

ELECTOR: Evaluator for long reads correction methods

Camille Marchet ^{1,3†}, Pierre Morisse ^{2†}, Lolita Lecompte ³,
Arnaud Lefebvre ², Thierry Lecroq ², Pierre Peterlongo ³ and Antoine Limasset ¹

¹Univ. Lille, CNRS, Inria, UMR 9189 - CRISTAL.

²Normandie Univ, UNIROUEN, LITIS, Rouen 76000, France.

³Univ Rennes, CNRS, Inria, IRISA - UMR 6074, F-35000 Rennes, France.

[†]these authors contributed equally to this work

Abstract

Motivation: In the last few years, the error rates of third generation sequencing data have been capped above 5%, including many insertions and deletions. Thereby, an increasing number of long reads correction methods have been proposed to reduce the noise in these sequences. Whether hybrid or self-correction methods, there exist multiple approaches to correct long reads. As the quality of the error correction has huge impacts on downstream processes, developing methods allowing to evaluate error correction tools with precise and reliable statistics is therefore a crucial need. Since error correction is often a resource bottleneck in long reads pipelines, a key feature of assessment methods is therefore to be efficient, in order to allow the fast comparison of different tools.

Results: We propose ELECTOR, a reliable and efficient tool to evaluate long reads correction, that enables the evaluation of hybrid and self-correction methods. Our tool provides a complete and relevant set of metrics to assess the read quality improvement after correction and scales to large datasets. ELECTOR is directly compatible with a wide range of state-of-the-art error correction tools, using whether simulated or real long reads. We show that ELECTOR displays a wider range of metrics than the state-of-the-art tool, LRCstats, and additionally importantly decreases the runtime needed for assessment on all the studied datasets.

Availability: ELECTOR is available at <https://github.com/kamimrcht/ELECTOR>.

Contact: camille.marchet@univ-lille.fr or pierre.morisse2@univ-rouen.fr

1 Introduction

1.1 Motivations

Pacific Biosciences (PB) and Oxford Nanopore Technologies (ONT) long reads, despite their high error rates and complex error profiles, were rapidly adopted for various applications. An increasing number of projects, especially for assembly or structural variant calling [1], indeed benefit from the long range information these reads provide. These reads display high error rates (from 9% to as much as 30%, according to technologies and libraries), that largely surpass those of Illumina reads. Given these high error rates, the first step of many applications is error correction. However, this stage can be a time bottleneck [1].

Moreover, contrary to Illumina, where the majority of errors are substitutions, long reads mainly contain insertions and deletions (indels) errors (ONT reads are more deletion-prone whereas PB reads contain more insertions). This combination of issues requires novel and specific algorithmic developments. To this extent, dozens of error correction methods directly targeting these long reads emerged in the last five years. A first range of error correction tools, called hybrid correctors, uses both short and long reads to perform error correction, relying on the important coverage and low error rate of the short reads in order to enhance long reads sequences. A second group of methods, called self-correctors, intends to correct long reads with the sole information contained in their sequences (see [2] for a review of correctors). Both paradigms include quite diverse algorithmic solutions, which makes it difficult to globally compare the correction results (in terms of throughput, quality and performances) without a proper benchmark.

In addition, the quality of the error correction has considerable impacts on downstream processes. Hence, it is interesting to know beforehand which corrector is best suited for a particular experimental design (coverage, read type, or genome, for instance). Developing methods allowing to evaluate error correction tools with precise and reliable statistics is therefore a crucial need.

Such evaluation methods should allow to perform reproducible and comprehensive benchmarks, and thus to efficiently identify which error correction method is best suited for a given case. They must be usable on datasets of various complexity (from bacteria to eukaryotes) in order to reproduce a wide variety of the scenarios that can be encountered. They also should be fast and lightweight, and should not be orders of magnitude more resource and time consuming than the actual correction methods they assess. This aspect is particularly critical when correction evaluators also stand in the perspective of new correction methods developments. They can help providing accurate and quick comparisons with state-of-the-art correctors. For developers as well as users, correction evaluators should describe with precision the correction method's behavior (i.e. quantity of corrected bases, introduced errors or read break ups, and throughput), in order to identify its potential pitfalls.

1.2 Previous works

Works introducing novel correction methods usually evaluate the quality of their tools based on how well the corrected long reads realign to the reference. Despite being interesting, this information remains incomplete. In particular, it is likely not to mention poor quality reads, or regions to which it is difficult to align. Inspired by earlier works by [4] and [5], La et al. introduced a new way to obtain metrics describing the quality of the error correction itself [6], that does not solely present the similarity between the aligned corrected reads and the reference genome. Relying on simulated data, they proposed the idea of a three way alignment between the reference genome, the uncorrected reads, and the corrected reads. They presented results on Pacific Biosciences data for hybrid error correction tools, by introducing LRCstats, an evaluation tool aiming at answering to the aforementioned problematics. With its three way alignment scheme, LRCstats provides reads' error rate before and after correction, as well as the detailed count of every type of error. However, only studying the

reads' error rate after correction is not a satisfying indication of the corrector's behavior. For instance, there is no clue about the putative insertions of new errors by the corrector, because its precision is not assessed. To overcome this issue, additional metrics such as precision (relevant corrected bases among all bases changed by the corrector), but also recall (correct bases that have been retrieved by the corrector among all bases to be corrected) should be given, in order to better understand the correction methods' pros and cons.

Moreover, LRCstats suffers from high resource consumption when processing large number of reads, i.e. when coverage or genome size are large. However, deep coverage is expected to help the correction of very long sequences [1]. Thus, the correction of such datasets must be assessed in a reasonable amount of time. Additionally, LRCstats's alignment scheme becomes limited when sequences to process grow longer. However, extremely long reads start to appear in recent works for larger genomes [7], and require correction as well.

