

Atropos: specific, sensitive, and speedy trimming of sequencing reads

John P Didion ^{Corresp., 1}, Marcel Martin ², Francis S Collins ¹

¹ National Human Genome Research Institute, National Institutes of Health, Bethesda, Maryland, United States

² Science for Life Laboratory, Department of Biochemistry and Biophysics, Stockholm University, Stockholm, Sweden

Corresponding Author: John P Didion

Email address: john.didion@nih.gov

A key step in the transformation of raw sequencing reads into biological insights is the trimming of adapter sequences and low-quality bases. Read trimming has been shown to increase the quality and reliability while decreasing the computational requirements of downstream analyses. Many read trimming software tools are available; however, no tool simultaneously provides the accuracy, computational efficiency, and feature set required to handle the types and volumes of data generated in modern sequencing-based experiments. Here we introduce Atropos and show that it trims reads with high sensitivity and specificity while maintaining leading-edge speed. Compared to other state-of-the-art read trimming tools, Atropos achieves a four-fold increase in trimming accuracy and a decrease in execution time of up to 40% (using 16 parallel execution threads). Furthermore, Atropos maintains high accuracy even when trimming data with elevated rates of sequencing errors. The accuracy, high performance, and broad feature set offered by Atropos makes it an appropriate choice for the pre-processing of most current-generation sequencing data sets. Availability. Atropos is open source and free software written in Python and available at <https://github.com/jdidion/atropos>.

Atropos: specific, sensitive, and speedy trimming of sequencing reads

John P Didion¹, Marcel Martin², and Francis S Collins¹

¹National Human Genome Research Institute, National Institutes of Health, Bethesda, MD

²Science for Life Laboratory, Department of Biochemistry and Biophysics, Stockholm University, Sweden

Corresponding author:

John P Didion, PhD¹

Email address: john.didion@nih.gov

ABSTRACT

Summary. A key step in the transformation of raw sequencing reads into biological insights is the trimming of adapter sequences and low-quality bases. Read trimming has been shown to increase the quality and reliability while decreasing the computational requirements of downstream analyses. Many read trimming software tools are available; however, no tool simultaneously provides the accuracy, computational efficiency, and feature set required to handle the types and volumes of data generated in modern sequencing-based experiments. Here we introduce Atropos and show that it trims reads with high sensitivity and specificity while maintaining leading-edge speed. Compared to other state-of-the-art read trimming tools, Atropos achieves a four-fold increase in trimming accuracy and a decrease in execution time of up to 40% (using 16 parallel execution threads). Furthermore, Atropos maintains high accuracy even when trimming data with elevated rates of sequencing errors. The accuracy, high performance, and broad feature set offered by Atropos makes it an appropriate choice for the pre-processing of most current-generation sequencing data sets. **Availability.** Atropos is open source and free software written in Python and available at <https://github.com/jdidion/atropos>.

1 INTRODUCTION

All current-generation sequencing technologies, including Illumina, SOLiD, PacBio, and Nanopore, require a library construction step that involves the introduction of short adapter sequences at the ends of the template DNA fragments. Depending on the sequencing platform and the fragment size distribution of the sequencing library, an often substantial fraction of reads will consist of both template and adapter sequences (Figure 1A). Additionally, the error rates of these sequencing technologies vary from ~0.1% on Illumina to 5% or more on long-read sequencing platforms. Error rates tend to be enriched at the ends of reads (where adapters are located), thus exacerbating the effects of adapter contamination. Adapter contamination and sequencing errors can lead to increased rates of misaligned and unaligned reads, which results in errors in downstream analysis including spurious variant calls (Del Fabbro et al., 2013; Sturm et al., 2016). Certain sequencing protocols may introduce other artifacts in sequencing reads. For example, some methylation sequencing (Methyl-Seq) protocols result in artificially methylated bases at the 3' ends of reads that can lead to inflated estimates of methylation levels (Bock, 2012).

Read trimming is an important step in the analysis pipeline to mitigate the effects of adapter contamination, sequencing errors, and other artifacts. The development of tools for read trimming is an active area of bioinformatics research, thus there are currently many options. In terms of adapter trimming, these tools fall into two general categories: 1) those that rely solely on matching the adapter sequence (*adapter-match trimming*) using semi-global alignment (which is the only option available for single-end reads; Figure 1B); and 2) those that leverage the overlap between paired-end reads to identify adapter starting positions (*insert-match trimming*; Figure 1C) (Sturm et al., 2016). Cutadapt (Martin, 2011) is a mature and feature-rich example of a tool that provides adapter-match trimming, while SeqPurge (Sturm et al., 2016) is a recent example of a highly accurate insert-match trimmer designed specifically

