

## ARTICLE

# Experiences with array-based sequence capture; toward clinical applications

Rowida Almomani<sup>1</sup>, Jaap van der Heijden<sup>1</sup>, Yavuz Ariyurek<sup>1,2</sup>, Yuching Lai<sup>2</sup>, Egbert Bakker<sup>1</sup>, Michiel van Galen<sup>1,2</sup>, Martijn H Breuning<sup>1</sup> and Johan T den Dunnen<sup>\*,1,2</sup>

Although sequencing of a human genome gradually becomes an option, zooming in on the region of interest remains attractive and cost saving. We performed array-based sequence capture using 385K Roche NimbleGen, Inc. arrays to zoom in on the protein-coding and immediate intron-flanking sequences of 112 genes, potentially involved in mental retardation and congenital malformation. Captured material was sequenced using Illumina technology. A data analysis pipeline was built that detects sequence variants, positions them in relation to the gene, checks for presence in databases (eg, db single-nucleotide polymorphism (SNP)) and predicts the potential consequences at the level of RNA splicing and protein translation. In the samples analyzed, all known variants were reliably detected, including pathogenic variants from control cases and SNPs derived from array experiments. Although overall coverage varied considerably, it was reproducible per region and facilitated the detection of large deletions and duplications (copy number variations), including a partial deletion in the *B3GALT1* gene from a patient sample. For ultimate diagnostic application, overall results need to be improved. Future arrays should contain probes from both DNA strands, and to obtain a more even coverage, one could add fewer probes from densely and more probes from sparsely covered regions.

*European Journal of Human Genetics* (2011) 19, 50–55; doi:10.1038/ejhg.2010.145; published online 24 November 2010

**Keywords:** capture array; heterogeneous disorders; sequencing

## INTRODUCTION

For many years, the amplification of target sequences by PCR, followed by Sanger sequencing, has been the gold standard for screening of variants in terms of both read length and accuracy of sequencing.<sup>1</sup> However, when it comes to conditions with highly heterogeneous etiology, a large number of different genes need to be screened for mutations. In such cases, gathering information becomes laborious, expensive and time-consuming. There are many examples of diseases that can be caused by mutations in many different genes, including mental retardation (MR),<sup>2</sup> Charcot–Marie–Tooth disease,<sup>3</sup> cardiomyopathy,<sup>4</sup> retinitis pigmentosa,<sup>5</sup> autism,<sup>6</sup> hearing loss<sup>7</sup> and congenital disorders of glycosylation.<sup>8</sup> Extensive resequencing of many disease-associated genes is required to explore, at the sequence and structural level, the genomic variation that might be involved in causing such diseases.

Several next-generation sequencing (NGS) platforms are now available and they have allowed the sequencing and analysis of large numbers of genes in one experiment,<sup>9–11</sup> and are able to generate a massive amount of sequence data and have considerably reduced the cost of DNA sequencing.<sup>12</sup> However, although NGS platforms have enormously increased throughput and have permitted whole-genome sequencing, high cost still prevents routine whole human genome resequencing projects. Therefore, zooming in on the region of interest is an attractive option. In addition, it circumvents the problem of identifying variants in genes for which the analyses were not intended (with associated ethical problems).

Microarray-based genomic selection combined with massively parallel high-throughput sequencing is the method of choice to analyze large numbers of genes in a more comprehensive and cost-effective manner.<sup>13–15</sup> We have used custom high-density microarrays (Roche NimbleGen, Inc., Madison, WI, USA) for the enrichment of 112 distinct genes potentially involved in MR and congenital malformation, followed by sequencing on the Illumina Genome Analyzer I platform (Illumina, San Diego, CA, USA).

The first aim of our study was to apply and validate the array-based enrichment method as an efficient and convenient strategy to capture any desired portion of the human genome. The second aim was to accelerate the detection of sequence and copy number variations (CNV) in the selected candidate genes with lower costs, especially for the genes that are potentially involved in MR.

## MATERIALS AND METHODS

### Sample selection and validation

Six DNA samples were used in this study, including two controls containing known pathogenic variants. Sample S-2 contains a known *MECP2* (OMIM 300005) pathogenic point mutation (c.538C>T); the second sample, patient S-6, carries a large deletion spanning exons 8–15 in one allele and a splice site mutation (c.660+1G) at the other allele of the *B3GALT1* (OMIM 610308) gene.

