



# NGmerge: merging paired-end reads via novel empirically-derived models of sequencing errors

## Citation

Gaspar, John M. "NGmerge: Merging Paired-end Reads via Novel Empirically-derived Models of Sequencing Errors." *Bmc Bioinformatics* 19, no. 1 (2018): 536.

## Published version

<https://doi.org/10.1186/s12859-018-2579-2>

## Link

<http://nrs.harvard.edu/urn-3:HUL.InstRepos:42661733>

## Terms of use

This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Other Posted Material (LAA), as set forth at

<https://harvardwiki.atlassian.net/wiki/external/NGY5NDE4ZjgzNTc5NDQzMGIzZWZhMGFIOWI2M2EwYTg>

## Accessibility

<https://accessibility.huit.harvard.edu/digital-accessibility-policy>

## Share Your Story

The Harvard community has made this article openly available.  
Please share how this access benefits you. [Submit a story](#)

1 **NGmerge: merging paired-end reads via novel empirically-derived**  
2 **models of sequencing errors**

3

4 John M. Gaspar\*

5

6 Informatics Group, Division of Science, Harvard University Faculty of Arts and  
7 Sciences, Cambridge MA, 02138

8

9 \*Corresponding author

10

11 Email address:

12 JMG: [jsh58@wildcats.unh.edu](mailto:jsh58@wildcats.unh.edu)

13

14

15

16

17

18

19

20 **Abstract**

21 **Background**

22 Advances in Illumina DNA sequencing technology have produced longer  
23 paired-end reads that increasingly have sequence overlaps. These reads can be merged  
24 into a single read that spans the full length of the original DNA fragment, allowing for  
25 error correction and accurate determination of read coverage. Extant merging  
26 programs utilize simplistic or unverified models for the selection of bases and quality  
27 scores for the overlapping region of merged reads.

28

29 **Results**

30 We first examined the baseline quality score - error rate relationship using  
31 sequence reads derived from PhiX. In contrast to numerous published reports, we  
32 found that the quality scores produced by Illumina were not substantially inflated  
33 above the theoretical values, once the reference genome was corrected for unreported  
34 sequence variants. The PhiX reads were then used to create empirical models of  
35 sequencing errors in overlapping regions of paired-end reads, and these models were  
36 incorporated into a novel merging program, NGmerge. We demonstrate that NGmerge  
37 corrects errors and ambiguous bases better than other merging programs, and that it  
38 assigns quality scores for merged bases that accurately reflect the error rates. Our

39 results also show that, contrary to published analyses, the sequencing errors of paired-  
40 end reads are not independent.

41

## 42 **Conclusions**

43 We provide a free and open-source program, NGmerge, that performs better  
44 than existing read merging programs. NGmerge is available on GitHub  
45 (<https://github.com/harvardinformatics/NGmerge>) under the MIT License; it is written  
46 in C and supported on Linux.

47

## 48 **Keywords**

49 high-throughput sequencing; Illumina paired-end sequencing; read merging;  
50 sequencing errors; quality scores; PhiX

51

## 52 **Background**

53 Among the high-throughput DNA sequencing technologies, the Solexa/Illumina  
54 platform [1] produces the greatest quantity of sequence data in a single run [2]. One  
55 unique attribute of this technology is its ability to generate sequence reads from both  
56 ends of a given DNA molecule. This provides many opportunities for biological

57 interpretation; for example, one can infer the full extent of a DNA molecule without  
58 sequencing its entirety, by aligning the paired-end reads to a reference sequence.

59

60 The output from an Illumina sequencing run is a set of FASTQ files, which  
61 contain read sequences and corresponding quality scores [3]. As first developed for  
62 Sanger sequencing, a quality score is determined from the probability that a given  
63 sequenced base is wrong, via the following equation [4]:

$$64 \quad Q_{\text{base}} = -10 \times \log_{10}( P(\text{base wrong}) ) \quad (1)$$

65 Thus, a base with a quality score of 40 should have a 1 in 10,000 chance of being wrong.

66 In sequence variant detection, many programs will consider only bases that achieve a  
67 minimum quality score, in order to reduce false positives [5], and in clinical variant  
68 detection, such as cancer diagnostics, published guidelines frequently incorporate such  
69 standards [6, 7]. However, numerous studies have shown that the raw quality scores  
70 produced by Illumina machines are inflated; that is, sequenced bases with a given  
71 quality score have higher error rates than expected from equation (1), especially at the  
72 high end of the scale [8-11].

