



HAL
open science

Design and validation of a metabolic disorder resequencing microarray (BRUM1).

Christopher Bruce, Matthew Smith, Fatima Rahman, Zhi-Feng Liu, Dominic
Mcmullan, Sarah Ball, Jane Hartley, Marian A Kroos, Lesley Heptinstall,
Arnold Jj Reuser, et al.

► **To cite this version:**

Christopher Bruce, Matthew Smith, Fatima Rahman, Zhi-Feng Liu, Dominic Mcmullan, et al.. Design and validation of a metabolic disorder resequencing microarray (BRUM1).. Human Mutation, 2010, 31 (7), pp.858. 10.1002/humu.21261 . hal-00552386

HAL Id: hal-00552386

<https://hal.science/hal-00552386>

Submitted on 6 Jan 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.



Design and validation of a metabolic disorder resequencing microarray (BRUM1).

Journal:	<i>Human Mutation</i>
Manuscript ID:	humu-2009-0430.R1
Wiley - Manuscript type:	Methods
Date Submitted by the Author:	23-Feb-2010
Complete List of Authors:	<p>Bruce, Christopher; The University of Birmingham, Clinical and Experimental Medicine</p> <p>Smith, Matthew; Birmingham Women's Hospital, West Midlands Regional Genetics Laboratory and Clinical Genetics Unit</p> <p>Rahman, Fatima; The University of Birmingham, Clinical and Experimental Medicine</p> <p>Liu, Zhi-feng; Nanjing Medical University, Department of Digestory</p> <p>McMullan, Dominic; Birmingham Women's Hospital, West Midlands Regional Genetics Laboratory and Clinical Genetics Unit</p> <p>Ball, Sarah; Birmingham Children's Hospital</p> <p>Hartley, Jane; The University of Birmingham, Clinical and Experimental Medicine</p> <p>Kroos, Marian; Erasmus MC, Clinical Genetics</p> <p>Heptinstall, Lesley; Royal Manchester Children's Hospital, Willink Biochemical Genetics Unit</p> <p>Reuser, Arnold; Erasmus MC, Clinical Genetics</p> <p>Rolfs, Arndt; Universität Rostock, Albrecht-Kossel-Institut für Neuroregeneration</p> <p>Hendriksz, Chris; Birmingham Children's Hospital, The Metabolic Unit; Birmingham Children's Hospital, The Metabolic Unit</p> <p>Kelly, Deirdre; Birmingham Children's Hospital, The Liver Unit</p> <p>Barrett, Timothy; The University of Birmingham, Clinical and Experimental Medicine</p> <p>Macdonald, Fiona; Birmingham Women's Hospital, West Midlands Regional Genetics Laboratory and Clinical Genetics Unit</p> <p>Maher, Eamonn; The University of Birmingham, Clinical and Experimental Medicine; Birmingham Women's Hospital, West Midlands Regional Genetics Laboratory and Clinical Genetics Unit</p> <p>Gissen, Paul; The University of Birmingham, Clinical and Experimental Medicine; Birmingham Children's Hospital, The Metabolic Unit</p>

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



For Peer Review

Design and validation of a metabolic disorder resequencing microarray (BRUM1).

CK Bruce¹, M Smith², F Rahman¹, ZF Liu³, DJ McMullan², S Ball⁴, J Hartley^{1,5}, MA Kroos⁶,
L Heptinstall⁷, AJJ Reuser⁶, A Rolfs⁸, C Hendriksz⁹, DA Kelly⁵, TG Barrett^{1,10}, F
MacDonald², ER Maher^{1,2}, P Gissen^{1,9}

¹School of Clinical and Experimental Medicine, The University of Birmingham;

²West Midlands Regional Genetics Service, Birmingham Women's Hospital;

³Department of Digestory, Nanjing Children's Hospital, Nanjing Medical University;

⁴Department of Clinical Chemistry, Birmingham Children's Hospital, Birmingham, UK;

⁵The Liver Unit, Birmingham Children's Hospital, Birmingham, UK;

⁶Department of Clinical Genetics, Erasmus MC, Rotterdam, The Netherlands;

⁷Royal Manchester Children's Hospital, Manchester, UK;

⁸Albrecht-Kossel-Institut für Neuroregeneration, Universität Rostock, Rostock, Germany;

⁹The Metabolic Unit, Birmingham Children's Hospital, Birmingham, UK;

¹⁰The Diabetes Unit, Birmingham Children's Hospital, Birmingham, UK.

Corresponding Author: Dr. P. Gissen.

Tel: +44 121 415 8529

Email: p.gissen@bham.ac.uk

Abstract

The molecular genetic diagnosis of inherited metabolic disorders is challenging. The diseases are rare and most show locus heterogeneity. Hence testing of the genes associated with IMDs is time consuming and often not easily available. We report a resequencing array that allows the simultaneous resequencing of up to 92 genes associated with IMDs.

To validate the array, DNA samples from 51 patients with 52 different known variants (including point variants, small insertion and deletions (indels)) in 7 genes (*CI4ORF133*, *GAA*, *NPC1*, *NPC2*, *VPS33B*, *WFS1*, *SLC19A2*) were amplified by PCR and hybridised to the array. A further patient cohort with 48 different mutations in *NPC1* were analysed blind. Out of 76 point variants, 73 were identified using automated software analysis followed by manual review. Ten insertion and deletion variants were detected in the extra tiling using mutation specific probes whilst 11 heterozygous deletions and 3 heterozygous insertions.

In summary, we identified 96% (95%CI 89-99%) of point variants added to the array, but the pickup rate reduced to 83% (95%CI 75-89%) when insertions/deletions were included. Whilst the methodology has strengths and weaknesses, application of this technique could expedite diagnosis in most patients with multi-locus IMDs.

Introduction

A number of inherited diseases present with a similar phenotype, but may be caused by mutations within different genes. For example, amongst the lysosomal storage disorders, patients with genetically distinct mucopolysaccharidoses have common clinical features (Neufeld and Muenzer, 2001). This also applies to the group of neuronal ceroidlipofuscinoses (Jalanko and Braulke, 2009). Further, the severe hepatic neurodegenerative disorder Niemann-Pick Type C is caused by mutations in two genes *NPC1* and *NPC2*. Although these forms are indistinguishable clinically, the management of patients with *NPC2* may include bone marrow transplantation as well as other therapeutic options available in this disease (Wraith et al., 2009). The neuromuscular and the mitochondrial diseases are two other groups of diseases in which the clinical phenotype leaves choices as to the underlying metabolic and causative genetic defects. Thus the availability of a rapid accurate molecular diagnostic platform could enhance acute clinical management as well as genetic counselling and prenatal diagnosis.

Current sequencing technologies employed in routine diagnostic laboratories rely on Sanger sequencing, and whilst providing the gold standard for the detection of small intragenic mutations, most diagnostic laboratories do not provide a service offering rapid analysis of multiple candidate genes (generally candidate genes are sequenced sequentially until a mutation is identified). In addition, there is a limit to how much sequence can be obtained from di-deoxysequencing; often, each exon has to be assayed as a single experiment. Second generation sequencing technologies can enable multiple gene testing in a single experiment, but it is not yet clear how suitable this approach is for routine testing in a clinical service laboratory. Microarray based resequencing is an evolving laboratory technique that offers the potential for rapid simultaneous mutation testing of multiple genes (Hacia, 1999). This technology has been used for pathogen identification (Malanoski et al.,

1
2
3 2006) bacterial genotyping (Corless et al., 2008; Zwick et al., 2008), mitochondrial DNA
4 sequencing (Hartmann et al., 2009), and human gene mutation identification (Denning et al.,
5
6 2007; Liu et al., 2007; Takahashi et al., 2008). One group has designed a resequencing chip,
7
8 the Jaundice Chip, that is being offered as a diagnostic test to investigate inherited syndromes
9
10 of intrahepatic cholestasis ((Liu et al., 2007)
11 <http://www.cincinnatichildrens.org/svc/alpha/m/molecular-genetics/jaundice-chip.htm>).
12
13
14
15
16

