



**HAL**  
open science

## Massive parallel amplicon sequencing of the breast cancer genes BRCA1&2: opportunities, challenges and limitations.

Kim de Leeneer, Jan Hellemans, Joachim de Schrijver, Machteld Baetens, Bruce Poppe, Wim van Criekinge, Anne Depaepe, Paul Coucke, Kathleen B.M. Claes

### ► To cite this version:

Kim de Leeneer, Jan Hellemans, Joachim de Schrijver, Machteld Baetens, Bruce Poppe, et al.. Massive parallel amplicon sequencing of the breast cancer genes BRCA1&2: opportunities, challenges and limitations.. Human Mutation, 2011, 32 (3), pp.335. 10.1002/humu.21428 . hal-00613912

**HAL Id: hal-00613912**

**<https://hal.science/hal-00613912>**

Submitted on 8 Aug 2011

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



**Massive parallel amplicon sequencing of the breast cancer genes *BRCA1&2* : opportunities, challenges and limitations.**

Journal:	<i>Human Mutation</i>
Manuscript ID:	humu-2010-0318.R1
Wiley - Manuscript type:	Methods
Date Submitted by the Author:	29-Oct-2010
Complete List of Authors:	De Leeneer, Kim; Center for Medical Genetics, Ghent University Hospital Hellemans, Jan; Center for Medical Genetics, Ghent University Hospital De Schrijver, Joachim; Laboratory for Bioinformatics and Computational Genomics, Ghent University Baetens, Machteld; Center for Medical Genetics, Ghent University Hospital Poppe, Bruce; Center for Medical Genetics, Ghent University Hospital Van Criekeing, Wim; Laboratory for Bioinformatics and Computational Genomics, Ghent University DePaepe, Anne; Center for Medical Genetics, Ghent University Hospital Coucke, Paul; Center for Medical Genetics, Ghent University Hospital Claes, Kathleen; Ghent University Hospital, Center for Medical Genetics
Key Words:	massive parallel sequencing, BRCA1/2, multiplex, barcoding, amplicon sequencing, breast cancer

SCHOLARONE™  
Manuscripts

1  
2  
3 **Massive parallel amplicon sequencing** of the breast cancer genes *BRCA1&2*:  
4 opportunities, challenges and limitations.  
5  
6

7 **Kim De Leeneer (1), Jan Hellemans (1), Joachim De Schrijver (2), Machteld Baetens (1),**  
8 **Bruce, Poppe (1), Wim Van Criekinge (2), Anne De Paepe (1), Paul Coucke (1),**  
9 **Kathleen Claes (1)\***  
10

11 1 Center for Medical Genetics, Ghent University Hospital, De Pintelaan 185, B-9000 Gent,  
12 Belgium  
13

14 2 Laboratory for Bioinformatics and Computational Genomics, Ghent University, B- 9000  
15 Gent, Belgium  
16  
17

18  
19  
20  
21 Corresponding author: Kathleen Claes, Ph.D.  
22 E-mail address: Kathleen.Claes@UGent.be  
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

1  
2  
3 ABSTRACT  
4  
5  
6  
7  
8

9 This study describes how the new massive parallel sequencing technology can be  
10 implemented in a diagnostic setting by proof of concept studies for the breast cancer  
11 susceptibility genes (*BRCA1&2*). Throughput was maximized by increasing uniformity in  
12 coverage. This was obtained by a multiplex approach, which outperformed pooling of  
13 singleplex PCRs. We evaluated the sensitivity by analysis of 133 distinct sequence variants; 3  
14 (2%) deletions or duplications in homopolymers of  $\geq 7$  nucleotides remained undetected,  
15 illustrating a limitation of pyrosequencing. Furthermore, other limitations like non random  
16 sequencing errors, pseudogene amplification and failure to detect of multi exon deletions are  
17 thoroughly described.  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

29  
30 Our workflow certainly has the potential for high throughput analysis of large genes in  
31 diagnostic settings, which is of great importance to meet the increasing expectations of  
32 genetic testing. Implementation of this approach will hopefully lead to a strong reduction in  
33 turnaround times, so a wider spectrum of at risk women will be able to benefit from  
34 therapeutic interventions for which knowledge of the DNA sequence is essential.  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46

47 Key Words: massive parallel sequencing, *BRCA1/2*, multiplex, barcoding, amplicon  
48 sequencing  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## INTRODUCTION

The development of high throughput or massive parallel sequencing (MPS) has opened many new research opportunities. Several platforms are released, of which the three major are the Genome Sequencer from Roche/454 Life Sciences, the Genome Analyzer from Illumina/Solexa and the SOLiD System from Applied Biosystems. These platforms differ in several ways, such as the technology applied, read length and the number of DNA molecules sequenced.

In the case of 454 sequencing, single DNA strands with 5' and 3' adaptor sequences are attached to beads and then clonally amplified by PCR in an oil-water emulsion. The beads are mixed with DNA polymerase and deposited in plates containing more than 1 million wells, with one bead per well. Nucleotides then flow sequentially over the wells and as each nucleotide is added to form complementary DNA strands, pyrophosphate is released and detected in a chemiluminescent flash (pyrosequencing chemistry).

Since the launch of MPS, an increasing number of applications have been published, amongst others *de novo* sequencing (Pearson, et al., 2007; Pol, et al., 2007; Velasco, et al., 2007), whole genome re-sequencing (Albert, et al., 2007; Korbel, et al., 2007), amplicon sequencing (e.g. exon re-sequencing, virus variant detection, DNA methylation) (Dahl, et al., 2007; Korshunova, et al., 2008; Pettersson, et al., 2008; Taylor, et al., 2007; Thomas, et al., 2006), miRNA and splice variant discovery (Ruby, et al., 2006; Yao, et al., 2007). However, the implementation of this high throughput sequencing in molecular diagnostics remains largely unexplored. Over the last decades, Sanger sequencing (Sanger, et al., 1977) has been the dominant DNA sequencing technology and the “gold standard” for DNA-based mutation

1  
2  
3 detection. However, due to cost limitations direct sequencing of large genes in diagnostics, is  
4  
5 often preceded by a mutation scanning technique, followed by characterization of the variant  
6  
7 with Sanger sequencing. Denaturing gradient gel electrophoresis (DGGE) (van der Hout, et  
8  
9 al., 2006), denaturing High-performance Liquid chromatography (dHPLC) (Liu, et al., 1998)  
10  
11 and High Resolution Melting Curve Analysis (HRMCA) (De Leeneer, et al., 2008; De  
12  
13 Leeneer, et al., 2009) are well known examples. For these pre-screening techniques,  
14  
15 sensitivities varying from 50-100% and specificities close to 100% were reported (Gerhardus,  
16  
17 et al., 2007). Some of these methods are laborious and cannot be fully automated since a  
18  
19 sequencing step is required to define the nature of the variant.  
20  
21  
22  
23  
24

25 We chose the BREast CAncer susceptibility genes **BRCA1 (MIM +113705)** and **BRCA2 (MIM**  
26  
27 **+600185)** to optimize high throughput amplicon sequencing, because of their size,  
28  
29 polymorphic character and lack of mutation hot spots. These genes require high throughput  
30  
31 screening as an increasing number of samples need to be tested within shorter turnaround  
32  
33 times. With the approval of targeted therapeutic agents like PARP inhibitors with selective  
34  
35 toxicity for tumors derived from germline carriers of mutations in *BRCA1* & *2*, expectations for  
36  
37 genetic testing keep increasing (Curtin, 2005). To obtain a cost efficient MPS strategy, an  
38  
39 easy set-up, and uniform distribution of coverage are required to maximize throughput. We  
40  
41 explored the challenges of optimizing the clinical use of MPS with two different approaches  
42  
43 and compared sensitivity and specificity of MPS with Sanger sequencing and prescreening  
44  
45 techniques currently used in molecular diagnostics. For our studies we chose the Roche  
46  
47 Genome Sequencer FLX (GS FLX) system because of the longer read lengths generated by  
48  
49 this instrument compared to the other 2 major platforms.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## MATERIAL AND METHODS

### DNA samples and sequence variants evaluated

In total 123 DNA samples isolated from blood were evaluated. To evaluate the specificity of our MPS set-up genomic DNA samples from 30 patients previously analysed for the complete coding region of *BRCA1/2* with other mutation detection techniques were used. Eleven of these patients were completely Sanger sequenced. Nineteen were analysed with High resolution melting curve analysis (HRMCA) followed by Sanger sequencing of all aberrant melting curves (De Leeneer, et al., 2008; De Leeneer, et al., 2009). For Sanger sequencing and HRMCA identical primers sets were used as for MPS.

Furthermore, 93 samples with previously characterized (Sanger sequencing) deletion or insertion variants were used as positive control samples to validate the detection capacities of MPS. These samples were only sequenced for the amplicon containing the mutation and other amplicons within the specific multiplex set. An overview of all variants evaluated is shown in [Supp. Table 1](#). In total forty-three 1-3 bp deletions, fourteen 1-3 bp insertions (of which 13 duplications), three combined indels and twenty-two deletions and an insertion larger than 3bp were sequenced. Twenty of these deletions and insertions are located in a homopolymeric region longer than 3bp.

Furthermore, 410 (40 unique) nucleotide substitutions were present in all samples analyzed. The majority (406) were frequent SNPs and 4 were nonsense mutations.

### PCR set-up

To cover all coding regions and splice sites, we used primer sets thoroughly validated by HRMCA (De Leeneer, et al., 2008; De Leeneer, et al., 2009) and Sanger sequencing. In run 1&2 we started from equimolar pools of singleplex PCRs, in run 3&4 multiplex PCRs were

1  
2  
3 generated prior to sequencing. To fuse the amplicon-specific primers with the adaptor-MID  
4 (Multiplex Identifiers) barcoded primers, 2 consecutive PCR rounds were used with a  
5 universal M13 tail as linker (Hellemans J et al, under review). A schematic representation of  
6 the principle and workflow of both approaches is shown in Figure 1.  
7  
8  
9  
10  
11

### 12 13 14 15 16 *Run 1 & 2*

17  
18 Per sample 111 singleplex PCRs with amplicon specific primers were performed with  
19 identical reaction conditions and primers as published before (De Leeneer, et al., 2008). PCRs  
20 were performed on the CFX384 (Bio-Rad) and RFU data (endpoint fluorescence) were used  
21 for subsequent normalization to obtain 12 equimolar pools of PCR products per patient.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

### 33 *Run 3 & 4*

34  
35 Sixteen multiplex reactions were optimised to cover the complete coding and splice site  
36 regions of *BRCA1* and *BRCA2*, except *BRCA1* exon 2. As this amplicon amplified less  
37 efficiently compared to the other PCRs, it was optimized as an individual reaction.  
38  
39  
40  
41  
42

43  
44 The specifications of each multiplex reaction are given in **Supplemental Table 2 & 3**.  
45 Multiplex PCR was performed in 20µl total volume. For sets 1,2,5,8-15 the amplification  
46 mixture included 2X Titanium buffer (Clontech), 3% DMSO (VWR International), 200µM of  
47 each dNTP (GE Healthcare), 1x Titanium Taq Polymerase (Clontech) and approximately  
48 100ng of DNA. In 5 reactions (set 3,4,6,7 and 16), 3 mM MgCl<sub>2</sub> (Invitrogen), 2X PCR buffer  
49 (Invitrogen) and 1.5U Platinum Taq polymerase was used instead of Titanium buffer and  
50 polymerase. Two touchdown PCR programs were used (abbreviated as Touch46 and Touch48  
51 in **Supp.** Table 2). The temperature cycling protocol consists of an initial denaturation step at  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 94°C for 2 minutes, followed by 12 cycles of denaturation at 94°C for 20 seconds, annealing  
4  
5 starting at 60 (58) °C for 20 seconds (decreasing 1°C per cycle) and an extension at 72°C  
6  
7 during 1 minute. This initial PCR is followed by 25 additional cycles of denaturation at 94°C  
8  
9 for 40 seconds, annealing at 48 (46) °C during 40 seconds and extension at 72°C for 30  
10  
11 seconds. Final extension was accomplished at 72°C for 10 minutes.  
12  
13