73

74 Advances in Illumina sequencing technology have given rise to reads of  
75 increasing length, such that the paired-end reads for a particular library may have  
76 substantial sequence overlaps. Since these overlapping regions do not represent

77 independent sequence data, it is possible to merge the reads into a single read spanning  
78 the full length of the original DNA fragment. This merging process allows for error  
79 correction and accurate determination of read coverage, and it has become increasingly  
80 appreciated in applications ranging from targeted variant resequencing [12] to  
81 metabarcoding (e.g., 16S/18S rRNA studies) [13].

82

83 An early merging program was fastq-join [14], which was the default option in  
84 the microbial ecology analysis package QIIME [15]. The latest version of that package,  
85 QIIME 2, uses the open-source VSEARCH [16]. A third merging program is PEAR [17],  
86 which has the significant advantage over the other two programs (and other programs  
87 such as FLASH [18] and CASPER [19]) of considering “dovetailed” alignments, in  
88 which one read’s 3’ end extends past its pair’s 5’ end (see Fig. 1B). The sequencing of  
89 DNA fragments that are shorter than the read lengths will result in reads that contain  
90 portions of sequencing adapters on their 3’ ends; such read pairs will not be merged by  
91 programs that fail to consider dovetailed alignments [17].

92

93 Another important difference among read merging programs is the method for  
94 assigning quality scores to the merged bases. Fastq-join and FLASH use a simplistic  
95 scheme in which, if the bases of the two reads match, the higher quality score is used for  
96 the merged base. Where the bases disagree, the base with the higher quality score is

97 selected, and the difference in quality scores becomes the merged base's score. PEAR  
 98 was designed with the reasoning that, if the bases of the original R1 and R2 reads agree,  
 99 then the quality score of the merged base should reflect the probability that both  
 100 original bases were wrong:

101  $Q_{\text{merged base}} = -10 \times \log_{10}( P(\text{R1 base wrong AND R2 base wrong}) )$  (2)

102  $= -10 \times \log_{10}( P(\text{R1 base wrong}) \times P(\text{R2 base wrong}) )$  (3)

103  $= -10 \times \log_{10}( P(\text{R1 base wrong}) ) + -10 \times \log_{10}( P(\text{R2 base wrong}) )$  (4)

104  $= Q_{R1} + Q_{R2}$  (5)

105 Thus, PEAR sums quality scores for matching bases. VSEARCH utilizes a more  
 106 sophisticated model developed by Edgar and Flyvbjerg [20], but the resulting scheme  
 107 for matching bases is nearly identical to that of PEAR (see Fig. S1). For example, a  
 108 merged base created from matching bases with quality scores of 40 would be assigned a  
 109 quality score of 40 by fastq-join, 80 by PEAR, and 85 by VSEARCH (ignoring the  
 110 artificial caps on quality scores placed by the programs).

111  
 112 None of these quality score profiles has been tested empirically, despite possible  
 113 shortcomings. For example, the profiles of PEAR and VSEARCH are based on equation  
 114 (1), which, as noted above, has been demonstrated to be inaccurate with Illumina  
 115 sequencing. Furthermore, both PEAR and VSEARCH were designed under the  
 116 assumption that sequencing errors in the two reads are independent; that is, in the

117 analysis above, equation (3) follows from equation (2) only if the two events (“R1 base  
118 wrong” and “R2 base wrong”) are independent. This assumption has, to our  
119 knowledge, never been verified.

120

121 In this manuscript, we first evaluate the baseline quality score - error rate  
122 relationship produced by Illumina machines using reads derived from the  
123 enterobacteria phage  $\Phi$ X174 (“PhiX”). This virus’ genome was the first DNA genome to  
124 be sequenced [21], and a library composed of fragments of PhiX DNA is routinely  
125 added to Illumina sequencing runs as a control. In addition to the wide accessibility of  
126 PhiX reads, the sizes of the fragments in this library are such that most PhiX-derived  
127 read pairs produced by longer Illumina runs will have sufficient overlaps that can be  
128 used to create quality score profiles for merged bases (Fig. S2). We have done this and  
129 incorporated these profiles into a novel merging program, NGmerge (Fig. 1). We  
130 demonstrate that NGmerge corrects errors and ambiguous bases (Ns) better than other  
131 merging programs, and produces merged reads whose quality scores accurately reflect  
132 the bases’ error rates.

133

## 134 **Results**

### 135 **Baseline error rates**



136 We began with nine Illumina sequencing runs that yielded 2×250bp paired-end  
137 reads, produced at Harvard University. After identifying reads that originated from  
138 PhiX, we calculated error rates for each quality score. Consistent with previous studies  
139 [8-11], the error rates were higher than expected based on equation (1) above; bases  
140 with quality score 40 had error rates an order of magnitude above the predicted  $1 \times 10^{-4}$ .