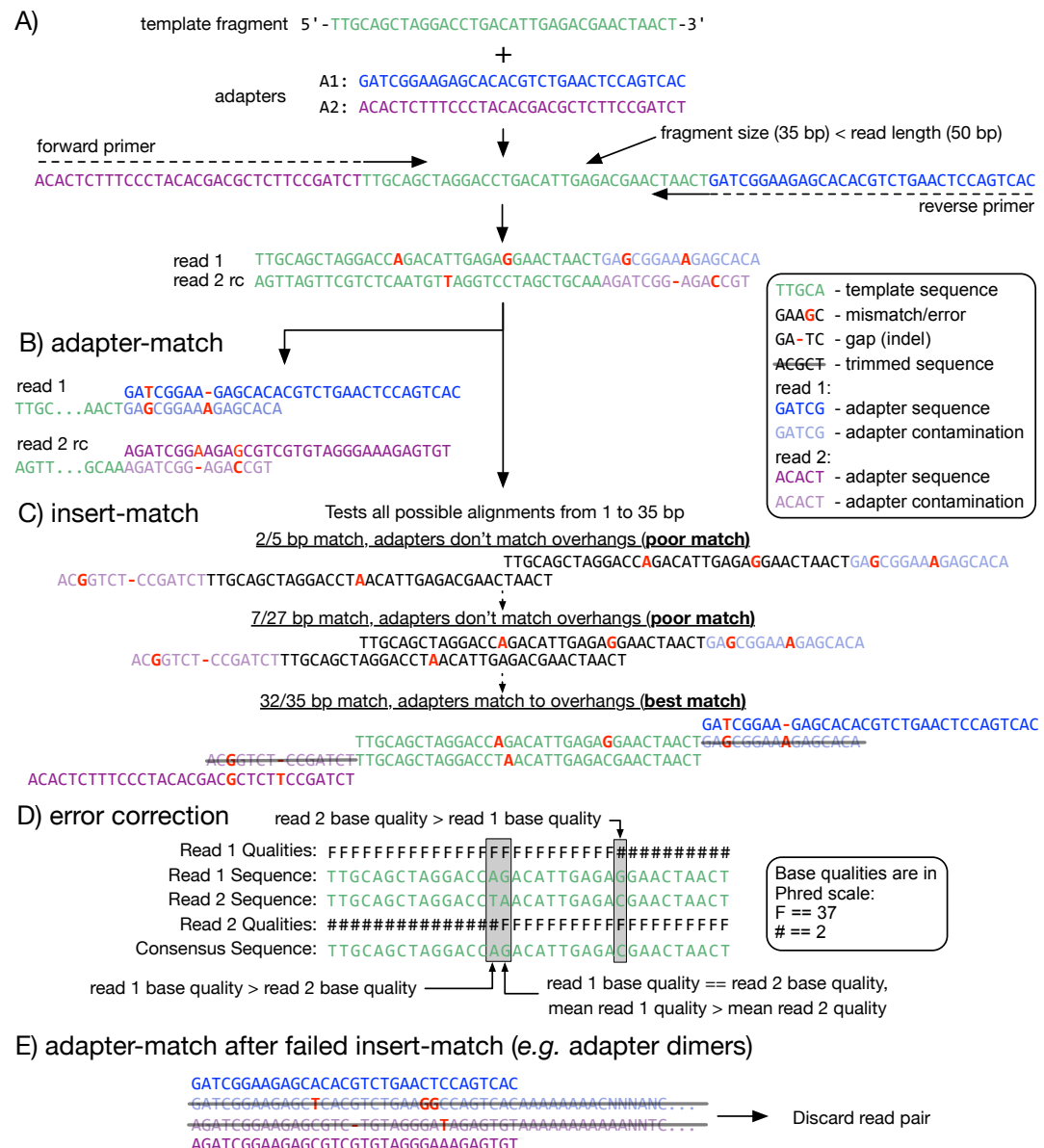


Figure 1. Adapter detection and trimming. A) When a fragment (or insert; green) is shorter than the read length, the read sequence will contain partial to full-length adapter sequences (blue and purple). B,C) Methods for detecting adapter contamination using semi-global alignment. Adapter-match (B) identifies the best alignment between each adapter and the end of its corresponding read. Insert-match (C) first identifies the best alignment between read 1 and the reverse-complement (rc) of read 2; if a valid alignment is found, then adapters are matched to the remaining overhangs. D) If a match is found, the overlapping inserts can be used for mutual error correction. The consensus base is the one with the highest quality, or, if the bases have equal quality, the one from the read with highest mean quality. E) If insert-match fails (for example, with an adapter dimer) adapter-match is performed. Reads that are too short after trimming are discarded.

47 for paired-end data. Additionally, hybrid tools are available that optimize their choice of read trimming
48 method based on the type of data. Skewer (Jiang et al., 2014) is a fast and accurate hybrid trimmer
49 that works with both single-end and paired-end data. However, choosing a read-trimming tool currently
50 requires a trade-off between feature set, efficiency, and accuracy. Furthermore, even state-of-the-art tools

51 still have a relatively high rate of over-trimming (removing usable template bases from reads) and/or
52 under-trimming (leaving low-quality and adapter-derived bases in the read sequence) (Sturm et al., 2016).

53 We sought to develop a read trimming tool that would combine the best aspects of currently available
54 software to provide high speed and accuracy while also offering a rich feature set. To accomplish this aim,
55 we used Cutadapt as a starting point, as it provides the broadest feature set of currently available tools and
56 is published under the MIT license, which allows modification and improvement of the code. We focused
57 on making three specific improvements to Cutadapt: 1) improve the accuracy of paired-end read trimming
58 by implementing an insert-match algorithm; 2) improve the performance by adding multiprocessing
59 support (as it is currently only able to use a single processor); and 3) add important additional features
60 such as automated trimming of Methyl-Seq reads and automated detection of adapter sequences in reads
61 where the experimental protocols are not known to the analyst. Because these modifications required
62 substantial changes to the Cutadapt code base, and because there are software tools that depend on the
63 current implementation of Cutadapt, we chose to “fork” the Cutadapt code base and develop our software,
64 Atropos, as a separate tool. Here, we show that we have accomplished our three aims. In addition to
65 extending the already rich set of features provided by the original Cutadapt tool, Atropos demonstrates
66 paired-end read trimming accuracy that is superior to other state-of-the-art tools, and it is among the fastest
67 read trimming tools when a moderate number of parallel execution threads are used (4). Furthermore,
68 Atropos achieves a performance increase that is roughly linear with the number of threads used, making it
69 the fastest tool when 8 or more threads are available.

70 2 MATERIALS AND METHODS

71 2.1 Implementation

72 Atropos is developed in Python and is available to install from GitHub or via the pip package manager
73 (see Data Availability).

74 2.1.1 Semi-global Alignment

75 Traditionally, the overlap between two sequences is detected by computing an optimal semi-global
76 alignment (Gusfield, 1997, Section 11.6.4), which is the same as global alignment except that neither
77 initial nor trailing gaps are penalized. This allows the sequences to shift relative to each other. An optimal
78 semi-global alignment maximizes the sum of alignment column scores, thus tending to favor longer over
79 short overlaps. Since score-based optimization is often not intuitively understood, the adapter alignment
80 algorithm uses edit operations instead, which has the advantage that it gives the user the ability to specify
81 a “maximum error rate” as a intuitive parameter. For a given alignment between read and adapter, the
82 error rate is computed as the number of edits (mismatches, insertions, deletions) divided by the length of
83 the matching part of the adapter. Minimizing the edit distance while at the same time not penalizing end
84 gaps would lead to optimal but meaningless zero-length overlaps; thus, a hybrid approach is chosen. The
85 adapter alignment algorithm computes edit distances for all allowed shifts of the adapter relative to the
86 read. Among those having an error rate not higher than the specified threshold, the shift (and therefore
87 alignment) with the highest number of matches is chosen.