The other four DNA samples were from patients with MR with an unknown cause. Single-nucleotide polymorphism (SNP) array data were available for two samples: S-7 with 250K Nsp Affymetrix and S-5 with 317K Illumina data. We used these data to validate the sequences obtained after capture-array and

<sup>1</sup>Center for Human and Clinical Genetics, Leiden University Medical Center, Leiden, The Netherlands and <sup>2</sup>Leiden Genome Technology Center, Leiden University Medical Center, Leiden, The Netherlands

\*Correspondence: Professor Dr JT den Dunnen, Center for Human and Clinical Genetics and Leiden Genome Technology Center, Leiden University Medical Center, Postzone S4-P, P.O. Box 9600, 2300 RC, Leiden, The Netherlands. Tel: +31 71 5269501; Fax: +31 71 526 8285; E-mail: ddunnen@HumGen.nl

Received 22 April 2010; revised 16 July 2010; accepted 20 July 2010; published online 24 November 2010

Illumina sequencing. Causative large deletions and duplications had been previously excluded by SNP array testing in S-3, S-5, S-7 and S-8.

### Exon array design

Microarrays with 385K probe capacity (Roche NimbleGen, Inc.) were used to capture all exons, the splice site and the immediately adjacent intron sequence of 112 human genes. On the basis of searches in OMIM and literature, we selected 112 human genes known to cause MR, either as part of a known syndrome or in isolation (Supplementary Table 1). Primary sequence data from all exons were extracted from NCBI's genome (Build 36). Microarrays were designed by Roche NimbleGen, Inc. with long oligonucleotide probes (54–99 nucleotides) that span each target region, overlapped and shifted on an average of seven bases.<sup>13</sup> The oligonucleotides were designed to achieve isothermal hybridization across the arrays capturing one strand only. All highly repetitive regions were excluded from the probe selection in order to avoid nonspecific capturing of genomic regions. Using all criteria listed, for 2% of the target sequences, no capture probe could be designed (note that, theoretically, these sequences can be covered partly through capture from directly flanking unique sequences). Four of the arrays were reused at least twice.

### Genomic DNA library preparation and target capture

The methods used for target capture, enrichments and elution followed previously described protocols with slight modifications (Roche NimbleGen, Inc.).<sup>16</sup> Genomic DNA (20–10 µg) was fragmented using a nebulizer or Bioruptor according to instructions from the manufacturer to yield fragments from 250–1000 bp (nebulization) or 250–600 bp (Bioruptor). Adapter oligonucleotides from Illumina (single reads) were ligated to the ends. After the ligation was completed, successful adapter ligation was confirmed by PCR. The DNA-adapter ligated fragments were then hybridized to the sequence capture microarray for 65 h. After hybridization and washing, the DNA fragments bound to the array were eluted, using 300 µl of the elution buffer (Qiagen, Valencia, CA, USA) on each array. A gasket (Agilent) was applied and placed on the thermal elution device (homemade) for 20 min at 95°C. We repeated this process once by adding 200 µl of elution buffer (Qiagen). DNA from each eluted sample was enriched by 18-cycle PCR using a high-fidelity polymerase and a single primer pair corresponding to the Illumina adapters ligated earlier.

### Check enrichments by qPCR

To verify successful hybridization capture, we performed qPCR (quantitative PCR) on DNA samples (S-2, S-3, S-5, S-7, S-6 and S-8) before and after array enrichment. The primers amplified five loci from *MBL2*, *DMD* and *BRCA1* (100 bp) as negative controls (no capture probes on the array) and four loci from *MECP2*, *CREBBP* and *NSD1* genes as positive controls (capture probes on the array) (Supplementary Table 2). All primers for qPCR were designed using Primer 3 (<http://frodo.wi.mit.edu/>).