17 Microarray based resequencing could be a promising technology for many molecular genetic
18 diagnostic laboratories as the equipment is often available in clinical and research
19 laboratories and can be used for clinical purposes. However, the mutation detection rate, the
20 reproducibility, and the accuracy of the resequencing technology have not been explored in
21 detail. Therefore we used a custom designed resequencing microarray to address these
22 questions and to elucidate areas of potential improvement for the resequencing technology
23 and microarray analysis. The novel array (Birmingham ReseqUencing Microarray version
24 1(BRUM1)) includes 92 genes involved in metabolic pathways (or encoding proteins
25 involved in various aspects of metabolism) and enables simultaneous resequencing of
26 multiple genes that may each cause similar phenotypes within a group of disorders (eg
27 hyperlipidaemia, glycogen storage disease and lysosomal storage diseases).
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

48 **Materials and Methods**

49 **DNA Samples**

50
51 A total of 51 genomic DNA samples harbouring disease causal gene mutations in
52
53
54 *NPC1*, *NPC2*, *VPS33B*, *C14ORF133*, *GAA*, *SLC19A2*, and *WFS1* (Supp. Table S1) were
55
56
57
58
59
60 obtained from the West Midlands Regional Genetics Service (Birmingham, UK),

1
2
3 Birmingham Children's Hospital (Birmingham, UK), Willink Biochemical Genetics Unit
4
5 (Manchester, UK), and the Erasmus Medical Center (Rotterdam, NL). Genes and sample
6
7 DNA were chosen based on the availability of DNA covering a wide range of mutation types.
8
9
10 DNA from a human cell line MRC5VA was used as a control.
11

12
13 A further cohort of 28 DNA samples that harboured 48 mutations detected by direct
14
15 sequencing within *NPCI* were obtained from Albrecht-Kossel-Institute for
16
17 Neuroregeneration (Rostock, Germany). Only the exon number was known prior to testing
18
19 the arrays but not the nature or precise location of the mutation.
20
21

22 23 *Microarray Resequencing Chip Design* 24

25
26 A total of 92 genes involved in various inherited disorders (Supp. Table S2) were
27
28 selected to be sequenced on a 300kb resequencing microarray (Affymetrix). All coding exons
29
30 with the addition of 25bp intronic sequence either side of the exon were selected. This allows
31
32 for full sequencing of the exon, plus 13bp of intronic sequence allowing splice site mutations
33
34 to be identified. The design of the array is such that there are a series of 25bp probes designed
35
36 against the supplied reference sequences. As the array simultaneously sequences both
37
38 forward and reverse strands, there are 8 probes for any given base; 4 for the forward and 4 for
39
40 the reverse. The central base differs between the four probes allowing all variants to be
41
42 identified. Sequencing commences at the 13thbp of the reference and completes 13thbp away
43
44 from the end of the reference sequence.
45
46
47
48
49

50
51 Small insertions/ deletions and indels, especially in a heterozygous state are known
52
53 limitations for resequencing arrays, but by including specific probes for known insertions and
54
55 deletions on the array, these limitations can be overcome (Karaman et al., 2005; Kothiyal et
56
57 al., 2009). Thus, 1343 known insertions/ deletions and indels (all of 5 or less nucleotides)
58
59
60

1
2
3 were identified from the Human Genome Mutation Database (www.hgmd.org) and included
4
5 on the array. Thus this array is divided into two sections, the first herein termed main tiling is
6
7 used to primarily detect single nucleotide variations, whilst the second section herein termed
8
9 extra tiling, is used to detect insertion/ deletion variations. The sequences of all the genes
10
11 were first analysed for repeat regions using repeat masker ([http://repeatmasker.org/cgi-](http://repeatmasker.org/cgi-bin/webrepeatmasker)
12
13 [bin/webrepeatmasker](http://repeatmasker.org/cgi-bin/webrepeatmasker)) and any repeats or low complexity regions of >25bp were excluded.
14
15 The sequences were submitted to Affymetrix for final array design and production. In total,
16
17 232 515bp can be sequenced on this array.
18
19
20
21
22

23 *Microarray resequencing.*

24
25
26 In order to efficiently test the arrays with 52 mutations, a combination of short and
27
28 long-range PCR amplifying individual exons of 7 genes (*C14ORF133*, *GAA*, *NPC1*, *NPC2*,
29
30 *SLC19A2*, *VPS33B* and *WFS1*) was undertaken. For exons without a known mutation, control
31
32 DNA was used. Each fragment was independently amplified in triplicate allowing
33
34 comparisons of the sequencing to be undertaken. Pooled PCR products were labelled PCR-A,
35
36 PCR-B or PCR-C. Triplicate hybridisations for each pool were labelled PCR-A1, PCR-A2
37
38 etc. PCR was undertaken in 50µl volumes using 50ng of DNA, 0.5µM each primer (Sigma),
39
40 0.5mM dNTPs (Bioline) and 2 units of Taq DNA polymerase (Biomix Red, Bioline UK) in
41
42 1x PCR buffer. Primers were designed using Exon Primer
43
44 (<http://ihg2.helmholtz-muenchen.de/ihg/ExonPrimer.html>; primers available on request) or
45
46 had been previously published (Gissen et al., 2004; Cullinane et al., 2009; Cullinane et al.,
47
48 2010). A 1kb control fragment was amplified using a template supplied by Affymetrix. The
49
50 size range of amplicons used in this test ranged from 156bp to 4500bp with 250pmoles of
51
52 each amplicon required for analysis. PCR product quantitation was undertaken using a
53
54 picogreen assay, pooling, fragmentation, labelling and hybridisation were performed
55
56
57
58
59
60

1
2
3 according to the GeneChip Custom Resequencing Array Protocol V2.1. Arrays were washed
4
5 and stained using a FS450 fluidics station before being scanned with a GCS3000 7G scanner.
6
7

8 *BRUM1 Data Analysis.*

9
10
11 Intensity files were produced using AGCC (Command Console V1.0) and processed
12
13 in GSeq 4.1 (Affymetrix) that uses the resequencing algorithm version 2. Several parameters
14
15 within this base calling can be altered which thus has an effect on the call rate and accuracy
16
17 of the base calls (see Cutler et al., (2001) and Di and Cawley, (2005) for full details of the
18
19 algorithm). Increasing the threshold on 2 parameters, the Quality Score Threshold (QST) and
20
21 Base Reliability Threshold (BRT) will decrease the number of calls, thus making the data
22
23 more stringent, whereas increasing the threshold on the sequence profile threshold and trace
24
25 threshold will increase the number of calls, thus making the data less stringent. A further two
26
27 parameters, the Max Signal to noise ratio and Modeltype have no effect on increasing or
28
29 decreasing the call rate. The Modeltype allows the choice between haploid and diploid data.
30
31 Haploid fits the data to 5 outcomes (A,G,C,T and N) whereas the diploid fits data to 11
32
33 outcomes (A,G,C,T,N,K,M,R,S,W, and Y) allowing heterozygous calls to be made. Altering
34
35 the weak signal and no signal thresholds has an indeterminate effect on base calling. For the
36
37 purpose of the analysis, base calling utilised the diploid model, whilst signal to noise ratio, no
38
39 signal threshold, sequence profile threshold and weak signal fold threshold were left at
40
41 default values. A range of quality score threshold, a measure of the reliability of the base call
42
43 (QST; 2,3,6,9,12,30) was used. As increasing the QST, makes the data more stringent a base
44
45 with a higher quality score is deemed to be more reliable and more accurate. The effect of
46
47 changing the base reliability threshold (BRT) was also assessed, using either 0 or 0.5 (if 50%
48
49 of samples no-called at this base, all samples would be no-called).
50
51
52
53
54
55
56
57
58
59
60

Di-deoxysequencing.

