

RESEARCH

Open Access



A comprehensive evaluation of long read error correction methods

Haowen Zhang¹, Chirag Jain¹ and Srinivas Aluru^{1,2*}

From 8th IEEE International Conference on Computational Advances in Bio and medical Sciences (ICCBS 2018) Las Vegas, NV, USA. 18-20 October 2018

Abstract

Background: Third-generation single molecule sequencing technologies can sequence long reads, which is advancing the frontiers of genomics research. However, their high error rates prohibit accurate and efficient downstream analysis. This difficulty has motivated the development of many long read error correction tools, which tackle this problem through sampling redundancy and/or leveraging accurate short reads of the same biological samples. Existing studies to assess these tools use simulated data sets, and are not sufficiently comprehensive in the range of software covered or diversity of evaluation measures used.

Results: In this paper, we present a categorization and review of long read error correction methods, and provide a comprehensive evaluation of the corresponding long read error correction tools. Leveraging recent real sequencing data, we establish benchmark data sets and set up evaluation criteria for a comparative assessment which includes quality of error correction as well as run-time and memory usage. We study how trimming and long read sequencing depth affect error correction in terms of length distribution and genome coverage post-correction, and the impact of error correction performance on an important application of long reads, genome assembly. We provide guidelines for practitioners for choosing among the available error correction tools and identify directions for future research.

Conclusions: Despite the high error rate of long reads, the state-of-the-art correction tools can achieve high correction quality. When short reads are available, the best hybrid methods outperform non-hybrid methods in terms of correction quality and computing resource usage. When choosing tools for use, practitioners are suggested to be careful with a few correction tools that discard reads, and check the effect of error correction tools on downstream analysis. Our evaluation code is available as open-source at <https://github.com/haowenz/LRECE>.

Keywords: Long read, Error correction, Benchmark, Evaluation

Background

Third-generation sequencing technologies produce long reads with average length of 10 Kbp or more that are orders of magnitudes longer than the short reads available through second-generation sequencing technologies (typically a few hundred bp). In fact, the longest read

length reported to date is > 1 million bp [1]. Longer lengths are attractive because they enable disambiguation of repetitive regions in a genome or a set of genomes. The impact of this valuable long-range information has already been demonstrated for *de novo* genome assembly [2–4], novel variant detection [5, 6], RNA-seq analysis [7], metagenomics [8], and epigenetics [9, 10].

The benefit of longer read lengths, however, comes with the major challenge of handling high error rates. Currently, there are two widely used third-generation single molecule sequencing platforms – Pacific Biosciences

*Correspondence: aluru@cc.gatech.edu

¹School of Computational Science and Engineering, Georgia Institute of Technology, Atlanta, GA 30332, USA

²Institute for Data Engineering and Science, Georgia Institute of Technology, Atlanta, GA 30332, USA



© The Author(s). 2020 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

(PacBio) and Oxford Nanopore Technologies (ONT). Both sequencing platforms are similar in terms of their high error rates (ranging from 10–20%) with most errors occurring due to insertions or deletions (indels); however the error distribution varies [4, 11, 12]. Pacbio sequencing errors appear to be randomly distributed over the sequence [13]. For ONT on the other hand, the error profile has been reported to be biased. For example, A to T and T to A substitutions are less frequent than other substitutions, and indels tend to occur in homopolymer regions [12, 14]. These error characteristics pose a challenge for long read data analyses, particularly for detecting correct read overlaps during genome assembly and variants at single base pair resolution, thus motivating the development of error correction methods.

Error correction algorithms are designed to identify and fix or remove sequencing errors, thereby benefiting resequencing or *de novo* sequencing analysis. In addition, the algorithms should be computationally efficient to handle increasing volumes of sequencing data, particularly in the case of large, complex genomes. Numerous error correction methodologies and software have been developed for short reads; we refer readers to [15] and [16] for a thorough review. Given the distinct characteristics of long reads, i.e., significantly higher error rates and lengths, specialized algorithms are needed to correct them. Till date, several error correction tools for long reads have been developed including PacBioToCA [17], LSC [18], ECTools [19], LoRDEC [20], proovread [21], NaS [22], Nanocorr [23], Jabba [24], CoLoRMap [25], LoRMA [26], HALC [27], FLAS [28], FMLRC [29], HG-CoLoR [30] and Hercules [31].

In addition, error correction modules have been developed as part of long read *de novo* assembly pipelines, such as Canu [32] and HGAP [33]. In the assembly pipeline, correction helps by increasing alignment identities of overlapping reads, which facilitates overlap detection and improves assembly. Many long read error correction tools require and make use of highly accurate short reads to correct long reads (accordingly referred to as hybrid methods). Others, referred to as non-hybrid methods, perform self-correction of long reads using overlap information among them.

A few review studies have showcased comparisons among rapidly evolving error correction algorithms to assess state-of-the-art. Laehnemann et al. [34] provide an introduction to error rates/profiles and a methodology overview of a few correction tools for various short and long read sequencing platforms, although no benchmark is included. A review and benchmark for hybrid methods is also available [35]. However, the study only used simulated reads and focused more on speed rather than correction accuracy. Besides, it does not include non-hybrid

methods in the assessment. More recently, LRCstats [36] was developed for evaluation of long read error correction software; however, it is restricted to benchmarking with simulated reads.

While benchmarking with simulated reads is useful, it fails to convey performance in real-world scenarios. Besides the base-level errors (i.e., indels and substitutions), real sequencing data sets also contain larger structural errors, e.g., chimeras [37]. However, state-of-the-art simulators (e.g., SimLoRD [38]) only generate reads with base-level errors rather than structural errors. Furthermore, Miclotte et al. [24] consistently observed worse performance when using real reads instead of simulated reads, suggesting that simulation may fail to match the characteristics of actual error distribution. Therefore, benchmarking with real data sets is important.

A more recent study by Fu et al. [39] was published during the final stage of our work, which like our survey is also aimed to characterize long read error correction performance using real data. In contrast to their review, we include both hybrid and non-hybrid methods rather than just the former. Our benchmark data includes data from the latest generation of long read sequencers (e.g., ONT R9 instead of ONT R7.3). We also shed light on the effect of long read sequencing depth to the overall performance, which is important to evaluate expected sequencing coverage and cost requirements to achieve good correction performance.

In this study, we establish benchmark data sets, present an evaluation methodology suitable to long reads, and carry out comprehensive evaluation of the quality and computational resource requirements of state-of-the-art long read correction software. We also study the effect of trimming and different sequencing depths on correction quality. To understand impact of error correction on downstream analysis, we perform assembly using corrected reads generated by various tools and assess quality of the resulting assemblies.

Overview of long read error correction methods

Hybrid methods

Hybrid methods take advantage of high accuracy of short reads (error rates often $< 1\%$) for correcting errors in long reads. An obvious requirement is that the same biological sample must be sequenced using both short read and long read technologies. Based on how these methods make use of short reads, we further divide them into two categories: *alignment-based* and *assembly-based*. The first category includes Hercules, CoLoRMap, Nanocorr, Nas, proovread, LSC and PacBioToCA, whereas HG-CoLoR, HALC, Jabba, LoRDEC, and ECTools are in the latter. The ideas underlying the methods are summarized below.

Short-read-alignment-based methods

As a first step, these methods align short reads to long reads using a variety of aligners, e.g. BLAST [40], Novoalign (<http://www.novocraft.com/products/novoalign/>). As long reads are usually error-prone, some alignments can be missed or biased. Thus, most of the tools in this category utilize various approaches to increase accuracy of alignments. Drawing upon the alignments, these methods use distinct approaches to generate corrected reads.

PacBioToCA: Consensus sequences for long reads are generated by multiple sequence alignment of short reads using AMOS consensus module [41].

LSC: Short reads and long reads are compressed using homopolymer compression (HC) transformation prior to alignment. Then error correction is performed at HC points, mismatches and indels by temporarily decompressing the aligned short reads and then generating consensus sequences. Finally, the corrected sequences are decompressed.

proovread: Similar to PacBioToCA and LSC, short reads are mapped to long reads, and the resulting alignments are used to call consensus. But its alignment parameters are carefully selected and adapted to the PacBio sequencing error profile. To further improve correction, the phred quality score and Shannon entropy value are calculated at each nucleotide for quality control and chimera detection, respectively. To reduce run time, an iterative correction strategy is employed. Three pre-correction steps are performed using increasing subsamples of short reads. In each step, the long read regions are masked to reduce alignment search space once they are corrected and covered by a sufficient number of short read alignments. In the final step, all short reads are mapped to the unmasked regions to make corrections.