1.3 Contribution

In order to cope with the identified limits of LRCstats, we propose ELECTOR, a new evaluation tool for long read error correction methods. ELECTOR can be used on simulated as well as real long read datasets, provided a reference genome is available for the sequenced species. It takes as input a reference genome in FASTA format, a set of corrected reads in FASTA format, and the corresponding uncorrected reads, either via a FASTA format file in the case of real data, or via the suite of files provided by the simulator in case of simulated data. ELECTOR provides a wider range of metrics than LRCstats, that assess the actual quality of the correction, such as recall, precision, and correct bases rate for each read. Such metrics have already been proposed in earlier works dedicated to short reads, such as ECTools [4]. However, ECTools' contribution is out of the scope of this work since algorithms to process short reads are different from those at stake in our case, due to the long reads high error rates and complex error profiles. ELECTOR also informs about typical difficulties long read correctors can encounter, such as homopolymers, and reads that have been trimmed, split or extended during the correction. In order to provide these additional metrics, the three way alignment paradigm used in LRCstats is replaced by a scalable multiple sequence alignment

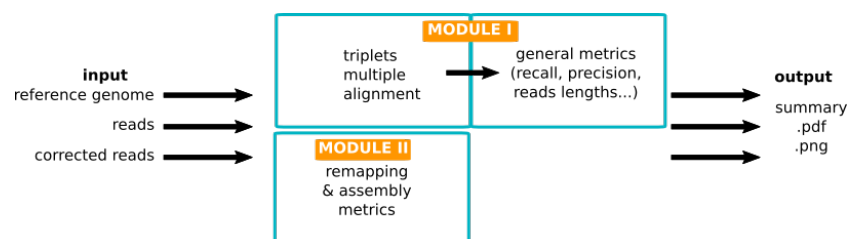


Figure 1: Overview of the ELECTOR’s pipeline. Input are the sequences at the different stages: without errors (from the reference genome), with errors (simulated or real reads) and corrected (after running a correction method). In a first module, a multiple sequence alignment of the three versions of each sequence is computed

, and the results are analyzed to provide correction quality measures. In a second module, reads are assembled using Minimap2 and Miniasm [3], and both the reads and the contigs are aligned to the reference genome, to provide remapping and assembly statistics. A text summary, plots and a pdf summary are output.

(MSA) in ELECTOR. This allows to compare three different versions of each read: the *uncorrected* version, as provided by the sequencing experiment or by the read simulator, the *corrected* version, as provided by the error correction method, and the *reference* version, which is a portion of the reference genome, representing a perfect version of the original read, on which no error would have been introduced. In addition, ELECTOR also performs and evaluates reads remapping and assembly, which were not assessed by the LRCstats pipeline.

In order to allow the multiple sequence alignment strategy to scale to ultra-long reads, that can reach lengths up to 1 million bp, and to large datasets, of several billion of bp, we also propose a novel heuristic that combines anchoring and partial order alignment. This way, we also propose a faster and more scalable evaluation pipeline than LRCstats.

For simulated reads, it is compatible with state-of-the-art long reads simulation tools, such as Nanosim [8] or SimLord [9], on which introduced errors are precisely known. Moreover, ELECTOR is meant to be a user friendly tool, that delivers its results through different output formats, such as graphics that can be directly integrated to the users’ projects. This tool was designed to be directly compatible with a wide range of state-of-the-art error correction tools, without requiring any preprocessing by the user. In particular, ELECTOR is compatible with the latest self-correction methods, and we thus present novel results on such tools, that were not tackled by LRCstats.

2 Material and methods

2.1 Input sequences

ELECTOR is implemented as a pipeline that is divided in two modules. An overview is shown in Figure 1. Input sequences are passed to the two modules independently. The full evaluation pipeline was initially designed for simulated long reads. This choice was motivated by the need to know the *reference* sequences (which are portions of the reference genome, representing perfect versions of the original reads, on which no error would have been introduced) in order to precisely control the results brought by the assessed correction method.

Our pipeline is compatible with long reads simulators SimLoRD and NanoSim. This means that when using long reads simulated with one of these two tools, the reference sequences are directly retrieved by ELECTOR, by parsing the files generated during the simulation. By using these state-of-the-art long reads simulation tools, we ensure to take as input sequences that closely simulate the actual characteristics of the long reads. However, other long reads simulation tools can also be used. In this case, the user must provide the *reference* sequences to ELECTOR itself. Further configuration of the simulation tools such as the error rate, or the long reads coverage, is the user’s call and has no impact on the ELECTOR pipeline. The genome used for the simulation, the files generated by the simulator, and the corrected reads, output by the desired correction method, are then provided as an input to our pipeline. For hybrid correction methods, whether the short reads are real

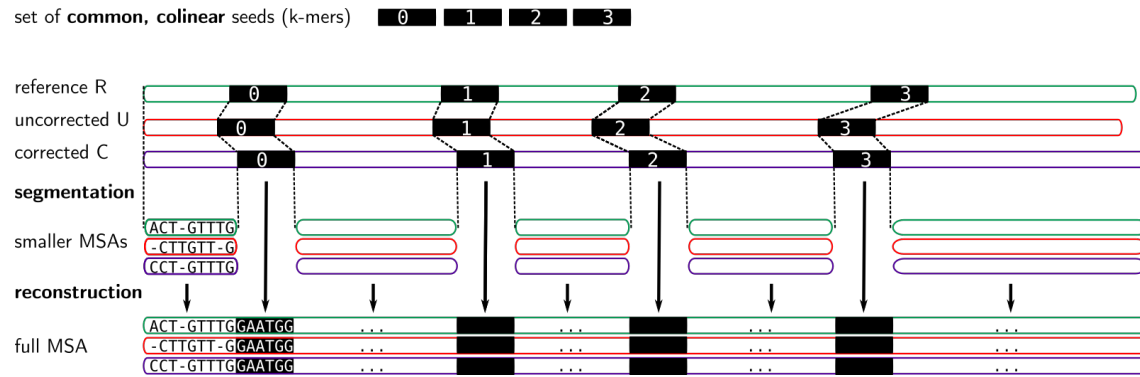


Figure 2: **Segmentation strategy to compute a multiple sequence alignment for a triplet of reference, uncorrected and corrected versions of a read.** Instead of computing a multiple alignment on the whole lengths of the sequences, we rather divide this problem into smaller multiple sequence alignments (MSA). As each version is different, in order to decide where to start and end the alignments, we find seed k -mers (in black) that are local exact matches between the three sequences. We thus compute local, separate MSAs, for subsequences bordered by seeds (or located at the extremities of the sequences). These multiple MSAs are then concatenated, along with the seed k -mers, in order to obtain a single, full MSA, of the whole length of the sequences.

or simulated has no impact on ELECTOR.

In the case of real data, the reads are also passed along with their corrected versions and with the genome to our pipeline. The *reference* sequences are then retrieved by aligning the *uncorrected* reads to the reference genome, using Minimap2 [10]. Only the best hit for each read is kept, and used to determine the corresponding reference sequence. In the case a read cannot align to the reference genome, and thus cannot produce a reference, the read is excluded from the analysis. Apart from that, the rest of the pipeline remains the same, although, as the alignment can soft-clip the reads' extremities, the computation of some metrics slightly vary. These differences are further detailed in Section 2.3.1. We then propose to compare the three different versions of each read (*reference*, *uncorrected*, and *corrected*) in a triplet multiple alignment. These three versions of each read undergo a multiple sequence alignment, in order to collect their differences and similarities at each position of the alignment.

