# Atria: An Ultra-fast and Accurate Trimmer for Adapter and Quality Trimming

Jiacheng Chuan [1,2], Aiguo Zhou[1,3], Lawrence Richard Hale [2], Miao He[4], Xiang Li [1,*]

[1] *Canadian Food Inspection Agency, Charlottetown, PE C1A5T1, Canada*

[2] *Department of Biology, University of Prince Edward Island, Charlottetown, PE C1A4P3, Canada*

[3] *Guangdong Laboratory for Lingnan Modern Agriculture, South China Agricultural University, Guangzhou 510642, China*

[4] *School of Life Sciences, Sun Yat-sen University, Guangzhou, Guangdong, 510275, China*

E-mail:

jiacheng.chuan@inspection.gc.ca (Chuan J)

aiguozhou@scau.edu.cn (Zhou A)

lhale@upei.ca (Hale L)

lsshem@mail.sysu.edu.cn (He M)

sean.li@inspection.gc.ca (Li X)

* Corresponding author: Xiang Li, sean.li@inspection.gc.ca

**Abstract**

22    **Background:** As Next Generation Sequencing takes a dominant role in terms of output

23    capacity and sequence length, adapters attached to the reads and low-quality bases

24    hinder the performance of downstream analysis directly and implicitly, such as

25    producing false-positive single nucleotide polymorphisms (SNP), and generating

26    fragmented assemblies. A fast trimming algorithm is in demand to remove adapters

27    precisely, especially in read tails with relatively low quality.

28    **Findings:** We present a trimming program named Atria. Atria matches the adapters in

29    paired reads and finds possible overlapped regions with a super-fast and carefully

30    designed byte-based matching algorithm ($O(n)$ time with $O(1)$ space). Atria also

31    implements multi-threading in both sequence processing and file compression and

32    supports single-end reads.

33    **Conclusions:** Atria performs favorably in various trimming and runtime benchmarks

34    of both simulated and real data with other cutting-edge trimmers. We also provide an

35    ultra-fast and lightweight byte-based matching algorithm. The algorithm can be used in

36    a broad range of short-sequence matching applications, such as primer search and seed

37    scanning before alignment.

38    **Availability & Implementation**: The Atria executables, source code, and benchmark

39    scripts are available at https://github.com/cihga39871/Atria under the MIT license.

40

41    **Research Area:** Software and Workflows

42    **Classifications:** Bioinformatics, Software Engineering

43

## Statement of Need

**Background**

46  Next generation sequencing (NGS) is a revolutionary new technology that produces

47  massive, high-resolution genome sequence data to facilitate a broad range of biological

48  applications. Illumina paired-end sequencing can read a DNA fragment from both ends

49  and generate accurate reads for downstream bioinformatics analysis, such as assembly,

50  resequencing, transcriptome profiling, variant calling, epigenome profiling, chromatin

51  interaction, and chromosomal rearrangements [1, 2].

52      In paired-end library preparation, adapter sequences are the technical sequences

53  ligated to both sides of inserts, which are the DNA fragments of interest. Then, DNA

54  molecules with adapters are sequenced from both ends of the inserts so paired-end reads

55  are generated. If insert sizes of paired-end reads are less than the read lengths, inserts

56  are reversely complementary, and adapters are sequenced after reading through the

57  inserts (**Fig. 1**). Thus, adapter contamination in the 3' end needs to be removed before

58  downstream analysis.

59      Cleaning adapters can therefore be achieved by searching adapter sequences and/or

60  aligning paired reads (**Fig. 1**). To date, some trimmers, such as AdapterRemoval [3],

61  Trim Galore [4], and Trimmomatic [5], use both types of information to clean adapters.

62  However, when the quality of sequencing reads decreases, the trimming process

63  employing both types of information is likely to give different trimming suggestions.

64    Trimmers thus face a bottleneck when working on trimming adapters at accurate

65    positions. Also, extremely short adapters at the low-quality 3' end are sometimes

66    difficult to detect. Thus, trade-offs between trimming truncated adapters, and retaining

67    inserts intact, become necessary.

68    These two issues hinder trimmers from cleaning adapter sequences and leaving

69    DNA inserts intact. To combat this, we launch Atria, an integrated trimming program

70    for NGS data. Atria uses a super-fast byte-based matching algorithm to detect adapters

71    and reverse complementary regions of paired reads, and integrates carefully designed

72    decision rules to infer true adapter positions. Thus, Atria can trim extremely short

73    adapter sequences at accurate positions and not over-trim reads without adapters (**Fig.**

74    **1**).

75    In addition to adapter trimming, Atria integrated a set of trimming and filtering

76    methods, such as consensus calling for overlapped regions, quality trimming,

77    homopolymer trimming, N trimming, hard clipping from both ends, and read

78    complexity filtration.

79

80    **Implementation**

81    The adapter finding algorithms used in Atria can be categorized in the following

82    portions: DNA encoding, matching algorithm, matching and scoring, decision rules,

83    consensus calling, quality trimming, and IO optimization (**Fig. 2**).

84

85  *DNA encoding*

86  The DNA encoding algorithm is developed based on BioSequences, a Julia package

87  from BioJulia [6]. The original BioSequences package encodes DNA bases A, C, G, T

88  as four-bit codes 0001, 0010, 0100, 1000, respectively. Extended codes are also

89  supported, such as N (1111), S (0110), and gap (0000). DNA sequences are encoded

90  and stored in a contiguous block of Random-Access Memory (RAM) as a dense array

91  of unsigned 64-bit integers (UInt64) (**Fig. 2A**).

92  Atria makes use of the property of dense arrays to extract sequences as unsigned

93  integers from available memory locations. When accessing the last several indices of a

94  sequence, the extraction is illegal because operating systems do not allow the loading

95  of data outside of sequence boundary. To solve the issue, Atria constructs a bit-safe

96  sequence array, which elongates the sequence boundary by appending a UInt64 to the

97  end of the original array, and setting all bits after the end of encoded DNA to 0 (**Fig.**

98  **2A**).

99  It is noticeable that the smallest addressable unit of memory is one byte (8 bits)

100 while each DNA is encoded in four bits, so only the even indices of sequence can be

101 directly extracted (defining indices start from 0) (**Fig. 2A**). The extraction from odd

102 indices requires extra operations, which is avoidable in many scenarios of a well-

103 designed algorithm.

104    We denote a UInt64 extracted from the memory position $n$ of sequence $a$ by $a_n$. $a_n$

105    is a 16-mer and represents the subsequence of $a$ indexed from $2n$ to $2n+15$, which is

106    denoted by $a[2n:2n+15]$ (**Fig. 2A**).

107

108    *Matching algorithm*

109    Given two sequences $a$ and $b$, we plan to match the 16-base-long head of $a$ to each

110    index of $b$. However, only the even indices of $b$ can be extracted from memory without

111    bitwise operations, so we prepared two UInt64 of $a$: $a_0$ and $a_-$. $a_0$ is the 16-mer UInt64

112    loaded from the position 0 of $a$, and $a_-$ can be computed from the following bitwise

113    operations: ($a_0 >> 4 \,|\, a_1 << 4$). In this way, $a_0$ represents the subsequence of $a$ indexed

114    from 0 to 15 ($a[0:15]$), and $a_-$ represents $a[1:16]$ (**Fig. 2B**).

115    In this way, the problem of matching the 16-base-long head of $a$ to each index of $b$

116    is converted to the problem of matching two 16-mers, $a_0$, and $a_-$, to each addressable

117    memory position of $b$. The latter requires less bitwise operations.

118    The number of mismatches $K$ is computed in the formula:

$$K_{an,bn} = 16 - count\_ones(an \,\&\, bn)$$

120    where *count_ones* counts the number of ones in the binary representation of the UInt64.

