



HAL
open science

Assessment of Copy Number Variation using the Illumina Infinium 1M SNP-array: A comparison of methodological approaches in the Spanish Bladder Cancer / EPICURO Study.

Gaëlle Marenne, Benjamín Rodríguez Santiago, Montserrat García-Closas,
Luis Alberto Pérez Jurado, Nathaniel Rothman, Daniel Rico, Guillermo Pita,
David G Pisano, Manolis Kogevinas, Debra T Silverman, et al.

► To cite this version:

Gaëlle Marenne, Benjamín Rodríguez Santiago, Montserrat García-Closas, Luis Alberto Pérez Jurado, Nathaniel Rothman, et al.. Assessment of Copy Number Variation using the Illumina Infinium 1M SNP-array: A comparison of methodological approaches in the Spanish Bladder Cancer / EPICURO Study.. Human Mutation, 2011, 32 (2), pp.240. 10.1002/humu.21398 . hal-00610793

HAL Id: hal-00610793

<https://hal.science/hal-00610793>

Submitted on 25 Jul 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.



Assessment of Copy Number Variation using the Illumina Infinium 1M SNP-array: A comparison of methodological approaches in the Spanish Bladder Cancer / EPICURO Study.

| | |
|-------------------------------|--|
| Journal: | <i>Human Mutation</i> |
| Manuscript ID: | humu-2010-0239.R1 |
| Wiley - Manuscript type: | Methods |
| Date Submitted by the Author: | 30-Sep-2010 |
| Complete List of Authors: | <p>Marenne, Gaëlle; Spanish National Cancer Research Centre Rodríguez Santiago, Benjamín; Departament de Ciències Experimentals i de la Salut, Universitat Pompeu Fabra</p> <p>García-Closas, Montserrat; Division of Cancer Epidemiology and Genetics, National Cancer Institute, Department of Health and Human Services</p> <p>Pérez Jurado, Luis; Universitat Pompeu Fabra, Ciències Experimentals i de la Salut; Hospital Vall d'Hebron, Program in Molecular Medicine and Genetics</p> <p>Rothman, Nathaniel; Division of Cancer Epidemiology and Genetics, National Cancer Institute, Department of Health and Human Services</p> <p>Rico, Daniel; Spanish National Cancer Research Centre</p> <p>Pita, Guillermo; Spanish National Cancer Research Centre</p> <p>Pisano, David; Spanish National Cancer Research Centre</p> <p>Kogevinas, Manolis; Centre for Research in Environmental Epidemiology</p> <p>Silverman, Debra; Division of Cancer Epidemiology and Genetics, National Cancer Institute, Department of Health and Human Services</p> <p>Valencia, Alfonso; Spanish National Cancer Research Centre</p> <p>Real, Francisco; Spanish National Cancer Research Centre</p> <p>Chanock, Stephen; National Cancer Institute, Pediatric Oncology Branch</p> <p>Génin, Emmanuelle; Inserm UMR-S946, Univ. Paris Diderot, Institut Universitaire d'Hématologie,</p> <p>Malats, Núria; Spanish National Cancer Research Centre</p> |
| Key Words: | Copy Number Variation, Genome Wide Association Study, Specificity, Sensitivity, Reliability, Accuracy, CNVpartition, |

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



SCHOLARONE™
Manuscripts

For Peer Review

1
2 **Title:** Assessment of Copy Number Variation using the Illumina Infinium 1M SNP-array: A
3
4 comparison of methodological approaches in the Spanish Bladder Cancer / EPICURO
5
6 Study.
7

8
9
10 **Authors:** Gaëlle Marenne(1, 2), Benjamín Rodríguez-Santiago(3, 4), Montserrat García
11
12 Closas(5), Luis Pérez-Jurado(3, 4, 6, 7), Nathaniel Rothman(5), Daniel Rico(1),
13
14 Guillermo Pita(1), David G. Pisano(1), Manolis Kogevinas(8, 9), Debra T
15
16 Silverman(5), Alfonso Valencia(1), Francisco X Real(1), Stephen Chanock*(5),
17
18 Emmanuelle Génin*(2), Núria Malats*(1). (* co-senior authors)
19

20
21
22 **Affiliations of authors:** (1) Centro Nacional de Investigaciones Oncológicas (CNIO) Madrid,
23
24 Spain; (2) Inserm UMR-S946, Univ. Paris Diderot, Institut Universitaire d'Hématologie, Paris,
25
26 France; (3) Departament de Ciències Experimentals i de la Salut, Universitat Pompeu Fabra,
27
28 Barcelona, Spain; (4) CIBER de Enfermedades Raras, CIBERER, E-08003 Barcelona, Spain;
29
30 (5) Division of Cancer Epidemiology and Genetics, National Cancer Institute, Department of
31
32 Health and Human Services, Bethesda, MD, USA; (6) Programa de Medicina Molecular i
33
34 Genètica, Hospital Universitari Vall d'Hebron, E-08035 Barcelona, Spain; (7) Department of
35
36 Genome Sciences, University of Washington, Seattle, WA 98195, United States; (8) Institut
37
38 Municipal d'Investigació Mèdica (IMIM-Hospital del Mar), Barcelona, Spain; (9) Centre for
39
40 Research in Environmental Epidemiology (CREAL), Barcelona, Spain.
41

42
43 **Corresponding author:** Núria Malats (nmalats@cnio.es).
44

45 **Running Title:** Accuracy study on CNV assessment
46

47 **Conflict of interest statement:** We declare we have no conflict of interest.
48
49
50

51 Deleted: 29/09/2010
52

53 | Accuracy Ms (30/09/2010)
54

55 1
56
57
58
59
60

ABSTRACT

High-throughput SNP-array technologies allow to investigate CNVs in genome-wide scans and specific calling algorithms have been developed to determine CNV location and copy number.

We report the results of a reliability analysis comparing data from 96 pairs of samples processed with CNVpartition, PennCNV and QuantiSNP for Infinium Illumina Human 1Million probe chip data. We also performed a validity assessment with **multiplex ligation-dependent probe amplification** (MLPA) as a reference standard.

The number of CNVs per individual varied according to the calling algorithm. Higher numbers of CNVs were detected in saliva than in blood DNA samples regardless of the algorithm used. All algorithms presented low agreement with mean Kappa Index (KI) <66. PennCNV was the most reliable algorithm (KI_w=98.96) when assessing the number of copies. The agreement observed in detecting CNV was higher in blood than in saliva samples. When comparing to MLPA, all algorithms identified poorly known copy aberrations (sensitivity=0.19-0.28). In contrast, specificity was very high (0.97-0.99). Once a CNV was detected, the number of copies was **truly** assessed (**sensitivity>0.62**).

Our results indicate that the current calling algorithms should be improved for high performance CNV analysis in genome-wide scans. Further refinement is required to assess CNVs as risk factors in complex diseases.

Key Words: Copy Number Variation, Genome Wide Association Study, Specificity, Sensitivity, Reliability, Accuracy, CNVpartition, PennCNV, QuantiSNP

Deleted: 29/09/2010

| Accuracy Ms (30/09/2010)

2

INTRODUCTION

Structural variations of the human genome emerge as novel major contributors to genetic diversity and disease susceptibility. Copy number variation (CNV) refers to deletions or duplications larger than 1kb (Feuk et al., 2006). It was estimated that 12% of the genome could be affected by such variants in comparison to 1-2% covered by single nucleotide polymorphisms (SNPs) (Redon et al., 2006); although a recent study provided a lower figure: 3.7% (Conrad et al., 2010). These large variations can overlap with genes and there is substantial evidence for correlation between CNVs and gene expression levels (Stranger et al., 2007). CNVs are also known to be involved both in mendelian disorders, such as Williams–Beuren Syndrome (deletion at chromosome region 7q11.23) or Charcot–Marie Tooth neuropathy Type 1A (duplications at chromosome region 17p11.2), and complex traits such as HIV infection and asthma, among others (Ionita-Laza et al., 2009).

Recently, efforts have been made to provide resources supporting studies of structural variation in human diseases such as the Database of Genomic Variation which annotates genomic coordinates along with estimated frequencies of the CNVs (Conrad et al., 2010; Iafrate et al., 2004; Redon et al., 2006). However, the cost and the complexity of CNV assessment have restricted CNV studies to a list of carefully selected candidate genes. The possibility to study CNVs at a genome-wide scale is now possible using high-throughput SNP-array technologies. The new-generation SNP-arrays, such as the Infinium Illumina Human 1Million probe chip and the Affymetrix 6.0 platform, allow a cost-effective detection of CNVs by interpreting allele intensities for each marker. These platforms also include monomorphic probes in regions of common CNVs that presented technical problems for SNP array design due to a lack of polymorphic probes or because of disruption from Mendelian inheritance and Hardy-Weinberg equilibrium. The Illumina 1 Million SNP-array works with Beadstudio software that provides the variables used to perform the CNV calling. Different

Deleted: 29/09/2010

| Accuracy Ms (30/09/2010)

3

1
2 algorithms can then be employed to locate CNVs by finding breakpoints and assessing the
3
4 number of copies present per individual. The most frequently-used algorithms for Illumina
5
6 data are CNVpartition – an Illumina developed plug-in –, PennCNV (Wang et al., 2007) and
7
8 QuantiSNP (Colella et al., 2007).

9
10 Several studies have successfully assessed the role of CNVs in complex diseases such as
11
12 asthma, autism, schizophrenia or cancer by applying high throughput analysis at genome-
13
14 wide level (Bae et al., 2008; Bassett et al., 2008; Blauw et al., 2008; Cronin et al., 2008;
15
16 Diskin et al., 2009; Friedman et al., 2006; Glessner et al., 2009; Greenway et al., 2009;
17
18 InternationalSchizophreniaConsortium, 2008; Ionita-Laza et al., 2008; Kathiresan et al., 2009;
19
20 Liu et al., 2009; Marshall et al., 2008; Matarin et al., 2008; Need et al., 2009; Sha et al., 2009;
21
22 Simon-Sanchez et al., 2008; Stefansson et al., 2008; Walsh et al., 2008; Weiss et al., 2008; Xu
23
24 et al., 2008; Yang et al., 2008). A review of these studies indicates that they have used a wide
25
26 range of methodologies, thus raising the issue of comparability of discovery rates. The rapid
27
28 development of technologies in this field has not been accompanied by a careful evaluation of
29
30 the software tools to assess disease risk association. In contrast to the nearly 100%
31
32 concordance observed for bi-allelic genotypes, a recent study reported very low agreement
33
34 estimates when the performance of different algorithms assessing CNV was compared using
35
36 HapMap data (Winchester et al., 2009).

37
38 Here, we report the results from reliability and validity analyses comparing three CNV calling
39
40 algorithms for Illumina 1M probe-array data (CNVpartition, PennCNV and QuantiSNP) using
41
42 multiplex ligation-dependent probe amplification (MLPA) as the gold-standard analysis. The
43
44 study was conducted on 96 duplicate samples from the Spanish Bladder Cancer Study. We
45
46 also assessed whether the source of DNA (blood or saliva) and the number and type of SNPs
47
48 considered in the CNV definition influenced the performance of the SNP calling algorithms.

51
52 Deleted: 29/09/2010

53
54 | Accuracy Ms (30/09/2010)

55
56
57
58
59
60 4

MATERIALS AND METHODS

Samples and genotyping data

Study subjects were recruited to the Spanish Bladder Cancer Study (SBCS)/EPICURO, conducted between 1998-2000. Individuals were from 5 different regions in Spain (Barcelona, Vallès/Bages, Alicante, Tenerife and Asturias). Leukocyte and saliva DNA were obtained as described elsewhere (Garcia-Closas et al., 2005). Genotyping was performed at the Core Genotyping Facility, National Cancer Institute, USA, using the Infinium Illumina Human 1M probe BeadChip containing 1,072,820 markers, among which 206,665 are in reported CNVs regions. For quality control reasons, 141 individuals were genotyped two to four times providing genetic data for 178 pairs out of 299 assays (Supp. Table S1).

Log R Ratio (LRR) and the B Allele Frequency (BAF) were exported from the normalized Illumina data through the Beadstudio software to perform CNV calling. LRR is the ratio between the observed and the expected probe intensity. The expected intensity is an interpolation of the mean intensities of the surrounding genotype clusters. BAF represents the proportion of B alleles in the genotype. A region without evidence of CNV should show a LRR around zero and three clusters of BAF of 0, 0.5 and 1 corresponding to the three genotypes AA, AB and BB, respectively (Supp. Figure S1). Individuals not fitting at least one of the CNV specific quality control metric recommended by PennCNV (Wang et al., 2007) were excluded from the analysis: $LRR - Standard\ Deviation > 0.28$, $0.45 > BAF - median > 0.55$, $BAF - drift > 0.002$, and $-0.04 > Wave\ Factor > 0.04$. After applying the abovementioned criteria, 92 individuals (90 duplicates and 2 triplicates) were suitable for this study, thus providing 96 pairs for comparison (90 from duplicate individuals and 6 from triplicate individuals) and 186 assays (90 individuals * 2 samples and 2 individuals * 3 samples). Among the duplicates there were 63 and 33 pairs from blood and saliva samples, respectively (Supp. Table S1).

Deleted: 29/09/2010

CNV calling

Three algorithms available for Illumina data were applied: CNVpartition, PennCNV (Wang et al., 2007) and QuantiSNP (Colella et al., 2007). CNVpartition was developed by Illumina and is available as a plug-in in the Beadstudio software. It is based on the assumption that the majority of CNV vary between 0 and 4 copies (i.e. AAAA, AAAB, AABB ...) thus yielding five options (homozygous deletion, heterozygous deletion, dizygous (normal state), trizygous (one extra copy), and tetrazygous (two extra copies). CNVpartition model LRR and BAF as simple bivariate Gaussian distributions for each of the fourteen possible copy genotypes. A preliminary copy number estimate is computed for each assayed locus by comparing its observed LRR and BAF to values predicted from each of the fourteen genotypes. Specifically, the likelihood of observing a given LRR and BAF under each of the fourteen models is computed and the number of copies is estimated by maximizing the likelihood. Once each probe is assigned a number of copies, breakpoints are determined by a partitioning method identifying regions where the estimated number of copies of the probes inside and outside the region is different. A confidence value is also provided to allow the filtering of the CNV and limit the number of false positive callings.