NaS: Like the other tools in this category, it first aligns short reads with long reads. However, only the stringently aligned short reads are found and kept as seed-reads. Then instead of calling consensus, similar short reads are retrieved with these seed-reads. Micro-assemblies of these short reads are performed to generate contigs, which are regarded as corrected reads. In other words, the long reads are only used as template to select seed-reads.

Nanocorr: It follows the same general approach as PacBioToCA and LSC, by aligning short reads to long reads and then calling consensus. But before the consensus step, a dynamic programming algorithm is utilized to select an optimal set of short read alignments that span each long read.

CoLoRMap: CoLoRMap does not directly call consensus. Instead, for each long read region, it runs a shortest

path algorithm to construct a sequence of overlapping short reads aligned to that region with minimum edit distance. Subsequently, the region is corrected by the constructed sequence. In addition, for each uncovered region (called gap) on long reads, any unmapped reads with corresponding mapped mates are retrieved and assembled locally to fill the gap.

Hercules: It first aligns short reads to long reads. Then unlike other tools, Hercules uses a machine learning-based algorithm. It creates a profile Hidden Markov Model (pHMM) template for each long read and then learns posterior transition and emission probabilities. Finally, the pHMM is decoded to get the corrected reads.

Short-read-assembly-based methods

These methods first perform assembly with short reads, e.g., generate contigs using an existent assembler, or only build the de Bruijn graph (DBG) based on them. Then the long reads are aligned to the assemblies, i.e., contigs/unitigs or a path in the DBG, and corrected. Algorithms for different tools in this category are summarized below.

ECTools: First, unitigs are generated from short reads using any available assembler and aligned to long reads. Afterwards, the alignments are filtered to select a set of unitigs which provide the best cover for each long read. Finally, differences in bases between each long read and its corresponding unitigs are identified and corrected.

LoRDEC: Unlike ECTools which generates assemblies, LoRDEC only builds a DBG of short reads. Subsequently, it traverses paths in the DBG to correct erroneous regions within each long read. The regions are replaced by the respective optimal paths which are regarded as the corrected sequence.

Jabba: It adopts a similar strategy as in LoRDEC, and builds a DBG of short reads followed by aligning long reads to the graph to correct them. The improvement is that Jabba employs a seed-and-extend strategy using maximal exact matches (MEMs) as seeds to accelerate the alignment.

HALC: Similar to ECTools, short reads are used to generate contigs as the first step. Unlike other methods which try to avoid ambiguous alignments [17, 42], HALC aligns long reads to the contigs with a relatively low identity requirement, thus allowing long reads to align with their similar repeats which might not be their true genomic origin. Then long reads and contigs are split according to the alignments so that every aligned region on read has its corresponding aligned contig region. A contig graph is constructed with the aligned contig regions as vertices. A weighted edge is added between two vertices if there

are adjacent aligned long read regions supporting it. The more regions support the edge, the lower is the weight assigned to it. Each long read is corrected by the path with minimum total weight in the graph. Furthermore, the corrected long read regions are refined by running LoRDEC, if they are aligned to similar repeats.

FMLRC: This software uses a DBG-based correction strategy similar to LoRDEC. However, the key difference in the algorithm is that it makes two passes of correction using DBGs with different k -mer sizes. The first pass does the majority of correction, while the second pass with a longer k -mer size corrects repetitive regions in the long reads. Note that a straightforward implementation of a DBG does not support dynamic adjustment of k -mer size. As a result, FMLRC uses FM-index to implicitly represent DBGs of arbitrary length k -mers.

HG-CoLoR: Similar to FMLRC, it avoids using a fixed k -mer size for the de Bruijn graph. Accordingly, it relies on a variable-order de Bruijn graph structure [43]. It also uses a seed-and-extend approach to align long reads to the graph. However, the seeds are found by aligning short reads to long reads rather than directly selecting them from the long reads.

Non-hybrid methods

These methods perform self-correction with long reads alone. They all contain a step to generate consensus sequences using overlap information. However, the respective methods vary in how they find the overlaps and generate consensus sequences. The details are as follows.

FLAS: It takes all-to-all long read overlaps computed using MECAT [44] as input, and clusters the reads that are aligned with each other. In case of ambiguous instances, i.e., the clusters that share the same reads, FLAS evaluates the overlaps by computing alignments using sensitive alignment parameters either to augment the clusters or discard the incorrect overlaps. The refined alignments are then used to correct the reads. To achieve better accuracy, it also corrects errors in the uncorrected regions of the long reads. Accordingly, it constructs a string graph using the corrected regions of long reads, and aligns the uncorrected ones to the graph for further correction.

LoRMA: By gradually increasing the k -mer size, LoRMA iteratively constructs DBGs using k -mers from long reads exceeding a specified frequency threshold, and runs LoRDEC to correct errors based on the respective DBGs. After that, a set of reads similar to each read termed *friends* are selected using the final DBG, which should be more accurate due to several rounds of corrections. Then, each read is corrected by the consensus sequence generated by its friends.

Canu error correction module: As a first step during the correction process, Canu computes all-versus-all overlap information among the reads using a modified version of MHAP [45]. It uses a filtering mechanism during the correction to favor true overlaps over the false ones that occur due to repetitive segments in genomes. The filtering heuristic ensures that each read contributes to correction of no more than D other reads, where D is the expected sequencing depth. Finally, a consensus sequence is generated for each read using its best set of overlaps by leveraging “falcon_sense” consensus module [3].

Results and discussion

Experimental setup

All tests were run on the Swarm cluster located at Georgia Institute of Technology. Each compute node in the cluster has dual Intel Xeon CPU E5-2680 v4 (2.40 GHz) processors equipped with a total of 28 cores and 256GB main memory.

Evaluated software

We evaluated 14 long read error correction programs in this study: Hercules, HG-CoLoR, FMLRC, HALC, CoLoRMap, Jabba, Nanocorr, proovread, LoRDEC, ECTools, LSC, FLAS, LoRMA and the error correction module in Canu. NaS was not included in the evaluation because it requires Newbler assembler which is no longer available from 454. PacBioToCA was also excluded since it is deprecated and no longer being maintained. The command line parameters were chosen based on user documentations of each software (Additional file 1 section “Versions and configurations”). The tools were configured to run exclusively on a single compute node and allowed to leverage all the 28 cores if multi-threading is supported. A cutoff on wall time was set to three days.

Performance on benchmark data sets

We evaluated the quality and computational resource requirements of each software on our benchmark data sets (“Methods” section). Results for the different data sets are shown in Tables 1, 2, 3, 4, 5 and 6. Because multiple factors are at play when considering accuracy, it is important to consider their collective influence in assessing quality of error correction. In what follows, we present a detailed and comparative discussion on correction accuracy, runtime and memory-usage. In addition, to guide error correction software users and future developers, we provide further insights into the strengths and limitations of various approaches that underpin the software. This includes evaluating their resilience to handle various sequencing depths, studying the effect of discarding or trimming input reads to gain higher accuracy, and impact on genome assembly.

