can be separated by sonication.

**[0218]** Following the separation step (e.g., digesting with one or more restriction enzymes), the sample can be partitioned into multiple partitions. Each of the plurality of partitions can comprise about 0, 1, 2 or several target polynucleotides. In some cases, each partition can have, on average, less than 5, 4, 3, 2, or 1 copies of a target nucleic acid per partition (e.g., droplet). In some cases, at least 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 125, 150, 175, or 200 droplets have zero copies of a target nucleic acid.

**[0219]** Often, target nucleic acid is amplified in the partitions. In some cases, the amplification comprises use of one or more TaqMan probes.

**[0220]** In another embodiment, the method further comprises the step of enumerating the number of partitions comprising a reference nucleic acid sequence. A reference nucleic acid sequence can be known to be present in a certain number of copies per genome and can be used to estimate the number of genome copies of a target nucleic acid sequence in a sample. Estimating the copy number can comprise comparing the number of partitions comprising the target sequence to the number of partitions comprising the reference nucleic acid sequence. In another instance, a CNV estimate is determined by a ratio of the concentration of target nucleic acid sequence to a reference sequence.

**[0221]** In another embodiment, the method further comprises the step of analyzing a second sample, wherein the second sample and the first sample are derived from the same sample (e.g., a nucleic acid sample is split to the first sample and the second sample). The method can further comprise not contacting the second sample with one or more restriction enzymes. In some cases, the method further comprises separating the second sample into a plurality of partitions. The method can further comprise enumerating the number of partitions of the second sample that comprise the target sequence. In another embodiment, the method further comprises enumerating the number of partitions of the second sample that comprise a reference sequence. In another embodiment, the method comprises estimating the copy number of the target sequence in the second sample. In another embodiment, estimating the copy number of the target sequence in the second sample comprises comparing the number of partitions from the second sample with the target sequence and the number of partitions from the second sample with the reference sequence.

**[0222]** The copy number of the target sequence from the first sample and the copy number of the target sequence in the second sample can be compared to determine whether the copy number of the target sequence in the second sample was

underestimated. The degree to which the copy number was underestimated may be indicative of whether interrogated copies were all on one chromosome or if at least one copy was on one homologous chromosome and at least one copy was on the other homologous chromosome. Values closer to one per diploid genome may indicate the first case, while values closer to two may indicate the second case.

**[0223]** Additional methods of determining copy number differences by amplification are described, e.g., in U.S. Patent Application Publication No. 20100203538. Methods for determining copy number variation are described in U.S. Pat. No. 6,180,349 and Taylor et al. (2008) *PLoS One* 3(9): e3179.

**[0224]** Copy number variations described herein can involve the loss or gain of nucleic acid sequence. Copy number variations can be inherited or can be caused by a de novo mutation. A CNV can be in one or more different classes. See, e.g. Redon et al. (2006) Global variation in copy number in the human genome. *Nature* 444 pp. 444-454. A CNV can result from a simple de novo deletion, from a simple de novo duplication, or from both a deletion and duplication. A CNV can result from combinations of multi-allelic variants. A CNV can be a complex CNV with de novo gain. A CNV can include about, or more than about 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, or 10 contiguous genes. A CNV can include about 1 to about 10, about 1 to about 5, about 1 to about 4, about 1 to about 3, about 1 to about 2, about 0 to about 10, about 0 to about 5, or about 0 to about 2 contiguous genes. A copy number variation can involve a gain or a loss of about, or more than about, 100, 500, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10,000, 20,000, 30,000, 40,000, 50,000, 60,000, 70,000, 80,000, 90,000, 100,000, 200,000, 500,000, 750,000, 1 million, 5 million, or 10 million base pairs. In some cases, a copy number variation can involve the gain or loss of about 1,000 to about 10,000,000, about 10,000 to about 10,000,000, about 100,000 to about 10,000,000, about 1,000 to about 100,000, or about 1,000 to about 10,000 base-pairs of nucleic acid sequence. A copy number variation can be a deletion, insertion, or duplication of a nucleic acid sequence. In some cases, a copy number variation can be a tandem duplication.

