

# Accurate detection of single nucleotide polymorphisms using nanopore sequencing

Espada Rocio<sup>1</sup>, Zarevski Nikola<sup>1,2</sup>, Drame-Maigne Adele<sup>1</sup>, and Rondelez Yannick<sup>1</sup>

<sup>1</sup>Gulliver, ESPCI Paris, PSL University, CNRS, 75005 Paris, France

<sup>2</sup>CRI Paris, Paris, France

## Abstract (350 words)

## 1 Background

**Background** Nanopore sequencing is a powerful single molecule DNA sequencing technology which provides a high throughput and long sequence reads. Nevertheless, its relatively high native error rate limits the direct detection of point mutations in individual reads of amplicon libraries, as these mutations are difficult to distinguish from the sequencing noise.

**Results** We propose a computational method to reduce noise in nanopore detection of point variations. Our approach uses the fact that all reads are expected to be very similar to a wild type sequence, for which we experimentally characterize the position-specific systematic sequencing error pattern. We then use this information to reweight, in individual reads from the variant library, the confidence given to nucleotides read that do not match the wild type. We tested this method on two sets of known variants of Klen Taq, where the true mutation rate was 3.3 mutations per kb, well below the sequencing noise. We observed that the actual mutations became more distinguishable from sequencing noise after correction. This approach can be used, for example to help the clustering of variants, or to decrease the number of reads necessary to call a consensus.

**Conclusions** The computational method is simple to implement and requires only a few thousands reads of the wild type sequence of interest, which can be easily obtained by multiplexing in a single minION run. The approach does not require any modification in the experimental protocol for sequencing and can be simply implemented downstream standard base calling.

## Keywords

minION, nanopore sequencing, next generation sequencing, amplicons, SNP detection, logistic regression.

Nanopore is a powerful technology for high throughput DNA sequencing, currently commercialized by Oxford Nanopore Technologies [1]. It provides sequences base calls reconstructed from conductivity records during the translocation of a single DNA molecule through a protein pore. This approach can be implemented in small devices, such as minION, which offers portability, real time sequencing, simple protocols, and a relatively low cost. A minION device can provide sequences for DNA strands of various lengths, from PCR products up to megabases genomic fragments, and provides a minimum of  $5 \cdot 10^9$  bases in one run. These characteristics make it an attractive device for, among other applications, sequencing amplicon libraries that are too long for other next generation sequencing technologies. Here we will discuss the case where the amplicon library has high diversity but low variability, i.e. it contains many different sequences but differing from each other by just a few point mutations. This is the case for example, when the genetic sample originate from a single ancestral sequence (the wild type) that has been submitted to error-prone replication. In these cases, although minION can provide full length read irrespective of the size of the DNA fragments, its relatively high error rate ( $\approx 5-10\%$ ) prevents the accurate detection of point genetic variation directly from individual reads [2, 3].

Several experimental approaches have been developed to mitigate these errors and provide more accurate sequences. One possibility is to read several times each amplicon, and then construct a consensus sequence. This has been done by creating sequence concatenates, for example using rolling circular amplification, [4, 5, 6] or via gene barcoding prior to amplification [7]. These methods usually reduce the number of different variants that can be studied, because a part of minION throughput is invested in reading duplicated reads, and is sensitive to bias occurring during the amplification steps [8]. Other approaches, used for genome assembly, combine long minION reads with shorter and more accurate reads obtained via other technologies, such as Illumina sequencing-by-synthesis approach [9].