1
2
3 To compare the microarray based sequencing with direct sequencing, surplus PCR
4 products were directly sequenced. PCR product were treated with ExoSap IT (GE
5 Healthcare) before being directly sequenced using BigDye 3.1 (Applied Biosystems)
6
7 Following precipitation and washing, sequencing products were HiDiformamide (Applied
8 Biosystems) and analysed using a DNA analyzer3730xl (Applied Biosystems). The resulting
9
10 sequence traces were analysed using Sequence Analysis 5.2.2 (Applied Biosystems).
11
12
13
14
15
16
17

18 *Statistics.*

19
20
21 In order to assess the reliability and reproducibility of the microarray, Fleiss Kappa
22 statistic was used to compare base calls between samples and also to compare the base calls
23 against di-deoxysequencing (Fleiss, 1971). Fleiss Kappa allows comparisons between
24 multiple raters (here each BRUM1 chip) to assess their agreement. This was implemented in
25 Microsoft Excel (King, 2004). Confidence intervals for proportions were calculated at
26 Measuring Usability (<http://www.measuringusability.com/wald.htm>) using the Adjusted
27 Wald method (Sauro and Lewis, 2005)
28
29
30
31
32
33
34
35
36
37
38
39
40
41

42 **Results**

43 *Chip Assembly.*

44
45
46 Analysing the sequences prior to BRUM1 array assembly, there were eight regions
47 that proved to have repetitive DNA sequences of >25bp. These were located in *RET*, *IDUA*,
48
49 *NAGLU*, *ARSB*, *PTPN11*, *EIF2AK3*, *APOB* and *PCSK9*. Repeats shorter than 25bp can be
50 safely tiled onto the array (GeneChip Custom Resequencing Array Design Guide). However,
51 this means that the above genes cannot be fully sequenced, resulting in 30bp of *RET* (CTG
52 repeat), 72bp of *IDUA* (GC rich low complexity), 147bp of *NAGLU* (GC rich, low
53
54
55
56
57
58
59
60

1
2
3 complexity), 42bp of *ARSB* (GC rich, low complexity), 31bp of *PTPN11* (CAAAA repeat),
4
5 75bp of *EIF2AK3* (CTG and CGG repeats), 38bp of *APOB* (CTG repeat) and 38bp of *PCSK9*
6
7 (CTG repeat) being excluded. Further during the design phase, all sequences were compared
8
9 to each other to predict cross hybridisation that resulted in 927 submitted sequences having a
10
11 homology of greater than 10% (range 10.07%-100%). All the sequences that demonstrated
12
13 homology were between the parental sequences and extra tiling included to detect insertion
14
15 and deletion variants or between the extra tiling itself. For example, NPC1#01 (NPC1 exon 1)
16
17 and NPC1#M21 (used to detect NPC1 46-47delTG) have 15% homology. However, this
18
19 observation was to be expected.
20
21
22
23
24

25 The final BRUM1 chip comprises 1.8million 25bp probes generated as 1.525 micron
26
27 features and has the capabilities of sequencing 232 515 bp in a single experiment.
28
29

30 *Call Rate and accuracy of the BRUM1.*

31
32
33

34 The call rate was assessed under several conditions with and without the BRT
35
36 (Table 1). The overall call rate for the BRUM1 array was low, ranging from 0.05% for a QST
37
38 of 30 and BRT of 0.5 to 31% for a QST of 2 and BRT of 0. However, only a portion of the
39
40 BRUM1 array was used, approximately 15 890 out of the total 232 515 which reflects in this
41
42 low call rate. Assessing the average call rate for the genes sequenced, an average call rate of
43
44 0.25% for a QST of 30 with a BRT of 0.5 was observed, rising to 94.25% for a QST of 2 and
45
46 BRT of 0. This equates to 15815bp that could not be assigned unambiguously at a QST of 30,
47
48 dropping to 912bp at a QST of 2.
49
50
51
52

53 The accuracy of base calling was better at a higher QST being 100% with a QST of
54
55 30. However, using a QST of 2, the accuracy of base calling dropped to 99.45%. This
56
57 indicates that 87bp out of the 15 890bp sequenced may be called incorrect. Using a QST of 2
58
59
60

1
2
3 without the BRT gave the highest call rate with only a marginal increase in incorrect base
4 calls. Therefore these settings were used in all subsequent analysis. Quality scores ranged
5
6 from 0-45.
7
8

9
10
11 *Reproducibility.*
12

13
14 In order to assess the reproducibility of the assay, triplicate PCR's and triplicate
15 hybridisation experiments were performed. Each hybridisation challenge comprised the same
16 fragments either from independent PCR's (A,B and C) or from the same PCR pool (A1, A2,
17 A3; B1, B2, B3; C1, C2, C3). Base calling was undertaken with the above settings and the
18 resulting sequences exported. Resulting sequences from arrays comprising either PCR-A,
19 PCR-B or PCR-C were aligned to the reference sequences and compared using Fleiss Kappa
20 statistics (Table 2). In total 15880bp were compared between each array for the triplicate
21 PCR pool. PCR-A had the lowest agreement with 91.7% bases in agreement (Kappa=0.895
22 CI 0.891-0.898). PCR-B and PCR-C were very similar with 96.8% (Kappa=0.959 CI 0.955-
23 0.962) and 97.3% (Kappa=0.966 CI 0.962-0.970) bases in agreement. Analysing the same
24 15 880bp between all nine arrays, 94.8% (Kappa=0.932 CI 0.933-0.935) were in agreement.
25
26 However, it must be noted that two fragments failed to hybridise efficiently from PCR-A,
27 which is reflected in the lower agreement between the repeats.
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45

46 *Consistent regions of no-calls.*
47
48

49 Certain probes, for example those containing increased numbers of Gs, can be
50 problematic for resequencing arrays resulting in decreased signal and an increased chance of
51 being no-called (Cutler et al., 2001), To assess whether there were any consistent regions of
52 no-calls, all fragments were aligned. Any base or region no-called in all 9 BRUM1s was
53 defined as a region of consistent no-call. Table 3 gives an overview of the consistent regions
54
55
56
57
58
59
60

1
2
3 of no-call. Interestingly, the most common base consistently no-called was T accounting for
4
5 19.5% of all no-called regions. The most common multiple base no-called region were
6
7 regions containing multiple C's accounting for 12.7% of these regions. Analysis of the
8
9 nucleotide content of the probes showed that probes with higher than average GC content
10
11 were more likely to be no-called. Further as the proportion of G bases increased in a single
12
13 probe, there was an increased chance of the base being no-called.
14
15
16

17
18 *Comparison of BRUM1 array based sequencing to di-deoxysequencing.*
19

20
21 Direct sequencing of the DNA used on the microarray was undertaken to determine
22
23 whether there was any difference when the two methods were compared. For the purpose of
24
25 automated analysis of di-deoxysequencing, a heterozygous call was made if the lower peak
26
27 height was >40% of the peak height of the main peak. A total of 8147bp from 58 different
28
29 fragments were compared to BRUM1 array sequencing with an overall agreement of 93.1%
30
31 (Kappa =0.91 CI 0.90-0.92). There were 39 discordant calls between the two methodologies,
32
33 385 were no-called in either methodology with 11 no-calls being common between both
34
35 methodologies. A further 30 were a heterozygous call in one methodology, but with a related
36
37 homozygous call in the other. There was an overall agreement of 93% between automated
38
39 analysis of di-deoxysequencing to the reference sequence.
40
41
42
43
44
45

46 However, following manual review of both di-deoxysequencing and BRUM1
47
48 sequencing, 112 bases were removed from the analysis due to low quality di-
49
50 deoxysequencing. These low quality bases were confined at the beginning or end of the di-
51
52 deoxysequencing. Comparison of the remaining 8035 bases demonstrated an agreement of
53
54 99.7% (Kappa 0.99 CI 0.98->0.99). There were no discordant calls, 5 bases were no-called in
55
56 either methodology and 22 bases were a heterozygous call in one methodology but related
57
58 homozygous in the other. Comparison of the reviewed di-deoxysequencing to the reference
59
60