Table 1 Experimental results for *E. coli* PacBio data set D1-P

Method	# Reads	# Bases (Mbp)	# Aligned reads	# Aligned bases (Mbp)	Maximum length (bp)	N50 (bp)	Genome fraction (%)	Alignment identity(%)	CPU time (hh:mm:ss)	Wall time (hh:mm:ss)	Memory usage (GB)
Original	85460	748.0	82886	688.0	44113	13990	100.000	86.8763	-	-	-
<i>Non-hybrid methods</i>											
FLAS	69327	632.3	68786	621.2	40117	13212	100.000	99.5959	09:47:50	00:56:45	4.9
LoRMA	330811	623.3	330715	623.0	22499	2441	100.000	99.6814	45:24:49	02:10:36	67.2
Canu	9283	168.1	9193	166.7	39693	20391	100.000	99.6970	07:47:33	00:27:14	6.0
<i>Short-read-assembly-based methods</i>											
HG-CoLoR	-	-	-	-	-	-	-	-	-	-	-
FMLRC	85260	706.5	83320	669.9	44084	13364	100.000	99.6983	03:05:06	00:30:07	9.8
HALC	85256	711.1	84030	661.7	44117	13399	100.000	99.4374	60:41:59	16:02:32	30.2
Jabba	77508	620.2	77508	619.7	41342	12557	99.258	99.9624	02:05:09	00:12:01	37.0
LoRDEC	85324	716.9	83507	665.9	44311	13491	100.000	98.4149	15:03:42	00:40:05	2.0
ECTools	55687	577.4	55687	575.7	39772	13583	100.000	99.8592	11:25:22	00:29:49	8.2
<i>Short-read-alignment-based methods</i>											
Hercules	-	-	-	-	-	-	-	-	-	>72:00:00	-
CoLoRMap	85674	730.7	83765	678.6	44113	13641	100.000	95.2930	31:35:16	02:53:33	34.9
Nanocorr	73368	504.9	73316	493.1	41079	10796	100.000	98.3257	1862:59:19	70:57:19	15.1
proofread	85367	720.2	83142	665.7	44113	13524	100.000	96.7250	71:17:14	12:21:53	53.9
LSC	-	-	-	-	-	-	-	-	-	>72:00:00	-

Note: HG-CoLoR reported an error when correcting this dataset

Table 2 Experimental results for *E. coli* ONT data set D1-O

Method	# Reads	# Bases (Mbp)	# Aligned reads	# Aligned bases (Mbp)	Maximum length (bp)	N50 (bp)	Genome fraction (%)	Alignment identity (%)	CPU time (hh:mm:ss)	Wall time (hh:mm:ss)	Memory usage (GB)
Original	163747	1481.5	163386	1454.4	131969	14895	100.000	81.3559	-	-	-
<i>Non-hybrid methods</i>											
FLAS	138472	1401.3	138458	1392.9	130497	14748	99.997	93.0176	20:27:50	01:56:52	8.0
LoRMA	595072	1433.5	595051	1432.5	31743	3333	99.924	96.6525	182:14:17	07:30:30	77.8
Canu	19335	226.2	19326	225.0	133168	38034	99.953	94.5969	17:14:11	00:50:04	6.7
<i>Short-read-assembly-based methods</i>											
HG-CoLoR	159856	1540.7	159854	1518.1	138002	15744	100.000	98.1308	231:20:30	44:41:19	13.8
FMLRC	163749	1555.4	163593	1546.3	137960	15687	100.000	99.6423	05:50:54	00:32:27	3.3
HALC	-	-	-	-	-	-	-	-	-	>72:00:00	-
Jabba	162970	1287.0	162970	1286.1	93923	12795	99.515	99.9557	02:51:05	00:10:33	37.1
LoRDEC	163838	1555.5	163722	1530.1	137887	15664	100.000	98.9920	32:35:27	01:12:37	2.2
ECTools	116868	1431.7	116868	1428.2	137863	16354	100.000	99.8116	19:44:40	00:46:51	8.1
<i>Short-read-alignment-based methods</i>											
Hercules	-	-	-	-	-	-	-	-	-	>72:00:00	-
CoLoRMap	164072	1518.3	163782	1495.7	134302	15180	100.000	89.2049	32:55:26	04:01:18	35.5
Nanocorr	-	-	-	-	-	-	-	-	-	>72:00:00	-
proofread	163815	1514.0	163481	1489.1	135798	15222	100.000	89.2071	104:33:09	18:35:46	47.8
LSC	-	-	-	-	-	-	-	-	-	>72:00:00	-

Table 3 Experimental results for yeast PacBio data set D2-P

Method	# Reads	# Bases (Mbp)	# Aligned reads	# Aligned bases (Mbp)	Maximum length (bp)	N50 (bp)	Genome fraction (%)	Alignment identity (%)	CPU time (hh:mm:ss)	Wall time (hh:mm:ss)	Memory usage (GB)
Original	239408	1462.7	235620	1332.6	35196	8656	99.976	87.2637	-	-	-
<i>Non-hybrid methods</i>											
FLAS	173187	1093.2	173046	1078.8	30046	8132	99.976	99.5777	11:46:31	01:15:40	7.9
LoRMA	650467	1142.0	650333	1141.4	18127	2323	99.951	99.7583	172:24:38	07:03:03	72.9
Canu	38228	453.2	38172	446.7	28748	12021	99.975	99.5864	15:18:34	00:50:12	6.5
<i>Short-read-assembly-based methods</i>											
HG-CoLoR	-	-	-	-	-	-	-	-	-	-	-
FMLRC	238706	1380.8	236883	1311.0	33658	8185	99.977	99.3889	07:52:17	00:28:55	5.5
HALC	238787	1395.4	238097	1287.6	34785	8270	99.976	99.0796	52:12:11	09:45:10	29.0
Jabba	202980	1087.2	202879	1086.6	30141	7847	95.627	99.9832	00:38:30	00:04:57	21.4
LoRDEC	238847	1405.0	237278	1297.1	34896	8326	99.978	97.9568	01:10:03	00:57:17	1.9
ECTools	130863	946.9	130832	943.1	28749	8412	99.810	99.7712	938:25:28	58:25:00	4.3
<i>Short-read-alignment-based methods</i>											
Hercules	239389	1460.3	235630	1330.4	35196	8644	99.976	87.6711	87:53:55	03:18:41	247.8
CoLoRMap	239309	1429.6	237135	1321.3	34850	8409	99.976	96.3912	18:44:48	03:07:34	37.3
Nanocorr	-	-	-	-	-	-	-	-	-	>72:00:00	-
proovread	238992	1412.4	236519	1298.0	35122	8369	99.978	97.9568	184:02:07	23:45:37	47.9
LSC	-	-	-	-	-	-	-	-	-	>72:00:00	-

Note: HG-CoLoR reported an error when correcting this dataset

Table 4 Experimental results for yeast ONT data set D2-O

Method	# Reads	# Bases (Mbp)	# Aligned reads	# Aligned bases (Mbp)	Maximum length (bp)	N50 (bp)	Genome fraction (%)	Alignment identity (%)	CPU time (hh:mm:ss)	Wall time (hh:mm:ss)	Memory usage (GB)
Original	118723	715.7	108463	638.1	55374	7003	99.976	86.1986	-	-	-
<i>Non-hybrid methods</i>											
FLAS	95606	585.6	95290	581.5	26592	6893	99.940	97.1699	07:42:10	07:42:10	4.4
LoRMA	398863	497.0	398350	495.2	16027	1439	99.485	98.4024	68:02:36	02:55:05	68.8
Canu	64829	475.1	64649	475.1	26895	7518	99.914	97.7710	12:31:04	00:37:53	9.0
<i>Short-read-assembly-based methods</i>											
HG-CoLoR	-	-	-	-	-	-	-	-	-	-	-
FMLRC	118701	713.7	111869	666.4	55374	6990	99.975	99.2529	03:35:44	00:17:21	2.2
HALC	118707	718.2	114071	647.9	55379	7025	99.976	98.8884	50:11:58	04:03:18	3.6
Jabba	99044	536.9	98631	535.9	28194	6730	95.400	99.9809	00:55:32	00:04:20	21.5
LoRDEC	118727	720.8	110606	647.8	55375	7049	99.976	96.9369	11:22:09	00:26:13	2.1
ECTools	81105	531.9	80843	529.3	26810	7071	99.314	99.7697	09:31:32	20:17:33	5.6
<i>Short-read-alignment-based methods</i>											
Hercules	118721	716.3	108467	638.9	55374	7008	99.976	87.2912	125:22:19	04:37:01	246.6
CoLoRMap	118774	722.0	108969	649.4	55374	7049	99.976	95.5851	11:01:38	01:34:52	27.8
Nanocorr	-	-	-	-	-	-	-	-	-	>72:00:00	-
proovread	118729	716.7	109057	643.4	55374	7007	99.976	96.3689	66:14:09	07:20:18	28.1
LSC	-	-	-	-	-	-	-	-	-	>72:00:00	-

