

1           **Accurate sample assignment in a multiplexed, ultra-sensitive, high-**  
2           **throughput sequencing assay for minimal residual disease**

3  
4 Jack Bartram<sup>1,2</sup>, Edward Mountjoy<sup>3</sup>, Tony Brooks<sup>4</sup>, Jeremy Hancock<sup>5</sup>, Helen Williamson<sup>5</sup>, Gary  
5 Wright<sup>2</sup>, John Moppett<sup>6</sup>, Nick Goulden<sup>2</sup>, Mike Hubank<sup>1,4\*</sup>

6 <sup>1</sup>Genetics and Genomic Medicine Programme, Institute of Child Health, University College London,  
7 UK; <sup>2</sup>Department of Haematology, Great Ormond Street Hospital for Children, London, UK; <sup>3</sup>School  
8 of Social and Community Medicine, University of Bristol, UK; <sup>4</sup>UCL Genomics, Institute of Child  
9 Health, University College London, UK; <sup>5</sup>Bristol Genetics Laboratory, Southmead Hospital, North  
10 Bristol NHS Trust, UK; <sup>6</sup>Department of Paediatric Haematology/Oncology, Royal Hospital for  
11 Children, Bristol, UK

12 **Short running head:** Accurate multiplexed high-throughput MRD

13 \***Correspondence:** Dr Mike Hubank, Genetics and Genomic Medicine Programme, Institute of Child  
14 Health, University College London, 30 Guilford Street, London, WC1N 1EH, Telephone: (+44)  
15 02079052266, Fax:       Email: m.hubank@ucl.ac.uk

16 **Text pages:** 16

17 **Text word count:** 4388

18 **Abstract word count:** 192

19 **Figures:** 5

20 **Tables:** 5

21 **Reference count:** 46

22 **Grant numbers:** GOSHCC W1069; Above & Beyond project 4/2012-13

23 **Conflicts of interest:** The authors have no conflicts of interest to declare

24 **ABSTRACT**

25 High throughput sequencing (HTS) (next generation sequencing) of the rearranged immunoglobulin  
26 and T-cell receptor genes promises to be cheaper and more sensitive than current methods for  
27 monitoring minimal residual disease (MRD) in patients with acute lymphoblastic leukemia. However,  
28 adoption of new approaches by clinical laboratories requires careful evaluation of all potential sources  
29 of error and the development of strategies to ensure the highest accuracy. Timely and efficient clinical  
30 use of HTS platforms will depend on combining multiple samples (multiplexing) in each sequencing  
31 run. Here we examine *immunoglobulin heavy chain* gene HTS on the Illumina MiSeq platform for  
32 MRD (HTS-MRD). We identify errors associated with multiplexing that could potentially impact on  
33 the accuracy of MRD analysis. We optimise a strategy combining high purity, sequence-optimised  
34 oligonucleotides, dual-indexing and an error-aware demultiplexing approach to minimise errors and  
35 maximise sensitivity. We present a probability-based demultiplexing pipeline, Error-Aware  
36 Demultiplexer (EAD) - that is suitable for all MiSeq sequencing strategies and accurately assigns  
37 samples to the correct identifier without excessive loss of data. Finally using controls quantified by  
38 digital PCR, we show that HTS-MRD can accurately detect as few as 1 in  $10^6$  copies of specific  
39 leukemic MRD.

## 40 **Introduction**

41 The accurate determination of minimal residual disease (MRD) during the early months of therapy in  
42 acute lymphoblastic leukemia (ALL), particularly childhood ALL, is well established as a biomarker  
43 for guiding therapy<sup>1-6</sup>. Current methods for MRD measurement – allele-specific real-time quantitative  
44 (ASO-RQ) PCR of clone-defining immunoglobulin (IG)/T-cell receptor (TCR) gene rearrangements in  
45 the patients' leukemia<sup>7, 8</sup> and flow cytometric (FC) tracking of leukemia associated  
46 immunophenotypes<sup>1, 9</sup> - are both expensive, time consuming and suffer from technical limitations.  
47 ASO-RQ PCR requires assays tailored to each individual patient and, depending on template  
48 availability and primer selection has a maximal sensitivity of  $1:1 \times 10^{-5}$  due to non-specific background  
49 amplification<sup>10</sup>. This prevents the identification of even lower risk patients who could benefit from  
50 safer protocols that further reduce treatment-related mortality. Standardisation of FC is difficult,  
51 requires experienced operators (especially pediatric samples), and inter-operator variation can lead to  
52 inconsistent reporting<sup>11</sup>. Finally, clonal architecture is dynamic. When disease relapse occurs, it can  
53 involve clones that were not identified, or only viewed as minor clones, at diagnosis and therefore  
54 were not tracked<sup>12</sup>.

55 Advances in high throughput sequencing (HTS) offer a potential solution to these problems. Highly  
56 parallel HTS can be employed to sequence the rearranged *VDJ* (*Variable, Diverse* and *Joining*) of the  
57 *immunoglobulin heavy chain (IGH)* genes, which encode the hypervariable CDR3 domain. Combined  
58 with the high capacity of HTS, this allows a single, clone-unbiased, and highly sensitive test to be  
59 applied to all patients, revealing persisting or evolving clones, potentially even if these were not the  
60 defining clones at presentation. Importantly, HTS generates exact nucleotide sequences for all clones  
61 which are unique to each leukemic clone can be traced through subsequent follow-up analysis.  
62 Several reports have demonstrated the potential advantages of HTS for the molecular characterization  
63 of haematological malignancies<sup>13-21</sup>.

64 In order to translate these advances, it is necessary to establish a high throughput sequencing MRD  
65 (HTS-MRD) workflow that is practical to operate in a clinical laboratory, cost effective, and

66 demonstrated to be rapid, accurate and reproducible<sup>22</sup>. Continuing development of HTS technologies  
67 has resulted in cheaper platforms, such as the MiSeq (Illumina), that operate at a capacity which more  
68 appropriately matches the demands of turnaround time and cost required in clinical practice<sup>23-25</sup>.  
69 Timely and efficient clinical use of HTS platforms depends on combining multiple samples in each  
70 sequencing run, or multiplexing. Indexing of samples with unique “index sequences” allows a wide  
71 degree of multiplexing to be achieved per run, maximising use of sequencing space. However the  
72 technical limitations of this multiplexing strategy require detailed investigation – in particular, how  
73 accurately are sequences matched with patient on the basis of the indexing? Small errors in MRD  
74 assignment could have serious clinical implications. Incorrectly stratifying a patient could result in  
75 under or over treating with exposure to unnecessary toxicity or increased risk of relapse. HTS-MRD  
76 experiments generate large datasets. Interpretation and reproducibility of these data are of crucial  
77 importance in the development of a reliable clinical assay<sup>26</sup>. Use of these technologies in clinical  
78 practice has been cautioned until they are fully validated<sup>27</sup>.

79 The potential for primer bias influencing MRD detection has been extensively modelled by other  
80 groups<sup>28</sup>, so we do not address this here. The amplification strategy used for the current ASO-RQ  
81 PCR assay is already clinically approved<sup>10</sup> and therefore provides the most practical basis for  
82 translation into a HTS-MRD assay. After evaluating alternatives, we chose the MiSeq as a suitable  
83 platform for delivery as it is able to sequence single read (unidirectional) libraries of sufficient length  
84 (minimum 150 nucleotides) and quality to identify clones in less than 24 hours. MiSeqDx has been  
85 FDA approved for diagnostic use in cystic fibrosis<sup>29</sup>.

86 Sequencing multiple indexed samples per run reduces costs and increases scalability; however it  
87 introduces the potential risk of assigning reads to the wrong patient sample (misassignment). We  
88 therefore systematically investigated potential sources of multiplexing error on the MiSeq that could  
89 reduce the sensitivity and accuracy of MRD identification. These include index cross-contamination,  
90 sequencing error, misassignment of indices, run-to-run carry over, and the accuracy of the  
91 demultiplexing algorithm.

92 We found significant problems with “off-the-shelf” solutions and workflows for multiplexed  
93 amplicon sequencing, with low but unacceptable levels of sample misassignment, which could  
94 potentially lead to false-positive calls in clinical use. We overcame these issues by applying a dual-  
95 indexing strategy similar to that described by Kircher *et al.*<sup>30</sup>, using high quality preparations of index  
96 oligonucleotides<sup>31</sup>, and by developing an informatics pipeline to filter out low quality sequencing  
97 reads and reduced quality index reads. Finally, we implement our workflow and demonstrate an  
98 accurate quantification strategy using a reference “spike-in” method quantified by digital PCR  
99 (dPCR) that potentially exceeds the accuracy of current approaches by at least ten-fold.

## 100 **Materials and Methods**

### 101 *Samples and cell lines*

102 Ethical approval was given (Research Ethics Committee reference 13/LO/1262) for use of  
103 appropriately consented material from patients with B-lineage ALL at Great Ormond Street Hospital  
104 for Children. Forty one pre-treatment and eight post-induction chemotherapy bone marrow (BM)  
105 samples were obtained. Pooled “normal” lymphocyte DNA came from the UK National Blood  
106 Service. The leukemic cell lines SUPB15, REH and TOM-1 were from DSMZ and BEL-1 was kindly  
107 donated by Dr RW Stam (Rotterdam, NL).

### 108 *Sample Preparation*

109 The mononuclear cell fraction of BM samples was isolated following centrifugation on Ficoll-  
110 Hypaque (density 1.077g/l). Authentication of cell lines was performed by short tandem repeat  
111 analysis using the PowerPlex-16 system (Promega). DNA was extracted according to standardised  
112 protocols<sup>32</sup> using QIAamp DNA MiniKit (Qiagen). DNA concentration was estimated using  
113 spectrophotometry (Nanodrop, Thermo Scientific), then accurately quantitated by RQ-PCR using  
114 albumin as a control/reference gene.

115

116

117 ***Oligonucleotide synthesis***

118 Single-index strategy indices were HPLC purified and synthesised (Sigma-Aldrich) as per Kozarewa  
119 *et al.*<sup>33</sup> We designed new dual-indices (Table 1) synthesised using the TruGrade process (Integrated  
120 DNA Technologies)<sup>31</sup>.

121 ***Digital PCR***

122 The CDR3-encoding regions of *IGH* genes of the cell lines described above were amplified and  
123 sequenced (HTS and Sanger) to ensure clonality and purity of sequence (Supplementary figures 1 and  
124 2). TaqMan assays were then designed for the unique CDR3 region. Reactions containing TaqMan  
125 Gene Expression Master Mix (ABI), GE sample loading reagent (Fluidigm), TaqMan assay and  
126 template DNA, were pipetted into the loading inlets of a 12.765 Digital Array (Fluidigm). The  
127 BioMark IFC controller MX (Fluidigm) was used to uniformly partition the reactions into the panels.  
128 Template molecules are partitioned throughout the panels with a high degree of randomness and  
129 independence<sup>34</sup>. Absolute copy number quantification of cell line “spike-in” *IGH* CDR3-encoding  
130 regions was performed by dPCR, using the BioMark Real-Time PCR System (Fluidigm). For each  
131 12.765 dPCR array template DNA was analysed in triplicate.