1
2
3 sequence there was an agreement of 98.9% (Kappa 0.99 CI 0.97->0.99). There were 8
4
5 discordant calls, 8 no calls and 21 heterozygous calls.
6
7

8
9 *Detection of Variation: single nucleotide variations.*
10

11
12 The main purpose of the resequencing microarray was to detect mutations. In order to
13 assess the ability of the array to detect mutations, we initially amplified the DNA known to
14 contain 39 point mutations/ polymorphisms, which was then hybridised to the array. Using
15 automated analysis, 32 out of 39 point variations were detected demonstrating a pickup rate
16 of 82.1% (95%CI 67.0-91.3%). Of the changes not picked up by automated base calling, 6
17 were manually called increasing the pickup rate to 97.4% (95%CI 85.6-99.9%). In addition, 4
18 false positive calls and 5 known polymorphisms were identified (Table 4).
19
20
21
22
23
24
25
26
27
28

29 *Detection of Variation: insertions and deletions.*
30

31
32
33 As insertions and deletions are known problems for resequencing arrays, the majority
34 of the known insertions and deletions were tiled separately on the array (Karaman et al.,
35 2005; Kothiyal et al., 2009). To test how well the technology detects insertions and deletions,
36 3 samples with known insertions and 10 samples with known deletions were amplified and
37 hybridised to the array. Analysing the main tiling, six deletions that were in a heterozygous
38 state and one deletion in a homozygous state were not detected. However, 3 other
39 homozygous deletions were detected and manifested as regions of no-call. Furthermore, 1
40 homozygous insertion manifested as a compound heterozygote but no other insertions were
41 identified. However, when analysing the extra tiling all insertions and 6 of the 10 deletions
42 were easily identified. Three heterozygous deletions were not identified. Overall, 48 out of 52
43 variations were detected representing a pickup rate of 92.3% (95%CI 84.0-97.2%)
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

Detection of Variation: Blind test of BRUM1 sequencing NPC1.

To assess how the array performed when the primary investigators had no prior knowledge of the mutations (except for exonic location), 28 NPC1 samples were analysed for 48 variations (37 were single nucleotide changes, 7 were novel heterozygous deletions, 3 were novel heterozygous insertions and 1 was an extra tiled deletion). Of these variations 35 were detected representing a pickup rate of 72.3% (95% CI 58.9-83.5%). Of these changes we could detect 35/37 (94.6%) of point mutations, however, we were unable to detect any insertion or deletions.

Discussion

This paper describes the development and validation of a novel microarray based resequencing platform that has been designed to sequence a number of genes involved in inherited metabolic disorders simultaneously. In order to validate the BRUM1 array, 7 genes were chosen based on the availability of samples with known mutations. A range of mutations were chosen such that a broad spectrum of mutation types would be tested, i.e. homozygous and heterozygous point mutations, homozygous insertions and heterozygous/homozygous deletions. Whilst this would give an indication of mutation pickup rate, and which mutations type would cause technical problems, it was also important to determine what GSEQ algorithm parameters gave the best base call with minimal impact on accuracy, and further investigate reproducibility and comparison to di-deoxysequencing.

Array design.

During the design phase of the array, it was noted that several sequences had repetitive or low complexity stretches of DNA that could not be safely tiled on the array and

1
2
3 thus had to be excluded. However, the removal of these sequences means that the genes in
4
5 question could not be completely analysed. In these patients with negative microarray
6
7 resequencing results could then be analysed by di-deoxysequencing of these fragments.
8
9

10 11 *Base Calling.* 12

13
14 As is evident from Table 1, altering the algorithm parameters drastically alters the call
15
16 rates. Every base is given a quality score based on a comparison of signal intensity across the
17
18 8 probes for a single base. The higher the score is, the more reliable the call is. Hence, with a
19
20 QST of 30, all bases are called with extremely high degree of accuracy. However, this also
21
22 meant that only 0.2% of bases were called. Using this parameter would mean extensive
23
24 operator input to manually base call. Given the number of bases sequenced in a single
25
26 experiment, this would prove impractical. Lowering the QST did increase the call rate. Using
27
28 a QST of 2 increased the call rate to 94%. This is in broad agreement with other custom
29
30 resequencing microarrays that demonstrate call rates of 93.5% (Liu et al., 2007) and 93.6%
31
32 (Denning et al., 2007). However, by lowering the QST, there is a concomitant reduction in
33
34 base call accuracy. Comparing the QST of 3 and 2, there is a reduction in base call accuracy
35
36 from 99.53% to 99.45%. This equates to the potential of 74 inaccurate bases rising to the
37
38 potential of 87 incorrect bases by changing the QST from 3 to 2.
39
40
41
42
43
44
45

46 Using automated analysis, false positive and negative results were observed. The false
47
48 negative base calling regions were primarily located within NPC1 and with the exception of 1
49
50 base, manifested as regions of no-call. Interestingly, 3 of these false negative bases increased
51
52 the length of a run of identical bases to 4 or more. All the false negatives manifesting as no-
53
54 calls were picked up during manual review of the no-called regions. However, one
55
56 heterozygous change could not be detected. Increasing the stringency of analysis and using a
57
58 QST of 6, this base was consistently called as wild type sequence. Analysing the peak heights
59
60

1
2
3 at this region, the base was called as a heterozygote from the forward sequencing, but
4
5 homozygous wild type called in the reverse (Figure 1). The base calling algorithm uses a set
6
7 of rules to assign a final base call. One rule states “If one call is homozygote, the other is
8
9 heterozygote and the homozygote allele is one of the two in the heterozygote call, assign
10
11 homozygote call and its quality score to that base” (Di and Cawley, 2005). Consequently, as
12
13 this base was wild type called, it would not have been picked up during manual review of the
14
15 no called regions and is a true false negative. Five false positive results were detected of
16
17 which only one could not be eliminated following manual review. Comparison of this
18
19 remaining false positive to di-deoxysequencing proved that it was a true false positive result.
20
21 The overall false positive rate was calculated to be <0.01% (Suojanen, 1999). Thus we felt
22
23 the algorithm parameters used were satisfactory and are in agreement with other studies
24
25 (Kothiyal et al., 2009).
26
27
28
29
30
31
32

33 It must also be noted, that only 9 BRUM1 arrays have been used to analyse these data.
34
35 The algorithm performs best with large number of samples (Di and Cawley, 2005). Thus one
36
37 would expect the base calling and accuracy to improve with greater numbers of BRUM1
38
39 arrays. Indeed analysing a single sample gave an average call rate of 85.4% whilst the same
40
41 BRUM1 array analysed as a batch increased the call rate to 94.8%. An improvement in call
42
43 rate was also observed by (Liu et al., 2007) when multiple samples were analysed as a batch.
44
45 Accuracy also improved marginally from 99.6% to 99.7%.
46
47
48
49

50 *Reproducibility.*

51
52

53 In order to reliably sequence DNA, it is important that the sequencing is reproducible.
54
55 In order to test this, three independent PCRs were performed for every single amplicon
56
57 tested. Further, to test cross chip reliability, identical products were loaded across three
58
59 BRUM1 arrays. This would give an indication of interchip variation. Overall there was a 94%
60

1
2
3 agreement between samples. This varied between 90.3% ($\kappa=0.877$, CI 0.874-0.881) to
4
5 97% ($\kappa=0.968$, CI 0.964-0.971) depending which combination of BRUM1 arrays used.
6
7 The likely reason for the 7% variation is that BRUM1-A3, giving the lowest agreement, also
8
9 had three problematic exon hybridisations. Further, the eighth exon of WFS1 proved to be
10
11 problematic for all BRUM1 arrays. Whilst the PCRs were very clean, the call rate (84%) for
12
13 this particular exon was not efficient. Assessing the signal intensities for this exon, it is
14
15 evident that saturation has been reached for a large number of the bases in either the forward
16
17 or reverse sequence. As a result, the base calling algorithm no-called 16% of these bases.
18
19
20
21