121    Let $k$ denote the user-defined number of mismatches allowed in the 16-mer

122    comparison of UInt64 $a_n$ and $b_n$ ($k = 2$ by default). After matching $a_0$ and $a_-$ to each

123    addressable memory position of $b$, if the minimum number of mismatches is not greater

124    than $k$, the smallest index of $b$ of the minimum mismatches is reported.

125    Therefore, the complexity of the matching algorithm is $O(n)$ time with $O(1)$ space,

126    so its speed is extremely fast. One limitation is that when computing the number of

127    mismatches of $a_n$ and $b_n$, and if they have ambiguous bases in the same indices, the

128    number of mismatches is underestimated. Another limitation is that the algorithm does

129    not handle indels. Those limitations are compensated in the design of adapter matching,

130    scoring, and decision rules.

131

132    *Matching and scoring*

133    We implement four pairs of matching to utilize properties of paired-end reads

134    thoroughly: (1) matching adapter 1 head to read 1, (2) matching adapter 2 head to read

135    2, (3) matching read 1 head to reverse complement of read 2 and (4) matching read 2

136    head to reverse complement of read 1 (**Fig. 2C**). If the maximum number of bases

137    matched of (1) and (2) is less than a user-defined cut-off (default is 9), (3) and (4) will

138    be performed with a loosed $k$ ($= k_{original} + 1$). If the largest number of matched bases of

139    the four matches is greater than the cut-off, and some matches do not meet the

140    requirement, we will re-run those matches with a loosed $k$ ($= k_{original} + 3$) at the insert

141    size indicated from the best match. If the new number of matched bases is greater than

142    the cut-off, the old match will be discarded.

143    The scoring system measures the matching reliability of the whole 16-mer rather

144    than each base. The Phred quality score $Q$ of each base is converted to the probability

145    $P$ of that the corresponding base being correct using the formula:

146
$$P = 1 - 10^{\left(-\frac{Q}{10}\right)}$$

147    Then, the average base quality $\bar{P}$ of 16-mer sub-sequence $a$ at the memory position

148    $n$ is computed:

149
$$\overline{P_{a_n}} = \frac{1}{16} \sum_{i=2n}^{2n+15} P_{a[i]}$$

150    Notably, if the read quality is too low, it would imply an invalid match. However,

151    in reality, invalid matches are filtered out by the kmer-based algorithm. To solve the

152    discordance, we limit the lower bound of $\bar{P}$ to 0.75 manually.

153    The matching score $S$ between $a_n$ and $b_m$ is defined as:

154
$$S_{a_n,b_m} = count\_ones(a_n \ \& \ b_m) \cdot \overline{P_{a_n}} \cdot \overline{P_{b_m}}$$

155    where *count_ones* counts the number of ones in the binary representation of the UInt64.

156    When sequence $a$ is a user-defined adapter, $\bar{P}_a = 1$ is used. Generally, the matching

157    score $S$ is ranged from 0 to 16.

158

159    **Pseudocode 1: Matching and scoring**

160    `r1_pos_adpt, r1_nmatch_adpt = match(adapter1, r1, k)`

161    `r2_pos_adpt, r2_nmatch_adpt = match(adapter1, r1, k)`

162    `k_extra = max(r1_nmatch_adpt, r2_nmatch_adpt) < 9 ? 1 : 0`

163    `r1_pos_pe, r1_nmatch_pe = match(reverse_complement(r2), r1, k + k_extra)`

164    `r2_pos_pe, r2_nmatch_pe = match(reverse_complement(r1), r2, k + k_extra)`

165

166    `max_nmatch = max(r1_nmatch_adpt, r2_nmatch_adpt, r1_nmatch_pe, r2_nmatch_pe)`

167    `max_pos = corresponding position of max_nmatch`

168    `if max_nmatch > 9`

```
169        for matches with any nmatch < 9
170            redo match with loosed k = k + 3 at max_pos
171            replace old results if nmatch > 9
172
173    r1_prob_adpt = average_16mer_quality(r1, r1_pos_adpt)
174    r2_prob_adpt = average_16mer_quality(r2, r2_pos_adpt)
175    r1_prob_head = average_16mer_quality(r1, 1)
176    r2_prob_head = average_16mer_quality(r2, 1)
177    r1_prob_pe = average_16mer_quality(r1, r1_pos_pe)
178    r2_prob_pe = average_16mer_quality(r2, r2_pos_pe)
179    r*_prob_* = 0.75 if any r*_prob_* < 0.75
180
181    r1_score_adpt = r1_nmatch_adpt * r1_prob_adpt
182    r2_score_adpt = r2_nmatch_adpt * r2_prob_adpt
183    r1_score_pe = r1_nmatch_pe * r1_prob_pe * r2_prob_head
184    r2_score_pe = r2_nmatch_pe * r2_prob_pe * r1_prob_head
185
```

*Decision rules*

This module infers correct adapter positions from the four pairs of matching described

in the previous section. It is illustrated and self-explanatory in **Fig. 2D**. First, in each

read, the adapter and paired-end matches are compared. The one with the higher

matching score is chosen. If both matches support the same adapter position, the

matching score of the read is the sum of adapter and paired-end matching scores. Then,

the matches of the two paired-end reads are compared using the same strategy. If one

read finds an ideal adapter (matching score > 10 by default) while the other read is too

194    short to check or the average base accuracy of its 16-mer is less than 0.6 (Phred $Q < 5$),

195    both reads will be trimmed. If the matching score of a given read pair is less than 10

196    (by default), the read pair will not be trimmed.

197    Other read pairs will be taken a further examination to reduce false positives, which

198    are usually adapter matches at read tails. A read tail is defined as the last several bases

199    (default is 12 bp) of each read. Reads are not trimmed if both statements are true: (1)

200    In any paired read, the adapter is found at the tail, but the paired-end match is not; (2)

201    In both paired reads, adapter and pair-end matches suggest different trimming positions.

202    Before the final trimming, one additional step is required for the accurate

203    positioning of adapter sequences. The previous steps usually assume the read 1 and 2

204    have the same length of insert sizes, but indel in reads usually lead to over or under trim

205    one base. To prevent this circumstance, Atria re-positions the adapter by matching one

206    adjacent base with the first four bp of adapter sequences. The position of the highest

207    number of bases matched is chosen to trim. This step is ignored when the inferred insert

208    size is greater than the read length minus three because, in this situation, the adapter

209    sequence is too short to check.

210

211    **Pseudocode 2: Decision rules**

```
212   function correct_insert_size(pos1, score1, pos2, score2)

213       if pos1 == pos2

214           return pos1, score1 + score2

215       else

216           score = max(score1, score2)
```

```
217          pos = corresponding pos of max score

218          return pos, score

219

220   r1_pos, r1_score = correct_insert_size(r1_pos_adpt, r1_score_adpt, r1_pos_pe,
221   r1_score_pe)

222   r2_pos, r2_score = correct_insert_size(r2_pos_adpt, r2_score_adpt, r2_pos_pe,
223   r2_score_pe)

224   r12_pos, r12_score = correct_insert_size(r1_pos, r1_score, r2_pos, r2_score)

225

226   if r1_pos != r2_pos

227       if r1_score > 10

228         r2_prob = average_16mer_quality(r2, r1_pos)

229          @goto "trim" if r2_prob < 0.6

230       elseif r2_score > 10

231         r1_prob = average_16mer_quality(r1, r2_pos)

232          @goto "trim" if r1_prob < 0.6

233

234   function check_read_tail(read)

235       E_adpt = whether adapter found at read tail

236       E_pe = whether pair-end match found at read tail

237       E = E_adpt & E_pe  # both matches in read tail

238       R = rx_pos_adpt == rx_pos_pe  # adapter and pair-end match at same position

239       return E, R

240

241   E1, R1 = check_read_tail(r1)

242   E2, R2 = check_read_tail(r2)

243   E = E1 | E2  # at least one read matching in read tail

244   R = R1 | R2  # at least one read matching at the same position
```

```
245    is_false_positive = E & !R

246

247    if r12_score > trim_score & !is_false_positive

248        @label "trim"

249        r1_pos_adjusted = adjacent_one_bp_check(r1, adapter1, r12_pos)

250        r2_pos_adjusted = adjacent_one_bp_check(r2, adapter2, r12_pos)

251        trim(r1, r1_pos_adjusted)

252        trim(r2, r2_pos_adjusted)

253
```

254    *Consensus calling*

255    In this module, the overlapped base pairs of read 1 and 2 are corrected to the

256    corresponding bases with higher quality scores. It has three steps, prediction,

257    assessment, and correction.