132 ***HTS strategy***

133 *IGH* genes were amplified by multiplex PCR using AmpliTaq Gold (ABI) in a 2-stage PCR (Figure 1  
134 and Supplementary table 1). In the first stage, *IGH* family primers were modified to contain partial  
135 MiSeq adaptor sequences. First stage products were purified using solid phase reversible  
136 immobilisation (SPRI) beads (Agencourt AMPure XP, Beckman Coulter), before fluorometer  
137 quantification (Qubit, Invitrogen) and DNA Bioanalyser (Agilent) quality assessment. The purified  
138 product formed the template in a second stage PCR using NEBNext High Fidelity master mix (New  
139 England Biolabs) in which sample specific indices and full MiSeq adaptor sequences were added. The  
140 indexed samples were again purified, quantified, and then normalised to create the sequencing library  
141 pool. Sequencing libraries were re-quantified by RQ-PCR using the KAPA library quantification kit  
142 for Illumina sequencing (Kapa Biosystems) or TaqMan Gene Expression Assay for Illumina Library

143 Quantification (Life Technologies). Accurate quantification of molecules bearing appropriately  
144 ligated Illumina adaptors was crucial to ensure optimal cluster density for sequencing. Sequencing  
145 mix with 5-10% PhiX (to offset low early cycle sequence complexity) was loaded onto the MiSeq  
146 following the Illumina protocol “Preparing Libraries for Sequencing on the MiSeq” (Illumina Inc.,  
147 San Diego, CA). We used a single-end read from *J* to *V* to ensure optimal quality over the CDR3-  
148 encoding region. Indexing reads were performed to identify the 8 base pair (bp) index sequences  
149 (single or dual-indexed). Sequencing runs performed in this study are listed in Table 2. The pre-  
150 processed (i.e. multiplexed bcl) data files discussed are available in controlled access format in the  
151 European Genome-phenome archive (EGAS00001001303). Custom bioinformatics pipelines were  
152 used to identify, cluster and annotate sequencing reads.

153 For evaluation of our workflow in clinical practice we used 5 randomly chosen real patient samples  
154 (MRD quantified for patient stratification by ASO RQ-PCR as gold standard). HTS workflows with  
155 and without the technical safeguards described in the paper, were then followed. Input sample was  
156 100 000 cell equivalent DNA. The HTS preparations also included separately indexed diagnostic  
157 sample of each patient at a one in ten dilution to further test the versatility of the workflow. MRD  
158 samples were “spiked” with 1, 10 and 100 copies of SUPB15, TOM-1, and REH *IGH* sequences for  
159 quantification purposes. To assess the potential limit of detection of HTS-MRD we also spiked the  
160 same quantity of cell line reference DNA into one million cell equivalents of pooled donor  
161 lymphocyte DNA.

## 162 **Results**

### 163 *Misassignment of indices*

164 We first examined the accuracy with which oligonucleotide indexing assigned reads to the correct  
165 sample using standard laboratory and analytical methods. Initially we tagged each sample with a  
166 single P7 index composed of 8-mer oligonucleotides chosen at random a panel of 96 described by  
167 Kozarewa *et al.*<sup>33</sup>. Sequencing runs were initially demultiplexed using the on-board MiSeq software  
168 which bins samples according to the 8-mer index. In addition to indices used for sample tagging in the

169 experiment, we also instructed the demultiplexing program to search for the remaining indices that  
170 had been synthesised, but were not included in the experiment.

171 In a typical experiment (A7BK7, Table 2), with total reads passing filter 18.59 million, 89.1% (16 538  
172 154) reads were assigned to indices corresponding to samples included in the experiment  
173 (Supplementary table 2), with 10.9% (2 036 018) undetermined reads not assigned to indices after  
174 demultiplexing (Supplementary table 3). Of the undetermined reads, 55% aligned to PhiX genome,  
175 added for quality control purposes and to improve cluster resolution, 41.7% to rearranged *IGH* reads;  
176 0.002% to non-rearranged *IGH* reads and 3.3% to other non-*IGH* reads, aligned elsewhere on the  
177 genome (Supplementary table 3, Supplementary figure 3). We found that 0.12% - 23,044 reads were  
178 misassigned to one of the 68 indices not used to tag samples in the run (Figure 2, Supplementary table  
179 3). Overall this suggested that > 1 in 1000 reads are misassigned by the standard on-board MiSeq  
180 demultiplexing pipeline. This could be caused by sequencing error or factory oligonucleotide cross-  
181 contamination.

### 182 ***Quality scores of misassigned reads***

183 To assess the extent of sequencing error we generated quality scores for the misassigned reads and  
184 their associated index reads using FASTQC - [www.bioinformatics.babraham.ac.uk/projects/fastqc/](http://www.bioinformatics.babraham.ac.uk/projects/fastqc/) -  
185 (representative examples from run A7ELC on Table 2 are shown in Figure 3). We discovered that the  
186 quality of misassigned sample and index reads (Figure 3D, F) is inferior to that of reads assigned to  
187 real samples (Figure 3A, C) and deteriorates with sequence length, resulting in poor mean quality  
188 scores (Figure 3E). For the index reads, average Phred quality score for correctly assigned reads was  
189 32.5 (range 30.4 – 35.2) (Figure 3C) compared to 19.3 (range 18.4 – 22.1) for misassigned index  
190 reads (Figure 3F). This gives a mean difference of 13.3 which was highly significant with  $p < 0.00001$   
191 (95% CI 11.5 – 15.2). These results indicate that misassignment stems at least in part from poor  
192 quality index reads and read quality filters are required to ensure the most accurate demultiplexing  
193 strategy.

194



195 ***Redesigned oligonucleotide indices***

196 Misassignment may also result from oligonucleotide cross-contamination during synthesis, as  
197 previously shown by *Quail et al*<sup>31</sup>. We designed a new set of 8-mer oligonucleotide indices with high  
198 Hamming distance to optimise maximal sequence difference, homopolymer length and GC/AT  
199 balance<sup>35, 36</sup>. These were then synthesised using the high purity TruGrade process which reduces the  
200 risk of factory cross-contamination of oligonucleotide stocks<sup>31</sup>. We adopted a dual-indexing approach  
201 previously shown to improve accurate demultiplexing<sup>30</sup>, increasing total index nucleotides to 16 bp  
202 (an 8 bp index on either end of the amplicon). Twenty-four i7 and sixteen i5 indices were designed  
203 (Table 1). In order to monitor frameshift errors, some oligonucleotides were designed with sequences  
204 which maintained a high Hamming distance compared to other sequences, but which were shifted by  
205 a single base position (Supplementary table 4). Initially we synthesised only indices i7 1-12 and i5 1-  
206 8.

207 A sequencing run was performed using this new indexing strategy (A7ELC, Table 2). Thirty-one  
208 samples were multiplexed and sequenced. All index combinations of the indices i7 1-12 and i5 1-8  
209 were entered on the sample sheet, giving a total comprised of 31 double-indexed samples, 65 index  
210 combinations synthesised but not used in the experiment, and 288 index combinations where  
211 oligonucleotides were not synthesised at all.

212 We found that for the majority of combinations, the new strategy eliminated the assignment of reads  
213 to indices not present in the sequencing mix (Table 3 and Supplementary table 5). This was the case  
214 regardless of whether the sequences were synthesised or not, suggesting that factory oligonucleotide  
215 cross-contamination is effectively eliminated by TruGrade synthesis. However, significant  
216 misassignment due to amplification or sequencing error remained a problem. As anticipated, we  
217 detected significant misassignment to indices where a frameshift had been introduced in the index.  
218 For example, even though the two index sequences have a high Hamming distance, frameshift  
219 between index i7\_05 AACTCCGC and index i7\_22 ACTCCGCA (Table 3) results in 0.6% of reads  
220 being misassigned. For subsequent experiments, we split the indexing oligonucleotides into groups

221 sharing compatible sequence combinations (Supplementary table 6), removing the potential for  
222 misassignment by frameshift error.

223 Our results show that the use of dual indices designed to maximise Hamming difference and minimise  
224 frameshift error reduces the risk of misassignment of samples in multiplexed MRD. However, it is  
225 apparent that even without frameshift error, up to 0.5% of reads are still assigned to non-existent  
226 index combinations for dual index reads (e.g. Table 3; i7\_04/i5\_06). This suggests similar errors must  
227 be present but hidden in the undetectable (mis)-assignment of reads to real indices using dual or single  
228 index strategies. This could clearly pose a concern for accurate MRD assessment.

### 229 *MiSeq on-board demultiplexing software*

230 On-board MiSeq Reporter software automatically demultiplexes MiSeq output from sequencing runs  
231 and converts binary base call (bcl) files to human readable text (fastq) files for each index on the  
232 sample sheet entered for the run. Clusters are assigned to a sample when the index sequence matches  
233 exactly but also permit assignment with a single mismatch per index read. We modified the process to  
234 assign only sample names to reads bearing a perfectly matched (i.e. 0 bp mismatch) index by  
235 demultiplexing raw sequencing output data using Illumina Consensus Assessment of Sequence and  
236 Variance (CASAVA) software version 1.8.2 (Illumina Inc., San Diego, CA) or bcl2fastq conversion  
237 software version 18.4 (Illumina Inc., San Diego, CA), allowing for no mismatches in the index  
238 sequence. The reanalysed data from run A7ELC is shown in Table 4 (and Supplementary table 7). As  
239 expected, there was a reduction in the misassigned reads (compare with Table 3), but also an  
240 unacceptable reduction in the assigned reads, with a correspondingly large increase in the  
241 undetermined bin. We noted that the quality statistics from the reads which remained misassigned  
242 were generally poor with  $\% \geq Q30$  lower for misassigned sequence reads. We therefore filtered the  
243 fastq files to remove any reads with  $\% Q30 < 70 \%$ ,  $< 80 \%$  and  $< 90 \%$  (Supplementary table 8). As  
244 expected we lost a greater proportion of reads from the misassigned group than from the correctly  
245 assigned reads but the misassigned reads were still present and a large number of reads had been  
246 discarded – potentially impacting on sensitivity. We concluded that current demultiplexing strategies

247 are not stringent enough for very high sensitivity applications, and do not take account of the quality  
248 of the index read.

#### 249 *Use of unique index combinations*

250 To further reduce the potential “hidden” misassignment, in the following run, A7FDO, we used only  
251 unique combinations of group 1 indices (Supplementary table 6). We sequenced 8 samples, >1.5  
252 million reads per sample, output data was demultiplexed using MiSeq on-board demultiplexer and  
253 through CASAVA allowing for no mismatches in index reads, reducing overall misassignment from  
254 0.2% to 0.05% (Supplementary tables 9 and 10). This shows that for MRD measurement where a  
255 high degree of accuracy (minimising misassignment) is required only dual unique indices should be  
256 used.

#### 257 *Demultiplexing based on quality of index reads – Error Aware Demultiplexer*

258 In view of the significant misassignment using the MiSeq on-board demultiplexer, the unacceptable  
259 loss of potentially informative sequences when increasing stringency to allow no mismatches, poor  
260 quality statistics of misassigned index reads, we developed our own demultiplexing pipeline. “Error-  
261 Aware Demultiplexer” (EAD) utilises base call quality scores of index reads produced during  
262 Illumina sequencing (open source available at: <https://github.com/edml/error-aware-demultiplexer>).  
263 The pipeline incorporates Phred scores to probabilistically match read indices to the sample identities  
264 during demultiplexing. Index similarity is assessed with the same algorithm used in the Illumina pair-  
265 end assembler PANDAseq<sup>37</sup>. The pipeline calculates the probability that the true index and the index  
266 read represent the same underlying sequence. For example, if two bases match and the quality of  
267 those bases are high then we have good evidence that they represent the same base. Probabilities are  
268 calculated for each base and multiplied together to get the probability that the two reads represent the  
269 same sequence. This results in up to an 80% reduction in misassignment compared to other  
270 approaches and importantly, with minimal subsequent loss of sequence reads (Table 5).