**[0225]** In another embodiment, CNV haplotypes can be estimated from fluorescent signals generated by real-time PCR or ddPCR of partitioned samples. Before the late stages of a real-time PCR or ddPCR experiment, when reagents can become limiting, a partition with a higher copy number of a target sequence can have a higher signal than a partition with a lower copy number of the target sequence. In one embodiment, a sample (e.g., a subsample of a sample used in a linkage experiment) can be partitioned, and PCR can be performed on the partitions (e.g., droplets). The mean fluorescence intensity of partitions can be determined as they undergo exponential amplification for a target and/or reference nucleic acid sequence. The mean intensity can correspond to the number of starting copies of the target. If multiple targets are linked along a single polynucleotide strand, the intensity in the partition (e.g., droplet) that captures this strand may be higher than that of a partition (e.g., droplet) that captures a strand with only a single copy of the target. Excess presence of positive droplets with higher mean amplitudes can suggest the presence of a haplotype with multiple CNV copies. Conversely, presence of positive droplets with only low mean amplitudes can suggest that only haplotypes with single CNV copies are present in the sample. In another embodiment, the number of cycles used to estimate CNV can be optimized based on the size of the partitions and the

amount of reagent in the partitions. For example, smaller partitions with lower amounts of reagent may require fewer amplification cycles than larger partitions that would be expected to have higher amounts of reagent.

**[0226]** The method can be useful because it can be used to analyze even target copies that are near each other on the polynucleotide, e.g., less than about 10, 9, 8, 7, 6, 5, 4, 5, 2, 1, 0.7, 0.5, 0.3, 0.2, 0.1, 0.05, or 0.01 megabases apart; or that are very near each other on the polynucleotide, e.g., less than about 10, 9, 8, 7, 6, 5, 4, 3, 2, or 1 kilobase apart. In some cases, the method is useful for analyzing target copies that are very close to each other on the polynucleotide, e.g., within about 1, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 200, 300, 400, 500, 600, 700, 800, 900, or 950 base pairs (bp's) apart. In some cases, the method is useful for analyzing target copies that are separated by zero (0) base pairs. In some cases, the method can be applied to identical, near identical, and completely different targets.

**[0227]** In some embodiments, the copy number of a target in a genome is about, more than about, less than about, or at least about 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 500, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10,000, 20,000, 30,000, 40,000, 50,000, 60,000, 70,000, 80,000, 90,000, or 100,000 copies per haploid or diploid genome. In some embodiments, the copy number of a target is about 2 to about 5, about 2 to about 10, about 2 to about 20, about 2 to about 30, about 2 to about 40, about 2 to about 50, about 2 to about 100, about 5 to about 10, about 5 to about 25, about 5 to about 50, about 5 to about 100, about 10 to about 20, about 10 to about 50, about 10 to about 100, about 25 to about 50, about 25 to about 75, about 25 to about 100, about 100 to about 200, about 100 to about 500, about 100 to about 1000, about 500 to about 1000, about 1000 to about 5000, about 1000 to about 10,000, about 10,000 to about 20,000, about 10,000 to about 50,000, about 10,000 to about 100,000, or about 50,000 to about 100,000 per haploid or diploid genome.

**[0228]** In some embodiments, CNVs can be analyzed by measuring amounts of a target and a reference in a single reaction using probes with one fluorescence dye for the target and another for the reference. In some embodiments, e.g., when the target copy number is high, the concentration (or amount) of the target can be higher than the concentration (or amount) of the reference. In that case, it can be challenging to measure both the target and the reference in a single digital reaction (e.g., digital PCR), because the dynamic range of digital PCR can be limited. For example, a target may be present at 10,000 copies in a genome, but a reference may be present at only two copies per genome.