14  
15  
16 Primer concentrations in one multiplex vary between 0.025µM and 0.8µM and were adjusted  
17  
18 to obtain equimolar quantities of each amplicon in one reaction (Supp. Table 3).  
19  
20

#### 21 22 23 24 *PCR with MID barcoded primers*

25  
26  
27 MID barcoded primers consist of (Figure 1A): (i) the required sequencing adaptor (A or B),  
28  
29 (ii) a 10-nucleotide long MID tag or barcode to identify the patient (MID sequences provided  
30  
31 by Roche/454, application note (CRF00104)) (iii) a universal M13-tail (forward primer:  
32  
33 cagcagcttgtaaaacgac and reverse primer: caggaacagctatgacc), identical with the M13 tail  
34  
35 used in the first PCR round.  
36  
37  
38

39  
40  
41 After the initial PCRs (singleplex or multiplex), all samples were diluted 1000 times and 1 µl  
42  
43 product was used as a template for a second PCR with MID barcoded primers. Total volume  
44  
45 of this reaction was 15µl. The amplification mixture included 1.5 mM MgCl<sub>2</sub> (Invitrogen), 1X  
46  
47 PCR buffer (Invitrogen) and 1.5U Platinum Taq (Invitrogen), 3% DMSO (VWR  
48  
49 International), 200µM of each dNTP (GE healthcare) and 0.2 µM of both forward and reverse  
50  
51 primer. Temperature cycling protocol consists out of following steps: 4 minutes at 94°C, 15  
52  
53 cycles of denaturation at 94°C for 30 seconds, annealing at 60°C for 30 seconds, extension at  
54  
55 72°C during 50 seconds, and final extension at 72° for 10 minutes.  
56  
57  
58  
59  
60

1  
2  
3 PCRs were performed on a CFX384 instrument (Bio-Rad). During optimization FAM labeled  
4  
5 MID primers were used to evaluate equimolarity between amplicons within one multiplex  
6  
7  
8 reaction and fluorescent peaks were separated on an ABI3730 capillary system.  
9

### 14 Sequencing runs and data analysis

17  
18 PCRs were normalized and equimolarly pooled in relation to the RFU data. This pool was  
19  
20 purified on a High Pure PCR Cleanup Micro kit (Roche). Fragment length of this total  
21  
22 amplicon pool was evaluated on the Bioanalyzer (Agilent) and compared to the theoretically  
23  
24 predicted pattern (Figure 2).  
25

26  
27  
28 Emulsion PCR and sequencing reactions on the GS-FLX (454- Roche) were performed  
29  
30 according to the manufacturer's instructions. The FASTA files were analysed with in house  
31  
32 developed variant interpretation pipeline (VIP) software version 1.3 (De Schrijver et al.,  
33  
34 2010) and with the commercially available Nextgene software (Softgenetics) version 2.0.  
35  
36 Reads were aligned against GenBank reference sequences, NC\_000017.10 (41196313-  
37  
38 41277467, complement BRCA1) and NC\_000013.10 (32889617-32973809;BRCA2) from  
39  
40  
41 build 18 of the Human Genome assembly.  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## RESULTS

We selected crucial characteristics like uniformity of coverage, sensitivity, specificity and throughput, to evaluate the utility of MPS in diagnostics. Furthermore, MPS should outperform current methods used in molecular diagnostics in terms of cost efficiency to make the implementation worthwhile.

The performance evaluation of the selected approach is based on the results of 4 proof of concept (PoC) experiments (Figure 1). In the first two experiments an identical singleplex approach was used for the analysis of 22 different patient samples.

To reduce workload we started in PoC 3&4 from a multiplex set-up. For PoC4 a strong optimization was performed compared to PoC3: primer concentrations within several multiplex sets were adjusted for amplicons poorly covered in PoC3. Furthermore, the composition of a few multiplexes was changed to obtain more uniform coverage.

### **Evaluation of uniformity in coverage distribution: singleplex vs. multiplex approach**

In a cost efficient test, the full capacity of the GS-FLX instrument is used. This can be achieved by pooling different patients and/or disorders in a single lane. To maximize sample size in a single experiment, a uniform distribution of coverage is required. This means that the difference in number of reads between the “less efficient” amplified fragments and the “best performing” fragments should be as small as possible. We evaluated a singleplex (PoC1&2) and a multiplex approach (PoC3&4). A summary on the number of amplicons sequenced, reads mapped and coverage distribution is shown for each PoC study in Table 1. By calculating the “fold difference to mean coverage”, we showed that starting from strongly optimized multiplex sets results in a more uniform distribution of coverage compared to the

1  
2  
3 equimolar pooling of singleplex sets (Figure 3). The “fold difference to mean coverage” can  
4  
5 be used as “spread correction factor” to calculate the average coverage one needs to aim for to  
6  
7 make sure that also for the less efficiently amplified fragments the minimum required  
8  
9 coverage is obtained. The smaller the fold difference to mean coverage, the larger the number  
10  
11 of samples that can be pooled in a single experiment, the more cost efficient a test will be.  
12  
13 Considering our best optimized multiplex sets (PoC4) we obtained a 3.16-fold difference to  
14  
15 mean coverage for 95% of the amplicons. This value can be used to calculate the average  
16  
17 coverage we need to aim for to obtain a predefined minimum coverage for at least 95% of the  
18  
19 amplicons (see below).  
20  
21  
22  
23  
24  
25  
26  
27

### 28 **Multiplex optimization**

29  
30  
31  
32 Optimizing the multiplex sets by fragment analysis using FAM labeled MID primers, turned  
33  
34 out to be a good strategy: on average, we found a nearly linear correlation between peak  
35  
36 height (=relative fluorescence) on capillary electrophoresis within sets and coverage obtained  
37  
38 after 454 sequencing as shown for 2 multiplex sets in Figure 4. The better results for the  
39  
40 multiplex approach indicate that the second PCR round attaching the MID barcodes,  
41  
42 introduces inequimolarities between the amplicons within the pool of singleplex PCRs. These  
43  
44 inequimolarities are further increased during the emulsion PCR.  
45  
46  
47  
48  
49  
50  
51  
52

### 53 **Calculation of the number of samples that can be pooled in a single run**

54  
55  
56 Knowing the fold difference to mean coverage, allows calculating the number of patients that  
57  
58 can be analyzed in a single standard GS-FLX run (calculation template available in  
59  
60 supplemental files of Hellemans et al. (under review):

- With a fold difference to mean coverage of 3.16, the required average coverage to obtain a minimum coverage of  $38 = 38 \times 3.16 = 120$ . We opted for a threshold of 38-fold coverage based on statistical analyses made by Hellemans et al. who found that 38-fold coverage is required to detect a particular heterozygous variant with a probability of 99.9% when only variants present in at least 25% of the reads are considered as possible true variants. Reasons for the 25% variant frequency: see “Sensitivity”.
- A standard GS-FLX run has 400 000 reads available; based on previous runs we found that the number of reads mapped is  $95\% = 400\,000 \times 95\% = 380\,000$
- 111 amplicons are required to cover the complete coding region and splice sites of *BRCA1&2* and considering a 5% safety margin to correct for possible run errors and differences in MID amplification efficiencies during sequencing:  $111 \times 120 + 5\% = 13986$  reads/patient are required
- therefore,  $380\,000 / 13986 = 27$  patients can be pooled in a standard GS-FLX run, with a maximum of 5% of the amplicons not meeting the 38x coverage threshold.
- With the Titanium chemistry (1,100,000 reads available), the number of patients can be increased to 74.

Figure 5 shows that for PoC4 on average 4 of 111 amplicons (96.4%) did not meet the 38X-treshold. However, based on the capillary electrophoresis panels, higher coverage was expected for these amplicons, indicating that some fragments are less efficiently amplified by the emulsion PCR or had reduced coverage due to experimental variation. Primer concentrations in the multiplex reactions can still be increased for *BRCA2* 11-16, *BRCA2* 11-18, *BRCA2* 11-19 and *BRCA2* 18-1 to improve coverage for these amplicons and may provide a possible solution.

## Sensitivity, specificity and filter settings

### *Sensitivity*

In total 503 (133 distinct) sequence variants, previously identified with Sanger sequencing were evaluated in our sample set (30 patients and 93 control samples) with VIP v1.3. (The sensitivity data obtained with NextGene v2.0 are described below.)

All 40 unique substitutions (SNP's, missense, nonsense, splice site mutations) were easily detected. Detection of deletions or insertions is more challenging. MPS analysis software is based on mapping of single stranded reads on a reference sequence; hence, reads lacking one or more nucleotides or containing an insertion will complicate this process. Furthermore, pyrosequencing has its limitations for correct basecalling in homopolymeric regions (Huse, et al., 2007). Therefore, we specifically selected 93 insertion-deletion (indel) mutations, of which 20 deletions and insertions were present in homopolymeric tracts longer than 3 bp, to thoroughly evaluate the limitations of the technology.

Of the 93 indels, 2 remained completely undetectable and 1 additional mutation was filtered out due to low quality scores (<30) and was present in less than 25 % of the reads. Table 2 shows an overview of all undetected variants and their flanking sequences. All undetected variants affect homopolymer stretches of 7 nucleotides.

In total 130/133 of all unique variants were detected with the VIP software resulting in a sensitivity of 98% (100% for substitutions). In general, sensitivity of MPS will be higher, since we introduced a sample bias by selecting variants in complex sequence regions.

### *Reference bias*

1  
2  
3 A priori, heterozygote variants are assumed to be present in 50% of the reads and  
4  
5 homozygous mutant samples in 100%. On average the heterozygous variants were present in  
6  
7 48.4% (95%CI: 47.7-49.1%; range 27-77%) of the reads. Homozygous mutant variants were  
8  
9 found with an average variant frequency of 99.0% (range: 93-100%). Based on these data we  
10  
11 conclude that the reference bias in mapping is minimal.  
12  
13  
14  
15  
16  
17  
18

### 19 *Specificity*

20  
21 MPS is more sensitive than Sanger sequencing in terms of random sequencing errors and  
22  
23 errors introduced by Taq polymerases, since the technology is based on sequencing of single,  
24  
25 clonally amplified molecules. It is challenging to filter out these errors in the data analysis.  
26  
27  
28

29 The VIP software ((version 1.3)(De Schrijver et al., 2010)) generates a list of variants  
30  
31 detected in more than 10% of mapped reads. By analyzing 30 patient samples for the  
32  
33 complete *BRCA1/2* coding sequence (3330 amplicons), this program generated a list of 5513  
34  
35 variants, of which only 443 are true variants, leaving 5070 false positives. Based on the  
36  
37 analysis of the positive controls, criteria were defined to distinguish false positive from true  
38  
39 variants. Reads defining a variant have to fulfill all of them, before being accepted as a  
40  
41 possible true variant that requires confirmation with Sanger sequencing.  
42  
43  
44  
45  
46

47 To filter out as many false positives as possible the following filters were applied in VIP:  
48  
49

- 50 1. Min. 38X coverage (cov) is required (Filter 1, cov >38x)
- 51
- 52 2. The variant needs to be present in at least 25% of the reads (Filter 2, AF > 25%)
- 53
- 54 3. At least in 1 direction (forward (F) or reverse (R)) a high quality score is required
- 55  
56 (Filter 3 Quality (Q) >30)
- 57  
58  
59  
60

