

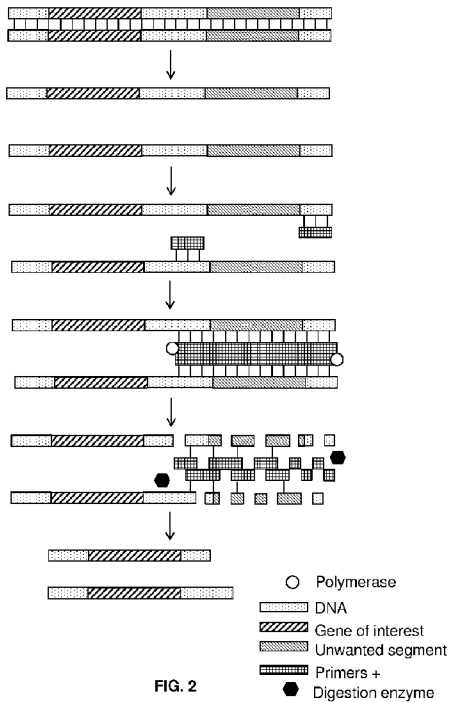


- (51) International Patent Classification:
C12Q 1/68 (2006.01)
- (21) International Application Number:
PCT/US2015/049132
- (22) International Filing Date:
9 September 2015 (09.09.2015)
- (25) Filing Language: English
- (26) Publication Language: English
- (30) Priority Data:
62/048,452 10 September 2014 (10.09.2014) US
- (71) Applicant: GOOD START GENETICS, INC. [US/US];
237 Putnam Avenue, Cambridge, MA 02139 (US).
- (72) Inventors: GOLE, Jeff; 237 Putnam Avenue, Cambridge,
MA 02139 (US). GORE, Athurva; 237 Putnam Avenue,
Cambridge, MA 02139 (US). UMBARGER, Mark; 5
Winchester Street, Unit 101, Brookline, MA 02446 (US).
- (74) Agents: MEYERS, Thomas, C. et al.; Brown Rudnick
LLP, One Financial Center, Boston, MA 02111 (US).

- (81) Designated States (unless otherwise indicated, for every kind of national protection available): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.
- (84) Designated States (unless otherwise indicated, for every kind of regional protection available): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, RU, TJ, TM), European (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

Published:
— with international search report (Art. 21(3))

(54) Title: METHODS FOR SELECTIVELY SUPPRESSING NON-TARGET SEQUENCES



(57) Abstract: The invention generally relates to negative selection of nucleic acids. The invention provides methods and systems that remove unwanted segments of nucleic acid in a sample so that a target gene or region of interest may be analyzed without interference from the unwanted segments. A sample is obtained that includes single-stranded nucleic acid with one or more unwanted segments. Complementary nucleic acid is added to the single-stranded nucleic acid to create a double-stranded region that includes the unwanted segment. The double-stranded region is then digested, leaving single-stranded nucleic acid that includes the target gene or region of interest. This allows paralogs, pseudogenes, repetitive elements, and other segments of the genome that may be similar to the target gene or region of interest to be removed from the sample.

WO 2016/040446 A1

1
2 METHODS FOR SELECTIVELY SUPPRESSING NON-TARGET SEQUENCES

3
4 Cross-Reference to Related Application

5 This application claims priority to, and the benefit of, U.S. Provisional Patent Application
6 Serial No. 62/048,452, filed September 10, 2014, the contents of which are incorporated by
7 reference.

8
9 Technical Field

10 The invention generally relates to negative selection of nucleic acids.

11
12 Background

13 The advent of high-throughput DNA sequencing has the potential to revolutionize
14 modern biology and transform diagnostic medicine. Instruments for next-generation sequencing
15 (NGS) continue to generate more data and become more inexpensive at a rate far outpacing
16 Moore's Law. However, the most popular sequencers have an extremely short read length,
17 limiting their ability to characterize any gene containing paralogous sequence or repetitive
18 elements. As nearly two thirds of the genome is highly repetitive and over 20,000 pseudogenic
19 regions exist, much of the genome is very difficult to characterize in a modern whole-genome
20 sequencing experiment. Unfortunately, for many genes of clinical interest, characterizing those
21 genes is made difficult by the presence of paralogs, pseudogenic homologs, and other segments
22 of the genome that may be similar to the gene of interest and thus stymie attempts to detect,
23 sequence, or isolate the gene of interest. As a result, despite the power of NGS instruments, some
24 disease-related genes and mutations, even where known, are difficult to detect.

25
26 Summary

27 The invention provides methods and systems that remove unwanted segments of nucleic
28 acid in a sample so that a target gene or region of interest may be analyzed without interference
29 from the unwanted segments. A sample is obtained that includes single-stranded nucleic acid
30 with one or more unwanted segments. Primers that are specific or preferentially bind to the
31 unwanted segment are hybridized to the single-stranded nucleic acid within the unwanted region

32 or in a non-repetitive section upstream of the unwanted region and extended by a polymerase to
33 create a double-stranded region that includes the unwanted segment. The double-stranded region
34 is then digested, leaving single-stranded nucleic acid that includes the target gene or region of
35 interest. This allows paralogs, pseudogenes, repetitive elements, and other segments of the
36 genome that may be similar to the target gene or region of interest to be removed from the
37 sample. The target gene or region of interest may thus be detected or characterized by analysis
38 without interference from the unwanted segments. This may provide an improved ability to
39 detect features such as disease-related genes and mutations, thus improving the clinical value of
40 NGS technologies.

41 Systems and methods of the invention may be used to remove unwanted regions from
42 genomic DNA (such as homologous genes, pseudogenes, or repetitive elements) prior to any
43 DNA-based experimental procedure, including but not limited to microarray hybridization,
44 quantitative or standard polymerase chain reaction, multiplex target capture, or DNA sequencing
45 (either targeted or shotgun). Systems and methods of the invention provide for the identification
46 of mutations in previously difficult-to-characterize genes, and therefore allow practitioners to
47 expand the number of genes included in a targeted or whole-genome sequencing assay.

48 In certain aspects, the invention provides a method of removing unwanted segments of a
49 nucleic acid from a sample. The method includes annealing a nucleic acid primer to a portion of
50 a single-stranded nucleic acid that flanks an unwanted segment of the nucleic acid, extending the
51 annealed primer in order to create a double-stranded region that includes the unwanted segment;
52 and digesting the double-stranded region, thereby removing the unwanted segment from the
53 nucleic acid.

54 The nucleic acid in the sample may include DNA, RNA, modified nucleic acids, or
55 combinations thereof. The method may include obtaining a sample from a subject and denaturing
56 double-stranded DNA in the sample. Denaturing can include the use of methods such as
57 exposing the sample to heat, a detergent, or an acidic or basic solution.

58 The primer may be annealed within the unwanted segment or within an area upstream of
59 the unwanted segment and extended. A pair or a number of primers may be used and primers that
60 flank the unwanted segment may be used. In certain embodiments, a plurality of primers are
61 annealed to a plurality of portions of that nucleic acid that flank an unwanted segment. The
62 primer or primers are preferably extended using a polymerase enzyme under conditions

63 sufficient to cause extension of the primer in a template-dependent manner. In some
64 embodiments, a primer or oligonucleotide is hybridized to the unwanted segment to create the
65 double-stranded region containing the unwanted segment without need for an extension step.

66 The double-stranded region is digested. This can include exposing the sample to an
67 enzyme that preferentially digests double-stranded nucleic acid such as certain double-stranded
68 endonucleases, restriction endonucleases, or nicking enzymes. After digestion, the enzyme may
69 be de-activated (e.g., by heat, chemicals, etc.). Digestion preferably results in intact genomic
70 DNA lacking one or more unwanted segment and that is compatible with a nucleic acid analysis
71 assay. Nucleic acid that is not digested may be analyzed by a nucleic acid analysis assay.

72 Assays suitable for analysis of the remaining un-digested nucleic acid may make use of
73 molecular inversion probe capture, hybrid capture, Haloplex, sequencing (e.g., Sanger
74 sequencing, NGS, or both), other methodologies, or combinations thereof. Where the unwanted
75 segment is a paralog, a pseudogene, or non-paralogous repetitive element, such elements may be
76 removed from the sample by methods of the invention.

77 In certain aspects, the invention provides a method of removing nucleic acid from a
78 sample. The method includes annealing at least one oligonucleotide to single-stranded DNA in a
79 sample, wherein the single-stranded DNA comprises target and non-target sequence. The
80 oligonucleotide may be annealed to the non-target sequence to create double-stranded DNA that
81 includes the non-target sequence or the oligonucleotide may be annealed elsewhere and extended
82 to create double-stranded DNA that includes the non-target sequence. The non-target sequence is
83 removed from the sample by digesting the double-stranded DNA. The target sequence may be
84 analyzed using, e.g., molecular inversion probes, microarray hybridization, multiplex ligation-
85 dependent probe amplification (MLPA), sequencing, fingerprinting techniques such as RFLP/
86 AFLP, chromatography, others, or combinations thereof. In some embodiments, the method
87 includes first obtaining the sample from a subject and denaturing double-stranded subject DNA
88 to produce the single-stranded DNA. Preferably, that single-stranded DNA consists essentially of
89 genomic DNA from the subject prior to the annealing of the oligo. The annealing may include
90 annealing a pair of oligonucleotides to the single-stranded DNA at sites that flank the non-target
91 sequence (i.e., to remove both strands of the unwanted segment or non-target sequence. In
92 certain embodiments, the target and non-target sequence are both located on at least one single
93 strand of the single-stranded DNA, and extending the at least one oligonucleotide and digesting

94 the double-stranded DNA results in removing the non-target sequence from the at least one
95 single strand of the single-stranded DNA.