**[0229]** In some embodiments, several different targets for the reference can be multiplexed with each being detectable using probes with the same fluorescent dye. (See e.g., FIG. 2) Often, these different reference targets represent different regions or loci within the same reference polynucleotide (e.g., gene); although, in some cases, different reference polynucleotides (e.g., genes) can be used. Use of multiple references can boost the counts of the reference and bring them closer to the counts of the target. In some embodiments, about, more than about, at least about, or less than about 2, 3, 4, 5, 6, 7, 8,

9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 500, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10,000, 20,000, 30,000, 40,000, 50,000, 60,000, 70,000, 80,000, 90,000, or 100,000 different references are used. In some embodiments, about 2 to about 5, about 2 to about 10, about 2 to about 20, about 2 to about 30, about 2 to about 40, about 2 to about 50, about 2 to about 100, about 5 to about 10, about 5 to about 25, about 5 to about 50, about 5 to about 100, about 10 to about 20, about 10 to about 50, about 10 to about 100, about 25 to about 50, about 25 to about 75, about 25 to about 100, about 100 to about 200, about 100 to about 500, about 100 to about 1000, about 500 to about 1000, about 1000 to about 5000, about 1000 to about 10,000, about 10,000 to about 20,000, about 10,000 to about 50,000, about 10,000 to about 100,000, or about 50,000 to about 100,000 different references are used. The reference can be any reference sequence described herein. Generally, the reference may be present at a different copy number than the target sequence. For example, the target may have copy number that is about, more than about, less than about, or at least about 1.5-fold, 2-fold, 2.5-fold, 3-fold, 3.5-fold, 4-fold, 4.5 fold, 5-fold, 6-fold, 7-fold, 8-fold, 9-fold, 10-fold, 20-fold, 50-fold, 75-fold, 100-fold, 200-fold, 500-fold, 700-fold, the copy number of the reference number. In other cases the copy number of the target is equal to that of the reference. In still other cases, the reference has a copy number that is about, more than about, less than about, or at least about 1.5-fold, 2-fold, 2.5-fold, 3-fold, 3.5-fold, 4-fold, 4.5 fold, 5-fold, 6-fold, 7-fold, 8-fold, 9-fold, 10-fold, 20-fold, 50-fold, 75-fold, or 100-fold the copy number of the target sequence.

**[0230]** In some embodiments, probes that anneal to each of the references can comprise the same label, e.g., fluorescent dye. Depending on the number of targets to be multiplexed, one can use universal probes, LNA probes, or ligation approaches. Any type of probe described herein can be used to multiplex references.

**[0231]** The methods described herein can be used to measure several gene expression targets in a single reaction. Several assays can be designed to target the lowest expressed gene and bring the measured counts closer to those of the higher expressed gene(s).

**[0232]** If the abundance of expression of two or more different targets on the same gene is being investigated, e.g., by converting mRNA to cDNA, a restriction digest on the cDNA can be performed in order to ensure that the different targets on a given gene end up in different partitions (e.g., droplets). Other methods of fragmenting nucleic acids described herein can be used to separate the targets.

**[0233]** The methods described herein can also apply to measuring viral load levels in a single reaction. A viral load can be measured by estimating the amount of virus in a bodily fluid. In some embodiments, determination of viral load can comprise PCR, reverse transcription PCR, or Nucleic Acid Sequence Based Amplification (NASBA) (transcription-based amplification system (TAS)). For example, PCR can be used to quantify integrated DNA (e.g., integrated into a chromosome of a cell). Reverse transcription PCR can be used to quantify viral RNA by converting it to cDNA. In some embodiments, NASBA is used to convert viral RNA into

DNA, and the DNA can be transcribed into RNA. NASBA can involve annealing a primer to the 3' end of an RNA template, reverse transcribing the RNA template, degrading the RNA template with RNase H, annealing a primer to the 5' end of the DNA strand, and using T7 RNA polymerase to produce a complementary RNA strand. The complementary RNA strand can be reused in the reaction cycle. In some embodiments, multiple references are used to such that the amount of viral nucleic acid and reference nucleic acid in a sample are within the dynamic range of the method used to determine the viral load. In some embodiments, probes used to detect the different references use the same label. In some embodiments, probes used to detect the references comprise different labels.