- 1  
2  
3 4. Variants in homopolymer stretches longer than 6 basepairs are non reliable calls  
4  
5 (Filter 4, Homopolymer (Hp) >6)  
6  
7

8  
9 An overview of the results is shown in Supplemental Table 3 & Figure 6. Application of these  
10 filters allowed obtaining a specificity of ~92 %. Specificity for each PoC separately is shown  
11 in Figure 7. Application of filters 1 and 2 resulted in the largest reduction of false positives:  
12 1432 data points were filtered out because of <38X coverage and another 2628 because of  
13 being present in less than 25% of the reads, resulting in a specificity of 70% (filter 1&2). Fine  
14 tuning occurred by applying filters 3 and 4 decreasing the remaining 1010 false positives to  
15 276. In total, 104 distinct false positive variants remained. The recurrence of some false  
16 positives can be attributed to the Taq polymerase used or **pseudogene amplification (Figure**  
17 **8)**. 89% (244/276) of the false positives found are indel variations in the close neighborhood  
18 of a homopolymer region.  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

#### 33 **Data analysis with a commercial available software package: Nextgene (Softgenetics)**

34  
35  
36 **The performance of our in house developed VIP software package was compared with the**  
37 **commercially available software Nextgene version 2.0 (Softgenetics). This software has**  
38 **intrinsic filters in terms of coverage and frequency of a particular variant being present in the**  
39 **reads. An overall variant score is calculated, which provides an empirical estimation of the**  
40 **likelihood that a given SNP is real and not an artifact of sequencing or alignment. This score**  
41 **is mainly based on the concept of Phred scores where quality scores are logarithmically linked**  
42 **to error probabilities. For example a quality score of 10, gives a chance of 1 out of 10 that the**  
43 **base is incorrectly called. Furthermore, sub scores are integrated in this general score, taking**  
44 **into account the allele frequency, forward and reverse balance and a homopolymeric score,**  
45 **which penalizes indels found in homopolymeric regions. The maximum value for this overall**  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 mutation score is 30. Filters and settings were defined based on our data obtained on the  
4  
5 analysis of the positive controls, as we did for VIP.  
6  
7

8  
9 We obtained an overall specificity of 84% (533 false positives in 30 patients), when all the  
10  
11 variants present in 25% of the reads at least 38X coverage were taken into account. This  
12  
13 specificity can be highly increased by applying an additional filter with the overall variant  
14  
15 score.  
16  
17

18  
19 Alle 93 indel variations and 40 substitutions were detected (100%), when no filter was set on  
20  
21 an overall variant score. To improve specificity, we included an overall variant score > 15  
22  
23 requirement (according to software developers recommendations), this resulted in a  
24  
25 specificity of 96%, but also in loss of detection of 3 variants present in homopolymeric  
26  
27 regions (shown in Table 2). *BRCA1* c.1010del was also not detected with the VIP software,  
28  
29 the remaining 2 are different compared to those missed with VIP, but they are also present in  
30  
31 homopolymeric regions >6. Hellemans et al. found that the vast majority of reads for  
32  
33 homopolymers up to the length of 6 can be correctly base called. However, for homopolymers  
34  
35 of 7 nucleotides or more, the number of correctly called reads decreases. Therefore, detection  
36  
37 of mutations in homopolymer tracts of 7 is challenging in in both programs evaluated.  
38  
39  
40  
41  
42  
43  
44  
45

#### 46 *Non random sequencing errors*

47  
48

49 Currently Sanger sequencing is considered as the gold standard. Since MPS involves the  
50  
51 sequencing of single clones, it can be more vulnerable to errors introduced by polymerases.  
52  
53 We did not use proofreading polymerase enzymes (except for the emulsion PCR). Since we  
54  
55 perform 3 PCRs (amplicon specific, MID and emulsion PCR) in our workflow, non specific  
56  
57 errors may occur. Performing our assays with proofreading Taq might lead to a reduction of  
58  
59  
60

1  
2  
3 some false positives. Only 11% (32) of all false positives were single nucleotide substitutions,  
4  
5 most of them can be explained by random sequencing or PCR artefacts and are only seen in a  
6  
7 single direction.  
8  
9

10  
11 For at least one variant we have strong evidence of a non random sequence error caused by  
12  
13 the Taq polymerase: *BRCA2* c.9502-44 G>T (exon 26) was found in about 36% (95% CI:30-  
14  
15 40% range: 19-61%) of all reads of all patients sequenced for this amplicon with Titanium  
16  
17 Taq polymerase. Sanger sequencing of exon 26 and splice sites for all of the patients with  
18  
19 Platinum Taq polymerase revealed only the G allele. Sanger sequencing of samples with  
20  
21 Titanium Taq polymerase clearly showed the G>T substitution in all samples (Figure 9).  
22  
23 These data clearly point at a possible role of the polymerase used for these non-random  
24  
25 sequencing errors.  
26  
27  
28  
29  
30  
31  
32  
33

#### 34 *Evaluation specificity of MID barcodes used*

35  
36

37 No other mutation detection technology allows pooling of several patients in a single lane.  
38  
39 MPS made this possible but the specificity of the MID barcodes used to distinguish individual  
40  
41 patients is crucial. MID tags are designed in such a way that 2 sequencing errors may occur,  
42  
43 without being defined as another MID. We analyzed the data for an experiment containing  
44  
45 only MIDs 1-5 and verified whether reads for MIDs 6-60 were generated. For MID10  
46  
47 approximately 900 reads were mapped (2 mismatches allowed), compared to about 14,000  
48  
49 reads for MIDs included in the experiment. Since these reads are randomly scattered over all  
50  
51 amplicons, chances on generating false positives are minimal, but not impossible for  
52  
53 amplicons with low coverage. Therefore, MID10 should be excluded when 2 mismatches are  
54  
55 allowed for mapping. MID10 is not present in the data when allowing only a single mismatch.  
56  
57  
58  
59  
60

1  
2  
3 Hereby, only a small fraction of the reads (maximum 2000 scattered over all amplicons and  
4 patients) are lost and the risk on false positives is strongly reduced. As this is of major  
5 importance for reliable analyses in diagnostic settings, this solution is preferred.  
6  
7  
8  
9

### 10 11 12 13 14 **Evaluation of multi-exon deletion detection** 15

16  
17 We evaluated the capacity to detect multi-exon deletions in both approaches. In PoC 2, a  
18 patient sample with a heterozygous deletion of *BRCA2* exon 1-4 was included and in PoC 4 a  
19 heterozygous deletion of *BRCA1* exon 18-19 was evaluated. By calculating the dosage  
20 quotient (DQ) (Goossens, et al., 2009), we obtained values not statistically different for the  
21 deleted exons compared to the non-deleted exons. To clearly distinguish a deleted exon, DQ  
22 values of 0.5 are expected. Normalizing the read counts by the average coverage of the  
23 reference patients for these amplicons, resulted in a DQ of 0.35 (*BRCA2* exon 2), DQ of 1  
24 (*BRCA2* exon 3) and 0.72 (*BRCA2* exon 4). For *BRCA1* exon 18 and 19, we found a dosage  
25 quotient of 0.9 for both exons. Therefore, differences in coverage for the deleted exons were  
26 in our experiments not significant and we failed to detect this type of mutations.  
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

## DISCUSSION

In this study we evaluated whether massive parallel amplicon sequencing on the GS-FLX is suitable for implementation in a diagnostic setting. We used a multiplex bar coded amplicon sequencing approach for *BRCA1* and *BRCA2* as an example. For these genes a thoroughly optimized and efficient mutation detection strategy was already previously optimized in our laboratory (De Leeneer, et al., 2008; De Leeneer, et al., 2009). Therefore, we critically evaluated the progress that could be made using amplicon sequencing on the GS-FLX. Different critical aspects such as uniformity of coverage, sensitivity and specificity and reliability of primers and consumables were considered. Our experiences will be very useful for the optimization of other genes.

The recommended approach for amplicon sequencing of multiple sequences is based on the use of fusion primers which start with an A or B sequence on which the pyrosequencing reaction is initiated, followed by a patient specific barcode (MID) and a target specific sequence at the 3' end. Although this fusion approach is very simple, some impractical issues, in terms of primer management and set-up, will arise when the complexity of the experiment (i.e. number of amplicons and number of patients) increases. Primer costs and workload can be strongly reduced by attaching sequencing adaptors and barcodes to the PCR product by ligation (Meyer, et al., 2007) or nested patch PCR (Varley and Mitra, 2008). In our study, a second PCR was used to attach the MID and sequencing adaptors to the amplicon. We preferred this approach because of its simplicity, lower cost and workload compared to some other workflows. Multiplexing PCR products allowed to further reduce workload and consumable cost, and to save patient material. We have chosen to multiplex about 10 amplicons in a single set: optimization of such sets can be obtained with a minimum of extra efforts and results in a reduction of the initial workload by tenfold. Higher degrees of

1  
2  
3 multiplexing would require significantly more efforts and would result in a relatively small  
4 additional increase in efficiency. Furthermore, it would hamper the evaluation of pools by  
5 capillary electrophoresis prior to sequencing, which turned out to be a cost effective  
6 optimization method. Our results showed a good correlation between peak height (=relative  
7 fluorescence) on capillary electrophoresis within sets and coverage obtained after 454  
8 sequencing, indicating a limited influence of the emulsion PCR. We currently have no other  
9 explanation than experimental variation for 4% of the PCR fragments that were less  
10 efficiently amplified by the emulsion PCR. A correlation between coverage and amplicon  
11 length, GC content or other sequence related characteristics could not be found. Although a  
12 homopolymeric tract is present in 2 of these amplicons (*BRCA2* 11.18 and 11.19), this cannot  
13 explain the lower coverage, since the presence of homopolymer stretches did not influence the  
14 amplification efficiency of other homopolymer-rich fragments.

15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35 By analysis of 30 patient samples for the complete *BRCA1/2* coding sequence with the VIP  
36 software, we obtained a list of 5513 variants present in at least 10% of the reads, of which  
37 only 443 are true variants, leaving 5070 false positives (specificity less than 10%). We are the  
38 first group defining filters to cope with MPS sequencing “errors” in a diagnostic setting. With  
39 this program we succeeded to obtain an overall specificity of 92%.