PennCNV and QuantiSNP are algorithms developed by academic teams and freely available (Colella et al., 2007; Wang et al., 2007). They are both based on a Hidden Markov Model (HMM) in which the number of gene copies is the hidden state and the LRR and the BAF are the two observed states that are considered independent of each other given the number of copies. A first-order HMM is considered where the number of copies at one probe depends on the number of copies at the previous probe. However, the two algorithms differ in their transition and emission probabilities. While transition probabilities depend on the distance between adjacent probes for both approaches, the probabilities for PennCNV are also state-specific, accounting for the fact that some state transition events (e.g., from normal state to

Deleted: 29/09/2010

| Accuracy Ms (30/09/2010)

6

1
2 heterozygous deletion) are more likely than others (e.g., from heterozygous deletion to
3 trizygous). Regarding the BAF emission probabilities, PennCNV uses a more sophisticated
4 model than QuantiSNP. Both algorithms provide a confidence value to filter CNVs. For
5
6 QuantiSNP, the confidence value is the Log Bayes Factor (LBF). All algorithms were used
7
8 with their default options and CNV calls from QuantiSNP with a LBF lower than 10 were
9
10 filtered out as recommended whereas no filter was applied on CNVpartition and PennCNV
11
12 calls.
13
14

15
16 Each of the 1,029,591 probes of the Illumina 1M array corresponding to the autosomal
17
18 chromosomes was assigned with an estimated number of copies if were included in a CNV
19
20 and with two copies otherwise. This procedure was applied to each of the 186 experiments
21
22 performed in this study and for each of the algorithms.
23

24 Reliability analysis

25
26 The calling agreement between duplicates was evaluated for each of the algorithms to
27
28 determine presence of CNV and number of copies. First, we assessed the agreement in
29
30 detecting the presence of an aberration by estimating the kappa index (KI) between
31
32 duplicates. KI compared the observed agreement against the agreement expected by chance in
33
34 all the probes (Cohen, 1960). For probes in which the algorithm was concordant in detecting
35
36 an aberration, we computed the agreement in assessing the number of copies by estimating
37
38 the weighted Kappa Index (KI_w). This was done by applying quadratic weights that decreased
39
40 while increasing differences in copy numbers (Supp. Figure S2). A total of 96 KI and KI_w
41
42 values were obtained for each algorithm. Summary statistics (mean, median, standard-
43
44 deviation, and quartiles) were computed and differences between algorithms were tested using
45
46 paired t-tests.

47
48 To further limit the number of false positive CNV callings from SNP-array platforms, Itsara
49
50 et al proposed to filter the called CNVs according to the type of aberration and the number of
51

Deleted: 29/09/2010

52
53
54 | Accuracy Ms (30/09/2010)

7

1 genotyped SNPs included in the CNV (Itsara et al., 2009). The LRR intensities were
2 transformed into standard normal measurements (Z-scores) and the B-deviation value for each
3 probe was estimated. Putative CNVs were classified into two categories (small and large)
4 according to a cut-off of 100 probes and 1 Mb length. Large CNVs were manually curated.
5 Small CNVs were subject to automated filtering. Homozygous deletions were required to
6 comply with: 1) ≥ 3 probes, median LRR Z-score ≤ -4 , and mean B-deviation ≥ 0.1 or 2) ≥ 3
7 probes and median LRR Z-score ≤ -8 . Heterozygous deletions were required to span ≥ 10
8 probes, have LRR Z-score ≤ -1.5 , and less than 10% of probes called as heterozygous. To
9 define duplications, the requirements were: ≥ 10 probes, LRR Z-score ≥ 1.5 , and B-deviation
10 among heterozygote probes ≥ 0.075 . The reliability of applying the Itsara's filter was
11 assessed, too.

12 We analyzed the calling agreement of paired samples depending on the DNA source by
13 stratifying the data according to whether the DNA was from blood (N=63) or saliva (N=33).
14 In addition, we assessed whether the number of SNPs included in each CNV influenced the
15 agreement rate by comparing the CNV calling performance between replicates by filtering for
16 the number of SNPs in the CNVs. The reliability results were plotted for the three algorithms
17 and the number of CNVs called according to the number of SNPs.

18 Select commercial SNP genotyping platforms contain monomorphic probes in regions of
19 known common CNVs to facilitate analysis, particularly when prior analyses in HapMap
20 indicated a substantial problem of fitness with Hardy Weinberg proportions. The overall
21 percentage of monomorphic probes in the 1M Illumina Infinium platform in autosomal
22 chromosomes is 1.4% (14,716/1,029,591). To test the impact of the type of probe
23 (monomorphic or polymorphic) on the reliability of the calling, we compared for these two
24 types of probes the ratio of concordant vs. discordant probes included in CNVs. We excluded
25 the regions with a concordant result for the absence of CNV because the density of the

Deleted: 29/09/2010

1
2 monomorphic probes in those regions was lower according to the design of the SNP-array,
3
4 hence not being comparable.
5

6 **Validity Study**

7
8 Multiplex ligation-dependent probe amplification (MLPA) assay is a standard laboratory
9
10 approach to assess differences in the number of alleles copies at a particular locus. It is based
11
12 on hybridization, specific probe ligation, amplification and capillary migration, and it was
13
14 used as the gold-standard method to assess the number of copies of a given sequence. **Regions**
15
16 **were selected for validation with MLPA if at least one algorithm detected a minimum of 8**
17
18 **individuals carrying a CNV to avoid performing experiments in regions where no CNV exist.**

19
20 Commercial probe mixes (kits P070 and P036 covering the selected regions (MRC-Holland
21
22 Amsterdam, The Netherlands) and custom designed probes (Supp. Table S2) were used.
23
24 MLPA reactions were carried out as described previously (Schouten et al., 2002) with slight
25
26 modifications when custom probes were used (Rodriguez-Santiago et al., 2009). The relative
27
28 peak height (RPH) method recommended by MRC-Holland was used to determine the copy
29
30 number status. Theoretically, heterozygous deletions and duplications showed a relative peak
31
32 height of approximately 0.5 and 1.5, respectively. Only blood samples were considered for
33
34 this analysis.

35
36 Leukocyte DNA from 56 individuals was analyzed twice by MLPA, providing a concordance
37
38 rate of 97.25%. Among the discordant assays, 10 showing a “non-calling” rate greater than
39
40 70% were re-analyzed. Since the results of four of them slightly improved after the 2nd MLPA
41
42 run they were included in the validity study and data were updated.

43
44 To assess the validity of each algorithm, sensitivity, specificity, and positive and negative
45
46 predictive values were computed by comparing CNV callings with MLPA data. Sensitivity
47
48 (SE) indicates the proportion of CNV identified by the algorithm over the total number of
49
50 existing CNV according to MLPA. Specificity (SP) is the proportion of the non-CNV by an

51
52 Deleted: 29/09/2010

53
54 | Accuracy Ms (30/09/2010)

55
56
57
58
59
60 9

1
2 algorithm over the true non-CNV number. Positive (PPV) and negative predictive values
3
4 (NPV) indicate the proportion of the true CNV and the true non-CNV over all CNV and non-
5
6 CNV regions each algorithm assigns, respectively. These estimates are given as proportions
7
8 with a 95%CI for the overall aberration assessment and for each type of CNV. The validity
9
10 analysis considered those probes and individuals that provided agreement in detecting CN
11
12 event according to each algorithm.

13
14 Statistical analyses were performed in R version 2.9.0 (<http://www.r-project.org>) with the
15
16 *epiR* package (Mark Stevenson, <http://epicentre.massey.ac.nz>). Significance was declared
17
18 when the p-value was smaller than 0.05.

21 22 RESULTS

23
24 The number of CNVs detected per individual varied substantially according to the calling
25
26 algorithm (Table 1). CNVpartition identified an average of 28.0 CNVs per individual whereas
27
28 the two algorithms based on the HMM, PennCNV and QuantiSNP, identified a median CNV
29
30 number of 58.5 and 56.0, respectively. The number of CNVs per individual detected in saliva
31
32 DNA was higher than in leukocyte DNA, regardless of the algorithm used (Table 1).

33 Reliability analysis

34
35 The SNP calling provided by the genotyping platform showed a very high agreement with a
36
37 mean Kappa Index (KI) of 99.99 (95%CI, 99.94 – 100) (Figure 1a). The distribution of this
38
39 KI was similar for experiments using blood or saliva DNA. Regarding CNV assessment in
40
41 duplicate samples, PennCNV, QuantiSNP, and CNVpartition presented a lower agreement
42
43 with mean KI values of 65.10, 63.09, and 57.24, respectively. The KI distribution based on
44
45 CNVpartition callings significantly differed from that based on PennCNV and QuantiSNP
46
47 callings ($p=2.68 \times 10^{-10}$ and $p=7.28 \times 10^{-5}$, respectively) (Figure 1b). Once a region of CNV was
48
49 detected, the algorithms also showed differences in the KI distribution when assessing the
50
51

Deleted: 29/09/2010

52
53 | Accuracy Ms (30/09/2010)

10

1
2 number of copies (Figure 1c). PennCNV appeared to be the most reliable algorithm with an
3 average KI_w (weighted KI) = 98.96 for the 96 pairs of replicates, and regardless the type of
4 CNV (gain or loss). However, QuantiSNP and CNVpartition performed differently and poorly
5 (Supp. Figure S3). This figure was significantly higher than those of CNVpartition
6 ($KI_w=94.55$, $p=5.18 \times 10^{-5}$) and QuantiSNP ($KI_w=92.88$, $p=7.43 \times 10^{-8}$). Applying the Itsara
7 filtering method, we did not observe an improvement of the agreement neither at the CNV
8 detection level nor at the level of copy number (Supp. Figure S4).

9
10 Regardless of the algorithm applied, the agreement observed in detecting CNV was always
11 higher in blood than in saliva samples (Figure 2), although the difference of the mean KI was
12 only significant for CNVpartition and PennCNV callings ($p=3.93 \times 10^{-7}$ and $p=8.16 \times 10^{-5}$,
13 respectively). The distribution of KI_w when assessing the number of copies, according to the
14 DNA source, was similar for all algorithms (data not shown).

15
16 The number of probes selected by each algorithm to identify CNVs varied widely: 1,742 for
17 CNVpartition, 2,361 for PennCNV, and 4,591 for QuantiSNP (Table 2). The percentage of
18 probes showing agreement for the presence of a CNV was significantly different for the three
19 algorithms: 37.7%, 50.7%, and 55.5% for CNVpartition, PennCNV, and QuantiSNP,
20 respectively, ($p=2.43 \times 10^{-35}$). The ratio between discordant/concordant probes was higher for
21 monomorphic than polymorphic probes: 2.17 vs. 1.61 for CNVpartition ($p=0.09$), 1.78 vs.
22 0.94 for PennCNV ($p=4.34 \times 10^{-4}$), and 1.51 vs. 0.72 for QuantiSNP ($p=1.31 \times 10^{-17}$).

23
24 The correlation between the calling agreement and the number of probes or the length of a
25 given CNV region is shown in Figure 3. A direct relationship between agreement and the
26 number of probes included in the CNVs was observed suggesting that reliability is greater for
27 CNVs containing more probes. This effect was observed for all algorithms but it was higher
28 for PennCNV. Our results also suggested that filtering CNVs by QuantiSNP for length, by

Deleted: 29/09/2010

PennCNV for length lower than 500 kb or by CNVpartition for length lower than 1Mb did not increase the reliability.

Validity analysis

Sensitivity (SE) and Specificity (SP) estimates for the presence and the type of CNV were estimated according to each algorithm (Figure 4). When considering the presence of CNVs (first line in Figure 4), we found that none of the algorithms used identified known CNV well ($0.19 \leq SE \leq 0.28$). In contrast, SP was very high ($0.97 \leq SP \leq 0.99$), indicating that algorithms rarely assigned a CNV in a region where it did not exist. QuantiSNP showed the best SE (0.28) with a SP of 0.97, similar to that of the other two algorithms. Nonetheless, the false positive (FP) calling rate for this algorithm (FP=34) was 2.8-fold higher compared to CNVpartition (FP=12), the latter showing the highest SP (0.99) and the lowest SE (0.19) (Supp. Table S3). PennCNV presented intermediate values of SE (0.23) and SP (0.98), yielding 22 false positive CNVs out of 1319 true “non-CNV”.

We also aimed at assessing whether copy number was well estimated when a CNV was identified. Since MLPA is prone to misclassify copy number states >3 , we classified CNVs in the following categories, instead: “duplications”, “homozygous deletions”, and “heterozygous deletions”; for specific purposes, we used the combined category “deletions” including both homozygous and heterozygous deletions. Once a CNV was identified, gene copy number was usually well estimated, the overall SEs for all types of CNVs being >0.62 . As expected, SP estimates remained very high ($SP > 0.87$). PennCNV and CNVpartition performed better than QuantiSNP, the latter showing the highest rates of FP and FN callings. QuantiSNP performed especially poorly when calling homozygous deletions ($SE=0.68$ and $SP=0.92$). When the Itsara filter was used, SE estimates were significantly decreased to values of 0.05, 0.07, and 0.08 for CNVpartition, PennCNV, and QuantiSNP, respectively; SP increased up to 0.997 for all algorithms (Supp. Table S3).

Deleted: 29/09/2010

DISCUSSION

In the past few years, the genomics community has begun to annotate a CNV genome wide map that provides better information on the contribution of structural genomic variation to genetic diversity in humans. SNP-array based-methods have allowed their association with disease susceptibility. However, the tools to carry out this task are still relatively rudimentary and the approach applied until now has mainly been based on reporting and validating individual CNVs located in candidate genes rather than assessing disease risk using genome wide analyses. This is primarily because of issues related to the accuracy of the available CNV calling algorithms. Which is, then, the most suitable method to identify CNVs for association studies using data from SNP-arrays?