The qPCR assays were performed in triplicate in the Lightcycler using 384-well plates (Roche NimbleGen, Inc.) in 10 µl total volume: 5 µl of 2× SYBR Green master Rox (Roche NimbleGen, Inc.), 0.25 µl of each primer (10 pmol/µl), 2 µl of DNA template and 2.5 µl of ultrapure water. The thermo-cycling protocol was carried out as follows: 10 min at 95°C, 45 cycles of 10 s at 95°C, 30 s at 60°C, 20 s at 72°C and 5 min at 72°C, followed by melting curve analysis in order to determine the specific and nonspecific amplified products and other artifacts that might interfere with CP values. To calculate the relative fold enrichment of the targeted regions, we compared amplification of the positive *versus* negative controls. The relative fold enrichment, *R*, was calculated using the values of ΔCP (ie, the difference between average CP of non-captured and average CP of captured samples) according to  $R = E^N$ , where *E* is the efficiency of the qPCR assay for a particular amplicon and *N* = ΔCP (crossing point).

### DNA sequencing

The eluted enriched DNA fragments were sequenced using the Illumina GAI platform at the Leiden Genome Technology Center (LGTC). Single-end sequencing of 36 or 50 nucleotides was performed following the instructions of the manufacturer.

### Reads mapping and data analysis

Sequence read mapping was carried out by ELAND and ELAND-extended programs, which were a part of the Illumina GAI data analysis package. Only reads of high-quality scores were mapped to the human reference genome (NCBI, BUILD 36.2), allowing up to two mismatches. We created different Perl scripts to extract and process data from the ELAND files. Coverage was calculated at the target level (gene–exons), the nucleotide level and at the per probe region. SNP calling was performed by searching for nucleotides discordant with the reference genome with a base call quality score of 30 (99.9% base call accuracy), a read depth of 8 or greater and the variant allele larger than 30% of the total coverage. Thereafter, all variants were checked for their presence in known databases, for example, dbSNP. Perl scripts were designed to predict the potential consequences at the level of RNA splicing and protein translation on the basis of Ensemble v.51. Furthermore, we designed a Perl script to facilitate detection of small deletions/insertions (up to three nucleotides). All Perl scripts are available on request.

### Sanger sequencing

A total of 21 variants detected by Illumina GAI analyzer were selected and confirmed by Sanger sequencing using the standard Sanger sequencing protocol at the Leiden Genome Technology Center (LGTC). The primer sequences (with M13 tail) used are shown in Supplementary Table 3.

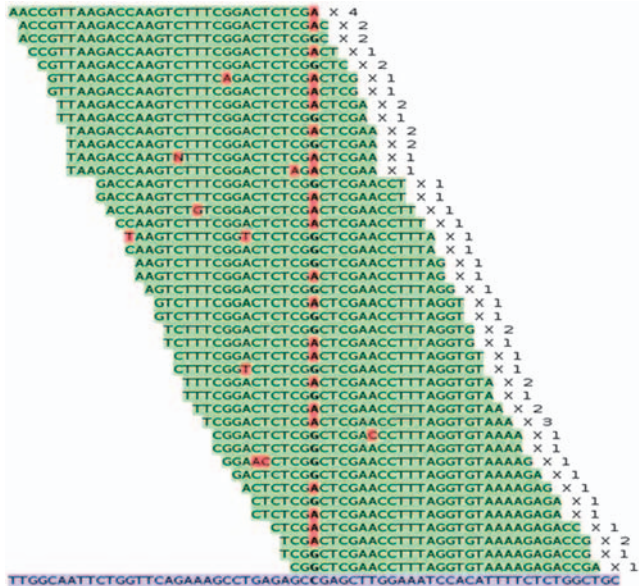
## RESULTS

The methodology used starts with fragmentation of the genomic DNA. Linker and primer addition can then be performed either before or after array-capture target enrichment. To facilitate limited amplification of the expected low-yield array elution, we decided to perform full Illumina sample preparation before array capture. Initially, experiments were conducted using 20 µg genomic DNA, later we reduced this to 10 µg. We used qPCR, comparing targeted (four positive controls) and non-targeted regions (five negative controls), to check successful array enrichment and to estimate the fold enrichment obtained (see Supplementary Tables 4 and 5 for examples). As enrichment varies significantly from locus to locus, we tested multiple loci to obtain an accurate estimate. Samples in which qPCR did not indicate clear enrichment (>100×) were discarded. The ultimate enrichments achieved varied from experiment to experiment with a tendency to increase over time, indicating that lab experience is an important aspect of the array capture technology. As the fold enrichments determined by qPCR correlate positively with the average sequence depth obtained, we conclude that qPCR provides an effective and cost-saving check for successful enrichment (examples are listed in Supplementary Tables 4 and 5).