96

97 Brief Description of the Drawings

98 FIG. 1 diagrams a method of removing unwanted segments of a nucleic acid.

99 FIG. 2 illustrates methods according to certain embodiments.

100 FIG. 3 gives a diagram of a system according to embodiments of the invention.

101

102 Detailed Description

103 To enable the characterization of difficult genomic regions using high-throughput short-
104 read sequencing, the invention provides methods for the removal of unwanted genomic regions
105 from a population of DNA molecules (e.g. genomic DNA). Most DNA-based techniques rely on
106 the amplification of specific regions of interest or sequencing library molecules in a positive
107 selection process (e.g. amplification utilizing primers that are unique to a single paralog).
108 Methods of the invention instead involve a negative selection technique that removes any
109 undesired analogous sequence, allowing application of standard high-throughput sequencing
110 techniques or other analyses to any difficult-to-characterize gene of interest.

111 Applicability of methods of the invention may be illustrated by reference to two
112 exemplary genes of interest for which direct high-throughput sequencing-based approaches are
113 currently insufficient. One gene is “glucosidase beta acid,” or GBA, which has been implicated
114 as causative in Gaucher disease. Currently, long-range polymerase chain reaction experiments
115 are required to characterize this gene, as a pseudogene with nearly identical sequence exists a
116 mere 15,000 base pairs away. By removing this pseudogenic region using the invention, GBA
117 can be characterized with high specificity, enabling construction of a genetic screen for Gaucher
118 disease. This gene is a suitable target for methods of the invention, as it is relatively small and
119 contains nearby unique flanking sequence.

120 An additional gene of interest is “survival of motor neuron 1,” or SMN1. This gene has
121 been implicated in spinal muscular atrophy. Currently, due to the presence of a paralogous gene
122 known as SMN2 that is 100,000 base pairs away from SMN1, characterization of SMN1 is
123 extremely challenging. By removing SMN2 using the invention, SMA could be screened for
124 with a high-throughput sequencing approach that would not require a complex statistical model.

125 Additionally, novel causative mutations in genes such as SMN1 could also be identified.
126 This gene is a suitable target for methods of the invention. It is of a suitable size and flanked by
127 highly repetitive regions.

128 FIG. 1 diagrams a method 101 of removing unwanted segments of a nucleic acid from a
129 sample according to embodiments of the invention. The method includes obtaining 105 a sample
130 that includes nucleic acid. An oligonucleotide is annealed 109 to an unwanted segment of the
131 nucleic acid or a portion of the nucleic acid that flanks an unwanted segment of the nucleic acid.
132 In embodiments in which the oligonucleotide flanks the unwanted segment, the oligonucleotide
133 is extended 113 to create a double-stranded region that includes the unwanted segment. In
134 embodiments in which the oligonucleotide is annealed to the unwanted segment, a double-
135 stranded region that includes the unwanted segment is created by virtue of the hybridization of
136 the oligonucleotide at that segment. The double-stranded region is digested 117, thus removing
137 the unwanted segment from the nucleic acid. This allows for a region or gene of interest to be
138 analyzed 121.

139 The sample that includes nucleic acid may be obtained 105 by any suitable method. The
140 sample may be obtained from a tissue or body fluid that is obtained in any clinically acceptable
141 manner. Body fluids may include mucous, blood, plasma, serum, serum derivatives, bile, blood,
142 maternal blood, phlegm, saliva, sweat, amniotic fluid, menstrual fluid, mammary fluid, follicular
143 fluid of the ovary, fallopian tube fluid, peritoneal fluid, urine, and cerebrospinal fluid (CSF),
144 such as lumbar or ventricular CSF. A sample may also be a fine needle aspirate or biopsied
145 tissue. A sample also may be media containing cells or biological material. Samples may also be
146 obtained from the environment (e.g., air, agricultural, water and soil) or may include research
147 samples (e.g., products of a nucleic acid amplification reaction, or purified genomic DNA, RNA,
148 proteins, etc.).

149 Isolation, extraction or derivation of genomic nucleic acids may be performed by
150 methods known in the art. Isolating nucleic acid from a biological sample generally includes
151 treating a biological sample in such a manner that genomic nucleic acids present in the sample
152 are extracted and made available for analysis. Generally, nucleic acids are extracted using
153 techniques such as those described in Green & Sambrook, 2012, *Molecular Cloning: A*
154 *Laboratory Manual* 4 edition, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY
155 (2028 pages), the contents of which are incorporated by reference herein. A kit may be used to

156 extract DNA from tissues and bodily fluids and certain such kits are commercially available
157 from, for example, BD Biosciences Clontech (Palo Alto, CA), Epicentre Technologies (Madison,
158 WI), Genra Systems, Inc. (Minneapolis, MN), and Qiagen Inc. (Valencia, CA). User guides that
159 describe protocols are usually included in such kits.

160 It may be preferable to lyse cells to isolate genomic nucleic acid. Cellular extracts can be
161 subjected to other steps to drive nucleic acid isolation toward completion by, e.g., differential
162 precipitation, column chromatography, extraction with organic solvents, filtration,
163 centrifugation, others, or any combination thereof. The genomic nucleic acid may be
164 resuspended in a solution or buffer such as water, Tris buffers, or other buffers. In certain
165 embodiments the genomic nucleic acid can be re-suspended in Qiagen DNA hydration solution,
166 or other Tris-based buffer of a pH of around 7.5.

167 Any nucleic acid may be analyzed using methods of the invention. Nucleic acids suitable
168 for use in aspects of the invention may include without limit genomic DNA, genomic RNA,
169 synthesized nucleic acids, whole or partial genome amplification product, and high molecular
170 weight nucleic acids, e.g. individual chromosomes. In certain embodiments, a sample is obtained
171 that includes double-stranded DNA, such as bulk genomic DNA from a subject, and the double-
172 stranded DNA is then denatured.

173 Double stranded nucleic acid may be denatured using any suitable method such as, for
174 example, through the use of heat, detergent incubation, or an acidic or basic solution.

175 FIG. 2 illustrates the progress of methods according to certain embodiments. As shown in
176 FIG. 2, methods may start with double stranded DNA (dotted shading if not otherwise hatched)
177 that contains a gene of interest (first angled hatching pattern) and a paralog of the gene of interest
178 (second angled hatching pattern). It will be appreciated that methods of the invention may
179 operate starting with any suitable nucleic acid such as double- or single-stranded DNA or RNA
180 or any combination thereof. The unwanted segment may be any sequence for which removal is
181 desired from the starting nucleic acid. For example, the unwanted segment may include a paralog
182 or homolog of a gene or region of interest; a pseudogene; or non-paralogous repetitive element.
183 As used herein, homolog refers to a gene related to a second gene by descent from a common
184 ancestral DNA sequence. Homolog describes the relationship between genes separated by the
185 event of speciation (i.e., orthology) or to the relationship between genes separated by the event of
186 genetic duplication (i.e., paralogy). Orthologs generally refers to genes in different species that

187 evolved from a common ancestral gene by speciation. Normally, orthologs retain the same
188 function in the course of evolution and paralogs are genes related by duplication within a
189 genome. See Fitch, 1970, Distinguishing homologs from analogous proteins, *Syst Biol* 19(2):99-
190 113 and Jensen, 2001, Orthologs and paralogs—we need to get it right, *Genome Biol* 2(8):1002-
191 1002.3. Pseudogenes include dysfunctional relatives of genes that have lost their protein-coding
192 ability or are otherwise no longer expressed in the cell. Methods of the invention may be used to
193 target a pseudogene that is present as a homolog to another gene or pseudogene within a sample
194 and methods of the invention may be used to target a pseudogene that is present even where no
195 known homologs of the pseudogene are suspected to also be present in the sample.

196 As illustrated in FIG. 2, the double-stranded DNA is denatured into its two
197 complementary strands prior to primer hybridization. Any suitable method may be used to
198 denature nucleic acid. Heat-based denaturing is a process by which double-stranded nucleic acid
199 unwinds and separates into single-stranded strands. Heat denaturation of a nucleic acid of an
200 unknown sequence typically uses a temperature high enough to ensure denaturation of even
201 nucleic acids having a very high GC content, e.g., 95° C-98° C in the absence of any chemical
202 denaturant. It is well within the abilities of one of ordinary skill in the art to optimize the
203 conditions (e.g., time, temperature, etc.) for denaturation of the nucleic acid. Temperatures
204 significantly lower than 95° C can also be used if the DNA contains nicks (and therefore sticky
205 overhangs of low T_m), sequence of sufficiently low T_m , or chemical additives such as betaine.

206 Denaturing nucleic acids with the use of pH is also well known in the art, and such
207 denaturation can be accomplished using any method known in the art such as introducing a
208 nucleic acid to high or low pH, low ionic strength, and/or heat, which disrupts base-pairing
209 causing a double-stranded helix to dissociate into single strands. For methods of pH-based
210 denaturation see, for example, Ageno et al., 1969, The alkaline denaturation of DNA, *Biophys J*
211 9:1281-1311.