Note: HG-CoLoR reported an error when correcting this dataset

Table 5 Experimental results for fruit fly PacBio data set D3-P

Method	# Reads	# Bases (Mbp)	# Aligned reads	# Aligned bases (Mbp)	Maximum length (bp)	N50 (bp)	Genome fraction (%)	Alignment identity (%)	CPU time (hh:mm:ss)	Wall time (hh:mm:ss)	Memory usage (GB)
Original	5366088	28797.8	1839681	16543.5	74735	15374	99.191	85.2734	-	-	-
<i>Non-hybrid methods</i>											
FLAS	1435682	14585.2	1428018	13574.1	43556	13550	98.915	98.8363	271:44:27	36:30:42	53.1
LoRMA	-	-	-	-	-	-	-	-	-	-	-
Canu	-	-	-	-	-	-	-	-	-	>72:00:00	-
<i>Short-read-assembly-based methods</i>											
HG-CoLoR	-	-	-	-	-	-	-	-	-	-	-
FMLRC	5246485	27354.6	2477890	16543.5	74735	14554	99.191	96.5284	327:37:22	13:49:04	31.2
HALC	4451474	21997.5	3434779	12793.3	74735	14349	99.178	96.8863	770:35:46	55:58:24	73.0
Jabba	35549	239.8	35505	239.1	37729	10461	65.616	99.9615	656:05:15	24:33:41	175.8
LoRDEC	5363998	28354.1	2056812	15636.9	74719	15078	99.200	92.2954	1011:52:27	36:19:18	5.9
ECTools	-	-	-	-	-	-	-	-	-	>72:00:00	-
<i>Short-read-alignment-based methods</i>											
Hercules	-	-	-	-	-	-	-	-	-	-	-
CoLoRMap	5366107	28891.6	1841822	14976.8	74735	15442	99.189	83.2580	495:11:17	64:52:25	189.4
Nanocorr	-	-	-	-	-	-	-	-	-	>72:00:00	-
proovread	-	-	-	-	-	-	-	-	-	>72:00:00	-
LSC	-	-	-	-	-	-	-	-	-	>72:00:00	-

Note: LoRMA, HG-CoLoR and Hercules reported errors when correcting this dataset

Table 6 Experimental results for fruit fly ONT data set D3-O

Method	# Reads	# Bases (Mbp)	# Aligned reads	# Aligned bases (Mbp)	Maximum length (bp)	N50 (bp)	Genome fraction (%)	Alignment identity (%)	CPU time (hh:mm:ss)	Wall time (hh:mm:ss)	Memory usage (GB)
Original	642255	4609.5	554083	3857.9	446050	11956	98.719	83.5921	-	-	-
<i>Non-hybrid methods</i>											
FLAS	423097	3507.6	422206	3402.6	64365	11517	97.588	95.3301	23:04:50	03:12:50	10.8
LoRMA	703097	615.5	682288	592.3	32644	865	30.338	98.1230	666:37:35	25:52:14	92.8
Canu	430082	3415.6	421475	3220.2	254967	12090	97.592	96.3739	88:51:10	04:36:20	20.2
<i>Short-read-assembly-based methods</i>											
HG-CoLoR	-	-	-	-	-	-	-	-	-	-	-
FMLRC	641945	4647.2	578290	3978.2	444605	12088	98.592	97.6010	47:45:17	03:06:05	31.2
HALC	643002	4668.5	611191	3955.7	451284	12115	98.616	97.6634	126:30:01	05:43:37	42.4
Jabba	494546	2878.2	494430	2876.3	72501	9305	83.166	99.9745	175:19:34	06:56:29	136.8
LoRDEC	642882	4655.9	567878	3921.1	447726	12079	98.691	94.0382	152:05:32	05:38:05	5.7
ECTools	-	-	-	-	-	-	-	-	-	>72:00:00	-
<i>Short-read-alignment-based methods</i>											
Hercules	642287	4612.8	554630	3859.4	449799	11966	98.713	83.9340	398:10:17	17:32:36	247.7
CoLoRMap	649041	4692.1	565881	3963.8	442948	12050	98.715	94.3361	160:00:22	16:07:18	57.3
Nanocorr	-	-	-	-	-	-	-	-	-	>72:00:00	-
proovread	-	-	-	-	-	-	-	-	-	>72:00:00	-
LSC	-	-	-	-	-	-	-	-	-	>72:00:00	-

Note: HG-CoLoR reported an error when correcting this dataset

Correction quality

We measure quality using the number of reads and total bases output in comparison with the input, the resulting alignment identity, fraction of the genome covered and read length distribution including maximum size and N50 length. From Tables 1, 2, 3, 4, 5 and 6, we gather that the best performing hybrid methods (e.g., FMLRC) are capable of correcting reads to achieve base-level accuracy in the high 90's. For the *E. coli* and yeast data sets, many of these programs achieve alignment identity > 99%. A crucial aspect to consider here is whether the high accuracy is achieved while preserving input read depth and N50. Few tools (e.g. Jabba) seem to attain high alignment identity at the cost of producing shorter reads and reduced depths because they choose to either discard uncorrected reads or trim the uncorrected regions. This may have a negative impact on downstream analyses. This trade-off is further discussed later in “[Effect of discarding reads during correction](#)” section.

Among the hybrid methods, a key observation is that short-read-assembly-based methods tend to show better performance than short-read-alignment-based methods. We provide the following explanation. Given that long reads are error-prone, short read alignment to long reads is more likely to be wrong (or ambiguous) than long read alignment to graph structures built using short reads. Errors in long reads can cause false positives in identifying the true positions where the respective short reads should align, which causes false correction later. For example, during the correction of D3-P, the alignment identity of corrected reads generated by CoLoRMap in fact decreased when compared to the uncorrected reads. The reason is that CoLoRMap uses BWA-mem [46] to map short reads, which is designed to report best mapping. However, due to the high error rates, the best mapping is not necessarily the true mapping. Large volume of erroneous long reads in D3-P can lead to many false alignments, which affected the correction process. On the other hand, long read lengths make it possible to have higher confidence when aligning them to paths in the graph. Therefore, in most of the experiments, assembly-based methods were able to produce reasonable correction.

Non-hybrid correction is more challenging as it relies solely on overlaps between erroneous long reads, yet the tools in this category yield competitive accuracy in many cases. However, non-hybrid methods may achieve lower alignment identity when the long reads are more erroneous. For example, the alignment identity of corrected reads generated by FLAS, LoRMA and Canu is lower than almost all hybrid methods for D1-O where the average alignment identity of uncorrected reads is only 81.36%.

Runtime and memory usage

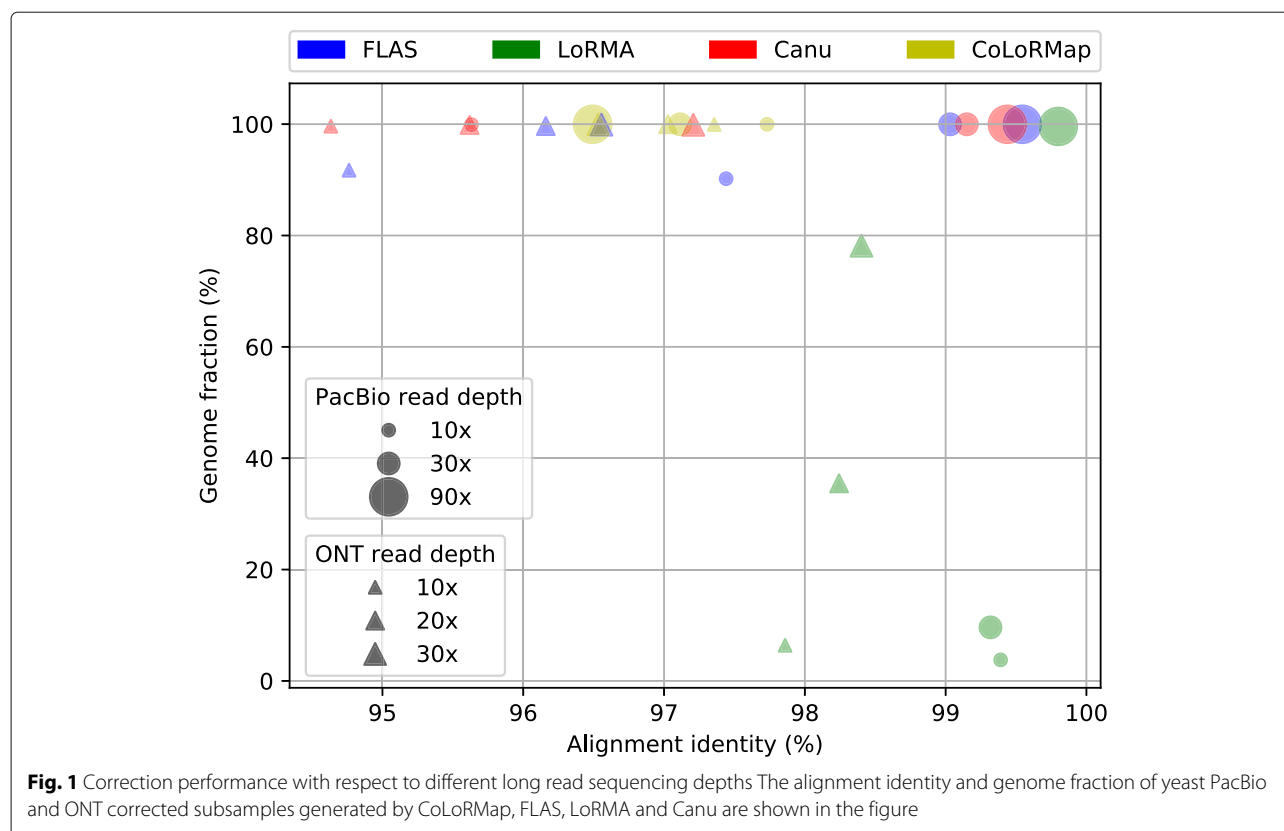
Scalability of the correction tools is an important aspect to consider in their evaluation. Slow speed or high memory usage makes it difficult to apply them to correct large data sets. Our results show that hybrid methods, in particular assembly-based methods, are much faster than the rest. For instance, Hercules and LSC failed to generate corrected reads in three days for D1-P, while most of the assembly-based tools finished the same job in less than one hour. Hercules, Nanocorr and LSC were unable to finish the correction of D2-O in three days, which was finished by FMLRC or LoRDEC in hours. Although proovread can complete the corrections of D2-P and D2-O, the run-time was 49.3 and 34.4 times longer, respectively, than run-time needed by FMLRC. Moreover, assembly-based methods, e.g., LoRDEC and FMLRC, used less memory in most of the experiments. Therefore, in terms of computational performance, users should give priority to short-read assembly-based methods over short-read alignment-based methods.