141

142 However, a closer look at the alignments revealed variants from the canonical  
143 PhiX reference genome. In all of the sequencing runs, the same five sequence variants  
144 were identified at a minimum 95% allele frequency, with most at greater than 99%  
145 (Table 1). No other variants were identified.

146

147 We modified the reference genome to incorporate the five observed variants.  
148 Once this was done, the error spectrum more closely matched the expected relationship  
149 (Fig. 2). For example, bases with a quality score of 40 had error rates whose average  
150 corresponded to a true value of 38.6. The major deviation was at the low end of the  
151 scale, where bases with quality scores of two had considerably lower error rates than  
152 expected (Fig. 2).

153

154 **Creating quality score profiles**

155           Again using the Harvard datasets, we computed error rates in regions where the  
156 paired reads overlapped each other, for each possible combination of the two reads'  
157 quality scores. We then converted these rates back into quality scores, using the  
158 baseline error rate already calculated (Fig. 2). In cases where the bases of the two reads  
159 agreed (Fig. 3A), which amounted to 96.7% of all overlapping bases, no combinations  
160 yielded scores below 25, even where both reads had low-quality bases. However, two  
161 high-quality matching bases did not produce substantially increased quality scores,  
162 with no combined scores rising above 40. This contrasts with the scoring schemes of  
163 VSEARCH and PEAR, which assign scores of up to 85 and 80, respectively (Fig. S1).

164

165           On the other hand, where the two reads' bases disagreed (Fig. 3B), the combined  
166 quality scores were lowest when the two original quality scores were similar to each  
167 other, and rose above 30 only when the two quality scores were at opposite ends of the  
168 scale. This is similar to the schemes of fastq-join and VSEARCH, whereas PEAR does  
169 not reduce quality scores for mismatches (Fig. S1).

170

171           These two quality score models were incorporated into our merging program,  
172 NGmerge, to be utilized in the creation of merged reads.

173

174   **Comparing merging programs**

175 To compare the performance of NGmerge against other merging programs, we  
176 opted not to use the same Harvard datasets, since they had trained NGmerge's models.  
177 Instead, we queried the Sequence Read Archive (SRA) for datasets that had sufficient  
178 read lengths and PhiX content to be usable for calculating error rates. We found 33 such  
179 datasets, containing a total of 1,386 sequencing runs (see Methods and Table S1 for  
180 details).

181

182 First, we determined that each of the SRA datasets had the same five PhiX  
183 genomic sequence variants previously identified (Table 1). In addition, the baseline  
184 error rates of the datasets' reads followed a similar trend to those of the training  
185 datasets. We then processed the datasets through each of the merging programs (since  
186 neither fastq-join nor VSEARCH consider dovetailed alignments, the datasets were  
187 preprocessed by NGmerge in adapter-removal mode (Fig. 1B) prior to analysis with  
188 those programs).

189

190 All of the programs reduced the total error rates in the overlapping regions of the  
191 reads of the SRA datasets (Fig. 4A, Table S2), with NGmerge producing the lowest rate,  
192 slightly better than fastq-join. The error rate after PEAR was more than twice those of  
193 the others, due to PEAR's more aggressive merging algorithm causing a higher starting  
194 value. Increasing the fraction mismatch parameter of NGmerge (-p 0.2) led to merging

195 results similar to those of PEAR, though with a lower final error rate (Fig. 4A, Table S2).  
196 However, we found that the program CASPER produced an even lower error rate than  
197 NGmerge on a subset of the datasets, due to CASPER's reliance on k-mer based  
198 contexts to resolve mismatches.

199

200 In addition to errors, sequence reads sometimes contain ambiguous bases (Ns),  
201 which can also complicate downstream analyses. Though ambiguous bases comprised  
202 just ~0.03% of overlapping bases in the SRA datasets, NGmerge's unique approach  
203 (counting them as neither matches nor mismatches during alignment) led to the  
204 correction of the most, twice the counts of fastq-join and VSEARCH, and 25% more than  
205 PEAR (Fig. 4B).

206

207 We further examined the quality scores produced by each merging program. In  
208 cases where the overlapping reads' bases matched (Fig. 5A), the error profile produced  
209 by NGmerge closely tracked that of the original reads. With the other three mergers,  
210 bases assigned quality scores below 28 had lower error rates than expected. However,  
211 both VSEARCH and PEAR greatly overstated the quality scores at higher values. For  
212 example, at a quality score of 80, VSEARCH and PEAR produced bases whose actual  
213 error rates were  $1.8 \times 10^{-4}$  and  $1.3 \times 10^{-4}$ , respectively, more than four orders of magnitude  
214 above the theoretical value of  $1 \times 10^{-8}$ .