### Sequence data

The custom arrays used contained 112 different human genes that are known to be or potentially involved in MR and congenital malformation. Samples were run on one channel of the Illumina GAI. For sequence analysis, we used only those QC-filtered reads that map back uniquely to the reference sequence (M0) or with one or two mismatches (M1, M2) (Figure 1). Using these settings, 85–92% of the targeted nucleotides were covered by at least eight reads (Table 1) and 94–98% by at least one read (note that for 2% of the targeted sequences, no probe could be designed, see M&M). Effectively, this means that for 78% of the targeted sequences on the array, coverage was sufficient (>20×) to detect any variants that were present.

Two of the samples had been previously analyzed using SNP arrays. The region selected using the capture array included 67 different SNPs that had been present on the SNP arrays. We observed a perfect agreement (100%) between array-based SNP calls and those obtained using NGS (67/67 variants) (Supplementary Table 6).



**Figure 1** Detection of sequence variants. A total of 32 nucleotide NGS reads (top, sequence mismatches in red) aligned with the genomic reference sequence (bottom). The center of the alignment shows a variant present in the heterozygous state. 'x n' behind the read indicates how many identical reads were obtained.

To determine our ability to detect pathogenic mutations, we included one sample from a female patient (S-2) harboring a dominant pathogenic point mutation in the *MECP2* gene, (c.538C>T) on the X chromosome. Our results clearly detected the change in the heterozygous state (Supplementary Table 7). Similarly, we detected a homozygous change in the *B3GALTL* gene in a Peter's Plus patient (c.660+1G>A, Supplementary Table 7, see below).

We next selected 21 variants detected in samples S-2, S-3, S-5, S-7 and S-8 and checked these by traditional Sanger sequencing. We were able to confirm 21 of the 21 variants, including their status being homozygous or heterozygous (Supplementary Table 7). The analysis of the variants found in all 112 genes of the patients did not reveal a clear cause of their MR Supplementary Table 8 and 9.

### CNV

Changes that cannot be easily detected using the sequence itself include deletions and duplications (CNVs). However, such variants can be expected to yield quantitative changes in coverage. To determine whether overall coverage can be used to detect quantitative changes, we first analyzed the 39 genes located on the X chromosome. Indeed, when coverage was normalized using autosomal genes (Figure 2a), samples from females showed a clearly higher X-chromosome coverage compared with male samples (Figure 2b). Furthermore, as expected, the gene on the Y chromosome (*NLGN4Y*) gave no coverage in the female sample (Figure 2b). To determine the sensitivity of our method for detecting smaller CNVs, we carefully analyzed a sample from a compound heterozygous patient (S-6) carrying a partial deletion (exons 8–15) and a splice site mutation (c.660+1G>A, intron 8) in the *B3GALTL* gene. The splice site mutation was evident as no wild-type sequence was present. The presence of a deletion emerged as, compared with other samples, we observed a significantly lower average coverage for the *B3GALTL* gene (53× versus 155×, 150×, 140×) (Figure 2c). In addition, although the splice site mutation in exon 8 was detected in the 'homozygous' state (similar to all nine variants downstream), we observed variants in the

first exons (1–7) also in heterozygous state (Supplementary Table 10). These data show that not only have we obtained an excellent specificity of the capture process but we have also been able to distinguish between male and female samples.

### DISCUSSION

Array-based genomic selection offers several advantages for large-scale targeted DNA isolation over other approaches such as PCR-based methods (long-range PCR or multiplexed short PCR),<sup>17–19</sup> selector technology<sup>20,21</sup> and BACs technology.<sup>22</sup> PCR-based methods become laborious, time-consuming and costly if hundreds to thousands of regions (exons) need to be amplified, especially if all the sequences are required. Furthermore, when PCRs are multiplexed, it becomes difficult to check successful amplification per fragment, the chance of obtaining artifacts increases and equimolar loading before sequencing becomes very difficult. New approaches for massive individual PCR have been introduced recently<sup>23</sup> but experiences with these are still limiting. Selector technology<sup>20,21</sup> seems attractive but it largely depends on proper in-house probe design, and experience thus far is very limited. Successful genomic selection using BACs has been demonstrated but has several limitations. As a BAC is the unit of selection, multiple BACs are required to isolate discontinuous regions of interest.