212 Nucleic acids can also be denatured via electro-chemical means, for example, by
213 applying a voltage to a nucleic acid within a solution by means of an electrode. Varying methods
214 of denaturing by applying a voltage are discussed in detail in U.S. Patent Nos. 6,197,508 and
215 U.S. Patent No. 5,993,611. After denaturation, unwanted segments can be targeted for removal.

216 Methods of the invention include targeting unwanted segments of nucleic acid for
217 removal. An unwanted segment of nucleic acid can be targeted for removal by making it into a

218 double-stranded segment. The unwanted segment can be made double-stranded by hybridizing a
219 complementary oligonucleotide to the unwanted segment, by hybridizing a complementary
220 oligonucleotide to a genomic segment flanking the unwanted segment and extending the
221 oligonucleotide, or a combination thereof (e.g., an oligonucleotide can be hybridized so that it
222 sits partially within the unwanted segment and then extended via methods described herein).

223 In certain embodiments, the oligonucleotide to be hybridized is a primer that is unique to
224 the unwanted segment. For example, methods may include using a primer that is unique to a
225 certain paralog or other element. The invention provides methods of making a primer and primer
226 extension reactions that are unique to a paralog or similar segment by including or using a primer
227 with a 3' end that terminates on a differentiating base (i.e., the 3'-most base or bases of the
228 primer may be complementary to a base or bases that appear only in association with the
229 segment (e.g., paralog) targeted for removal.

230 In some embodiments, double stranded DNA is created by hybridization alone (e.g.,
231 rather than by using oligonucleotide primer with polymerase extension). One or more long
232 segments of nucleic acid complementary to the unwanted segments could be used. For example,
233 long segments of synthetic DNA could be used. The segments of complementary nucleic acid
234 could have any suitable length such as, for example, tens of bases, hundreds of bases, length of
235 an exon, length of a gene, etc. Use of one or more long segments of nucleic acid complementary
236 to the unwanted segments (e.g., followed by digestion of dsDNA) may provide for enrichment
237 of, for example, target relative to non-target.

238 As noted above, the recognition site for the oligonucleotide, primer, or complementary
239 nucleic acid may flank the unwanted segment, lie within the unwanted segment, or both.
240 Additionally, methods may include using one or any suitable number of oligonucleotides or
241 primers to target an unwanted segment or segments of nucleic acid.

242 In the non-limiting, illustrative embodiment shown in FIG. 2, primers (cross-hatching
243 pattern) are annealed to unique genomic segments flanking the paralogous region. The primer
244 may be annealed at any suitable location. For example, it may be preferable to anneal any of the
245 one or more primers to a portion within 50 or fewer bases from the unwanted segment, although
246 it may not be necessary to anneal the primers within 50 bases of the unwanted region. As shown
247 in FIG. 2, primers are annealed at locations that flank the unwanted segment, i.e., each primer of
248 a pair hybridizes to its target strand in a region that flanks the 5' end of the unwanted segment. In

249 this way, extension of the primers will result in most or all of the unwanted segments being
250 present in exclusively double-stranded form, whereas the desired region(s) should remain in a
251 primarily single-stranded state.

252 In certain embodiments, polymerase (drawn as an open circle in FIG. 2) is used to
253 perform second-strand synthesis over the paralogous region. Extending the annealed primer
254 creates a double-stranded region that includes the unwanted segment. The primer is extended
255 using a polymerase enzyme under conditions sufficient to cause extension of the primer in a
256 template-dependent manner. Suitable polymerase enzymes include phi29, Bst, Exo-minus E.
257 Coli Polymerase I, Taq Polymerase, and T7 Polymerase I.

258 An enzymatic digestion (the digestion enzyme is represented by a darkened hexagon in
259 FIG. 2) is then used to degrade only the double-stranded paralogous region, leaving behind the
260 gene of interest. Any suitable digestion platform may be employed such as, for example,
261 dsDNase, fragmentase, a non-specific nicking enzyme such as a modified Vvn, restriction
262 enzymes such as MspJI and FspEI, and a combination of USER plus T7 endonuclease I.

263 Thermo Scientific dsDNase is an engineered shrimp DNase designed for rapid and safe
264 removal of contaminating genomic DNA from RNA samples. It is an endonuclease that cleaves
265 phosphodiester bonds in DNA to yield oligonucleotides with 5'-phosphate and 3'-hydroxyl
266 termini. Highly specific activity towards double-stranded DNA ensures that RNA and single-
267 stranded DNA such as cDNA and primers are not cleaved. dsDNase is easily inactivated by
268 moderate heat treatment (55°C). Thermo Scientific dsDNase is available from Thermo Fisher
269 Scientific, Inc. (Waltham, MA).

270 Fragmentase includes the enzyme sold under the trademark NEBNEXT dsDNA
271 fragmentase by New England Biolabs (Ipswich, MA). NEBNEXT dsDNA fragmentase generates
272 dsDNA breaks in a time-dependent manner to yield 50–1,000 bp DNA fragments depending on
273 reaction time. NEBNext dsDNA Fragmentase contains two enzymes, one randomly generates
274 nicks on dsDNA and the other recognizes the nicked site and cuts the opposite DNA strand
275 across from the nick, producing dsDNA breaks. The resulting DNA fragments contain short
276 overhangs, 5'-phosphates, and 3'-hydroxyl groups. The random nicking activity of NEBNext
277 dsDNA Fragmentase has been confirmed by preparing libraries for next-generation sequencing.
278 A comparison of the sequencing results between genomic DNA (gDNA) prepared with
279 NEBNext dsDNA fragmentase and with mechanical shearing demonstrates that the NEBNext

280 dsDNA Fragmentase does not introduce any detectable bias during the sequencing library
281 preparation and no difference in sequence coverage is observed using the two methods

282 The *Vibrio vulnificus* nuclease, Vvn, is a non-specific periplasmic nuclease capable of
283 digesting DNA and RNA. It has been suggested that Vvn hydrolyzes DNA by a general single-
284 metal ion mechanism. See Li, et al., 2003, DNA binding and cleavage by the periplasmic
285 nuclease Vvn: a novel structure with a known active site, EMBO J 22(15):4014-4025.

286 MspJI is a modification dependent endonuclease that recognizes certain methylation
287 patterns. The most common epigenetic modifications found in eukaryotic organisms are
288 methylation marks at CpG or CHG sites. A subset of these modified sites are recognized and
289 cleaved by MspJI. MspJI is available from New England Biolabs. T7 Endonuclease I recognizes
290 and cleaves non-perfectly matched DNA, cruciform DNA structures, Holliday structures or
291 junctions, hetero-duplex DNA and more slowly, nicked double-stranded DNA. The cleavage site
292 is at the first, second or third phosphodiester bond that is 5' to the mismatch. The protein is the
293 product of T7 gene 3. Any other suitable enzyme for digesting the target unwanted segments
294 may be used.

295 The added enzymes may then be deactivated using an irreversible heat or chemical
296 treatment, leaving genomic DNA lacking an intact undesired region(s) yet still compatible with
297 any downstream assay (e.g. molecular inversion probe capture or any other library construction
298 methodology).

299 The digesting step results in intact genomic DNA lacking one or more unwanted segment
300 and that is compatible with a nucleic acid analysis assay. This DNA can then be utilized for any
301 downstream assay. Downstream assays may include molecular inversion capture, sequencing,
302 others, or a combination thereof.

303 Methods of the invention can be used to negatively select out pseudogenic regions from
304 the genome. Methods of the invention can be combined with a genetic test, screening, or other
305 assay in order to screen patients for mutations in a gene (e.g., GBA, SMN1, or other genes
306 containing paralogous regions). Some background may be found in published international
307 patent application WO 2013/191775, to Nugen Technologies, Inc.

308 After removing the unwanted segment from the nucleic acid, the sample may be enriched
309 for genes of interest using methods known in the art, such as hybrid capture. Methods suitable
310 for use may be found discussed in U.S. Pat. 8,529,744; U.S. Pat. 7,985,716; U.S. Pat. 7,666,593;

311 and U.S. Pat. 6,613,516. As will be described in more detail below, a preferable capture method
312 uses molecular inversion probes.

313 Nucleic acids, including genomic nucleic acids, can be fragmented using any of a variety
314 of methods, such as mechanical fragmenting, chemical fragmenting, and enzymatic fragmenting.
315 Methods of nucleic acid fragmentation are known in the art and include, but are not limited to,
316 DNase digestion, sonication, mechanical shearing, and the like. U.S. Pub 2005/0112590 provides
317 a general overview of various methods of fragmenting known in the art.