40  
41  
42  
43  
44  
45  
46  
47  
48 Filter 2 (AF >25%) is based on the minimal allele frequency found (27%) by analysis of 93  
49 distinct indel variants and should avoid random PCR errors. The majority of false positives  
50 were generated by sequencing of homopolymeric regions, a known complication for  
51 pyrosequencing, resulting in undercalls and overcalls in homopolymeric stretches (i.e. one  
52 nucleotide missing or one nucleotide added compared to the reference sequence). Filter 3  
53 (Q>30) and 4 (Hp>6) were applied to reduce these numbers, since homopolymeric stretches

1  
2  
3 influenced the Q value assigned to a specific read. These results were compared to those  
4  
5 obtained with the commercially available software Nextgene version 2.0. Application of  
6  
7 filters 1 and 2 (AF >25% and coverage >38 x), resulted in a specificity of 84%. Using an  
8  
9 overall mutation score filter of >15, allowed to increase the specificity to 96% but resulted in  
10  
11 loss of detection of some mutations present in homopolymeric regions of >6 nucleotides.  
12  
13  
14

15  
16 Our results provide evidence that some Taq polymerases introduce non random sequencing  
17  
18 errors. For example *BRCA2* c.9502-44 G>T was found in almost all the patients analysed for  
19  
20 the relevant amplicon when amplified with Titanium Taq polymerase (ClonTech), but not  
21  
22 with Platinum Taq polymerase (Invitrogen). Since these errors will be reproducible in every  
23  
24 run in almost all patients, they can be considered as true false positives and can be ignored in  
25  
26 the long term.  
27  
28  
29

30  
31 Furthermore, the specificity of the barcodes used to identify patients was evaluated to confirm  
32  
33 that false positives could not have been generated by aligning reads to a wrong MID. Despite  
34  
35 the fact that no patients with MID10 were included in one of our experiments, reads for  
36  
37 MID10 were detected when allowing 2 mismatches for the MID. To avoid possible false  
38  
39 positives for amplicons with lower coverage we excluded MID10 from further experiments,  
40  
41 since 10 misaligned reads for a given amplicon are sufficient to generate a variant with at a  
42  
43 frequency of > 25% if only 38 fold coverage is obtained. In future experiments, data will be  
44  
45 mapped allowing only a single mismatch as hereby the fraction of reads lost is minimal and  
46  
47 misalignments to an incorrect MID will be avoided.  
48  
49  
50

51  
52  
53  
54  
55  
56 Homopolymeric stretches turned out to be major sources for false positive and false negative  
57  
58 variants. For an efficient workflow a high specificity is required and our study shows that the  
59  
60

1  
2  
3 number of false positives is strongly reduced by the application of predefined filters.  
4  
5 However, it was impossible to define adequate filters in both software programs evaluated  
6  
7 without loss of detection of some variants in homopolymeric regions longer than 6  
8  
9 nucleotides. The coding region of *BRCA1/2* contains 7 homopolymer regions with 7 and 3  
10  
11 homopolymer regions with 8 nucleotides, in 11 of our amplicons. Until homopolymer  
12  
13 analysis improves, these eleven amplicons are currently analysed with HRMCA in our setting,  
14  
15 increasing detection capacity to 100% for all mutations evaluated (De Leeneer, et al., 2008).  
16  
17 Sequencing technologies not based on pyrosequencing may outperform 454 sequencing for  
18  
19 detection of variants in homopolymer regions.  
20  
21  
22  
23  
24

25 Breast cancer diagnostics was earlier this year evaluated by deep sequencing on the GAIIX  
26  
27 instrument (Illumina) by Morgan, et al. (starting from long-range PCRs) and by Walsh, et al.  
28  
29 (DNA capturing by hybridization in solution to custom-designed cRNA oligonucleotide  
30  
31 baits). Both reported detection of all variants/mutations evaluated, but the number of variants  
32  
33 evaluated was much smaller and insertions/deletions were maximum 19bp in length. The  
34  
35 deletion c.5503\_5564del62 evaluated in our study may have remained undetected with this  
36  
37 sequencing technology with read lengths of  $2 \times 76$ -bp paired-end reads (Walsh, et al.) or 51  
38  
39 bp. (Morgan, et al.).  
40  
41  
42  
43  
44

45 Walsh, et al. did not report any false positives after filtering out variants present in less than  
46  
47 15% of the reads. With an average coverage of 1286 (range 781-1854) a reduction in false  
48  
49 positives is indeed expected, however, such high coverage largely increases the cost per  
50  
51 sample. We calculated consumable costs and labor time for our MPS approach, and found that  
52  
53 our MPS set-up by pooling 74 patients in a single run costs about 345 EUR per sample  
54  
55 (consumable cost: 232 EUR), which becomes cost competitive with our HRMCA approach  
56  
57 and is much smaller than Sanger sequencing. The current development of commercial user-  
58  
59  
60

1  
2  
3 friendly software for data analysis, allows MPS to outperform prescreening techniques used in  
4  
5 combination with Sanger sequencing. Automation of the workflow will even further decrease  
6  
7  
8 the workload.  
9

10 Turnaround times for genetic testing need to be strongly reduced to meet increasing  
11  
12 expectations. Pooling large numbers of patients in a single run will only be useful in a  
13  
14 diagnostic setting if different disorders can be pooled in a single run. This requires the  
15  
16 development of uniform workflows for different genetic tests.  
17

18  
19 In some populations, large intragenic founder deletions represent an important fraction of the  
20  
21 *BRCA1/2* mutation spectrum. Our MPS set-up allowed to detect point mutations with high  
22  
23 sensitivity but turned out to be unreliable for the identification of large exon (or multi-exon)  
24  
25 deletions. The application of 3 consecutive PCR rounds prior to sequencing most likely  
26  
27 explains why deviations from diploidy remained undetected. Additionally we aimed for an  
28  
29 average coverage of 120 (to obtain minimally 38), to allow a cost efficient test by pooling of a  
30  
31 large number of samples in a single lane. Studies successfully reporting the detection of large  
32  
33 intragenic rearrangements worked with much higher read depths, allowing more reliable  
34  
35 quantifications of copy numbers (Goossens, et al., 2009; Walsh, et al.) For a sensitive analysis  
36  
37 an additional technique for the detection of copy number variations needs to be included in  
38  
39 the mutation detection strategy.  
40  
41  
42  
43  
44  
45

46 In conclusion, we developed an efficient workflow for high throughput *BRCA1/2* amplicon  
47  
48 sequencing. Sensitivity and specificity of MPS amplicon sequencing is high and can be  
49  
50 further increased by supplementing MPS assays to overcome issues related to homopolymeric  
51  
52 regions. In terms of throughput, diagnostic testing can be highly accelerated and MPS  
53  
54 facilitates offering genetic analyses to more at-risk patients. Considering cost efficiency MPS  
55  
56 outperforms all other mutation screening techniques, but there is a shift from “wet” lab work  
57  
58  
59  
60



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

towards data analysis. In our opinion, Sanger sequencing should still be used for confirmation of deleterious variants in diagnostics.

For Peer Review

## ACKNOWLEDGEMENTS

This project was realized with the funding of an Emmanuel van der Schueren scholarship of the Flemish foundation against cancer to Kim De Leeneer. This research was supported by grant 1.5.150.07 from the Fund for Scientific Research Flanders (FWO) to Kathleen Claes and by GOA grant BOF10/GOA/019 (Ghent University). Bruce Poppe is a senior clinical investigator from FWO. This study was supported by a StepStone grant from the Industrial Research Fund from Ghent University; the Roche 454 GS-FLX instrument is part of the NXTGNT infrastructure, funded by a Hercules grant (middle heavy infrastructure). NXTGNT (initiated and supervised by Sofie Bekaert, Jo Vandesompele, Dieter Deforce, Philippe van Nieuwerburgh, Wim Van Criekinge, Jan Hellemans, Paul Coucke) is a genome analysis platform from Ghent University.

## REFERENCES

- 1  
2  
3  
4  
5  
6  
7 Albert TJ, Molla MN, Muzny DM, Nazareth L, Wheeler D, Song X, Richmond TA, Middle CM, Rodesch  
8 MJ, Packard CJ and others. 2007. Direct selection of human genomic loci by microarray  
9 hybridization. *Nat Methods* 4(11):903-5.
- 10 Curtin NJ. 2005. PARP inhibitors for cancer therapy. *Expert Rev Mol Med* 7(4):1-20.
- 11 Dahl F, Stenberg J, Fredriksson S, Welch K, Zhang M, Nilsson M, Bicknell D, Bodmer WF, Davis RW, Ji  
12 H. 2007. Multigene amplification and massively parallel sequencing for cancer mutation  
13 discovery. *Proc Natl Acad Sci U S A* 104(22):9387-92.
- 14 De Leeneer K, Coene I, Poppe B, De Paepe A, Claes K. 2008. Rapid and sensitive detection of BRCA1/2  
15 mutations in a diagnostic setting: comparison of two high-resolution melting platforms. *Clin  
16 Chem* 54(6):982-9.
- 17 De Leeneer K, Coene I, Poppe B, De Paepe A, Claes K. 2009. Genotyping of frequent BRCA1/2 SNPs  
18 with unlabeled probes: a supplement to HRMCA mutation scanning, allowing the strong  
19 reduction of sequencing burden. *J Mol Diagn* 11(5):415-9.
- 20 De Schrijver JM, De Leeneer K, Lefever S, Sabbe N, Pattyn F, Van Nieuwerburgh F, Coucke P, Deforce  
21 D, Vandesompele J, Bekaert S and others. Analysing 454 amplicon resequencing experiments  
22 using the modular and database oriented Variant Identification Pipeline. *BMC Bioinformatics*  
23 11:269.
- 24 Gerhardus A, Schleberger H, Schlegelberger B, Gadzicki D. 2007. Diagnostic accuracy of methods for  
25 the detection of BRCA1 and BRCA2 mutations: a systematic review. *Eur J Hum Genet*  
26 15(6):619-27.
- 27 Goossens D, Moens LN, Nelis E, Lenaerts AS, Glassee W, Kalbe A, Frey B, Kopal G, De Jonghe P, De  
28 Rijk P and others. 2009. Simultaneous mutation and copy number variation (CNV) detection  
29 by multiplex PCR-based GS-FLX sequencing. *Hum Mutat* 30(3):472-6.
- 30 Hellemans J, De Leeneer K, De Schrijver J, Clemente L, Baetens M, Lefever S, De Keulenaer S, Claes K,  
31 Pattyn F, De Wilde B, Coucke P, Vandesompele J. . Massively parallel sequencing of PCR  
32 products, the road to next generation molecular diagnostics. Under review.
- 33 Huse SM, Huber JA, Morrison HG, Sogin ML, Welch DM. 2007. Accuracy and quality of massively  
34 parallel DNA pyrosequencing. *Genome Biol* 8(7):R143.
- 35 Korb J, Urban AE, Affourtit JP, Godwin B, Grubert F, Simons JF, Kim PM, Palejev D, Carriero NJ, Du  
36 L and others. 2007. Paired-end mapping reveals extensive structural variation in the human  
37 genome. *Science* 318(5849):420-6.
- 38 Korshunova Y, Maloney RK, Lakey N, Citek RW, Bacher B, Budiman A, Ordway JM, McCombie WR,  
39 Leon J, Jeddloh JA and others. 2008. Massively parallel bisulphite pyrosequencing reveals  
40 the molecular complexity of breast cancer-associated cytosine-methylation patterns  
41 obtained from tissue and serum DNA. *Genome Res* 18(1):19-29.
- 42 Liu W, Smith DI, Reztzigel KJ, Thibodeau SN, James CD. 1998. Denaturing high performance liquid  
43 chromatography (DHPLC) used in the detection of germline and somatic mutations. *Nucleic  
44 Acids Res* 26(6):1396-400.
- 45 Meyer M, Stenzel U, Myles S, Prufer K, Hofreiter M. 2007. Targeted high-throughput sequencing of  
46 tagged nucleic acid samples. *Nucleic Acids Res* 35(15):e97.
- 47 Morgan JE, Carr IM, Sheridan E, Chu CE, Hayward B, Camm N, Lindsay HA, Mattocks CJ, Markham AF,  
48 Bonthron DT and others. Genetic diagnosis of familial breast cancer using clonal sequencing.  
49 *Hum Mutat* 31(4):484-91.
- 50 Pearson BM, Gaskin DJ, Segers RP, Wells JM, Nuijten PJ, van Vliet AH. 2007. The complete genome  
51 sequence of *Campylobacter jejuni* strain 81116 (NCTC11828). *J Bacteriol* 189(22):8402-3.
- 52 Pettersson E, Zajac P, Stahl PL, Jacobsson JA, Fredriksson R, Marcus C, Schioth HB, Lundeberg J,  
53 Ahmadian A. 2008. Allelotyping by massively parallel pyrosequencing of SNP-carrying  
54 trinucleotide threads. *Hum Mutat* 29(2):323-9.
- 55  
56  
57  
58  
59  
60