**[0234]** In some embodiments, multiplexing can also be useful for evening out biological variation where a reference varies in copy number from individual to individual. By averaging across multiple targets and/or reference sequences, the impact of the variation can be reduced. This method can be used, e.g., for diagnostic tests, including those used for measuring copy number alterations.

**[0235]** In some embodiments, a reference sequence that is present at two copies per diploid genome can be used, e.g., a housekeeping gene (e.g., a gene that is required for the maintenance of basic cellular function). Dividing the concentration or amount of the target by the concentration or amount of the reference can yield an estimate of the number of target copies per genome.

**[0236]** A housekeeping gene that can be used as reference in the methods described herein can include a gene that encodes a transcription factor, a transcription repressor, an RNA splicing gene, a translation factor, tRNA synthetase, RNA binding protein, ribosomal protein, RNA polymerase, protein processing protein, heat shock protein, histone, cell cycle regulator, apoptosis regulator, oncogene, DNA repair/replication gene, carbohydrate metabolism regulator, citric acid cycle regulator, lipid metabolism regulator, amino acid metabolism regulator, nucleotide synthesis regulator, NADH dehydrogenase, cytochrome C oxidase, ATPase, mitochondrial protein, lysosomal protein, proteosomal protein, ribonuclease, oxidase/reductase, cytoskeletal protein, cell adhesion protein, channel or transporter, receptor, kinase, growth factor, tissue necrosis factor, etc. Specific examples of housekeeping genes that can be used in the methods described include, e.g., HSP90, Beta-actin, tRNA, rRNA, ATF4, RPP30, and RPL3.

**[0237]** A single copy reference nucleic acid (e.g., gene) can be used to determine copy number variation. Multi-copy reference nucleic acids (e.g., genes) can be used to determine copy number to expand the dynamic range. For example, the multi-copy reference gene can comprise about, or more than about 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 500, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10,000, 20,000, 30,000, 40,000, 50,000, 60,000, 70,000, 80,000, 90,000, or 100,000 copies in a genome.

**[0238]** Ranges can be expressed herein as from "about" one particular value, and/or to "about" another particular value. When such a range is expressed, another embodiment includes from the one particular value and/or to the other

particular value. Similarly, when values are expressed as approximations, by use of the antecedent “about,” it will be understood that the particular value forms another embodiment. It will be further understood that the endpoints of each of the ranges are significant both in relation to the other endpoint, and independently of the other endpoint. The term “about” as used herein refers to a range that is 15% plus or minus from a stated numerical value within the context of the particular usage. For example, about 10 would include a range from 8.5 to 11.5.

## EXAMPLES

### Example 1

**[0239]** One thousand genome equivalents (about 6 ng of DNA) can be sequenced at 100× depth. From 1000 genome equivalents, there can be 2,000 copies of every (normal copy number) target. Steps can be taken such that for every locus, a large majority of fragments end up in separate partitions. Steps can also be taken to ensure that most of the 2,000 fragments are tagged with a unique barcode.

**[0240]** The first goal can be accomplished by increasing the number of partitions. For example, with 100,000 partitions, only about 0.5% of fragments at a particular locus from different chromosomes are expected to end up in the same partition. Note that many such cases will be readily identified by the appearance of distinct alleles from heterozygous SNPs with the same barcode as well as by increased coverage of the locus by a barcode.

**[0241]** In order to ensure that most fragments are tagged with distinct barcodes, a large number of different barcodes can be used, and an approach that distributes barcodes so that any given partition is furnished with a small number (preferably one) of barcode-containing droplets can be used. The distribution can be random so that some partitions receive zero barcodes, some one, some multiple. Thus, for 100,000 partitions 100,000 barcoding droplets can be supplied. In this case, it is anticipated that 37% of the partitions will receive no adaptors and will thus be unavailable for sequencing. The number of barcoding droplets can be increased if sample preservation is a goal. 37% of the partitions can be barcoded with a single barcode and up to 25% can be coded with potentially different barcodes. In the case above, 740 fragments will be unavailable for sequencing, 740 will be sequestered in with their own barcodes and 500 will be sequestered with multiple barcodes. Ideally all of the  $740 \times 1 + 360 \times 2 + \dots = 2,000$  barcodes in the partitions associated with a particular fragment would be unique. If there are 10,000 different barcode types, then more than 80% of the fragments would be uniquely tagged.