318 Genomic nucleic acids can be fragmented into uniform fragments or randomly
319 fragmented. In certain aspects, nucleic acids are fragmented to form fragments having a fragment
320 length of about 5 kilobases or 100 kilobases. Desired fragment length and ranges of fragment
321 lengths can be adjusted depending on the type of nucleic acid targets one seeks to capture and the
322 design and type of probes such as molecular inversion probes (MIPs) that will be used. Chemical
323 fragmentation of genomic nucleic acids can be achieved using methods such as a hydrolysis
324 reaction or by altering temperature or pH. Nucleic acid may be fragmented by heating a nucleic
325 acid immersed in a buffer system at a certain temperature for a certain period of time to initiate
326 hydrolysis and thus fragment the nucleic acid. The pH of the buffer system, duration of heating,
327 and temperature can be varied to achieve a desired fragmentation of the nucleic acid. Mechanical
328 shearing of nucleic acids into fragments can be used e.g., by hydro-shearing, trituration through a
329 needle, and sonication. The nucleic acid can also be sheared via nebulization, hydro-shearing,
330 sonication, or others. See U.S. Pat. 6,719,449; U.S. Pat. 6,948,843; and U.S. Pat. 6,235,501.
331 Nucleic acid may be fragmented enzymatically. Enzymatic fragmenting, also known as
332 enzymatic cleavage, cuts nucleic acids into fragments using enzymes, such as endonucleases,
333 exonucleases, ribozymes, and DNazymes. Varying enzymatic fragmenting techniques are well-
334 known in the art. Additionally, DNA may be denatured again as needed after the digestion and
335 any other sample prep steps. For example, during a fragmentation step, ssDNA may anneal to
336 form dsDNA and it may be desirable to again denature the dsDNA. In certain embodiments, the
337 sample nucleic acid is captured or targeted using any suitable capture method or assay such as
338 hybridization capture or capture by probes such as MIPs.

339 MIPs, or molecular inversion probes, can be used to detect or amplify particular nucleic
340 acid sequences in complex mixtures. Use of molecular inversion probes has been demonstrated
341 for detection of single nucleotide polymorphisms (Hardenbol et al., 2005, Highly multiplexed

342 molecular inversion probe genotyping: over 10,000 targeted SNPs genotyped in a single tube
343 assay, *Genome Res* 15:269-75) and for preparative amplification of large sets of exons (Porreca
344 et al., 2007, Multiplex amplification of large sets of human exons, *Nat Methods* 4:931-6,
345 Krishnakumar et al., 2008, A comprehensive assay for targeted multiplex amplification of human
346 DNA sequences, *PNAS* 105:9296-301). One of the main benefits of the method is in its capacity
347 for a high degree of multiplexing, because generally thousands of targets may be captured in a
348 single reaction containing thousands of probes.

349 In certain embodiments, molecular inversion probes include a universal portion flanked
350 by two unique targeting arms. The targeting arms are designed to hybridize immediately
351 upstream and downstream of a specific target sequence located on a genomic nucleic acid
352 fragment. The molecular inversion probes are introduced to nucleic acid fragments to perform
353 capture of target sequences located on the fragments. According to the invention, fragmenting
354 aids in capture of target nucleic acid by molecular inversion probes. As described in greater
355 detail herein, after capture of the target sequence (e.g., locus) of interest, the captured target may
356 further be subjected to an enzymatic gap-filling and ligation step, such that a copy of the target
357 sequence is incorporated into a circle. Capture efficiency of the MIP to the target sequence on
358 the nucleic acid fragment can be improved by lengthening the hybridization and gap-filing
359 incubation periods. (See, e.g., Turner et al., 2009, Massively parallel exon capture and library-
360 free resequencing across 16 genomes, *Nature Methods* 6:315-316.)

361 A library of molecular inversion probes may be created and used in capturing DNA of
362 genomic regions of interests (e.g., SMN1, SMN2, control DNA). The library includes a plurality
363 of oligonucleotide probes capable of capturing one or more genomic regions of interest (e.g.,
364 SMN1, SMN2 and control loci) within the samples to be tested.

365 The result of MIP capture as described above is a library of circular target probes, which
366 then can be processed in a variety of ways. Adaptors for sequencing may be attached during
367 common linker-mediated PCR, resulting in a library with non-random, fixed starting points for
368 sequencing. For preparation of a shotgun library, a common linker-mediated PCR is performed
369 on the circle target probes, and the post-capture amplicons are linearly concatenated, sheared,
370 and attached to adaptors for sequencing. Methods for shearing the linear concatenated captured
371 targets can include any of the methods disclosed for fragmenting nucleic acids discussed above.

372 In certain aspects, performing a hydrolysis reaction on the captured amplicons in the presence of
373 heat is the desired method of shearing for library production.

374 In some embodiments, the amount of target nucleic acid and probe used for each reaction
375 is normalized to avoid any observed differences being caused by differences in concentrations or
376 ratios. In some embodiments, in order to normalize genomic DNA and probe, the genomic DNA
377 concentration is read using a standard spectrophotometer or by fluorescence (e.g., using a
378 fluorescent intercalating dye). The probe concentration may be determined experimentally or
379 using information specified by the probe manufacturer.

380 Similarly, once a locus has been captured, it may be amplified and/or sequenced in a
381 reaction involving one or more primers. The amount of primer added for each reaction can range
382 from 0.1 pmol to 1 nmol, 0.15 pmol to 1.5 nmol (for example around 1.5 pmol). However, other
383 amounts (e.g., lower, higher, or intermediate amounts) may be used.

384 A targeting arm may be designed to hybridize (e.g., be complementary) to either strand of
385 a genetic locus of interest if the nucleic acid being analyzed is DNA (e.g., genomic DNA). For
386 MIP probes, whichever strand is selected for one targeting arm will be used for the other one. In
387 the context of RNA analysis, a targeting arm should be designed to hybridize to the transcribed
388 RNA. It also should be appreciated that MIP probes referred to herein as “capturing” a target
389 sequence are actually capturing it by template-based synthesis rather than by capturing the actual
390 target molecule (other than for example in the initial stage when the arms hybridize to it or in the
391 sense that the target molecule can remain bound to the extended MIP product until it is denatured
392 or otherwise removed).

393 A targeting arm may include a sequence that is complementary to one allele or mutation
394 (e.g., a SNP or other polymorphism, a mutation, etc.) so that the probe will preferentially
395 hybridize (and capture) target nucleic acids having that allele or mutation. Sequence tags (also
396 referred to as barcodes) may be designed to be unique in that they do not appear at other
397 positions within a probe or a family of probes and they also do not appear within the sequences
398 being targeted. Uniformity and reproducibility can be increased by designing multiple probes per
399 target, such that each base in the target is captured by more than one probe.

400 The length of a capture molecule on a nucleic acid fragment (e.g., a target nucleic acid or
401 sub-region thereof) may be selected based upon multiple considerations. For example, where
402 analysis of a target involves sequencing, e.g., with a next-generation sequencer, the target length

403 should typically match the sequencing read-length so that shotgun library construction is not
404 necessary. However, it should be appreciated that captured nucleic acids may be sequenced using
405 any suitable sequencing technique as aspects of the invention are not limited in this respect.

406 It is also to be appreciated that some target nucleic acids on a nucleic acid fragment are
407 too large to be captured with one probe. Consequently, it may be helpful to capture multiple sub-
408 regions of a target nucleic acid in order to analyze the full target.

409 Methods of the invention also provide for combining the method of fragmenting the
410 nucleic acid prior to capture with other MIP capture techniques that are designed to increase
411 target uniformity, reproducibility, and specificity. Other MIP capture techniques are shown in
412 U.S. Pub. 2012/0165202, incorporated by reference.

413 Multiple probes, e.g., MIPs, can be used to amplify each target nucleic acid. In some
414 embodiments, the set of probes for a given target can be designed to ‘tile’ across the target,
415 capturing the target as a series of shorter sub targets. In some embodiments, where a set of
416 probes for a given target is designed to ‘tile’ across the target, some probes in the set capture
417 flanking non-target sequence). Alternately, the set can be designed to ‘stagger’ the exact
418 positions of the hybridization regions flanking the target, capturing the full target (and in some
419 cases capturing flanking non-target sequence) with multiple probes having different targeting
420 arms, obviating the need for tiling. The particular approach chosen will depend on the nature of
421 the target set. For example, if small regions are to be captured, a staggered-end approach might
422 be appropriate, whereas if longer regions are desired, tiling might be chosen. In all cases, the
423 amount of bias-tolerance for probes targeting pathological loci can be adjusted by changing the
424 number of different MIPs used to capture a given molecule.

425 Probes for MIP capture reactions may be synthesized on programmable microarrays
426 because of the large number of sequences required. Because of the low synthesis yields of these
427 methods, a subsequent amplification step is required to produce sufficient probe for the MIP
428 amplification reaction. The combination of multiplex oligonucleotide synthesis and pooled
429 amplification results in uneven synthesis error rates and representational biases. By synthesizing
430 multiple probes for each target, variation from these sources may be averaged out because not all
431 probes for a given target will have the same error rates and biases.

432 Using methods described herein, a single copy of a specific target nucleic acid may be
433 amplified to a level that can be sequenced. Further, the amplified segments created by an

434 amplification process such as PCR may be, themselves, efficient templates for subsequent PCR
435 amplifications.

436 Amplification or sequencing adapters or barcodes, or a combination thereof, may be
437 attached to the fragmented nucleic acid. Such molecules may be commercially obtained, such as
438 from Integrated DNA Technologies (Coralville, IA). In certain embodiments, such sequences are
439 attached to the template nucleic acid molecule with an enzyme such as a ligase. Suitable ligases
440 include T4 DNA ligase and T4 RNA ligase, available commercially from New England Biolabs
441 (Ipswich, MA). The ligation may be blunt ended or via use of complementary overhanging ends.
442 In certain embodiments, following fragmentation, the ends of the fragments may be repaired,
443 trimmed (e.g. using an exonuclease), or filled (e.g., using a polymerase and dNTPs) to form
444 blunt ends. In some embodiments, end repair is performed to generate blunt end 5'
445 phosphorylated nucleic acid ends using commercial kits, such as those available from Epicentre
446 Biotechnologies (Madison, WI). Upon generating blunt ends, the ends may be treated with a
447 polymerase and dATP to form a template independent addition to the 3'-end and the 5'-end of
448 the fragments, thus producing a single A overhanging. This single A can guide ligation of
449 fragments with a single T overhanging from the 5'-end in a method referred to as T-A cloning.
450 Alternatively, because the possible combination of overhangs left by the restriction enzymes are
451 known after a restriction digestion, the ends may be left as-is, i.e., ragged ends. In certain
452 embodiments double stranded oligonucleotides with complementary overhanging ends are used.