74 In this paper, we propose a computational protocol to im-  
75 prove variant detection in individual reads from libraries  
76 for which a reference gene is known, using standard 1D  
77 protocol minION sequencing. We base our method on two  
78 observations made during the sequencing of many (identical)  
79 copies of the parent sequence. First, the confidence or  
80 quality scores ( $Q_{\text{score}}$ ) assigned by the base calling process  
81 to each nucleotide are usually low when a wrong nucleotide  
82 is assigned (Suppl. figure S2), as expected. Second, the er-  
83 rors are not homogeneously distributed, and they are more  
84 frequent in some positions of the DNA (Suppl. figure S1).  
85 These observations suggest that it should be possible to re-  
86 duce the non-random part of the sequencing errors, using  
87 the information contained in the ( $Q_{\text{score}}$ ). The method we  
88 propose has two steps: the first one uses the reference reads  
89 to build a statistical model of the error pattern. Here we  
90 used a position and nucleotide-specific logistic regression.  
91 In the second step, this information is used to re-analyze  
92 minION base calls for the variant library and to update the  
93 confidence value of each nucleotide read in this dataset.  
94 We tested our method using the gene of Klen Taq DNA  
95 polymerase, a truncated variant of the well-known Taq  
96 polymerase. We trained our model using approximately  
97 6000 reads of the wild type and applied it on two toy li-  
98 braries containing 10 known variants with 2 to 9 point  
99 mutations. Our computational correction reduced the se-  
100 quencing noise, allowing a better identification of true  
101 point mutations. The corrected confidence values were  
102 used to clusterize sequencing reads of the same variant  
103 with an accuracy of 96% leaving aside barcodes or other  
104 physical links. Furthermore, in the context of point mu-  
105 tant libraries, the corrected confidence values reduce the  
106 number of reads needed to obtain an exact consensus se-  
107 quence by at least 5 times, with good performance ob-  
108 tained using only typically 10 reads.

## 109 2 Method

110 Like other sequencing approaches, nanopore data analysis  
111 pipeline provide reads where each base is associated with a  
112 confidence values ( $Q_{\text{score}}$ ). This number reflects the prob-  
113 ability that the assigned nucleotide at that position is the  
114 correct one  $p_{\text{right}}$ , via  $p_{\text{right}} = 1 - 10^{-\frac{Q_{\text{score}}}{10}}$ .  
115 We first characterized the errors and  $Q_{\text{score}}$  distributions  
116 on more than 6000 reads of the wild type Klen Taq gene  
117 (length 1665 nucleotides), for which we have a ground  
118 truth sequence (see Suppl. table S1). In this data set,  
119 we can confidently attribute mismatches between the read  
120 nucleotide and the wild type as sequencing errors. At  
121 each position, and for each of the three non wild type nu-  
122 cleotides and deletions, we quantified the proportion of  
123 errors made by minION and guppy base-caller or, equiv-  
124 alently, the proportion of correct reads (which we denom-  
125 inate  $n_{\text{correct}}$ ), and grouped them according to the  $Q_{\text{score}}$

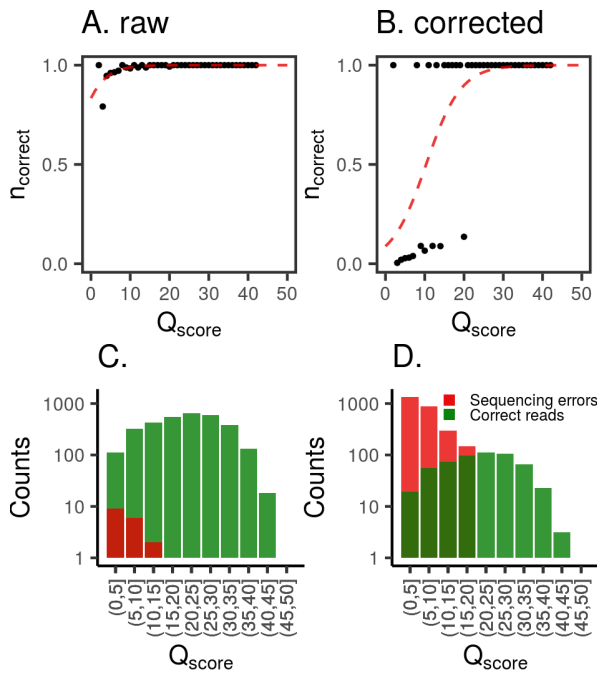
126 provided by the basecaller.  $n_{\text{correct}}$  is generally smaller for  
127 low  $Q_{\text{score}}$  (figure 1A and S4), reflecting that most errors  
128 appear where guppy base-calling confidence is low (figure  
129 1C). However, the pattern strongly differs base-to-base and  
130 position-to-position (see examples in Suppl. fig. S4). We  
131 fitted this relation position and base-wise by a logistic re-  
132 gression, which provided a classifier able to convert the re-  
133 ported  $Q_{\text{score}}$  to the probability that this read is indeed cor-  
134 rect ( $p_{\text{right}}$ ). This was made separately for the sequences  
135 read for sense and antisense strands. To include deletions  
136 in this analysis (which do not have a  $Q_{\text{score}}$  assigned by the  
137 basecaller), we fixed their confidence value as the average  
138 of the  $Q_{\text{score}}$  of their nearest neighbors in the nucleotide  
139 sequence. This decision was inspired by the observation  
140 that  $Q_{\text{score}}$  are correlated between consecutive nucleotides  
141 (Fig. S3). Insertions were ignored. All together, we ob-  
142 tained 13320=1665·4·2 regressions, one for each position  
143 and non-wild type nucleotide/deletion, in the forward and  
144 reverse sense of sequencing.