**[0242]** If the number of genome equivalents is lower then fewer partitions and barcodes could be used.

**[0243]** Note that perfection is not necessary for this application, because only a small subset of SNPs from any given genomic location can be captured to yield phasing information. It can be acceptable if a substantial fraction of fragments is not informative.

**[0244]** One can attain greater efficiency of sample processing if each partition is supplied with a barcode in a controlled manner. For example, sample containing partitions and barcode containing partitions can be merged using droplet merging technology from RAINDANCE™ (RAINSTORM™). Droplet merging can be performed using a microfluidic circuit similar to FLUIDIGM's array designs. If it can be guar-

anteed that a given partition receives precisely one ADF, fewer ADFs and fewer ADF types can be used.

**[0245]** A microfluidic chip can be used in an analogous manner for partitioning. Sample partitions can be supplied with their own barcodes via a two-dimensional arrangement of channels as described above. A large number of unique barcodes can be readily supplied by combining vertical and horizontal barcodes.

**[0246]** While preferred embodiments of the methods, compositions, systems, and kits described herein have been shown and described herein, it will be obvious to those skilled in the art that such embodiments are provided by way of example only. Numerous variations, changes, and substitutions will now occur to those skilled in the art without departing from the methods, compositions, systems, and kits described herein. It should be understood that various alternatives to the embodiments of the methods, compositions, systems, and kits described herein may be employed in practicing the methods, compositions, systems, and kits. It is intended that the following claims define the scope of the methods, compositions, systems, and kits within the scope of these claims and their equivalents be covered thereby.

What is claimed is:

1. A method comprising:

- a. subdividing a plurality of adaptors into a plurality of first partitions, wherein each of said first partitions has on average a first volume and wherein said adaptors comprise unique barcodes;
- b. subdividing a sample comprising multiple polynucleotides into a plurality of second partitions, wherein each of said second partitions has on average a second volume, wherein said second volume is greater than said first volume;
- c. merging at least one of said first partitions with at least one of said second partitions to form a merged partition; and
- d. tagging one of said multiple polynucleotides, or fragment thereof, with at least one of said adaptors.

2. A method comprising:

- a. subdividing a plurality of adaptors into a plurality of first partitions, wherein each of said first partitions has on average a first volume and wherein said adaptors comprise unique barcodes;
- b. subdividing a sample comprising multiple polynucleotides into a plurality of second partitions, wherein each of said second partitions has on average a second volume, wherein said second volume is less than said first volume;
- c. merging at least one of said first partitions with at least one of said second partitions to form a merged partition; and
- d. tagging one of said multiple polynucleotides, or fragment thereof, with at least one of said adaptors.

3. The method of claim 1 or 2, wherein said first partitions are droplets.

4. The method of claim 3, wherein said second partitions are droplets.

5. The method of claim 3, wherein said droplets are within an immiscible fluid.

6. The method of claim 1 or 2, wherein said polynucleotides are genomic DNA.

7. The method of claim 4, wherein said tagging comprises merging at least one first droplet comprising one of said adaptors with at least one second droplet comprising one of said polynucleotides.

8. The method of claim 1, wherein said first partitions are first droplets and said second partitions are second droplets; and wherein, prior to said merging, said at least one second droplet comprises said at least one first droplet.

9. The method of claim 1, wherein said first partitions are first droplets and said second partitions are second droplets; and wherein, prior to said merging, said at least one second droplet does not comprise said at least one first droplet.

10. The method of claim 2, wherein said first partitions are first droplets and said second partitions are second droplets; and wherein, prior to said merging, said at least one first droplet comprises said at least one second droplet.

11. The method of claim 2, wherein said first partitions are first droplets and said second partitions are second droplets; and wherein, prior to said merging, said at least one first droplet does not comprise said at least one second droplet.