88 We summarize the algorithm here; see (Martin, 2013, Section 2.2) for details. Let s be the adapter, t the
89 read and $m = |s|$, $n = |t|$. Adapter alignment computes edit distances $D(i, j)$ between the i -length prefix of
90 s and the j -length prefix of t for all $i = 0, \dots, m$ and $j = 0, \dots, n$ with the standard dynamic-programming
91 (DP) recurrence

$$D(i, j) = \min\{D(i-1, j-1) + [s_i \neq t_j], D(i-1, j), D(i, j-1)\} \quad (1)$$

92 The base cases are $D(i, 0) = 0$ or $D(i, 0) = i$ and $D(0, j) = 0$ or $D(0, j) = j$, depending on the adapter
93 type, allowing to skip a prefix of s and/or t at no cost. The algorithm additionally keeps track of $M(i, j)$,
94 which is the number of matches between the prefixes of s and t , and of the “origin” $O(i, j)$, which is
95 the number of skipped characters in t in the optimal alignment (if negative, characters in s are skipped
96 instead). All three DP matrices D, M, O are filled in at the same time, after which the cells of the bottom
97 row ($i = m$) are inspected. They represent possible end positions of the adapter sequence within the
98 read. For each position j , the error rate is computed from $D(m, j)$ and $O(m, j)$, and positions with a too
99 high error rate are discarded. If positions remain, the one with the highest number of matches $M(m, j)$

is returned as the position J of the adapter sequence. Together with the start of the adapter sequence at $O(m, J)$, the adapter sequence can then be removed from the read.

Observing that no backtrace within the DP matrix is required, the actual implementation keeps only a single column of the matrices in memory for better cache locality. Significant runtime improvements are achieved by employing the optimization described by Ukkonen (Ukkonen, 1985) of stopping the computation of a column as soon as the costs are too high and provably cannot decrease for the remainder of the column. When the user supplies an anchored adapter and disables insertions and deletions (indels) at the same time, the algorithm also switches to a much simpler variant that computes only the Hamming distance between the adapter and a prefix or suffix of the read.

2.1.2 Insert Match Algorithm

For each read pair, the insert-match algorithm uses the same semi-global alignment algorithm described above (with indels disabled) to find all possible alignments between the first read and the reverse complement of the second read that satisfy user-specified specificity thresholds (Figure 1C). Specificity is determined by the combination of up to three user-configurable thresholds: 1) minimum number of overlapping bases, 2) maximum number of mismatch bases, and 3) random mismatch probability (Sturm et al., 2016). The probability of a random match at k bases out of the n bases being compared is computed using the binomial distribution:

$$P = \sum_{i=k}^n \frac{n!}{i!(n-i)!} p^i (1-p)^{n-i} \quad (2)$$

The candidate alignments are tested in order of decreasing length until one is found in which the overhanging sequences on either end match the user-specified adapter sequences. Comparison between the adapter and overhang sequences is done using a constrained adapter-match approach. Briefly, starting at the end of the insert overlap, a pairwise comparison is made between the adapter and the read at each possible offset. The offset that best satisfies the user-configurable specificity thresholds (the same three described above) is taken to be the location of the adapter sequence, and all bases from that position to the 3' end of the read are removed. If an adapter is only found in one of the two reads, then the same offset is used to trim both reads, under the assumption that the location of the adapter sequence must be symmetric across the read pair.

Optionally, the overlapping inserts can be used for mutual error correction (Figure 1D). Where the aligned inserts have mismatches, the base with the highest quality score is chosen as the consensus. When the bases have equal quality, there is an option to leave the bases unchanged, convert them both to N, or to choose the base from the read with the highest mean quality as the consensus. There are additional options to 1) completely overwrite one read in the pair if its quality is very poor, and/or 2) merge the overlapping read pair into a single read, which avoids double-counting overlapping read pairs in read depth-based analyses.

If no insert match is found, or if an adapter is not found in an overhang, then an unconstrained adapter-match approach is attempted separately in each read (Figure 1E).

2.1.3 Parallel processing

The performance improvements in Atropos relative to Cutadapt and other read trimming tools are based in two observations: 1) each read (or read pair) is trimmed separately, and thus trimming can be parallelized across multiple processor cores, and 2) a significant fraction of the execution time is spent decompressing input files and re-compressing results. Compression of sequencing data is increasingly becoming necessary due to the large volumes of data generated in sequencing experiments.

To address the first bottleneck, we implemented a parallel processing pipeline based on the Python multiprocessing module. Briefly, a single thread is dedicated to a "reader" process that loads reads (or read pairs) from input file(s), with support for a variety of data formats and automatic decompression of compressed data. Reads are loaded in batches, and each batch is added to an in-memory queue. A user-specified number of "worker" threads (which is constrained by the number of processing cores available on the user's system) extract batches from the queue and perform trimming and filtering operations on the reads in the same manner as Cutadapt. Atropos addresses the second bottleneck by offering a choice of three modes for writing the results to disk. The first two modes involve adding the results to a second in-memory queue, from which a dedicated "writer" process extracts batches and performs

the serialized write operation. These two modes differ in how the trimmed reads are compressed – in “worker compression” mode, each worker is responsible for compressing the results using the Python gzip module prior to placing the results on the queue, whereas in “writer compression” mode, the writer process performs compression using the much faster system-level gzip program. The choice between these two modes is selected automatically based on the number of worker threads used, with worker compression becoming faster than writer compression when at least 8 threads are available. The third output mode, called “parallel writing,” does not use a dedicated writer process (and thus an additional worker process can be used in its place). Instead, each worker process writes its results to a separate file. This can dramatically reduce the execution time of the program (~50% reduction in our experiments; see Results) and is generally compatible with downstream analysis since many mapping and assembly tools accept multiple input files (and for those that don’t, gzipped files can be safely concatenated without needing to be decompressed and recompressed). An additional speed-up is gained by recognizing that the reader process often finishes loading data well before the worker processes finish processing it; thus, an additional worker thread is started as soon as the reader process completes.

2.1.4 Adapter detection

Often, details of sequencing library construction are not fully communicated from the individual or facility that generated the library to the individual(s) performing data analysis. Manual determination of sequencing adapters and other potential library contaminants can be a tedious and error-prone task. Thus, we implemented in Atropos a new ‘detect’ command that automatically identifies adapters/contaminants from a sample of read sequences. First, a profile is built of k -mers (where k is a fixed number of consecutive nucleotides, defaulting to $k = 12$) within N read sequences (where N defaults to 10,000). When at least 8 consecutive A bases are detected, those bases along with all subsequent bases (in the 3’ direction) are first trimmed, as that pattern is a strong indicator that the sequencer scanned past the end of the template (i.e. the length of the fragment + adapter is less than the read length; Figure 1E). Additionally, low-complexity reads are excluded, where complexity $X(S)$ is defined as follows. Let $C(i, S)$ be the number of elements of a nucleotide sequence $S = s_1, \dots, s_n$, that are nucleotide $i \in A, C, G, T$.

$$X(S) = - \sum \frac{C(i, S) * \log(C(i, S))}{\log(2)} \quad (3)$$

Sequences with $X(S) < 1.0$ are defined as low-complexity. All remaining k -mers are counted, and each k -mer is linked to all of the sequences from which it originated. This process continues iteratively for increasing values of k , with only those read sequences linked to high-abundance k -mers in the previous iteration being used to build the k -mer profile in the next iteration. k -mer K is considered high-abundance when:

$$|K| > \frac{N * (l - k + 1) * O}{4^k} \quad (4)$$

where l is the read length and $O = 100$ by default. Finally, high-abundance k -mers of all lengths are merged to eliminate shorter sequences that are fully contained in longer sequences. Optionally, the high-abundance k -mers are matched to a list of known adapters/contaminants. Atropos reports to the user an ordered list of up to 20 of the most likely contaminants.