- 1  
2  
3 Pol A, Heijmans K, Harhangi HR, Tedesco D, Jetten MS, Op den Camp HJ. 2007. Methanotrophy below  
4 pH 1 by a new *Verrucomicrobia* species. *Nature* 450(7171):874-8.  
5  
6 Ruby JG, Jan C, Player C, Axtell MJ, Lee W, Nusbaum C, Ge H, Bartel DP. 2006. Large-scale sequencing  
7 reveals 21U-RNAs and additional microRNAs and endogenous siRNAs in *C. elegans*. *Cell*  
8 127(6):1193-207.  
9  
10 Sanger F, Nicklen S, Coulson AR. 1977. DNA sequencing with chain-terminating inhibitors. *Proc Natl*  
11 *Acad Sci U S A* 74(12):5463-7.  
12  
13 Taylor KH, Kramer RS, Davis JW, Guo J, Duff DJ, Xu D, Caldwell CW, Shi H. 2007. Ultradeep bisulfite  
14 sequencing analysis of DNA methylation patterns in multiple gene promoters by 454  
15 sequencing. *Cancer Res* 67(18):8511-8.  
16  
17 Thomas RK, Nickerson E, Simons JF, Janne PA, Tengs T, Yuza Y, Garraway LA, LaFramboise T, Lee JC,  
18 Shah K and others. 2006. Sensitive mutation detection in heterogeneous cancer specimens  
19 by massively parallel picoliter reactor sequencing. *Nat Med* 12(7):852-5.  
20  
21 van der Hout AH, van den Ouweland AM, van der Lijjt RB, Gille HJ, Bodmer D, Bruggenwirth H,  
22 Mulder IM, van der Vlies P, Elfferich P, Huisman MT and others. 2006. A DGGE system for  
23 comprehensive mutation screening of BRCA1 and BRCA2: application in a Dutch cancer clinic  
24 setting. *Hum Mutat* 27(7):654-66.  
25  
26 Varley KE, Mitra RD. 2008. Nested Patch PCR enables highly multiplexed mutation discovery in  
27 candidate genes. *Genome Res* 18(11):1844-50.  
28  
29 Velasco R, Zharkikh A, Troggio M, Cartwright DA, Cestaro A, Pruss D, Pindo M, Fitzgerald LM, Vezzulli  
30 S, Reid J and others. 2007. A high quality draft consensus sequence of the genome of a  
31 heterozygous grapevine variety. *PLoS One* 2(12):e1326.  
32  
33 Walsh T, Lee MK, Casadei S, Thornton AM, Stray SM, Pennil C, Nord AS, Mandell JB, Swisher EM, King  
34 MC. Detection of inherited mutations for breast and ovarian cancer using genomic capture  
35 and massively parallel sequencing. *Proc Natl Acad Sci U S A* 107(28):12629-33.  
36  
37 Yao Y, Guo G, Ni Z, Sunkar R, Du J, Zhu JK, Sun Q. 2007. Cloning and characterization of microRNAs  
38 from wheat (*Triticum aestivum* L.). *Genome Biol* 8(6):R96.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## FIGURE LEGENDS

**Figure 1: Schematic overview of the two set-up and different workflows used in the 4 PoC studies.**

**Panel A:** Amplicon specific primers are fused with an universal M13tail. This target specific PCR product is amplified in a second PCR with primers, consisting of an M13tail, MID-tag (patient specific) and at the end an A or B adaptor

**Panel B** Upper panel: Schematic representation of the approach used in PoC 1 & 2. 111 PCR reactions per patient are performed and equimolarly pooled in 12 pools, followed by a second PCR to attach MID primers and sequencing adaptors. All products are equimolarly pooled prior to emulsion PCR and sequencing.

Lower panel: Approach used in PoC 3 & 4; 16 multiplex reactions were optimized containing the 111 amplicons to be amplified for each patient. By a second PCR round MID barcode and sequencing adaptors were attached. PCR products are equimolarly pooled prior to emulsion PCR and sequencing.

**Figure 2: Comparison of fragment length of the total amplicon pool in theory and in practice.**

**Left panel:** Theoretical profile of the amplicon length of the total amplicon pool for one patient when all amplicons are equimolarly pooled together. Amplicon length range is 233bp to 437 bp, with an average amplicon length of 335bp.

**Right panel:** Amplicon length profile obtained with the Bioanalyzer (Agilent) after column purification of the pool. Comparison of the two patterns is used as a quality control prior to emulsion PCR.

**Figure 3: Evolution of distribution of coverage uniformity in the 4 PoC experiments.**

In total we performed 4 GS-FLX runs. In our first two experiments (green and orange curve), equimolar pooling of simplex PCR reactions was used. The third (blue) and the fourth (purple) experiment were prepared according to the multiplexing protocol.

**Panel A:** A distribution plot of the coverage is shown. The experiments with the multiplex approach (PoC 3&4) clearly show a more uniform distribution of coverage compared to the pooled simplex runs (PoC 1&2). 38 fold coverage threshold is depicted by the dotted line. Failure rate in PoC 3 was higher because of poor amplification of some multiplex sets.

**Panel B:** Fold difference to the mean coverage is shown in function of the fraction of amplicons. The dotted lines depict the factor where 90 or 95% fraction of all amplicons is

1  
2  
3 considered. Our last experiment (PoC4-purple curve) clearly shows the best result where the  
4 smallest difference to mean coverage was obtained. 95% of all amplicons are sequenced with  
5 a spread correction factor of 3.16 (Table 1) (X-axis is shown in log<sub>10</sub> scale).  
6  
7  
8  
9

#### 10 **Figure 4: Correlation between relative fluorescence and coverage within multiplex sets**

11  
12 The correlation between the relative fluorescence seen on capillary electrophoresis and  
13 coverage is shown.  
14

15  
16 Equimolarity of amplicons within a multiplex set was verified on capillary electrophoresis.  
17  
18  
19

#### 20 **Figure 5: Distribution of coverage for each amplicon in BRCA1/2 (PoC 4)**

21  
22 Coverage for each BRCA1 and BRCA2 amplicon is shown, with the line indicating the  
23 threshold of 38 fold coverage. In total 111 amplicons are needed to cover the complete coding  
24 region of both genes. In BRCA2, 4 amplicons did on average not reach the 38x coverage  
25 threshold; further optimisation for these amplicons is required.  
26  
27  
28  
29  
30

#### 31 **Figure 6: Overview of data generated**

32  
33 Data were analysed with an in house developed software program (VIP). Combined forward  
34 and reverse allele frequencies are plotted against coverage. Variants detected in 30 patients  
35 and 93 positive control samples are shown. Grey data points are the false positives filtered out  
36 when filters are applied. Green data points are true variants (503). The green outlier at almost  
37 80% AF is deletion of 62 nucleotides, probably preferential amplification of the shorter allele  
38 has occurred. 80% of the true variants have an allele frequency in the 40-60 range. Red data  
39 points are the remaining false positives (276).  
40  
41  
42  
43  
44  
45

#### 46 **Figure 7: Specificity of 4 PoC studies analysed with VIP software**

47  
48 Specificity (%) of each PoC study and in total is plotted, for every filter applied in the VIP  
49 software. Filter 1 and 2 clearly result in the largest increase in specificity. The total specificity  
50 found is ~92%.  
51  
52  
53  
54

#### 55 **Figure 8: Pseudogene amplification by MPS**

56  
57 A part of the sequence of BRCA1 exon 2 is shown in both panels. In BRCA1 exon 2 for each  
58 sample 10-15 variants (in 25-50% of the reads, blue boxes in upper panel) were detected  
59 with MSP, but were not observed Sanger sequencing (lower panel, blue boxes represent the  
60

1  
2  
3 possible pseudogene nucleotides). Blasting of these sequences, showed co-amplification of  
4 BRCA1P1 (90% analogy with BRCA1 exon 2) by the reverse primer.  
5  
6  
7  
8

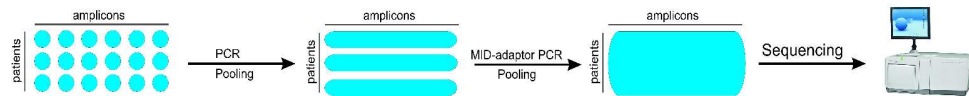
9 **Figure 9: Sanger sequencing results for *BRCA2* c.9502-44G>T**

10  
11 In the upper panel the *BRCA2* a Sanger sequence flanking the G allele at position c.9502-44 is  
12 shown for a sample amplified with Platinum Taq polymerase. In the lower panel, the same  
13 Sanger sequence for the same sample is shown after amplification with Titanium Taq  
14 polymerase. In this sample we see a clear reduction of the G allele and the replacement by a T  
15 allele.  
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

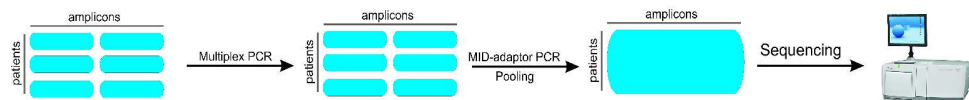
For Peer Review

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

PoC experiments 1 & 2



PoC experiments 3 & 4

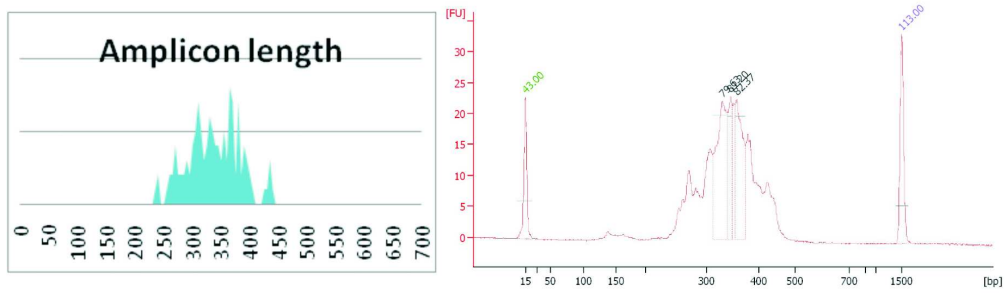


678x215mm (300 x 300 DPI)

Or Peer Review

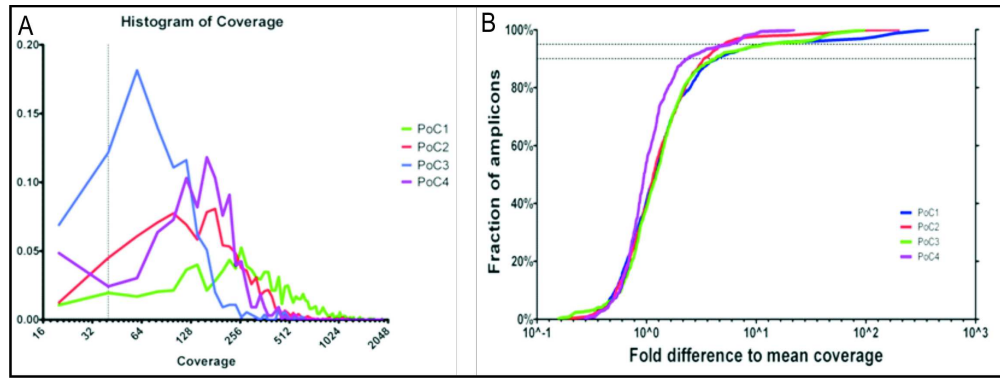


1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



389x122mm (600 x 600 DPI)