12. The method of claim 1, wherein said second volume is at least two times said volume of said first volume.

13. The method of claim 2, wherein said first volume is at least two times said volume of said second volume.

14. The method of claim 7, further comprising modifying the temperature of said droplets.

15. The method of claim 7, wherein said merging comprises use of a controller such that each of said first droplets merges with each of said second droplets.

16. The method of claim 7, wherein said merging comprises randomly merging droplets comprising polynucleotides with droplets comprising adaptors.

17. The method of claim 1, further comprising pooling said adaptor-tagged polynucleotides, or fragments thereof.

18. The method of claim 1, further comprising analyzing said adaptor-tagged polynucleotides, or fragments thereof.

19. The method of claim 18, wherein said analyzing comprises sequencing said adaptor-tagged polynucleotides, or fragments thereof.

20. The method of claim 1, further comprising determining whether said adaptor-tagged polynucleotides, or fragments thereof, were located in the same partition.

21. The method of claim 19, further comprising estimating the likelihood that any two sequence reads generated by said sequencing came from the same or different partitions.

22. The method of claim 6, wherein said sample is partitioned so that it is unlikely that a given partition comprises two or more polynucleotides, or fragments thereof, from the same locus but from different chromosomes.

23. The method of claim 6, wherein said genomic DNA is high molecular-weight DNA.

24. The method of claim 1 or 2, further comprising fragmenting said polynucleotides within said second partitions to form polynucleotide fragments.

25. The method of claim 24, wherein said polynucleotide fragments are generated by fragmenting said polynucleotides with an endonuclease.

26. The method of claim 1, wherein said polynucleotides are tagged by ligating said adaptors to said polynucleotides within a plurality of said merged partitions.

27. The method of claim 1, wherein said tagging is accomplished using transposons.

28. The method of claim 1, wherein said tagged polynucleotides are amplified.

29. The method of claim 28, wherein said amplification comprises a polymerase chain reaction.

30. The method of claim 1, wherein said polynucleotides are amplified before tagging.

31. The method of claim 1 or 2, wherein each of said first partitions comprises, on average, less than five adaptors.

32. The method of claim 1 or 2, wherein each of said second partitions comprises, on average, less than five of said multiple polynucleotides.

33. The method of claim 1, wherein said subdividing of said sample comprises emulsifying or mixing said sample with said second partitions.

34. The method of claim 2, wherein said subdividing of said plurality of adaptors comprises emulsifying or mixing said plurality of adaptors with said second partitions.

35. The method of claim 30, wherein said amplification is multiple-displacement amplification.

36. A method comprising

a. partitioning organelles into a plurality of partitions, wherein each partition comprises on average less than five organelles per partition;

b. lysing said extracellular organelles in the plurality of partitions, wherein the lysing releases RNA from said organelles;

c. generating tagged cDNA from said released RNA in said plurality of partitions with adaptors comprising a barcode, wherein each partition in the plurality of partitions comprises adaptors with a unique barcode.

37. The method of claim 36, wherein said organelles are extracellular organelles.

38. The method of claim 36, wherein said organelles are exosomes.

39. The method of claim 36, wherein said generating tagged cDNA comprises reverse transcription of said released RNA with partition-specific barcoded primers.

40. The method of claim 36, further comprising sequencing said tagged cDNA.

41. The method of claim 36, further comprising determining if said tagged cDNA is from the same organelle.

42. A method comprising

a. partitioning microorganisms into a plurality of partitions,

b. obtaining polynucleotides from the microorganisms in the plurality of partitions; and

c. tagging the polynucleotides in the plurality of partitions with adaptors comprising a barcode, wherein each partition in the plurality of partitions comprises adaptors with a unique barcode.

43. The method of claim 42, wherein each of said partitions comprises, on average, less than five microorganisms.

44. The method of claim 42, further comprising sequencing the tagged polynucleotides.

45. The method of claim 42, further comprising determining if the tagged polynucleotide fragments are from the same partition.

\* \* \* \* \*