2.1.5 Error Rate Estimation

Quality and adapter trimming is sensitive to the choice of several parameters. For example, relative to data sets with typical rates of sequencing error, data sets with higher error-rates require higher thresholds for mismatches and/or random-match probability during insert- and adapter-matching to perform with the same level of sensitivity. Thus, we implemented in Atropos a new ‘error’ command that provides an estimate of the error rate in each input file. The error command gives the choice between two algorithms: 1) averaging all base qualities across a sample of reads, which is fast but likely overestimates the true rate of sequencing error (Dohm et al., 2008; DePristo et al., 2011); and 2) the shadow regression method proposed by Wang et al. (Victoria Wang et al., 2012), which more accurately estimates error rates at the cost of reduced speed and greater memory usage.

2.1.6 Shared Cutadapt and Atropos Improvements

In addition to improvements in the semi-global alignment algorithm above, Atropos also benefits from the following improvements that were made to Cutadapt subsequent to the publication of Martin *et al.* (2011), but prior to the Atropos fork, and are therefore features in both programs.

- Adapters can now be *anchored*, which removes their ability to occur anywhere within a read. An anchored 5' adapter thus matches only if it is a prefix of the read, and a 3' adapter only if it is a suffix of the read. This is useful, for example, when one or both sequencing adapters are known to be ligated directly to a PCR primer.
- *Linked adapters* combine a 5' with a 3' adapter. Trimming multiple adapters from each read was also supported previously, but linked adapters make it possible to require that one of them is a 5' adapter and one a 3' one.
- *IUPAC ambiguity codes* are fully supported. Thus, adapter sequences containing characters such as N (matching any nucleotide), H (A, C, or T), Y (C or T) work as expected. They are useful when adapters contain barcodes or random nucleotides. The nucleotides and ambiguity codes are internally represented as patterns of four bits, in which each set bit corresponds to an allowed nucleotide. Comparisons are thus simple "binary and" operations, resulting in no runtime overhead.
- *Paired-end data* can be trimmed with sequences specified for the forward and reverse reads independently. Read pairs are guaranteed to remain in sync. Even *interleaved* data (paired-end reads in a single file) is accepted.
- Quality trimming can now work in a *NextSeq-specific* mode in which spurious runs of high-quality G nucleotides at the 3' end of a read are correctly trimmed. NextSeq instruments use "dark" or "black" cycles for G nucleotides, making them unable to distinguish between regular G and reaching the end of the fragment.
- Other additions include support for *trimming a fixed number of bases* from a read, support for files compressed using the bzip2 and lzma algorithms, and improved filtering options.

2.2 Benchmarks

Data Set	Error Rate*	Read Length	Total Read Pairs	Reads w/ Adapters**	Adapter Bases**
Simulated 1	0.20%	125	781,923	325,982	12,447,262
Simulated 2	0.60%	125	780,899	325,105	12,416,861
Simulated 3	1.20%	125	782,237	325,860	12,464,235
GM12878 WGBS	2.79%	125	1,000,000	57,130	3,082,003
K562 mRNA-seq	4.31%	75	6,100,265	14,384	749,451

Table 1. Description of data sets. For the real data sets, * actual error rates are unknown – we estimate error rates from base qualities over a sample of 10,000 read pairs; and ** total adapter content is unknown – we provide the number of reads containing exact matches for the first 35 adapter bases, and the number of adapter bases present.

We evaluated both the speed and the accuracy of Atropos relative to other state-of-the-art read trimming tools using both simulated and real-world data (Table 1). As trimming of single-end reads is unchanged from the original Cutadapt method and is also decreasing in relevance as most current experiments use paired-end data, we focused our benchmarking on trimming of paired-end reads. Sturm et al. demonstrate that Skewer (Jiang et al., 2014) and SeqPurge (Sturm et al., 2016) stand out as having superior performance in paired-end read trimming; thus, we chose to benchmark Atropos against these tools. We also compared the new insert-match algorithm against the adapter-match algorithm that is used by Cutadapt, and which can be enabled in Atropos using a command line option.

To simulate paired-end read data, we use the ART simulator (Huang et al., 2012) that was modified by Jiang et al. to add adapter sequences to the ends of simulated fragments. ART simulates reads based on empirically derived quality profiles specific to each sequencing platform. A quality profile consists

of distributions of quality scores for each nucleotide at each read position, expressed as read counts aggregated from multiple sequencing experiments, where quality scores are in Phred scale ($-10\log_{10}(e)$, where e is the probability that the base call is erroneous). We developed an *R* script to artificially inflate the error rates in an ART profile to a user-defined level. For each row in the profile with quality score bins $e_1..e_n$ and corresponding read counts $r_1..r_n$, the overall error rate can be computed as:

$$E = \frac{\sum_{i=1}^n e_i r_i}{\sum_{i=1}^n r_i} \quad (5)$$

We use the *R* function *optim* with the variable metric ("BFGS") algorithm to optimize a function that adds an equal number of counts C to each Phred-score bin in the distribution and then compares the overall error rate to the user-desired error rate E' :

$$f(C, E') = \frac{\sum_{i=1}^n e_i (r_i + C)}{\sum_{i=1}^n (r_i + C)} - E' \quad (6)$$

We simulated ~800k 125 bp paired-end reads using the Illumina 2500 profile at error rates that were low/typical (~0.2%, the unmodified profile), intermediate (~0.6%), and high (~1.2%). We evaluated the accuracy of the tools on the simulated data by comparing each trimmed read pair to the known template sequence. We counted the frequency of following outcomes: the fragment does not contain adapters but is trimmed anyway ("wrongly trimmed"), the fragment is under-trimmed, or the fragment is over-trimmed. We also counted the total number of under- and over-trimmed bases.