258        In the prediction step, Atria makes a preliminary estimate of whether a read pair

259    contains an overlapped region. If adapters are trimmed and the remaining lengths of

260    read 1 and 2 are the same, the prediction passes. If no adapter can be trimmed, two

261    additional matching and scoring are required. The head of the reverse complement of

262    read 2 is matched to read 1, and the head of the reverse complement of read 1 is matched

263    to read 2. If the two matches reach a consensus, the prediction passes. Otherwise, the

264    prediction fails and consensus calling is skipped.

265        In the assessment step, Atria compares the whole overlapped region using a similar

266    matching algorithm, except that ambiguous bases (N, 1111) are converted to gaps (0000)

267    before matching. If the ratio of mismatch is greater than a user-defined value (28% by

268    default), the assessment fails, and consensus calling is skipped.

269        In the correction step, each base pair in the overlapped region is corrected to the

270    corresponding base with the highest quality score.

271

272    *Quality trimming*

273    Atria implements a traditional sliding window algorithm to remove the low-quality tail.

274    The sliding window scans from the front of the read and computes the average Phred

275    quality score of the sliding window. If the average quality is less than a given threshold,

276    the read tail is removed.

277

278    *IO optimization*

279    Atria spends more time on reading and writing than matching and trimming, so the key

280    to reducing runtime is to optimize IO usage. Considering that a large amount of RAM

281    is easily accessible nowadays, Atria trades increased RAM usage with decreased time.

282    A large block of memory is allocated for reading input files, which is then wrapped and

283    encoded to FASTQ objects parallelly using multi-threading. On the contrary, in the

284    writing process, Atria unboxes and decodes FASTQ objects to string vectors in parallel

285    and writes sequentially to files. In addition, pigz (parallel gzip) and pbzip2 (parallel

286    bzip2) are called for compression and decompression when needed [7, 8]. Atria also

287    support running with a single thread.

288

## Comparison to related work

**The performance of adapter trimming on a simulated dataset**

289

290

291   We simulated 8.9 G bases with 100 bp paired-end reads from the *Arabidopsis thaliana*

292   reference genome using the Skewer modified ART, a public NGS read simulator to

293   allow adapters in the reads [9, 10]. The simulation profile was trained from a 101 bp

294   paired-end public dataset SRR330569, and the 33 bp adapter pair used in read

295   simulation is AGATCGGAAGAGCACACGTCTGAACTCCAGTCA and

296   AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT [11].

297   Atria v3.0.0 was benchmarked with cutting-edge and popular trimmers, including

298   AdapterRemoval v2.3.1 [3], Skewer v0.2.2 [10], Fastp v0.21.0 [12], Ktrim v1.2.1 [13],

299   Atropos v1.1.29 [14], SeqPurge v2012_12 [15], Trim Galore v0.6.5 [4] and

300   Trimmomatic v0.39 [5]. Only adapter trimming was used, and other trimming and

301   filtration were disabled. Detailed command line arguments are listed in **Table S1**. Each

302   trimming software was running on an idle Ubuntu 19.10 server with a 32-thread Intel

303   i9-9960X Central Processing Unit (CPU) @ 3.10 GHz, 128 gigabyte (GB) DDR4-3200

304   RAM, and a 2 terabyte (TB) Samsung 970 EVO Solid State Drive (SSD) (sequential

305   reads and writes up to 3.5 and 2.5 TB/s).

306   The trimming performance was evaluated based on the following metrics: positive

307   predictive value (PPV), as the fraction of the number of correctly trimmed reads to all

308   trimmed reads; sensitivity, as the fraction of the number of correctly trimmed reads to

309      the reads with adapters; specificity, as the fraction of the number of untrimmed reads

310      without adapters to all reads without adapters; and Matthew's correlation coefficient

311      (MCC) measuring overall quality of pattern recognition, as

312
$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FN)(TN + FP)}}$$

313      where TP is the number of reads trimmed correctly, TN is the number of untrimmed

314      reads without adapters, FP is the number of over-trimmed reads, and FN is the number

315      of under-trimmed reads [3, 10].

316

317      **Table 1    Adapter trimming performance on the 8.9 G bases with 100 bp paired-**

318      **end simulated data**

| Trimmer | PPV (%) | Sensitivity (%) | Specificity (%) | MCC (%) |
|---|---|---|---|---|
| Atria | 99.35 | 99.81 | 99.82 | 99.51 |
| AdapterRemoval | 99.42 | 99.94 | 99.83 | 99.61 |
| Atropos | 99.57 | 97.34 | 99.88 | 98.00 |
| Fastp | 98.73 | 99.58 | 99.61 | 98.92 |
| Ktrim | 91.51 | 85.85 | 98.84 | 87.84 |
| SeqPurge | 57.92 | 99.80 | 76.84 | 66.62 |
| Skewer | 99.58 | 99.53 | 99.88 | 99.44 |
| Trim Galore | 40.05 | 82.98 | 62.39 | 38.96 |
| Trimmomatic | 99.29 | 57.86 | 99.88 | 71.05 |

319

320      The adapter trimming performance is shown in **Table 1**. AdapterRemoval, Atria

321      and Skewer were the top-class adapter trimmers in terms of MCC (99.61%, 99.51%,

322      99.44%, respectively) (**Table 1**). Fastp (98.92%) and Atropos (98.00%) were in the

323      second tier (**Table 1**). Ktrim obtained a good specificity (98.84%) but sacrificed its

324     sensitivity (85.85%), and Trimmomatic achieved an exceptional specificity (99.88%)

325     by trading off its sensitivity (57.86%) (**Table 1**).

326        To compare speed and efficiency, elapsed time (wall time) and average CPU

327     consumption of each trimmer were recorded in different threading (1-32 threads) for

328     uncompressed and gzip compressed data formats (**Fig. 3, Table S1**). Efficiency was

329     defined as the fraction of processing speed to the percent of CPU utilized, so it was a

330     better measurement, especially in CPU-intensive scenarios, such as running on a server

331     with a job scheduling system or trimming multiple samples at the same time. Ktrim and

332     Atria were two of the fastest trimmers in terms of speed and efficiency, from one to 16

333     threads (**Fig. 3, Table S1**). For uncompressed data, Trimmomatic was faster than Atria

334     using 8-32 threads, but its real CPU usage was much greater than Atria (**Fig. 3, Table**

335     **S1**). The speed and efficiency of AdapterRemoval and Skewer were generally 2-4 times

336     less than Atria, and Atropos was the slowest one (**Fig. 3, Table S1**). SeqPurge did not

337     support the output of uncompressed data, so it was only tested in the compressed

338     benchmark.