271

272

273 ***Reducing run-to-run carryover***

274 Template molecules may remain in the MiSeq fluidics system, even after a standard wash program,  
275 and can be washed onto the flow cell in subsequent runs. MiSeq instruments maintained according to  
276 standard Illumina recommendations typically have sample carryover rates of  $\leq 0.1\%$ . If the same  
277 indexing strategy is applied in a subsequent run, a 0.1% carryover rate could potentially cause errors  
278 in clinical interpretation of MRD. We added two non-human samples to the sequencing pool (A7ELC,  
279 Table 2) and then performed a sequencing run (A7FDO, Table 2) with the same indices but different  
280 samples 3 weeks later and after a further 2 different runs on the MiSeq (and therefore three standard  
281 washes). Carryover was detected at a rate of 0.002%. Performing the wash recipe recommended by  
282 Illumina for highly sensitive applications “Technical support: Reducing Run-to-Run Carryover on the  
283 MiSeq Using Dilute Sodium Hypochlorite Solution” (Illumina Inc., San Diego, CA), completely  
284 eradicated the carryover (data not shown). We observed that this wash recipe can sometimes cause  
285 lower cluster densities, presumably related to small amounts of sodium hypochlorite washed onto the  
286 flow cell.

287 ***Accurate quantification using a dPCR calculated reference “spike-in” – evidence for one in a***  
288 ***million cells detection sensitivity***

289 HTS-MRD is theoretically limited by the number of cells input. To accurately quantify sensitivity, we  
290 performed a “spike-in” experiment. While other groups have used plasmids or synthetic templates<sup>18</sup>,  
291 <sup>28</sup>, we used a pre-determined quantity of reference *IGH* DNA target derived from B-cell leukemia  
292 lines with unique *IGH* clonotypes extracted, prepared and accurately quantified using dPCR (data not  
293 shown). Dilutions ranging from  $1 \times 10^{-4}$  to  $1 \times 10^{-6}$  were created by adding known quantities (1, 10 and  
294 100 copies) of SUPB15, TOM-1, and REH *IGH* sequences into one million cell equivalents of pooled  
295 donor lymphocyte DNA. We then amplified the DNA of all one million cells and sequenced the  
296 products using the HTS strategy described above. The sequencing run (ABG7Y, Table 2) was  
297 demultiplexed using EAD resulting in over 2.5 million reads per sample. We found using our  
298 workflow HTS-MRD achieved linear amplification, with  $R^2 > 0.9998$ , good reproducibility (Figure 4)

299 and the ability to detect the spike-in cell line *IGH* target down to one copy in a million normal cell  
300 equivalents.

### 301 *Application to clinical MRD samples*

302 To compare the effectiveness of our workflow with a standard HTS approach, we prepared libraries  
303 from five clinical samples previously scored for MRD by the gold standard ASO-PCR assay and  
304 classified into risk categories based on a 0.01% threshold. Each sample was indexed and sequenced  
305 together with  $10^{-1}$  dilutions of its own diagnostic sample tagged with a different index.  
306 Demultiplexing and pre-processing errors in standard HTS workflows results in false positive risk  
307 classification due to the misassignment of diagnostic sample to patient MRD (Figure 5A). In the clinic  
308 this would lead to overtreatment of patients. The scenario was then repeated using the improved  
309 workflow with EAD. The samples were now all correctly classified (Figure 5B). This experiment  
310 demonstrates the potential clinical consequences of misclassification, and the power of an improved  
311 workflow to prevent this occurring even with a one in ten dilution of diagnostic sample present in the  
312 material.

### 313 **Discussion**

314 Several groups have demonstrated that HTS of the rearranged *IGH* gene is a potentially sensitive  
315 method for MRD detection in patients with ALL<sup>15, 16, 18, 38</sup>. However, adoption of new approaches in  
316 clinical laboratories requires careful evaluation of all potential sources of error and the development  
317 of strategies to ensure the highest accuracy<sup>27</sup>.

318 Several important sources of potential error impact on workflow choice for HTS-MRD. Differential  
319 hybridisation kinetics of oligonucleotide primers can introduce significant biases that alter the  
320 composition of sequence libraries prepared by multiplex PCR<sup>28</sup>. However, for EuroMRD approved  
321 centres<sup>10</sup>, the accepted clinical diagnostic assay for MRD is PCR based using consensus primers and it  
322 therefore seemed reasonable to adopt the same primer sets as the basis for the proposed HTS-MRD  
323 test. The Illumina MiSeq is subject to characteristic base-calling errors, but these are significantly  
324 lower than current competing systems. *Kennedy et al.*<sup>39</sup> describe a ligation and capture based assay

325 which overcome HTS errors, but the method is not suitable for introducing molecular indices into the  
326 specific locus for IG/TCR genes. Current techniques used for sensitive HTS-MRD are all amplicon-  
327 based approaches.

328 Multiplexing of indexed samples means that HTS-MRD could be an economical clinical assay.  
329 However, our analysis identified low, but clearly detectable and clinically relevant levels of sample  
330 misassignment using the standard MiSeq demultiplexing approach. We therefore developed a strategy  
331 that reduces misassignment to the absolute minimum while maintaining maximal sensitivity.

332 We adopted TruGrade-synthesised oligonucleotide primers designed with a high Hamming distance  
333 and screened to avoid frameshift error. *Quail et al.* reported contamination rates for purification by  
334 HPLC or PAGE purification of 0.56% and 0.34%, however with TruGrade this reduced to just  
335 0.03%<sup>31</sup>. Also, for high sensitivity applications, unique combination dual-indexing, which identifies  
336 the sample origin of each sequence twice, independently, during demultiplexing, is superior to single-  
337 index multiplex sequencing<sup>30</sup>.

338 Although misassigned reads are broadly of low quality, filtering based on read quality alone is an  
339 inefficient method for improving accuracy. A better method is to filter on index read quality, so we  
340 developed a custom demultiplexing pipeline (EAD) that uses a probabilistic approach to remove  
341 unreliable index reads while optimising retention of informative sequences. EAD out-performs  
342 standard demultiplexing software, including the on-board MiSeq demultiplexer, producing high  
343 quality data, with up to 80% reduction in read misassignment. Despite reduced tolerance for  
344 inaccurate index reads, EAD achieves up to 10% more allocated reads (Table 5) than other  
345 demultiplexing methods. EAD has potential application in any assay requiring highly accurate  
346 demultiplexing (e.g. pathogen detection and single cell applications). We also confirmed run-to-run  
347 carryover in the MiSeq fluidic pathway and demonstrated that it is essential to perform high-  
348 stringency post-run washes of the MiSeq to avoid the risk of contamination. Together, our workflow  
349 reduces misassignment to less than 0.05% with no loss of potentially usable data, resulting in high  
350 quality and accurately assigned data from multiplexed sequencing experiments.

351 We have highlighted the importance of eliminating all avoidable risk in MRD analysis, and  
352 appropriate quality control measures could provide some safeguards. *Seitz et al.*<sup>40</sup>, for example,  
353 describe a novel method to prevent carry-over contaminations during similar two-step PCR protocols.  
354 It will also be good practice that no diagnostic sample should be sequenced together with its follow-  
355 up, that separate laboratory areas be used to prepare follow-up sample libraries, and that index  
356 combinations be alternated. Despite this, errors may still occur, chiefly associated with mis-  
357 identification of samples during multiplexing.

358 Firstly there is a significant risk that false-positive results may occur if two patients share exactly the  
359 same clonal sequence. Wu *et al.* evaluated this by HTS in post treatment samples and estimated the  
360 chance in B cells was 0.1% and 0.72% at 1 cell in 1 00 000 and 1 000 000 resolution respectively<sup>16</sup>.  
361 Next, the spike-in (cell lines, synthetic templates or plasmids) used for quantification (as in all current  
362 *IGH* HTS-MRD assays) will potentially all have the same clonal *IGH* sequence in each sample on the  
363 same run, risking inaccurate quantification with consequent potential for misclassification.  
364 Furthermore, simple human error in sample processing could result in misdiagnosis, or even where  
365 detected, the cost of a repeat run. We clearly demonstrate that a  $10^{-1}$  diluted diagnostic sample can  
366 result in misclassification with the standard multiplexing approach, with the likelihood presumably  
367 diminishing with clone concentration. In the clinical setting, cost-effective diagnostic sequencing will  
368 depend on a high degree of multiplexing with optimal use of sequencing space on cheaper, higher  
369 capacity sequencers. In this scenario, cross-checking and excluding potentially contaminating highly  
370 expressed clones before sequencing would be inconvenient and reduce efficiency. It is therefore  
371 notable that our optimised approach with EAD reduces misassignment errors from all these scenarios  
372 and allows correct classification even with a 10% diagnostic sample present.

373 We have also introduced highly accurate and biologically relevant sequencing reference controls for  
374 HTS-MRD. Using dPCR, we absolutely quantified the number of molecules of biologically relevant  
375 *IGH* controls present in each preparation. This is more accurate than conventional RQ-PCR or  
376 estimations based on molecular weight and concentration. The entire workflow was then applied to  
377 demonstrate sensitivity of at least one in a million cell equivalents, improving sensitivity by at least

378 ten fold compared to the current ASO-RQ MRD approach. A major advantage of such sensitivity,  
379 combined with confident sample assignment in multiplexing would be to allow application to  
380 peripheral blood rather than BM, with obvious advantages for patients (especially children), as well as  
381 potentially giving a more representative picture of MRD<sup>41-46</sup>.

382 The improvements described represent a step towards the rigorous validation required to produce a  
383 robust clinical HTS-MRD assay. Large prospective head-to-head comparison with current methods  
384 are needed to prove if the increased sensitivity and broader view of IG repertoire that HTS can  
385 achieve can replace the current burdensome and expensive techniques. We also propose EAD as a  
386 flexible method applicable beyond MRD detection to any multiplexing approaches where high  
387 accuracy of assignment is paramount.

### 388 **Acknowledgments**

389 This work was supported by Great Ormond Street Hospital Children's Charity. We wish to give  
390 special thanks to the Smurfit family for the generous support of this MRD research. The work  
391 described forms part of the EuroClonality-NGS consortium to standardise workflows. Finally we  
392 would like to acknowledge Above and Beyond, Charitable Trustees of University Hospitals Bristol  
393 for their funding which supported the original design of the MRD analysis pipeline.

394 Authorship contribution: J.B., M.H. and N.G. designed the research. J.B., G.W., and T.B. performed  
395 experiments; E.M. designed bioinformatics pipelines with help of J.B., H.W., J.H., J.M. J.B. and  
396 M.H.; J.B. and M.H. analysed results and J.B. made the figures; J.B. and M.H. wrote the paper. E.M.,  
397 T.B., G.W., H.W., J.H., N.G. and J.M. contributed to the writing of the manuscript.