22 *No-Called regions.*

23
24
25
26 Comparisons between the samples demonstrated several regions of consistent no-call.
27
28 These were mainly confined to areas surrounding repeat regions and runs of a single base, in
29
30 particular C. Analysing the GC content of the no-called probes, there is a skew towards
31
32 probes having a higher GC content. Further, assessing the individual nucleotide content of the
33
34 probes, there was a decrease in fluorescence for those with an increased number of C bases.
35
36 This phenomenon has been observed previously (Cutler et al., 2001), but we have no
37
38 explanation for this observation. Despite these no-called regions, when assessing the no-
39
40 called regions manually, it is evident that the majority of these regions can be manually
41
42 called. Interestingly, the most common no-called single base was a T. It is important to note,
43
44 that despite regions of no-call, mutations that reside in these no-called areas can be detected.
45
46 However, if the mutation increases the length of a run of bases that are no-called, the mutated
47
48 base is also likely to be no-called. But, if no-called, these changes are likely to get picked up
49
50 during manual review of the sequences.
51
52
53
54
55

56
57
58 From our data, an average of 55 (0.35%) no-calls remained following manual review.
59
60 The majority of these no-calls were confined to *WFS1*. We are exploring recommendations

1
2
3 by (Kothiyal et al., 2009) to tile each sequence in triplicate but altering the interrogation base
4
5 of each set of probes to assess whether this tiling strategy will decrease the no-calls observed.
6
7
8 At present however, if a base is no-called and cannot be assigned through manual review and
9
10 no mutations are identified elsewhere in the gene, this sequence will need to be di-
11
12 deoxysequenced to ascertain the true nature of this base.
13
14

15 16 *Variation detection.* 17

18
19 The main goal of the custom resequencing array is the ability to detect mutations in
20
21 patients with a suspected disease to pinpoint the cause and give a better diagnosis and more
22
23 informed prognosis. To determine how well the BRUM1 array faired at detecting sequence
24
25 variations, a broad spectrum of changes were tested. A pick up rate of 97.4% (95%CI 85.6-
26
27 99.9%) was found for point mutations, 100% for known tiled deletions, 25% for unknown or
28
29 not tiled deletions and 100% for known tiled insertions. For the point variations, 18 changes
30
31 were in a heterozygous state, whilst 20 were in a homozygous state. As expected, deletions
32
33 were the main problem for the array. Only 1 out of 4 untiled deletions were detected in the
34
35 main tiling, and this was in a homozygous state (Supp. Figure S1a). Of the remaining three
36
37 deletions, one of which was a large 12bp deletion that was felt sufficiently large to be
38
39 detected by other means. The other two deletions had not been included in the HGMD since
40
41 they were unknown at the time of the array design. Furthermore, analysing the peak intensity
42
43 at these deletion sites in the main tiling, no obvious reduction in signal intensity was
44
45 observed, thus these deletions would have been missed (Supp. Figure S1b).
46
47
48
49
50
51

52
53 In our blind test, an overall pickup rate of 72.9% was observed. Whilst this pickup
54
55 rate is low, stratifying the mutation types, we noticed all 35 detected variations were single
56
57 nucleotide variations suggesting a pickup rate of 94.3% (95% CI 81.4-99.4%) for this
58
59 variation type. However, this also meant that no insertion or deletion mutations were
60

1
2
3 identified. After the analysis was undertaken, the insertion and deletion mutations were
4
5 compared against the extra tiling and it transpires that only 1 of the deletions had been
6
7 present in the HGMD at the time of the array design and had been tiled. All other insertions
8
9 and deletions blind tested had not been tiled and must therefore be classed as novel. It would
10
11 be expected that any novel heterozygous insertion and deletion would be missed by the array
12
13 technology. Indeed in our blind analysis, we failed to detect any heterozygous insertions or
14
15 deletions. However, novel homozygous deletions and insertions should be detected. We did
16
17 detect an untiled 5bp deletion in *C14ORF133* that was in a homozygous state. Deletions will
18
19 manifest as no-called regions, whilst insertions may appear as compound heterozygous
20
21 samples with heterozygous calls a few base pairs apart. Consequently, consideration must be
22
23 given during the design phase of the microarray resequencing chip to include all known
24
25 insertions and deletions irrespective of the size of the change to allow adequate detection.
26
27 Further, it would be prudent to update the resequencing chip with newer insertion and
28
29 deletions when any re-order is performed.
30
31
32
33
34
35

36 37 *Comparison to di-deoxysequencing.*

38
39
40 To see how the BRUM1 array compares to traditional di-deoxysequencing, surplus
41
42 products from microarray analysis were sequenced by direct sequencing. Using automated
43
44 analysis, 93% of sequenced bases were in agreement with the BRUM1 array sequencing.
45
46 There were 39 bases that were truly discordant between the two technologies. Comparison of
47
48 the BRUM1 sequencing and di-deoxysequencing to the reference sequence indicated that
49
50 these discordant bases were concordant using BRUM1 sequencing, but discordant with di-
51
52 deoxysequencing indicating that the di-deoxysequencing was wrong. Interestingly, these
53
54 bases were confined to the start of the sequencing traces when di-deoxysequencing doesn't
55
56 perform well. The primary cause of this was the primers used were not designed with di-
57
58
59
60

1
2
3 deoxysequencing as the primary analysis technique. Thus the location of the primers may be
4
5 too close to the start of the sequence of interest to obtain clean sequence. There were also
6
7 heterozygous calls made in the di-deoxysequencing that were homozygous calls in BRUM1
8
9 array sequencing. These incorrect heterozygous calls were due to influence of neighbouring
10
11 bases raising the secondary peak height above the threshold limit.
12
13
14
15

16 Following manual review of the data, the agreement between the di-deoxysequencing
17
18 and BRUM1, and di-deoxysequencing and reference sequence increased to 99% and 98%
19
20 respectively. Whilst no discordancy was observed between BRUM1 and di-deoxysequencing,
21
22 discordancy was observed when di-deoxysequencing was compared against the reference
23
24 sequence. However this was to be expected as 28 mutations resided in these fragments
25
26 sequenced.
27
28
29
30
31
32
33

34 In conclusion, we performed a stringent analysis of the BRUM1 array using DNA of
35
36 patients with known mutations and found that it had an overall 96% pick up rate for single
37
38 base changes, 50% for homozygous insertions and deletions and 39% pick up rate for
39
40 heterozygous insertion/ deletions. In general we acknowledge that the current resequencing
41
42 technology will not be able to detect a significant proportion of novel insertions and
43
44 deletions. These problems can be partly overcome by incorporating all known insertion and
45
46 deletions into the extra tiling and regularly updating the extra tiling with recently identified
47
48 insertion/ deletions. Clinically, this could mean that a patient with novel heterozygous
49
50 insertion/ deletions would be missed. BRUM1 array probe analysis was undertaken using
51
52 GSeq4.1. A number of other genes have been sequenced using this array (data not shown),
53
54 with comparable call rates and accuracy and positive identification of known mutations. Thus
55
56 we feel that the statistics generated would translate to the rest of the array. There are other
57
58
59
60

1
2
3 software packages available to analyse Affymetrix resequencing arrays that are currently
4
5 being assessed for performance (Schroeder et al., 2009).
6
7

8
9 We feel the BRUM1 array can be recommended as a cost-effective, fast, transferrable
10 and reliable screening tool in patients whose phenotype can be caused by mutations in one of
11 several genes. We acknowledge that before the technology can be successfully implemented
12 as a screening tool further evaluation needs to be performed in the clinical setting. This is
13 currently being undertaken. Our estimates suggest that the price of the additional gene
14 sequencing performed on the same array will not be significantly increased compared with
15 the basic price. This price and time reduction is largely achieved by using long range PCR to
16 amplify the genomic DNA fragments. In addition, biochemical testing such as tissue enzyme
17 diagnosis could be avoided if pathogenic mutations are detected. Therefore sequencing of
18 several genes per patient and also using the array for several patients (in whom different
19 genes will be sequenced) will make this method extremely cost-effective. For example we
20 estimate that screening 7 genes by microarray resequencing will be significantly cheaper and
21 faster than sequencing two genes by the conventional di-deoxysequencing methods. Thus,
22 this methodology can be used in clinical practice but it is extremely important to bear in mind
23 the limitations in relation to the detection of novel insertions and deletions. Further, regular
24 updates to the extra tiling to facilitate the detection of recently identified insertion and
25 deletions would be advantageous.
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 **Acknowledgements.** 52