145 With this information in hand, we then looked at mis-  
146 matches in reads originated from mutated variant se-  
147 quences. For each of them, we can now use the correspond-  
148 ing logistic regression to compute  $p_{\text{right}}$ , the probability  
149 of actually being a mutation according to the associated  
150  $Q_{\text{score}}$ . The higher the value of  $p_{\text{right}}$ , the more likely it is  
151 a true mutation.

152 However, the training set uses wild type sequences, which  
153 do not contain true mutations and hence our naive classi-  
154 fier is heavily biased against classifying mutations, as rep-  
155 resented in figures 1A and 1C. This is not representa-  
156 tive of the proportion of errors/correct reads present in  
157 the mismatches of a set of mutants, even if the muta-  
158 tion rate is low. To adapt the classifier, we need an *a*  
159 *priori* expectation of mutations ( $p_{\text{prior-right}}$ ), that should  
160 come from independent information. In the present case,  
161 the variant sequences originates from an error prone PCR  
162 (ePCR) process, for which we possess an estimate of the  
163 mutation frequency. We also obtain the *a priori* expect-  
164 ation of an observed mismatch to be a sequencing error  
165 ( $p_{\text{prior-error}}$ ) as the sequencing error rate at that position  
166 in the wild type set. We therefore compute  $p_{\text{prior-right}}$  and  
167  $p_{\text{prior-error}}$ , which we use to reweight the wild type set be-  
168 fore fitting. This shifts the logistic regression towards the  
169 higher  $Q_{\text{score}}$ , allowing the classifier to accept a number of  
170 observed mismatches consistent with the prior expectation  
171 (figure 1B). We then proceed to re-score the mismatches  
172 observed in the mutant library, taking into account posi-  
173 tion, nucleotide and sense of the strand.

## 174 3 Results

175 To test our method, we built two toy libraries of Klen  
176 Taq mutants and sequenced them using the minION de-  
177 vice (see experimental methods). The first one consists



**Figure 1:** Example of logistic regression over reads of a known wild type sequence to re-scale  $Q_{\text{score}}$  at one particular position and nucleotide. A: black dots are the proportion of correct reads against the  $Q_{\text{score}}$  reported by guppy base-caller. Red line is the logistic regression performed over this data. C: Distribution of correct reads (green) and sequencing errors (red). B and D: Same plots as A and C, when each read is weighted according to the prior probabilities.

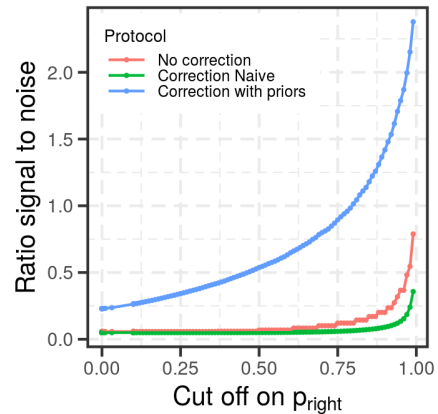
of 7 known variants uniquely barcoded before sequencing, thus we know exactly what mutations to expect in each read. The second toy library is a mixture of three known variants. In this case we know which mutations to expect in the library but we do not have a mark on each sequencing read. The complete list of mutations is available in Suppl. Table S2.

We used the logistic regressions performed on the wild type data to convert  $Q_{\text{score}}$  into a probability of being a true mutation the  $p_{\text{right}}$  of each nucleotide, as described in section Method.

### 3.1 Signal to noise ratio

MinION's high error rate makes it difficult to distinguish actual mutations from sequencing errors in single reads, especially when the mutation rate is low. A straightforward procedure to filter errors is to only trust read position which have a high probability of being correct. In this section, we compare how it performs to carry this process using the raw  $p_{\text{right}}$  returned by guppy base-caller on minION reads, the  $p_{\text{right}}$  after fitting with a logistic regression (or correction naive), or  $p_{\text{right}}$  after fitting with a logistic regression using prior probabilities.