339        When trimming compressed data, the speed of AdapterRemoval, Skewer, Fastp,

340     Atropos and Trimmomatic kept constant when the number of threads increased from 4

341     to 32, because they failed to utilize more than four CPU in the IO process, while Atria

342     and Trim Galore did not have the limitation (**Fig. 3, Table S1**). Atria was faster than

343     Trim Galore, and the efficiency of Atria was constantly two to three times greater than

344     Trim Galore (**Fig. 3, Table S1**). SeqPurge showed strange speed curves; when

345    assigning a single thread to SeqPurge, the average CPU usage was 300%, and the speed

346    and average CPU usage dropped when assigning 8 to 32 threads (**Fig. 3, Table S1**). In

347    addition, Ktrim did not support output compressed files, so we ignored it. In general,

348    Atria was the fastest trimmer when trimming compressed files.

349

350    **The detailed statistics of adapter trimming accuracy on a simulated dataset**

351    The previous portion benchmarks on a whole dataset. This section evaluates trimming

352    accuracy regarding different read properties, including adapter presence or absence,

353    base error, and adapter length. To achieve the goal, Atria integrates a benchmarking

354    toolkit for read simulation and trimming analysis.

355        The read simulation method was inspired by how sequencers read DNA. First, an

356    original DNA fragment (insert) with a given original insert size is simulated base by

357    base. Adenine, thymine, cytosine, and guanine are randomly chosen repetitively. Then,

358    the insert and adapter sequences are copied base by base with an error profile, which

359    simulates the procedure of sequencing by synthesis. The error profile defines

360    substitution rate, insertion rate, and deletion rate.

361        Twenty-one million read pairs were simulated with a uniform read length (100 bp),

362    different error profiles, adapter length, and original insert sizes. The baseline error

363    profile comprises a 0.1% substitution rate, 0.001% insertion rate, and 0.001% deletion

364    rate, inspired by an Illumina error profile analysis [16]. 1x, 2x, 3x, 4x, and 5x baseline

365    error profile, 16, 20, 24, 28, and 33 adapter lengths, and 66 to 120 even insert sizes are

366    chosen. In this way, the reads with the least insert size have full lengths of adapters.

367    The reads with 66-98 original insert sizes contain adapters, and the reads with 100-120

368    original insert sizes are free from adapter contamination, except for few reads with a

369    100 bp insert size containing indels. Therefore, in each condition combination, 30

370    thousand read pairs were simulated to avoid random errors. The reads were trimmed

371    with the same method described in the last section.

372       The average trimming performance among different conditions is shown in **Fig. 4**

373    **A**. When adapters are present, Atria trims 99.9% adapters accurately, and SeqPurge,

374    Fastp, and Atropos follow closely with an accuracy of 99.7% (**Fig. 4 A1**). When

375    adapters are absent, AdapterRemoval, Skewer, Trimmomatic, Atropos, and Atria

376    successfully leave 100.0% reads intact, and Fastp falls behind with 99.8% accuracy

377    (**Fig. 4 A2**).

378       **Fig. 4 B** illustrates the trimming accuracy on different read error profiles. When

379    adapters are present, the accuracy of all trimmers drops as error rates increase (**Fig. 4**

380    **B1**). Atria keeps the highest accuracy from 100.0% to 99.9%, and is almost not affected

381    by different error rates (**Fig. 4 B1**). The accuracy of SeqPurge, Fastp, and Atropos

382    decrease from 99.9% to 99.6%, 99.5%, and 99.4%, respectively (**Fig. 4 B1**). With no

383    adapter present in reads, the accuracy is hardly influenced by error profiles (**Fig. 4 B2**),

384    so the performance is similar to **Fig. 4 A2**. In addition, adapter lengths ranging from 16

385    to 33 bp are not relevant to most trimmers' accuracy, including Atria (**Fig. 4 C**).

386

387    **The performance of adapter trimming on real sequencing dataset**

388    *RNA-Seq paired-end dataset (SRR330569)*

389    SRR330569 is a real RNA-Seq dataset sequenced from *Drosophila simulans* with 5.46

390    G bases and 2 x 101 bp read length. It contains 38 bp adapter sequences

391    AGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCG                         and

392    AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGTAGAT in read 1 and read 2,

393    respectively. Adapter trimming was performed by different trimmers without other

394    trimming or filtering methods. Then, a sliding-window based quality trimming was

395    performed to remove low-quality tails (sliding window size = 5 and average Q score ≥

396    15). The adapter-trimmed reads and adapter-and-quality-trimmed reads were mapped

397    to the *Drosophila simulans* genome version 2.02 from FlyBase using Hisat2 v2.2.1,

398    respectively [17, 18]. Mapping statistics were collected using SAMTools Stat v1.10

399    [19]. Skewer did consensus calling after adapter trimming, and no option was provided

400    to disable it. To achieve benchmark parity, Skewer was compared to Atria with

401    consensus calling enabled, and other trimmers were compared to Atria without

402    consensus calling. Time trimming was recorded in accordance with a common scenario:

403    inputs were gzip-compressed and trimmed with eight threads, and outputs were also

404    gzip-compressed to reduce massive disk use. All tested trimmers worked in the scenario

405    except that Ktrim could not output gzip files (**Table 2**).

406         Atria was the fastest program to process and output compressed data in terms of

407    wall time (**Table 2**). It also achieved the highest number of reads mapped and paired,

408    and the percent of properly paired reads with or without quality trimming. Generally,

409    higher base mapped is accompanied with higher error rate in the mapping process, so

410    it is important to interpret the two metrics together. Atria had the lowest mapping error

411    rate of 8.1833‰ and the forth highest base mapped (**Table 2**). The trimmers

412    (AdapterRemoval, Fastp, and Atropos) with the highest three error rates has the highest

413    base mapped (**Table 2**). Our program generally improved more than 5% compared to

414    other trimmers for the data without quality trimming (**Table 2**). The mapping statistics

415    of data without quality trimming were generally worse than with quality trimming

416    except for Atria. Specifically, the properly paired rates of other trimmers without

417    quality trimming were 0.5 to 4% less than with quality trimming (**Table 2**). Quality

418    trimming also increased the number of mapped and paired reads and reduced the

419    number of unmapped reads (**Table 2**).