The early comparisons have focused on evaluations using simulations or data from a few HapMap or CEPH samples (Kidd et al., 2008; Korbelt et al., 2007; Redon et al., 2006; Winchester et al., 2009). Here we provide, for the first time, a direct comparison of the accuracy (reliability and validity) of 3 CNV calling algorithms (PennCNV, QuantiSNP, and CNVpartition) using MLPA as a gold standard and therefore eliminating some of the concerns for the validity when using simulation or resequencing data. We also investigated a more stable platform, Illumina Infinium 1M array that may not suffer from the same clustering biases as the former ones.

The algorithms used displayed wide variation in the number of CNV events. Overall, we conclude that the reproducibility of the algorithms is less than optimal. Our results indicate that PennCNV and QuantiSNP are more reliable in detecting CNVs than CNVpartition. Yet, the agreement achieved with these algorithms was much lower (mean KI ranged 57-65) than that observed for SNP calling (KI=99.99). Winchester et al, reported a moderate overlap between PennCNV and QuantiSNP, ranging from 58-78% for the NA15510 CEPH sample

Deleted: 29/09/2010

1
2 (Winchester et al., 2009). One explanation for the unsatisfactory concordance in experimental
3 replicates for CNV detection and breakpoint identification relates to the different signal to
4 noise tolerance for SNP genotyping and CNV assessment. While the background signal of
5 SNP-arrays does not significantly affect SNP genotyping, it may affect CNV assessment due
6 to the need of different normalization approaches for the latter (Curtis et al., 2009; Winchester
7 et al., 2009).

8
9
10 Importantly, the three tools used performed poorly regarding their sensitivity to detect CNVs
11 when using MLPA experimental results as the gold standard, the percentage of missed CNV
12 ranging from 72-81%. Therefore, improved sensitivity of algorithms is a must in order to use
13 genome wide chip data for CNV detection and disease association studies. When the analysis
14 was restricted to concordant CNVs according to the applied algorithms, these estimated
15 adequately gene copy number. This result supports the notion of performing a two-stage
16 calling to increase accuracy. That is, to assess first the identification of CNVs and second, to
17 characterize those already detected.

18
19
20 Another important finding of our work relates to the source of DNA. Many studies have
21 shown that buccal cell and blood DNA provide similar calling rates for SNP. By contrast, we
22 found that leukocyte DNA is more reliable for CNV detection and that buccal cell DNA
23 yields a higher CNV calling rate. These findings are compatible with the idea that the
24 abundance of bacterial DNA in buccal samples can interfere with the performance of
25 genotyping bi-alleles as well, notably demonstrated by the higher discordance rates and lower
26 completion rates. Furthermore, while tissue-related differences in genome architecture leading
27 to variation in the number of CNVs may be real, other technical explanations such as DNA
28 quality should also be considered. In the Spanish Bladder Cancer/EPICURO Study, saliva
29 was obtained after a buccal rinse with Listerine® as a fixative. Saliva was then frozen until
30 DNA extraction. This simple and costless procedure yielded substantial amounts of DNA and
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

Deleted: 29/09/2010

1
2 allowed accurate SNP genotyping using TaqMan assays as well as Illumina technology. For
3
4 the latter, the calling agreement for leukocyte and buccal DNA was 99.99%. In the absence of
5
6 other studies providing similar information, caution is needed when analyzing buccal cell
7
8 DNA and new methodological studies specifically addressing these issues are needed.

9
10 Select commercial SNP-array platforms have included monomorphic probes to improve
11
12 coverage of CNV analyses. We have analyzed whether monomorphic and polymorphic
13
14 probes performed differently in assessing CNV. Surprisingly, we observed that, regardless of
15
16 the algorithm used, CNVs showing discordance between duplicates contained a higher
17
18 proportion of monomorphic probes than CNVs that were concordant. The difference was
19
20 greater for QuantiSNP. Hence, our findings indicate that polymorphic probes deliver more
21
22 robust information than monomorphic probes, at least using the current CNV calling tools.

23
24 Alternatively, it is possible that monomorphic probes may concentrate in a small number of
25
26 large CNVs being difficult to call since they are not homogeneously distributed across the
27
28 genome and are placed in those regions suspected of harbouring CN changes (Iafate et al.,
29
30 2004; Redon et al., 2006). Nevertheless, there is no evidence that CNVs in these regions are
31
32 larger than those elsewhere.

33
34 Despite the limitations described above, SNP-arrays offer important advantages over other
35
36 techniques to assess CNV at a genome wide level, including the possibility of analyzing a
37
38 large number of samples because of their relatively low cost and the small amount of DNA
39
40 required. CNV detection largely depends on the coverage of the platform. The low reliability
41
42 that we have observed may be partially due to the fact that the localization of the CNV
43
44 breakpoints depends on the position of the markers. While the Illumina 1M platform is one of
45
46 the densest arrays offering a genome wide coverage, the average distance between two probes
47
48 is around 3kb, larger than the smallest CNVs which are defined as having 1kb length. We
49
50 have found that the average distance between surrounding probes was greater for discordant

Deleted: 29/09/2010

1
2 than for concordant CN events. This effect was stronger for PennCNV and QuantiSNP than
3
4 for CNVpartition (results not shown). Small CNVs containing a small number of probes were
5
6 less reliable than large CNVs that are generally called based on more probes. Furthermore,
7
8 because the algorithms discard CNVs containing <3 probes, there was also an inherited
9
10 disadvantage to small CNVs as compared to larger ones. By applying the filter proposed by
11
12 Itsara *et al* (Itsara et al., 2009), agreement did not improve while sensitivity decreased
13
14 dramatically.

15
16 The relatively poor agreement between algorithms increases the heterogeneity in CNV
17
18 detection, raising the chance of false positive results in association studies. Furthermore,
19
20 current algorithms lack sensitivity for CNV identification, mainly when they are small. To
21
22 partially overcome this limitation, some authors have proposed to use the normalized intensity
23
24 obtained from the SNP-arrays, without performing the calling, and compare its distribution at
25
26 the individual probe level between cases and controls (Ionita-Laza et al., 2009; McCarroll and
27
28 Altshuler, 2007). Although this strategy has not been formally evaluated and power is
29
30 probably limited because of lack of biological meaning, it constitutes an alternative
31
32 exploratory approach to assess association of CNVs and phenotypes. Others have suggested
33
34 performing the calling and the association test simultaneously to take into account the
35
36 uncertainty of the calling in the test (Barnes et al., 2008; Gonzalez et al., 2009). However,
37
38 these methods require a priori definition of CNVs.

39
40 We used MLPA as the gold standard technique to estimate sensitivity and specificity of the
41
42 algorithms used. MLPA is reproducible, allows the detection of small differences in gene
43
44 copy number, requires low amounts of DNA, can be applied for mid-throughput studies, and
45
46 has a low cost. Among its limitations are the fact that it only detects CNVs in
47
48 targeted/selected genes and the results are bound to be affected by sequence polymorphisms
49
50 and by the occurrence of gene copy number changes in mosaicism. Despite careful probe

Deleted: 29/09/2010

1
2 design, we cannot rule out that an incomplete overlapping between probes and CNVs may
3
4 contribute to the low sensitivity for CNV detection found.

5
6 The algorithms used here are those that model both LRR and BAF to assess CNV, a practice
7
8 that allows the correction for bias effects and minimizes noise in the intensity measures (Yau
9
10 and Holmes, 2008). In addition, these algorithms are widely applied for CNV assessment
11
12 using Illumina derived data. Other CNV calling softwares are also available, such as Circular
13
14 Binary Segmentation (Olshen et al., 2004), GADA originally developed for array-CGH data
15
16 and adapted for SNP-array (Pique-Regi et al., 2008), DchipSNP (Lin et al., 2004), Tri Typer
17
18 (Franke et al., 2008) and SCIMM (Cooper et al., 2008). However, they do not jointly
19
20 incorporate both LRR and BAF information, their strengths and weaknesses have been
21
22 reviewed elsewhere (Winchester et al., 2009). Nevertheless, none of them has proven to be
23
24 superior to the ones used here. Winchester et al (Winchester et al., 2009) reported that
25
26 QuantiSNP yielded a higher number of events when measuring CNV in the NA15510 CEPH
27
28 sample in our study, QuantiSNP and PennCNV provided a similar mean number of CN
29
30 changes that was higher than that provided by CNVpartition. Recently, Dellinger et al
31
32 reported a comparison of 7 algorithms, including QuantiSNP, CNVpartition and PennCNV on
33
34 simulation studies on the basis of genotyped data by Affymetrix 6.0. The authors compared
35
36 sensitivity and specificity of the algorithms with CNV described in external databases (DGV,
37
38 HapMap Asian and HapMap confirmed) and concluded that QuantiSNP performed better than
39
40 the other algorithms (Dellinger et al., 2010).

41
42 Nevertheless, the current CNV calling algorithms do not yet provide stable, high quality calls
43
44 comparable to those in common usage for SNP calling algorithms. In particular, the
45
46 sensitivity is extremely low. Small/common CNVs may be less detectable because the
47
48 cumulative likelihood of CNV versus normal copy for a limited number of markers suffers
49
50 from a low signal-to-noise ratio. In order to improve this sensitivity in regions of known

Deleted: 29/09/2010

1
2 CNVs, some authors have proposed to look at some specific markers located within these
3
4 regions and use reported deletion and duplication frequencies as prior probabilities in the
5
6 calling. Such models are implemented in two widely used approaches, namely Canary (Korn
7
8 et al., 2008) and PennCNV-validation packages in which they have been shown to
9
10 substantially increase the sensitivity of calling CNV in these known regions. Efforts are also
11
12 made to improve technologies such as CGH-arrays and (Park et al., 2010) and next generation
13
14 sequencing. Hopefully, these will improve the detection of rare or novel CNVs in the near
15
16 future.

17
18 In conclusion, there is a need for better assays and tools to identify CNVs at the genome wide
19
20 level and test for their association with disease in large samples of cases and controls. The
21
22 main current limitations are the low reliability and sensitivity. Sensitivity showed differences
23
24 according to the algorithm applied and the type of change. The use of leukocyte DNA,
25
26 polymorphic probes, and a high number of probes per CNV should contribute to increase
27
28 reliability and PennCNV algorithm yield higher concordance rates.

29
30 The annotation of large CNVs across the genome has opened a new scenario to explore
31
32 genetic variation and its association with complex diseases and traits. While a few studies
33
34 support a major contribution of CNV to disease, there is an urgent need to develop and refine
35
36 better techniques and algorithms to assess CNVs at a genome wide level as disease-
37
38 predisposing variants.

41 ACKNOWLEDGEMENTS

42
43 We thank Juan Cruz Cigudosa, Ramón Díaz-Uriarte, Gonzalo Gómez, Kevin Jacobs, Kristel
44
45 Van Steen, and Marc Zindel for scientific sound comments and for technical support.

Deleted: 29/09/2010

We also acknowledge the support provided by Adonina Tardón, Alfredo Carrato, Consol Serra, Reina García-Closas, Josep Lloreta, Montserrat Torà, Gemma Castaño, María Salas, and Francisco Fernández, physicians, field workers, and lab technicians during the study.

This work was partially supported by the Fondo de Investigación Sanitaria, Spain (G03/174, PI061614, FI09/00205), Asociación Española Contra el Cáncer (AECC), Fundació Marató de TV3, Red Temática de Investigación Cooperativa en Cáncer (RTICC), Spain; by the Intramural Research Program of the Division of Cancer Epidemiology and Genetics, National Cancer Institute, USA; and by Egide-PHRC Picasso travel grant.

REFERENCE

Bae JS, Cheong HS, Kim JO, Lee SO, Kim EM, Lee HW, Kim S, Kim JW, Cui T, Inoue I, Shin HD. 2008. Identification of SNP markers for common CNV regions and association analysis of risk of subarachnoid aneurysmal hemorrhage in Japanese population. *Biochem Biophys Res Commun* 373:593-6.

Barnes C, Plagnol V, Fitzgerald T, Redon R, Marchini J, Clayton D, Hurles ME. 2008. A robust statistical method for case-control association testing with copy number variation. *Nat Genet* 40:1245-52.

Bassett AS, Marshall CR, Lionel AC, Chow EW, Scherer SW. 2008. Copy number variations and risk for schizophrenia in 22q11.2 deletion syndrome. *Hum Mol Genet* 17:4045-53.

Blauw HM, Veldink JH, van Es MA, van Vught PW, Saris CG, van der Zwaag B, Franke L, Burbach JP, Wokke JH, Ophoff RA, van den Berg LH. 2008. Copy-number variation in sporadic amyotrophic lateral sclerosis: a genome-wide screen. *Lancet Neurol* 7:319-26.

Cohen J. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement* 20:37-46.

Colella S, Yau C, Taylor JM, Mirza G, Butler H, Clouston P, Bassett AS, Seller A, Holmes CC, Ragoussis J. 2007. QuantiSNP: an Objective Bayes Hidden-Markov Model to detect and accurately map copy number variation using SNP genotyping data. *Nucleic Acids Res* 35:2013-25.

Conrad DF, Pinto D, Redon R, Feuk L, Gokcumen O, Zhang Y, Aerts J, Andrews TD, Barnes C, Campbell P, Fitzgerald T, Hu M, Ihm CH, Kristiansson K, Macarthur DG, Macdonald JR,

Deleted: 29/09/2010

1
2 Onyiah I, Pang AW, Robson S, Stirrups K, Valsesia A, Walter K, Wei J, Tyler-Smith C,
3 Carter NP, Lee C, Scherer SW, Hurles ME. 2010. Origins and functional impact of copy
4 number variation in the human genome. *Nature* 464:704-12.

5
6
7 Cooper GM, Zerr T, Kidd JM, Eichler EE, Nickerson DA. 2008. Systematic assessment of
8 copy number variant detection via genome-wide SNP genotyping. *Nat Genet* 40:1199-203.