We used the library containing seven known mutants to



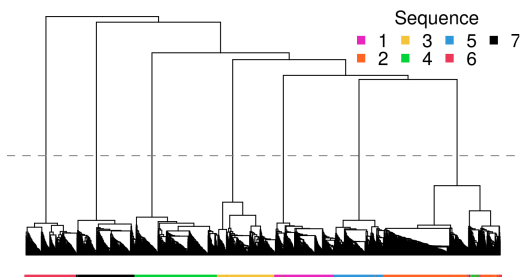
**Figure 2:** Signal to noise ratio as a function of the cut-off on the quality values  $p_{\text{right}}$ , for the three methods analyzed: using  $p_{\text{right}}$  provided by minION (red), using  $p_{\text{right}}$  after naive correction (green) and  $p_{\text{right}}$  after correction with priors (blue). The last method improves the signal to noise ratio for any  $p_{\text{right}}$  cut off.

quantify the signal to noise ratio. In this example each sequence is barcoded, so we know if each mismatch is a sequencing error or an actual mutation. We defined *signal* as the number of mismatches known to be mutations with a  $p_{\text{right}}$  higher than the threshold (true positives), and *noise* as the number of those mismatches known to be non mutated (false positives). The counts are weighted by  $p_{\text{right}}$  for each nucleotide. Results are shown in figure 2. For all thresholds, the proposed correction has a higher signal to noise ratio, thus facilitating the identification of actual mutations. This remains true when no cut off is applied (cut off = 0). Results are similar if we count number of mismatches over the threshold without weighting them by  $p_{\text{right}}$  (Suppl. fig S5). We noticed that the signal to noise ratio improvement is mainly driven by the reduction of the false positive rate, i.e. reduction of sequencing noise. Nevertheless, the true positive rate also decreases faster (Suppl. fig S6). This should be taken into account if the goal is to detect mutations which are poorly represented in the sequenced set.

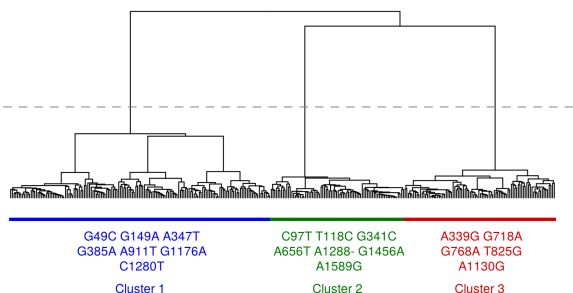
### 3.2 Sequences clustering

To show that the reduction of noise is relevant to characterize an ensemble of reads into sequence clusters, we used the corrected probabilities  $p_{\text{right}}$  to group reads by weighted sequence similarity. We defined the dissimilarity between two reads  $s_1$  and  $s_2$  as a modified hamming distance:

$$d(s_1, s_2) = \sum_{n=\text{position}} f(s_1^n, s_2^n) \quad (1)$$



**Figure 3:** Hierarchical clustering of reads of a toy library constituted by 7 mutants of the Klen Taq gene, each of them experimentally identified by a barcode. The colors below each leaf indicate to which variant the read corresponds to. The clustering groups reads of the same variant.



**Figure 4:** Hierarchical clustering of reads of a toy library constituted by 3 variants. Clusters are marked in colors. After clustering, we calculated the consensus sequence within each of the four groups detected. Detected mutations with respect to the wild type gene are annotated below. Each cluster correspond to a known variant.

where

$$f(s_1^n, s_2^n) = \begin{cases} 0 & s_1^n = \text{wt} \text{ and } s_2^n = \text{wt} \\ +p_{\text{right}}(s_1^n)p_{\text{right}}(s_2^n) & s_1^n \neq s_2^n \\ -p_{\text{right}}(s_1^n)p_{\text{right}}(s_2^n) & s_1^n = s_2^n \neq \text{wt} \end{cases} \quad (2)$$

where  $p_{\text{right}}(s_i^n)$  is  $p_{\text{right}}$  assigned to sequence  $i$  at position  $n$ . We summed zero if both sequences are equal to wild type (wt), we summed the product of corrected probabilities of both sequences if they are different (one mutation present in only one of them), and we subtracted the product of corrected probabilities of both sequences if they are the same and differ from wild type (both reads display the same mutation). Then we used R implementation of Ward's minimum variance method to clusterize the reads. We evaluated this procedure on the toy library of 7 variants. The resulting dendrogram is displayed in figure 3, where each leaf is a different read, and the color of the point below represents which Klen Taq variant it is. When cutting the dendrogram at a fixed height to obtain seven clusters, we can find a prevailing variant on each of them, which allows us to use this as a classification method.