215

216           Where the original reads' bases did not match (Fig. 5B), NGmerge slightly  
217 underestimated the quality scores throughout most of the score range. Fastq-join and  
218 VSEARCH followed similar paths, going above the baseline profile only at the ends of  
219 the quality score range. The merged bases produced by PEAR had far higher error rates  
220 than expected throughout the quality score range.

221

## 222 **Discussion**

223           Of the merging programs analyzed, NGmerge has the best performance. It  
224 considers dovetailed alignments and thus does not require a separate adapter-removal  
225 step prior to merging reads; this more than compensates for its slightly worse run-time  
226 compared to VSEARCH (Note S1). Furthermore, NGmerge produces lower error rates  
227 and corrects more Ns than the other programs. We note that NGmerge's method for  
228 resolving mismatched bases may be improved by implementing a context-based  
229 scheme like that of CASPER. However, such an approach may have difficulty  
230 distinguishing sequencing errors from true biological variants in real samples; this is an  
231 area for further research.

232

233           In addition, NGmerge creates merged reads whose quality scores accurately  
234 reflect the bases' error rates, unlike the other merging programs. It is noteworthy that

235 the quality scores that deviate the most from the expected error rates are produced by  
236 VSEARCH and PEAR in merging matching bases (Fig. 5A). As explained above, these  
237 programs' quality score calculations are based on the assumption of the independence  
238 of sequencing errors in paired reads [17, 20]. Our results demonstrate that this  
239 assumption is false. Therefore, the models produced by Edgar and Flyvbjerg [20] are  
240 invalid.

241  
242 One reason for the lack of independence of errors in paired-end sequencing  
243 stems from the beginning of a sequencing run, during first-strand synthesis (Fig. 6).  
244 Since the original DNA fragment is denatured after it is copied, any errors made during  
245 this step will be propagated throughout the cluster that is formed during bridge  
246 amplification. Thus, both paired reads will contain the errors, but the erroneous bases  
247 will not have reduced quality scores.

248  
249 Because of our reliance on reads derived from PhiX, NGmerge's quality score  
250 profiles were tested only on datasets generated by MiSeq and HiSeq instruments. Thus,  
251 they may not work as well with other Illumina platforms, such as the NextSeq.  
252 NGmerge provides the option to forgo its default quality score profiles and instead to  
253 utilize calculations similar to those of fastq-join (and FLASH), which, though simplistic,

254 are conservative over most of the score ranges. A third option is for the user to supply  
255 custom matrices of quality score profiles to NGmerge.

256

257         Although a number of studies have concluded that Illumina’s quality scores are  
258 substantially inflated, our results contradict this notion. Inaccuracies in reference  
259 sequences are a persistent problem that adversely affect error rate calculations [22], and  
260 in fact that proved to be the case here. Once the PhiX reference genome was corrected  
261 to account for the five sequence variants, the calculated error rates closely followed the  
262 expected relationship shown in equation (1). It is important to note that, in general,  
263 errors occurring during the library preparation process (e.g. PCR amplification) can be  
264 misconstrued as sequencing errors, leading to specious conclusions [23]. This is another  
265 reason why unamplified PhiX remains an enduring control in Illumina sequencing  
266 applications.

267

## 268 **Conclusions**

269         We have examined errors produced by Illumina sequencing technology via reads  
270 derived from PhiX. We have found that variants from the canonical PhiX reference  
271 genome account for most of the discrepancy between the actual and theoretical  
272 relationships between quality scores and error rates. Furthermore, in the course of  
273 developing empirical models for error rates of paired-end sequence reads, we have

274 demonstrated the fallacy of the assumption that has been repeatedly made, both  
275 implicitly and explicitly, that errors in such reads are independent.

276

277         Finally, we have described a free and open-source program, NGmerge, that  
278 merges paired-end sequence reads, thus correcting errors and ambiguous bases, and  
279 assigning quality scores that are consistent with the measured error rates. The program  
280 can also be run in an alternative mode simply to remove contaminating sequencing  
281 adapters. Complete descriptions of the usage and options of NGmerge are found on the  
282 homepage of the software (<https://github.com/harvardinformatics/NGmerge>) and in the  
283 accompanying UserGuide. The program is written in C and is parallelized with  
284 OpenMP 4.0.