9
10 Cronin S, Blauw HM, Veldink JH, van Es MA, Ophoff RA, Bradley DG, van den Berg LH,
11 Hardiman O. 2008. Analysis of genome-wide copy number variation in Irish and Dutch ALS
12 populations. *Hum Mol Genet* 17:3392-8.

13
14
15 Curtis C, Lynch AG, Dunning MJ, Spiteri I, Marioni JC, Hadfield J, Chin SF, Brenton JD,
16 Tavare S, Caldas C. 2009. The pitfalls of platform comparison: DNA copy number array
17 technologies assessed. *BMC Genomics* 10:588.

18
19
20 Dellinger AE, Saw SM, Goh LK, Seielstad M, Young TL, Li YJ. 2010. Comparative analyses
21 of seven algorithms for copy number variant identification from single nucleotide
22 polymorphism arrays. *Nucleic Acids Res* 38:e105.

23
24
25 Diskin SJ, Hou C, Glessner JT, Attiyeh EF, Laudenslager M, Bosse K, Cole K, Mosse YP,
26 Wood A, Lynch JE, Pecor K, Diamond M, Winter C, Wang K, Kim C, Geiger EA, McGrady
27 PW, Blakemore AI, London WB, Shaikh TH, Bradfield J, Grant SF, Li H, Devoto M,
28 Rappaport ER, Hakonarson H, Maris JM. 2009. Copy number variation at 1q21.1 associated
29 with neuroblastoma. *Nature* 459:987-91.

30
31
32 Feuk L, Carson AR, Scherer SW. 2006. Structural variation in the human genome. *Nat Rev*
33 *Genet* 7:85-97.

34
35
36 Franke L, de Kovel CG, Aulchenko YS, Trynka G, Zhernakova A, Hunt KA, Blauw HM, van
37 den Berg LH, Ophoff R, Deloukas P, van Heel DA, Wijmenga C. 2008. Detection,
38 imputation, and association analysis of small deletions and null alleles on oligonucleotide
39 arrays. *Am J Hum Genet* 82:1316-33.

40
41
42 Friedman JM, Baross A, Delaney AD, Ally A, Arbour L, Armstrong L, Asano J, Bailey DK,
43 Barber S, Birch P, Brown-John M, Cao M, Chan S, Charest DL, Farnoud N, Fernandes N,
44 Flibotte S, Go A, Gibson WT, Holt RA, Jones SJ, Kennedy GC, Krzywinski M, Langlois S,
45 Li HI, McGillivray BC, Nayar T, Pugh TJ, Rajcan-Separovic E, Schein JE, Schnerch A,
46 Siddiqui A, Van Allen MI, Wilson G, Yong SL, Zahir F, Eydoux P, Marra MA. 2006.

47
48
49
50
51 Deleted: 29/09/2010

52
53 | Accuracy Ms (30/09/2010)

54 20

1
2 Oligonucleotide microarray analysis of genomic imbalance in children with mental
3 retardation. *Am J Hum Genet* 79:500-13.
4

5 Garcia-Closas M, Malats N, Silverman D, Dosemeci M, Kogevinas M, Hein DW, Tardon A,
6 Serra C, Carrato A, Garcia-Closas R, Lloreta J, Castano-Vinyals G, Yeager M, Welch R,
7 Chanock S, Chatterjee N, Wacholder S, Samanic C, Tora M, Fernandez F, Real FX, Rothman
8 N. 2005. NAT2 slow acetylation, GSTM1 null genotype, and risk of bladder cancer: Results
9 from the Spanish Bladder Cancer Study and meta-analyses. *Lancet* 366:649-659.
10

11 Glessner JT, Wang K, Cai G, Korvatska O, Kim CE, Wood S, Zhang H, Estes A, Brune CW,
12 Bradfield JP, Imielinski M, Frackelton EC, Reichert J, Crawford EL, Munson J, Sleiman PM,
13 Chiavacci R, Annaiah K, Thomas K, Hou C, Glaberson W, Flory J, Otieno F, Garris M,
14 Soorya L, Klei L, Piven J, Meyer KJ, Anagnostou E, Sakurai T, Game RM, Rudd DS,
15 Zurawiecki D, McDougle CJ, Davis LK, Miller J, Posey DJ, Michaels S, Kolevzon A,
16 Silverman JM, Bernier R, Levy SE, Schultz RT, Dawson G, Owley T, McMahon WM,
17 Wassink TH, Sweeney JA, Nurnberger JI, Coon H, Sutcliffe JS, Minshew NJ, Grant SF,
18 Bucan M, Cook EH, Buxbaum JD, Devlin B, Schellenberg GD, Hakonarson H. 2009. Autism
19 genome-wide copy number variation reveals ubiquitin and neuronal genes. *Nature* 459:569-
20 73.
21

22 Gonzalez JR, Subirana I, Escaramis G, Peraza S, Caceres A, Estivill X, Armengol L. 2009.
23 Accounting for uncertainty when assessing association between copy number and disease: a
24 latent class model. *BMC Bioinformatics* 10:172.
25

26 Greenway SC, Pereira AC, Lin JC, DePalma SR, Israel SJ, Mesquita SM, Ergul E, Conta JH,
27 Korn JM, McCarroll SA, Gorham JM, Gabriel S, Altshuler DM, Quintanilla-Dieck Mde L,
28 Artunduaga MA, Eavey RD, Plenge RM, Shadick NA, Weinblatt ME, De Jager PL, Hafler
29 DA, Breitbart RE, Seidman JG, Seidman CE. 2009. De novo copy number variants identify
30 new genes and loci in isolated sporadic tetralogy of Fallot. *Nat Genet* 41:931-5.
31

32 Iafrate AJ, Feuk L, Rivera MN, Listewnik ML, Donahoe PK, Qi Y, Scherer SW, Lee C. 2004.
33 Detection of large-scale variation in the human genome. *Nat Genet* 36:949-51.
34

35 InternationalSchizophreniaConsortium. 2008. Rare chromosomal deletions and duplications
36 increase risk of schizophrenia. *Nature* 455:237-41.
37

38 Ionita-Laza I, Perry GH, Raby BA, Klanderma B, Lee C, Laird NM, Weiss ST, Lange C.
39 2008. On the analysis of copy-number variations in genome-wide association studies: a
40 translation of the family-based association test. *Genet Epidemiol* 32:273-84.
41

Deleted: 29/09/2010

1
2 Ionita-Laza I, Rogers AJ, Lange C, Raby BA, Lee C. 2009. Genetic association analysis of
3 copy-number variation (CNV) in human disease pathogenesis. *Genomics* 93:22-6.

4
5 Itsara A, Cooper GM, Baker C, Girirajan S, Li J, Absher D, Krauss RM, Myers RM, Ridker
6 PM, Chasman DI, Mefford H, Ying P, Nickerson DA, Eichler EE. 2009. Population analysis
7 of large copy number variants and hotspots of human genetic disease. *Am J Hum Genet*
8 84:148-61.

9
10
11 Kathiresan S, Voight BF, Purcell S, Musunuru K, Ardissino D, Mannucci PM, Anand S,
12 Engert JC, Samani NJ, Schunkert H, Erdmann J, Reilly MP, Rader DJ, Morgan T, Spertus JA,
13 Stoll M, Girelli D, McKeown PP, Patterson CC, Siscovick DS, O'Donnell CJ, Elosua R,
14 Peltonen L, Salomaa V, Schwartz SM, Melander O, Altshuler D, Ardissino D, Merlini PA,
15 Berzuini C, Bernardinelli L, Peyvandi F, Tubaro M, Celli P, Ferrario M, Fetiveau R,
16 Marziliano N, Casari G, Galli M, Ribichini F, Rossi M, Bernardi F, Zonzin P, Piazza A,
17 Mannucci PM, Schwartz SM, Siscovick DS, Yee J, Friedlander Y, Elosua R, Marrugat J,
18 Lucas G, Subirana I, Sala J, Ramos R, Kathiresan S, Meigs JB, Williams G, Nathan DM,
19 MacRae CA, O'Donnell CJ, Salomaa V, Havulinna AS, Peltonen L, Melander O, Berglund G,
20 Voight BF, Kathiresan S, Hirschhorn JN, Asselta R, Duga S, Spreafico M, Musunuru K, Daly
21 MJ, Purcell S, Voight BF, Purcell S, Nemes J, Korn JM, McCarroll SA, Schwartz SM, Yee
22 J, Kathiresan S, Lucas G, Subirana I, Elosua R, Surti A, Guiducci C, Gianniny L, Mirel D,
23 Parkin M, Burt N, Gabriel SB, Samani NJ, Thompson JR, Braund PS, Wright BJ, Balmforth
24 AJ, Ball SG, Hall AS, Schunkert H, Erdmann J, Linsel-Nitschke P, Lieb W, Ziegler A, Konig
25 I, Hengstenberg C, Fischer M, Stark K, Grosshennig A, Preuss M, Wichmann HE, Schreiber
26 S, Schunkert H, Samani NJ, Erdmann J, Ouwehand W, Hengstenberg C, Deloukas P, Scholz
27 M, Cambien F, Reilly MP, Li M, Chen Z, Wilensky R, Matthai W, Qasim A, Hakonarson
28 HH, Devaney J, Burnett MS, Pichard AD, Kent KM, Satler L, Lindsay JM, Waksman R,
29 Epstein SE, Rader DJ, Scheffold T, Berger K, Stoll M, Huge A, Girelli D, Martinelli N,
30 Olivieri O, Corrocher R, Morgan T, Spertus JA, McKeown P, Patterson CC, Schunkert H,
31 Erdmann E, Linsel-Nitschke P, Lieb W, Ziegler A, Konig IR, Hengstenberg C, Fischer M,
32 Stark K, Grosshennig A, Preuss M, Wichmann HE, Schreiber S, Holm H, Thorleifsson G,
33 Thorsteinsdottir U, Stefansson K, Engert JC, Do R, Xie C, Anand S, Kathiresan S, Ardissino
34 D, Mannucci PM, Siscovick D, O'Donnell CJ, Samani NJ, Melander O, Elosua R, Peltonen L,
35 Salomaa V, Schwartz SM, Altshuler D. 2009. Genome-wide association of early-onset
36 myocardial infarction with single nucleotide polymorphisms and copy number variants. *Nat*
37 *Genet* 41:334-41.

Deleted: 29/09/2010

51
52
53 | Accuracy Ms (30/09/2010)

22

1
2 Kidd JM, Cooper GM, Donahue WF, Hayden HS, Sampas N, Graves T, Hansen N, Teague B,
3 Alkan C, Antonacci F, Haugen E, Zerr T, Yamada NA, Tsang P, Newman TL, Tuzun E,
4 Cheng Z, Ebling HM, Tusneem N, David R, Gillett W, Phelps KA, Weaver M, Saranga D,
5 Brand A, Tao W, Gustafson E, McKernan K, Chen L, Malig M, Smith JD, Korn JM,
6 McCarroll SA, Altshuler DA, Peiffer DA, Dorschner M, Stamatoyannopoulos J, Schwartz D,
7 Nickerson DA, Mullikin JC, Wilson RK, Bruhn L, Olson MV, Kaul R, Smith DR, Eichler EE.
8 2008. Mapping and sequencing of structural variation from eight human genomes. *Nature*
9 453:56-64.

10
11
12
13
14 Korbelt JO, Urban AE, Grubert F, Du J, Royce TE, Starr P, Zhong G, Emanuel BS, Weissman
15 SM, Snyder M, Gerstein MB. 2007. Systematic prediction and validation of breakpoints
16 associated with copy-number variants in the human genome. *Proc Natl Acad Sci U S A*
17 104:10110-5.

18
19
20
21 Korn JM, Kuruvilla FG, McCarroll SA, Wysoker A, Nemesh J, Cawley S, Hubbell E, Veitch
22 J, Collins PJ, Darvishi K, Lee C, Nizzari MM, Gabriel SB, Purcell S, Daly MJ, Altshuler D.
23 2008. Integrated genotype calling and association analysis of SNPs, common copy number
24 polymorphisms and rare CNVs. *Nat Genet* 40:1253-60.

25
26
27 Lin M, Wei LJ, Sellers WR, Lieberfarb M, Wong WH, Li C. 2004. dChipSNP: significance
28 curve and clustering of SNP-array-based loss-of-heterozygosity data. *Bioinformatics* 20:1233-
29 40.

30
31
32 Liu W, Sun J, Li G, Zhu Y, Zhang S, Kim ST, Sun J, Wiklund F, Wiley K, Isaacs SD, Stattin
33 P, Xu J, Duggan D, Carpten JD, Isaacs WB, Gronberg H, Zheng SL, Chang BL. 2009.
34 Association of a germ-line copy number variation at 2p24.3 and risk for aggressive prostate
35 cancer. *Cancer Res* 69:2176-9.

36
37
38 Marshall CR, Noor A, Vincent JB, Lionel AC, Feuk L, Skaug J, Shago M, Moessner R, Pinto
39 D, Ren Y, Thiruvahindrapuram B, Fiebig A, Schreiber S, Friedman J, Ketelaars CE, Vos YJ,
40 Ficiocioglu C, Kirkpatrick S, Nicolson R, Sloman L, Summers A, Gibbons CA, Teebi A,
41 Chitayat D, Weksberg R, Thompson A, Vardy C, Crosbie V, Luscombe S, Baatjes R,
42 Zwaigenbaum L, Roberts W, Fernandez B, Szatmari P, Scherer SW. 2008. Structural
43 variation of chromosomes in autism spectrum disorder. *Am J Hum Genet* 82:477-88.

44
45
46
47 Matarin M, Simon-Sanchez J, Fung HC, Scholz S, Gibbs JR, Hernandez DG, Crews C,
48 Britton A, Wavrant De Vrieze F, Brott TG, Brown RD, Jr., Worrall BB, Silliman S, Case LD,

Deleted: 29/09/2010

1
2 Hardy JA, Rich SS, Meschia JF, Singleton AB. 2008. Structural genomic variation in
3 ischemic stroke. *Neurogenetics* 9:101-8.