Out of 3538 reads analyzed, 3381 were correctly clustered (96%). This was not true when using guppy base-caller's  $p_{\text{right}}$  values, in which only 63% were correctly classified (see suppl. fig S7).

We also evaluated the toy library of three mixed variants (named variants 8, 9 and 10, see sequences at suppl. table S2). We measured dissimilarities between reads and obtained the dendrogram in figure 4. We cut the dendrogram in a fixed height to obtain three clusters denoted in colors in the figure. We calculated the consensus sequence for the reads belonging to the same cluster. In the figure, we annotated the mutations for each consensus. We found that all three match one of the known variant: cluster 1 is variant 8, cluster 3 is variant 9 and cluster 2 is variant 10 (in this last case, the consensus sequence adds a deletion at position 1288).

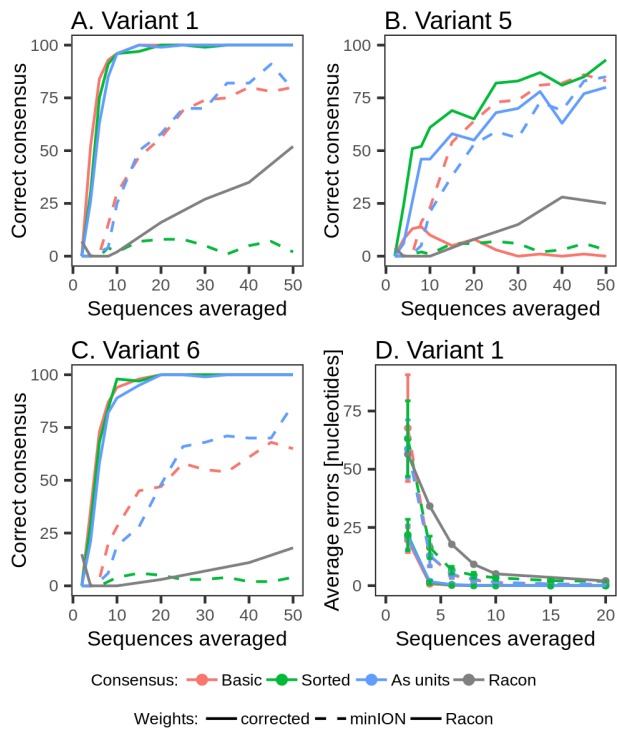
### 3.3 Consensus calculation

One strategy for lowering noise in nanopore sequencing consists in computing a variant consensus sequence (VCS), i.e. reading several times a variant and use these reads to calculate the most likely sequence [4, 10]. Usually, average on more reads leads to a more accurate VCS. We hypothesized that the corrected weights should help to converge faster on the actual sequence, as they increase the confidence on mutations over the sequencing errors. We tested this hypothesis over the set of 7 known Klen Taq variants sequenced using minION.

We computed the VCS over a set of aligned reads of the same Klen Taq variant (using the fact that these reads were barcoded). For each position, we calculated the weighted frequency of each nucleotide and deletions. The weights used are the  $p_{\text{right}}$  returned by guppy base-caller, or the ones obtained in the logistic regression. The nucleotide (or deletion) with the highest weighted count is kept as the consensus in that position.

We computed the consensus on a set of 2 to 50 sequences drawn randomly, and repeated 100 times for each set size. In figures 5 and S9 (red curves) we show how many times the VCS matches exactly the actual sequence. In all cases where there are no deletions (variants 1-4 and 7), or where there is a deletion in a non-homopolymer position (variant 6), the convergence on the actual sequence is faster when using the corrected weights: perfect consensus are more than 95% of attempts starting from typically 10 sequences. When using raw  $Q_{\text{score}}$  provided by guppy base-caller this performance is not obtained in sets of up to 50 sequences. Results are different for variant 5 which has a deletion within a homopolymer ('GG' at positions 489 and 490), making it a challenging target. The basic VCS matches performed worse with corrected scores and including more sequences even reduces the correct percentage (5B, red solid and dashed curves). We noticed that this is produced by the misrepresentation of the deletion. The VCS is com-





**Figure 5:** A-C Number of times the VCS matches the actual sequence according to the number of reads used to compute it. We used different strategies for computing the consensus: In red, no pre-processing of sequences. In green, sorting homopolymers to move deletions towards the end. In blue, treating homopolymers as a unit. Panel A show the results for variant 1 (no deletions), panel B for variant 5 (which has a true deletion within a homopolymer), and panel C for variant 6 (which has one deletion in a non-homopolymeric position). Panel D shows the average error per consensus sequence according to the number of sequences used. In all cases we compare the use of guppy base-caller weights (dashed line) and weights corrected by the logistic regression method (solid line). In grey, results using racon software.