453 In certain embodiments, one or more bar code is attached to each, any, or all of the
454 fragments. A bar code sequence generally includes certain features that make the sequence useful
455 in sequencing reactions. The bar code sequences are designed such that each sequence is
456 correlated to a particular portion of nucleic acid, allowing sequence reads to be correlated back to
457 the portion from which they came. Methods of designing sets of bar code sequences is shown for
458 example in U.S. Pat. 6,235,475, the contents of which are incorporated by reference herein in
459 their entirety. In certain embodiments, the bar code sequences range from about 5 nucleotides to
460 about 15 nucleotides. In a particular embodiment, the bar code sequences range from about 4
461 nucleotides to about 7 nucleotides. In certain embodiments, the bar code sequences are attached
462 to the template nucleic acid molecule, e.g., with an enzyme. The enzyme may be a ligase or a
463 polymerase, as discussed above. Attaching bar code sequences to nucleic acid templates is
464 shown in U.S. Pub. 2008/0081330 and U.S. Pub. 2011/0301042, the content of each of which is

465 incorporated by reference herein in its entirety. Methods for designing sets of bar code sequences
466 and other methods for attaching bar code sequences are shown in U.S. Pats. 6,138,077;
467 6,352,828; 5,636,400; 6,172,214; 6,235,475; 7,393,665; 7,544,473; 5,846,719; 5,695,934;
468 5,604,097; 6,150,516; 7,537,897; 6,172,218; and 5,863,722, the content of each of which is
469 incorporated by reference herein in its entirety. After any processing steps (e.g., obtaining,
470 isolating, fragmenting, amplification, or barcoding), nucleic acid can be sequenced.

471 Sequencing may be by any method known in the art. DNA sequencing techniques include
472 classic dideoxy sequencing reactions (Sanger method) using labeled terminators or primers and
473 gel separation in slab or capillary, sequencing by synthesis using reversibly terminated labeled
474 nucleotides, pyrosequencing, 454 sequencing, Illumina/Solexa sequencing, allele specific
475 hybridization to a library of labeled oligonucleotide probes, sequencing by synthesis using allele
476 specific hybridization to a library of labeled clones that is followed by ligation, real time
477 monitoring of the incorporation of labeled nucleotides during a polymerization step, polony
478 sequencing, and SOLiD sequencing. Separated molecules may be sequenced by sequential or
479 single extension reactions using polymerases or ligases as well as by single or sequential
480 differential hybridizations with libraries of probes.

481 A sequencing technique that can be used includes, for example, Illumina sequencing.
482 Illumina sequencing is based on the amplification of DNA on a solid surface using fold-back
483 PCR and anchored primers. Genomic DNA is fragmented, and adapters are added to the 5' and 3'
484 ends of the fragments. DNA fragments that are attached to the surface of flow cell channels are
485 extended and bridge amplified. The fragments become double stranded, and the double stranded
486 molecules are denatured. Multiple cycles of the solid-phase amplification followed by
487 denaturation can create several million clusters of approximately 1,000 copies of single-stranded
488 DNA molecules of the same template in each channel of the flow cell. Primers, DNA polymerase
489 and four fluorophore-labeled, reversibly terminating nucleotides are used to perform sequential
490 sequencing. After nucleotide incorporation, a laser is used to excite the fluorophores, and an
491 image is captured and the identity of the first base is recorded. The 3' terminators and
492 fluorophores from each incorporated base are removed and the incorporation, detection and
493 identification steps are repeated. Sequencing according to this technology is described in U.S.
494 Pat. 7,960,120; U.S. Pat. 7,835,871; U.S. Pat. 7,232,656; U.S. Pat. 7,598,035; U.S. Pat.
495 6,911,345; U.S. Pat. 6,833,246; U.S. Pat. 6,828,100; U.S. Pat. 6,306,597; U.S. Pat. 6,210,891;

496 U.S. Pub. 2011/0009278; U.S. Pub. 2007/0114362; U.S. Pub. 2006/0292611; and U.S. Pub.
497 2006/0024681, each of which are incorporated by reference in their entirety.

498 Sequencing generates a plurality of reads. Reads generally include sequences of
499 nucleotide data wherein read length may be associated with sequencing technology. For
500 example, the single-molecule real-time (SMRT) sequencing technology of Pacific Bio produces
501 reads thousands of base-pairs in length. For 454 pyrosequencing, read length may be about 700
502 bp in length. In some embodiments, reads are less than about 500 bases in length, or less than
503 about 150 bases in length, or less than about 90 bases in length. In certain embodiments, reads
504 are between about 80 and about 90 bases, e.g., about 85 bases in length. In some embodiments,
505 these are very short reads, i.e., less than about 50 or about 30 bases in length.

506 The sequence reads may be analyzed to characterize the target gene or region of interest.
507 For example, mutations can be “called” (i.e., identified and reported), a haplotype for the sample
508 may be reported, or other analyses may be performed. Mutation calling is described in U.S. Pub.
509 2013/0268474. In some embodiments, an analysis may include determining copy number states
510 of genomic regions of interest. A set of sequence reads can be analyzed by any suitable method
511 known in the art. For example, in some embodiments, sequence reads are analyzed by hardware
512 or software provided as part of a sequence instrument. In some embodiments, individual
513 sequence reads are reviewed by sight (e.g., on a computer monitor). A computer program may be
514 written that pulls an observed genotype from individual reads. In certain embodiments, analyzing
515 the reads includes assembling the sequence reads and then genotyping the assembled reads.

516 Sequence assembly can be done by methods known in the art including reference-based
517 assemblies, de novo assemblies, assembly by alignment, or combination methods. Assembly can
518 include methods described in U.S. Pat. 8,209,130 titled Sequence Assembly by Porecca and
519 Kennedy, the contents of each of which are hereby incorporated by reference in their entirety for
520 all purposes. In some embodiments, sequence assembly uses the low coverage sequence
521 assembly software (LOCAS) tool described by Klein, et al., in LOCAS-A low coverage
522 sequence assembly tool for re-sequencing projects, PLoS One 6(8) article 23455 (2011), the
523 contents of which are hereby incorporated by reference in their entirety. Sequence assembly is
524 described in U.S. Pat. 8,165,821; U.S. Pat. 7,809,509; U.S. Pat. 6,223,128; U.S. Pub.
525 2011/0257889; and U.S. Pub. 2009/0318310, the contents of each of which are hereby
526 incorporated by reference in their entirety.

527 Functions described above such as sequence read analysis or assembly can be
528 implemented using systems of the invention that include software, hardware, firmware,
529 hardwiring, or combinations of any of these.

530 FIG. 3 gives a diagram of a system 301 according to embodiments of the invention.
531 System 301 may include an analysis instrument 303 which may be, for example, a sequencing
532 instrument (e.g., a HiSeq 2500 or a MiSeq by Illumina). Instrument 303 includes a data
533 acquisition module 305 to obtain results data such as sequence read data. Instrument 303 may
534 optionally include or be operably coupled to its own, e.g., dedicated, analysis computer 333
535 (including an input/output mechanism, one or more processor, and memory). Additionally or
536 alternatively, instrument 303 may be operably coupled to a server 313 or computer 349 (e.g.,
537 laptop, desktop, or tablet) via a network 309.

538 Computer 349 includes one or more processors and memory as well as an input/output
539 mechanism. Where methods of the invention employ a client/server architecture, steps of
540 methods of the invention may be performed using the server 313, which includes one or more of
541 processors and memory, capable of obtaining data, instructions, etc., or providing results via an
542 interface module or providing results as a file. The server 313 may be engaged over the network
543 309 by the computer 349 or the terminal 367, or the server 313 may be directly connected to the
544 terminal 367, which can include one or more processors and memory, as well as an input/output
545 mechanism.

546 In system 301, each computer preferably includes at least one processor coupled to a
547 memory and at least one input/output (I/O) mechanism.

548 A processor will generally include a chip, such as a single core or multi-core chip, to
549 provide a central processing unit (CPU). A process may be provided by a chip from Intel or
550 AMD.

551 Memory can include one or more machine-readable devices on which is stored one or
552 more sets of instructions (e.g., software) which, when executed by the processor(s) of any one of
553 the disclosed computers can accomplish some or all of the methodologies or functions described
554 herein. The software may also reside, completely or at least partially, within the main memory
555 and/or within the processor during execution thereof by the computer system. Preferably, each
556 computer includes a non-transitory memory such as a solid state drive, flash drive, disk drive,
557 hard drive, etc. While the machine-readable devices can in an exemplary embodiment be a single

558 medium, the term “machine-readable device” should be taken to include a single medium or
559 multiple media (e.g., a centralized or distributed database, and/or associated caches and servers)
560 that store the one or more sets of instructions and/or data. These terms shall also be taken to
561 include any medium or media that are capable of storing, encoding, or holding a set of
562 instructions for execution by the machine and that cause the machine to perform any one or more
563 of the methodologies of the present invention. These terms shall accordingly be taken to include,
564 but not be limited to one or more solid-state memories (e.g., subscriber identity module (SIM)
565 card, secure digital card (SD card), micro SD card, or solid-state drive (SSD)), optical and
566 magnetic media, and/or any other tangible storage medium or media.