4
5 McCarroll SA, Altshuler DM. 2007. Copy-number variation and association studies of human
6 disease. *Nat Genet* 39:S37-42.

7
8 Need AC, Ge D, Weale ME, Maia J, Feng S, Heinzen EL, Shianna KV, Yoon W,
9 Kasperaviciute D, Gennarelli M, Strittmatter WJ, Bonvicini C, Rossi G, Jayathilake K, Cola
10 PA, McEvoy JP, Keefe RS, Fisher EM, St Jean PL, Giegling I, Hartmann AM, Moller HJ,
11 Ruppert A, Fraser G, Crombie C, Middleton LT, St Clair D, Roses AD, Muglia P, Francks C,
12 Rujescu D, Meltzer HY, Goldstein DB. 2009. A genome-wide investigation of SNPs and
13 CNVs in schizophrenia. *PLoS Genet* 5:e1000373.

14
15 Olshen AB, Venkatraman ES, Lucito R, Wigler M. 2004. Circular binary segmentation for the
16 analysis of array-based DNA copy number data. *Biostatistics* 5:557-72.

17
18 Park H, Kim JI, Ju YS, Gokcumen O, Mills RE, Kim S, Lee S, Suh D, Hong D, Kang HP,
19 Yoo YJ, Shin JY, Kim HJ, Yavartanoo M, Chang YW, Ha JS, Chong W, Hwang GR,
20 Darvishi K, Kim H, Yang SJ, Yang KS, Kim H, Hurles ME, Scherer SW, Carter NP, Tyler-
21 Smith C, Lee C, Seo JS. 2010. Discovery of common Asian copy number variants using
22 integrated high-resolution array CGH and massively parallel DNA sequencing. *Nat Genet*
23 42:400-5.

24
25 Pique-Regi R, Monso-Varona J, Ortega A, Seeger RC, Triche TJ, Asgharzadeh S. 2008.
26 Sparse representation and Bayesian detection of genome copy number alterations from
27 microarray data. *Bioinformatics* 24:309-18.

28
29 Redon R, Ishikawa S, Fitch KR, Feuk L, Perry GH, Andrews TD, Fiegler H, Shapero MH,
30 Carson AR, Chen W, Cho EK, Dallaire S, Freeman JL, Gonzalez JR, Gratacos M, Huang J,
31 Kalaitzopoulos D, Komura D, MacDonald JR, Marshall CR, Mei R, Montgomery L,
32 Nishimura K, Okamura K, Shen F, Somerville MJ, Tchinda J, Valsesia A, Woodwark C,
33 Yang F, Zhang J, Zerjal T, Zhang J, Armengol L, Conrad DF, Estivill X, Tyler-Smith C,
34 Carter NP, Aburatani H, Lee C, Jones KW, Scherer SW, Hurles ME. 2006. Global variation
35 in copy number in the human genome. *Nature* 444:444-54.

36
37 Rodriguez-Santiago B, Brunet A, Sobrino B, Serra-Juhe C, Flores R, Armengol L, Vilella E,
38 Gabau E, Guitart M, Guillamat R, Martorell L, Valero J, Gutierrez-Zotes A, Labad A,
39 Carracedo A, Estivill X, Perez-Jurado LA. 2009. Association of common copy number
40
41
42
43
44
45

46 Deleted: 29/09/2010

47
48
49
50
51
52
53
54 | Accuracy Ms (30/09/2010)

24

1
2 variants at the glutathione S-transferase genes and rare novel genomic changes with
3 schizophrenia. *Mol Psychiatry*.

4
5 Schouten JP, McElgunn CJ, Waaijer R, Zwiijnenburg D, Diepvens F, Pals G. 2002. Relative
6 quantification of 40 nucleic acid sequences by multiplex ligation-dependent probe
7 amplification. *Nucleic Acids Res* 30:e57.

8
9
10 Sha BY, Yang TL, Zhao LJ, Chen XD, Guo Y, Chen Y, Pan F, Zhang ZX, Dong SS, Xu XH,
11 Deng HW. 2009. Genome-wide association study suggested copy number variation may be
12 associated with body mass index in the Chinese population. *J Hum Genet* 54:199-202.

13
14
15 Simon-Sanchez J, Scholz S, Matarin Mdel M, Fung HC, Hernandez D, Gibbs JR, Britton A,
16 Hardy J, Singleton A. 2008. Genomewide SNP assay reveals mutations underlying Parkinson
17 disease. *Hum Mutat* 29:315-22.

18
19
20 Stefansson H, Rujescu D, Cichon S, Pietilainen OP, Ingason A, Steinberg S, Fossdal R,
21 Sigurdsson E, Sigmundsson T, Buizer-Voskamp JE, Hansen T, Jakobsen KD, Muglia P,
22 Francks C, Matthews PM, Gylfason A, Halldorsson BV, Gudbjartsson D, Thorgeirsson TE,
23 Sigurdsson A, Jonasdottir A, Jonasdottir A, Bjornsson A, Mattiasdottir S, Blondal T,
24 Haraldsson M, Magnusdottir BB, Giegling I, Moller HJ, Hartmann A, Shianna KV, Ge D,
25 Need AC, Crombie C, Fraser G, Walker N, Lonnqvist J, Suvisaari J, Tuulio-Henriksson A,
26 Paunio T, Toulopoulou T, Bramon E, Di Forti M, Murray R, Ruggeri M, Vassos E, Tosato S,
27 Walshe M, Li T, Vasilescu C, Muhleisen TW, Wang AG, Ullum H, Djurovic S, Melle I,
28 Olesen J, Kiemeny LA, Franke B, Sabatti C, Freimer NB, Gulcher JR, Thorsteinsdottir U,
29 Kong A, Andreassen OA, Ophoff RA, Georgi A, Rietschel M, Werge T, Petursson H,
30 Goldstein DB, Nothen MM, Peltonen L, Collier DA, St Clair D, Stefansson K, Kahn RS,
31 Linszen DH, van Os J, Wiersma D, Bruggeman R, Cahn W, de Haan L, Krabbendam L,
32 Myin-Germeys I. 2008. Large recurrent microdeletions associated with schizophrenia. *Nature*
33 455:232-6.

34
35
36
37
38
39
40
41 Stranger BE, Forrest MS, Dunning M, Ingle CE, Beazley C, Thorne N, Redon R, Bird CP, de
42 Grassi A, Lee C, Tyler-Smith C, Carter N, Scherer SW, Tavare S, Deloukas P, Hurles ME,
43 Dermitzakis ET. 2007. Relative impact of nucleotide and copy number variation on gene
44 expression phenotypes. *Science* 315:848-53.

45
46
47
48
49
50
51
52 Walsh T, McClellan JM, McCarthy SE, Addington AM, Pierce SB, Cooper GM, Nord AS,
53 Kusenda M, Malhotra D, Bhandari A, Stray SM, Rippey CF, Roccanova P, Makarov V,
54 Lakshmi B, Findling RL, Sikich L, Stromberg T, Merriman B, Gogtay N, Butler P, Eckstrand

Deleted: 29/09/2010

1
2 K, Noory L, Gochman P, Long R, Chen Z, Davis S, Baker C, Eichler EE, Meltzer PS, Nelson
3 SF, Singleton AB, Lee MK, Rapoport JL, King MC, Sebat J. 2008. Rare structural variants
4 disrupt multiple genes in neurodevelopmental pathways in schizophrenia. *Science* 320:539-
5 43.
6
7

8 Wang K, Li M, Hadley D, Liu R, Glessner J, Grant SF, Hakonarson H, Bucan M. 2007.
9 PennCNV: an integrated hidden Markov model designed for high-resolution copy number
10 variation detection in whole-genome SNP genotyping data. *Genome Res* 17:1665-74.
11

12 Weiss LA, Shen Y, Korn JM, Arking DE, Miller DT, Fossdal R, Saemundsen E, Stefansson
13 H, Ferreira MA, Green T, Platt OS, Ruderfer DM, Walsh CA, Altshuler D, Chakravarti A,
14 Tanzi RE, Stefansson K, Santangelo SL, Gusella JF, Sklar P, Wu BL, Daly MJ. 2008.
15 Association between microdeletion and microduplication at 16p11.2 and autism. *N Engl J*
16 *Med* 358:667-75.
17

18 Winchester L, Yau C, Ragoussis J. 2009. Comparing CNV detection methods for SNP arrays.
19 *Brief Funct Genomic Proteomic* 8:353-66.
20

21 Xu B, Roos JL, Levy S, van Rensburg EJ, Gogos JA, Karayiorgou M. 2008. Strong
22 association of de novo copy number mutations with sporadic schizophrenia. *Nat Genet*
23 40:880-5.
24

25 Yang TL, Chen XD, Guo Y, Lei SF, Wang JT, Zhou Q, Pan F, Chen Y, Zhang ZX, Dong SS,
26 Xu XH, Yan H, Liu X, Qiu C, Zhu XZ, Chen T, Li M, Zhang H, Zhang L, Drees BM,
27 Hamilton JJ, Pappasian CJ, Recker RR, Song XP, Cheng J, Deng HW. 2008. Genome-wide
28 copy-number-variation study identified a susceptibility gene, UGT2B17, for osteoporosis. *Am*
29 *J Hum Genet* 83:663-74.
30

31 Yau C, Holmes CC. 2008. CNV discovery using SNP genotyping arrays. *Cytogenet Genome*
32 *Res* 123:307-12.
33
34
35
36

37 38 39 40 41 42 **FIGURE LEGENDS**

43
44 **Figure 1.** Box plots of the distribution of kappa index estimates comparing duplicated pairs
45 for A) the SNP callings, B) the detection of CNVs according to the different algorithms, and
46 C) the number of copies assigned by the different algorithms in the regions where a CNV was
47 detected.
48
49
50
51

Deleted: 29/09/2010

1
2 **Figure 2.** Box plots of the distribution of kappa indexes comparing the callings on duplicated
3 samples by the different algorithms depending on the source of DNA.
4

5
6 **Figure 3.** Average Kappa Index for the agreement in detecting CNVs (first row) and median
7 number of CNVs across the 92 individuals (second row) for each algorithm while filtering the
8 called CNVs according the number of probes in the CNV (first column) and the length of the
9 CNV (second column).
10
11
12

13
14 **Figure 4.** Sensitivity (SE) and Specificity (SP) estimates for the presence and for the type-
15 specific CNV according to each algorithm.
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

Deleted: 29/09/2010

Table 1: Median number of CNVs detected in the 92 individuals included in this study. The results are displayed according to the algorithm applied and the source of DNA. One of the replicates was randomly selected to obtain these estimates.

| Algorithm | Source of DNA | Number of Copies | | | | Total |
|--------------|---------------|------------------|------|----|---|-------|
| | | 0 | 1 | 3 | 4 | |
| CNVpartition | All | 10 | 10 | 8 | 1 | 28 |
| | Blood | 8 | 10 | 6 | 1 | 25 |
| | Saliva | 14 | 12 | 13 | 2 | 51 |
| PennCNV | All | 5 | 31.5 | 23 | 2 | 58.5 |
| | Blood | 5 | 28 | 19 | 1 | 53 |
| | Saliva | 6 | 40 | 32 | 2 | 101 |
| QuantiSNP | All | 18.5 | 24 | 9 | 2 | 56 |
| | Blood | 18 | 22 | 8 | 1 | 51 |
| | Saliva | 20 | 30 | 12 | 4 | 90 |

Table 2: Distribution of probes in the two agreement categories (disagree and agree on calling CNV) for each of the algorithms. Results are displayed for all (All), monomorphic (Mono) and polymorphic (Poly) probes.

| | CNVpartition | | | PennCNV | | | QuantiSNP | | |
|-----------------------------|--------------|--------|--------|---------|-------|--------|-----------|--------|--------|
| | All | Mono | Poly | All | Mono | Poly | All | Mono | Poly |
| Disagree | 1085 | 113 | 972 | 1165 | 89 | 1076 | 2044 | 385 | 1659 |
| | 100% | 10.43% | 89.57% | 1 | 7.63% | 92.37% | 100% | 18.83% | 81.17% |
| Agree in calling CNV | 657 | 52 | 605 | 1196 | 50 | 1146 | 2547 | 255 | 2292 |
| | 100% | 7.97% | 92.03% | 1 | 4.16% | 95.84% | 100% | 10.00% | 90.00% |
| ratio Disagree/Agree | | 2.17 | 1.61 | | 1.78 | 0.94 | | 1.51 | 0.72 |

Or Peer Review

Supplementary Table S1. Number of Individuals, assays, and pairs analyzed (before CNV criteria) and considered in the accuracy study (after CNV criteria) and according to DNA source.

| Overall | | | Blood | | | Saliva | | | Blood / Saliva | | |
|----------------------------|--------|-------|-------------|--------|-------|-------------|--------|-------|----------------|----------|-------|
| Individuals | Assays | Pairs | Individuals | Assays | Pairs | Individuals | Assays | Pairs | Individuals | Assays | Pairs |
| Before CNV criteria | | | | | | | | | | | |
| 141 | 299 | 178 | 71 | 142 | 71 | 66 | 146 | 97 | 4 | 11 | 10 |
| 127 dup | | | 71 dup | | | 55 dup | | | 1 dup | 5 Blood | 1 B/B |
| 11 trip | | | | | | 8 trip | | | 3 trip | 6 Saliva | 2 S/S |
| 3 quadrip | | | | | | 3 quadrip | | | | | 7 B/S |
| After CNV criteria | | | | | | | | | | | |
| 92 | 186 | 96 | 63 | 126 | 63 | 29 | 60 | 33 | - | - | - |
| 90 dup | | | 63 dup | | | 27 dup | | | | | |
| 2 trip | | | | | | 2 trip | | | | | |

Assays are count by summing all duplicate, triplicate and quadruplicate samples
Pairs refer to the by-two comparisons provided by duplicate (2), triplicate (3) and quadruplicate (6) samples.