In this study, we have tested array-based sequence capture to determine the sequence of 112 genes potentially involved in MR. We show that array-based sequence capture technology is an efficient, quick and reliable method for the parallel sequencing of a range of genes of interest. Known variants (array-based calls) for 67 SNPs matched perfectly with those obtained using NGS Supplementary Table 6. Two positive controls with known pathogenic changes in the *MECP2* gene (sample S-2) and *B3GALTL* gene (sample S-6) were readily detected. In addition, 21/21 selected variants found in the five samples analyzed could be confirmed using Sanger sequencing (Supplementary Table 7). Sequence coverage of the nucleotide of interest is critical for reliably detecting sequence changes. If coverage is too low, both false positives (caused by sequence errors) and false negatives (if only one allele from a heterozygous sample is observed) will occur.

The coverage we obtained differs significantly not only between targeted genomic regions (genes) but also between different samples (Supplementary Table 1, Figure 2a). As the overall methodology is rather complex, particularly the collection of the hybridized array-enriched DNA sequences, the difference between samples is most probably influenced by technical factors such as variations in hybridization, washing conditions and potential reuse of the capture array. Furthermore, coverage is influenced by array design, including probe sequence (melting temperature, GC content), probe density and spacing (Supplementary Table 1). Our data show that AT-rich regions (>55%), regions with an overall low probe density (<3) and small exons (on average 90bp) yield a low coverage, which also varies significantly between experiments. For a second-generation capture array, the results obtained could be used to change the probe density, that is, decreased in well-covered and increased in low-covered regions.

Our data show that longer reads (50bp) improve accuracy and selectivity of read mapping to the reference genome, which influenced the SNP calling by having less false positives and slightly better coverage.

As CNVs (deletions/duplications) are a significant cause in the etiology of MR,<sup>24</sup> we tested the feasibility of detecting large CNVs using array capture and NGS. Our results indicate that, if coverage is sufficiently high, array capture can also be used to detect such

**Table 1** Sequence summary results of the different array-capture experiments performed

Sample ID, sex	Total reads $\times 10^3$	Reads passing QC filter $\times 10^3$	Total number of reads mapped $\times 10^3$	MM0 reads $\times 10^3$	MM1 reads $\times 10^3$	MM2 reads $\times 10^3$	Coverage per nucleotide	% of Nucleotides were covered $\geq 8$ times	% of nucleotides were covered 0 times	Read length	Array reused
S-2, F	6.744	4.804	2.428	1.359	691	378	138	87.11	6.22	50	No
S-3, M	7.305	5.354	2.176	1.225	618	333	100	90.71	4.49	50	No
S-5, M	10.43	7.237	5.576	4.935	499	142	120	92.42	2.09	32	No
S-7, M	15.771	6.112	4.719	3.885	638	196	100	91.13	2.70	32	No
S-6, M	12.154	6.575	6.575	5.914	486	174	99	99.24	7.08	32	Yes, 2nd time
S-8, F	11.077	3.531	3.531	2.301	736	485	44	85.38	4.43	49	Yes, 3rd time

Abbreviations: F, female; M, male; MM# reads, number of reads with # mismatches to the reference sequence; QC, quality control.

quantitative changes. Our array contained one gene from the Y chromosome that gave no coverage in females (Figure 2b), whereas the 39 X-linked genes when compared with the 69 autosomal genes yielded overall 50% lower coverage in male samples (Figure 2b). Another example derives from a sample containing a partial *B3GALT1* gene deletion on one allele (exons 8–15) and a splice site mutation on the other allele (c.660+1G>A). Although coverage over the entire gene seems reduced (experimental variation/coincidence), coverage for the second half of the gene clearly drops below that of normal (Figure 2d). An algorithm for detecting local deviations from the average coverage is currently under development.