Among the non-hybrid methods, LoRMA's memory usage was generally the highest among all the tools, and was slower than assembly-based methods. However, Canu showed superior scalability. Owing to a fast long read overlap detection algorithm using MinHash [45], Canu was able to compute long read overlaps and used them to correct the reads in reasonable time, which is comparable to most of the hybrid methods. The memory footprint of Canu was also lower than many hybrid-methods. However, Canu did not finish the correction of D3-P in three days probably because this data set is too large to compute pairwise overlaps. FLAS showed performance comparable to Canu as FLAS also leverages the fast overlap computation method in MECAT [44].

Effect of long read sequencing depth on error correction

Requiring high sequencing coverage for effective error correction can impact both cost and time consumed during sequencing and analysis. The relative cost per base pair using third-generation sequencing is still several folds higher when compared to the latest Illumina sequencers [1]. Accordingly, we study how varying long read sequencing depth affects correction quality, while keeping the short read data set fixed. We conducted this experiment using data sets D2-P and D2-O with various depth levels obtained using random sub-sampling. The details of the subsamples are summarized in Additional file 1 Table S1. The output behavior of the correction tools is shown in Additional file 1 Tables S2-S9.

Figure 1 shows the correction performance of four tools— CoLoRMap, LoRMA, Canu and FLAS that were significantly affected by the long read sequencing depth. For corrected reads generated by hybrid methods, no significant change on the metrics was observed except



those generated by CoLoRMap. The alignment identity of its corrected reads increased with decreased sequencing depth. This observation is consistent with the experimental results reported by its authors. Similarly, CoLoRMap did not perform well on large data sets such as D3-P as large data sets increase the risk of false positive alignments.

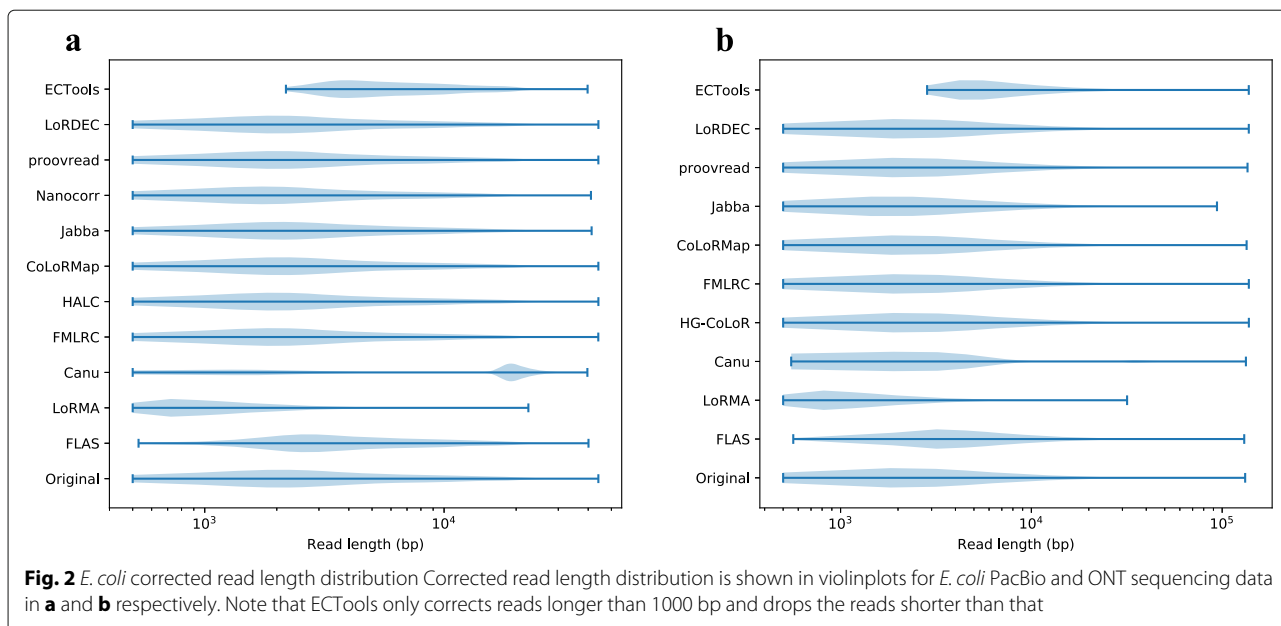
On the other hand, the performance of non-hybrid methods deteriorated significantly when sequencing depth was decreased. As non-hybrid methods leverage overlap information to correct errors, they require sufficient long read coverage to make true correction. The genome fraction covered by corrected reads produced by LoRMA with subsamples of D2-P decreased from 99.59% to 82.97% when sequencing depth dropped from 90x to 60x, and further decreased to 9.61%, 5.39% and 3.78% for 30x, 20x and 10x respectively, implying loss of many long reads after correction. The alignment identities were still greater than 99% using all subsamples because LoRMA trimmed the uncorrected regions. For corrected reads generated by Canu, no significant change on genome fraction was observed. But the alignment identity dropped from above 99% to 97.03% and 95.63% for 20x and 10x sequencing depths, respectively. FLAS showed similar performance but genome fraction for 10x sequencing depth was only 90.20% lower than the 99.92% achieved

by Canu, which indicates FLAS drops some reads when sequencing depth is low.

Effect of discarding reads during correction

Many correction tools opt for discarding input reads or regions within reads that they fail to correct. As a result, the reported alignment identity is high (> 99%), but much fewer number of bases survive after correction. This effect is more pronounced in corrected reads generated by Jabba, ECTools, and LoRMA. They either trim uncorrected regions at sequence ends, or even in the middle, to avoid errors in the final output which eventually yields high alignment identity. However, aggressive trimming also makes the correction lossy and may influence downstream analysis because long range information is lost if the reads are shortened or broken into smaller pieces. Therefore, users should be conservative in trimming and turn it off when necessary. One good practice is to keep the uncorrected regions and let downstream analysis tools perform the trimming, e.g. overlap-based trimming after read correction in Canu.

A direct implication of discarding or trimming reads is the change of read length distribution. Figure 2 shows the original and corrected read length distributions. Among all the tools, HG-CoLoR, FMLRC, HALC, CoLoRMAP, LoRDEC and proovread can maintain a similar read



length distribution after correction whereas Nanocorr, Jabba and ECTools lost many long reads after correction due to their trimming step. Nanocorr drops a long read when there is no short read aligning to it. This procedure can remove many error-prone long reads, which leads to a higher alignment identity after correction. However, the fraction of discarded reads in many cases is found to be significant. For example, a mere 376.3 million bp cumulative length of sequences survived out of 1462.7 million bp data set, after correction of D2-P. ECTools also generated only 946.9 million corrected bases using this data set. Canu changed the read length distribution significantly after correction although due to a different reason. Canu estimates the read length after correction and tries to keep the longest 40x reads for subsequent assembly. FLAS kept most of the reads with short length while losing many reads with long length.