291 puted position by position, and the aligner used (LAST)  
 292 does not always represent the deletion in the same position  
 293 (sometimes is 'G-' and sometimes is '-G'). Thus the pres-  
 294 ence of a deletion is averaged out. We tested two strategies  
 295 to address this issue: on the one hand we tried to sort all  
 296 homopolymers so that the deletions are always at the end  
 297 before computing the VCS. Results are shown in green  
 298 curves in fig. 5). In this case, the results for variant 5 im-  
 299 proved, reaching a performance slightly better than using  
 300  $Q_{\text{score}}$  provided by guppy base caller, while similar perfor-  
 301 mance for all other variants compared to the basic consen-  
 302 sus calculation. Notably, this strategy is counter-productive  
 303 for the VCS calculation using guppy raw scores in all vari-  
 304 ants (green dashed lines). On the other hand, we tried  
 305 evaluating the homopolymers as a unit instead of position  
 306 by position (blue curves in fig. 5). Here again more, the  
 307 results for variant 5 improved, reaching a performance as

308 good as using  $Q_{\text{score}}$  from guppy base-caller, without this  
 309 affecting other variants. For a deeper discussion on this  
 310 point, refer to the Supplementary material.

311 We also evaluated how different the VCS are from true  
 312 sequences. In figures 5D and S10, we plotted the mean  
 313 number of wrong nucleotides in the VCS. When using the  
 314 corrected weights, it is enough to use 6 reads to obtain a  
 315 mean error of less than one nucleotide in the VCS, (in all  
 316 cases except variant 5). This number increases up to 10  
 317 sequences when using raw  $Q_{\text{score}}$  as weights. Also, there  
 318 are no significant changes among the different consensus  
 319 strategies used.

320 As a reference, we computed consensus sequences using  
 321 Racon, a state of the art tool [10] (grey curves in fig. 5, S9  
 322 and S10). In all cases, we observed that the convergence  
 323 for Racon requires more reads than the methods presented  
 324 in this manuscript. However, we note that Racon does not  
 325 use a reference sequence, which makes direct comparison  
 326 difficult.

## 327 4 Conclusions

328 Single molecule MinION sequencing comes with a rela-  
 329 tively high error rate, which limits some applications such  
 330 as the analysis of libraries containing many different, but  
 331 genetically similar, sequences. This is for example the case  
 332 for amplicon libraries used in directed evolution. The ap-  
 333 proach that we propose here leverage the fact that the er-  
 334 rors are partly systematic, as previously noted [11]. We  
 335 therefore accumulate many reads from the reference gene  
 336 to build a sequence specific error model that can locally  
 337 correct for the sequencing biases. Applying this procedure  
 338 on two toy libraries of the Klen Taq gene with an average  
 339 mutation rate of 3.3 bases/kb, we showed that correcting  
 340 the confidence values provides large increase in the signal  
 341 to noise ratio. This can be used to cluster reads and com-  
 342 pute accurate consensus with a minimized burden on the  
 343 sequencing throughput.

344 An important ingredient of our approach is a correct prior  
 345 for the number of mutations in the train and test set. Here,  
 346 our reference sequence was assumed to be perfect, and we  
 347 could precisely evaluate the average mutation rate in the  
 348 test set, because it originated from a controlled experi-  
 349 mental mutagenesis protocol. In other situations it would  
 350 be possible to use short read sequencing, for example Il-  
 351 lumina, to evaluate this number. If the full sequence is  
 352 submitted to short read high quality sequencing, it would  
 353 even be possible to obtain more precise priors, for exam-  
 354 ple specific to each position and nucleotide. Our approach  
 355 would then be used to associate these mutations on single  
 356 long reads.