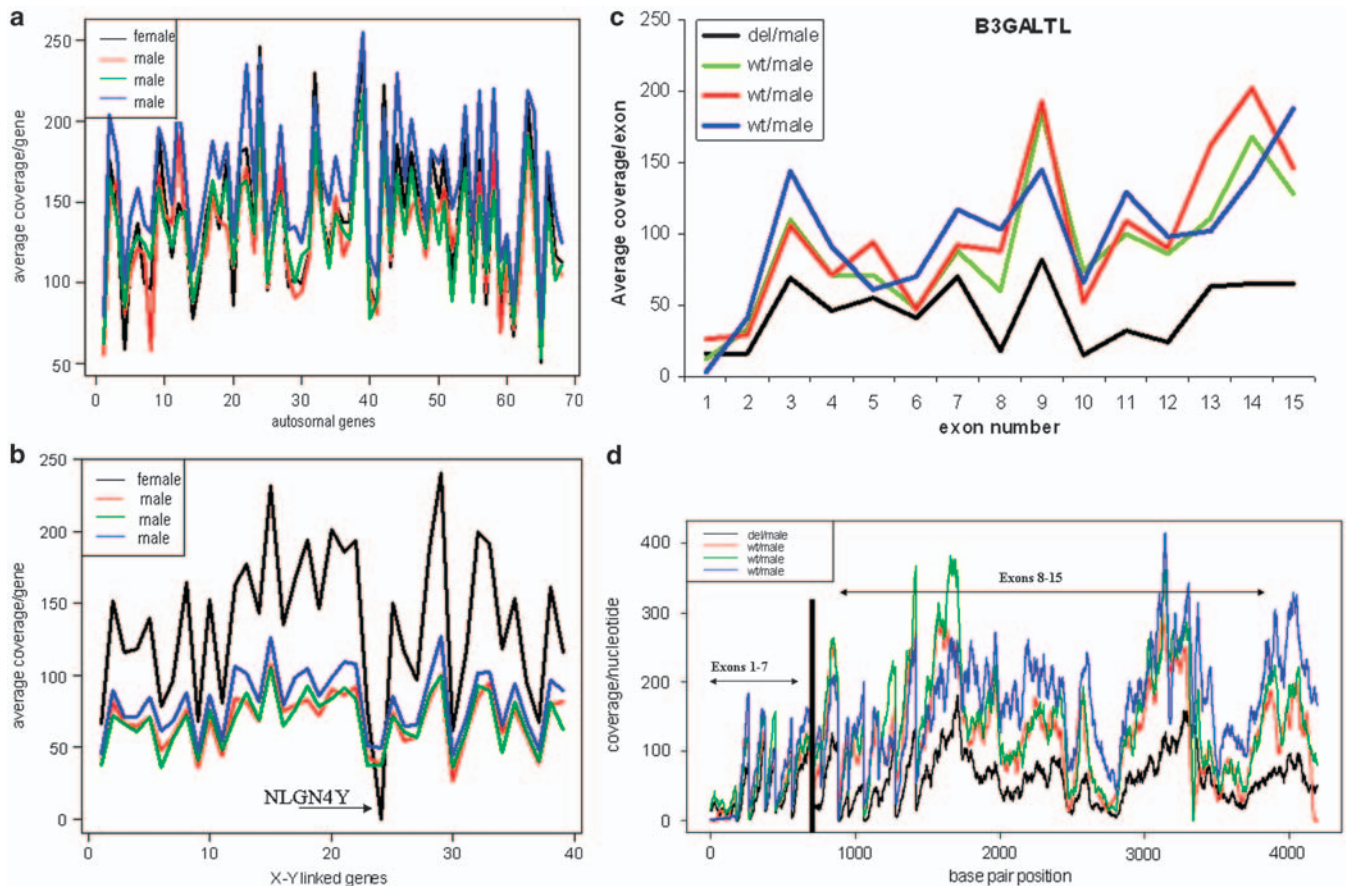
Regarding probe design (performed by Roche NimbleGen, Inc.), it should be noted that all array probes are from one strand (coding DNA strand) and thus DNA molecules from only the non-coding strand are captured. This has several consequences. First, the sequence obtained is from one strand only, whereas for diagnostic applications, quality assurance requires that sequences be obtained in forward and reverse orientation. Sequencing this one strand in both directions is partly fooling oneself. Second, we observed that the sequences obtained relative to the array probes extend in a 5' but not in a 3' direction. The most probable cause for the latter is steric hindrance during array hybridization, preventing non-hybridizing tails at the surface side of the array. When capture probes are attached with their 3' ends, this has consequences for probe design at the edges of the targeted regions; on the 5' side, coverage will be significantly better than on the 3' side. Both effects could be overcome simply by reversing the probe sequence of every other nucleotide on the array. Theoretically, this would also mean that the overall yield of enriched DNA would double, as both strands from the sample will be captured.