2.2 Scalable triplet multiple alignment

2.2.1 Principle

For each of the three versions of a read, the triplet multiple alignment module computes a multiple sequence alignment (MSA) using a partial order align-

ment algorithm, starting with the *reference* sequence, then adding the *corrected* version, and finally the *uncorrected* version. This step yields a multiple alignment matrix that is output in pseudo FASTA (PIR) format for each triplet. The triplet multiple alignment is computed using an implementation of partial order alignment graphs [11]. Partial order alignment graphs are used as structures containing the information of the multiple aligned sequences. In this method a directed acyclic graph (DAG) contains the previous multiple sequence alignment result. Successive nucleotides from the sequences are stored in vertices, and each new sequence is aligned to this DAG in a generalization of the Needleman-Wunsch algorithm. Paths in the graph represent the successive alignments.

However, such a procedure can be time-consuming when applied to noisy long reads. Thus, we propose a novel multiple sequence alignment heuristic, that we implemented for ELECTOR's purpose, and describe below. Although initially developed for ELECTOR, this multiple sequence alignment heuristic presents an interest that goes beyond the scope of this work. In particular, it would be interesting to generalize it, in order to perform scalable multiple sequence alignment for a larger number of long, noisy sequences.

2.2.2 Segmentation strategy for the MSA

Very long reads induce long MSA running times, that rise according to their length. Moreover, they also imply more errors and branches in the graph, which further increases the computation duration.

In order to reduce the time footprint of our approach, we propose a segmentation strategy. It consists in dividing the triplet multiple alignment into several smaller multiple sequence alignments. Drawing inspiration from MUMmer's [12] and Minimap's [3] longest increasing subsequence approaches, we divide the multiple alignment problem into instances of short windows separated by regions of exact matches. See Figure 2 for an example. If we were able to bound the size of the windows we could guarantee an asymptotic time linear to the read length. In practice our implementation can produce large windows, but we observe a running time almost linear in the size of the reads, as shown in our experimental results, in Section 3

The windows are computed as follows. For each triplet, we compute seed k -mers that have the following properties: 1-they appear in each of the three versions of the sequence, 2- they are not repeated across any of the versions of the sequence, 3-they are not overlapping in any of the versions of the sequence. Using dynamic programming, the longest seed k -mers subsequence \mathcal{S} common to the three sequences is computed. Pairs of successive seed k -mers from \mathcal{S} delineate windows. Thus, we align the subsequences triplets from windows independently, as described in the previous paragraph, using subsequently smaller alignment matrices. Then, the multiple small MSAs are concatenated, along with the seed k -mers, to obtain a single MSA of the whole length of the triplet.

The size of these seed k -mers is adapted according to the current observed error rates [3, 13], *i.e.* 9 to 15 nucleotides. As it is difficult to *a priori* set a k -mer size, we designed a quick iterative strategy that tries several values of k , in order to choose the most suitable for a given triplet.

To avoid computing metrics on poorly corrected reads we filter out corrected reads which length is below 1% of the *reference* length (l being a parameter set to 10 by default) or reads for which an insufficient number of seeds k -mers were found. These two types of filtered reads are tagged and reported apart in ELECTOR's summary to inform the user about their numbers.

2.2.3 Handle reads of different sizes in the segmentation strategy

In the case of a trimmed/split read, the *corrected* version is shortened in comparison to the two other versions, and a part of the *reference* version is thus missing in the *corrected* version. A prefix and/or a suffix of the *reference* can be missing depending on the case. Different scenarios are outlined in Figure 4. In the case of a missing prefix, the first window selected by the segmentation strategy will contain a very large prefix from the *reference* and *uncorrected* versions and a very small sequence from the *corrected* version. This is due to the fact that we only use anchors shared among the three sequences. As the *corrected* subsequence is substantially shorter than the *uncorrected* and *reference* subsequences, it would be irrelevant to compute a MSA between those three sequences. Furthermore computing a MSA on two very large sequences is extremely expensive. To cope with this problem, we detect such cases by checking the sequences length. If we detect a very large first window and *reference* and *uncorrected* sequences are longer than *corrected*, we use a segmentation scheme only with k -mers from *reference* and *uncorrected*, and only align these two prefixes. This way, we are able to efficiently compute a MSA when the corrected reads do not cover all the original region, avoiding to run a MSA on large/unrelated sequences. The procedure is symmetrical for a missing suffix in the *corrected* sequence (see Figure 4 for an example with a suffix). This procedure is extremely important for correctors that output numerous split reads, which would induce extremely long runtime due to large sequence MSA computations described before.

2.3 Inference of quality assessment metrics from MSA

2.3.1 Classification of corrected reads

We report different categories of *corrected* reads in ELECTOR.

“Regular” reads are neither trimmed/split nor extended. Figure 4 shows how we deduce trimmed/split/extended categories from the MSA result.

Split/trimmed reads (two first scenarios in Figure 4). They are reads composed of fragments that come from a single original read that could only be corrected on one or several distinct parts that are

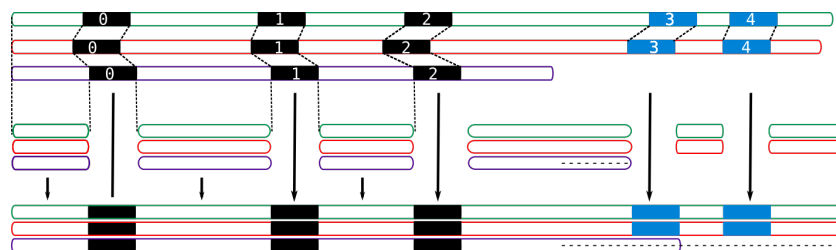


Figure 3: **Segmentation strategy when the corrected read is smaller.** The *corrected* read is shortened on its right end. In order to avoid passing subsequences starting from seed 2 to the end of each sequence to the MSA module, which would be costly to compute, we perform a second segmentation strategy. This allows us to retrieve a new set of seeds (gray seeds 3 and 4). This new set of seeds divides the remaining subsequences (suffixes in this case) in *reference* and *uncorrected* into windows on which we compute MSA separately. The full MSA is reconstructed by concatenation, and dots are added to complete the *corrected* MSA line.

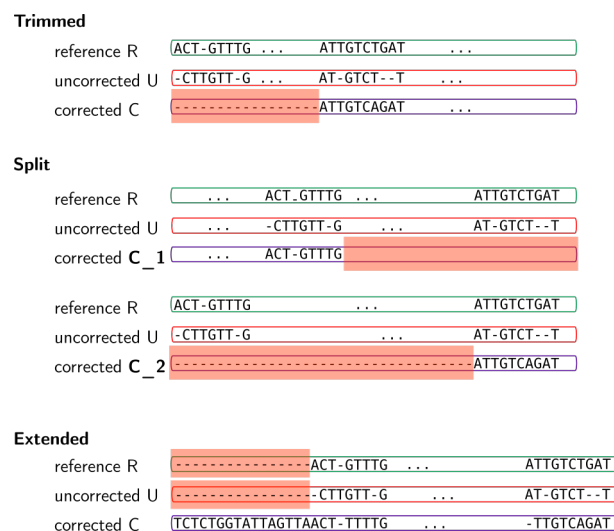


Figure 4: **Three scenarios of corrected read categories in MSA result.** Trimmed/split reads have a corrected version with missing left or right part. Extended corrected reads have a corrected version with a longer left or right part, which is not present in the two other versions.

reported apart. We collect all fragments that come from a single initial read and report how many reads were split. For each trimmed or split corrected read, we report the total uncorrected length of its associated *reference* read (*i.e.* the length that is covered by no fragment).