420

421    **Table 2    Performance of trimmers on real data**

| Metric | Trimming and consensus | | Trimming only | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Atria | Skewer | Atria | AR | Atropos | Fastp | Ktrim* | SeqPurge | Trim Galore | Trimmomatic |
| **Low-quality dataset (SRR330569, RNA, Hisat2 mapping)** | | | | | | | | | | |
| Elapsed time (min:sec)* | 2:38 | 9:19 | **2:32** | 11:29 | 10:08 | 9:17 | 1:34 + GZ | 3:53 | 3:39 | 9:38 |
| **No quality trimming** | | | | | | | | | | |
| Reads mapped and paired | 26,126,804 | 24,694,330 | **25,781,268** | 24,559,060 | 24,505,656 | 24,545,646 | 24,196,658 | 24,240,072 | 24,046,542 | 22,797,620 |
| Reads unmapped | 27,276,761 | 28,254,804 | **27,379,299** | 28,338,455 | 28,410,022 | 28,350,294 | 28,747,248 | 28,591,442 | 28,647,873 | 29,649,287 |
| Properly paired reads (%) | 48.3 | 45.6 | 47.6 | 45.3 | 45.2 | 45.3 | 44.5 | 44.7 | 44.2 | 38.3 |
| Base mapped (cigar) | 2,387,212,225 | 2,354,164,041 | 2,322,436,545 | **2,346,791,204** | 2,341,355,822 | 2,344,847,438 | 2,316,915,673 | 2,304,776,514 | 2,317,321,510 | 2,237,846,534 |
| Error rate (‰) | 7.3952 | 9.5897 | **8.1833** | 9.8902 | 9.8536 | 9.8683 | 9.7920 | 9.7904 | 9.6994 | 9.3106 |
| **With quality trimming** | | | | | | | | | | |
| Reads mapped and paired | 25,942,092 | 25,787,464 | **25,728,206** | 25,721,788 | 25,714,956 | 25,725,480 | 25,473,670 | 25,364,392 | 25,654,498 | 24,744,754 |
| Reads unmapped | 27,245,720 | 27,364,827 | **27,361,655** | 27,369,773 | 27,373,854 | 27,360,527 | 27,556,820 | 27,736,292 | 27,400,932 | 28,064,739 |
| Properly paired reads (%) | 47.9 | 47.6 | **47.5** | 47.5 | 47.5 | 47.5 | 46.9 | 46.8 | 47.3 | 42.3 |
| Base mapped (cigar) | 2,317,238,536 | 2,316,981,456 | 2,302,740,463 | **2,304,584,269** | 2,304,325,743 | 2,304,437,244 | 2,292,034,762 | 2,263,465,110 | 2,297,815,439 | 2,246,076,618 |
| Error rate (‰) | 7.1114 | 7.7882 | 7.8902 | 7.9160 | 7.9141 | 7.9149 | 7.9059 | 7.8787 | 7.8921 | **7.5649** |
| **High-quality dataset (ERR4695159, cell-free DNA, Bowtie2 mapping)** | | | | | | | | | | |
| Elapsed time (min:sec)* | 3:08 | 11:34 | **3:03** | 13:48 | 13:41 | 11:29 | 1:41 + GZ | 4:05 | 4:34 | 11:44 |
| **No quality trimming** | | | | | | | | | | |
| Reads mapped and paired | 54,367,548 | 54,287,616 | 54,324,964 | 54,319,438 | 54,299,922 | 54,322,088 | 53,087,420 | **54,446,104** | 54,218,344 | 54,128,760 |
| Reads unmapped | 1,094,145 | 1,016,244 | 1,119,103 | 989,335 | 1,002,745 | **978,665** | 2,317,968 | 999,099 | 1,041,005 | 1,094,752 |
| Properly paired reads (%) | 96.8 | 96.7 | 96.7 | 96.7 | 96.7 | 96.7 | 94.1 | **97.0** | 96.4 | 88.6 |
| Base mapped (cigar) | 7,703,820,585 | 7,700,134,673 | 7,700,298,217 | 7,701,482,302 | 7,700,749,164 | 7,699,298,008 | 7,607,799,306 | 7,512,839,360 | 7,677,087,845 | **7,720,352,493** |
| Error rate (‰) | 3.3082 | 3.8388 | 3.8724 | 3.8771 | 3.8834 | 3.8564 | 4.3239 | **3.8007** | 3.9173 | 6.1984 |
| **With quality trimming** | | | | | | | | | | |
| Reads mapped and paired | 54,553,566 | 54,526,276 | 54,546,192 | 54,541,948 | 54,539,502 | 54,549,462 | 53,335,674 | **54,608,308** | 54,482,002 | 54,403,982 |
| Reads unmapped | 965,447 | 984,845 | 967,917 | 970,869 | 973,217 | **826,424** | 2,136,081 | 890,884 | 999,918 | 914,003 |
| Properly paired reads (%) | 97.0 | 97.0 | 97.0 | 97.0 | 97.0 | 97.0 | 94.4 | **97.1** | 96.8 | 89.0 |
| Base mapped (cigar) | 7,653,879,312 | 7,649,380,218 | 7,646,989,362 | 7,647,893,624 | 7,648,184,196 | 7,646,574,606 | 7,556,468,882 | 7,461,588,482 | 7,625,484,706 | **7,668,777,971** |
| Error rate (‰) | 2.9547 | 3.2535 | 3.2678 | 3.2698 | 3.2792 | 3.2634 | 3.7183 | **3.2117** | 3.3109 | 5.5798 |

*Note*: AR = AdapterRemoval. In the trimming-only benchmark, bold and underline formats indicate the first and second trimmers (including tie) in terms of each metric, respectively. *Elapsed time (wall time) is benchmarked based on trimming and output gzip files with 8 threads, except that Ktrim cannot output gzip files (marked with time + GZ).

*Genome-wide human cell-free DNA dataset (ERR4695159)*

Generally, plasma cell-free DNA is short in length [20], and trimming is extremely important in medical diagnosis. Here, we chose a human genome-wide cell-free DNA dataset ERR4695159. It has 8.4 G bases with 2 x 150 bp read length with 33 bp adapter sequences AGATCGGAAGAGCACACGTCTGAACTCCAGTCA in read 1 and AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT in read 2. The benchmark workflow was the same as the RNA-Seq analysis, except that the clean reads were

436     mapped to the human reference genome hg38 (GRCh38.p13) using Bowtie2 v2.3.5.1

437     [21].

438     Atria was also the fastest trimmer in the scenario (3 min 3 s) (**Table 2**). SeqPurge

439     and Trim Galore finished the task in more than 4 minutes, while others spent more than

440     11 minutes (**Table 2**).

441     In adapter-trimming-only statistics, SeqPurge had the highest mapped and paired

442     reads (54,446,104) and the highest properly paired reads (97.0%) (**Table 2**). Atria

443     followed with 54,324,964 mapped and paired reads. Atria, AdapterRemoval, Atropos,

444     and Fastp all had 96.7% properly paired reads (**Table 2**).

445     With quality trimming, the overall performance increased, and properly paired

446     reads were closer; SeqPurge had 97.1% properly paired reads, with Atria,

447     AdapterRemoval, Atropos, and Fastp close behind at 97.0% (**Table 2**). Only 89.0% of

448     reads were properly paired with Trimmomatic (**Table 2**).

449

450     **Discussion**

451     Atria performs favorably with other cutting-edge adapter trimmers in accuracy,

452     robustness, speed, and efficiency. Its performance is ascribed to the byte-based

453     matching algorithm. The design concept of the algorithm is to minimize any

454     unnecessary CPU operations by taking advantage of the data structure of dense arrays.

455     Matrix-based algorithms, such as the Needleman-Wunsch algorithm and the Smith-

456     Waterman algorithm, allocate and update a matrix and perform base-to-base

457    comparison [22, 23]. They report every mismatch and gap between two sequences

458    while Atria skips this step since it is focussed on the start positions of successful

459    matches. Despite that, the matching algorithm used in Atria is able to identify mismatch

460    loci when needed.

461    The byte-based matching algorithm is lightweight and designed for short sequence

462    scanning. Each DNA is encoded in four bits and stores continuously in RAM. A sub-

463    sequence can be extracted as an unsigned integer from a given memory position. For

464    example, a 64-bit unsigned integer (UInt) represents a 16-mer, and a 128-bit UInt

465    represents a 32-mer. The comparison between two sub-sequences is completed within

466    the accumulator register, a CPU unit for arithmetic or logical operation. It does not

467    require addressing or updating a scoring matrix from RAM. When comparing a short

468    sequence, such as an adapter, to a long sequence, such as the read, the 16-mer of the

469    short sequence is compared to every position of the long sequence. Hence, the byte-

470    based matching algorithm has $O(n)$ expected time complexity and $O(1)$ space

471    complexity in adapter matching, where n is the length of the long sequence.

472    The algorithm also has its limitations. It only reports the number of matched bases

473    and does not report the positions of mismatches, so it cannot be used for sequence

474    alignment solely. Besides, the algorithm does not handle insert and deletion. However,

475    the average indel rate of Illumina library is $10^{-6}$ to $10^{-5}$ [16], and the low indel rate is

476    almost negligible in real data analysis. In addition, Atria does four pairs of matches in

477    different locations to compensate for the limitation. If one match is failed because of

478    indel, other matches will suggest the real adapter positions.

479         In the runtime benchmark, we compared how trimmers performed using extremely

480    high CPU cores. In general, efficiency marginally decreased as CPU usage increased

481    due to the trimmers' parallel implementation and the inevitable cost of multi-threading,

482    such as task scheduling and context switching. In addition, IO could be the main

483    bottleneck for most hard disk drives and some solid-state drives. Thus, if the system IO

484    reaches a bottleneck, an efficiency plateau would be expected sooner.