We also benchmarked the tools on two real-world data sets. First, we sampled ~1M read pairs from a whole-genome bisulfite sequencing (WGBS) library generated from the GM12878 cell line. Second, we used 6.1M paired-end mRNA-seq reads generated from the K562 cell line. Both of these data sets were generated by the ENCODE project (ENCODE Project Consortium, 2012). Since the genomic origins of the templates are not known *a priori*, we instead compared the read trimming tools based on improvement in the results of mapping the trimmed versus untrimmed reads. We used STAR (Dobin et al., 2013) to map the mRNA-seq reads to GRCh37, and we used bwa-meth (Pedersen et al., 2014) to map the WGBS reads to the bisulfite-converted GRCh37. We also compared the results of only adapter trimming to the results of adapter trimming plus additional quality trimming using a minimum quality threshold of 20 (Phred-scale).

Although sequence analysis is sometimes performed using a desktop computer, analysis of the volumes of data currently being generated increasingly requires the use of high-throughput computing facilities ("clusters"). The hardware architecture of a cluster is often different from that of a desktop computer. Most importantly, storage in a cluster is typically centralized and accessed by the compute nodes via high-speed networking. Such an architecture inevitably adds latency to file reading and writing operations ("I/O"). Cluster nodes also typically have more processing cores and memory available than desktop computers. This means that the performance of software with intensive I/O usage (such as read trimming) is likely to be quite different on a desktop versus a cluster. To examine the impact of these architectural differences, we ran all benchmarks on both a desktop computer (a Mac Pro) having a 3.7 GHz quad-core Intel Xeon E5 processor and 32 GB RAM, and on a cluster node having 4 quad-core 2.2 GHz AMD Opteron processors and 128 GB memory, and with all data being read from and written to network-accessible storage over a 1 Gbit ethernet connection.

Finally, we provide scripts in our GitHub repository that can be used to re-run our analysis, and that also enable other tool developers to benchmark their software against Atropos.

3 RESULTS

3.1 Simulated Data

3.1.1 Performance

On a desktop computer with 4 processing cores, we found that SeqPurge had the fastest overall execution time, followed closely by Skewer and Atropos in parallel write mode (Table 2).

As expected, execution times on a cluster node using 4 threads were approximately 70% greater than those observed on a desktop computer (Table 3). We expect that much of this disparity is due to the

Program	Execution Time (sec.)	
	Min	Max
Atropos (adapter + worker compr.)	185.74	195.83
Atropos (adapter + writer compr.)	55.25	57.10
Atropos (adapter + parallel write)	30.67	32.06
Atropos (insert + worker compr.)	187.04	199.11
Atropos (insert + writer compr.)	56.01	57.96
Atropos (insert + parallel write)	35.23	41.17
SeqPurge	26.72	27.14
Skewer	29.49	30.10

Table 2. Execution time on simulated datasets for programs running on desktop with 4 parallel processes (threads). Each program was executed multiple times, and Atropos was run with combinations of alignment algorithm (insert-match or adapter-match) and output mode (writer-compression or parallel-write). The minimum and maximum execution times for each program are shown, with the lowest overall execution time in bold.

increased latency involved in network-based I/O on the cluster, although some may also be explained by CPU differences (3.7 GHz Intel on the desktop versus 2.2 GHz on the cluster node).

When increasing the number of parallel execution threads from 4 to 8 and 16, Atropos achieves a roughly linear decrease in execution time, and the execution time of SeqPurge decreases to a lesser extent. Interestingly, the execution time of Skewer actually increases with an increasing number of threads. With 8 and 16 threads, Atropos in parallel-write mode is the fastest of the tools.

On the other hand, Atropos uses substantially more memory than SeqPurge or Skewer (Table 3). We expect this is partially due to overhead of automatic memory management in Python compared to C++ (in which SeqPurge and Skewer are implemented), but in larger part results from Atropos' use of in-memory queues to communicate between parallel processes. For all three programs, memory usage is independent of the number of reads processed. We note that Atropos provides parameters to limit memory usage (although typically at the expense of reduced speed).

3.1.2 Accuracy

We found that the four trimming algorithms had different biases toward under- and over-trimming (Table 4). Across the three sequencing error rates, Skewer had the lowest frequency of wrongly trimming reads while SeqPurge had the highest. The Atropos adapter-match algorithm and Skewer had similarly low frequencies of over-trimming, while the Atropos insert-match algorithm and SeqPurge had similarly low frequencies of under-trimming. Overall, the Atropos insert-match algorithm and SeqPurge demonstrated the lowest error rates (0.01%).

In terms of numbers of over- and under-trimmed bases, the Atropos insert-match algorithm and SeqPurge clearly had the best performance (Table 4) at all sequencing error rates. The two algorithms had similarly low numbers of under-trimmed bases, but the Atropos insert-match algorithm had a lower number of over-trimmed bases, giving it the lowest overall error rate (0.001%). On the other hand, Skewer and the Atropos adapter-match algorithm left substantial numbers of under-trimmed bases, resulting in about 10-fold higher overall error rates.

Additionally, we found that all tools discarded very similar numbers of reads (~1.8%) that were below the minimum length threshold of 25 bp after trimming. These were reads with short insert sizes, which have a high rate of spurious mapping, and thus it is common practice to discard them.

3.2 Real Data

We first tested Atropos' adapter detection module on the real data sets. Using the first 10,000 reads in each pair of FASTQ files, Atropos correctly detected the exact sequences of the adapters used in constructing each library. For 3 of the 4 adapters, the sequences were found in a list of known contaminants (WGBS read 1: "TruSeq Adapter, Index 7"; WGBS read 2 and RNA-seq read 2: "TruSeq Universal Adapter"); the RNA-seq read 1 adapter appears to have a custom-designed sequence.