398 **Conflict-of-interest disclosure:** The authors declare no competing financial interests.

399 **Correspondence:** Dr Mike Hubank, Genetics and Genomic Medicine Programme, UCL Institute of  
400 Child Health, London. E-mail: m.hubank@ucl.ac.uk

401



402 **References**

- 403 [1] Campana D: Role of minimal residual disease monitoring in adult and pediatric acute  
404 lymphoblastic leukemia. *Hematology/oncology clinics of North America* 2009, 23:1083-98, vii.
- 405 [2] Yeoh AE, Ariffin H, Chai EL, Kwok CS, Chan YH, Ponnudurai K, Campana D, Tan PL, Chan MY, Kham  
406 SK, Chong LA, Tan AM, Lin HP, Quah TC: Minimal residual disease-guided treatment deintensification  
407 for children with acute lymphoblastic leukemia: results from the Malaysia-Singapore acute  
408 lymphoblastic leukemia 2003 study. *Journal of clinical oncology : official journal of the American*  
409 *Society of Clinical Oncology* 2012, 30:2384-92.
- 410 [3] Campana D: Minimal residual disease monitoring in childhood acute lymphoblastic leukemia.  
411 *Current opinion in hematology* 2012, 19:313-8.
- 412 [4] Vora A, Goulden N, Wade R, Mitchell C, Hancock J, Hough R, Rowntree C, Richards S: Treatment  
413 reduction for children and young adults with low-risk acute lymphoblastic leukaemia defined by  
414 minimal residual disease (UKALL 2003): a randomised controlled trial. *The lancet oncology* 2013,  
415 14:199-209.
- 416 [5] Conter V, Bartram CR, Valsecchi MG, Schrauder A, Panzer-Grumayer R, Moricke A, Arico M,  
417 Zimmermann M, Mann G, De Rossi G, Stanulla M, Locatelli F, Basso G, Niggli F, Barisone E, Henze G,  
418 Ludwig WD, Haas OA, Cazzaniga G, Koehler R, Silvestri D, Bradtke J, Parasole R, Beier R, van Dongen  
419 JJ, Biondi A, Schrappe M: Molecular response to treatment redefines all prognostic factors in  
420 children and adolescents with B-cell precursor acute lymphoblastic leukemia: results in 3184  
421 patients of the AIEOP-BFM ALL 2000 study. *Blood* 2010, 115:3206-14.
- 422 [6] Yamaji K, Okamoto T, Yokota S, Watanabe A, Horikoshi Y, Asami K, Kikuta A, Hyakuna N, Saikawa  
423 Y, Ueyama J, Watanabe T, Okada M, Taga T, Kanegane H, Kogawa K, Chin M, Iwai A, Matsushita T,  
424 Shimomura Y, Hori T, Tsurusawa M: Minimal residual disease-based augmented therapy in childhood  
425 acute lymphoblastic leukemia: a report from the Japanese Childhood Cancer and Leukemia Study  
426 Group. *Pediatric blood & cancer* 2010, 55:1287-95.

427 [7] Goulden NJ, Knechtli CJ, Garland RJ, Langlands K, Hancock JP, Potter MN, Steward CG, Oakhill A:  
428 Minimal residual disease analysis for the prediction of relapse in children with standard-risk acute  
429 lymphoblastic leukaemia. *British journal of haematology* 1998, 100:235-44.

430 [8] van Dongen JJ, Seriu T, Panzer-Grumayer ER, Biondi A, Pongers-Willemse MJ, Corral L, Stolz F,  
431 Schrappe M, Masera G, Kamps WA, Gadner H, van Wering ER, Ludwig WD, Basso G, de Bruijn MA,  
432 Cazzaniga G, Hettinger K, van der Does-van den Berg A, Hop WC, Riehm H, Bartram CR: Prognostic  
433 value of minimal residual disease in acute lymphoblastic leukaemia in childhood. *Lancet* 1998,  
434 352:1731-8.

435 [9] Gaipa G, Basso G, Biondi A, Campana D: Detection of minimal residual disease in pediatric acute  
436 lymphoblastic leukemia. *Cytometry Part B, Clinical cytometry* 2013.

437 [10] van der Velden VH, Cazzaniga G, Schrauder A, Hancock J, Bader P, Panzer-Grumayer ER, Flohr T,  
438 Sutton R, Cave H, Madsen HO, Cayuela JM, Trka J, Eckert C, Foroni L, Zur Stadt U, Beldjord K, Raff T,  
439 van der Schoot CE, van Dongen JJ: Analysis of minimal residual disease by Ig/TCR gene  
440 rearrangements: guidelines for interpretation of real-time quantitative PCR data. *Leukemia* 2007,  
441 21:604-11.

442 [11] Kalina T, Flores-Montero J, van der Velden VH, Martin-Ayuso M, Bottcher S, Ritgen M, Almeida J,  
443 Lhermitte L, Asnafi V, Mendonca A, de Tute R, Cullen M, Sedek L, Vidriales MB, Perez JJ, te Marvelde  
444 JG, Mejstrikova E, Hrusak O, Szczepanski T, van Dongen JJ, Orfao A: EuroFlow standardization of flow  
445 cytometer instrument settings and immunophenotyping protocols. *Leukemia* 2012, 26:1986-2010.

446 [12] Szczepanski T, Willemse MJ, Brinkhof B, van Wering ER, van der Burg M, van Dongen JJ:  
447 Comparative analysis of Ig and TCR gene rearrangements at diagnosis and at relapse of childhood  
448 precursor-B-ALL provides improved strategies for selection of stable PCR targets for monitoring of  
449 minimal residual disease. *Blood* 2002, 99:2315-23.

450 [13] Kohlmann A, Grossmann V, Haferlach T: Integration of next-generation sequencing into clinical  
451 practice: are we there yet? *Seminars in oncology* 2012, 39:26-36.

452 [14] Logan AC, Zhang B, Narasimhan B, Carlton V, Zheng J, Moorhead M, Krampf MR, Jones CD,  
453 Waqar AN, Faham M, Zehnder JL, Miklos DB: Minimal residual disease quantification using  
454 consensus primers and high-throughput IGH sequencing predicts post-transplant relapse in chronic  
455 lymphocytic leukemia. *Leukemia* 2013, 27:1659-65.

456 [15] Ladetto M, Bruggemann M, Monitillo L, Ferrero S, Pepin F, Drandi D, Barbero D, Palumbo A,  
457 Passera R, Boccadoro M, Ritgen M, Gokbuget N, Zheng J, Carlton V, Trautmann H, Faham M, Pott C:  
458 Next-generation sequencing and real-time quantitative PCR for minimal residual disease detection in  
459 B-cell disorders. *Leukemia* 2014, 28:1299-307.

460 [16] Wu D, Emerson RO, Sherwood A, Loh ML, Angiolillo A, Howie B, Vogt J, Rieder M, Kirsch I,  
461 Carlson C, Williamson D, Wood BL, Robins H: Detection of Minimal Residual Disease in B  
462 Lymphoblastic Leukemia by High-Throughput Sequencing of IGH. *Clinical cancer research : an official*  
463 *journal of the American Association for Cancer Research* 2014, 20:4540-8.

464 [17] Logan AC, Vashi N, Faham M, Carlton V, Kong K, Buno I, Zheng J, Moorhead M, Klinger M, Zhang  
465 B, Waqar A, Zehnder JL, Miklos DB: Immunoglobulin and T cell receptor gene high-throughput  
466 sequencing quantifies minimal residual disease in acute lymphoblastic leukemia and predicts post-  
467 transplantation relapse and survival. *Biology of blood and marrow transplantation : journal of the*  
468 *American Society for Blood and Marrow Transplantation* 2014, 20:1307-13.

469 [18] Faham M, Zheng J, Moorhead M, Carlton VE, Stow P, Coustan-Smith E, Pui CH, Campana D:  
470 Deep-sequencing approach for minimal residual disease detection in acute lymphoblastic leukemia.  
471 *Blood* 2012, 120:5173-80.

472 [19] Logan AC, Gao H, Wang C, Sahaf B, Jones CD, Marshall EL, Buno I, Armstrong R, Fire AZ,  
473 Weinberg KI, Mindrinos M, Zehnder JL, Boyd SD, Xiao W, Davis RW, Miklos DB: High-throughput VDJ  
474 sequencing for quantification of minimal residual disease in chronic lymphocytic leukemia and  
475 immune reconstitution assessment. *Proceedings of the National Academy of Sciences of the United*  
476 *States of America* 2011, 108:21194-9.

477 [20] Martinez-Lopez J, Lahuerta JJ, Pepin F, Gonzalez M, Barrio S, Ayala R, Puig N, Montalban MA,  
478 Paiva B, Weng L, Jimenez C, Sopena M, Moorhead M, Cedena T, Rapado I, Mateos MV, Rosinol L,  
479 Oriol A, Blanchard MJ, Martinez R, Blade J, San Miguel J, Faham M, Garcia-Sanz R: Prognostic value  
480 of deep sequencing method for minimal residual disease detection in multiple myeloma. *Blood*  
481 2014, 123:3073-9.

482 [21] Wu D, Sherwood A, Fromm JR, Winter SS, Dunsmore KP, Loh ML, Greisman HA, Sabath DE,  
483 Wood BL, Robins H: High-throughput sequencing detects minimal residual disease in acute T  
484 lymphoblastic leukemia. *Science translational medicine* 2012, 4:134ra63.

485 [22] Mattocks CJ, Morris MA, Matthijs G, Swinnen E, Corveleyn A, Dequeker E, Muller CR, Pratt V,  
486 Wallace A: A standardized framework for the validation and verification of clinical molecular genetic  
487 tests. *European journal of human genetics : EJHG* 2010, 18:1276-88.

488 [23] Gullapalli RR, Desai KV, Santana-Santos L, Kant JA, Becich MJ: Next generation sequencing in  
489 clinical medicine: Challenges and lessons for pathology and biomedical informatics. *Journal of*  
490 *pathology informatics* 2012, 3:40.

491 [24] Williams ES, Hegde M: Implementing genomic medicine in pathology. *Advances in anatomic*  
492 *pathology* 2013, 20:238-44.

493 [25] Xuan J, Yu Y, Qing T, Guo L, Shi L: Next-generation sequencing in the clinic: promises and  
494 challenges. *Cancer letters* 2013, 340:284-95.

495 [26] Nekrutenko A, Taylor J: Next-generation sequencing data interpretation: enhancing  
496 reproducibility and accessibility. *Nature reviews Genetics* 2012, 13:667-72.

497 [27] van Dongen JJ, van der Velden VH, Bruggemann M, Orfao A: Minimal residual disease (MRD)  
498 diagnostics in acute lymphoblastic leukemia (ALL): need for sensitive, fast and standardized  
499 technologies. *Blood* 2015.

500 [28] Carlson CS, Emerson RO, Sherwood AM, Desmarais C, Chung MW, Parsons JM, Steen MS,  
501 LaMadrid-Herrmannsfeldt MA, Williamson DW, Livingston RJ, Wu D, Wood BL, Rieder MJ, Robins H:

502 Using synthetic templates to design an unbiased multiplex PCR assay. *Nature communications* 2013,  
503 4:2680.

504 [29] Collins FS, Hamburg MA: First FDA authorization for next-generation sequencer. *The New*  
505 *England journal of medicine* 2013, 369:2369-71.