485

486    **Conclusions**

487    We introduce not only Atria, a cutting-edge trimming software for sequence data, but

488    also the ultra-fast and lightweight byte-based matching algorithm. The algorithm can

489    be used in a broad range of short-sequence matching applications, such as primer search

490    and seed scanning before alignment. Atria is implemented in Julia, a programming

491    language designed specifically for high performance. The source code, executables, and

492    benchmark scripts are available on Atria's Github page [24].

493

494    **Availability and requirements**

495    Project name: Atria

496    Project home page: https://github.com/cihga39871/Atria

497    Operating system(s): Linux, OSX

498    Programming language: Julia

499    Other requirements: Julia v1.4, Pigz v2.4 or higher, Pbzip v1.1.13 or higher

500    License: MIT

501    Research Resource Identification Initiative ID: SCR_021313

502

## Data Availability

504    The datasets SRR330569, and ERR4695159 analyzed during the current study are

505    available in the Sequence Read Archive from the National Center for Biotechnology

506    Information [11, 25].

507    The Atria source codes, releases, documents, and benchmark scripts can be

508    downloaded from Atria's Github page [24].

509

## Abbreviations

511    CPU: Central processing unit; DNA: Deoxyribonucleic acid; GB: Gigabyte; MCC:

512    Matthew's correlation coefficient; NGS: Next-generation sequencing; PPV: Positive

513    predictive value; RAM: Random-access memory; RNA: Ribonucleic acid; SNP: Single

514    nucleotide polymorphism; SSD: Solid-state drive; TB: Terabyte; UInt: Unsigned

515    integer; UInt64: Unsigned 64-bit integer; WGS: Whole-genome sequencing.

516

## Competing interests

518    The authors declare that they have no competing interests.

519

## Funding

523

## Authors' contributions

525 JC developed Atria software, performed benchmark experiments under the supervision

526 of XL. Both XL and LH serve as co-supervisors and participates in the design of the

527 study. MH contributes to the optimization of the algorithm. AZ participates in

528 benchmark validation. JC, LH, and XL drafted the manuscript. All authors read and

529 approved the final version of the manuscript.

530

## Acknowledgments

536

## References

538 1. Schluth-Bolard C, Diguet F, Chatron N, Rollat-Farnier PA, Bardel C, Afenjar
539 A, et al. Whole genome paired-end sequencing elucidates functional and
540 phenotypic consequences of balanced chromosomal rearrangement in patients

541  with developmental disorders. J Med Genet. 2019;56 8:526-35.
542  doi:10.1136/jmedgenet-2018-105778.

543  2. Tan G, Opitz L, Schlapbach R and Rehrauer H. Long fragments achieve lower
544  base quality in Illumina paired-end sequencing. Sci Rep. 2019;9 1:2856.
545  doi:10.1038/s41598-019-39076-7.

546  3. Schubert M, Lindgreen S and Orlando L. AdapterRemoval v2: rapid adapter
547  trimming, identification, and read merging. BMC Res Notes. 2016;9 1:88.
548  doi:10.1186/s13104-016-1900-2.

549  4. Krueger F. Trim galore. A wrapper tool around Cutadapt and FastQC to
550  consistently apply quality and adapter trimming to FastQ files. 2015;516:517.

551  5. Bolger AM, Lohse M and Usadel B. Trimmomatic: a flexible trimmer for
552  Illumina sequence data. Bioinformatics. 2014;30 15:2114-20.
553  doi:10.1093/bioinformatics/btu170.

554  6. BioJulia/BioSequences.jl: Biological sequences for the Julia language.
555  https://github.com/BioJulia/BioSequences.jl. Accessed 1 Dec 2020.

556  7. Pigz - Parallel gzip. https://zlib.net/pigz/. Accessed 1 Dec 2020.

557  8. Parallel BZIP2 (PBZIP2). http://compression.ca/pbzip2/. Accessed 1 Feb 2021.

558  9. Huang W, Li L, Myers JR and Marth GT. ART: a next-generation sequencing
559  read simulator. Bioinformatics. 2012;28 4:593-4.
560  doi:10.1093/bioinformatics/btr708.

561  10. Jiang H, Lei R, Ding SW and Zhu S. Skewer: a fast and accurate adapter trimmer
562  for next-generation sequencing paired-end reads. BMC Bioinformatics.
563  2014;15 1:182. doi:10.1186/1471-2105-15-182.

564  11. Barrett T, Clark K, Gevorgyan R, Gorelenkov V, Gribov E, Karsch-Mizrachi I,
565  et al. BioProject and BioSample databases at NCBI: facilitating capture and
566  organization of metadata. Nucleic Acids Res. 2012;40 D1:D57-D63.

567  12. Chen S, Zhou Y, Chen Y and Gu J. fastp: an ultra-fast all-in-one FASTQ
568  preprocessor. Bioinformatics. 2018;34 17:i884-i90.
569  doi:10.1093/bioinformatics/bty560.

570  13. Sun K. Ktrim: an extra-fast and accurate adapter- and quality-trimmer for
571  sequencing data. Bioinformatics. 2020;36 11:3561-2.
572  doi:10.1093/bioinformatics/btaa171.

573  14. Didion JP, Martin M and Collins FS. Atropos: specific, sensitive, and speedy
574  trimming of sequencing reads. PeerJ. 2017;5:e3720. doi:10.7717/peerj.3720.

575  15. Sturm M, Schroeder C and Bauer P. SeqPurge: highly-sensitive adapter
576  trimming for paired-end NGS data. BMC Bioinformatics. 2016;17 1:208.
577  doi:10.1186/s12859-016-1069-7.

578  16. Schirmer M, D'Amore R, Ijaz UZ, Hall N and Quince C. Illumina error profiles:
579  resolving fine-scale variation in metagenomic sequencing data. BMC
580  Bioinformatics. 2016;17 1:125. doi:10.1186/s12859-016-0976-y.

581   17.   Thurmond J, Goodman JL, Strelets VB, Attrill H, Gramates LS, Marygold SJ,
582         et al. FlyBase 2.0: the next generation. Nucleic Acids Res. 2019;47 D1:D759-
583         D65. doi:10.1093/nar/gky1003.
584   18.   Kim D, Paggi JM, Park C, Bennett C and Salzberg SL. Graph-based genome
585         alignment and genotyping with HISAT2 and HISAT-genotype. Nat Biotechnol.
586         2019;37 8:907-15. doi:10.1038/s41587-019-0201-4.
587   19.   Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The
588         Sequence Alignment/Map format and SAMtools. Bioinformatics. 2009;25
589         16:2078-9. doi:10.1093/bioinformatics/btp352.
590   20.   Sun K, Jiang P, Wong AIC, Cheng YKY, Cheng SH, Zhang H, et al. Size-tagged
591         preferred ends in maternal plasma DNA shed light on the production
592         mechanism and show utility in noninvasive prenatal testing. Proceedings of the
593         National Academy of Sciences. 2018;115 22:E5106-E14.
594   21.   Langmead B, Wilks C, Antonescu V and Charles R. Scaling read aligners to
595         hundreds of threads on general-purpose processors. Bioinformatics. 2019;35
596         3:421-32. doi:10.1093/bioinformatics/bty648.
597   22.   Needleman SB and Wunsch CD. A general method applicable to the search for
598         similarities in the amino acid sequence of two proteins. J Mol Biol. 1970;48
599         3:443-53. doi:10.1016/0022-2836(70)90057-4.
600   23.   Smith TF and Waterman MS. Identification of common molecular
601         subsequences. J Mol Biol. 1981;147 1:195-7. doi:10.1016/0022-
602         2836(81)90087-5.
603   24.   cihga39871/Atria: An ultra-fast and accurate NGS adapter and quality trimmer.
604         https://github.com/cihga39871/Atria. Accessed 31 Mar 2021.
605   25.   Sequence Read Archive from the National Center for Biotechnology
606         Information. https://www.ncbi.nlm.nih.gov/sra/. Accessed 15 Jan 2021.