357 An underlying assumption of our method is that the dis-  
 358 tribution of  $Q_{\text{score}}$  observed at a particular position for  
 359 the wild type base appropriately reflects the distribution

of  $Q_{\text{score}}$  that would be observed for a variant base at that position. This approximation is necessary since the error model is build from a single sequence and hence has a single “true” base per position. Fortunately, the difference in  $Q_{\text{score}}$  distributions for “true” versus “error” seems large enough for our method to perform well within that approximation. Moreover, when working with large variant libraries, it would in principle to use an iterative process to complete the error model using observation made on the variants.

The protocol proposed needs to characterize the sequencing errors done on appropriate reference sequence. As such it is limited to analyze variants which are close neighbors of the reference, and where mutations can be considered to be independent. We did not try to adapt the method to detect alterations beyond point replacements or deletions, which may require a more complex analysis pipeline. Encouragingly, variant 6 contained two contiguous mutations, and was properly analysed by our consensus approach. We also note that, in the case of highly diverse sequences, and if the library can be sequenced at sufficient depth, it become easier to cluster sequences and compute direct consensus. Finally, our approach provides large improvement of signal to noise at very little experimental effort or throughput reduction. This is because only the reference DNA needs to be sequenced many times, and this can be done simultaneously with the libraries, using standard barcoding protocols. There is no other modifications to the experimental protocol, and the computational error correction process can be simply added to any analysis pipeline after base calling.

## 5 Experimental methods

### Samples preparation

Klen Taq wild type gene was amplified using a high fidelity PCR (Q5 polymerase from NEB) from a stored plasmid. See DNA sequence in Suppl. table S1.

Mutants were obtained via error prone PCR (ePCR) using Agilent’s kit GeneMorph II. We started from 1.1 nM of dam-methylated DNA. We used primers GGGAT-TATTCTTTGGCGCTCAGCCAAT and ACCATGCGTCT-GCTGCATGAAT. Thermocycling was performed as follows: 95°C for 2 min, followed by 25 cycles of [95°C for 30sec + 65°C for 30sec + 72°C for 2min] and a final extension at 72°C for 10 min. We digested the product with DpnI (NEB) and purified it using columns (Macheray-Nagel). We put the mutagenized genes in a pIVEX vector via Gibson assembly (NEB Hi-Fi DNA assembly) using 125ng of gene DNA, 100 ng of vector in a 2:1 insert:vector molar ratio and incubated for 15min at 50°C. We purified and concentrated DNA with a Zymo Research kit. We transformed the product into chemocompetent KRX bacteria. We spread them on a Petri dish with LB and Ampicillin. We incubated overnight and picked some clones randomly. We verified the presence of the plasmid via colony PCR (using Dream-

Taq polymerase from Thermofisher). 10 positive clones were grown overnight in liquid LB with antibiotic and mini-prepped to obtain the plasmid DNA. A fraction was used for high quality sequencing (Sanger sequencing), and another fraction used for minION sequencing.

### minION sequencing

We used the prepared DNA of each clone and the wild type gene for amplification by PCR with Q5 polymerase (NEB) using primers which included the minION barcodes adapters: ACTTGCCCTGTCGCTCTATCTTCAGTGT-GCTGGAATTCGCCCTTTTA and TTTCTGTTGGTGCT-GATATTGCAGACCACAACGGTTTCCCTCTAGAAATA.

Thermocycling was performed as follows: 98°C for 30sec, 23 cycles of [98°C for 10sec + 59°C for 30sec + 72°C for 1min], final extension at 72°C for 2min.

We digested the product with Dpn1 (NEB), gel purified it using Macheray-Nagel kit. We proceeded following standard minION protocols. We used one barcode for wild type, one for each of the mutants 1 to 7, and one extra barcode for mutants 8-10. We then mixed the DNA and followed standard protocol for preparing and loading the DNA into a flowcell.

We used kits EXP-PCB001 for barcoding, SQK-LSK108 for ligation, EXP-LLB001 for flowcell loading, and minION flowcell version was R9.4/FLO-MIN106, thus the sequencing was 1D.

### minION reads pre-processing

minION raw data was base called and demultiplexed using ONT Guppy version 3.3.3. We obtained 140197 reads from which 45635 (33%) did not match any barcode. Each sequence was pairwise aligned to wild type Klen Taq sequence using LAST [12]. Those sequences with low alignment score were discarded, keeping 32507 at the end (23% from original reads).

### Racon consensus

We selected a subset of sequences on the fastq file after base calling using an R script, we aligned them to the wild type sequence using minimap2 2.17-r941 [13] (default parameters), and computed the VCS with racon v1.4.11 [10] (default parameters).