Extended reads are reads that have a subsequence at their left and/or right end that was not present in the *reference* sequence (last scenario in Figure 4).

These reads can be chimeras from the correction step. However they can also be reads that were over-corrected by a graph-based correction method, that kept on traversing the graph after reaching the *uncorrected* reads' extremities. We do not compute quality assessment metrics on the extended regions, but we report the number of extended reads, as well as their mean extension size, with respect to the *reference* reads.

Soft-clipped reads are reads from a real dataset for which the extremities were soft clipped during the alignment to the reference genome. This category can only arise when processing real data, as we only retrieve *reference* reads by aligning the *uncorrected* reads to the reference genome in this case. For such reads, we do not compute quality assessment metrics on the soft clipped regions, as they could not be properly aligned to the reference genome, and were therefore not used to determine the *reference* read.

Bad quality reads are reads which quality is so bad they were removed before the MSA step. As mentioned before, these are the reads for which an insufficient number of seed k -mers were found. We only report their number as no metric can be computed, since they are not aligned.

2.3.2 Recall, precision, error rate

Once the MSA is computed, we have a base-wise information of the differences and similarities in nucleotide content for each of the three versions of a sequence. Insertions or deletions are represented by

a "." in the deleted parts, and by the corresponding nucleotide (A,C,T or G) in the inserted parts. Let us denote $nt(R, p_i)$, $nt(C, p_i)$, $nt(U, p_i)$ the characters of *reference*, *corrected* and *uncorrected* versions in $\{A, C, G, T, .\}$, at position p_i ($0 \leq i < N$), in a MSA of size N . Figure 5 shows how recall and precision are computed. The set of positions to correct \mathcal{P} contains positions p_i such as $nt(R, p_i) \neq nt(U, p_i)$. The set of existing position in the corrected version \mathcal{E} is defined by including any position p_x from the *corrected* version that is not counted in a trimmed/split/extended region. The processed positions set \mathcal{C} is defined as $\mathcal{P} \cup \{p_j/nt(C, p_j) \neq nt(R, p_j)\} \cap \mathcal{E}$. The correct positions set \mathcal{Co} is defined as $\mathcal{C} \cap \{p_j/nt(C, p_j) = nt(R, p_j)\}$. The recall, precision and error rate are computed as such:

$$Recall = \frac{card(\mathcal{C} \cap \mathcal{P})}{card(\mathcal{P})} \quad (1)$$

$$Precision = \frac{card(\mathcal{Co} \cap \mathcal{C})}{card(\mathcal{C})} \quad (2)$$

$$Error\ rate = 1 - \frac{card(\mathcal{Co})}{\sum_{i=0}^{c-1} i} \quad (3)$$

with c the length of the corrected read.

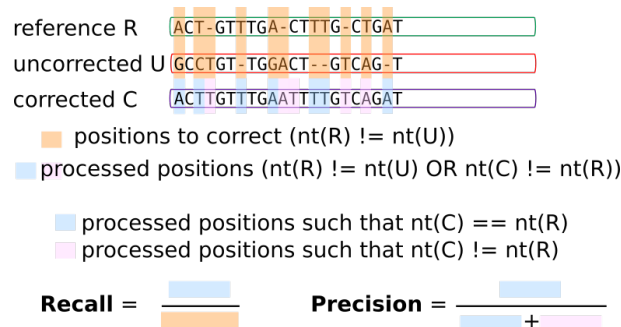


Figure 5: Compute recall and precision using triple base-wise comparison at each MSA's position. $nt(R)$ (respectively $nt(U)$, $nt(C)$) represents the character in *reference* (respectively *uncorrected*, *corrected*) line of the MSA at a given position.

2.3.3 Additional metrics

ELECTOR's first module also provides information on the number of split or trimmed corrected reads, and on the mean missing size of these reads, as well as on the number of extended reads, and on the

mean extension size of these reads. The size distribution of sequences before and after correction is presented graphically. In the case of split reads, we report the length of each split fragment in the distribution. The %GC is also output, as well as the insertion/deletion/substitution counts, before and after correction. Oxford Nanopore reads are known to be more error-prone than PacBio reads in homopolymers. Thus, we propose metrics to examine these particular regions. We show the ratio of homopolymer sizes in the *corrected* version over the *reference* version. The closer it is to one, the better the corrector overcame possible non-systematic errors in ONT reads. More details on the computation of these metrics are shown in Supplementary Material (Supplementary Figure 3).

2.3.4 Remapping of corrected reads

We perform remapping of the corrected reads to the reference genome using BWA-MEM [14], due to the high quality of the long reads after error correction. We report the percentage of aligned reads, the average identity of the alignments, as well as the genome coverage, *i.e.* the percentage of bases of the reference genome to which at least a nucleotide aligned.

2.3.5 Post-correction assembly metrics

We perform the assembly of the corrected reads using Miniasm [3], as we mainly seek to develop a pipeline providing fast results. We acknowledge that assemblers such as Smartdenovo [15] or Canu [16] are more sensitive, but as they display much larger runtimes, Miniasm provides a satisfying compromise.

As for the metrics of the assembly, we output the overall number of contigs, the number of contigs that could be aligned, the number of breakpoints of the aligned contigs, and the NGA50 and NGA75 sizes of the aligned contigs. The alignment of the contigs is also performed with BWA-MEM [14], and the computation of the different metrics is performed by parsing the generated SAM file.

3 Results

3.1 Validation of segmentation strategy for MSA

In order to validate our segmentation strategy for MSA, we show to which extent its results differ from

the classic MSA approach. We expect that recall, precision and correct base rate hardly differ, thus showing that both behaviors produce very similar results. Conversely, we expect an important gain in time with our segmentation strategy compared to the original algorithm. We thus compared multiple alignment results obtained with our strategy to results obtained with the regular implementation of the partial order alignment on multiple datasets of different read lengths, which affects the runtime of the alignments. Results are presented in Table 1.

They show that our segmentation strategy and the regular approach only differ by a few digits in the presented metrics for all the experiments. However, using segmentation, a substantial gain in time is achieved. On the largest dataset, the runtime is hence reduced from several days to roughly an hour. Moreover, while the classic MSA strategy runtime raises with respect to the read length, our approach no longer suffers from this drawback.