506 [30] Kircher M, Sawyer S, Meyer M: Double indexing overcomes inaccuracies in multiplex sequencing  
507 on the Illumina platform. *Nucleic acids research* 2012, 40:e3.

508 [31] Quail MA, Smith M, Jackson D, Leonard S, Skelly T, Swerdlow HP, Gu Y, Ellis P: SASI-Seq: sample  
509 assurance Spike-Ins, and highly differentiating 384 barcoding for Illumina sequencing. *BMC genomics*  
510 2014, 15:110.

511 [32] Verhagen OJ, Wijkhuijs AJ, van der Sluijs-Gelling AJ, Szczepanski T, van der Linden-Schrever BE,  
512 Pongers-Willems MJ, van Wering ER, van Dongen JJ, van der Schoot CE: Suitable DNA isolation  
513 method for the detection of minimal residual disease by PCR techniques. *Leukemia* 1999, 13:1298-9.

514 [33] Kozarewa I, Turner DJ: 96-plex molecular barcoding for the Illumina Genome Analyzer. *Methods*  
515 *Mol Biol* 2011, 733:279-98.

516 [34] Bhat S, Herrmann J, Armishaw P, Corbisier P, Emslie KR: Single molecule detection in nanofluidic  
517 digital array enables accurate measurement of DNA copy number. *Analytical and bioanalytical*  
518 *chemistry* 2009, 394:457-67.

519 [35] Mir K, Neuhaus K, Bossert M, Schober S: Short barcodes for next generation sequencing. *PloS*  
520 *one* 2013, 8:e82933.

521 [36] Bystrykh LV: Generalized DNA barcode design based on Hamming codes. *PloS one* 2012,  
522 7:e36852.

523 [37] Masella AP, Bartram AK, Truszkowski JM, Brown DG, Neufeld JD: PANDAseq: paired-end  
524 assembler for illumina sequences. *BMC bioinformatics* 2012, 13:31.

525 [38] Giraud M, Salson M, Duez M, Villenet C, Quief S, Caillaud A, Gardel N, Roumier C, Preudhomme  
526 C, Figeac M: Fast multiclonal clusterization of V(D)J recombinations from high-throughput  
527 sequencing. *BMC genomics* 2014, 15:409.

528 [39] Kennedy SR, Schmitt MW, Fox EJ: Detecting ultralow-frequency mutations by Duplex  
529 Sequencing. 2014, 9:2586-606.

530 [40] Seitz V, Schaper S, Droge A, Lenze D, Hummel M, Hennig S: A new method to prevent carry-over  
531 contaminations in two-step PCR NGS library preparations. Nucleic acids research 2015.

532 [41] Brisco MJ, Sykes PJ, Hughes E, Dolman G, Neoh SH, Peng LM, Toogood I, Morley AA: Monitoring  
533 minimal residual disease in peripheral blood in B-lineage acute lymphoblastic leukaemia. British  
534 journal of haematology 1997, 99:314-9.

535 [42] Brisco MJ, Sykes PJ, Hughes E, Story CJ, Rice MS, Schwarzer AP, Morley AA: Molecular relapse can  
536 be detected in blood in a sensitive and timely fashion in B-lineage acute lymphoblastic leukemia.  
537 Leukemia 2001, 15:1801-2.

538 [43] Coustan-Smith E, Sancho J, Hancock ML, Razzouk BI, Ribeiro RC, Rivera GK, Rubnitz JE, Sandlund  
539 JT, Pui CH, Campana D: Use of peripheral blood instead of bone marrow to monitor residual disease  
540 in children with acute lymphoblastic leukemia. Blood 2002, 100:2399-402.

541 [44] Helgestad J, Rosthoj S, Johansen P, Varming K, Ostergaard E: Bone marrow aspiration technique  
542 may have an impact on therapy stratification in children with acute lymphoblastic leukaemia.  
543 Pediatric blood & cancer 2011, 57:224-6.

544 [45] Martin H, Atta J, Bruecher J, Elsner S, Schardt C, Stadler M, von Melchner H, Hoelzer D: In  
545 patients with BCR-ABL-positive ALL in CR peripheral blood contains less residual disease than bone  
546 marrow: implications for autologous BMT. Annals of hematology 1994, 68:85-7.

547 [46] van der Velden VH, Jacobs DC, Wijkhuijs AJ, Comans-Bitter WM, Willemse MJ, Hahlen K, Kamps  
548 WA, van Wering ER, van Dongen JJ: Minimal residual disease levels in bone marrow and peripheral  
549 blood are comparable in children with T cell acute lymphoblastic leukemia (ALL), but not in  
550 precursor-B-ALL. Leukemia 2002, 16:1432-6.

551

552

553 **Tables**

554 Table 1. Oligonucleotide designs for dual-indexing. These are used in the second stage of a nested  
 555 PCR. First stage primers have adaptors<sup>†</sup> at each end of amplicon which are complementary to the  
 556 indexing oligonucleotide.

Index	Oligonucleotide Sequence	Index in oligonucleotide	Index read sequence
i7_01	5'-CAAGCAGAAGACGGCATAACGAGATTCTAGCTAGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	TCTAGCTA	TAGCTAGA
i7_02	5'-CAAGCAGAAGACGGCATAACGAGATCTAGCTATGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	CTAGCTAT	ATAGCTAG
i7_03	5'-CAAGCAGAAGACGGCATAACGAGATAGGTTGGCGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	AGGTTGGC	GCCAACT
i7_04	5'-CAAGCAGAAGACGGCATAACGAGATGACCAACGGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	GACCAACG	CGTTGGTC
i7_05	5'-CAAGCAGAAGACGGCATAACGAGATGCGGAGTTGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	GCGGAGTT	AACTCCGC
i7_06	5'-CAAGCAGAAGACGGCATAACGAGATGTGCCATAGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	GTGCCATA	TATGGCAC
i7_07	5'-CAAGCAGAAGACGGCATAACGAGATTAATGTCCGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	TAATGTCC	GGACATTA
i7_08	5'-CAAGCAGAAGACGGCATAACGAGATCGAAGGACGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	CGAAGGAC	GTCTTCG
i7_09	5'-CAAGCAGAAGACGGCATAACGAGATAATGTCTGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	AATGTCT	AGGACATT
i7_10	5'-CAAGCAGAAGACGGCATAACGAGATAGAACATTGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	AGAACATT	AATGTTCT
i7_11	5'-CAAGCAGAAGACGGCATAACGAGATTGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	TGTCAGTC	GAAGTACA
i7_12	5'-CAAGCAGAAGACGGCATAACGAGATCACCGTTGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	CACCGCTT	AAGCGGTG
i7_13	5'-CAAGCAGAAGACGGCATAACGAGATCAGACGAGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	CAGACGCA	TGCGTCTG
i7_14	5'-CAAGCAGAAGACGGCATAACGAGATGCTACTAGGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	GCTACTAG	CTAGTAGC
i7_15	5'-CAAGCAGAAGACGGCATAACGAGATGTCAGTCTGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	GTCAGTCT	AGACTGAC
i7_16	5'-CAAGCAGAAGACGGCATAACGAGATTTACCCGCGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	TTCACCGC	GCGGTGAA
i7_17	5'-CAAGCAGAAGACGGCATAACGAGATGGTCTAATGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	GGTCTAAT	ATTAGACC
i7_18	5'-CAAGCAGAAGACGGCATAACGAGATACCTGGATGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	ACCTGGAT	ATCCAGGT
i7_19	5'-CAAGCAGAAGACGGCATAACGAGATAGCGACAGGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	AGCGACAG	CTGTGCT
i7_20	5'-CAAGCAGAAGACGGCATAACGAGATATAGGCTCGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	ATAGGCTC	GAGCCTAT
i7_21	5'-CAAGCAGAAGACGGCATAACGAGATTAGAACATGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	TAGAACAT	ATGTTCTA
i7_22	5'-CAAGCAGAAGACGGCATAACGAGATTGCGGAGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	TGCGGAGT	ACTCCGCA
i7_23	5'-CAAGCAGAAGACGGCATAACGAGATTGCGGAGGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	TTGCGGAG	CTCCGCAA
i7_24	5'-CAAGCAGAAGACGGCATAACGAGATTAGAACAGTGACTGGAGTTTCAGACGTGTGCTCTTCCGATC*T-3'	TTAGAACA	TGTTCTAA
i5_01	5'-AATGATACGGCGACCACCGAGATCTACACCACTTGAGACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	CACTTGAG	CACTTGAG
i5_02	5'-AATGATACGGCGACCACCGAGATCTACACGTTACCGAACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	GTTACCGA	GTTACCGA
i5_03	5'-AATGATACGGCGACCACCGAGATCTACACTGACGACTACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	TGACGACT	TGACGACT
i5_04	5'-AATGATACGGCGACCACCGAGATCTACACACGATTACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	ACGGATTC	ACGGATTC
i5_05	5'-AATGATACGGCGACCACCGAGATCTACACCCATAGGAACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	CCATAGGA	CCATAGGA
i5_06	5'-AATGATACGGCGACCACCGAGATCTACACTGGAAGGCACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	TGGAAGGC	TGGAAGGC
i5_07	5'-AATGATACGGCGACCACCGAGATCTACACGCATCATGACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	GCATCATG	GCATCATG
i5_08	5'-AATGATACGGCGACCACCGAGATCTACACAGCGTGAACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	AGCGGTGA	AGCGGTGA
i5_09	5'-AATGATACGGCGACCACCGAGATCTACACAGTTACCGACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	AGTTACCG	AGTTACCG
i5_10	5'-AATGATACGGCGACCACCGAGATCTACACCATGCATAACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	CATGCATA	CATGCATA
i5_11	5'-AATGATACGGCGACCACCGAGATCTACACACATGCATAACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	ACATGCAT	ACATGCAT
i5_12	5'-AATGATACGGCGACCACCGAGATCTACACCATAGGACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	ACCATAGG	ACCATAGG
i5_13	5'-AATGATACGGCGACCACCGAGATCTACACTCCAGGTAACACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	TCCAGGTA	TCCAGGTA
i5_14	5'-AATGATACGGCGACCACCGAGATCTACACCTTAATTGACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	CTTAATTG	CTTAATTG
i5_15	5'-AATGATACGGCGACCACCGAGATCTACACCGATTCAACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	CGGATTCA	CGGATTCA
i5_16	5'-AATGATACGGCGACCACCGAGATCTACACTTAGACCAACTCTTTCCCTACACGACGCTCTTCCGATC*T-3'	TTAGACCA	TTAGACCA

557 <sup>†</sup>First stage primers including partial adaptor sequences for - Forward (i7) primers:  
 558 5'-GTGACTGGAGTTTCAGACGTGTGCTCTTCCGATCTCTGGGTGCGACAGGCCCTGGACAA-3'  
 559 5'-GTGACTGGAGTTTCAGACGTGTGCTCTTCCGATCTTGGATCCGTCAGCCCCAGGGAAGG-3'  
 560 5'-GTGACTGGAGTTTCAGACGTGTGCTCTTCCGATCTGGTCCGCCAGGCTCCAGGAA-3'  
 561 5'-GTGACTGGAGTTTCAGACGTGTGCTCTTCCGATCTTGGATCCGCCAGGCCAGGGAAGG-3'  
 562 5'-GTGACTGGAGTTTCAGACGTGTGCTCTTCCGATCTGGGTGCGCCAGATGCCCGGGAAGG-3'