567 A computer of the invention will generally include one or more I/O device such as, for
568 example, one or more of a video display unit (e.g., a liquid crystal display (LCD) or a cathode
569 ray tube (CRT)), an alphanumeric input device (e.g., a keyboard), a cursor control device (e.g., a
570 mouse), a disk drive unit, a signal generation device (e.g., a speaker), a touchscreen, an
571 accelerometer, a microphone, a cellular radio frequency antenna, and a network interface device,
572 which can be, for example, a network interface card (NIC), Wi-Fi card, or cellular modem.

573 Any of the software can be physically located at various positions, including being
574 distributed such that portions of the functions are implemented at different physical locations.

575 System 301 or components of system 301 may be used to perform methods described
576 herein. Instructions for any method step may be stored in memory and a processor may execute
577 those instructions. System 301 or components of system 301 may be used for the analysis of
578 genomic sequences or sequence reads (e.g., sequence assembly or variant calling).

579 In certain embodiments, as part of the analysis and determination of copy number states
580 and subsequent identification of copy number variation, the sequence read counts for genomic
581 regions of interest are normalized based on internal controls. In particular, an intra-sample
582 normalization is performed to control for variable sequencing depths between samples. The
583 sequence read counts for each genomic region of interest within a sample will be normalized
584 according to the total read count across all control references within the sample.

585 After normalizing read counts for both the genomic regions of interest and control
586 references, copy number states may be determined. In one embodiment, the normalized values
587 for each sample of interest will be compared to the normalized values for a control sample. A
588 ratio, for example, may be generated based on the comparison, wherein the ratio is indicative of

589 copy number and further determinative of any copy number variation. In the event that the
590 determined copy number of a genomic region of interest of a particular sample falls within a
591 tolerable level (as determined by ratio between test and control samples), it can be determined
592 that genomic region of interest does not present copy number variation and thus the patient is at
593 low risk for being a carrier of a condition or disease associated with such. In the event that the
594 determined copy number of a genomic region of interest of a particular sample falls outside of a
595 tolerable level, it can be determined that genomic region of interest does present copy number
596 variation and thus the patient is at risk for being a carrier of a condition or disease associated
597 with such.

598

599

Incorporation by Reference

600 References and citations to other documents, such as patents, patent applications, patent
601 publications, journals, books, papers, web contents, have been made throughout this disclosure.
602 All such documents are hereby incorporated herein by reference in their entirety for all purposes.

603

604

605

606

Equivalents

607 Various modifications of the invention and many further embodiments thereof, in
608 addition to those shown and described herein, will become apparent to those skilled in the art
609 from the full contents of this document, including references to the scientific and patent literature
610 cited herein. The subject matter herein contains important information, exemplification and
611 guidance that can be adapted to the practice of this invention in its various embodiments and
612 equivalents thereof.

613

614

Examples

615 Example 1: Determination of Copy Number State of SMN1

616 Approximately 28 samples are collected to determine carrier status with respect to spinal
617 muscular atrophy (SMA). Genomic DNA is extracted from whole human blood using a Gentra
618 Puregene Blood Kit and following the Puregene protocol for DNA Purification from Whole
619 Blood (Qiagen). Of the 28 samples, there is 1 water negative control and 7 control DNA samples

620 and 20 test samples. Each of the control samples includes two or more genomic regions of
621 interest (e.g. loci) having known (or stable) copy numbers. Control samples 1-4 each include
622 control loci and survival motor neuron genes (SMN), including telomeric SMN (SMN1) and
623 centromeric SMN (SMN2) genes. There are a total of 17 control loci, 5 SMN1, and 5 SMN2, all
624 of which have a known copy number of 2. Control sample 5 includes 17 control loci, each
625 having a known copy number of 2, and 5 SMN1, each having a known copy number of 0.
626 Control sample 6 includes 17 control loci, each having a known copy number of 2, and 5 SMN1,
627 each having a known copy number of 1. Control sample 7 includes 17 control loci, each having a
628 known copy number of 2, and 5 SMN1, each having a known copy number of 3 or more.

629 Samples are processed via method 101 to remove copies of SMN2. The sample is treated
630 (e.g., heated) to denature genomic dsDNA. Primers specific to SMN2 that are complementary to
631 regions flanking the SMN2 sequence are introduced. The primers are annealed to the ssDNA in
632 the regions flanking the unwanted SMN2 segment. The annealed primers are then extended using
633 a polymerase in a template-dependent manner to make double-stranded any single-stranded
634 instance of SMN2 present in any sample. A double-stranded endonuclease is introduced and
635 allowed to digest all dsDNA, thus digesting any segments that include SMN2. This stage of
636 processing of the sample is completed by inactivating the ds endonuclease and the remaining
637 DNA is analyzed for SMN1 by MIP capture and sequencing.

638 The processed samples are then fragmented and/or denatured in preparation for
639 hybridization with molecular inversion probes. The genomic DNA of each sample is
640 fragmented/denatured by any known method or technique sufficient to fragment genomic DNA.

641 Once it is isolated, MIP capture probes are hybridized to the fragmented genomic DNA
642 in each sample by introducing capture probe mix into each sample well. In particular, the capture
643 probe mix will generally include a plurality of SMA molecular inversion probes that are capable
644 of binding to one or more of the genomic regions of interest (e.g., SMN1) or the control DNA. A
645 library of molecular inversion probes is generated. The library may include a variety of different
646 probe configurations. For example, one or more probes are capable of hybridizing specifically to
647 the control loci and one or more probes are capable of hybridizing only to SMN1. Of those
648 probes specific to SMN1, some are capable of producing sequences specific to that paralog while
649 some are not capable of producing paralog-specific sequences. The library may also include one
650 or more probes capable of hybridizing nonspecifically to both SMN1 and SMN2. However, since

651 SMN2 segments are removed from the sample via methods of the invention, copies of SMN2
652 will not interfere with analysis of SMN1.

653 Diluted probes are introduced to the isolated fragmented genomic DNA in each sample
654 and the isolated whole genomic DNA is incubated in the diluted probe mix to promote
655 hybridization. The time and temperature for incubation may be based on any known
656 hybridization protocol, sufficient to result in hybridization of the probes to the DNA. After
657 capture of the genomic region of interest (e.g., SMN1) the captured region is subjected to an
658 enzymatic gap-filling and ligation step, in accordance with any known methods or techniques,
659 including those generally described herein. The captured material may further be purified.

660 The purified captured DNA is then amplified by any known amplification methods or
661 techniques. In one embodiment, the purified captured DNA is amplified using barcode-based
662 PCR. The resulting barcodes PCRs for each sample are then combined into a master pool and
663 quantified.

664 After PCR, portions of the PCR reactions for each sample are pooled and purified, then
665 quantified. In particular, the PCR reactions for all samples are pooled in equal volumes into one
666 master pool. The master sample pool is then purified via a PCR cleanup protocol according to
667 manufacturer's instructions. The purified pool is then run on a microfluidics-based platform for
668 sizing, quantification and quality control of DNA, RNA, proteins and cells. In particular, the
669 purified pool and control samples (pre-purification) are run on an Agilent Bioanalyzer for the
670 detection and quantification of SMN1 probe products.

671 Next, the sample pool is prepared for sequencing. In a preferred embodiment, Illumina
672 sequencing techniques are used. Prior to sequencing, the sample pool is reduced to 2 nM by
673 diluting with 1X TE. Template DNA for cluster generation is prepared by combining 10 micro-
674 Liter of 0.1 N NaOH with 10 micro-Liter of 2 nM DNA library (sample pool) and incubating
675 said mixture at room temperature for 5 min. The mixture is then mixed with 980 micro-Liter of
676 HT1 buffer (Illumina), thereby reducing the denatured library to a concentration of 20 pM. This
677 mixture is then mixed (e.g., inversion) and pulse centrifuged. Next, 225 micro-Liter of the 20 pM
678 library is mixed with 775 micro-Liter of HT1 buffer to reduce the library pool to a concentration
679 of 4.5 pM. The library pool having a concentration of 4.5 pM is used for on-board clustering in
680 the sequencing.

681 The sequencing is carried out on the HiSeq 2500/1500 system sold by Illumina, Inc. (San
682 Diego, CA). Sequencing is carried out with the TruSeq Rapid PE Cluster Kit and TruSeq Rapid
683 SBS 200 cycle kit (Illumina) and in accordance with manufacturer's instructions. In addition to
684 the reagents and mixes included within the kits, additional reagents are prepared for genomic
685 read sequencing primers and reverse barcode sequencing primers.

686 The library pool undergoes sequencing under paired-end, dual-index run conditions.
687 Sequencing generates a plurality of reads. Reads generally include sequences of nucleotide data
688 less than about 150 bases in length, or less than about 90 bases in length. After obtaining
689 sequence reads, they are further processed as described in U.S. Pat. 8,209,130.