3.2 Validation on synthetic datasets

3.2.1 Datasets

We present results of ELECTOR on several simulated datasets of several species (*A. baylyi*, *E. coli*, *S. cerevisiae*, *C. elegans*). Further details on each dataset are given in Supplementary Table 1.

3.2.2 ELECTOR general results

We display ELECTOR’s results using reads corrected by the following tools: HALC [17], HG-CoLoR [18], LoRDEC [19], Canu, Daccord [20], MECAT [21] and LoRMA [22]. As previously detailed in Section 2.3.2, the first module of ELECTOR computes general metrics: mean recall, precision, correct base rate and other additional metrics. For the three first results, a graphic representation of their distribution is also made available.

A subset of the results output by this first module of ELECTOR are presented, for the six correction methods, and for the *E.coli*, *S. cerevisiae*, and *C. elegans* datasets, in Table 3. The complete results, including all the metrics output by ELECTOR are presented in Supplementary Material Tables. These results show that ELECTOR reported recalls from 95 to almost 100% and precisions from 94 to more than 99% for all tools. These results are consistent with results presented in the different tools’ publications

on datasets from the same species. Others metrics output by ELECTOR’s first module are further detailed in Section 3.2.3.

The second module of ELECTOR performs remapping of the reads on the reference and assembles them, thus providing additional information that are not available through LRCstats. An output example of this module is given in Table 6.

3.2.3 Comparison to state-of-the-art

Other works Recently, several benchmark analysis were proposed for long reads (comparison of hybrid correction methods [23], comparison of hybrid and self correction methods [24], analysis of long read correction on transcriptomic reads [25]). In this work, we rather focus on the methodological basis to efficiently perform and reproduce such benchmarks, than highlight pros and cons of correction methods. Though, we present corrector’s result in order to illustrate the validity of our approach. The presented results are in accordance with those reported in the other works. In the rest of the result section, we report comparisons to the only other automated assessment tool for long reads correction: LRCstats.

Comparison of the metrics displayed by ELECTOR and LRCstats are shown in Table 2, for a hybrid and a self-correction method, respectively HALC and Canu, on the *S. cerevisiae* dataset. This table shows the results obtained for each tool, and the supplementary metrics offered by ELECTOR that are not displayed by LRCstats. The complete results provided by LRCstats and ELECTOR, for each correction tool, and on each dataset, are presented in Supplementary Material Tables 2-3.

Common metrics Both LRCstats and ELECTOR compute metrics for *uncorrected* and *corrected* versions of the reads. The first result to notice is that the error rate announced in *uncorrected* sequences can differ from one correction method to the other, both for ELECTOR and LRCstats. This is explained by the fact that HALC and Canu do not correct the same set of reads. As a result, the corresponding *uncorrected* reads used to compute the error rates are not the same either.

Second, as it can be seen from the throughput metrics, ELECTOR and LRCstats do not process the same quantity of reads. This is due to the fact that they rely on different rules to exclude reads that are too difficult to process in the alignment schemes, but also because of their respective behaviors toward split

Experiment	Recall (%)	Precision (%)	Correct bases (%)	Time
"1k" MSA	99.712	98.996	98.980	2h 05 min
"1k" segmentation +MSA	99.769	98.992	98.979	28 min
"10k" MSA	99.921	99.781	99.794	20 h 50 min
"10k" segmentation + MSA	99.921	99.795	99.793	29 min
"100k" MSA	99.913	99.925	99.956	8 days 18 h 38 min
"100k" segmentation +MSA	99.924	99.903	99.902	1 h 11 min

Table 1: Comparison of the two multiple alignment strategies on a simulated *E. coli* datasets. The reads were simulated with a given length, a 10% error rate and a coverage of 100x. The reads were corrected with Canu with default parameters.

reads, which are not aligned and reported the same way. LRCstats concatenates the different parts of a split read before aligning the concatenation, even if a missing zone can exist between two fragments. This behavior can complicate the alignment task and introduce a bias in the output metrics. On the contrary, ELECTOR processes the different fragments separately before reconstituting the whole alignment, and takes into account missing parts. These differences thus have an impact on the metrics displayed for corrected reads. ELECTOR's throughput is a little smaller than LRCstats', however reads uncaptured in the throughput are directly reported in precise categories in ELECTOR (very short reads and low quality reads), while they are lost in LRCstats' output.

Different alignment strategies in both tools also have impacts on the results, which explains the differences seen in indels and substitution counts. However, ELECTOR and LRCstats globally report the same trends of two successful corrections that decreased the error rates.

ELECTOR exclusive metrics ELECTOR's novel metrics point out important differences between the two correction methods, such as the high quantity of trimmed and/or split reads when using HALC in comparison to Canu.

3.3 Performances comparison

We compared LRCstats and ELECTOR's runtime and memory consumption on several datasets in Table 4 and Table 5. The datasets are chosen to represent different factors of pitfalls. We thus vary genome sizes and read throughput, with datasets presented in Supplementary Table 1, as well as reads size distribution, with new *E. coli* datasets composed of reads of length 1 kbp, 10 kbp, 100 kbp, and 1 Mbp. By assessing ELECTOR on reads of such length, we mainly wish to demonstrate its scalability. ELECTOR's runtime and memory peaks are computed for the two in-

Metric	Uncorrected		Corrected by HALC	
	ELECTOR	LRCstats	ELECTOR	LRCstats
Throughput	238,309,333	237,655,341	212,266,193	214,152,119
Error Rate	0.1403	0.1751	0.0042	0.0023
Insertions	28,772,841	32,589,970	100,874	215,507
Deletions	5,235,890	8,991,984	1,035,978	120,743
Substitutions	4,058,953	1,633,123	198,853	221,646
Recall	-	-	0.9997	✖
Precision	-	-	0.9959	✖
Trimmed/split	-	-	12,043	✖
Mean missing size	-	-	577.5	✖
Extension	-	-	71	✖
Mean extension size	-	-	53.2	✖
Low quality	-	-	160	✖
Small reads	-	-	3436	✖
%GC	38.2	✖	38.2	✖
Metric	Uncorrected		Corrected by Canu	
	ELECTOR	LRCstats	ELECTOR	LRCstats
Throughput	244,560,743	244,633,066	229,555,492	229,825,812
Error Rate	0.1425	0.1781	0.0506	0.0694
Insertions	30,090,583	34,105,075	12,252,413	12,942,568
Deletions	5,483,119	9,489,618	2,574,320	3,134,365
Substitutions	4,375,017	1,748,302	2,197,172	1,591,650
Recall	-	-	0.9515	✖
Precision	-	-	0.9495	✖
Trimmed/split	-	-	2,216	✖
Mean missing size	-	-	35.1	✖
Extension	-	-	178	✖
Mean extension size	-	-	30.7	✖
Low quality	-	-	43.0	✖
Small reads	-	-	0.0	✖
%GC	38.2	✖	38.7	✖

Table 2: Comparative results of ELECTOR's and LRCStats' outputs on the *S. cerevisiae* dataset using a hybrid corrector (HALC) and a self corrector (Canu). A dash in the Uncorrected columns indicates that the metric is not computed for the *uncorrected* reads. A cross indicates that the assessment tool does not provide the metric.