563 5'-GTGACTGGAGTTCAGACGTGTGCTCTCCGATCTTGGATCAGGCAGTCCCCATCGAGAG-3'  
564 5'-GTGACTGGAGTTCAGACGTGTGCTCTCCGATCTTTGGGTGCGACAGGCCCTGGACAA-3'  
565 Reverse primer (i5):  
566 5'-ACACTCTTCCCTACACGACGCTCTCCGATCTTACCTGAGGAGACGGTGACC-3'  
567 \*Indicates addition of phosphorothioated DNA bases



568 Table 2. Sequencing runs performed on the Illumina MiSeq

<b>MiSeq Run</b>	<b>Unique run name</b>	<b>Number of samples</b>	<b>Indexing strategy</b>	<b>MiSeq kit version</b>	<b>Read length</b>	<b>Total reads - million</b>	<b>Cluster density (K/mm<sup>2</sup>)</b>	<b>Reads passing filter - million</b>	<b>Phred quality (%<math>\geq</math>Q30)</b>
1	A478Y	36	Single, Kozarewa	2	300	22.4	1211	17.38	71.0
2	A5U9G	34	Single, Kozarewa	2	300	22.12	1094	18.69	74.4
3	A6PKD	20	Single, Kozarewa	2	300	20.91	1006	17.55	78.6
4	A6FMV	23	Single, Kozarewa	3	167	12.11	465	11.70	95.9
5	A7BK7	28	Single, Kozarewa	2	208	22.35	1063	18.59	79.2
6	A7ELC	31	Dual, TruGrade	3	151	32.83	1267	22.81	80.3
7	A7FDO	8	Dual, TruGrade	3	151	23.98	865	18.77	82.1
8	A72MP	8	Dual, TruGrade	2	300	23.01	1067	20.17	70.1
9	A7B79	21	Dual, TruGrade	2	300	24.45	1175	19.66	72.5
10	A8FPT	12	Dual, TruGrade	2	300	21.66	1020	19.53	78.6
11	ABALR	12	Dual, TruGrade	2	300	23.44	1120	18.33	72.6
12	ABAJL	37	Dual, TruGrade	2	300	24.19	1170	21.77	69.3
13	ABG7Y	6	Dual, TruGrade	3	151	24.29	892	22.48	85

569

570

571 Table 3. Assignment of reads from dual-indexed run A7ELC using the MiSeq on-board  
572 demultiplexer. Index combinations corresponding to actual samples are highlighted in bold.  
573 Misassigned reads (non-bold print) are shown for all possible dual-index combinations where at least  
574 one index is assigned (complete information in Supplementary table 5). Index 1 has the prefix i5 and  
575 index 2 has the prefix i7. An example of frameshift causing misassignment occurs between indices  
576 i7\_05 and i7\_22 (underlined text) resulting in misassignment of 0.6% of reads.

577

		i7 indices															
		i7_01	i7_02	i7_03	i7_04	<u>i7_05</u>	i7_06	i7_07	i7_08	i7_09	i7_14	i7_17	i7_19	i7_20	<u>i7_22</u>	i7_23	i7_24
i5 indices	i5_01	<b>532594</b>	<b>553987</b>	<b>606447</b>	2105	<b><u>715211</u></b>	<b>586926</b>	3581	<b>500618</b>	0	2	5	2	4	<b><u>2642</u></b>	140	0
	i5_02	<b>483605</b>	<b>420978</b>	3407	<b>563981</b>	<u>2130</u>	<b>694559</b>	<b>457959</b>	1129	153	4	7	3	0	<b><u>2</u></b>	0	3
	i5_03	4306	<b>413460</b>	<b>853442</b>	<b>912389</b>	<b><u>542539</u></b>	2107	<b>785862</b>	<b>433723</b>	234	0	6	0	16	<b><u>1851</u></b>	103	0
	i5_04	<b>869143</b>	2301	<b>820461</b>	<b>485673</b>	<b><u>708010</u></b>	<b>548037</b>	3374	<b>501739</b>	2	0	0	0	10	<b><u>2470</u></b>	162	0
	i5_05	5541	<b>398969</b>	4490	<b>794442</b>	<b><u>608924</u></b>	3922	<b>631108</b>	1442	170	0	3	3	0	<b><u>1949</u></b>	97	2
	i5_06	<b>1410284</b>	3156	<b>807859</b>	2536	<u>2938</u>	<b>731447</b>	<b>613956</b>	1927	189	0	0	0	14	<b><u>5</u></b>	0	0
	i5_09	163	131	0	186	<u>0</u>	223	123	0	0	0	0	0	0	<b><u>0</u></b>	0	0
	i5_10	0	2	0	0	<u>0</u>	0	0	2	0	0	0	0	0	<b><u>0</u></b>	0	0
	i5_12	1	68	0	145	<u>130</u>	0	127	1	0	0	0	0	0	<b><u>0</u></b>	0	0
	i5_14	0	0	1	0	<u>0</u>	0	0	0	0	0	0	0	0	<b><u>0</u></b>	0	0
	i5_15	803	1	835	443	<u>718</u>	556	6	543	0	0	0	0	0	<b><u>5</u></b>	0	0
	i5_16	4	6	8	2	<u>7</u>	2	9	2	0	0	0	0	0	<b><u>0</u></b>	0	0
	Undetermined reads 3E+06																

578 Table 4. Assignment of reads from dual-indexed run A7ELC demultiplexed allowing for no  
579 mismatches in index sequences. Index combinations for samples included in the run are shown in  
580 bold print. Misassigned reads (non-bold print) are shown for all possible dual-index combinations  
581 where at least one index is assigned (complete information in Supplementary table 7). Index 1 has the  
582 prefix i5 and index 2 has the prefix i7.  
583

		i7 indices									
		<i>i7_01</i>	<i>i7_02</i>	<i>i7_03</i>	<i>i7_04</i>	<i>i7_05</i>	<i>i7_06</i>	<i>i7_07</i>	<i>i7_08</i>	<i>i7_22</i>	<i>i7_23</i>
i5 indices	<i>i5_01</i>	<b>484180</b>	<b>505456</b>	<b>549076</b>	402	<b>654781</b>	<b>535588</b>	661	<b>449768</b>	1689	12
	<i>i5_02</i>	<b>436266</b>	<b>387331</b>	741	<b>511368</b>	440	<b>629466</b>	<b>415014</b>	241	0	0
	<i>i5_03</i>	1138	<b>376242</b>	<b>790058</b>	<b>893669</b>	<b>489608</b>	472	<b>714305</b>	<b>394937</b>	1177	10
	<i>i5_04</i>	<b>809651</b>	454	<b>738566</b>	<b>439508</b>	<b>690965</b>	<b>490128</b>	683	<b>450862</b>	1715	2
	<i>i5_05</i>	1218	<b>358552</b>	970	<b>722159</b>	<b>555629</b>	637	<b>569741</b>	250	1239	10
	<i>i5_06</i>	<b>1291695</b>	535	<b>733096</b>	525	614	<b>660168</b>	<b>552237</b>	360	2	0
	<i>i5_09</i>	0	0	0	0	0	1	0	0	0	0
	<i>i5_12</i>	0	1	0	1	1	0	1	0	0	0
	<i>i5_15</i>	561	0	583	298	522	375	0	365	4	0
	<i>i5_16</i>	0	0	0	0	0	0	0	0	0	0
	Undetermined reads 5E+06										

584 Table 5. Reduction of misassignment due to index read sequence error using the Error Aware  
585 Demultiplexing strategy. Numbers of reads assigned, with percentage change in reads assigned  
586 compared to standard MiSeq on-board demultiplexing, allowing up to 1 base pair mismatch (in  
587 brackets); Reads assigned using standard methods allowing for 0 base pair mismatches in index reads  
588 are on upper rows (0 bp) of each index combination and Error Aware Demultiplexer read assignments  
589 are on the lower rows (EAD) of each index combination. The bold text indicate sample index  
590 combinations that were present in sequenced samples, the non-bold are index combinations not  
591 included in the sequencing run. Data from sequencing run A7FDO. Index 1 has the prefix i5 and  
592 index 2 has the prefix i7.

		i7 indices								
		<i>i7_01</i>	<i>i7_03</i>	<i>i7_04</i>	<i>i7_05</i>	<i>i7_06</i>	<i>i7_07</i>	<i>i7_08</i>	<i>i7_10</i>	
i5 indices	<i>i5_01</i>	0 bp	<b>1633722</b> <b>(-8.6)</b>	44 (-76.2)	76 (-43.3)	132 (-68.3)	109 (-72.5)	118 (-66.9)	86 (-74.9)	108 (-73.3)
		EAD	<b>1761443</b> <b>(-1.4)</b>	89 (-51.9)	63 (-53.0)	102 (-75.5)	140 (-64.6)	110 (-69.2)	93 (-72.9)	85 (-79.0)
	<i>i5_02</i>	0 bp	154 (-50.2)	<b>1714456</b> <b>(-8.8)</b>	121 (-46.5)	152 (-72.0)	169 (-50.4)	165 (-76.9)	129 (-81.8)	156 (-54.5)
		EAD	124 (-59.9)	<b>1921136</b> <b>(-2.2)</b>	87 (-61.5)	127 (-76.6)	177 (-48.1)	174 (-75.6)	129 (-81.8)	181 (-47.2)
	<i>i5_03</i>	0 bp	50 (-60.0)	79 (-42.8)	<b>1506573</b> <b>(-8.5)</b>	92 (-85.8)	59 (-67.6)	68 (-89.4)	61 (-69.3)	65 (-85.2)
		EAD	37 (-70.4)	73 (-47.1)	<b>1620923</b> <b>(-1.5)</b>	63 (-90.3)	36 (-80.2)	55 (-91.4)	62 (-68.8)	55 (-87.5)
	<i>i5_04</i>	0 bp	86 (-54.0)	84 (-45.5)	60 (-37.5)	<b>1667394</b> <b>(-9.0)</b>	75 (-58.3)	61 (-72.9)	80 (-80.5)	90 (-88.9)
		EAD	71 (-62.0)	76 (-50.6)	52 (-45.8)	<b>1800250</b> <b>(-1.7)</b>	64 (-64.4)	62 (-72.4)	73 (-82.2)	101 (-87.6)
	<i>i5_05</i>	0 bp	185 (-82.8)	171 (-60.7)	150 (-66.7)	254 (-87.3)	<b>1748758</b> <b>(-9.0)</b>	158 (-82.3)	118 (-78.9)	215 (-88.6)
		EAD	194 (-81.9)	189 (-56.6)	115 (-74.4)	168 (-91.6)	<b>1909784</b> <b>(-0.6)</b>	169 (-81.1)	110 (-80.3)	180 (-90.5)
	<i>i5_06</i>	0 bp	185 (-64.4)	193 (-73.9)	154 (-79.8)	185 (-80)	170 (-57.5)	<b>1924357</b> <b>(-8.2)</b>	160 (-60.7)	161 (-55.6)
		EAD	159 (-69.4)	160 (-78.3)	93 (-87.8)	173 (-81.3)	179 (-55.3)	<b>2083579</b> <b>(-0.6)</b>	127 (-68.8)	164 (-54.8)
	<i>i5_07</i>	0 bp	152 (-70.9)	222 (-79.7)	130 (-60.2)	198 (-76.1)	167 (-70.1)	137 (-67.0)	<b>1586001</b> <b>(-9.1)</b>	172 (-69.4)
		EAD	138 (-73.6)	206 (-81.1)	147 (-55.0)	174 (-79.0)	157 (-71.9)	134 (-67.7)	<b>1765359</b> <b>(-1.1)</b>	152 (-73.0)
	<i>i5_08</i>	0 bp	44 (-86.3)	57 (-87.8)	40 (-79.9)	98 (-93.4)	60 (-92.2)	55 (-85.7)	31 (-88.0)	<b>1781666</b> <b>(-8.8)</b>
		EAD	32 (-90.0)	58 (-87.6)	38 (-80.9)	74 (-95.0)	46 (-94.1)	47 (-87.8)	30 (-88.4)	<b>1979346</b> <b>(-1.3)</b>

593 **Figure Legends**

594 **Figure 1. Nested PCR and final library product for *IGH* VDJ amplicon sequencing on the**  
595 **Illumina MiSeq.** (A) Family variable heavy chain segment (VH) and joining heavy chain segment  
596 (JH) primers with complementary partial index sequences are used to amplify the VDJ junction of the  
597 rearranged *IGH* gene, containing the hypervariable region (purple). (B) In the second stage, index  
598 sequences (blue) and Illumina P5 and P7 platform adaptors sequences are added. The final amplicon  
599 construct with sequencing strategy is shown in (C). One or two index reads were used depending on  
600 indexing strategy, with single-end sequencing read from the JH (P5) end of the amplicon.