690 Read counts for a genomic region of interest are normalized with respect to an internal
691 control DNA. Normalized read counts are compared to the internal control DNA, thereby
692 obtaining a ratio. A copy number state of the genomic region of interest is determined based on
693 the comparison, specifically the ratio.

694 The plurality of reads generated by the sequencing method described above are analyzed
695 to determine copy number states, and ultimately copy number variation, in any of the genomic
696 regions of interest (e.g., SMN1) that would necessarily indicate the presence of an autosomal
697 recessive trait in which copy number variation is diagnostic (e.g., spinal muscular atrophy).
698 Analysis of the read counts is carried out using Illumina's HiSeq BclConverter software. Files
699 (e.g. qSeq files) may be generated for both the genomic and barcode reads. In particular, in
700 accordance with one method of the present invention, genomic read data for each sample is split
701 based upon the barcode reads, which yields separate FASTQ files for each sample.

702 Based on the ratios, loci copy numbers may be called as follows: a ratio of <0.1 will be
703 called a copy number state of 0; a ratio between 0.1 and 0.8 will be called a copy number state of
704 1; a ratio between 0.8 and 1.25 will be called a copy number state of 2; and a ratio of >1.25 will
705 be called a copy number state of 3+.

706 The determined copy numbers can then be used to determine the carrier status of an
707 individual from which the sample was obtained (i.e. whether the patient is a carrier of the
708 disease). In particular, if the copy number state is determined to vary from the normal copy state
709 (e.g., CN is 0, 1 or 3+), it is indicative the condition (e.g., carrier of SMA).

710

711

What is claimed is:

- 712
- 713
- 714
- 715 1. A method of removing unwanted segments of a nucleic acid from a sample, the method
- 716 comprising:
- 717 obtaining a single-stranded nucleic acid that contains an unwanted segment;
- 718 adding complementary nucleic acid to create a double-stranded region that contains the
- 719 unwanted segment; and
- 720 digesting the double-stranded region, thereby removing the unwanted segment from the
- 721 nucleic acid.
- 722
- 723 2. The method of claim 1, wherein adding the complementary nucleic acid comprises annealing
- 724 an oligonucleotide to the unwanted segment, thereby creating the double-stranded region.
- 725
- 726 3. The method of claim 2, wherein the annealing step comprises annealing a plurality of primers
- 727 to a plurality of portions of the nucleic acid that flank an unwanted segment.
- 728
- 729 4. The method of claim 1, wherein adding the complementary nucleic acid comprises:
- 730 annealing an oligonucleotide to a portion of the single-stranded nucleic acid that flanks
- 731 the unwanted segment; and
- 732 extending the annealed oligonucleotide to create the double-stranded region.
- 733
- 734 5. The method of claim 4, wherein the annealing step comprises annealing a plurality of primers
- 735 to a plurality of portions of the nucleic acid that flank an unwanted segment.
- 736
- 737 6. The method of claim 4, wherein the extending step is conducted using a polymerase enzyme
- 738 under conditions sufficient to cause extension of the primer in a template-dependent manner.
- 739
- 740 7. The method of claim 1, wherein the digesting step comprising exposing the sample to an
- 741 enzyme that preferentially digests double-stranded nucleic acid.
- 742

743 8. The method of claim 7, wherein the enzyme is selected from double-stranded endonucleases,
744 restriction endonucleases, and nicking enzymes.

745

746 9. The method of claim 8, further comprising the step of deactivating the enzyme.

747

748 10. The method of claim 1, wherein the nucleic acid is selected from DNA, RNA, and modified
749 nucleic acids.

750

751 11. The method of claim 1, further comprising the step of analyzing nucleic acid remaining after
752 the digesting step.

753

754 12. The method of claim 1, wherein the digesting step results in intact genomic DNA lacking one
755 or more unwanted segment and that is compatible with a nucleic acid analysis assay.

756

757 13. The method of claim 12, wherein the assay comprises molecular inversion probe capture.

758

759 14. The method of claim 13, wherein the assay further comprises sequencing.

760

761 15. The method of claim 14, wherein the sequencing is selected from Sanger sequencing and
762 Next Generation Sequencing.

763

764 16. The method of claim 1, wherein the unwanted segment is a paralog, a pseudogene, or non-
765 paralogous repetitive element.

766

767 17. The method of claim 1, further comprising the step of obtaining a sample from a subject and
768 denaturing double-stranded DNA in the sample.

769

770 18. The method of claim 17, wherein the denaturing step comprises exposing the sample to heat,
771 a detergent, or a basic solution.