### Availability

R code for the analysis is available at <https://github.com/rociocespci/minIONLogReg>.

### Acknowledgements

This project has received funding from the European Union’s Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No 845976 and from the European Research Council (ERC, Consolidator Grant No. 647275 ProFF).

### Conflicts of interest

The authors declare no conflicts of interest.

## References

- 461 [1] <https://nanoporetech.com/>.
- 462 [2] Fritz J Sedlazeck, Philipp Rescheneder, Moritz Smolka,  
463 Han Fang, Maria Nattestad, Arndt von Haeseler, and  
464 Michael C Schatz. Accurate detection of complex struc-  
465 tural variations using single-molecule sequencing. *Nature*  
466 *methods*, 15(6):461–468, 2018.
- 467 [3] Liang Gong, Chee-Hong Wong, Wei-Chung Cheng,  
468 Harianto Tjong, Francesca Menghi, Chew Yee Ngan, Edi-  
469 son T Liu, and Chia-Lin Wei. Picky comprehensively de-  
470 tects high-resolution structural variants in nanopore long  
471 reads. *Nature methods*, 15(6):455–460, 2018.
- 472 [4] Chenhao Li, Kern Rei Chng, Esther Jia Hui Boey, Amanda  
473 Hui Qi Ng, Andreas Wilm, and Niranjan Nagarajan. Inc-  
474 seq: accurate single molecule reads using nanopore se-  
475 quencing. *Gigascience*, 5(1):s13742–016, 2016.
- 476 [5] Roger Volden, Theron Palmer, Ashley Byrne, Charles  
477 Cole, Robert J Schmitz, Richard E Green, and Christo-  
478 pher Vollmers. Improving nanopore read accuracy with  
479 the r2c2 method enables the sequencing of highly multi-  
480 plexed full-length single-cell cDNA. *Proceedings of the Na-*  
481 *tional Academy of Sciences*, 115(39):9726–9731, 2018.
- 482 [6] Szymon T Calus, Umer Z Ijaz, and Ameet J Pinto.  
483 Nanoampli-seq: a workflow for amplicon sequencing for  
484 mixed microbial communities on the nanopore sequencing  
485 platform. *Gigascience*, 7(12):giy140, 2018.
- 486 [7] Søren M Karst, Ryan M Ziels, Rasmus H Kirkegaard,  
487 Emil A Sørensen, Daniel McDonald, Qiyun Zhu, Rob  
488 Knight, and Mads Albertsen. Enabling high-accuracy  
489 long-read amplicon sequences using unique molecular iden-  
490 tifiers with nanopore or pacbio sequencing. *bioRxiv*, page  
491 645903, 2020.
- 492 [8] Marc A Sze and Patrick D Schloss. The impact of dna  
493 polymerase and number of rounds of amplification in pcr  
494 on 16s rrna gene sequence data. *mSphere*, 4(3):e00163–19,  
495 2019.
- 496 [9] Sara Goodwin, James Gurtowski, Scott Ethe-Sayers, Pan-  
497 chajanya Deshpande, Michael C Schatz, and W Richard  
498 McCombie. Oxford nanopore sequencing, hybrid error cor-  
499 rection, and de novo assembly of a eukaryotic genome.  
500 *Genome research*, 25(11):1750–1756, 2015.
- 501 [10] Robert Vaser, Ivan Sovic, Niranjan Nagarajan, and Mile  
502 Sikic. Fast and accurate de novo genome assembly from  
503 long uncorrected reads. *Genome research*, 27(5):737–746,  
504 2017.
- 505 [11] Raga Krishnakumar, Anupama Sinha, Sara W Bird,  
506 Harikrishnan Jayamohan, Harrison S Edwards, Joseph S  
507 Schoeniger, Kamlesh D Patel, Steven S Branda, and  
508 Michael S Bartsch. Systematic and stochastic influences on  
509 the performance of the minion nanopore sequencer across  
510 a range of nucleotide bias. *Scientific reports*, 8(1):1–13,  
511 2018.
- 512 [12] Szymon M Kielbasa, Raymond Wan, Kengo Sato, Paul  
513 Horton, and Martin C Frith. Adaptive seeds tame genomic  
514 sequence comparison. *Genome research*, 21(3):487–493,  
515 2011.
- [13] Heng Li. Minimap2: pairwise alignment for nucleotide  
517 sequences. *Bioinformatics*, 34(18):3094–3100, 2018. 518