601 **Figure 2. Misassignment of sequences from sequencing run A7BK7.** Indexed reads corresponding  
602 to samples are shown in red. Reads assigned to indices not included in the sequencing run are shown  
603 in blue. Eight bp indices used as per *Kowenza et al.*<sup>33</sup>

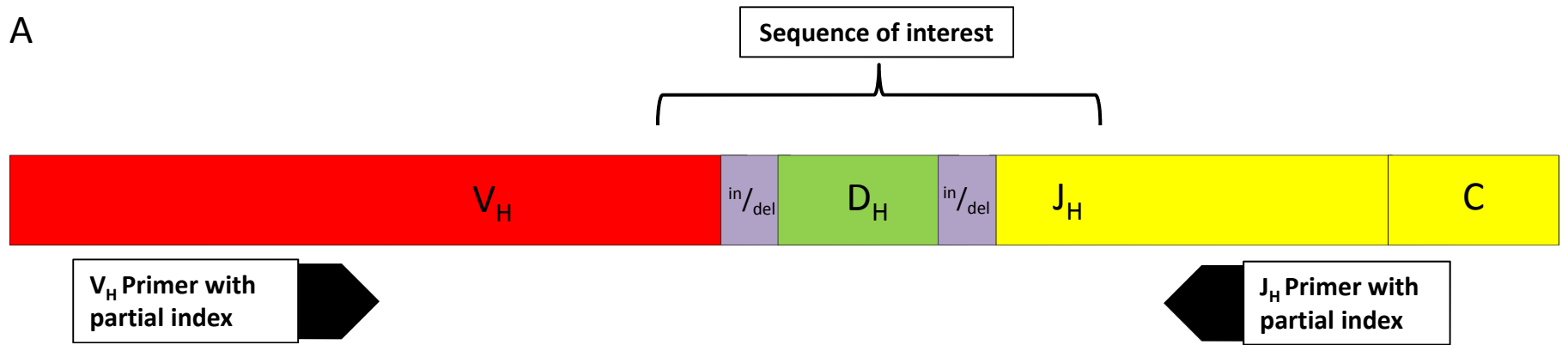
604 **Figure 3. Comparison of quality statistics between reads assigned to real samples (“true”) and**  
605 **reads “misassigned” to indices not included in the run (examples using data from run A7ELC).**  
606 (A-C) Base quality of sequencing reads shown by position (A) mean quality score distribution for all  
607 sequences (B) and index read quality (C) from a representative “true” indexed sample. (D-F) Base  
608 quality of sequencing reads shown by position (D) mean quality score distribution for all sequences  
609 (E) and index read quality (F) from a representative “misassigned” sample. Pooled libraries were  
610 sequenced (150 bases, single-end) on an Illumina MiSeq.

611 **Figure 4. Accurate quantification of MRD down to 1 in 1 million cells using a spike in control**  
612 **quantified by digital PCR.** Three separate serial dilutions of a cell lines; (A) SUPB15, (B) TOM-1  
613 and (C) REH, spiked into to 1 million cells equivalent of pooled normal lymphocyte DNA. Samples  
614 sequenced on Illumina MiSeq (> 2 million reads per sample), indexed using HTS strategy described  
615 above and demultiplexed using Error Aware Demultiplexer.

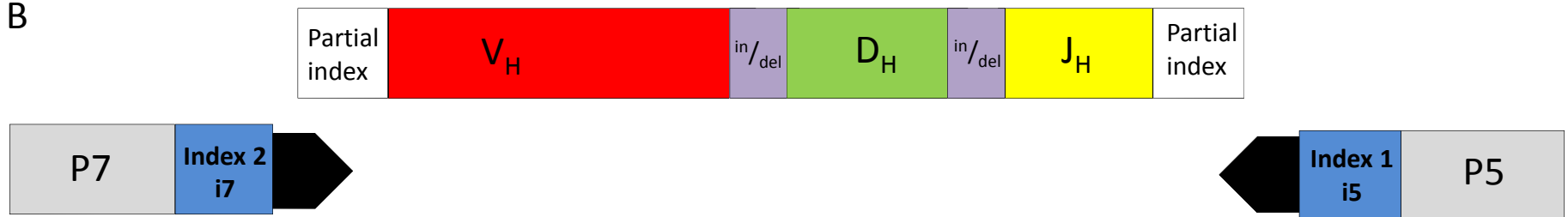
616 **Figure 5. MRD analysis performed in 5 patients at the end of induction chemotherapy for**  
617 **childhood B-ALL (day 28 MRD).** The current gold standard RQ-PCR technique is compared to a  
618 standard HTS work flow (A) and an optimised workflow using Error Aware Demultiplexer (B). The

619 sequencing run included a diluted ( $10^{-1}$ ) diagnostic sample for each patient to simulate an increased  
620 probability of misassignment. MRD samples were also “spiked” with cell line DNA for quantification  
621 purposes. Errors introduced during standard multiplexed sequencing result in mis-diagnosis in  
622 patients N001 and N005 due to misassignment (A). These incorrect calls due to misassignment are  
623 removed (B) using the workflow presented.

A



B



C

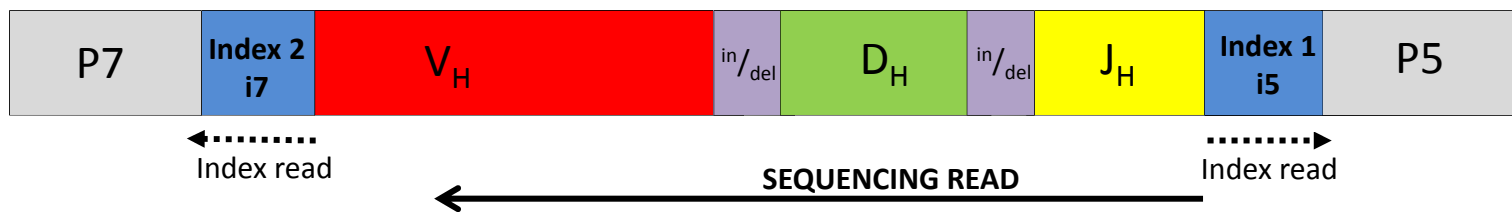
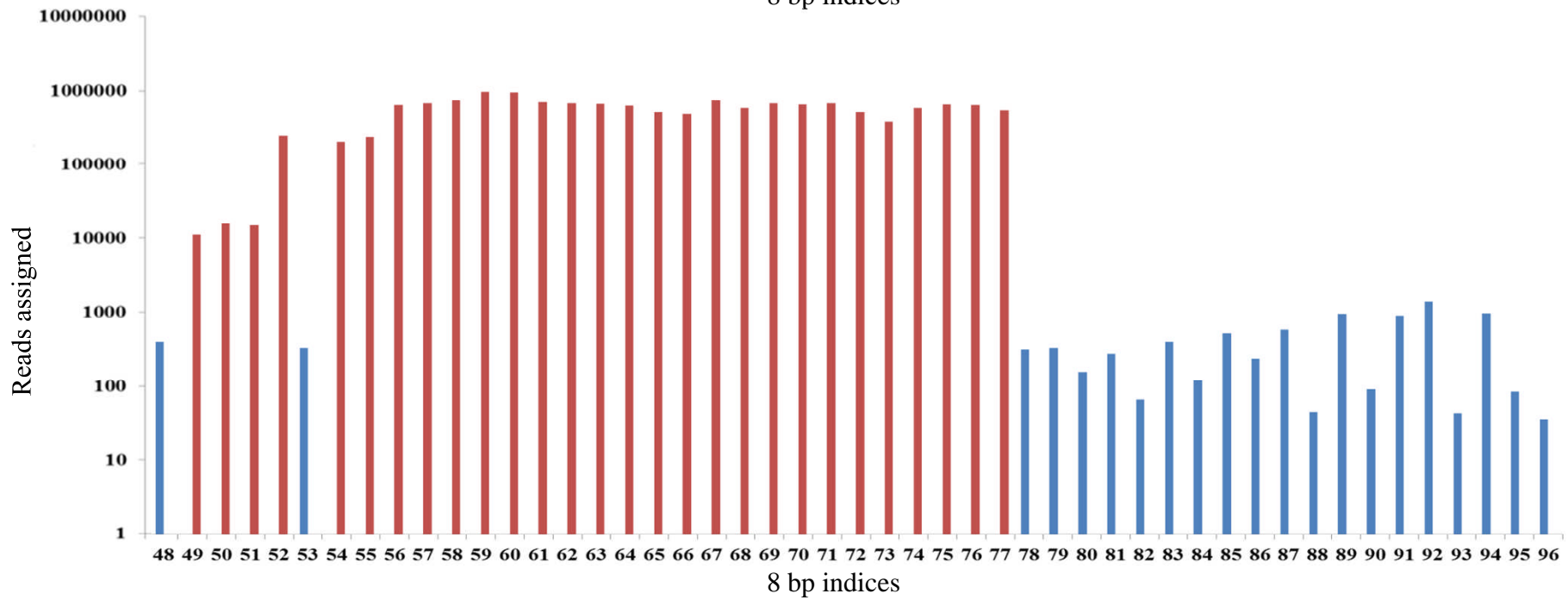
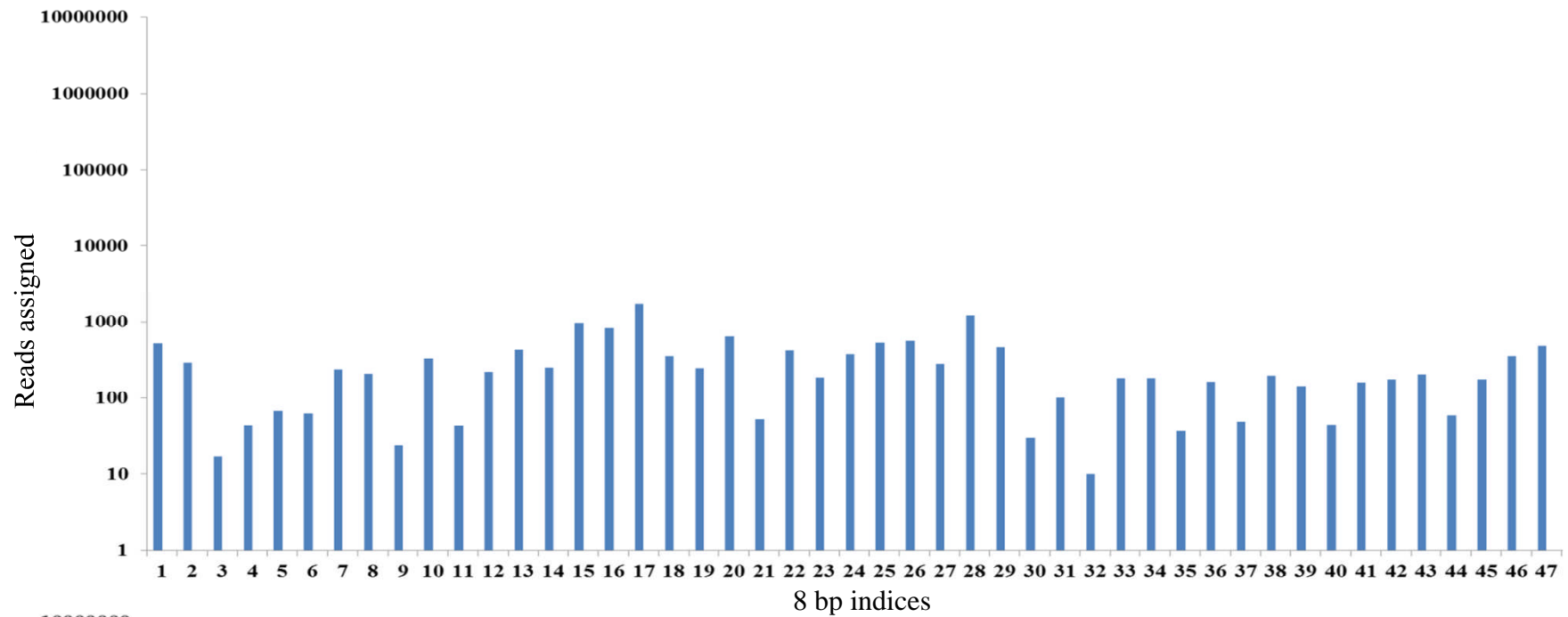