Metric	HALC		HG-CoLoR		LoRDEC		CANU		Daccord		MECAT	
	uncorrected	corrected	uncorrected	corrected	uncorrected	corrected	uncorrected	corrected	uncorrected	corrected	uncorrected	corrected
<i>E. coli</i>												
Throughput	91,950,978	81,199,351	93,003,632	84,089,814	89,077,682	77,969,503	91,933,413	86,443,218	92,936,636	83,773,362	80,380,557	58,979,203
Error Rate	0.1415	0.0015	0.1428	0.0007	0.1384	0.0015	0.1432	0.0524	0.1433	0.004	0.1332	0.0052
Recall	-	0.9999	-	1.0	-	0.9999	-	0.9495	-	0.9988	-	0.9983
Precision	-	0.9985	-	0.9993	-	0.9986	-	0.9476	-	0.9961	-	0.9949
<i>S. cerevisiae</i>												
Throughput	238,309,333	212,266,193	245,700,616	219,744,436	196,676,910	188,228,237	244,560,743	229,555,492	246,455,883	222,050,951	217,284,712	162,057,920
Error Rate	0.1403	0.0042	0.1414	0.003	0.1325	0.0054	0.1425	0.0506	0.1426	0.0054	0.1339	0.0066
Recall	-	0.9997	-	0.9999	-	0.9995	-	0.9515	-	0.9986	-	0.998
Precision	-	0.9959	-	0.9971	-	0.9947	-	0.9495	-	0.9946	-	0.9936
<i>C. elegans</i>												
Throughput	1,731,103,921	1,588,220,052	1,988,381,391	1,726,223,265	1,299,187,175	1,154,508,245	1,997,798,872	1,873,188,109	-	-	1,270,739,795	870,965,775
Error Rate	0.1377	0.0153	0.1397	0.0065	0.1242	0.0126	0.1427	0.0496	-	-	0.1199	0.0065
Recall	-	0.9989	-	0.9997	-	0.9989	-	0.9527	-	-	-	0.9982
Precision	-	0.985	-	0.9936	-	0.9875	-	0.9505	-	-	-	0.9936

Table 3: Examples of the main statistics reported by ELECTOR on simulated datasets. A dash in a row/column indicates that the metric is not computable from the data. On the *C. elegans* dataset, Daccord could not be run, and reported an error.

dependent modules of the ELECTOR pipeline: the read triplets multiple alignment step, allowing to access general metrics, and the remapping and assembly step. The first module is comparable to the LRCstats pipeline, as both perform similar operations. However, LRCstats does not include modules to perform read remapping and assembly. We thus present this second module of ELECTOR apart, in Table 6, where we display metrics, runtime and memory consumption for the assessment of the HALC corrected reads on the *S. cerevisiae* dataset.

We first assess in Table 5 the performances of both tools on simulated data with different read length. We observe that the runtime and memory consumption of LRCstats grow with the read length, and we were unable to run it on the 100 kbp and 1 Mbp datasets, even on a cluster node with 250 GB of RAM. This behavior may be a problem to work on very long reads. Due to its segmentation strategy ELECTOR is able to handle larger reads, up to one megabase length. Furthermore its time and memory consumption do not raise much with the read length. This shows that ELECTOR is able to scale to extremely long reads. Considering the availability and usefulness of such very long reads library, we believe that this ability to efficiently handle long sequences is one of the main advantages of ELECTOR.

In the experiment describing the runtimes of ELECTOR and LRCstats, presented in Table 4, we also added, for the perspective, the runtime of the correction method itself. Interestingly we observe that LRCstats is often more time consuming than the correction method, which is not desirable. ELECTOR presents reduced runtime, showing that it could be used to mitigate that heavy analysis time.

3.4 Assessment of real data

3.4.1 Validation of ELECTOR's real data mode

In order to validate the real data mode of ELECTOR, we ran the following experiment. We used a simulated dataset, and assessed its correction using the two different ELECTOR mode, simulated and real data. First, we ran it classically, by providing the simulation files as input, so ELECTOR could retrieve the actual *reference* reads by parsing the files. Second, we ran it by only providing the FASTA file of simulated reads as input, so ELECTOR had to retrieve the *reference* reads by aligning the uncorrected long reads to the reference genome, as if they were actual real long reads. We ran this experiment on the *S. cerevisiae* dataset, and to further validate ELECTOR's behavior on real data, we assessed correction of both a hybrid corrector, HALC, and a self-corrector, Canu.

Results of these experiments are shown in Table 7. We observe that ELECTOR's results in simulated and real mode are consistent. In particular, the recalls and precisions are very similar. The same trend appears as for the comparison to LRCstats: the two modes show some differences in the input uncorrected reads (as shown by the throughputs), that have an impact on the differences observed between their results. This is due to the bias induced by the additional alignment step that is required in the real mode. The main differences that appear in the metrics occur on metrics that are highly dependent on the alignment results, such as the number of trimmed, split and extended reads, and the sizes of these events; as well as indels counts.

Method	HALC	HG-CoLoR	LoRDEC	Canu	Daccord	MECAT
<i>A. baylyi</i>						
Corrector	22min	47min	6min	10min	20min	43sec
LRCstats	3h50	3h52	3h38	3h10	3h59	2h02
ELECTOR	<u>14min</u>	<u>10min</u>	<u>45min</u>	<u>9min</u>	<u>12min</u>	<u>9min</u>
<i>E. coli</i>						
Corrector	24min	45min	8min	12min	27min	52sec
LRCstats	4h58	5h02	4h37	4h05	4h20	2h30
ELECTOR	28min	<u>13min</u>	<u>1h17</u>	<u>11min</u>	<u>12min</u>	<u>11min</u>
<i>S. cerevisiae</i>						
Corrector	1h19	4h32	28min	31min	1h15	2min
LRCstats	10h56	12h26	12h14	10h46	12h04	6h59
ELECTOR	1h55	<u>1h07</u>	<u>4h59</u>	<u>32min</u>	<u>44min</u>	<u>32min</u>
<i>C. elegans</i>						
Corrector	5h59	88h56	6h01	4h33	-	22min
LRCstats	83h29	81h05	70h00	85h08	-	-
ELECTOR	32h35	<u>10h30</u>	<u>29h48</u>	<u>4h19</u>	-	3h12

Table 4: Runtimes of ELECTOR and LRCstats on different datasets and different correctors. Both tools were launched with 9 threads. The runtimes of the correctors are also included as a matter of comparison. The fastest assessment method is shown in bold for each case. When the assessment method is also quicker than the correction method itself, it is underlined. Daccord could be run on the *C. elegans* dataset, and reported an error, and LRCstats crashed on the *C. elegans* dataset corrected by Canu.