Effect of error correction on genome assembly

Error correction of long reads remains a useful pre-processing stage for reliable construction of overlap graphs during genome assembly. We examined how well the accuracy of error correction correlates with the quality of genome assembly performed using corrected reads. To do so, we conducted an experiment to compute genome assembly using corrected PacBio and ONT reads of *E. coli*, i.e., corrected reads for D1-P and D1-O. Assembly was computed using Canu with its error correction module turned off, and assembly quality was assessed using QAST [47]; the results are shown in Tables 7 and 8.

Considering the assemblies generated using corrected PacBio reads (Table 7), NGA50 of about 3 Mbp was obtained when using reads generated by FLAS, Canu,

FMLRC, HALC, Nanocorr, LoRDEC or ECTools. When using corrected ONT reads (Table 8), assemblies generated using reads corrected by Canu, HG-CoLoR, FMLRC, LoRDEC and ECTools have NGA50 near 3 Mbp. In contrast, assemblies generated using reads corrected by Jabba and LoRMA showed lower NGA50 in both cases. Their trimming procedure possibly led to the loss of some long range information, thereby causing lower continuity in assembly. Post error correction, alignment identity of corrected reads needs to be sufficiently high to identify true overlaps during assembly. We observe that NGA50 of assemblies generated using reads corrected by CoLoRMap and proovread is low, as the corrected reads generated by these two tools have low alignment identity (e.g., < 90% for D1-O, Table 2).

We also examined the frequency of mismatches and indels in the assemblies. For data set D1-P, corrected reads generated by HALC and ECTools produced assemblies with > 500 mismatches, significantly higher than the other tools. However, alignment identity of their corrected reads was either competitive with, or superior to, what is produced by other tools. Notably, both HALC and ECTools use assembled contigs from short reads to do error correction. Mis-assemblies of short reads, especially in repetitive and low-complexity regions, may cause false corrections, which leads to errors during assembly [29]. Corrected reads produced by FMLRC achieved the least number of errors in assembly. Meanwhile, its alignment identity was also the highest among the methods which avoid trimming. Therefore, higher alignment identity of corrected reads is an important but not a sufficient criteria to minimize errors in genome assemblies.

Table 7 Results of genome assembly computed using corrected reads of D1-P

Method	# contigs	NGA50 (bp)	Largest contigs (bp)	Total length (bp)	Genome fraction (%)	# misassemblies	# mismatches	# indels (<=5bp)	# indels (>5bp)	Indel length (bp)
<i>Non-hybrid methods</i>										
FLAS	2	3996362	4681650	4689583	99.998	4	4	162	0	167
LoRMA	14	696878	2501146	4663900	99.938	4	75	4181	6	4295
Canu	1	3976437	4670120	4670120	99.998	4	7	92	0	95
<i>Short-read-assembly-based methods</i>										
FMLRC	9	3821409	4657352	4831908	99.998	8	1	4	0	5
HALC	25	2947777	4682714	5388722	99.983	8	541	35	8	257
Jabba	58	138874	398327	4623296	97.273	1	172	32	3	167
LoRDEC	2	3996441	4681757	4703690	99.998	4	66	18	2	55
ECTools	19	3548731	4657296	5154324	99.974	4	592	80	2	188
<i>Short-read-alignment-based methods</i>										
CoLoRMap	86	1217587	1448649	5700143	99.998	4	42	3	7	478
Nanocorr	18	3095077	4646253	4931697	99.998	5	65	34	2	157
proovread	2	1686030	4626702	4666724	99.656	4	76	97	1	176

Non-hybrid methods such as LoRMA, Canu and FLAS produced more indels than mismatches in their assemblies while most of the hybrid methods showed the opposite behavior. These observations suggest that existing self-correction methods are not good at handling indels when compared to hybrid methods. Consequently, *de novo* long read assemblers that use self-correction methods typically rely on post-processing ‘polishing’ stages, using signal-level data from long read instruments [2, 33] or alternate sequencing technologies [48].

Conclusions and future directions

In this work, we established benchmark data sets and evaluation methods for comprehensive assessment of long read error correction software. Our results suggest that hybrid methods aided by short accurate reads can achieve better correction quality, especially when handling low coverage-depth long reads, compared with non-hybrid methods. Within the hybrid methods, assembly-based methods are superior to alignment-based methods in terms of scalability to large data sets. Besides, better

Table 8 Results of genome assembly computed using corrected reads of D1-O

Method	# contigs	NGA50 (bp)	Largest contigs (bp)	Total length (bp)	Genome fraction (%)	# misassemblies	# mismatches	# indels (<=5bp)	# indels (>5bp)	Indel length (bp)
<i>Non-hybrid methods</i>										
FLAS	1	1283465	4561925	4561925	99.834	4	4785	70062	1310	122264
LoRMA	14	726649	1239048	4583602	99.650	3	23275	57650	37	73578
Canu	1	3335496	4601279	4601279	99.914	2	14810	56511	283	78931
<i>Short-read-assembly-based methods</i>										
HG-CoLoR	32	3924167	4634988	5555776	99.845	28	393	82	8	246
FMLRC	9	4325756	4718452	4974808	99.874	6	6	7	4	242
Jabba	57	105474	311624	4460218	95.838	0	117	22	5	179
LoRDEC	57	3492326	4044623	5389657	99.800	10	402	162	18	448
ECTools	2	2891718	4733248	4797686	99.885	3	632	267	16	682
<i>Short-read-alignment-based methods</i>										
CoLoRMap	55	206971	501963	6007440	94.790	4	11074	14392	206	24183
proovread	99	75162	225568	4613429	98.737	14	276	41	1	75

performance on correction such as preserving higher proportion of input bases and high alignment identity often leads to better performance in downstream applications such as genome assembly. But the tools with superior correction performance should be further tested in the context of applications of interest.

Users can also select tools according to our experimental results for their specific expectations. When speed is a concern, assembly-based hybrid methods are preferred whenever short reads are available. Besides, hybrid methods are more immune to low long read sequencing depth than non-hybrid methods. Thus, users are recommended to choose hybrid methods when long read sequencing depth is relatively low. In cases where indel errors may cause a serious negative impact on downstream analyses, hybrid methods should be preferred over non-hybrid ones if short reads are available.

FMLRC outperformed other hybrid methods in almost all the experiments. For non-hybrid methods, Canu and FLAS showed better performance over LoRMA. Hence, these three are recommended as default when users want to avoid laborious tests on all the error correction tools.

For future work, better self-correction algorithms are expected to avoid hybrid sequencing, thus reducing experimental labor on short read sequencing preparation. In addition, most of the correction algorithms run for days to correct errors in the sequencing of even moderately large and complex genomes like the fruit fly. These algorithms will spend much more time on correcting larger sequencing data sets of human, and become a bottleneck in sequencing data analysis. Therefore, more efficient or parallel correction algorithms should be developed to ease the computational burden. Furthermore, none of the hybrid tools makes use of paired-end information in their correction, except CoLoRMap. But the use of paired-end reads in CoLoRMap did not improve correction performance significantly according to previous studies. Paired-end reads have already been used to resolve repeats and remove entanglements in de Bruijn graphs [49]. Since many error correction tools build de Bruijn graphs to correct long reads, the paired-end information may also be able to improve error correction.

Most of the published error correction tools focus on correction of long DNA reads sequenced from a single genome, which also served as the motivation for our review. Long read sequencing is increasingly gaining traction in transcriptomics and metagenomics applications. It is not clear whether the existing tools can be leveraged or extended to work effectively in such scenarios, and is an active area of research [50].

Methods

We selected data sets from recent publicly accessible genome sequencing experiments. For benchmarking

the different programs, our experiments used genome sequences from multiple species and different sequencing platforms with recent chemistry, e.g., R9 for ONT or P6-C4/P5-C3 for PacBio. We describe our evaluation criteria and use it for a comprehensive assessment of the correction methods/software.