Program	Execution Time (Min Max sec.)					
	4 Threads	8 Threads		16 Threads		
Atropos (adapter + worker compr.)	161.33	163.83	97.93	108.04	53.92	65.39
Atropos (adapter + writer compr.)	84.41	88.11	112.77	114.73	112.27	157.59
Atropos (adapter + parallel write)	51.73	63.61	33.71	46.77	32.88	68.63
Atropos (insert + worker compr.)	171.34	186.00	105.32	110.24	57.40	99.96
Atropos (insert + writer compr.)	87.43	88.59	111.01	114.81	110.53	123.94
Atropos (insert + parallel write)	64.12	75.55	41.22	51.80	34.24	37.24
SeqPurge	43.46	51.36	42.01	46.06	35.69	38.65
Skewer	43.10	55.20	59.00	66.34	54.36	67.79
Program	Memory Usage (Min Max Gb)					
	4 Threads	8 Threads		16 Threads		
Atropos (worker compr.)	1.019	1.5	0.528	0.730	0.760	1.291
Atropos (writer compr.)	1.884	1.992	1.737	1.951	1.336	1.777
Atropos (parallel write)	1.273	1.680	0.448	1.138	0.571	0.705
SeqPurge	0.013	0.013	0.013	0.013	0.013	0.015
Skewer	0.008	0.008	0.009	0.009	0.013	0.014

Table 3. Top: Execution time on simulated datasets for programs running on a cluster with 4, 8, or 16 parallel processes (threads). Each program was executed multiple times, and Atropos was run with all combinations of alignment algorithm (insert-match or adapter-match) and output mode (worker-compression, writer-compression or parallel-write). The minimum and maximum execution times for each program are shown, and notable results are highlighted in bold. Bottom: The minimum and maximum memory usage for each program across all executions.

Program	Reads				Bases		
	Wrongly Trimmed	Over-trimmed	Under-trimmed	Total Error	Over-trimmed	Under-trimmed	Total Error
Error rate 0.2%							
Atropos							
adapter-match	51 (0.01%)	1 (0.00%)	28,991 (9.72%)	3.71%	490	102,133	0.054%
insert-match	134 (0.03%)	26 (0.01%)	27 (0.01%)	0.01%	1,036	66	0.001%
SeqPurge	2,639 (0.55%)	37 (0.01%)	29 (0.01%)	0.01%	7,708	29	0.004%
Skewer	10 (0.00%)	7 (0.00%)	283 (0.09%)	0.04%	21	22,029	0.012%
Error rate 0.6%							
Atropos							
adapter-match	72 (0.01%)	6 (0.00%)	28,843 (9.68%)	3.69%	733	101,839	0.054%
insert-match	168 (0.03%)	14 (0.00%)	31 (0.01%)	0.01%	1,434	62	0.001%
SeqPurge	2,595 (0.54%)	26 (0.01%)	38 (0.01%)	0.01%	6,701	41	0.004%
Skewer	6 (0.00%)	3 (0.00%)	413 (0.14%)	0.05%	9	36,685	0.019%
Error rate 1.2%							
Atropos							
adapter-match	76 (0.02%)	5 (0.00%)	30,152 (10.12%)	3.86%	721	117,027	0.062%
insert-match	175 (0.04%)	12 (0.00%)	46 (0.02%)	0.01%	1,440	110	0.001%
SeqPurge	2,652 (0.55%)	24 (0.01%)	34 (0.01%)	0.01%	7,628	36	0.004%
Skewer	7 (0.00%)	5 (0.00%)	644 (0.22%)	0.08%	12	67,066	0.035%

Table 4. Trimming accuracy on simulated data with three different base-call error rates.

3.2.1 Performance

We performed adapter trimming on the real datasets in a cluster environments using 16 parallel cores. Again, we found that Atropos in parallel-write mode performed the best (Table 5). For the WGBS dataset, SeqPurge, Skewer, and Atropos in worker-compression mode were ~50% slower than Atropos in parallel-write mode. For the mRNA-seq data set, Atropos in writer-compression mode was nearly as fast

Program	Trimming Time (sec.)		Trimming Mem. (Gb)		Mapping Time (min.)	
	Q0	Q20	Q0	Q20	Q0	Q20
WGBS Data						
Untrimmed reads					37:59	
Atropos (worker compr.)	60.71	61.84	1.665	0.839	32:10	28:46
Atropos (writer compr.)	85.38	91.06	1.026	0.842		
Atropos (parallel write)	38.28	39.64	0.664	1.440		
SeqPurge	56.20	60.54	0.011	0.012	33:00	25:22
Skewer	69.34	69.83	0.014	0.014	31:11	26:20
mRNA-seq Data						
Untrimmed reads					03:34	
Atropos (worker compr.)	294.24		0.800		03:23	
Atropos (writer compr.)	257.17		0.586			
Atropos (parallel write)	254.01		0.466			
SeqPurge	351.46		0.011		02:14	
Skewer	279.08		0.014		02:20	

Table 5. Execution and read mapping times for programs trimming real data on a cluster with 16 threads. Each program was executed with no additional quality trimming (Q0); for WGBS data, we also performed quality trimming at a minimum base quality of 20 (Q20). Atropos was run with the insert-match algorithm. SeqPurge and Atropos were run with error correction enabled (Skewer performs error correction by default).

as parallel-write mode, while skewer was ~10% slower. We also performed read mapping on the cluster with 16 cores (Table 5). In all cases, we observed that the dataset with the least effective improvement over untrimmed reads (Figures 2 and 3) had the fastest mapping time, indicating that mapping time is not a valid proxy for overall trimming tool performance.

3.2.2 Effectiveness

We assessed read trimming effectiveness in practical terms. For the WGBS data, we computed the number of trimmed reads mapped at at given quality (MAPQ) cutoff, relative to the number of untrimmed reads mapped at that cutoff. We found that trimming by Atropos resulted in the greatest increase in number of mapped reads at all quality cutoffs (Figure 2). We found that trimming with SeqPurge resulted in similar, but smaller, gains in mapping quality, while trimming with Skewer resulted in the smallest gains in quality.

We also found that additional quality trimming in addition to adapter trimming has a substantial negative effect on read mapping, at least for bisulfite reads mapped using bwa-meth. Quality trimming by Skewer had the least negative effect on mapping quality of the three programs, and quality trimming by SeqPurge had the greatest negative effect on mapping quality.

For the mRNA-seq data, we additionally compared each alignment to GENCODE (v19) gene annotations (Harrow et al., 2012) to determine the number of reads mapped to expressed regions of the genome. We found that trimming with Atropos resulted an equal or greater number of mapped reads aligned to expressed regions compared to the other tools at all MAPQ thresholds (Figure 3).

4 CONCLUSIONS

Our results demonstrate that adapter trimming tools are approaching optimal accuracy, at least for the (currently) most common type of data – paired-end short reads with 3' adapters. On synthetic data with varying error rates, Atropos (using our new insert-match algorithm) and SeqPurge both demonstrated overall error rates of 0.01% at the read level, and Atropos has the lowest base-level error rate of 0.001%.

On real WGBS and mRNA-seq data, we found that adapter trimming with Atropos resulted in the greatest increase in read mapping quality. We also found that stringent quality trimming has a negative effect on WGBS read mapping quality, at least when using bwa-meth as the alignment tool. For reads trimmed with a quality threshold of 20, all mapping statistics were worse than those for untrimmed reads.