To save costs, we have reused the arrays up to three times by hybridizing different samples. The danger of this approach is of course contamination, if hybridized DNA from a previous experiment is not eluted completely. Indeed, in some experiments, we observed low-level contamination, for example, through heterozygous calls from X-chromosome sequences in male samples. It should be noted, however that cross-contamination can be easily controlled when samples containing differently tagged linkers are used in subsequent experiments.

Using the current design, low coverage was obtained mainly at the edges of the regions targeted, especially the 3' side (see above), that is, direct gene flanking or intronic regions. Although coverage varied widely, 78% of all regions targeted and present on the array were covered effectively by the sequence obtained. Note that there is a clear correlation between fragment size of the genomic DNA used and the coverage, the larger the fragment size used the lower the target coverage achieved, as more flanking DNA is captured. Especially for array-based capture, because of the steric hindrance described, this effect will be significant near the array-attached end of a probe-targeted region. Assuming that second-generation capture arrays will be more effective (ie, complete and with even coverage) and sequence power will improve further, it should soon be possible to sequence-tag, mix and simultaneously analyze different samples in one experiment, giving a significant cost reduction.

Recently in-solution capture was presented as an alternative to array-based capture.<sup>25</sup> Besides advantages of simplicity, a reduced workload and a potential for automation, when attempted, in-solution capture will not show the effect of steric hindrance we observed. However, capturing both strands would be complicated by the fact that capture probes will hybridize with each other. Initial experiences in our lab with in-solution capture were successful and for future projects we will change to this approach.





**Figure 2** Average coverage obtained for different genes in four different samples. (a) Shows average coverage of 69 autosomal genes from four different samples. (b) Shows average coverage of 39 genes located on X and one gene (*NLGN4Y*) located on the Y chromosome; a female sample exhibited an absence of hybridization in the captured array, with no coverage in the regions corresponding to the *NLGN4Y*. The female sample shows a higher average coverage per gene for all genes located on X-chromosome compared with male samples. (c) Lower average coverage of *B3GALT1* gene in a male patient sample with a known large deletion compared with three wild-type male samples. (d) Coverage per nucleotide/position for the whole *B3GALT1* gene: the patient sample shows lower coverage for the second half (exons 8–15) compared with wild type samples. del=deletion, wt=wild type.

Overall, we conclude that array-based sequence capture followed by NGS offers a versatile tool for successfully selecting sequences of interest from a total human genome. The approach will be especially helpful in speeding up the identification of the pathogenic mutation(s) in diseases in which the genomic region to be scanned is large. Our results indicate that the methodology can still be improved, in particular, with respect to probe design, obtaining a more even coverage of the targeted regions. On the basis of initial experiences and publications, we expect that array capture will be quickly replaced by in-solution capture. Ultimately, the cost of this approach is determined by the minimal coverage, which in turn determines the sensitivity required for the detection of potential sequence variants.

#### CONFLICT OF INTEREST

The authors declare no conflict of interest.

#### ACKNOWLEDGEMENTS

We thank the Leiden Genome Technology Center (LGTC), in particular Sophie Greve-Onderwater, Matthew Hestand and Rolf Vossen, for their expert technical assistance; Antoinette Gijbbers for sharing the SNP data; and Kamlesh Madan for critical reading of the paper. The research leading to these results has

received funding from the European Community's Seventh Framework Program (FP7/2007-2013) under Grant agreements 223026 (NMD-chip) and 223143 (the TechGene).