607

608   **Figure legends**

609   **Figure 1   Overview of Atria workflow**

610   **Figure 2   Adapter trimming algorithms**

611   **Figure 3   Benchmark of adapter-trimming speed for uncompressed and**

612   **compressed files on different threading options**

613   The 8.9 G bases simulated paired-end data with a 100 bp read length was trimmed in

614   both uncompressed and compressed format using up to 32 threads. Speed is the ratio of

615 the number of bases to elapsed time (wall time). SeqPurge does not support

616 uncompressed outputs, so it is not shown in the uncompressed benchmark. In the

617 trimming for compressed data, the speed of AdapterRemoval, Skewer, Fastp, Atropos,

618 and Trimmomatic kept constant when the number of threads increased from 4 to 32, so

619 we only benchmark those trimmers using 1, 2, and 4 threads. Ktrim does not support

620 output compressed files, so it is not shown in the compressed benchmark.

621 **Figure 4   Adapter trimming accuracy on adapter presence and absence, different**

622 **base errors, and adapter lengths**

623 A1, B1, and C1 are statistics for reads with adapter contamination, while A2, B2, C2

624 for reads without adapters. A1 and A2 show the accumulated rates of accurate trim, one

625 bp over trim, one bp under trim, multiple bp over trim, and multiple bp under trim. In

626 A1, the accuracy of Trimmomatic is 41.0%. In A2, the accuracy of SeqPurge is 78.8%,

627 the accuracy of Trim Galore is 68.3%. B1 and B2 show the trimming accuracy on

628 different error profiles. In B1, the accuracy of Trimmomatic drops from 41.9% to

629 40.1%. In B2, the accuracy of SeqPurge is 78.8%, and the accuracy of Trim Galore is

630 68.2 - 68.3%. C1 and C2 show the trimming accuracy on different adapter lengths. In

631 C1, the accuracy of Trimmomatic is 0.0% at 16 bp adapter length, 50.7% to 51.6% at

632 adapter lengths from 20 to 33 bp. In C2, the accuracy of SeqPurge ranges from 78.7%

633 at 16 bp to 78.9% at 33 bp, and the accuracy of Trim Galore ranges in 68.2 - 68.3%

634 from 16 to 33 bp.

635

## Supplementary material

**Table S1   Trimming speed on the 8.9 G bases 100 bp paired-end simulated data**

Atria (consensus) does both adapter trimming and paired-end consensus call (base correction of overlapped regions). In the trimming for uncompressed data, SeqPurge does not support uncompressed outputs, so it is not shown in the uncompressed benchmark. Fastp does not support 32 threads, so only 1-16 threads were tested. In the trimming for compressed data, the speed of AdapterRemoval, Skewer, Fastp, and Trimmomatic kept constant when the number of threads increased from 4 to 32, so we only benchmarked those trimmers using 1, 2, and 4 threads. Atropos was too slow to trim compressed data, and Ktrim did not support compressed outputs, so they are not shown in the compressed benchmark.

**Table 2   Performance of trimmers on real data (larger than A4)**

| | Trimming and consensus | | Trimming only | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Metric | Atria | Skewer | Atria | AR | Atropos | Fastp | Ktrim* | SeqPurge | Trim Galore | Trimmomatic |
| **Low-quality dataset (SRR330569, RNA, Hisat2 mapping)** | | | | | | | | | | |
| Elapsed time (min:sec)* | 2:38 | 9:19 | **2:32** | 11:29 | 10:08 | 9:17 | 1:34 + GZ | 3:53 | <u>3:39</u> | 9:38 |

**No quality trimming**

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Reads mapped and paired | 26,126,804 | 24,694,330 | **25,781,268** | 24,559,060 | 24,505,656 | 24,545,646 | 24,196,658 | 24,240,072 | 24,046,542 | 22,797,620 |
| Reads unmapped | 27,276,761 | 28,254,804 | **27,379,299** | 28,338,455 | 28,410,022 | 28,350,294 | 28,747,248 | 28,591,442 | 28,647,873 | 29,649,287 |
| Properly paired reads (%) | 48.3 | 45.6 | **47.6** | 45.3 | 45.2 | 45.3 | 44.5 | 44.7 | 44.2 | 38.3 |
| Base mapped (cigar) | 2,387,212,225 | 2,354,164,041 | 2,322,436,545 | **2,346,791,204** | 2,341,355,822 | 2,344,847,438 | 2,316,915,673 | 2,304,776,514 | 2,317,321,510 | 2,237,846,534 |
| Error rate (‰) | 7.3952 | 9.5897 | **8.1833** | 9.8902 | 9.8536 | 9.8683 | 9.7920 | 9.7904 | 9.6994 | 9.3106 |

**With quality trimming**

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Reads mapped and paired | 25,942,092 | 25,787,464 | **25,728,206** | 25,721,788 | 25,714,956 | 25,725,480 | 25,473,670 | 25,364,392 | 25,654,498 | 24,744,754 |
| Reads unmapped | 27,245,720 | 27,364,827 | **27,361,655** | 27,369,773 | 27,373,854 | 27,360,527 | 27,556,820 | 27,736,292 | 27,400,932 | 28,064,739 |
| Properly paired reads (%) | 47.9 | 47.6 | **47.5** | **47.5** | **47.5** | **47.5** | 46.9 | 46.8 | 47.3 | 42.3 |
| Base mapped (cigar) | 2,317,238,536 | 2,316,981,456 | 2,302,740,463 | **2,304,584,269** | 2,304,325,743 | 2,304,437,244 | 2,292,034,762 | 2,263,465,110 | 2,297,815,439 | 2,246,076,618 |
| Error rate (‰) | 7.1114 | 7.7882 | 7.8902 | 7.9160 | 7.9141 | 7.9149 | 7.9059 | 7.8787 | 7.8921 | **7.5649** |

**High-quality dataset (ERR4695159, cell-free DNA, Bowtie2 mapping)**

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Elapsed time (min:sec)* | 3:08 | 11:34 | **3:03** | 13:48 | 13:41 | 11:29 | 1:41 + GZ | 4:05 | 4:34 | 11:44 |

**No quality trimming**

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Reads mapped and paired | 54,367,548 | 54,287,616 | 54,324,964 | 54,319,438 | 54,299,922 | 54,322,088 | 53,087,420 | **54,446,104** | 54,218,344 | 54,128,760 |
| Reads unmapped | 1,094,145 | 1,016,244 | 1,119,103 | 989,335 | 1,002,745 | 978,665 | 2,317,968 | 999,099 | 1,041,005 | 1,094,752 |
| Properly paired reads (%) | 96.8 | 96.7 | 96.7 | 96.7 | 96.7 | 96.7 | 94.1 | **97.0** | 96.4 | 88.6 |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Base mapped (cigar) | 7,703,820,585 | 7,700,134,673 | 7,700,298,217 | <u>7,701,482,302</u> | 7,700,749,164 | 7,699,298,008 | 7,607,799,306 | 7,512,839,360 | 7,677,087,845 | **7,720,352,493** |
| Error rate (‰) | 3.3082 | 3.8388 | 3.8724 | 3.8771 | 3.8834 | <u>3.8564</u> | 4.3239 | **3.8007** | 3.9173 | 6.1984 |
| **With quality trimming** | | | | | | | | | | |
| Reads mapped and paired | 54,553,566 | 54,526,276 | 54,546,192 | 54,541,948 | 54,539,502 | <u>54,549,462</u> | 53,335,674 | **54,608,308** | 54,482,002 | 54,403,982 |
| Reads unmapped | 965,447 | 984,845 | 967,917 | 970,869 | 973,217 | **826,424** | 2,136,081 | <u>890,884</u> | 999,918 | 914,003 |
| Properly paired reads (%) | 97.0 | 97.0 | <u>97.0</u> | <u>97.0</u> | <u>97.0</u> | <u>97.0</u> | 94.4 | **97.1** | 96.8 | 89.0 |
| Base mapped (cigar) | 7,653,879,312 | 7,649,380,218 | 7,646,989,362 | 7,647,893,624 | <u>7,648,184,196</u> | 7,646,574,606 | 7,556,468,882 | 7,461,588,482 | 7,625,484,706 | **7,668,777,971** |
| Error rate (‰) | 2.9547 | 3.2535 | 3.2678 | 3.2698 | 3.2792 | <u>3.2634</u> | 3.7183 | **3.2117** | 3.3109 | 5.5798 |