Tool	Genome	Read length	Memory (MB)	Elapsed time	CPU time
LRC	<i>E. coli</i>	1k	1,803	42 min	8 h 18 min
ELECTOR	<i>E. coli</i>	1k	1,030	12 min	28 min
LRC	<i>E. coli</i>	10k	13,484	4 h 51 min	2 days 22 h 38 min
ELECTOR	<i>E. coli</i>	10k	3,091	13 min	29 min
ELECTOR	<i>E. coli</i>	100k	12,231	28 min	1 h 11 min
ELECTOR	<i>E. coli</i>	1M	24,881	2 h 44 min	11 h 05 min

Table 5: Performance results obtained from a simulated coverage of 100x of the respective reference genome, on a 20 cores cluster node equipped with 250 GB of RAM.

Remapping	
Aligned reads (%)	100
Average identity (%)	99.54
Genome coverage (%)	99.01
Assembly	
Number of contigs	165
Aligned contigs	165
Number of breakpoints	14
NGA50	94,358
NGA75	55,321
Time	5 min
Memory (MB)	2,759

Table 6: Metrics reported by the remapping and assembly module of ELECTOR, on the *S. cerevisiae* dataset, corrected with HALC.

3.4.2 Results on a real human dataset

In order to demonstrate ELECTOR’s results in a realistic scenario for large genomes, we assess the correction of a real human dataset. We report results, as well as runtime of the assessment, in Table 8. The reads were corrected with MECAT before running ELECTOR. Using 20 threads, we were able to obtain the results for the 650,771 corrected reads of the dataset in less than 19 hours. Results reported by ELECTOR show that MECAT is able to correct human reads with a 20% error rate with more than 90% of recall and precision.

4 Discussion and perspectives

We described and demonstrated ELECTOR's heuristics for multiple sequence alignment and its important speedup in comparison to LRCstats. While showing same trends in the results they display for correctors, we have pointed that LRCstats and ELECTOR have differences in their common metrics. This is explained by a set of different choices and heuristics between the two tools. Regarding LRCstats, we fulfilled our speed-up and scalability objective, being able to report results faster than LRCstats, and better scaling to both very long reads and large datasets.

ELECTOR is designed to work with simulated reads, and is expected to give the best results when working with it. The real data mode uses a prior alignment of the reads on a reference genome to retrieve the *reference* versions of the reads. We demonstrated that ELECTOR's main metrics in its real data mode remain similar to what would be obtained in the simulated mode. However, we can point out two limitations of ELECTOR: first, even if the data can come from an actual sequencing experiment, a reference genome needs to exist for the sequenced species in order to retrieve the *reference* reads, and thus perform evaluation. Second, we encourage users to be very cautious about ELECTOR's results with real data when looking at trimmed/split reads numbers and their sizes, since these metrics highly depend on the results of the alignment to the reference.

In ELECTOR we propose an efficient segmentation heuristic for multiple sequence alignment. We adapted this task for the original and specific long read application. There is a current interest for segmented multiple alignment scheme [26]. However, these methods are not specifically designed for noisy long reads. In such a perspective, a generalization of our segmentation strategy for long reads multiple alignments of any size would be very interesting.

A future application is the assessment of correction methods directly targeted at RNA long reads sequencing. As shown in a recent study [25], RNA long reads have specific requirements that are not met by current methods, which calls for new correctors in the future. ELECTOR could be coupled with a reference transcriptome or a RNA long read simulator, although, currently, such simulation software does not exist to our knowledge.

Metric	Uncorrected		Corrected by Halc	
	Simulated	Real	Simulated	Real
Throughput	238,309,333	238,119,170	212,266,193	212,141,319
Error Rate	0.1403	0.1449	0.0042	0.0104
Recall	-	-	0.9997	0.9938
Precision	-	-	0.9959	0.9897
Insertions	28,772,841	26,796,500	100,874	90,737
Deletions	5,235,890	5,042,365	1,035,978	1,490,680
Substitutions	4,058,953	3,682,863	198,853	182,590
Trimmed/split	-	-	12,043	13,320
Mean missing size	-	-	577.5	896.0
Extension	-	-	71.0	39.0
Mean extension size	-	-	53.2	72.0
Low quality	-	-	160.0	152.0
Small reads	-	-	3436.0	3438.0
%GC	-	-	38.2	38.2
Metric	Uncorrected		Corrected by Canu	
	Simulated	Real	Simulated	Real
Throughput	244,560,743	244,402,568	229,555,492	229,403,697
Error Rate	0.1425	0.1442	0.0506	0.052
Recall	-	-	0.9515	0.9499
Precision	-	-	0.9495	0.9481
Insertions	30,090,583	28,452,967	12,252,413	10,965,458
Deletions	5,483,119	5,800,286	2,574,320	2,916,564
Substitutions	4,375,017	4,081,445	2,197,172	1,940,888
Trimmed/split	-	-	2216.0	4943.0
Mean missing size	-	-	35.1	74.7
Extension	-	-	178.0	169.0
Mean extension size	-	-	30.7	31.9
Low quality	-	-	43.0	42.0
Small reads	-	-	0.0	0.0
%GC	-	-	38.7	38.7

Table 7: Comparative results output by ELECTOR, using simulated and real modes on the same *S. cerevisiae* dataset, using a hybrid corrector (HALC) and a self corrector (Canu).

	uncorrected	corrected with MECAT
Throughput	5,605,157,590	5,451,767,836
Recall(%)	-	92.70
Precision(%)	-	91.50
Error rate	0.1974	0.0861
Average correct bases rate	0.8026	0.9139
Number of trimmed/split reads	-	570,635
Mean missing size in trimmed/split reads	-	362.0
Number of over-corrected reads by extension	-	275
Mean extension size in over-corrected reads	-	62.4
%GC	41.6	41.1
Small reads	-	356
Low quality corrected reads	-	4,279
Insertions	247,953,086	10,144,736
Deletions	746,165,024	473,239,036
Substitutions	162,822,923	7,521,389
Homopolymer ratio	-	0.7570
Runtime	-	18 h 27 min

Table 8: ELECTOR’s results for MECAT correction of a real human dataset. ELECTOR assessed 650,771 reads. Small reads are corrected reads which length is lower than 10.0% of the original read. Homopolymer ratio is the ratio of homopolymer sizes in *corrected* vs *reference*. We reported the wallclock time of the run, using 30 threads.