Peer Review

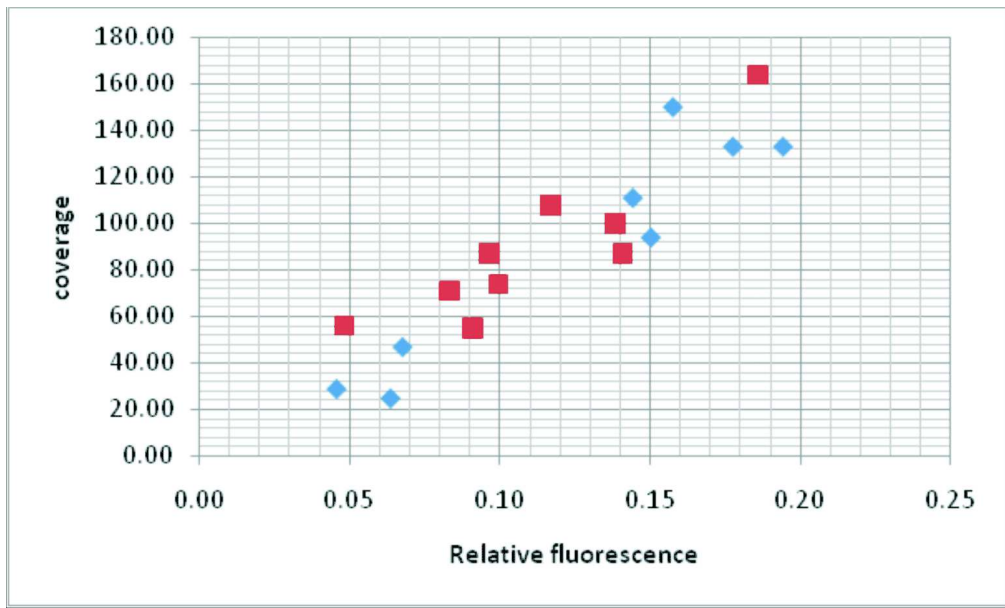


195x72mm (300 x 300 DPI)

Peer Review

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

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



130x78mm (300 x 300 DPI)

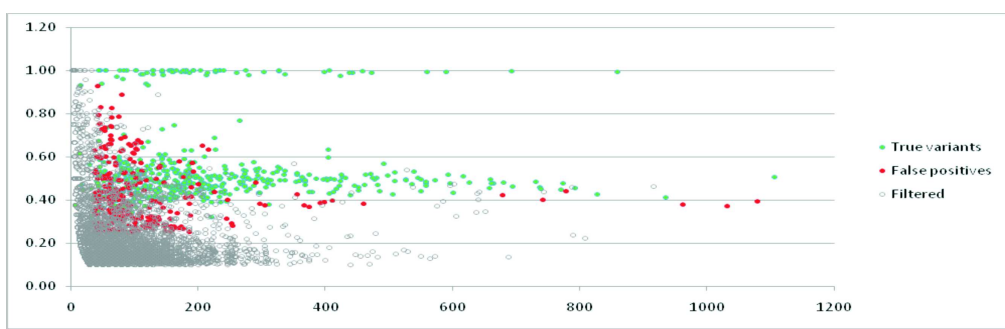
Review

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



299x182mm (300 x 300 DPI)

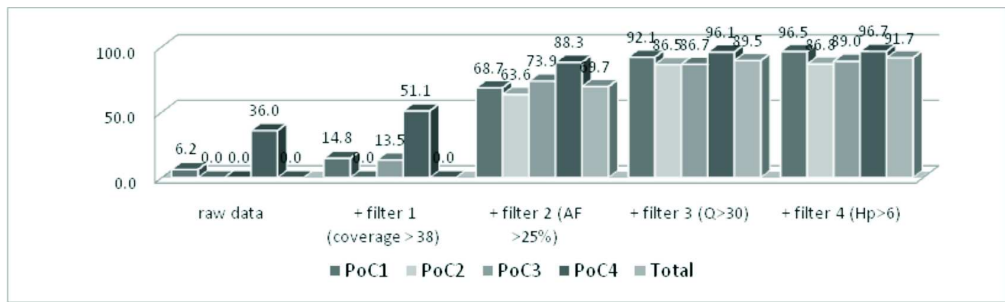
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



225x71mm (300 x 300 DPI)

ur Peer Review

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

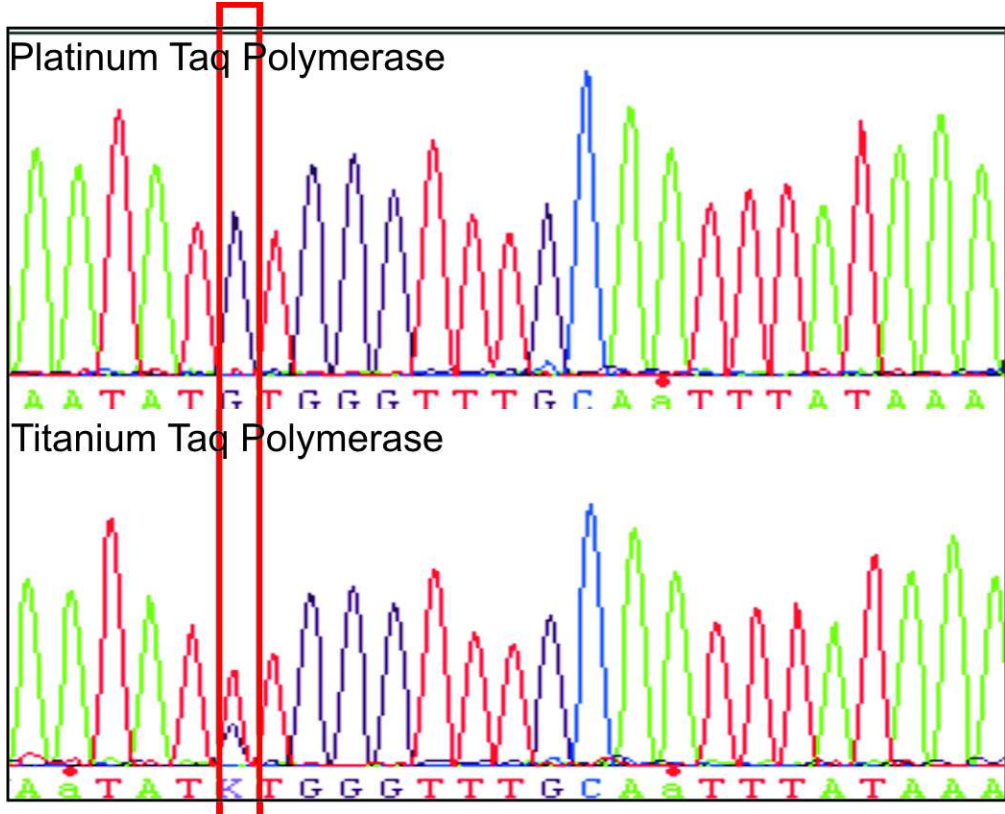


286x84mm (300 x 300 DPI)

or Peer Review



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



*BRCA2* c.9502-44 G>T

88x77mm (300 x 300 DPI)

ew



Table 1: Overview of characteristics of the 4 PoC experiments

	Singleplex pooling		Multiplex pooling	
	PoC 1	PoC 2	PoC 3	PoC 4
<b>SET UP</b>				
used run reads	100%	100%	~30%*	~25%*
mapped BRCA1/2 reads	515,916	250,884	78,937	55,574
amplicons sequenced	1221	1221	555	333
patients sequenced	11	11	5	3
<b>COVERAGE</b>				
min/average/max (per amplicon)	1/348/1870	1/191/1076	1/144/931	1/168/559
standard deviation	264	140	122	90
Variation coefficient	0.76	0.73	0.84	0.53
Fold difference to mean coverage 90%/95%	4.48/11.86	3.38/5.11	4.2/12.08	2.15/3.16
# amplicons <38 fold coverage (%)	107 (8.76)	84 (6.88)	54 (9.73)	21 (6.31)

\* only 30% and 25% of the reads were used in these experiments, since in the same runs pooling different disorders in one experiment was evaluated

Table 2: Examples of variants analysed in homopolymeric regions  $\geq 5$ 

Variant (c.)	flanking sequences	Change of homopolymer length	coverage of nucleotide	Quality	VF (%)
<b>Results with VIP software</b>					
<i>Undetected variants</i>					
<i>BRCA1</i> c.1016dup	ACTCCAGCACAG <u>AAAAAAAA</u> AGGTAGATCTGAATGCTGATC	7→8	38	n.a	<10%
<i>BRCA2</i> c.994del	CTAGMAAGACTAGG <u>AAAAAAAA</u> TTTTCCATGARGCAAACGCT	7→6	219	n.a	<10%
<i>BRCA1</i> c.1010del	ACTCCAGCACAG <u>AAAAAAAA</u> AGGTAGATCTGAATGCTGATCC	7→6	39	25	23%
<i>Detected variants in homopolymeric regions</i>					
<i>BRCA1</i> c.3329dup	ATCCTGAAATA <u>AAAAAAAA</u> AGCAAGAATATGAAGAAGTAGTTC	6→7	99	31	47%
<i>BRCA2</i> c.5577_5580del	ACATGAAACA <u>TTAAAA</u> AGTGAAAGACATATTTACAGAC	4→6	105	34	50%
<i>BRCA1</i> c.2989_2990dup	AAAACATAATGTAAG <u>AAAAA</u> ATCTGCTAGAGGAAAACCTT	5→7	167	32	40%
<b>Results with Nextgene software</b>					
<i>Examples of variants with low mutation score</i>					
<i>BRCA2</i> c.6351dup	GAAGATCA <u>AAAAAAAA</u> CACTAGTTTT	7→8	221	10	33%
<i>BRCA1</i> c.1010del	ACTCCAGCACAG <u>AAAAAAAA</u> AGGTAGATCTGAATGCTGATCC	7→6	20	12	50%
<i>BRCA1</i> c.1961delA	AGAGATAAAG <u>AAAAAAAA</u> AGTACAACCAA	8→7	76	11	52%

## Supp T1: Overview of all sequence variants evaluated on NGS.

<b>substitutions</b>		
<i>DNA level systematic nomenclature (BIC)</i>	<i>Protein level</i>	<i>Suggested classification</i>
<b>BRCA1</b>		
c.134+3A>C (IVS3+3A>C)	non-coding	Mutation
c.1067A>G (1186A>G )	p.Gln356Arg	Polymorphism
c.1456 T>C (1575T>C )	p.Phe486Leu	Polymorphism
c.1487G>A (1606G>A)	p.Arg496His	Polymorphism
c.2077G>A (2196G>A)	p.Asp693Asn	Polymorphism
c.2082C>T (2201C>T)	p.Ser694Ser	Polymorphism
c.2311T>C (2430T>C)	p.Leu771Leu	Polymorphism
c.2612C>T (2731C>T)	p.Pro871Leu	Polymorphism
c.3113A>G (3232A>G)	p.Glu1038Gly	Polymorphism
c.3179A>C (3298A>C)	p.Glu1060Ala	Polymorphism
c.3548A>G (3667A>G)	p.Lys1183Arg	Polymorphism
c.3661G>T (3780G>T)	p.Glu1221X	Mutation
c.3841C>T (3960C>T)	p.Gln1281X	Mutation
c.4308T>C (4427T>C )	p.Ser1436Ser	Polymorphism
c.4867A>G (4956A>G)	p.Ser1613Gly	Polymorphism
c.4956G>A (5075G>A)	p.Met1652Ile	Polymorphism
c.4987-53 C>T (IVS17-53C>T)	non-coding	Polymorphism
c.536A>G (655A>G)	p.Val179Cys	Polymorphism
<b>BRCA2</b>		
c.10110G>A (10338G>A )	p.Arg3370Arg	polymorphism
c.1113C>A (1342C>A)	p.His372Asn	polymorphism