772

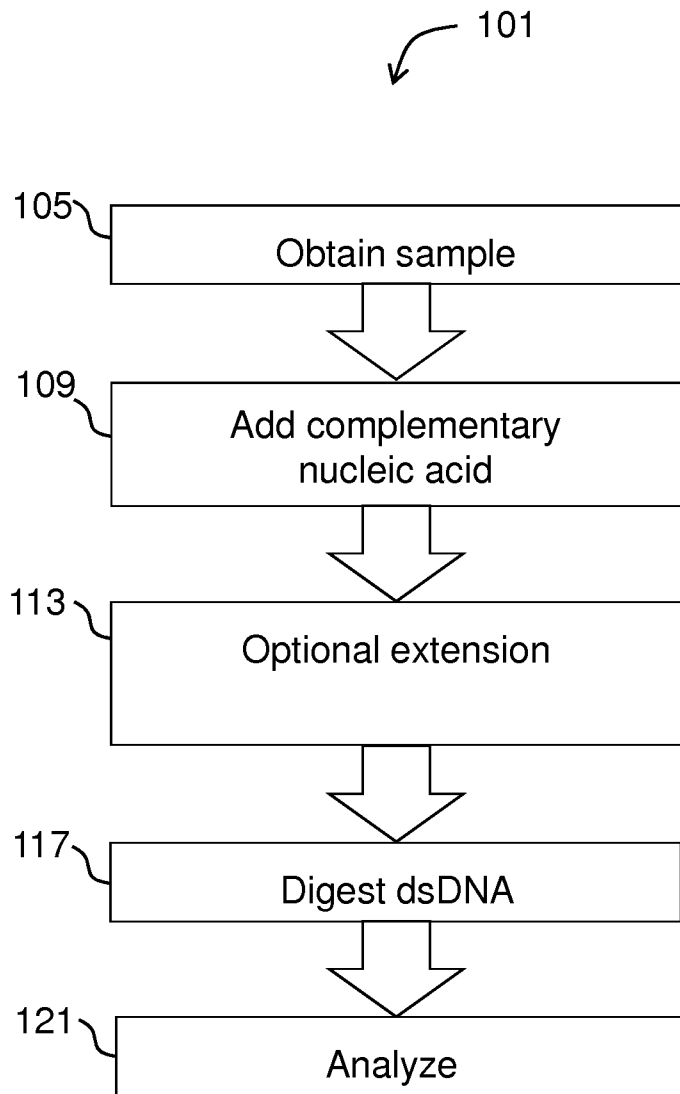


FIG. 1

2/3

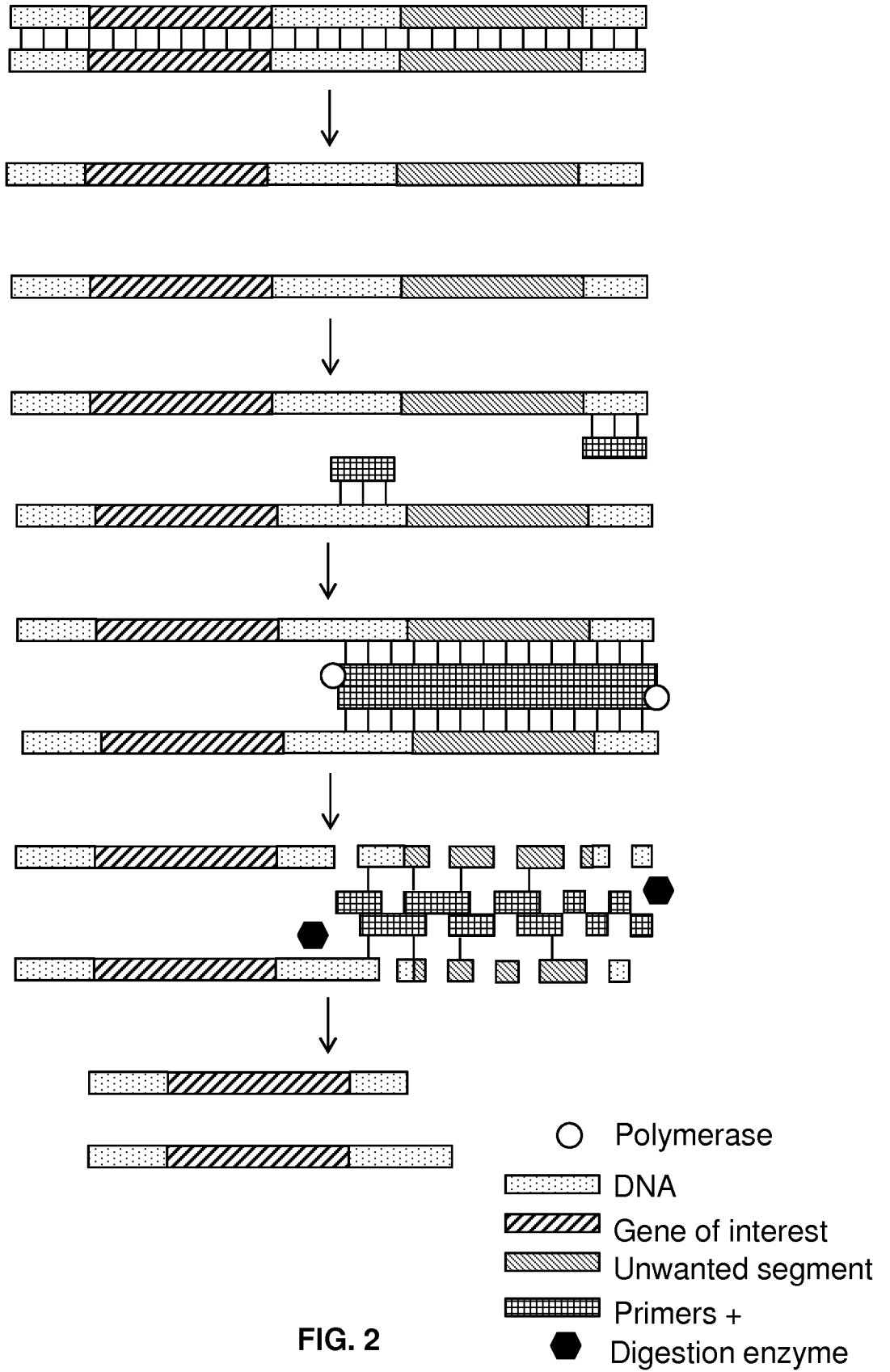


FIG. 2

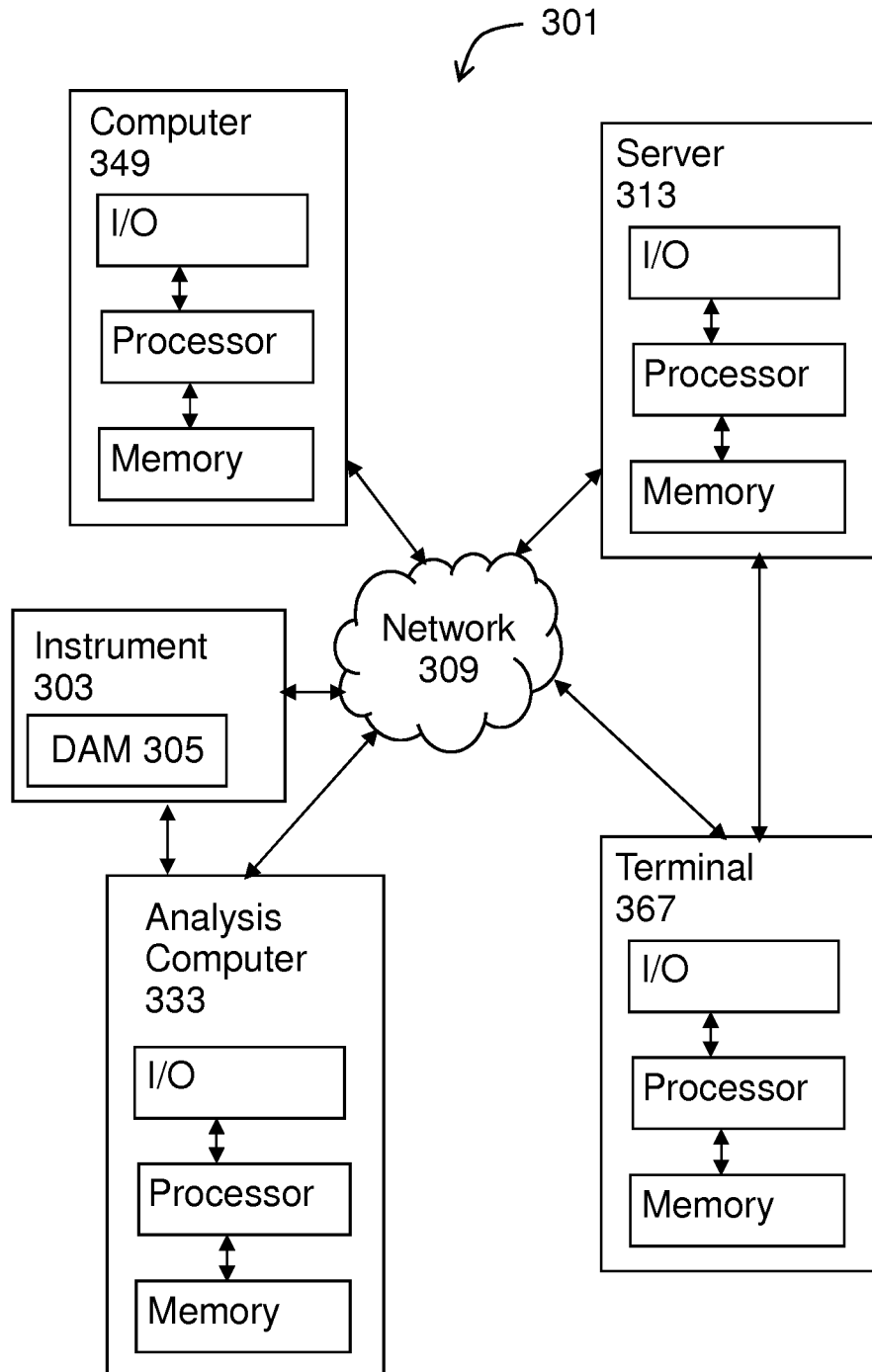


FIG. 3

INTERNATIONAL SEARCH REPORT

International application No PCT/US2015/049132

A. CLASSIFICATION OF SUBJECT MATTER INV. C12Q1/68 ADD.		
According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED Minimum documentation searched (classification system followed by classification symbols) C12Q		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) EPO-Internal, BIOSIS, EMBASE, FSTA, WPI Data		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	S. J. GREEN ET AL: "Suicide Polymerase Endonuclease Restriction, a Novel Technique for Enhancing PCR Amplification of Minor DNA Templates", APPLIED AND ENVIRONMENTAL MICROBIOLOGY, vol. 71, no. 8, 1 August 2005 (2005-08-01), pages 4721-4727, XP055229646, US ISSN: 0099-2240, DOI: 10.1128/AEM.71.8.4721-4727.2005	1-12, 16-18
Y	p. 4721, right-hand col., last para.; p. 4723, right-hand col., last para. - p. 4724, right-hand col., 1st para.; fig. 1, 2 ----- -/--	13-15
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents :		
"A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed	"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family	
Date of the actual completion of the international search	Date of mailing of the international search report	
23 November 2015	02/12/2015	
Name and mailing address of the ISA/ European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Fax: (+31-70) 340-3016	Authorized officer Ripaud, Leslie	

INTERNATIONAL SEARCH REPORT

International application No
PCT/US2015/049132

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	STUART K ARCHER ET AL: "Selective and flexible depletion of problematic sequences from RNA-seq libraries at the cDNA stage", BMC GENOMICS, BIOMED CENTRAL LTD, LONDON, UK, vol. 15, no. 1, 26 May 2014 (2014-05-26), page 401, XP021187323, ISSN: 1471-2164, DOI: 10.1186/1471-2164-15-401	1,2, 7-11,16, 17
Y	p. 2, para. bridging left- to right-hand col.; p. 3, left-hand col., last para. - right-hand col., last para.; fig. 1	13-15
X	US 2005/003369 A1 (CHRISTIANS FREDERICK C [US] ET AL) 6 January 2005 (2005-01-06)	1,2, 7-11, 16-18
Y	para. 5-8, 47, 51-55; fig. 1A	13-15
X	WO 2013/191775 A2 (NUGEN TECHNOLOGIES INC [US]; ARMOUR CHRISTOPHER [US]; AMORESE DOUG [US] 27 December 2013 (2013-12-27) cited in the application	1-12, 16-18
Y	para. 5, 6, 12, 202; fig. 1	13-15
Y	PORRECA GREGORY J ET AL: "Multiplex amplification of large sets of human exons", NATURE METHODS, NATURE PUBLISHING GROUP, GB, vol. 4, no. 11, 1 November 2007 (2007-11-01), pages 931-936, XP002525088, ISSN: 1548-7091, DOI: 10.1038/NMETH1110 cited in the application abstract; fig. 1	13-15
X,P	WO 2015/119941 A2 (IGENOMX INTERNAT GENOMICS CORP [US]) 13 August 2015 (2015-08-13) the whole document	1-18
A	DONNA J. E. HOUSLEY ET AL: "SNP Discovery and Haplotype Analysis in the Segmentally Duplicated DRD5 Coding Region", ANNALS OF HUMAN GENETICS, vol. 73, no. 3, 1 May 2009 (2009-05-01), pages 274-282, XP055230245, GB ISSN: 0003-4800, DOI: 10.1111/j.1469-1809.2009.00513.x abstract	1-18

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No PCT/US2015/049132

Patent document cited in search report		Publication date		Patent family member(s)	Publication date
US 2005003369	A1	06-01-2005	US	2005003369 A1	06-01-2005
			US	2005123987 A1	09-06-2005

WO 2013191775	A2	27-12-2013	CA	2877094 A1	27-12-2013
			CN	104619894 A	13-05-2015
			EP	2861787 A2	22-04-2015
			GB	2518078 A	11-03-2015
			JP	2015521468 A	30-07-2015
			US	2015299767 A1	22-10-2015
			WO	2013191775 A2	27-12-2013

WO 2015119941	A2	13-08-2015	NONE		