285

## 286 **Methods**

### 287 **NGmerge design**

288         NGmerge operates on paired-end reads in two distinct modes, “stitch” and  
289 “adapter-removal” (Fig. 1). In either mode, NGmerge tests all possible gapless  
290 alignments of a pair of reads in attempting to find an optimal alignment. By default,  
291 NGmerge requires that a valid alignment have a minimum overlap of 20bp and a  
292 maximum of 10% mismatches in the overlap region (-m 20 -p 0.1). If multiple valid  
293 alignments are found, the one with the lowest fraction mismatch is selected as the



294 optimum. In all of these calculations, ambiguous bases (Ns) are considered neither  
295 matches nor mismatches. When the '-d' option is set, or in adapter-removal mode,  
296 NGmerge will also attempt to align the reads in a dovetailed configuration (such as that  
297 shown in Fig. 1B), with 3' overhangs corresponding to contaminating sequencing  
298 adapters that will be removed.

299

300 In stitch mode, NGmerge forms a single merged read that spans between the 5'  
301 ends of the two original reads. The bases and quality scores of any non-overlapping  
302 regions are copied into the new read. For the overlapping region, if the bases of the R1  
303 and R2 reads match, that base is used for the merged read, with the quality score  
304 determined from the "match" matrix (see "Harvard datasets" below). Where the bases  
305 disagree, the base with the higher quality score is selected, and the "mismatch" matrix  
306 yields the merged quality score.

307

### 308 **Merging programs**

309 NGmerge (v0.2) was run with default alignment parameters, requiring a  
310 minimum overlap of 20bp and a maximum of 10% mismatches. With the SRA datasets,  
311 it was also run allowing 20% mismatches (-p 0.2). The '-d' option was set to allow for  
312 dovetailed alignments and the automatic removal of sequencing adapters.

313

314 Fastq-join (v1.01.759) [14] was run with alignment parameters analogous to those  
315 of the defaults of NGmerge (-m 20 -p 10). Because fastq-join does not allow for  
316 dovetailed alignments, adapters were removed from the reads with NGmerge prior to  
317 analysis with fastq-join.

318

319 VSEARCH (v2.6.2) [16], like fastq-join, does not consider dovetailed alignments,  
320 so it was also given reads from which adapters were removed with NGmerge. The  
321 minimum overlap length was increased to 20bp (--fastq\_minovlen 20). The maximum  
322 number of mismatches was greatly increased (--fastq\_maxdiffs 30); even so, VSEARCH  
323 still analyzed the fewest reads (Table S1). The cap on output quality scores was  
324 increased from the default value of 41 (--fastq\_qmaxout 85).

325

326 PEAR (v0.9.10) [17] was run with a 20bp minimum overlap (-v 20) and a  
327 maximum p-value threshold of 0.0001 (there is no fraction mismatch parameter). The  
328 cap on output quality scores was increased from the default value of 40 (-c 80).

329

330 With NGmerge, the arguments '-j <file> -b' were specified so that the program  
331 would produce a file listing overlap mismatches and Ns, for later error counting. With  
332 the other merging programs, a custom Python script (findDiffs.py) reconstructed the  
333 alignments and determined the overlap mismatches.

334

335 Further details of these programs' approaches toward read merging, along with  
336 an illustrative example, are provided in Note S2.

337

### 338 **Calculation of error rates**

339 The 5386-bp genome of the enterobacteria phage  $\Phi$ X174, *sensu lato*, was retrieved  
340 from NCBI (accession NC\_001422.1). The reads of each of the datasets were aligned to  
341 this genome using Bowtie2 [24], as described below. Pileup files were created from the  
342 alignment files using SAMtools (v1.5) mpileup (-B -Q 0 -d 1e9) [25], and variants were  
343 called with VarScan (v2.4.1) pileup2snp (--min-var-freq 0.15) [5].

344

345 The downloaded PhiX genome was modified to incorporate the five variants  
346 observed in the datasets (Table 1). Furthermore, because the PhiX genome is circular, a  
347 fragment corresponding to the first 1kb (including the variants at positions 587 and 833)  
348 was appended to the end of the genome. This produced a final reference genome of  
349 6386bp that was used in all further analyses.

350

351 Reads were aligned to the modified PhiX reference genome using Bowtie2  
352 (v2.3.2). The parameters of the program were modified to increase the strictness of  
353 accepted alignments, specifically by increasing the minimum score threshold (--score-

354 min L,0,-0.2) and increasing gap penalties (--rdg 5,15 --rfg 5,15). The size of allowed  
355 fragments was increased to 1kb (-X 1000), and the effort parameters were adjusted (--  
356 very-sensitive).

357

358 For the analyses of unmerged paired-end reads, contaminating sequencing  
359 adapters were removed by NGmerge ('-a' mode) prior to alignment. Only properly-  
360 paired alignments ('samtools view -f 0x2') were used to calculate error rates.