Benchmark data sets

Our benchmark includes resequencing data from three reference genomes – *Escherichia coli* K-12 MG1655 (*E. coli*), *Saccharomyces cerevisiae* S288C (yeast), and *Drosophila melanogaster* ISO1 (fruit fly). The biggest hurdle when benchmarking with real data is the absence of ground truth (i.e., perfectly corrected reads). However, the availability of reference genomes of these strains enables us to evaluate the output of correction software in a reliable manner using the reference. Essentially, differences in a corrected read with respect to the reference imply uncorrected errors. A summary of the selected read data sets is listed in Table 9. We leveraged publicly available high coverage read data sets of the selected genomes available from all three platforms – Illumina (for short reads), Pacbio, and ONT. In addition, some of these samples were sequenced using multiple protocols, yielding reads of varying quality. This enabled us to do a thorough comparison among error correction software across various error rates and error profiles.

To conduct performance evaluation under different sequencing depths, yeast sequencing reads (D2-P and D2-O) were subsampled randomly using Seqtk (<https://github.com/lh3/seqtk>). Subsamples with average depth of 10x, 20x and 30x were generated for ONT reads. In addition, 10x, 20x, 30x, 60x and 90x PacBio read subsamples were generated from D2-P. Details of these subsamples are available in Additional file 1 Table S1.

Evaluation methodology

Our evaluation method takes corrected reads and a reference genome as input. The corrected reads were filtered using a user defined length (default 500). Reads which were too short to include in downstream analysis were dropped during the filtration. Filtered reads were aligned to the reference genome using Minimap2 [51] (using “-ax map-pb” and “-ax map-ont” parameters for PacBio and ONT reads respectively). The primary alignment for each read was used in the evaluation.

In an ideal scenario, an error correction software should take each erroneous long read and produce the error-free version of it, preserving each read and its full length. To assess how close to the ideal one can get, measures such as error rate post-correction or percentage of errors removed (termed *gain*; see [15]) can be utilized. However, long read error correction programs do not operate in this fashion. They may completely discard some reads

Table 9 Details of the benchmark data sets

Data set	Sequencing specification	Sequencing NCBI accession	Sequencing depth ^a	Read length (bp) ^b	Number of reads	Reference genome	Genome length (Mbp)	Reference NCBI accession
D1-I	Illumina Miseq	- ^c	373x	2×151	2×5 729 470			
D1-P	Pacbio P6C4	- ^d	161x	13 982	87 217	<i>E. coli</i> K-12 MG1655	4.6	NC_000913.3
D1-O	MinION R9 1D	- ^e	319x	14 891	164 472			
D2-I	Illumina Miseq	ERR1938683	81x	2×150	2×3 318 467			
D2-P	Pacbio P6C4	PRJEB7245	120x	8 656	239 408	<i>S. cerevisiae</i> S288c	12.2	GCF_000146045.2
D2-O	MinION R9 2D	ERP016443	59x	7 001	119 955			
D3-I	Illumina Nextseq	SRX3676782	44x	2×151	2×20 619 401			
D3-P	Pacbio P5C3	SRX499318	204x	15 132	6 864 972	<i>D. melanogaster</i> ISO1	143.7	GCF_000001215.4
D3-O	MinION R9.5 1D	SRX3676783	32x	11 934	663 784			

^aSequencing depth is estimated using the sequencing data and reference genome size

^bN50 is reported for PacBio or ONT reads, since their lengths vary

^cDownloaded from Illumina at https://ftp://webdata.webdata@ussd-ftp.illumina.com/Data/SequencingRuns/MG1655/MiSeq_Ecoli_MG1655_110721_PF_R1.fastq.gz and https://ftp://webdata.webdata@ussd-ftp.illumina.com/Data/SequencingRuns/MG1655/MiSeq_Ecoli_MG1655_110721_PF_R2.fastq.gz

^dDownloaded from PacBio at <https://github.com/PacificBiosciences/DevNet/wiki/E.-coli-Bacterial-Assembly>

^eDownloaded from Loman Labs at https://s3.climb.ac.uk/nanopore/E_coli_K12_1D_R9.2_SpotON_2.pass.fasta

or choose to split an input read into multiple reads when the high error rate cannot be reckoned with. In addition, short read assembly based error correction programs use long read alignments to de Bruijn graphs, and produce sequences corresponding to the aligned de Bruijn graph paths as output reads instead. Though original reads may not be fully preserved, all that matters for effective use of error correction software is that its output consists of sufficient number of high quality long reads that reflect adequate read lengths, sequencing depth, and coverage of the genome. Accordingly, our evaluation methodology reflects such assessment.

We measure the number of reads and total bases output by each error correction software, along with the number of aligned reads and total number of aligned bases extracted from alignment statistics, because they together reveal the effectiveness of correction. Besides, statistics which convey read length distribution such as maximum length and N50 were calculated to assess effect of the correction process on read lengths. Fraction of bases on the genome covered by output reads is also reported to assess if there are regions of the genome that lost coverage or suffered significant deterioration in coverage depth post-correction. Any significant drop on these metrics can be a potential sign of information loss during the correction. Finally, alignment identity is calculated by the number of total matched bases divided by the total alignment length. Tools which achieve maximum alignment identity with minimum loss of information are desirable.

As part of this study, we provide an evaluation tool to automatically generate the evaluation statistics of corrected reads mentioned above (<https://github.com/haowenz/LRECE>). We include a wrapper script which can

run state-of-the-art error correction software on a grid engine given any input data from user. Using the script, two types of evaluations can be conducted; users can either evaluate the performance on a list of tools with their own data to find a suitable tool for their studies, or they can run any correction tool with the benchmark data and compare it with other state-of-the-art tools.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12864-020-07227-0>.

Additional file 1: Settings of the error correction tools and correction results on subsamples of yeast PacBio and ONT sequencing data.

Abbreviations

PacBio: Pacific bioscience; ONT: Oxford Nanopore technologies; pHMM: Profile Hidden Markov Model; DBG: de Bruijn graph; MEM: Maximal exact matches

Acknowledgements

Not applicable.

About this supplement

This article has been published as part of BMC Genomics Volume 21 Supplement 6, 2020: Selected articles from the 8th IEEE International Conference on Computational Advances in Bio and medical Sciences (ICCABS 2018): genomics. The full contents of the supplement are available online at <https://bmcbgenomics.biomedcentral.com/articles/supplements/volume-21-supplement-6>.

Authors' contributions

HZ, CJ and SA designed the study. HZ performed the experiments and wrote the paper. CJ and SA revised the paper. All authors read and approved the final manuscript.

Funding

This work was supported in part by the National Science Foundation under CCF-1718479. Publication costs are funded by the National Science Foundation under CCF-1718479. The funding body played no role in the design of the study and collection, analysis, interpretation of data or in writing the manuscript.

Availability of data and materials

E. coli Illumina reads are available at https://ftp://webdata.webdata@ussd-ftp.illumina.com/Data/SequencingRuns/MG1655/MiSeq_Ecoli_MG1655_110721_P1_P2.fastq.gz and https://ftp://webdata.webdata@ussd-ftp.illumina.com/Data/SequencingRuns/MG1655/MiSeq_Ecoli_MG1655_110721_P1_P2.fastq.gz. *E. coli* PacBio reads are available at <https://github.com/PacificBiosciences/DevNet/wiki/E.-coli-Bacterial-Assembly>. *E. coli* ONT reads are available at https://s3.climb.ac.uk/nanopore/E_coli_K12_1D_R9.2_SpotON_2.pass.fasta. *E. coli* reference genome is available with NCBI accession number NC_000913.3. Yeast reference genome and sequencing data are available with NCBI accession number GCF_000146045.2, ERR1938683, PRJEB7245 and ERP016443. Fruit fly reference genome and sequencing data are available with NCBI accession number GCF_000001215.4, SRX3676782, SRX499318 and SRX3676783.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Received: 8 November 2020 Accepted: 12 November 2020