1			
2			
3			
4	c.125A>G (353A>G)	p.Tyr42Cys	polymorphism
5	c.1-25G>A (203G>A)	non-coding	polymorphism
6	c.1365A>G (1593A>G)	p.Ser455Ser	polymorphism
7	c.2229 T>C (2457T>C)	p.His743His	polymorphism
8	c.2971A>G (3199A>G)	p.Asn991Asp	polymorphism
9			
10	c.3396A>G (3624A>G)	p.Lys1132Lys	polymorphism
11	c.3807T>C (4035T>C)	p.Val1269Val	polymorphism
12	c.3851G>A (4079G>A)	p.Ser1284Asn	polymorphism
13			
14	c.516+1G>A (IVS6+1G>A)	non-coding	Mutation
15	c.5645C>A (5873C>A)	p.Ser1882X	Mutation
16	c.5744C>T (5972C>T)	p.Thr1915Met	polymorphism
17	c.68+62T>G (IVS2+62T>G)	non-coding	polymorphism
18	c.681+56C>T (IVS8+56C>T)	non-coding	polymorphism
19			
20	c.7057G>C (7285G>C)	p.Gly2353Arg	polymorphism
21	c.7242A>G (7470A>G)	p.Ser2414Ser	polymorphism
22			
23	c.7806-14T>C (IVS16-14T>C)	non-coding	polymorphism
24	c.8182G>A (8410G>A)	p.Val2728Ile	polymorphism
25	c.865A>C 1093A>C	p.Asn289His	polymorphism
26			
27	c.9257-16T>C (IVS24-16T>C)	non-coding	polymorphism
28	c.9976A>T (10204A>T)	p.Lys3326X	polymorphism
29			

### indels

30			
31			
32	<b>BRCA1</b>		
33	c.562-58delT (IVS8-58delT)	non-coding	Polymorphism
34	c.1010delA (1129delA)	p.Glu337fs	Mutation
35	c.1016dupA (1135insA)	p.Lys339fs	Mutation
36			
37	c.1072delC (1191delC)	p.Leu358fs	Mutation
38	c.1121delC (1240delC)	p.Thr374fs	Mutation
39			
40	c.1287dupA (1406insA)	p.Asp430fs	Mutation
41			
42			
43			
44			
45			
46			
47			

1			
2			
3	<b>c.1292dupT (1411insT)</b>	p.Leu431fs	Mutation
4	<b>c.1319delT (1438delT)</b>	p.Leu440fs	Mutation
5	<b>c.1504_1508del5 (1623del5)</b>	p.Leu502fs	Mutation
6	<b>c.1881_1884del4 (2000del4)</b>	p.Val627fs	Mutation
7	<b>c.1961delA (2080delA)</b>	p.Lys654fs	Mutation
8	<b>c.2019delA (2138delA)</b>	p.Glu673fs	Mutation
9	<b>c.2197del5 (2316del5)</b>	p.Glu733fs	Mutation
10	<b>c.2210delC (2329delC)</b>	p.Thr737fs	Mutation
11	<b>c.2212_2215del4 (2331del4)</b>	p.Val738fs	Mutation
12	<b>c.232delA (351delA)</b>	p.Arg78fs	Mutation
13	<b>c.2380dupG (2478insG)</b>	p.Glu787fs	Mutation
14	<b>c.2405_2406delITG (2524delITG)</b>	p.Val802fs	Mutation
15	<b>c.2646_2648delITGC (2765delITGC)</b>	p.Cys882del	Mutation
16	<b>c.2685_2686delAA (2804delAA)</b>	p.Gln895fs	Mutation
17	<b>c.2689insA (2809insA)</b>	p.Pro897fs	Mutation
18	<b>c.2726delA (2845delA)</b>	p.Asn909fs	Mutation
19	<b>c.2727_2730del4 (2846del4)</b>	p.Asn909fs	Mutation
20	<b>c.2728delC (2847delC)</b>	p.Gln910fs	Mutation
21	<b>c.2764_2767del4 (2883delACAG)</b>	p.Thr922fs	Mutation
22	<b>c.2934T&gt;G (3053T&gt;G)</b>	p.Tyr978X	Mutation
23	<b>c.2989_2990dupA (3109insAA)</b>	p.Asn997fs	Mutation
24	<b>c.3329dupA (3448insA)</b>	p.Lys1110fs	Mutation
25	<b>c.3481_3491del11 (3600del11)</b>	p.Glu1161fs	Mutation
26	<b>c.3485delA (3604delA)</b>	p.Asp1162fs	Mutation
27	<b>c.3494_3495delTT (3613-3614delTT)</b>	p.Phe1169fs	Mutation
28	<b>c.3549AG&gt;T (3668AG&gt;T)</b>	p.Lys1183fs	Mutation
29	<b>c.3756_3759del4 (3875del4)</b>	p.Leu1252fs	Mutation
30	<b>c.3770_3771delAG (3889delAG)</b>	p.Glu1257fs	Mutation
31	<b>c.3820dupG (3939insG)</b>	p.Val1274fs	Mutation
32	<b>c.3891_3893delTTC (4010delTTC)</b>	p.Ser1297fs	Mutation
33			
34			
35			
36			
37			
38			
39			
40			
41			
42			
43			
44			
45			
46			
47			

c.4165_4166delAG (4284delAG)	p.Ser1389fs	Mutation
c.4391_4393delCTAinsTT (4510delCTAinsTT)	p.Pro1464fs	Mutation
c.4416delTTinsG (4535delTTinsG)	p.Leu1472fs	Mutation
c.4435delG (4554delG)	p.Val1479fs	Mutation
c.4575_4585del11 (4694del11)	p.Gln1525fs	Mutation
c.493del2 (612delCT)	p.Thr164fs	Mutation
c.5030_5033del4 (5149del4)	p.Thr1677fs	Mutation
c.5137delG (5256delG)	p.Val1713fs	Mutation
c.5191+2delT (IVS19+2delT)	non-coding	Mutation
c.5266dupC (5382insC)	p.Gln1756fs	Mutation
c.5329dupC (5448insC)	p.Thr1777fs	Mutation
c.5360_5361delGTinsAG (5479_5480delGTinsAG)	p.Cys1787fs	Mutation
c.5503_5564del62 (5622del62)	p.Arg1835fs	Mutation
<b>BRCA2</b>		
c.462_463delAA (690delAA)	p.Gln154fs	Mutation
c.1310_1313del4 (1538del4)	p.Lys437fs	Mutation
c.1389_1390delAG (1617delAG)	p.Thr463fs	Mutation
c.1705delC (1933delC)	p.Gln569fs	Mutation
c.2150delG (2378delG)	p.Cys717fs	Mutation
c.2584_2590del7 (2812del7)	p.Lys862fs	Mutation
c.2806_2809del4 (3034del4)	p.Lys936fs	Mutation
c.2957dupA (3185insA)	p.Asn986fs	Mutation
c.3269delT (3497delT)	p.Met1090fs	Mutation
c.3453delC (3681delC)	p.Ile1151fs	Mutation
c.3847_3848delGT (4075delGT)	p.Val1283fs	Mutation
c.3866_3867delAA (4094delAA)	p.Lys1289fs	Mutation
c.4171delG (4399delG)	p.Glu1391fs	Mutation
c.4435ins4 (4763 INS4)	p.Ser14797fs	Mutation

1			
2			
3	<b>c.4449delA (4677delA)</b>	p.Thr1483fs	Mutation
4	<b>c.4456_4459del4 (4684del4)</b>	p.Val1486fs	Mutation
5			
6	<b>c.4480dupA (4708insA)</b>	p.Ser1494fs	Mutation
7	<b>c.469_470delAA (697delAA)</b>	p.Lys157fs	Mutation
8	<b>c.4936_4939del4 (5164del4)</b>	p.Glu1646fs	Mutation
9			
10	<b>c.4940delCA (5168delCA)</b>	p.Thr1647fs	Mutation
11	<b>c.5131delG (5359delG)</b>	p.Val1711fs	Mutation
12	<b>c.5180delA (5408delA)</b>	p.Glu1727fs	Mutation
13			
14	<b>c.5213_5216del4 (5441delCTTA)</b>	p.Thr1738fs	Mutation
15	<b>c.5314delC (4542delC)</b>	p.Val1438fs	Mutation
16	<b>c.5350_5351delAA (5578delAA)</b>	p.Asn1784fs	Mutation
17	<b>c.5577_5580del4 (5805del4)</b>	p.Ile1859fs	Mutation
18	<b>c.5595_5596delAT (5823delAT)</b>	p.Ile1865fs	Mutation
19			
20	<b>c.5681dupA (5909insA)</b>	p.Tyr1894fs	Mutation
21	<b>c.5722_5723delCT (5950delCT)</b>	p.Leu1908fs	Mutation
22	<b>c.5771_5774del4 (5999del4)</b>	p.Ile1924fs	Mutation
23			
24	<b>c.5964delT (6174delT)</b>	p.Ser1982fs	Mutation
25	<b>c.6270_6271delTA (6498delTA)</b>	p.His2090fs	Mutation
26	<b>c.6280_6281delTT (6503delTT)</b>	p.Leu292fs	Mutation
27	<b>c.6280_6286del7 (6508del7)</b>	p.Tyr2094fs	Mutation
28	<b>c.634_635delAG (862delAG)</b>	p.Arg212fs	Mutation
29	<b>c.6351dupA (5579insA)</b>	p.Asn1784fs	Mutation
30	<b>c.6445delAT (6673delAT)</b>	p.Ile2149fs	Mutation
31	<b>c.6591_6592delTG (6819delTG)</b>	p.Thr2197fs	Mutation
32	<b>c.6603_6604delTG (6831delTG)</b>	p.Ser2201fs	Mutation
33	<b>c.6644_6647del4 (6872del4)</b>	p.Tyr2215fs	Mutation
34	<b>c.8904delC (9132delC)</b>	p.Val2969fs	Mutation
35	<b>c.9099_9100delTC (9327delTC)</b>	p.Thr3033fs	Mutation
36	<b>c.9458delG (9686delG)</b>	p.Gly3153fs	Mutation
37	<b>c.994delA (1222delA)</b>	p.Ile332fs	Mutation
38			
39			
40			
41			
42			
43			
44			
45			
46			
47			

- GenBank reference sequences NT\_010755.15 (4920610–5001764; BRCA1) and NT\_024524.13 (13869617– 13953809; BRCA2) cDNA numbering according to reference sequence NM\_007294.3 (BRCA1) and NM\_000059.3 (BRCA2)
- Nucleotide numbering reflects cDNA numbering with +1 corresponding to the A of the ATG translation initiation codon in the reference sequence. The initiation codon is codon 1.