Figure 1.





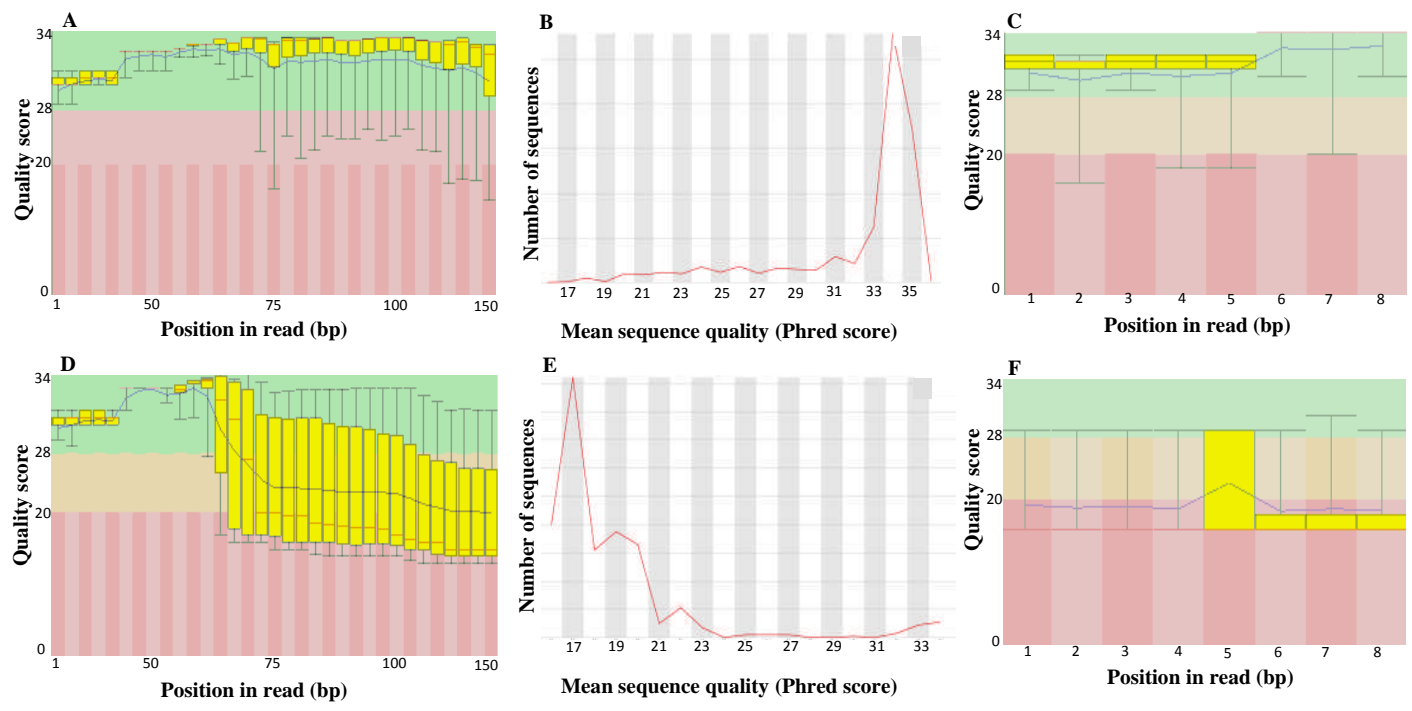


Figure 3.

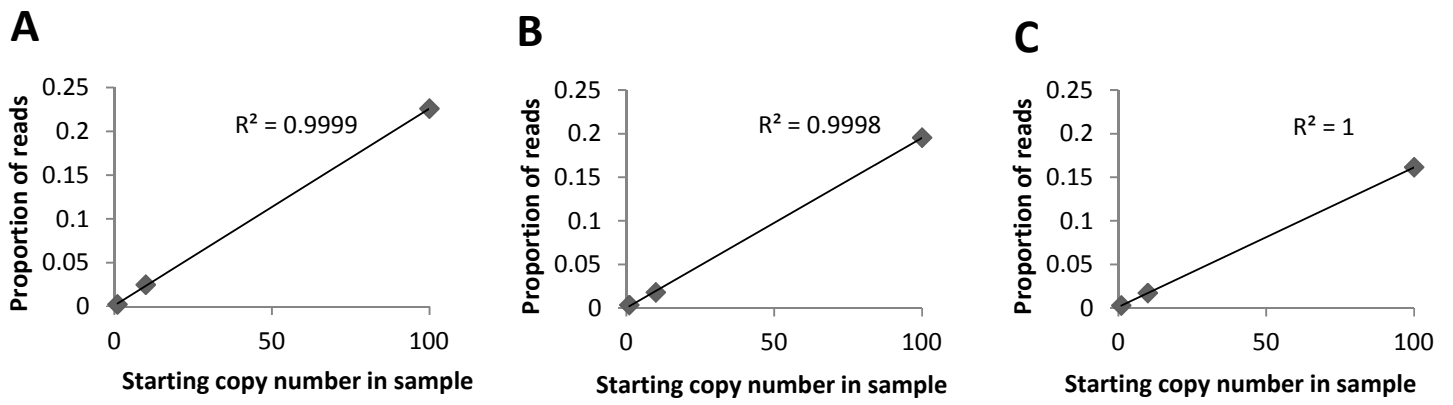


Figure 4.

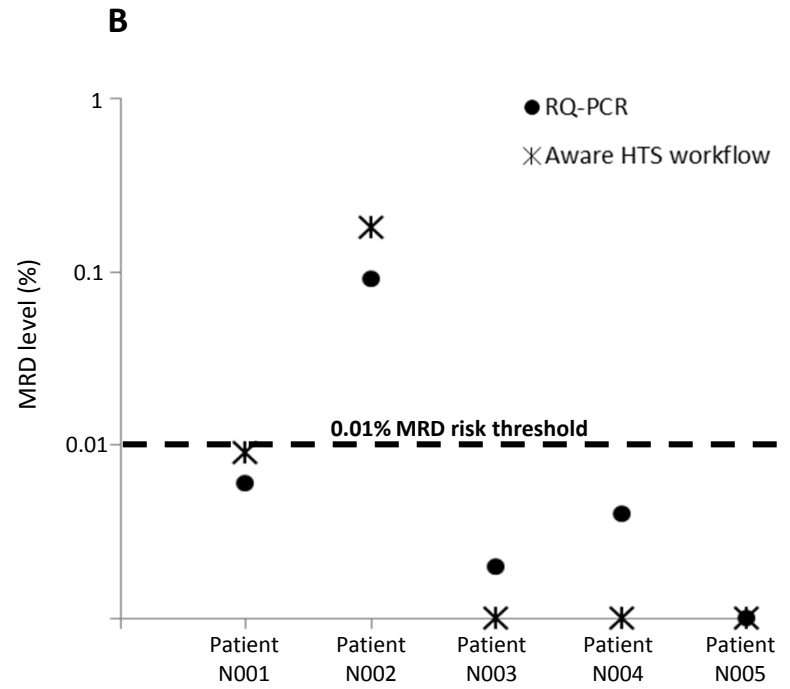
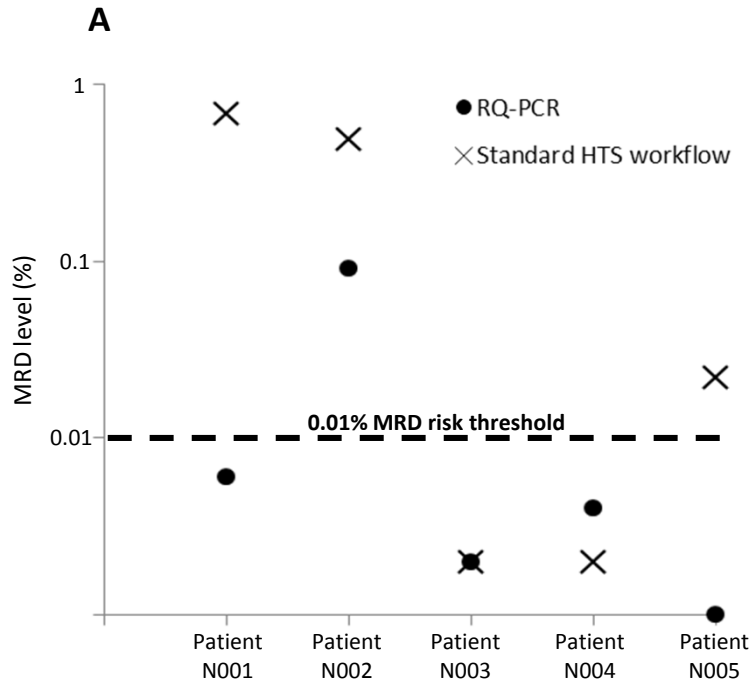


Figure 5.