Published: 21 December 2020

References

- Sedlazeck FJ, Lee H, Darby CA, Schatz MC. Piercing the dark matter: bioinformatics of long-range sequencing and mapping. *Nat Rev Genet*. 20181. <https://doi.org/10.1038/s41576-018-0003-4>.
- Loman NJ, Quick J, Simpson JT. A complete bacterial genome assembled de novo using only nanopore sequencing data. *Nat Methods*. 2015;12(8):733.
- Chin CS, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, et al. Phased diploid genome assembly with single-molecule real-time sequencing. *Nat Methods*. 2016;13(12):1050.
- Jain M, Koren S, Miga KH, Quick J, Rand AC, Sasani TA, et al. Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nat Biotechnol*. 2018;36(4):338.
- Sedlazeck FJ, Rescheneder P, Smolka M, Fang H, Nattestad M, Von Haeseler A, Schatz MC. Accurate detection of complex structural variations using single-molecule sequencing. *Nat Methods*. 2018;15(6):461–468. Nature Publishing Group.
- Chaisson MJ, Huddleston J, Dennis MY, Sudmant PH, Malig M, Hormozdiari F, et al. Resolving the complexity of the human genome using single-molecule sequencing. *Nature*. 2015;517(7536):608.
- Gordon SP, Tseng E, Salamov A, Zhang J, Meng X, Zhao Z, et al. Widespread polycistronic transcripts in fungi revealed by single-molecule mRNA sequencing. *PLoS ONE*. 2015;10(7):e0132628.
- Dilthey A, Jain C, Koren S, Phillippy A. MetaMaps-Strain-level metagenomic assignment and compositional estimation for long reads. *bioRxiv*. 2018372474. <https://doi.org/10.1101/372474>.
- Rand AC, Jain M, Eizenga JM, Musselman-Brown A, Olsen HE, Akeson M, et al. Mapping DNA methylation with high-throughput nanopore sequencing. *Nat Methods*. 2017;14(4):411.
- Simpson JT, Workman RE, Zuzarte P, David M, Dursi L, Timp W. Detecting DNA cytosine methylation using nanopore sequencing. *Nat Methods*. 2017;14(4):407.
- Carneiro MO, Russ C, Ross MG, Gabriel SB, Nusbaum C, DePristo MA. Pacific biosciences sequencing technology for genotyping and variation discovery in human data. *BMC Genomics*. 2012;13(1):375.
- Jain M, Fiddes IT, Miga KH, Olsen HE, Paten B, Akeson M. Improved data analysis for the MinION nanopore sequencer. *Nat Methods*. 2015;12(4):351.
- Korlach J, Biosciences P. Understanding accuracy in SMRT sequencing. Technical report. 2013. www.pacb.com.
- Ashton PM, Nair S, Dallman T, Rubino S, Rabsch W, Mwaigwisya S, et al. MinION nanopore sequencing identifies the position and structure of a bacterial antibiotic resistance island. *Nat Biotechnol*. 2015;33(3):296.
- Yang X, Chockalingam SP, Aluru S. A survey of error-correction methods for next-generation sequencing. *Brief Bioinforma*. 2012;14(1):56–66.
- Alic AS, Ruzafa D, Dopazo J, Blanquer I. Objective review of de novo stand-alone error correction methods for NGS data. *Wiley Interdiscip Rev Comput Mol Sci*. 2016;6(2):111–46.
- Koren S, Schatz MC, Walenz BP, Martin J, Howard JT, Ganapathy G, et al. Hybrid error correction and de novo assembly of single-molecule sequencing reads. *Nat Biotechnol*. 2012;30(7):693.
- Au KF, Underwood JG, Lee L, Wong WH. Improving PacBio long read accuracy by short read alignment. *PLoS ONE*. 2012;7(10):e46679.
- Lee H, Gurtowski J, Yoo S, Marcus S, McCombie WR, Schatz M. Error correction and assembly complexity of single molecule sequencing reads. *BioRxiv*. 2014006395. <https://doi.org/10.1101/006395>.
- Salmela L, Rivals E. LoRDEC: accurate and efficient long read error correction. *Bioinformatics*. 2014;30(24):3506–14.
- Hackl T, Hedrich R, Schultz J, Förster F. proovread: large-scale high-accuracy PacBio correction through iterative short read consensus. *Bioinformatics*. 2014;30(21):3004–11.
- Madoui MA, Engelen S, Cruaud C, Belsler C, Bertrand L, Alberti A, et al. Genome assembly using Nanopore-guided long and error-free DNA reads. *BMC Genomics*. 2015;16(1):327.
- Goodwin S, Gurtowski J, Ethe-Sayers S, Deshpande P, Schatz MC, McCombie WR. Oxford Nanopore sequencing, hybrid error correction, and de novo assembly of a eukaryotic genome. *Genome Res*. 2015;25(11):1750–6.
- Miclote G, Heydari M, Demeester P, Rombauts S, Van de Peer Y, Audenaert P, et al. Jabba: hybrid error correction for long sequencing reads. *Algorithm Mol Biol*. 2016;11(1):10.
- Haghshenas E, Hach F, Sahinalp SC, Chauve C. Colormap: Correcting long reads by mapping short reads. *Bioinformatics*. 2016;32(17):i545–51.
- Salmela L, Walve R, Rivals E, Ukkonen E. Accurate self-correction of errors in long reads using de Bruijn graphs. *Bioinformatics*. 2016;33(6):799–806.
- Bao E, Lan L. HALC: High throughput algorithm for long read error correction. *BMC Bioinformatics*. 2017;18(1):204.
- Bao E, Xie F, Song C, Dandan S. FLAS: fast and high-throughput algorithm for PacBio long-read self-correction. *Bioinformatics*. 2019;35(20):3953–60. Oxford University Press.
- Wang JR, Holt J, McMillan L, Jones CD. FMLRC: Hybrid long read error correction using an FM-index. *BMC Bioinformatics*. 2018;19(1):50.
- Morisse P, Lecroq T, Lefebvre A, Berger B. Hybrid correction of highly noisy long reads using a variable-order de Bruijn graph. *Bioinformatics*. 2018. <https://doi.org/10.1093/bioinformatics/bty521>.
- Firtina C, Bar-Joseph Z, Alkan C, Cicek AE. Hercules: a profile HMM-based hybrid error correction algorithm for long reads. *Nucleic Acids Res*. 2018;46(21):e125–.
- Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res*. 2017;27(5):722–36.
- Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al. Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. *Nat Methods*. 2013;10(6):563.
- Laehnmann D, Borkhardt A, McHardy AC. Denoising DNA deep sequencing data high-throughput sequencing errors and their correction. *Brief Bioinforma*. 2015;17(1):154–79.
- Mahmoud M, Zywicki M, Twardowski T, Karlowski WM. Efficiency of PacBio long read correction by 2nd generation Illumina sequencing. *Genomics*. 2017. <https://doi.org/10.1016/j.ygeno.2017.12.011>.
- La S, Haghshenas E, Chauve C. LRCstats, a tool for evaluating long reads correction methods. *Bioinformatics*. 2017;33(22):3652–4.
- Fichot EB, Norman RS. Microbial phylogenetic profiling with the Pacific Biosciences sequencing platform. *Microbiome*. 2013;1(1):10.
- Stöcker BK, Köster J, Rahmann S. SimLoRD: simulation of long read data. *Bioinformatics*. 2016;32(17):2704–6.
- Fu S, Wang A, Au KF. A comparative evaluation of hybrid error correction methods for error-prone long reads. *Genome Biol*. 2019;20(1):26.
- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol*. 1990;215(3):403–10.
- Pop M, Phillippy A, Delcher AL, Salzberg SL. Comparative genome assembly. *Brief Bioinforma*. 2004;5(3):237–48.
- Yang X, Dorman KS, Aluru S. Reptile: representative tiling for short read error correction. *Bioinformatics*. 2010;26(20):2526–33.

43. Kowalski T, Grabowski S, Deorowicz S. Indexing arbitrary-length k-mers in sequencing reads. *PLoS ONE*. 2015;10(7):e0133198.
44. Xiao CL, Chen Y, Xie SQ, Chen KN, Wang Y, Han Y, et al. MECAT: fast mapping, error correction, and de novo assembly for single-molecule sequencing reads. *nature methods*. 2017;14(11):1072.
45. Berlin K, Koren S, Chin CS, Drake JP, Landolin JM, Phillippy AM. Assembling large genomes with single-molecule sequencing and locality-sensitive hashing. *Nat Biotechnol*. 2015;33(6):623.
46. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv preprint arXiv*. 2013:13033997.
47. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics*. 2013;29(8):1072–5.
48. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS ONE*. 2014;9(11):e112963.
49. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol*. 2012;19(5):455–77.
50. de Lima LIS, Marchet C, Caboche S, Da Silva C, Istace B, Aury JM, et al. Comparative assessment of long-read error-correction software applied to RNA-sequencing data. *bioRxiv*. 2018:476622.
51. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. 2018;1:7.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