Supplementary Table S2. MLPA probes considered in the MLPA analysis.

| Probe | Chromosome | Band | Start | End |
|---------------------------|-------------------|-------------|--------------|-------------|
| SKI | 1 | 1p36.33 | 2,150,969 | 2,151,029 |
| IL1B | 2 | 2q13 | 113,306,801 | 113,306,852 |
| A_14_P103008 | 2 | 2q37.3 | 242,228,984 | 242,229,042 |
| PLCD1 | 3 | 3p22.3 | 38,026,650 | 38,026,709 |
| Chr3_46771035 | 3 | 3p21.31 | 46,781,196 | 46,781,253 |
| Chr4_69231671 | 4 | 4q13.2 | 69,109,638 | 69,109,698 |
| PCDHA9 | 5 | 5q31.1 | 140,208,267 | 140,208,335 |
| DOM3Z | 6 | 6p21.32 | 32,047,183 | 32,047,228 |
| HLA-DRB5 | 6 | 6p21.32 | 32,593,310 | 32,593,379 |
| FZD9 | 7 | 7q11.23 | 72,294,840 | 72,294,901 |
| Chr8_39356595 | 8 | 8p11.23 | 39,401,744 | 39,401,802 |
| RXRa | 9 | 9q34.2 | 136,453,357 | 136,453,414 |
| NOTCH1 | 9 | 9q34.3 | 138,523,724 | 138,523,783 |
| PPYR1 | 10 | 10q11.22 | 46,507,740 | 46,507,809 |
| ADAM8 | 10 | 10q26.3 | 134,933,411 | 134,933,468 |
| HRAS | 11 | 11p15.5 | 523,758 | 523,813 |
| A_14_P114204 | 11 | 11q13.1 | 66,952,984 | 66,953,039 |
| OR4K2 | 14 | 14q11.2 | 19,414,387 | 19,414,452 |
| Chr16_32481309 | 16 | 16p11.2 | 32,516,918 | 32,516,977 |
| chr17_415_A | 17 | 17q21.31 | 41,539,152 | 41,539,211 |
| chr17_42061812_42110026_B | 17 | 17q21.31 | 41,889,427 | 41,889,486 |
| NSF | 17 | 17q21.32 | 42,166,492 | 42,166,551 |
| STK11 | 19 | 19p13.3 | 1,171,375 | 1,171,442 |
| ENm007_1 | 19 | 19q13.42 | 59,427,206 | 59,427,263 |
| ENm007_2 | 19 | 19q13.42 | 59,968,534 | 59,968,593 |
| A_14_P105195 | 20 | 20q11.21 | 30,111,471 | 30,111,530 |
| GSTT1 | 22 | 22q11.23 | 22,706,190 | 22,706,250 |
| Chr22_22690592 | 22 | 22q11.23 | 22,709,442 | 22,709,496 |
| Chr22_Pop_1 | 22 | 22q13.1 | 37,684,655 | 37,684,714 |

Supplementary Table S3. Validity estimates for blood samples comparing the calling results with those obtained using MLPA as a reference. The estimates and their 95% confidence intervals (CI) for sensitivity (SE), specificity (SP), positive predictive value (VPP) and negative predictive value (VPN) are displayed according to the algorithms and the different types of aberrations with and without filtering using the Itsara *et al.* criteria.

| Steps | CNV type | | CNVpartition | | | | PennCNV | | | | QuantiSNP | | | |
|-----------|-----------------------------|------------|--------------|---------------|----------------------|---------------|-----------|---------------|----------------------|---------------|-----------|---------------|----------------------|---------------|
| | | | No filter | | Itsara et al. filter | | No filter | | Itsara et al. filter | | No filter | | Itsara et al. filter | |
| | | | Est. | 95% CI | Est. | 95% CI | Est. | 95% CI | Est. | 95% CI | Est. | 95% CI | Est. | 95% CI |
| 1 | CNV | SE | 0.19 | [0.14 - 0.23] | 0.05 | [0.03 - 0.08] | 0.23 | [0.18 - 0.28] | 0.07 | [0.05 - 0.11] | 0.28 | [0.23 - 0.33] | 0.08 | [0.05 - 0.11] |
| | | SP | 0.99 | [0.98 - 1.00] | 1.00 | [1.00 - 1.00] | 0.98 | [0.97 - 0.99] | 1.00 | [0.99 - 1.00] | 0.97 | [0.96 - 0.98] | 1.00 | [0.99 - 1.00] |
| | | VPP | 0.83 | [0.73 - 0.91] | 0.95 | [0.75 - 1.00] | 0.76 | [0.66 - 0.85] | 0.86 | [0.68 - 0.96] | 0.71 | [0.62 - 0.79] | 0.90 | [0.73 - 0.98] |
| | | VPN | 0.83 | [0.81 - 0.85] | 0.80 | [0.78 - 0.82] | 0.84 | [0.83 - 0.86] | 0.81 | [0.79 - 0.83] | 0.86 | [0.84 - 0.87] | 0.81 | [0.79 - 0.83] |
| 2a | Deletion* | SE | 0.97 | [0.86 - 1.00] | 1.00 | [0.73 - 1.00] | 0.95 | [0.84 - 0.99] | 1.00 | [0.79 - 1.00] | 0.98 | [0.91 - 1.00] | 1.00 | [0.81 - 1.00] |
| | | SP | 1.00 | [0.79 - 1.00] | 1.00 | [0.09 - 1.00] | 1.00 | [0.83 - 1.00] | 1.00 | [0.09 - 1.00] | 0.92 | [0.73 - 0.99] | 1.00 | [0.01 - 1.00] |
| | | VPP | 1.00 | [0.86 - 1.00] | 1.00 | [0.73 - 1.00] | 1.00 | [0.87 - 1.00] | 1.00 | [0.79 - 1.00] | 0.97 | [0.88 - 1.00] | 1.00 | [0.81 - 1.00] |
| | | VPN | 0.96 | [0.79 - 1.00] | 1.00 | [0.09 - 1.00] | 0.94 | [0.79 - 0.99] | 1.00 | [0.09 - 1.00] | 0.96 | [0.78 - 1.00] | 1.00 | [0.01 - 1.00] |
| 2b | Homozygous deletion* | SE | 1.00 | [0.83 - 1.00] | 1.00 | [0.66 - 1.00] | 0.86 | [0.57 - 0.98] | 1.00 | [0.64 - 1.00] | 0.68 | [0.45 - 0.86] | 1.00 | [0.66 - 1.00] |
| | | SP | 0.94 | [0.79 - 1.00] | 1.00 | [0.42 - 1.00] | 1.00 | [0.91 - 1.00] | 1.00 | [0.66 - 1.00] | 0.92 | [0.82 - 1.00] | 1.00 | [0.68 - 1.00] |

Supplementary Material – Accuracy Ms (30/09/2010)

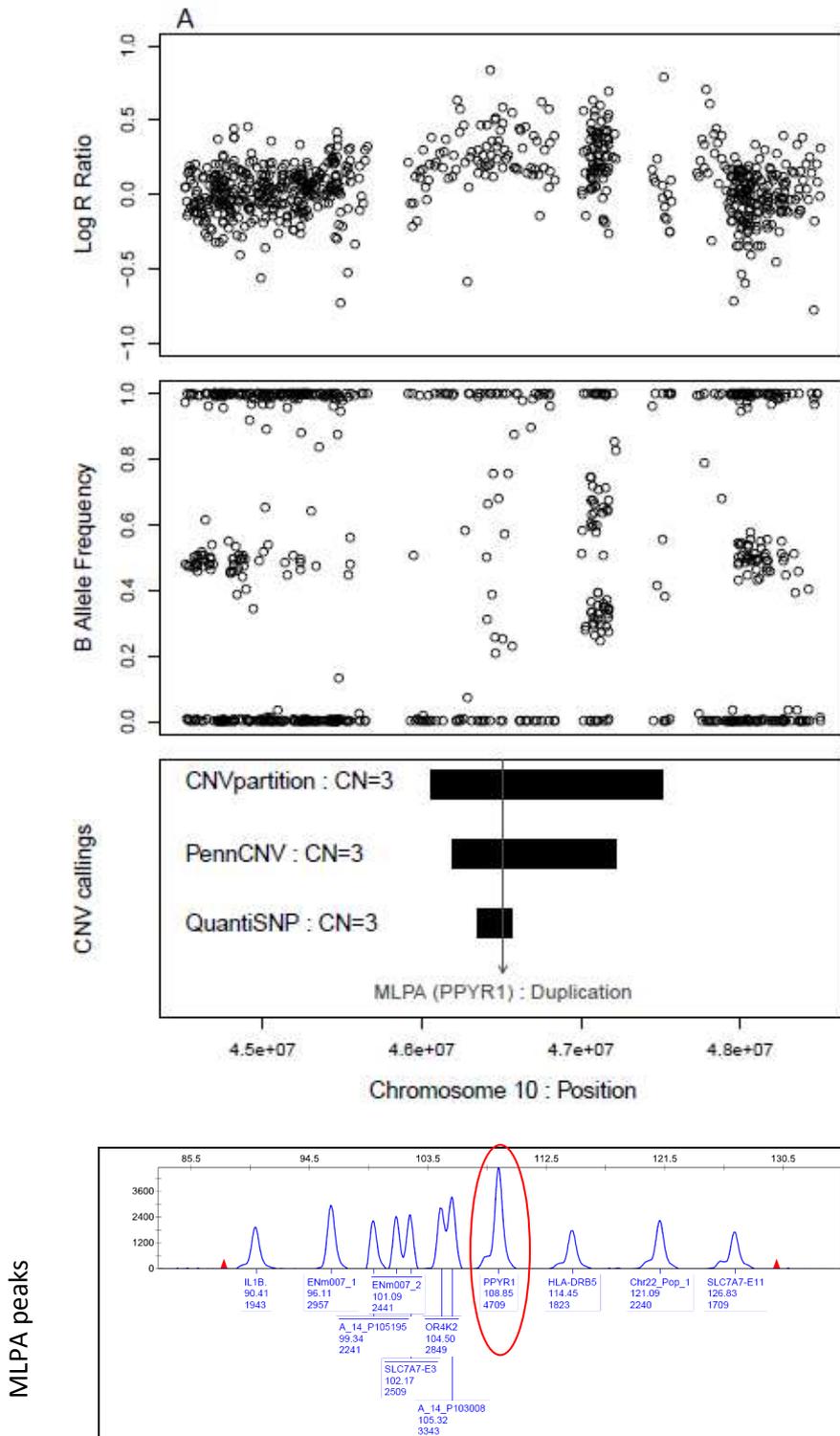
3

| | | | | | | | | | | | | | | |
|-----------|-----------------------------------|------------|-------|------------------|-------|------------------|-------|------------------|-------|------------------|-------|------------------|-------|------------------|
| | | | 0.99] | | 1.00] | | 1.00] | | 1.00] | | 0.97] | | 1.00] | |
| | | VPP | 0.94 | [0.79 - 0.99] | 1.00 | [0.66 - 1.00] | 1.00 | [0.64 - 1.00] | 1.00 | [0.64 - 1.00] | 0.75 | [0.51 - 0.91] | 1.00 | [0.66 - 1.00] |
| | | VPN | 1.00 | [0.83 - 1.00] | 1.00 | [0.42 - 1.00] | 0.97 | [0.88 - 1.00] | 1.00 | [0.66 - 1.00] | 0.89 | [0.78 - 0.95] | 1.00 | [0.68 - 1.00] |
| 2c | Heterozygous deletion* | SE | 0.63 | [0.24 - 0.91] | 1.00 | [0.28 - 1.00] | 0.93 | [0.76 - 0.99] | 1.00 | [0.62 - 1.00] | 0.92 | [0.78 - 0.98] | 1.00 | [0.66 - 1.00] |
| | | SP | 1.00 | [0.9 - 1.00] | 1.00 | [0.7 - 1.00] | 0.95 | [0.84 - 0.99] | 1.00 | [0.68 - 1.00] | 0.87 | [0.74 - 0.95] | 1.00 | [0.68 - 1.00] |
| | | VPP | 1.00 | [0.36 - 1.00] | 1.00 | [0.28 - 1.00] | 0.93 | [0.76 - 0.99] | 1.00 | [0.62 - 1.00] | 0.85 | [0.7 - 0.94] | 1.00 | [0.66 - 1.00] |
| | | VPN | 0.95 | [0.85 - 0.99] | 1.00 | [0.7 - 1.00] | 0.95 | [0.84 - 0.99] | 1.00 | [0.68 - 1.00] | 0.93 | [0.81 - 0.99] | 1.00 | [0.68 - 1.00] |
| 2d | Duplication* | SE | 1.00 | [0.79 - 1.00] | 1.00 | [0.09 - 1.00] | 1.00 | [0.83 - 1.00] | 1.00 | [0.09 - 1.00] | 0.92 | [0.73 - 0.99] | 1.00 | [0.01 - 1.00] |
| | | SP | 0.97 | [0.86 - 1.00] | 1.00 | [0.73 - 1.00] | 0.95 | [0.84 - 0.99] | 1.00 | [0.79 - 1.00] | 0.98 | [0.91 - 1.00] | 1.00 | [0.81 - 1.00] |
| | | VPP | 0.96 | [0.79 - 1.00] | 1.00 | [0.09 - 1.00] | 0.94 | [0.79 - 0.99] | 1.00 | [0.09 - 1.00] | 0.96 | [0.78 - 1.00] | 1.00 | [0.01 - 1.00] |
| | | VPN | 1.00 | [0.86 - 1.00] | 1.00 | [0.73 - 1.00] | 1.00 | [0.87 - 1.00] | 1.00 | [0.79 - 1.00] | 0.97 | [0.88 - 1.00] | 1.00 | [0.81 - 1.00] |