361

362 Error rates were calculated by quality score for the alignments in the SAM  
363 alignment files by a custom Python script (countErrors.py). When analyzing SAM files  
364 of merged reads, the script was provided the original read length(s) and a list of  
365 merging mismatches and Ns, in order to further categorize the errors based on the  
366 nucleotides in the R1 and R2 original reads, into matches, mismatches, and Ns (Fig. S3).  
367 The list of mismatches and Ns was produced by NGmerge or findDiffs.py, as described  
368 above.

369

370 In order to create the error profiles of NGmerge, we used the script  
371 countErrors3D.py, which tallied errors based on both original reads' quality scores.

372

373           The three custom Python scripts are freely available on GitHub  
374 (<https://github.com/jsh58/NGmerge/tree/master/scripts/>).

375

## 376 **Harvard datasets**

377           Nine sequencing datasets produced by the Bauer Core Facility of Harvard  
378 University, Faculty of Arts and Sciences Division of Science, between January 2016 and  
379 May 2017 were analyzed. Each sequencing run was produced on the Illumina HiSeq  
380 2500 platform, yielding 2×250bp paired-end reads. The reads placed into the  
381 “undetermined” bins were examined, a total of 553.0 million read pairs.

382

383           The paired-end reads were aligned to the modified PhiX genome after adapter-  
384 trimming with NGmerge, as described above. A LOESS regression function relating the  
385 quality scores to the logarithm (base 10) of the error rates was calculated in R (v3.4.1).  
386 This formed the baseline error profile for subsequent analyses.

387

388           To create the quality score profiles of NGmerge, the same reads were processed  
389 with NGmerge in stitch mode, allowing dovetailed alignments (-d). The merged reads  
390 were aligned to PhiX, and error rates were calculated for each combination of the  
391 quality scores of the R1 and R2 reads with countErrors3D.py. The “match” table,  
392 consisting of error rates for the locations where the bases of the two reads agreed, was

393 edited to exclude values derived from fewer than 1000 counts, and values of zero were  
394 given a pseudo-error count of 0.5 (yielding a rate of 0.5 / base count). These same edits  
395 were made to the “mismatch” table (error rates where the bases of the two reads  
396 disagreed), except that the minimum count threshold was lowered to 100 because of the  
397 reduced number of counts. Then, for each table, a two-dimensional LOESS regression  
398 function (relating both quality scores to the log (base 10) of the error rates) and  
399 predicted error rates were calculated in R. These error rates were then transformed  
400 back into quality scores using the baseline error profile calculated for the original  
401 paired-end reads. The resulting “match” and “mismatch” matrices were incorporated  
402 into NGmerge as the default quality score profiles.

403

#### 404 **SRA datasets**

405 The Sequence Read Archive (SRA) of NCBI ([www.ncbi.nlm.nih.gov/sra](http://www.ncbi.nlm.nih.gov/sra)) was  
406 queried for datasets containing paired-end reads that were minimum 2×250bp in length.  
407 The sequencing runs of over 160 SRA studies were examined, though some were  
408 eliminated immediately for various reasons (misabeled as paired-end; actual read  
409 lengths shorter than stated; reads already trimmed). The remaining datasets’ reads  
410 were adapter-trimmed with NGmerge and aligned to the PhiX genome. Those with at  
411 least 10,000 read pairs aligning to PhiX in a properly-paired configuration were further

412 analyzed. The details of these 33 datasets, which contained 1,386 sequencing runs and a  
413 total of 2.25 billion read pairs, are provided in Table S1.

414

415 The reads of the 33 SRA datasets were analyzed in a similar fashion to the  
416 Harvard datasets. The baseline error rates were calculated from the original reads, and  
417 error rates were also determined after processing the reads with each of the merging  
418 programs. For each set of error rates, LOESS regression functions were computed,  
419 relating the quality scores to the log (base 10) of the error rates.

420

## 421 **List of abbreviations**

422 PCR: polymerase chain reaction; SRA: Sequence Read Archive

423

## 424 **Declarations**

### 425 **Ethics approval and consent to participate**

426 Not applicable

427

### 428 **Consent for publication**

429 Not applicable

430

### 431 **Availability of data and materials**

432           The datasets analyzed in the current study, listed in Table S1, are available in the  
433 NCBI Sequence Read Archive, <https://www.ncbi.nlm.nih.gov/sra>.