- 1 Sanger F, Nicklen S, Coulson AR: DNA sequencing with chain-terminating inhibitors. *Proc Natl Acad Sci USA* 1977; **74**: 5463–5467.
- 2 Chelly J, Khelifaoui M, Francis F, Chérif B, Bienvu T: Genetics and pathophysiology of mental retardation. *Eur J Hum Genet* 2006; **14**: 701–713.
- 3 Szigeti K, Lupski JR: Charcot-Marie-Tooth disease. *Eur J Hum Genet* 2009; **17**: 703–710.
- 4 Paul M, Zumhagen S, Stallmeyer B, Koopmann M, Spieker T, Schulze-Bahr E: Genes causing inherited forms of cardiomyopathies. A current compendium. *Herz* 2009; **34**: 98–109.
- 5 Hartong DT, Berson EL, Dryja TP: Retinitis pigmentosa. *Lancet* 2006; **368**: 1795–1809.
- 6 Muhle R, Trentacoste SV, Rapin I: The genetics of autism. *Pediatrics* 2004; **113**: 472–486.
- 7 Hilgert N, Smith RJ, Van Camp G: Forty-six genes causing nonsyndromic hearing impairment: which ones should be analyzed in DNA diagnostics? *Mutat Res* 2009; **681**: 189–196.
- 8 Freeze H: Genetic defects in the human glycome. *Nat Rev Genet* 2006; **7**: 537–551.
- 9 Bonetta L: Genome sequencing in the fast lane. *Nat Methods* 2006; **3**: 141–147.
- 10 von Bubnoff A: Next-generation sequencing: the race is on. *Cell* 2008; **132**: 721–723.
- 11 Schuster SC: Next-generation sequencing transforms today's biology. *Nat Methods* 2008; **5**: 16–18.

- 12 Shendure J, Ji H: Next-generation DNA sequencing. *Nat Biotechnol* 2008; **26**: 1135–1145.
- 13 Albert TJ, Molla MN, Muzny DM *et al*: Direct selection of human genomic loci by microarray hybridization. *Nat Methods* 2007; **4**: 903–905.
- 14 Okou DT, Steinberg KM, Middle C, Cutler DJ, Albert TJ, Zwick ME: Microarray-based genomic selection for high-throughput resequencing. *Nat Methods* 2007; **4**: 907–909.
- 15 Hodges E, Xuan Z, Balija V *et al*: Genome-wide *in situ* exon capture for selective resequencing. *Nat Genet* 2007; **39**: 1522–1527.
- 16 Roche NimbleGen: *NimbleGen services user's guides: sequence capture service*.[http://www.nimblegen.com/products/lit/SeqCap\\_UsersGuide\\_Service\\_v3p0.pdf](http://www.nimblegen.com/products/lit/SeqCap_UsersGuide_Service_v3p0.pdf).
- 17 Edwards MC, Gibbs RA: Multiplex PCR: advantages, development, and applications. *PCR Methods Appl* 1994; **3**: S65–S75.
- 18 Markoulatos P, Siafakas N, Moncany M: Multiplex polymerase chain reaction: a practical approach. *J Clin Lab Anal* 2002; **16**: 47–51.
- 19 Cutler DJ, Zwick ME, Carrasquillo MM *et al*: High-throughput variation detection and genotyping using microarrays. *Genome Res* 2001; **11**: 1913–1925.
- 20 Dahl F, Gullberg M, Stenberg J, Landegren U, Nilsson M: Multiplex amplification enabled by selective circularization of large sets of genomic DNA fragments. *Nucleic Acids Res* 2005; **33**: 71.
- 21 Dahl F, Stenberg J, Fredriksson S *et al*: Multigene amplification and massively parallel sequencing for cancer mutation discovery. *Proc Natl Acad Sci USA* 2007; **104**: 9387–9392.
- 22 Bashiardes S, Veile R, Helms C, Mardis ER, Bowcock AM, Lovett M: Direct genomic selection. *Nat Methods* 2005; **2**: 63–69.
- 23 Tewhey R, Warner JB, Nakano M *et al*: Microdroplet-based PCR enrichment for large-scale targeted sequencing. *Nat Biotechnol* 2009; **27**: 1025–1031.
- 24 Shaw-Smith C, Redon R, Rickman L *et al*: Microarray based comparative genomic hybridisation (array-CGH) detects submicroscopic chromosomal deletions and duplications in patients with learning disability/mental retardation and dysmorphic features. *J Med Genet* 2004; **41**: 241–248.
- 25 Gnirke A, Melnikov A, Maguire J *et al*: Solution hybrid selection with ultra-long oligonucleotides for massively parallel targeted sequencing. *Nat Biotechnol* 2009; **27**: 182–189.



This work is licensed under the Creative Commons Attribution-NonCommercial-No Derivative Works 3.0 Unported Licence. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/3.0/>

Supplementary information accompanies the paper on European Journal of Human Genetics website (<http://www.nature.com/ejhg>)