5 Conclusion

We presented ELECTOR, a tool that enables the evaluation of self and hybrid correction methods, and that can be used in the conception of a benchmark. ELECTOR provides a wide variety of metrics that include base-wise computation of recall, precision, error rate of corrected reads as well as indels, substitutions and homopolymers correction. In particular, recall and precision allow to spot correction methods specific pitfalls, that remain unclear when only looking at the error rates of the corrected reads. These results are presented in a text summary and in pdf and png versions, which allows users to easily integrate them in their reports.

We used ELECTOR on a representative list of state-of-the-art hybrid and self correctors, ran on reads from small bacterial genomes to large mammal genomes. With the applications on large genomes and ever increasing lengths of long reads in mind, we designed ELECTOR to be time-saving and scalable. We shown that for most datasets and correctors, ELECTOR runs in a similar amount of time as the corrector itself. We also demonstrated that it can scale to novel ultra long reads.

Thus, it represents a major improvement in comparison to LRCstats. Since ELECTOR is based on multiple sequence alignment, we adapted this strategy to our scaling objectives. We proposed an inno-

vative and promising algorithm of segmentation for multiple sequence alignment of noisy long reads. This procedure drastically reduces the time footprint of the multiple alignment step in ELECTOR. We believe it could be generalized for broad applications implying multiple sequence alignment of long, noisy sequences.

References

- [1] Fritz J Sedlazeck, Hayan Lee, Charlotte A Darby, and Michael C Schatz. Piercing the dark matter: bioinformatics of long-range sequencing and mapping. *Nature Reviews Genetics*, page 1, 2018.
- [2] David Laehnemann, Arndt Borkhardt, and Alice Carolyn McHardy. Denoising DNA deep sequencing data-high-throughput sequencing errors and their correction. *Briefings in bioinformatics*, 17(1):154–179, 2015.
- [3] Heng Li. Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences. *Bioinformatics*, 32(14):2103–2110, 2016.
- [4] Xiao Yang, Sriram P Chockalingam, and Srinivas Aluru. A survey of error-correction methods for next-generation sequencing. *Briefings in bioinformatics*, 14(1):56–66, 2012.

- [5] Giles Miclotte, Mahdi Heydari, Piet Demeester, Stephane Rombauts, Yves Van de Peer, Pieter Audenaert, and Jan Fostier. Jabba: hybrid error correction for long sequencing reads. *Algorithms for Molecular Biology*, 11(1):10, 2016.
- [6] Sean La, Ehsan Haghshenas, and Cedric Chauve. LRCstats, a tool for evaluating long reads correction methods. *Bioinformatics*, 33(22):3652–3654, 2017.
- [7] Miten Jain, Sergey Koren, Karen H Miga, Josh Quick, Arthur C Rand, Thomas A Sasani, John R Tyson, Andrew D Beggs, Alexander T Dilthey, Ian T Fiddes, et al. Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nature biotechnology*, 36(4):338, 2018.
- [8] Chen Yang, Justin Chu, René L Warren, and Inanç Birol. NanoSim: nanopore sequence read simulator based on statistical characterization. *GigaScience*, 6(4):1–6, 2017.
- [9] Bianca K Stöcker, Johannes Köster, and Sven Rahmann. Simlord: Simulation of long read data. *Bioinformatics*, 32(17):2704–2706, 2016.
- [10] Heng Li. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, 1:7, 2018.
- [11] Christopher Lee, Catherine Grasso, and Mark F Sharlow. Multiple sequence alignment using partial order graphs. *Bioinformatics*, 18(3):452–464, 2002.
- [12] Arthur L. Delcher, Steven L Salzberg, and Adam M. Phillippy. Using MUMmer to identify similar regions in large sequence sets. *Current Protocols in Bioinformatics*, Chapter 10, 2 2003.
- [13] Mark J Chaisson and Glenn Tesler. Mapping single molecule sequencing reads using basic local alignment with successive refinement (blasr): application and theory. *BMC bioinformatics*, 13(1):238, 2012.
- [14] Heng Li. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv preprint arXiv:1303.3997*, 2013.
- [15] J Ruan. Smartdenovo: Ultra-fast de novo assembler using long noisy reads, 2017.
- [16] Sergey Koren, Brian P Walenz, Konstantin Berlin, Jason R Miller, Nicholas H Bergman, and Adam M Phillippy. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome research*, pages gr-215087, 2017.
- [17] Ergude Bao and Lingxiao Lan. HALC: High throughput algorithm for long read error correction. *BMC bioinformatics*, 18(1):204, 2017.
- [18] Pierre Morisse, Thierry Lecroq, and Arnaud Lefebvre. Hybrid correction of highly noisy long reads using a variable-order de Bruijn graph. *Bioinformatics*, 34(24):4213–4222, 2018.
- [19] Leena Salmela and Eric Rivals. LoRDEC: accurate and efficient long read error correction. *Bioinformatics*, 30(24):3506–3514, 2014.
- [20] German Tischler and Eugene W Myers. Non hybrid long read consensus using local de Bruijn graph assembly. *bioRxiv*, page 106252, 2017.
- [21] Chuan-Le Xiao, Ying Chen, Shang-Qian Xie, Kai-Ning Chen, Yan Wang, Yue Han, Feng Luo, and Zhi Xie. MECAT: fast mapping, error correction, and de novo assembly for single-molecule sequencing reads. *Nature Methods*, 14(11):1072, 2017.
- [22] Leena Salmela, Riku Walve, Eric Rivals, and Esko Ukkonen. Accurate self-correction of errors in long reads using de Bruijn graphs. *Bioinformatics*, 33(6):799–806, 2016.
- [23] Shuhua Fu, Anqi Wang, and Kin Fai Au. A comparative evaluation of hybrid error correction methods for error-prone long reads. *Genome Biology*, 20(1):26, Feb 2019.
- [24] Haowen Zhang, Chirag Jain, and Srinivas Aluru. A comprehensive evaluation of long read error correction methods. *bioRxiv*, 2019.
- [25] Leandro Ishi Soares de Lima, Camille Marchet, Segolene Caboche, Corinne Da Silva, Benjamin Istace, Jean-Marc Aury, Helene Touzet, and Rayan Chikhi. Comparative assessment of long-read error-correction software applied to rna-sequencing data. *bioRxiv*, page 476622, 2018.
- [26] Edgar Garriga Nogales, Paolo Di Tommaso, Cedrik Magis, Ionas Erb, Hafid Laayouni, Fyodor Kondrashov, Evan Floden, and Cedric

Notredame. Fast and accurate large multiple sequence alignments using root-to-leave regressive computation. *bioRxiv*, page 490235, 2018.