53
54 We wish to thank Birmingham Children's Hospital Research Foundation, The
55 Children's Liver Fund and WellChild for funding this research. We would also like to
56
57
58
59
60

1
2
3 acknowledge educational grants from Actelion and Ferring Pharmaceuticals. PG is a GSK
4
5 Clinician Scientist Fellow.
6
7

8 9 **References**

- 10
11
12
13
14
15 Corless CE, Kaczmarek E, Borrow R, Guiver M. 2008. Molecular characterization of
16
17 *Neisseria meningitidis* isolates using a resequencing DNA microarray. *J Mol Diagn*
18
19 10:265-271.
20
21
22
23 Cullinane AR, Straatman-Iwanowska A, Zaucker A, Wakabayash Y, Bruce C, Guanmei L,
24
25 Rahman F, Gurakan F, Utine E, Ozkan TB, Denecke J, Vukovic J, Rocco MD,
26
27 Mandel H, Cangul H, Matthews RP, Thomas SG, Rappoport J, Arias IM, Wolburg H,
28
29 Knisely AS, Kelly DA, Mueller F, Maher ER, Gissen P. 2010. Mutations in *VIPAR*
30
31 cause an arthrogyrosis, renal dysfunction, and cholestasis syndrome phenotype with
32
33 defects in epithelial polarisation. *Nat Genet* In Press.
34
35
36
37
38 Cullinane AR, Straatman-Iwanowska A, Zaucker A, Wakabayash Y, Bruce C, Rahman F,
39
40 Rappoport J, Arias IM, Wolburg H, Knisely AS, Kelly DA, Mueller F, Maher ER,
41
42 and Gissen P. 2009. Mutations in polarin (*PLRN*) cause an ARC syndrome phenotype
43
44 and defects in epithelial polarisation and apical junction complex formation (Abstract
45
46 290). 59th Annual Meeting of The American Society of Human Genetics. Honolulu,
47
48 Hawaii.
49
50
51
52
53
54 Cutler DJ, Zwick ME, Carrasquillo MM, Yohn CT, Tobin KP, Kashuk C, Mathews DJ, Shah
55
56 NA, Eichler EE, Warrington JA, Chakravarti A. 2001. High-throughput variation
57
58 detection and genotyping using microarrays. *Genome Res* 11:1913-1925.
59
60

- 1
2
3 Denning L, Anderson JA, Davis R, Gregg JP, Kuzdenyi J, Maselli RA. 2007. High
4 throughput genetic analysis of congenital myasthenic syndromes using resequencing
5 microarrays. PLoS One 2:e918.
6
7
8
9
10
11 Di X and Cawley S. 2005. Alternative base calling method for resequencing microarrays.
12 Conf Proc IEEE Eng Med Biol Soc. 2809 2812.
13
14
15
16
17 Fleiss JL. 1971. Measuring nominal scale agreement among many raters. Psychological
18 Bulletin 76:378-382.
19
20
21
22
23 Gissen P, Johnson CA, Morgan NV, Stapelbroek JM, Forsheew T, Cooper WN, McKiernan
24 PJ, Klomp LW, Morris AA, Wraith JE, McClean P, Lynch SA, Thompson RJ, Lo B,
25 Quarrell OW, Di RM, Trembath RC, Mandel H, Wali S, Karet FE, Knisely AS,
26 Houwen RH, Kelly DA, Maher ER. 2004. Mutations in VPS33B, encoding a
27 regulator of SNARE-dependent membrane fusion, cause arthrogryposis-renal
28 dysfunction-cholestasis (ARC) syndrome. Nat Genet 36:400-404.
29
30
31
32
33
34
35
36
37
38 Hacia JG. 1999. Resequencing and mutational analysis using oligonucleotide microarrays.
39 Nat Genet 21:42-47.
40
41
42
43
44 Hartmann A, Thieme M, Nanduri LK, Stempfl T, Moehle C, Kivisild T, Oefner PJ. 2009.
45 Validation of microarray-based resequencing of 93 worldwide mitochondrial
46 genomes. Hum Mutat 30:115-122.
47
48
49
50
51
52 Jalanko A, Braulke T. 2009. Neuronal ceroid lipofuscinoses. Biochim Biophys Acta
53 1793:697-709.
54
55
56
57
58
59
60

- 1
2
3 Karaman MW, Groshen S, Lee CC, Pike BL, Hacia JG. 2005. Comparisons of substitution,
4 insertion and deletion probes for resequencing and mutational analysis using
5 oligonucleotide microarrays. *Nucleic Acids Res* 33:e33.
6
7
8
9
10
11 King JE. 2004. Software solutions for obtaining a kappa-type statistic for use with multiple
12 raters. Annual Meeting of the Southwest Educational Research Association. Dallas,
13 Texas.
14
15
16
17
18
19 Kothiyal P, Cox S, Ebert J, Aronow BJ, Greinwald JH, Rehm HL. 2009. An overview of
20 custom array sequencing. *Curr Protoc Hum Genet* 61:17.1-17.11.
21
22
23
24
25 Liu C, Aronow BJ, Jegga AG, Wang N, Miethke A, Mourya R, Bezerra JA. 2007. Novel
26 resequencing chip customized to diagnose mutations in patients with inherited
27 syndromes of intrahepatic cholestasis. *Gastroenterology* 132:119-126.
28
29
30
31
32
33 Malanoski AP, Lin B, Wang Z, Schnur JM, Stenger DA. 2006. Automated identification of
34 multiple micro-organisms from resequencing DNA microarrays. *Nucleic Acids Res*
35 34:5300-5311.
36
37
38
39
40
41 Neufeld EF, Muenzer J. 2001. The Mucopolysaccharidoses. In: Scriver CR, Beaudet AL, Sly
42 WS, Valle D, editors. *The metabolic and molecular bases of inherited disease*. New
43 York: Mcgraw-Hill. p 3241-3251.
44
45
46
47
48
49 Sauro J and Lewis JR. 2005. Estimating completion rates from small samples using binomial
50 confidence intervals: comparisons and recommendations. PROCEEDINGS of the
51 HUMAN FACTORS AND ERGONOMICS SOCIETY 49th ANNUAL MEETING.
52 2100 2104.
53
54
55
56
57
58
59
60

1
2
3 Schroeder C, Stutzmann F, Weber BH, Riess O, Bonin M. 2009. High-throughput
4 resequencing in the diagnosis of BRCA1/2 mutations using oligonucleotide
5 resequencing microarrays. *Breast Cancer Res Treat.*
6
7
8
9

10
11 Suojanen JN. 1999. False false positive rates. *N Engl J Med* 341:131.
12
13

14
15 Takahashi Y, Seki N, Ishiura H, Mitsui J, Matsukawa T, Kishino A, Onodera O, Aoki M,
16
17 Shimosawa N, Murayama S, Itoyama Y, Suzuki Y, Sobue G, Nishizawa M, Goto J,
18
19 Tsuji S. 2008. Development of a high-throughput microarray-based resequencing
20 system for neurological disorders and its application to molecular genetics of
21 amyotrophic lateral sclerosis. *Arch Neurol* 65:1326-1332.
22
23
24
25
26

27
28 Wraith JE, Baumgartner MR, Bembi B, Covanis A, Levade T, Mengel E, Pineda M, Sedel F,
29
30 Topcu M, Vanier MT, Widner H, Wijburg FA, Patterson MC. 2009.
31
32 Recommendations on the diagnosis and management of Niemann-Pick disease type
33 C. *Mol Genet Metab* 98:152-165.
34
35
36

37
38 Zwick ME, Kiley MP, Stewart AC, Mateczun A, Read TD. 2008. Genotyping of *Bacillus*
39
40 *cereus* strains by microarray-based resequencing. *PLoS One* 3:e2513.
41
42
43
44
45

46 **Figure Legends**

47
48 Figure 1: False negative call in *NPCI*.