434

### 435 **Competing interests**

436           The author declares that he has no competing interests.

437

### 438 **Funding**

439           This work was not supported by any specific funding.

440

### 441 **Authors' contributions**

442           JMG wrote the software, analyzed the data, and wrote the paper.

443

### 444 **Acknowledgements**

445           The author thanks Dr. Timothy Sackton of Harvard University and Dr. W. Kelley  
446 Thomas of the University of New Hampshire for helpful discussions regarding software



447 design and testing. Dr. Thomas and the Hubbard Genome Center also provided  
448 preliminary datasets for testing. The author thanks Harvard University, George  
449 Washington University, and Dr. Jeremy Goecks for funding support. The computations  
450 in this paper were run on the Odyssey cluster supported by the FAS Division of Science,  
451 Research Computing Group at Harvard University.

452

## 453 **References**

454 [1] Bentley DR, Balasubramanian S, Swerdlow HP, Smith GP, Milton J, Brown CG, et al.  
455 Accurate whole human genome sequencing using reversible terminator chemistry.  
456 Nature. 2008;456:53-9.

457

458 [2] Reuter JA, Spacek DV, Snyder MP. High-throughput sequencing technologies. Mol  
459 Cell. 2015;58:586-97.

460

461 [3] Cock PJ, Fields CJ, Goto N, Heuer ML, Rice PM. The Sanger FASTQ file format for  
462 sequences with quality scores, and the Solexa/Illumina FASTQ variants. Nucleic Acids  
463 Res. 2010;38:1767-71.

464

465 [4] Ewing B, Green P. Base-calling of automated sequencer traces using phred. II. Error  
466 probabilities. Genome Res. 1998;8:186-94.

467

468 [5] Koboldt DC, Zhang Q, Larson DE, Shen D, McLellan MD, Lin L, et al. VarScan 2:  
469 somatic mutation and copy number alteration discovery in cancer by exome  
470 sequencing. *Genome Res.* 2012;22:568-76.

471

472 [6] Pant S, Weiner R, Marton MJ. Navigating the rapids: the development of regulated  
473 next-generation sequencing-based clinical trial assays and companion diagnostics. *Front*  
474 *Oncol.* 2014;4:78.

475

476 [7] Strom SP. Current practices and guidelines for clinical next-generation sequencing  
477 oncology testing. *Cancer Biol Med.* 2016;13:3-11.

478

479 [8] Dohm JC, Lottaz C, Borodina T, Himmelbauer H. Substantial biases in ultra-short  
480 read data sets from high-throughput DNA sequencing. *Nucleic Acids Res.* 2008;36:e105.

481

482 [9] Li R, Li Y, Fang X, Yang H, Wang J, Kristiansen K, et al. SNP detection for massively  
483 parallel whole-genome resequencing. *Genome Res.* 2009;19:1124-32.

484

485 [10] DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A  
486 framework for variation discovery and genotyping using next-generation DNA  
487 sequencing data. *Nat Genet.* 2011;43:491-8.  
488

489 [11] Manley LJ, Ma D, Levine SS. Monitoring Error Rates In Illumina Sequencing. *J*  
490 *Biomol Tech.* 2016;27:125-128.  
491

492 [12] Chiara M, Pavesi G. Evaluation of Quality Assessment Protocols for High  
493 Throughput Genome Resequencing Data. *Front Genet.* 2017;8:94.  
494

495 [13] Lagier JC, Khelafia S, Alou MT, Ndongo S, Dione N, Hugon P, et al. Culture of  
496 previously uncultured members of the human gut microbiota by culturomics. *Nat*  
497 *Microbiol.* 2016;1:16203.  
498

499 [14] Aronesty E. Comparison of Sequencing Utility Programs. *Open Bioinforma J.*  
500 2013;7:1-8.  
501

502 [15] Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, et  
503 al. QIIME allows analysis of high-throughput community sequencing data. *Nat*  
504 *Methods.* 2010;7:335-6.

505

506 [16] Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: a versatile open  
507 source tool for metagenomics. *PeerJ*. 2016;4:e2584.

508

509 [17] Zhang J, Kobert K, Flouri T, Stamatakis A. PEAR: a fast and accurate Illumina  
510 Paired-End reAd mergeR. *Bioinformatics*. 2014;30:614-20.

511

512 [18] Magoč T, Salzberg SL. FLASH: fast length adjustment of short reads to improve  
513 genome assemblies. *Bioinformatics*. 2011;27:2957-63.

514

515 [19] Kwon S, Lee B, Yoon S. CASPER: context-aware scheme for paired-end reads from  
516 high-throughput amplicon sequencing. *BMC Bioinformatics*. 2014;15 Suppl 9:S10.

517

518 [20] Edgar RC, Flyvbjerg H. Error filtering, pair assembly and error correction for next-  
519 generation sequencing reads. *Bioinformatics*. 2015;31:3476-82.