## Supp T2: Specifications of reaction conditions and composition of multiplex reactions

	SET 1	SET 2	SET 3	SET 4	SET 5	SET 6	SET 7	SET 8
<b>AMPLICONS</b>	BRCA2_10_4	BRCA1_21	BRCA1_23	BRCA1_10	BRCA1_15_2	BRCA2_27_3	BRCA1_11_1	BRCA1_11_14
	BRCA2_9	BRCA1_22	BRCA1_17	BRCA2_14_3	BRCA2_23_2	BRCA1_18	BRCA1_19	BRCA2_27_2
	BRCA1_11_12	BRCA2_14_2	BRCA1_3	BRCA1_5	BRCA1_6	BRCA2_18_2	BRCA1_15_1	BRCA1_11_8
	BRCA2_11_15	BRCA1_11_2	BRCA2_11_10	BRCA1_24	BRCA2_7	BRCA2_11_12	BRCA2_22_1	BRCA2_15
	BRCA1_8	BRCA1_9	BRCA1_11_15	BRCA2_11_3	BRCA2_25_1	BRCA1_11_18	BRCA2_10_3	BRCA2_11_24
	BRCA2_22_2	BRCA1_11_4	BRCA2_25_2	BRCA2_10_2		BRCA1_11_6	BRCA1_11_13	BRCA2_11_28
	BRCA1_16_1	BRCA1_16_2	BRCA1_13	BRCA2_11_11			BRCA1_12	BRCA2_20
	BRCA2_11_4	BRCA1_11_16		BRCA1_11_21			BRCA2_27_4	BRCA2_11_1
BRCA2_11_21			BRCA2_11_17				BRCA2_10_6	
<b>PCR program</b>	<i>Touch60</i>	<i>Touch60</i>	<i>Touch58</i>	<i>Touch58</i>	<i>Touch60</i>	<i>Touch58</i>	<i>Touch58</i>	<i>Touch60</i>
<b>Taq polymerase</b>	<i>Titanium Taq</i>	<i>Titanium Taq</i>	<i>PlatinumTaq</i>	<i>PlatinumTaq</i>	<i>Titanium Taq</i>	<i>PlatinumTaq</i>	<i>PlatinumTaq</i>	<i>Titanium Taq</i>
	SET 9	SET 10	SET 11	SET 12	SET 13	SET 14	SET 15	SET 16
<b>AMPLICONS</b>	BRCA1_11_17	BRCA1_11_20	BRCA2_8	BRCA2_12	BRCA2_18_1	BRCA2_11_6	BRCA2_21	BRCA2_2
	BRCA2_27_5	BRCA1_14	BRCA2_10_1	BRCA1_11_7	BRCA2_17	BRCA2_11_20	BRCA2_13	BRCA1_7
	BRCA1_20	BRCA2_11_13	BRCA2_11_7	BRCA2_27_1	BRCA2_3	BRCA1_11_19	BRCA2_11_27	BRCA2_4
	BRCA2_11_14	BRCA1_11_10	BRCA1_11_5	BRCA2_11_9	BRCA2_16	BRCA2_26	BRCA1_11_11	BRCA2_11_23



	BRCA2_23_1	BRCA2_11_22	BRCA2_11_26	BRCA2_10_7	BRCA2_19	BRCA2_11_19	BRCA2_11_25	
	BRCA2_10_5	BRCA2_11_8		BRCA2_14_1		BRCA2_11_5	BRCA1_11_3	
	BRCA1_11_9	BRCA2_5/6				BRCA2_24		
						BRCA2_11_16		
<b>PCR program</b>	<i>Touch60</i>	<i>Touch60</i>	<i>Touch60</i>	<i>Touch60</i>	<i>Touch60</i>	<i>Touch58</i>	<i>Touch60</i>	<i>Touch58</i>
<b>Taq polymerase</b>	<i>Titanium Taq</i>	<i>Titanium Taq</i>	<i>Titanium Taq</i>	<i>Titanium Taq</i>	<i>Titanium Taq</i>	<i>Titanium Taq</i>	<i>Titanium Taq</i>	<i>PlatinumTaq</i>

Supp T3: Primer concentrations in all multiplex reactions

SET 1			SET 2		
concentration	primers/amplicon		concentration	primers/amplicon	
0.005 $\mu$ M	BRCA2_10_4		0.009 $\mu$ M	BRCA1_21	
0.002 $\mu$ M	BRCA2_9		0.001 $\mu$ M	BRCA1_22	
0.004 $\mu$ M	BRCA1_11_12		0.007 $\mu$ M	BRCA2_14_2	
0.016 $\mu$ M	BRCA2_11_15		0.002 $\mu$ M	BRCA1_11_2	
0.016 $\mu$ M	BRCA1_8		0.024 $\mu$ M	BRCA1_9	
0.016 $\mu$ M	BRCA2_22_2		0.004 $\mu$ M	BRCA1_11_4	
0.008 $\mu$ M	BRCA1_16_1		0.008 $\mu$ M	BRCA1_16_2	

---

<b>0.016</b> $\mu\text{M}$	BRCA2_11_4	0.016 $\mu\text{M}$	BRCA1_11_16
----------------------------	------------	---------------------	-------------

<b>0.024</b> $\mu\text{M}$	BRCA2_11_21
----------------------------	-------------

**SET 3****SET 4**

<b>concentration</b>	<b>primers/amplicon</b>
----------------------	-------------------------

<b>concentration</b>	<b>primers/amplicon</b>
----------------------	-------------------------

<b>0.003</b> $\mu\text{M}$	BRCA1_23
----------------------------	----------

0.006 $\mu\text{M}$	BRCA1_10
---------------------	----------

<b>0.004</b> $\mu\text{M}$	BRCA1_17
----------------------------	----------

0.011 $\mu\text{M}$	BRCA2_14_3
---------------------	------------

<b>0.008</b> $\mu\text{M}$	BRCA1_3
----------------------------	---------

0.032 $\mu\text{M}$	BRCA1_5
---------------------	---------

<b>0.008</b> $\mu\text{M}$	BRCA2_11_10
----------------------------	-------------

0.003 $\mu\text{M}$	BRCA1_24
---------------------	----------

<b>0.008</b> $\mu\text{M}$	BRCA1_11_15
----------------------------	-------------

0.012 $\mu\text{M}$	BRCA2_11_3
---------------------	------------

<b>0.008</b> $\mu\text{M}$	BRCA2_25_2
----------------------------	------------

0.024 $\mu\text{M}$	BRCA2_10_2
---------------------	------------

<b>0.012</b> $\mu\text{M}$	BRCA1_13
----------------------------	----------

0.024 $\mu\text{M}$	BRCA2_11_11
---------------------	-------------

0.008 $\mu\text{M}$	BRCA1_11_21
---------------------	-------------

0.024 $\mu\text{M}$	BRCA2_11_17
---------------------	-------------

**SET 5****SET 6**

<b>concentration</b>	<b>primers/amplicon</b>
----------------------	-------------------------

---

<b>concentration</b>	<b>primers/amplicon</b>
----------------------	-------------------------

---

0.004	μM	BRCA1_15_2	0.007	μM	BRCA2_27_3
0.004	μM	BRCA2_23_2	0.002	μM	BRCA1_18
0.004	μM	BRCA1_6	0.003	μM	BRCA2_18_2
0.016	μM	BRCA2_7	0.008	μM	BRCA2_11_12
0.003	μM	BRCA2_25_1	0.008	μM	BRCA1_11_18
			0.024	μM	BRCA1_11_6

## SET 7

## SET 8

concentration primers/amplicon

concentration primers/amplicon

0.016	μM	BRCA1_11_1	0.007	μM	BRCA1_11_14
0.004	μM	BRCA1_19	0.008	μM	BRCA2_27_2
0.008	μM	BRCA1_15_1	0.008	μM	BRCA1_11_8
0.008	μM	BRCA2_22_1	0.005	μM	BRCA2_15
0.016	μM	BRCA2_10_3	0.004	μM	BRCA2_11_24
0.008	μM	BRCA1_11_13	0.016	μM	BRCA2_11_28
0.008	μM	BRCA1_12	0.024	μM	BRCA2_20
0.032	μM	BRCA2_27_4	0.016	μM	BRCA2_11_1

---

0.024  $\mu$ M BRCA2\_10\_6

**SET 9****SET 10**

**concentration**   **primers/amplicon**

**concentration**   **primers/amplicon**

**0.006**  $\mu$ M BRCA1\_11\_17

0.006  $\mu$ M BRCA1\_11\_20

**0.008**  $\mu$ M BRCA2\_27\_5

0.002  $\mu$ M BRCA1\_14

**0.003**  $\mu$ M BRCA1\_20

0.008  $\mu$ M BRCA2\_11\_13

**0.008**  $\mu$ M BRCA2\_11\_14

0.014  $\mu$ M BRCA1\_11\_10

**0.008**  $\mu$ M BRCA2\_23\_1

0.024  $\mu$ M BRCA2\_11\_22

**0.024**  $\mu$ M BRCA2\_10\_5

0.024  $\mu$ M BRCA2\_11\_8

**0.008**  $\mu$ M BRCA1\_11\_9

0.024  $\mu$ M BRCA2\_5en6

**SET 11****SET 12**

**concentration**   **primers/amplicon**

**concentration**   **primers/amplicon**

**0.006**  $\mu$ M BRCA2\_8

0.016  $\mu$ M BRCA2\_12

**0.024**  $\mu$ M BRCA2\_10\_1

0.016  $\mu$ M BRCA1\_11\_7

**0.008**  $\mu$ M BRCA2\_11\_7

0.008  $\mu$ M BRCA2\_27\_1

---

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

<b>0.008</b>	$\mu\text{M}$	BRCA1_11_5	0.008	$\mu\text{M}$	BRCA2_11_9
<b>0.008</b>	$\mu\text{M}$	BRCA2_11_26	0.016	$\mu\text{M}$	BRCA2_10_7
			0.008	$\mu\text{M}$	BRCA2_14_1
<b>SET 13</b>			<b>SET 14</b>		
<b>concentration</b>		<b>primers/amplicon</b>	<b>concentration</b>		<b>primers/amplicon</b>
<b>0.002</b>	$\mu\text{M}$	BRCA2_18_1	0.008	$\mu\text{M}$	BRCA2_11_6
<b>0.024</b>	$\mu\text{M}$	BRCA2_17	0.016	$\mu\text{M}$	BRCA2_11_20
<b>0.004</b>	$\mu\text{M}$	BRCA2_3	0.008	$\mu\text{M}$	BRCA1_11_19
<b>0.024</b>	$\mu\text{M}$	BRCA2_16	0.016	$\mu\text{M}$	BRCA2_26
<b>0.008</b>	$\mu\text{M}$	BRCA2_19	0.016	$\mu\text{M}$	BRCA2_11_19
			0.008	$\mu\text{M}$	BRCA2_11_5
			0.016	$\mu\text{M}$	BRCA2_24
			0.024	$\mu\text{M}$	BRCA2_11_16
<b>SET 15</b>			<b>SET 16</b>		
<b>concentration</b>		<b>primers/amplicon</b>	<b>concentration</b>		<b>primers/amplicon</b>

<b>0.008</b>	μM	BRCA2_21	0.008	μM	BRCA2_2
<b>0.008</b>	μM	BRCA2_13	0.008	μM	BRCA1_7
<b>0.01</b>	μM	BRCA2_11_27	0.008	μM	BRCA2_4
<b>0.01</b>	μM	BRCA1_11_11	0.024	μM	BRCA2_11_23
<b>0.008</b>	μM	BRCA2_11_25			
<b>0.016</b>	μM	BRCA1_11_3			

	Sensitivity					
	Detected variants		False negatives (coverage)		True False negatives	
	#	unique	#	unique		
<b>PoC 1</b>	146	30	16	14	0	0
<b>PoC 2</b>	160	40	12	5	1	1
<b>PoC 3</b>	43	29	12	8	0	0
<b>PoC 4</b>	204	115	2	2	2	2

Supp T4A: Summary of sensitivity data

(VIP)

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

<b>Total</b>	553	144	42	29	3	3
--------------	-----	-----	----	----	---	---

For Peer Review

Supp T4B: Summary of specificity data (VIP)

<b>Specificity (false positives)</b>										
	raw data		Filter 1 (coverage > 38)		Filter 2 (AF >25%)		Filter 3 (Q>30)		Filter 4 (Hp>7)	
	#	unique	#	unique	#	unique	#	unique	#	unique
<b>PoC 1</b>	1145	205	1040	168	382	48	96	27	43	14
<b>PoC 2</b>	2948	693	1955	474	444	112	165	59	161	55
<b>PoC 3</b>	764	377	480	248	145	76	74	36	61	28
<b>PoC 4</b>	213	177	163	140	39	47	13	19	11	7
<b>Total</b>	5070	1452	3638	1030	1010	283	348	141	276	104

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

For Peer Review