49
50 Forward strand shows a heterozygous call with peaks of C and T being equal and prominent.

51
52 Reverse strand shows a homozygous calls with a peak of A.
53
54
55
56
57
58
59
60

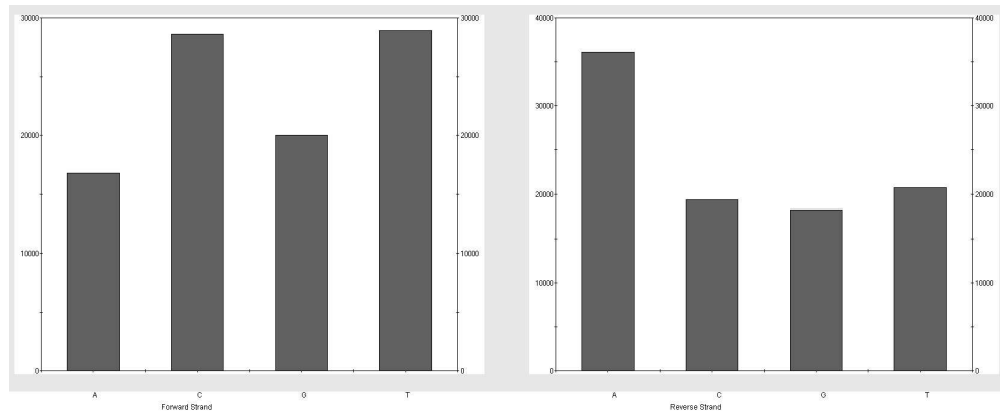


Figure 1: False negative call in NPC1.
Forward strand shows a heterozygous call with peaks of C and T being equal and prominent.
Reverse strand shows a homozygous calls with a peak of A.

140x57mm (300 x 300 DPI)

Peer Review

QST	Average Call Rate		Average Accuracy	
	BRT0	BRT 0.5	BRT0	BRT0.5
2	94.25% (914)	93.55% (1025)	99.45% (87)	99.51% (78)
3	90.70% (1478)	89.42%(1681)	99.53% (75)	99.58% (67)
6	74.52% (4049)	71.71% (4495)	99.72% (44)	99.76% (38)
9	53.67% (7362)	49.01% (8102)	99.82% (29)	99.85% (24)
12	33.83% (10 514)	29.23% (11 245)	99.89% (17)	99.91% (14)
30	0.42% (15 223)	0.22% (15 540)	100% (0)	100% (0)

Table 1 Call rate and accuracy of the sequencing.

Figures in parenthesis refer to the actual numbers of no-called bases for Average Call rate or the accurate numbers of incorrect bases for Average Accuracy; QST =quality score threshold, BRT=base reliability threshold.

For Peer Review

	PCR A	PCR B	PCR C	All Chips
pA	0.917	0.968	0.973	0.948
Kappa	0.895	0.959	0.966	0.934
LCI	0.891	0.955	0.962	0.933
UCI	0.898	0.962	0.970	0.935
No of bases different between samples	1319	508	429	826

Table 2: Reproducibility

The output from each BRUM1 array was compared within each PCR and between all samples. Fleiss Kappa was used to compare the outputs. pA is the level of agreement between the samples, whilst LCI and UCI are the lower and upper 95% confidence intervals.

For Peer Review

No Called base/ region	A	G	C	T	AT	>2CC	>2GG	Other*
Frequency	71	58	64	89	58	58	17	41
% of all no-calls	15.6%	12.7%	14.0%	19.2%	12.7%	12.7%	3.7%	9.0%

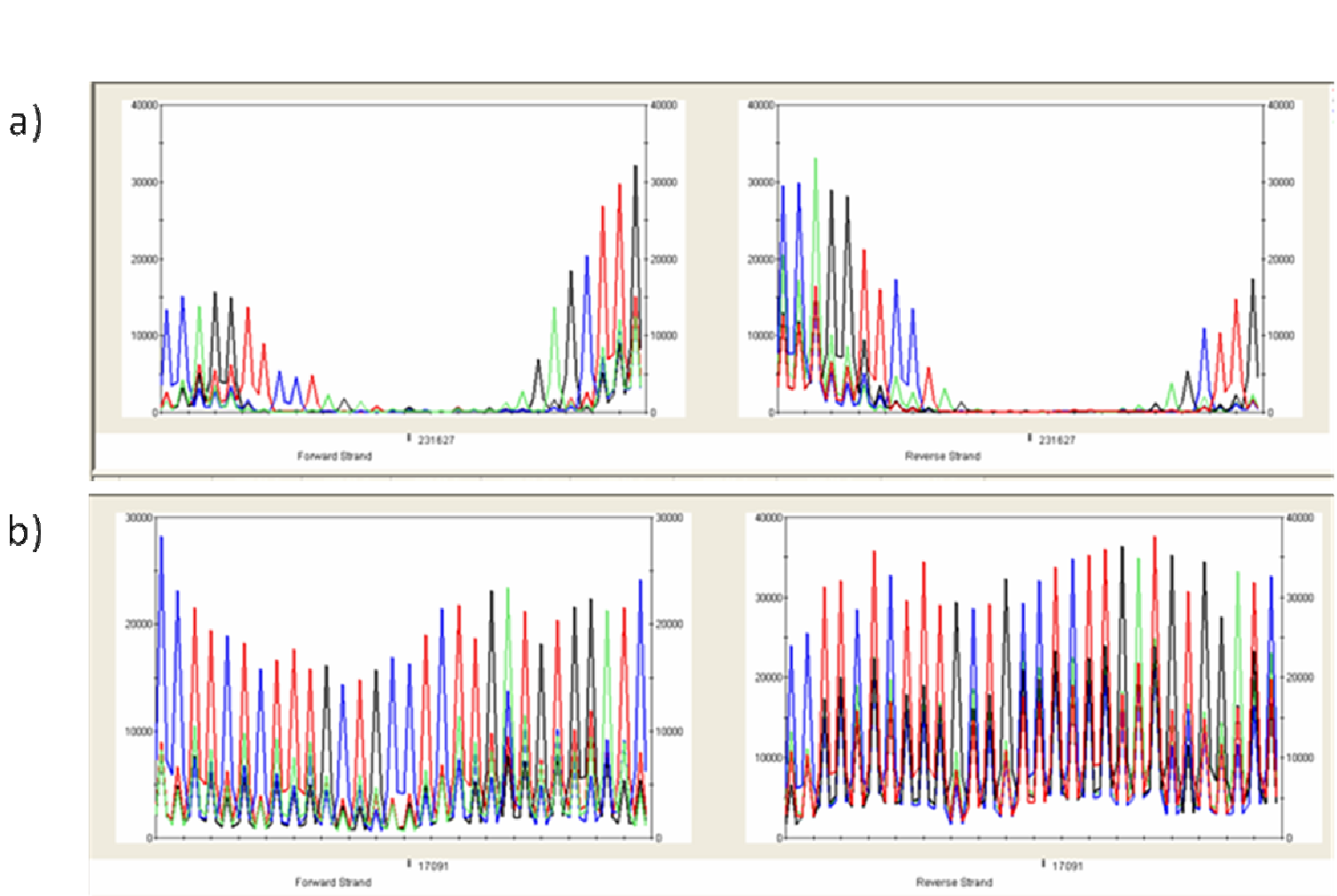
Table 3: Consistent regions of No-Call

*Other refers to no-called regions observed less than twice.