653

Adapter Trimming

Consensus Calling of Overlapped Region

Quality Trimming

ADAPTER1

READ1 INSERT

READ2 INSERT

ADAPTER2

READ1 INSERT

READ2 INSERT

READ1 INSERT

READ2 INSERT

**A** DNA Encoding

Encoding a DNA sequence to a dense array

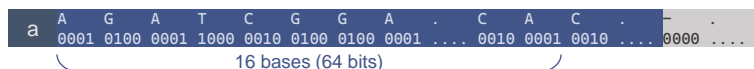| A | C | G | T | N | A | C | G | T | N | T | – | – | – | – | – | – | . | – |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
0001 0100 1000 1111 0001 0010 0100 1000 1111 1000 | 0000 0000 0000 0000 0000 0000 .... 0000

DNA encoding region — At least 64 extra bits set to 0 to prevent from overflow

Extracting a 16-mer subsequence as a 64-bit unsigned integer (UInt64)

Extract from **memory position 0** (16 bases extracted, sequence indices [0:15], valid indices [0:10])
0001 0010 0100 1000 1111 0001 0010 0100 1000 1111 1000 0000 0000 0000 0000 0000

8-bit

Extract from **memory position 1** (16 bases extracted, sequence indices [2:17], valid indices [2:10])
0100 1000 1111 0001 0010 0100 1000 1111 1000 0000 0000 0000 0000 0000 0000 0000
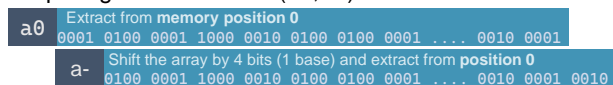
8 bits (1 byte, 2 bases) is the smallest addressable unit of memory in many computer architectures

**B** Matching Algorithm

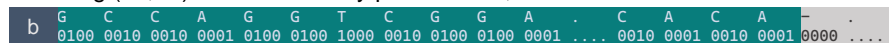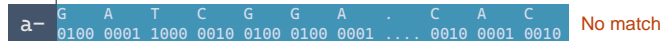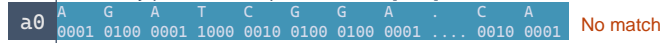Given two sequences a and b, matching the head of a to each memory position of b

| a | A | G | A | T | C | G | G | A | . | C | A | C | . | – | . |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
0001 0100 0001 1000 0010 0100 0100 0001 .... 0010 0001 0010 .... 0000 ....

16 bases (64 bits)

Preparing the heads of a (a0, a-)

a0 — Extract from **memory position 0**
0001 0100 0001 1000 0010 0100 0100 0001 .... 0010 0001

a- — Shift the array by 4 bits (1 base) and extract from **position 0**
0100 0001 1000 0010 0100 0100 0001 .... 0010 0001 0010

Matching (a0, a-) to each memory position of b, and the best index of b is returned

| b | G | C | C | A | G | G | T | C | G | G | A | . | C | A | C | A | – | . |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
0100 0010 0010 0001 0100 0100 1000 0010 0100 0100 0001 .... 0010 0001 0010 0001 0000 ....

b: memory position 0, sequence indices [0:15]

a0 | A G A T C G G A . C A
0001 0100 0001 1000 0010 0100 0100 0001 .... 0010 0001 — No match

a– | G A T C G G A . C A C
0100 0001 1000 0010 0100 0100 0001 .... 0010 0001 0010 — No match

b: memory position 1, sequence indices [2:17]

a0 | A G A T C G G A . C A
0001 0100 0001 1000 0010 0100 0100 0001 .... 0010 0001 — No match

a– | G A T C G G A . C A C
0100 0001 1000 0010 0100 0100 0001 .... 0010 0001 0010 — No match

b: memory position 2, sequence indices [4:19]

a0 | A G A T C G G A . C A
0001 0100 0001 1000 0010 0100 0100 0001 .... 0010 0001 — No match

a– | G A T C G G A . C A C
0100 0001 1000 0010 0100 0100 0001 .... 0010 0001 0010 — Match (1 error)
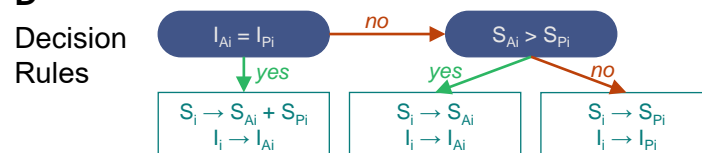
b: memory position 3, sequence indices [6:21] … b: memory position n, sequence indices [2n:2n+15]

**C** Matching & Scoring
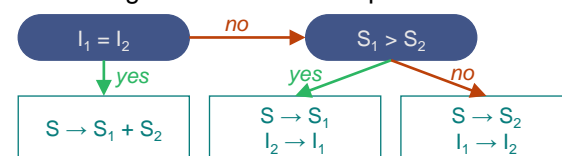
ADAPTER1

READ1 INSERT

READ2 INSERT

ADAPTER2

**C** Matching & Scoring (cont.)

(1) Matching adapter 1 to read 1
Insert size $I_{A1}$, matching score $S_{A1}$

(2) Matching read 1 to reverse complement of read 2
Insert size $I_{P1}$, matching score $S_{P1}$

(3) Matching adapter 2 to read 2
Insert size $I_{A2}$, matching score $S_{A2}$

(4) Matching read 2 to reverse complement of read 1
Insert size $I_{P1}$, matching score $S_{P2}$

**D** Decision Rules

Correcting insert size in each read. For i in 1, 2:

$I_{Ai} = I_{Pi}$ —no→ $S_{Ai} > S_{Pi}$

yes ↓

$S_i \to S_{Ai} + S_{Pi}$
$I_i \to I_{Ai}$

yes →
$S_i \to S_{Ai}$
$I_i \to I_{Ai}$

no →
$S_i \to S_{Pi}$
$I_i \to I_{Pi}$

Correcting insert size from its paired read.

$I_1 = I_2$ —no→ $S_1 > S_2$

yes ↓

$S \to S_1 + S_2$

yes →
$S \to S_1$
$I_2 \to I_1$

no →
$S \to S_2$
$I_1 \to I_2$

Checking false positive and trim.

$L_{1/2}$ as the valid sequence length of read 1 or 2.
$L_{tail}$ as the user-defined tail length.
**E** as whether a match is found at the end of read (Bool).
**R** as whether the adapter and insert matches suggest the same insert size (Bool).
$S_{trim}$ as the minimum score to trim adapters.

$E_{A1} \to (I_{A1} > L_1 - L_{tail})$
$E_{P1} \to (I_{P1} < L_1 - L_{tail})$
$E_1 \to (E_{A1} \& E_{P1})$
$R_1 \to (I_{A1} == I_{P1})$

$E_{A2} \to (I_{A2} > L_2 - L_{tail})$
$E_{P2} \to (I_{P2} < L_2 - L_{tail})$
$E_2 \to (E_{A2} \& E_{P2})$
$R_2 \to (I_{A2} == I_{P2})$

$S > S_{trim}$ —no→ NO TRIM

yes ↓

(E1 | E2) & !(R1 | R2) —yes→ NO TRIM

no ↓

TRIM