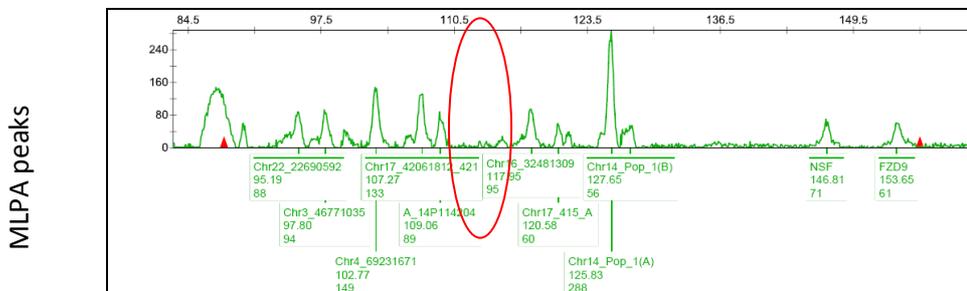
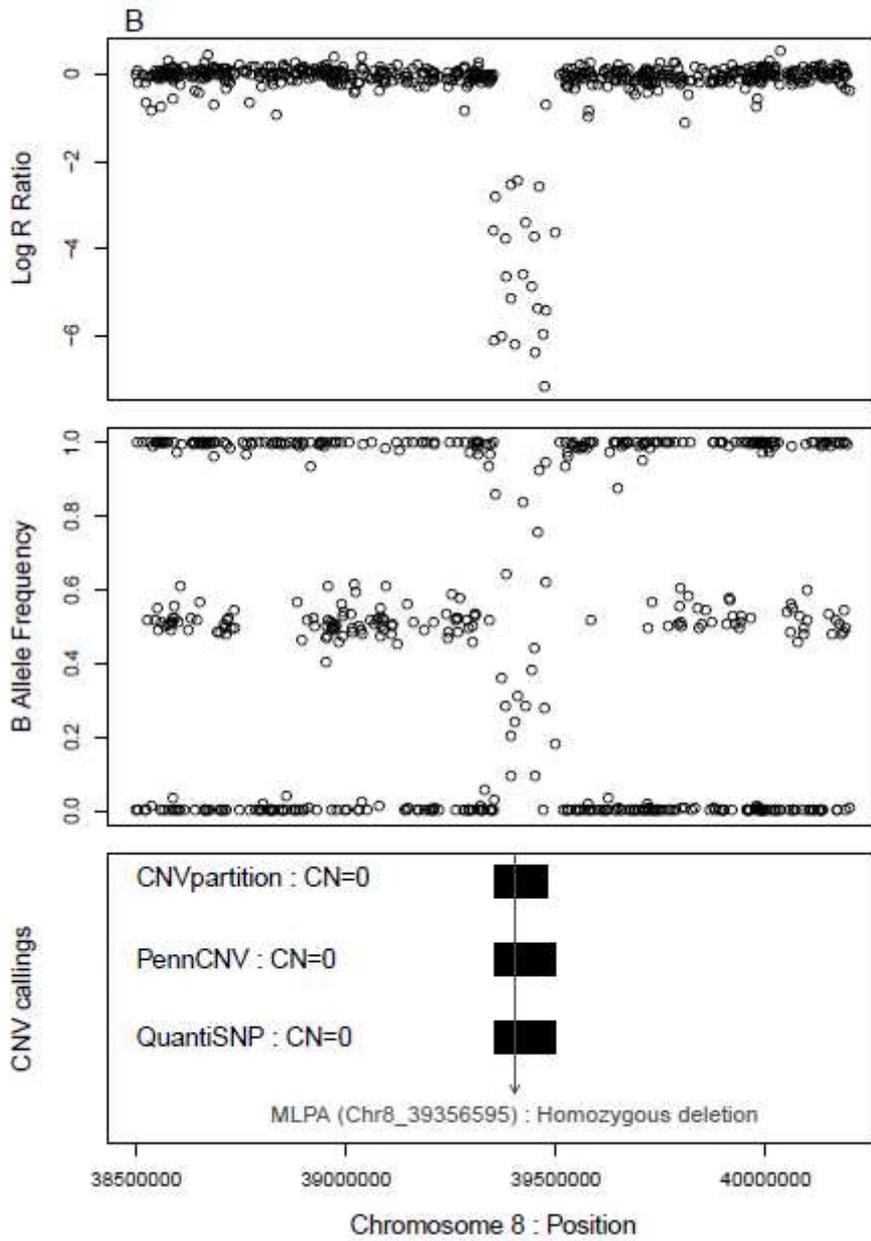
*Estimates for each CNV type were calculated only for these true positive CNVs identified in step 1.

Supplementary Figure S1: Log R Ratio (LRR), B Allele Frequency (BAF), algorithm and MLPA callings and MLPA peaks for A) a true positive duplication, B) a true positive homozygous deletion, C) a false negative heterozygous deletion and D) a false positive duplication. MLPA peaks are shown for the considering individual and for various probes used for validation.

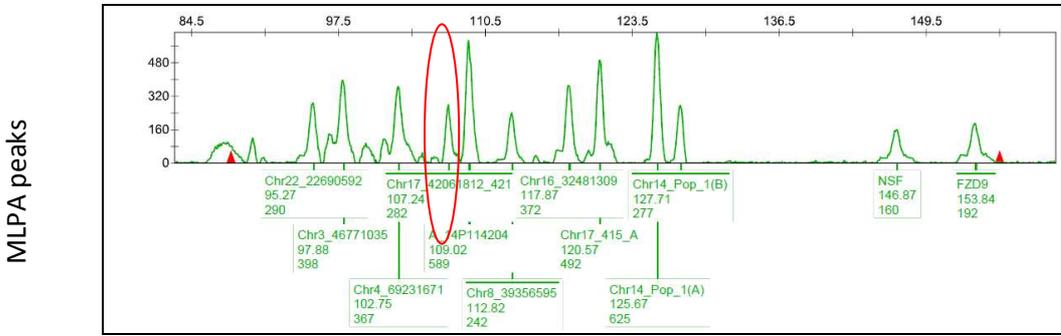
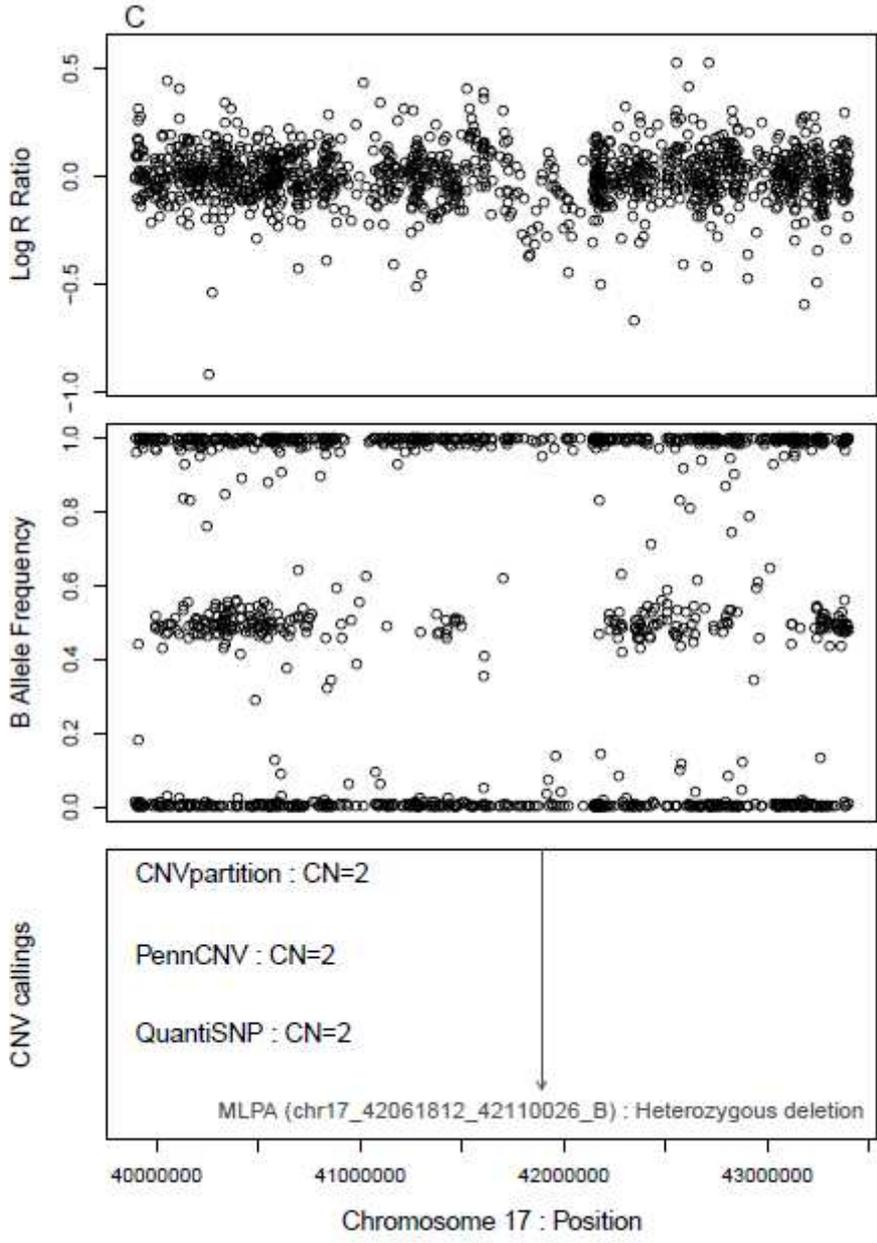
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



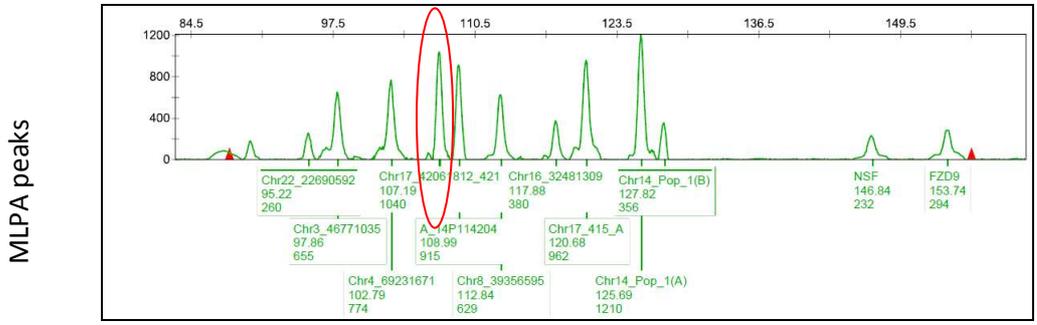
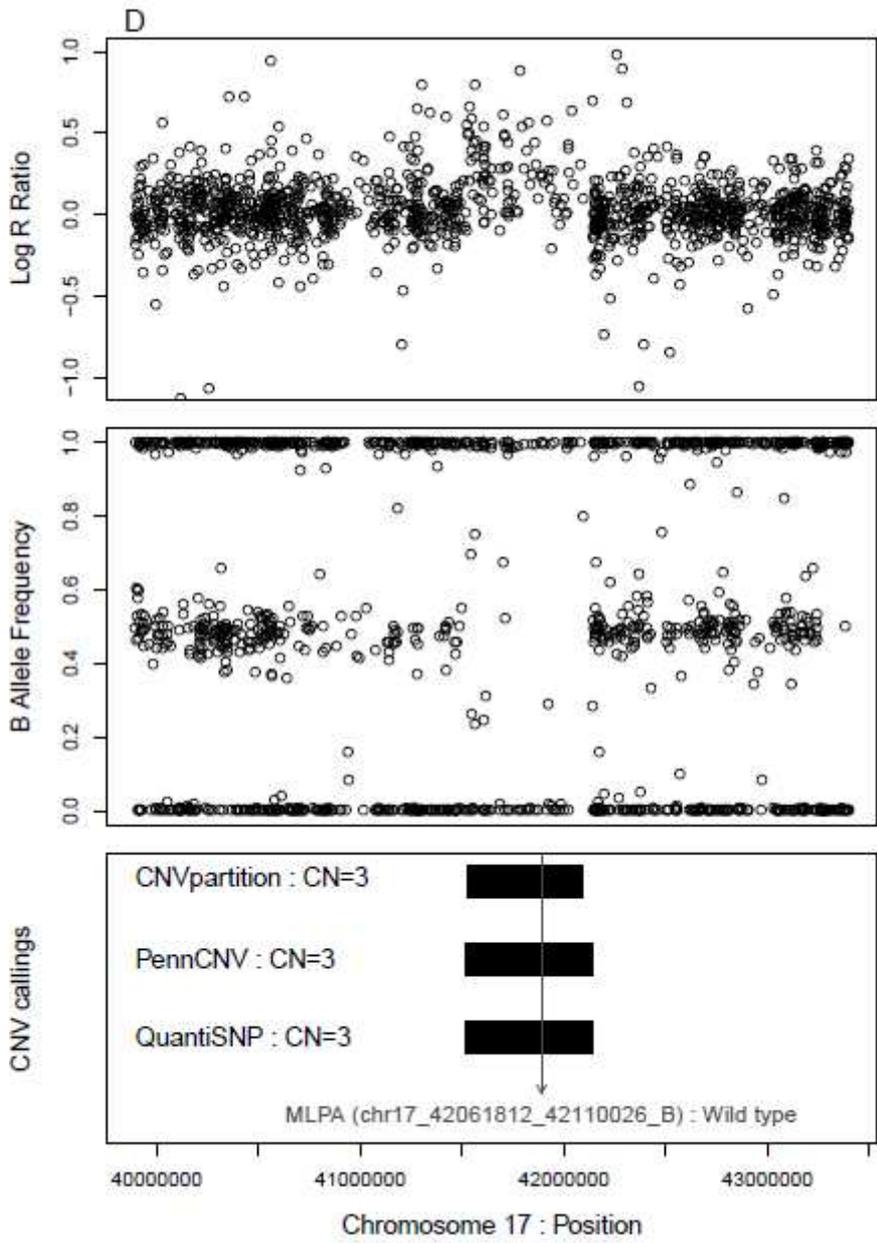
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