	Point Changes	Insertions	Deletions
Number of changes analysed	39 (76)	3 (6)	10 (18)
Number of changes detected	32 (65)	0 (0)	4 (4)
False Positives	3 (3)	-	-
False Positive rate	<0.01% (<0.01%)	-	-
False Negatives	7 (9)	3 (6)	5 (13)
False Positive following manual review	1 (3)	-	-
False negative following manual review	1 (3)	0 (3)	3 (11)
Novel changes	5		
Pickup Rate	97% (96%)	100% (50%)	70% (39%)
95% CI	82-99% (88.5-99.1)	47-100% (19-81%)	39-90% (20-61%)

Table 4: Summary of Changes

False positives were not calculated for insertions and deletions. Any no-called base has to be classed as a false positive deletion and thus needs manual review. Figures in parenthesis are combined data including additional *NPCI* variations used in the blind study. Novel changes were known polymorphisms that were detected but were not expected.



35 Supp. Figure S1: Deletion detection: 30bp of sequence incorporating two known deletions were
36 captured in GSeq4.1. When deletions are in a homozygous state (panel a; 5bp deletion in
37 *C14ORF133*) there is a drop in signal intensity on both strands indicating a deletion. However,
38 when deletions are in a heterozygous state (panel b; a 12bp heterozygous deletion in *VPS33B*) no
39 reduction in signal intensity is observed.
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Gene		Gene		Gene		Gene		Gene		Gene		Gene	
C14ORF133	Patient	GAA	Patient	NPC1	Patien	NPC2	Patient	SLC19A2	Patient	VPS33B	Patient	WFS1	Patient
1	ARC-18	2	GAA-33	1	MRC5VA	1	NPC-9	1	SLC-1	1	MRC5VA	2	MRC5VA
2	MRC5VA	3	GAA-24	2	MRC5VA	2	MRC5VA	2	SLC-2	2	ARC-8	3	WFS1-1
3	MRC5VA	4	GAA-36	3	MRC5VA	3	MRC5VA	3	MRC5VA	3	ARC-21	4	MRC5VA
4	MRC5VA	5	MRC5VA	4	MRC5VA	4	MRC5VA	4	MRC5VA	4	ARC-9	5	WFS1-2
5	MRC5VA	6	GAA-31	5	MRC5VA	5	MRC5VA	5	MRC5VA	5	ARC-10	6	WFS1-3
6	MRC5VA	7	GAA-3	6	MRC5VA			6	MRC5VA	6	ARC-11	7	MRC5VA
7	ARC-2	8	MRC5VA	7	MRC5VA					7	ARC-12	8a	WFS1-4
8	MRC5VA	9	GAA-15	8	NPC-5					8	ACR-13	8b	WFS1-4
9	ARC-3	10	GAA-21	9	NPC-2					9	MRC5VA	8c	WFS1-5
10	ARC-4	11	GAA-21	10	MRC5VA					10	MRC5VA	8d	WFS1-6
11	MRC5VA	12	GAA-16	11	MRC5VA					11	MRC5VA	8e	MRC5VA
12	ARC-5	13	GAA-27	12	NPC-20					12	ARC-14		
13	MRC5VA	14	GAA-18	13	NPC-20					13	ARC-15		
14	MRC5VA	15	GAA-8	14	NPC-21					14	MRC5VA		
15	MRC5VA	16	MRC5VA	15	MRC5VA					15	ARC16		
16	ARC-6	17	MRC5VA	16	MRC5VA					16	ARC16		
17	ARC-7	18	GAA-14	17	NPC-20					17	MRC5VA		
18	MRC5VA	19	GAA-4	18	NPC-20					18	ARC-17		
19	MRC5VA			19	NPC-7					19	ARC-1		
				20	NPC-6					20	ARC-19		
				21	NPC-6					21	ARC-20		
				22	NPC-12					22	ARC-20		
				23	NPC-11					23	MRC5VA		
				24	MRC5VA								
				25	MRC5VA								

Supp Table S2: Source of DNA used for the array

Lysosomal Storage	Lysosomal Storage	Glycogen Storage	Cholestasis	Iron-deposition/
<i>NPC1</i> (607623) ENST00000269228	<i>NAGLU</i> (609701)	<i>G6PC</i> (232200)	<i>VPS33B</i> (608552) ENST00000333371	<i>PANK2</i> (606157)
<i>NPC2</i> (601015) ENST00000238633	<i>HGSNAT</i> (610453)	<i>SLC37A4</i> (232220)	<i>ATP8B1</i> (602397)	<i>PLA2G6</i> (603604)
<i>PSAP</i> (607939)	<i>GNS</i> (607664)	<i>SLC2A1</i> (138140)	<i>ABCB11</i> (603201)	
<i>SUMF1</i> (176801)	<i>GALNS</i> (612222)	<i>SLC2A2</i> (138160)	<i>ABCB4</i> (171060)	
<i>SMPD1</i> (607608)	<i>ARSB</i> (611542)	<i>AGL</i> (610860)	<i>JAG1</i> (601920)	
<i>GBA</i> (606463)	<i>GUSB</i> (611499)	<i>GBE1</i> (607839)	<i>ATP7B</i> (606882)	
<i>GLB1</i> (611458)	<i>HYALI</i> (607071)	<i>PYGM</i> (608455)	<i>SLC25A13</i> (603859)	
<i>HEXB</i> (606873)	<i>DYM</i> (607461)	<i>PYGL</i> (608455)	<i>C14ORF133</i> ENST00000327028	
<i>HEXA</i> (606689)	<i>AGA</i> (208400)	<i>PFKM</i> (610681)		
<i>ASAHI</i> (228000)	<i>MAN2B1</i> (609458)	<i>PHKA2</i> (306000)		
<i>GNPTG</i> (607838)	<i>MANBA</i> (609489)	<i>PHKB</i> (172490)		
<i>GNPTAB</i> (607840)	<i>ABHD5</i> (604780)	<i>PHKG2</i> (172471)		
<i>MCOLN1</i> (605248)	<i>PNPLA2</i> (609059)	<i>GYS2</i> (138571)		
<i>SLC17A5</i> (604322)	<i>LAMP2</i> (309060)	<i>PRKAG2</i> (602743)		
<i>ARSA</i> (607574)	<i>NEU1</i> (608272)	<i>ALDOC</i> (103870)		
<i>GLA</i> (300644)				
<i>GAA</i> (606800) ENST00000302262				
<i>IDUA</i> (252800)				
<i>IDS</i> (309900)				
<i>SGSH</i> (605270)				

Supp. Table S2; Genes included in the array. Figures in parenthesis are OMIM reference numbers. ENST numbers are Ensembl accession numbers (www.ensembl.org).

Endocrine Malignancy	Growth	Diabetes	Obesity	Lipid Disorders
<i>VHL</i> (608537)	<i>GHI</i> (139250)	<i>WFS1</i> (606201) ENST0000022670	<i>LEP</i> (164160)	<i>LDLR</i> (606945)
<i>SDHB</i> (185470)	<i>IGF1</i> (147440)	<i>CISD2</i> (604928)	<i>LEPR</i> (601007)	<i>APOB</i> (107730)
<i>SDHC</i> (602413)	<i>IGF2</i> (147470)	<i>SLC19A2</i> (603941) ENST00000236137	<i>MC3R</i> (155540)	<i>OLR1</i> (602601)
<i>SDHD</i> (602690)	<i>PTPN11</i> (176876)	<i>EIF2AK3</i> (604032)	<i>GHRL</i> (605353)	<i>PCSK9</i> (607786)
<i>RET</i> (164761)	<i>RAF1</i> (164760)		<i>CD36</i> (173510)	<i>LDLRAP1</i> (600073)
				<i>APOE</i> (107741)
				<i>USF1</i> (191523)
				<i>CETP</i> (118470)
				<i>ABCA1</i> (600046)
				<i>ABCG5</i> (605459)
				<i>ABCG8</i> (605460)
				<i>CYP27A1</i> (606530)
				<i>MTP</i> (157147)

Supp. Table S2 continued.