520

521 [21] Sanger F, Coulson AR, Friedmann T, Air GM, Barrell BG, Brown NL, et al. The  
522 nucleotide sequence of bacteriophage phiX174. *J Mol Biol*. 1978;125:225-46.

523

524 [22] Gaspar JM, Thomas WK. FlowClus: efficiently filtering and denoising  
525 pyrosequenced amplicons. BMC Bioinformatics. 2015;16:105.

526

527 [23] Eren AM, Vineis JH, Morrison HG, Sogin ML. A filtering method to generate high  
528 quality short reads using illumina paired-end technology. PLoS One. 2013;8:e66643.

529

530 [24] Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat  
531 Methods. 2012;9:357-9.

532

533 [25] Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence  
534 Alignment/Map format and SAMtools. Bioinformatics. 2009;25:2078-9.

535

536

## 537 **Figures**

538 **Figure 1. Analysis modes of NGmerge.** The diagrams show the paired-end reads (R1,  
539 R2) derived from sequencing DNA fragments (white boxes) with sequencing adapters  
540 (gray boxes) on either end. **A:** In the default mode (“stitch”), NGmerge combines  
541 paired-end reads that overlap into a single read that spans the full length of the original  
542 DNA fragment. **B:** The alternative “adapter-removal” mode returns the original reads  
543 as pairs, removing the 3’ overhangs of those reads whose optimal alignment has this

544 characteristic.

545

546 **Figure 2. Error rate - quality score relationship.** The quality scores of the original  
547 paired-end reads in the Harvard datasets followed a nearly linear relationship with the  
548 log of the error rates, consistent with expectations.

549

550 **Figure 3. Quality score profiles of NGmerge.** **A:** A plot of the quality scores  
551 corresponding to the error rates calculated for each combination of the two reads'  
552 quality scores, for cases where the bases matched. **B:** Same as A, but for cases where the  
553 bases of the paired reads did not match.

554

555 **Figure 4. Errors and Ns corrected by the merging programs.** **A:** Error rates in the  
556 paired reads' overlap regions, before and after the application of the merging programs.  
557 Note that the "Before" error rates vary because different merging programs analyze  
558 slightly different sets of reads (see Table S1). **B:** Total number of Ns corrected by each  
559 of the merging programs.

560

561 **Figure 5. Error rate - quality score profiles produced by the merging programs.** **A:**  
562 Comparison of the profiles when the overlapping bases of the reads matched. The

563 black line represents the baseline error rate - quality score profile of the original reads.

564 **B:** Comparison of the profiles when the overlapping bases did not match.

565

566 **Figure 6. First-strand synthesis on the flow cell.** A single-stranded DNA fragment to  
567 be sequenced anneals to an oligonucleotide that is covalently attached to the flow cell  
568 surface. The primer is extended to copy the DNA fragment, which is then removed by  
569 denaturation [1].

570

## 571 **Tables**

572 **Table 1. Variants in the PhiX genome.**

Position (1-based)	NCBI base	iGenomes base	Observed base
587	G	A	A
833	G	A	A
2731	A	G	G
2793	C	T	C
2811	C	T	T
3133	C	C	T

573 Four of the five observed variants are in the version of the PhiX genome provided by

574 Illumina on its iGenomes website (<https://support.illumina.com/sequencing/>

575 [sequencing\\_software/igenome.html](https://support.illumina.com/sequencing/sequencing_software/igenome.html); retrieved Nov. 2017), and it lists an additional

576 variant at position 2793 that was not observed in any sequencing run. Nevertheless, the

577 iGenomes version is considered current by Illumina (personal communication, Nov.  
578 2017).

579

## 580 **Supplementary Information**

581 **Figure S1. Quality score profiles of the merging programs**

582 **Figure S2. PhiX library fragment lengths and paired-end read overlaps**

583 **Figure S3. Error rate calculation**

584 **Table S2. Error rates before and after merging**

585 **Note S1. Benchmarking of merging programs**

586 **Note S2. Methods for producing a merged read**

587 **Note S3. Additional notes on the merging programs**

588 **Note S4. NGmerge-PEAR disagreement**

589 **Table S1. Details of the 33 datasets downloaded from SRA.** The “Number of  
590 Sequencing Runs” gives the number of individual SRA “run” datasets that appeared to  
591 have been derived from the same sequencing run and thus were concatenated.  
592 Multiple datasets from the same SRA study that were not combined were given  
593 separate designations (“\_0”, “\_1”, etc.). The “Instrument” column is based on the  
594 metadata given in each SRA dataset and may not be accurate.