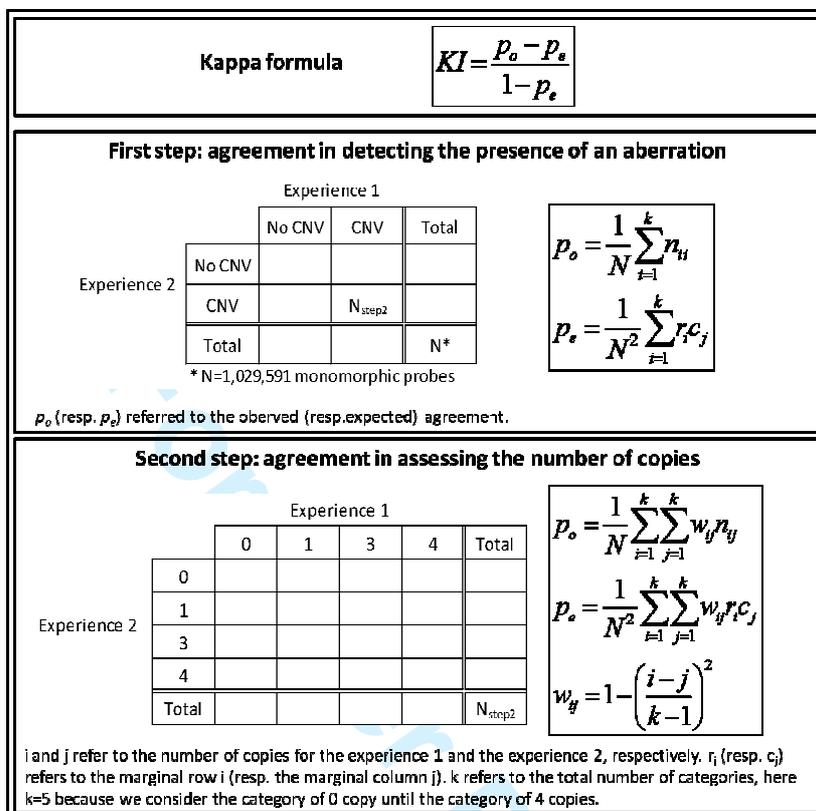
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



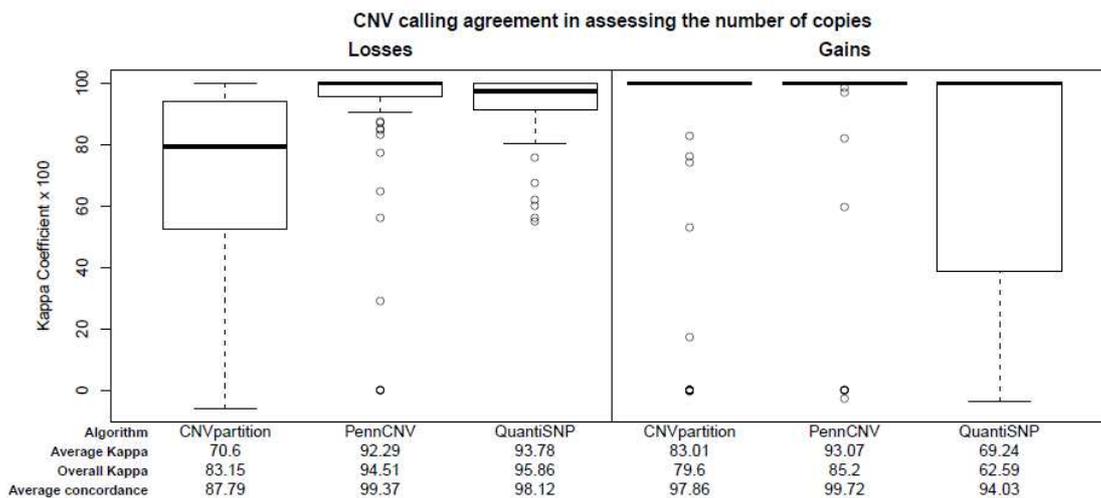
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



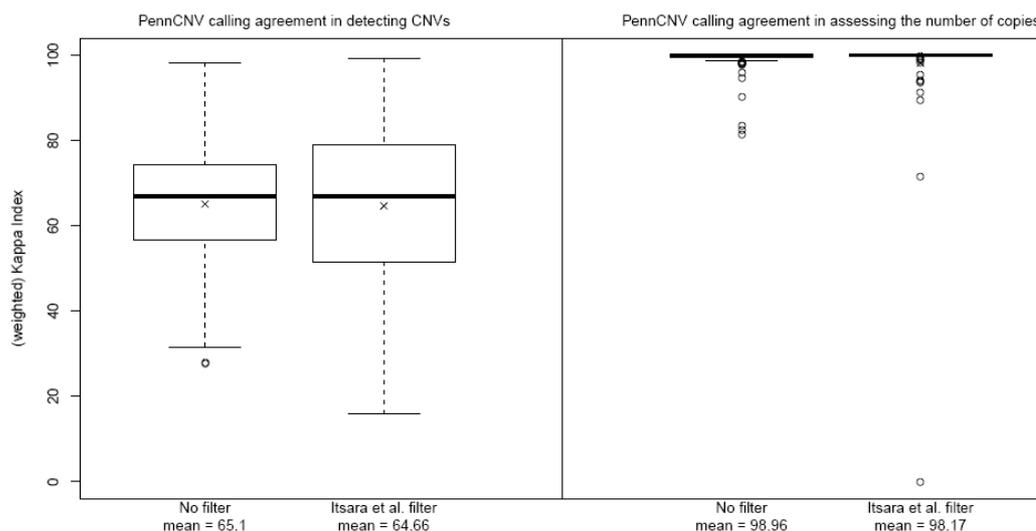
Supplementary Figure S2: Detail of the kappa calculation for the two-step agreement on calling CNVs.



Supplementary Figure S3: Agreement on assessing the number of copies once the type of CNV (loss or gain) was concordant for both replicates. For each type of CNV and each algorithm, we computed 1) the Kappa coefficient for each pair of duplicate and we provided the average Kappa across the 96 pairs, 2) a overall Kappa coefficient computed over all the 96 pairs of replicates and concordant probes, and 3) the classic concordance rate for each pair of duplicate and we provided the average concordance across the 96 pairs.

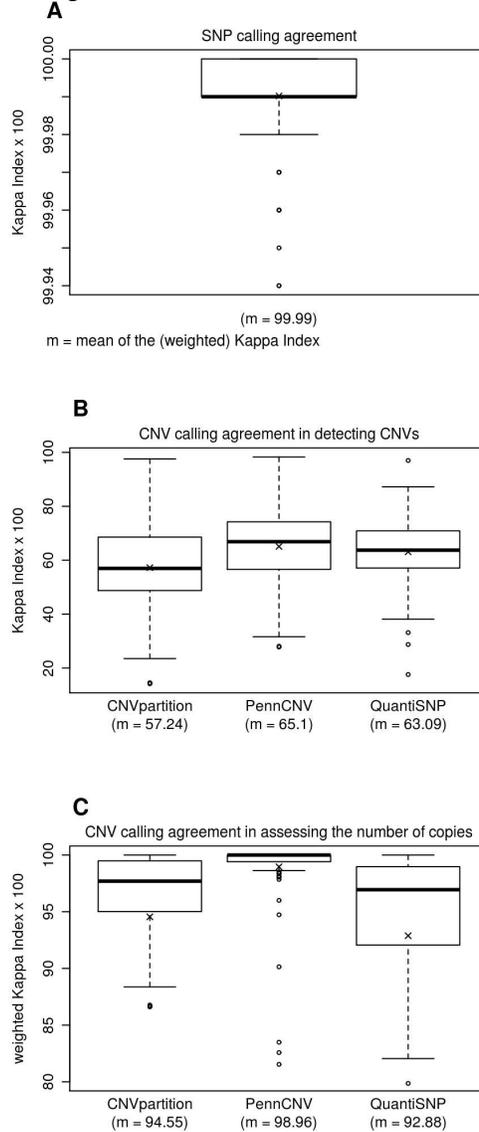


Supplementary Figure S4: Impact of the filtering on PennCNV calling agreement. Box plots before and after filtering for the distribution of A) Kappa Index estimates for CNV detection on duplicated samples, and B) weighted Kappa Index estimates for copy-number assessment when a CNV was detected.



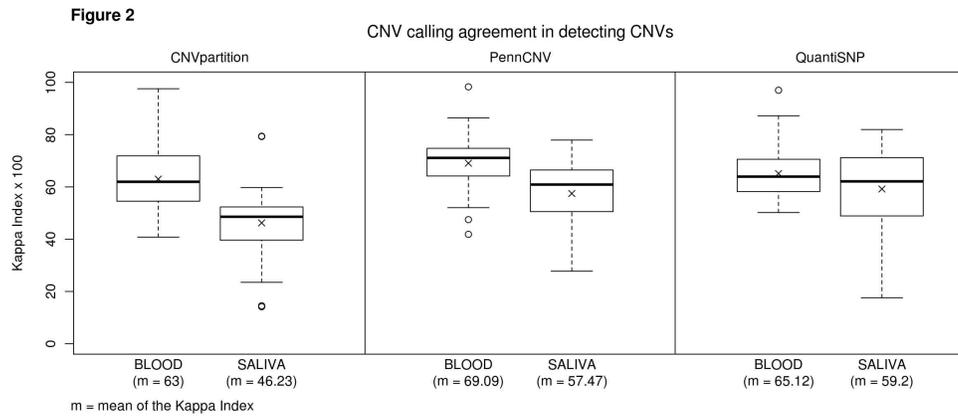
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 1



Box plots of the distribution of kappa index estimates comparing duplicated pairs for A) the SNP callings, B) the detection of CNVs according to the different algorithms, and C) the number of copies assigned by the different algorithms in the regions where a CNV was detected.
114x266mm (200 x 200 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

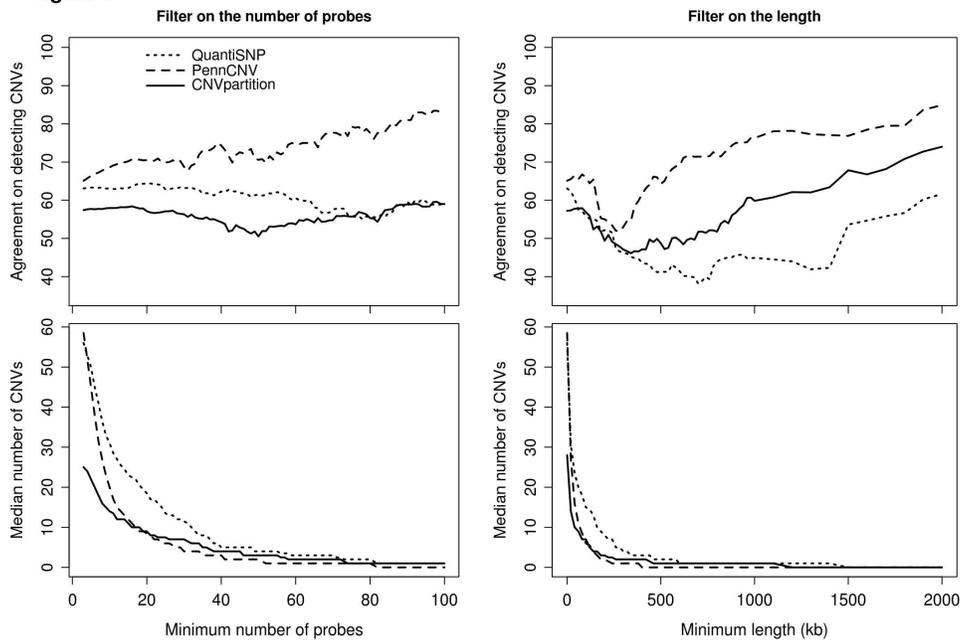


Box plots of the distribution of kappa indexes comparing the callings on duplicated samples by the different algorithms depending on the source of DNA.
304x133mm (200 x 200 DPI)

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

Figure 3

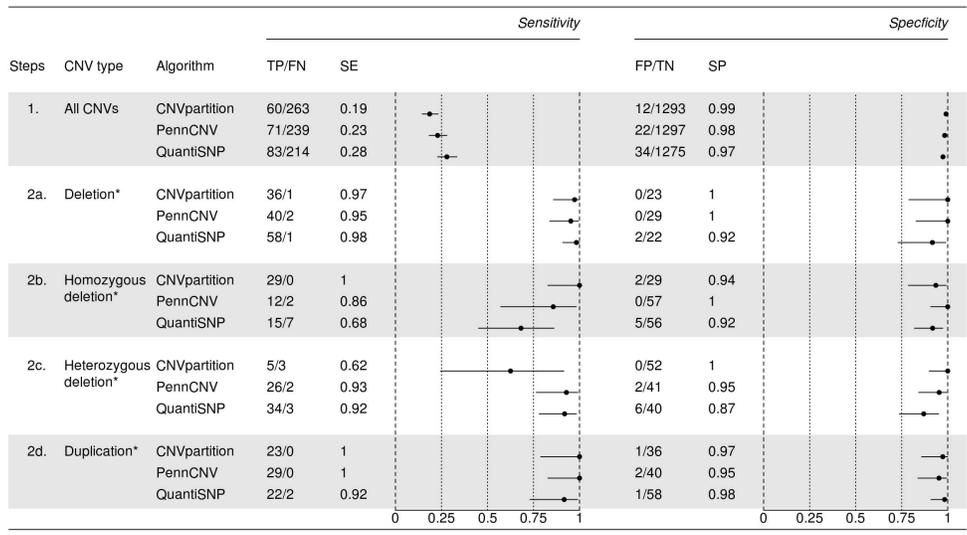


Average Kappa Index for the agreement in detecting CNVs (first row) and median number of CNVs across the 92 individuals (second row) for each algorithm while filtering the called CNVs according the number of probes in the CNV (first column) and the length of the CNV (second column).
279x190mm (200 x 200 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

Figure 4



*Estimates for each CNV type were calculated only for these true positive CNVs identified in step 1.

Sensitivity (SE) and Specificity (SP) estimates for the presence and for the type-specific CNV according to each algorithm.
304x190mm (200 x 200 DPI)

Review