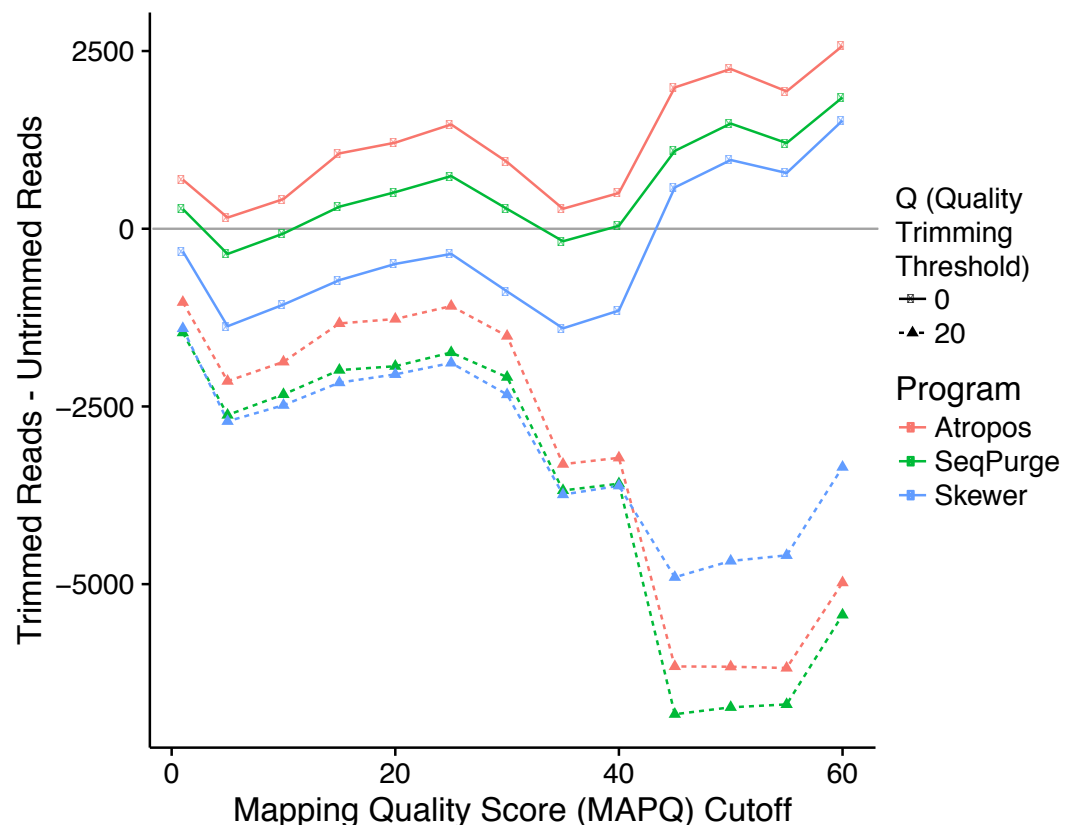


Figure 2. Atropos trimming best improves mapping of real WGBS sequencing reads. Reads were adapter-trimmed with all three programs, both without additional quality trimming (Q=0) and with quality trimming at a threshold of Q=20. We mapped both untrimmed and trimmed reads to the genome. For each MAPQ cutoff on the x-axis (1..60 at intervals of 5), the number of reads with MAPQ greater than or equal to the cutoff less the number of untrimmed reads with MAPQ greater than or equal to the cutoff is shown on the y-axis for each program.

In terms of performance, SeqPurge had the best performance of the three tools when only 4 threads were available, while Atropos and SeqPurge had superior performance on our cluster environment when there were at least 8 threads available. In the later case, Atropos in parallel-write mode had the best overall performance on all datasets, with the trade-off that parallel-write mode produces a larger number of data files, which may make analyses of large projects more complicated to manage. SeqPurge performed better than Skewer and Atropos in worker- and writer-compression mode on the smaller and lower-error WGBS dataset, while Skewer performed better than SeqPurge and Atropos in worker-compression mode on the larger and higher-error mRNA-seq dataset. Atropos' memory requirements were the highest among the three programs (0.5-2 GB versus 10-20 Mb), but well within the capabilities of most modern computer systems.

In summary, our results show that Atropos offers the best combination of accuracy and performance of the tools that we evaluated. Furthermore, Atropos has the richest feature set of the three tools, including Methyl-Seq-specific trimming options, automated adapter detection, estimation of sequencing error, and support for data generated by many sequencing methods (ABI SOLiD, Illumina NextSeq, mate-pair libraries, and single-end sequencing). Finally, Atropos is easily installable on any system with Python 3.3+.

5 DATA AVAILABILITY

- Atropos can be installed from the Python Package Index (pypi) using the pip tool: 'pip install atropos'.

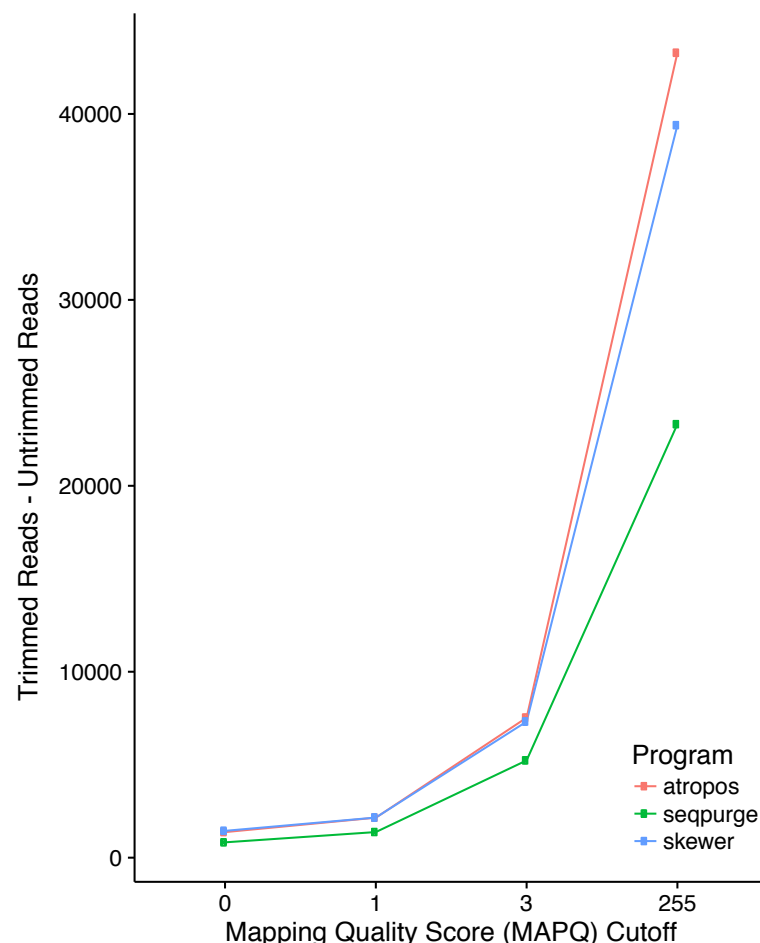


Figure 3. Atropos trimming results in the greatest increase in mRNA-seq reads mapped to expressed regions. Reads were adapter-trimmed with all three program without additional quality trimming. We mapped both untrimmed and trimmed reads to the genome using STAR. When parameter `outSAMmultNmax = 2`, STAR produces only four MAPQ values: 255=unique alignment, 3=two alignments with similar but unequal score; 1=two alignments with equal score; and 0=unmapped. For each of the four MAPQ cutoffs on the x-axis, the number of reads with MAPQ greater than or equal to the cutoff less the number of untrimmed reads with MAPQ greater than or equal to the cutoff is shown on the y-axis for each program.

- The Atropos source code, including all scripts needed to execute the analyses in this paper, are available at <https://github.com/jddidion/atropos>. The portions of Atropos developed by JPD are a work of the US government, and thus all copyright is waived under a CC0 1.0 Universal Public Domain Dedication (<https://creativecommons.org/publicdomain/zero/1.0/>).
- The simulated data sets are also available in the Atropos GitHub repository.
- The real K562 mRNA-seq data (accession SRR521458) is available from the NCBI Sequence Read Archive: <https://trace.ncbi.nlm.nih.gov/Traces/sra/?run=SRR521458>.
- The real GM12878 WGBS data (accession ENCLB794YYH) is available from the ENCODE project website: <https://www.encodeproject.org/experiments/ENCSR890UQO/>.
- We used human reference genome GRCh37, downloaded from ftp://ftp.ncbi.nlm.nih.gov/genomes/Homo_sapiens/ARCHIVE/BUILD.37.3.

- We used GENCODE v19 annotations, downloaded from ftp://ftp.sanger.ac.uk/pub/genocode/Gencode_human/release_19/genocode.v19.annotation.gtf.gz.

6 ACKNOWLEDGEMENTS

We thank Jim Mullikin for helpful comments on an earlier version of this manuscript. We also thank the users of Cutadapt and Atropos that have contributed bug reports and improvements. JPD and FSC are funded by the NIH intramural program. Additionally, JPD is funded by the American Diabetes Association (1-17-PDF-100). MM is supported by a grant from the Knut and Alice Wallenberg Foundation to the Wallenberg Advanced Bioinformatics Infrastructure.

REFERENCES

- Bock, C. (2012). Analysing and interpreting DNA methylation data. *Nature Reviews Genetics*, 13(10):705–719.
- Del Fabbro, C., Scalabrin, S., Morgante, M., and Giorgi, F. M. (2013). An extensive evaluation of read trimming effects on Illumina NGS data analysis. *PLoS One*, 8(12):e85024–None.
- DePristo, M. A., Banks, E., Poplin, R., Garimella, K. V., Maguire, J. R., Hartl, C., Philippakis, A. A., del Angel, G., Rivas, M. A., Hanna, M., McKenna, A., Fennell, T. J., Kernytsky, A. M., Sivachenko, A. Y., Cibulskis, K., Gabriel, S. B., Altshuler, D., and Daly, M. J. (2011). A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nature Genetics*, 43(5):491–498.
- Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M., and Gingeras, T. R. (2013). STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*, 29(1):15–21.
- Dohm, J. C., Lottaz, C., Borodina, T., and Himmelbauer, H. (2008). Substantial biases in ultra-short read data sets from high-throughput DNA sequencing. *Nucleic Acids Research*, 36(16):e105–e105.
- ENCODE Project Consortium (2012). An integrated encyclopedia of DNA elements in the human genome. *Nature*, 489(7414):57–74.
- Gusfield, D. (1997). *Algorithms on Strings, Trees and Sequences*. Cambridge University Press.
- Harrow, J., Frankish, A., Gonzalez, J. M., Tapanari, E., Diekhans, M., Kokocinski, F., Aken, B. L., Barrell, D., Zadissa, A., Searle, S., Barnes, I., Bignell, A., Boychenko, V., Hunt, T., Kay, M., Mukherjee, G., Rajan, J., Despacio-Reyes, G., Saunders, G., Steward, C., Harte, R., Lin, M., Howald, C., Tanzer, A., Derrien, T., Chrast, J., Walters, N., Balasubramanian, S., Pei, B., Tress, M., Rodriguez, J. M., Ezkurdia, I., van Baren, J., Brent, M., Haussler, D., Kellis, M., Valencia, A., Reymond, A., Gerstein, M., Guigó, R., and Hubbard, T. J. (2012). GENCODE: the reference human genome annotation for The ENCODE Project. *Genome research*, 22(9):1760–1774.
- Huang, W., Li, L., Myers, J. R., and Marth, G. T. (2012). ART: a next-generation sequencing read simulator. *Bioinformatics*, 28(4):593–594.
- Jiang, H., Lei, R., Ding, S.-W., and Zhu, S. (2014). Skewer: a fast and accurate adapter trimmer for next-generation sequencing paired-end reads. *BMC Bioinformatics*, 15:182–None.
- Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMB-net.journal*, 17(1):pp. 10–12.
- Martin, M. (2013). *Algorithms and Tools for the Analysis of High-Throughput DNA Sequencing Data*. PhD thesis, TU Dortmund.
- Pedersen, B. S., Eyring, K., De, S., Yang, I. V., and Schwartz, D. A. (2014). Fast and accurate alignment of long bisulfite-seq reads. *arXiv:1401.1129 [q-bio]*. arXiv: 1401.1129.
- Sturm, M., Schroeder, C., and Bauer, P. (2016). SeqPurge: highly-sensitive adapter trimming for paired-end NGS data. *BMC Bioinformatics*, 17:208.
- Ukkonen, E. (1985). Finding approximate patterns in strings. *Journal of Algorithms*, 6(1):132–137.
- Victoria Wang, X., Blades, N., Ding, J., Sultana, R., and Parmigiani, G. (2012). Estimation of sequencing error rates in short reads. *BMC Bioinformatics*, 13:185.