1 # Title Page:

2 MQF and buffered MQF: Quotient filters for efficient storage of k-mers with their counts and

3 metadata

4 ● Moustafa Shokrof

5 ○ Department of Computer Science, University of California, Davis, CA, USA

6 ● C. Titus Brown:

7 ○ Department of Population Health and Reproduction, School of Veterinary

8 Medicine, University of California, Davis, CA, USA

9 ● Tamer A. Mansour (corresponding author)

10 ○ Department of Clinical Pathology, School of Medicine, University of Mansoura,

11 Mansoura, Egypt

12 ○ Department of Population Health and Reproduction, School of Veterinary

13 Medicine, University of California, Davis, CA, USA

14

15

16


17
18
19
20
21
22
23
24

# Abstract

## **Background**

Specialized data structures are required for online algorithms to efficiently handle large sequencing datasets. The counting quotient filter (CQF), a compact hashtable, can efficiently store k-mers with a skewed distribution.

## **Result**

Here, we present the mixed-counters quotient filter (MQF) as a new variant of the CQF with novel counting and labeling systems. The new counting system adapts to a wider range of data distributions for increased space efficiency and is faster than the CQF for insertions and queries in most of the tested scenarios. A buffered version of the MQF can offload storage to disk, trading speed of insertions and queries for a significant memory reduction. The labeling system provides a flexible framework for assigning labels to member items while maintaining good data locality and a concise memory representation. These labels serve as a minimal perfect hash function but are ~10 fold faster than BBhash, with no need to re-analyze the original data for further insertions or deletions.

## **Conclusion**

The MQF is a flexible and efficient data structure that extends our ability to work with high throughput sequencing data.

# Keywords

Compact hash tables, k-mers, debruijn graphs, NGS, inexact data structures.

# Background

48    Online algorithms effectively support streaming analysis of large data sets, which is

49    important for analyzing data sets with large volume and high velocity(1). Approximate data

50    structures are commonly used in online algorithms to provide better average space and time

51    efficiency (2). For example, the Bloom filter supports approximate set membership queries with

52    a predefined false positive rate (FPR) (3). The count-min sketch (CMS) is similar to Bloom filters

53    and can be used to count items with a tunable rate of overestimation. However, there are a

54    number of problems with Bloom filters and the CMS - in particular, they do not support data

55    locality.

56    The Counting Quotient Filter (CQF) is a more efficient data structure that serves similar

57    purposes with better efficiency for skewed distributions and much better data locality(4). The

58    CQF is a recent variant of quotient filters that tracks the count of its items using a variable size

59    counter. As a compact hashtable, CQF can perform in either probabilistic or exact modes and

60    supports deletes, merges, and resizing.

61    Analysis of k-mers in biological sequencing data sets is an ongoing challenge(5). K-mers

62    in raw sequencing data often have a high Zipfian distribution, and the CQF was built to minimize

63    memory requirements for counting such items. However, this advantage deteriorates in

64    applications that require frequent random access to the data structure, and where the k-mer

65    count distribution may change in response to different sampling approaches, library preparation

66    and/or sequencing technologies. For example, k-mer frequency across 1000s of RNAseq

67    experiments shows different patterns of abundant k-mers (6).

68    Data structures like CMS (7) and CQF (4) also do not natively support associating k-

69    mers with multiple values, which can be useful for coloring in De Bruijn graphs as well as other

70    features (8). Classical hash tables are designed to associate their keys with a generic data type

71    but they are expensive memory-wise (9). Minimal Perfect Hash Functions (MPHFs) can provide

72   a more compact solution by mapping each k-mer into a unique integer. These integers can then

73   be used as indices for the k-mers to label them in other data structures (10). An implementation

74   capable of handling large scale datasets with fast performance requires ~3 bits per element

75   (11). However, such a concise representation comes with a high false-positive rate on queries

76   for non-existent items. Moreover, unlike hashtables, MPHF does not support insertions or

77   deletions thus any change in the k-mer set would require rehashing of the original dataset.

78          In this paper, we introduce the mixed-counters quotient filter (MQF), a modified version

79   of the CQF with a new encoding scheme and labeling system supporting high data locality. We

80   further show how Buffered MQF can be used to scale MQF to solid-state disks. We compare

81   between MQF and the CQF, CMS, and MPHF data structures regarding memory efficiency,

82   speed performance, and applicability to specific data analysis challenges. We further do a direct

83   comparison of the CMS to MQF in the khmer software package for sequencing data analysis, to

84   showcase the benefits of MQF is in real world applications.
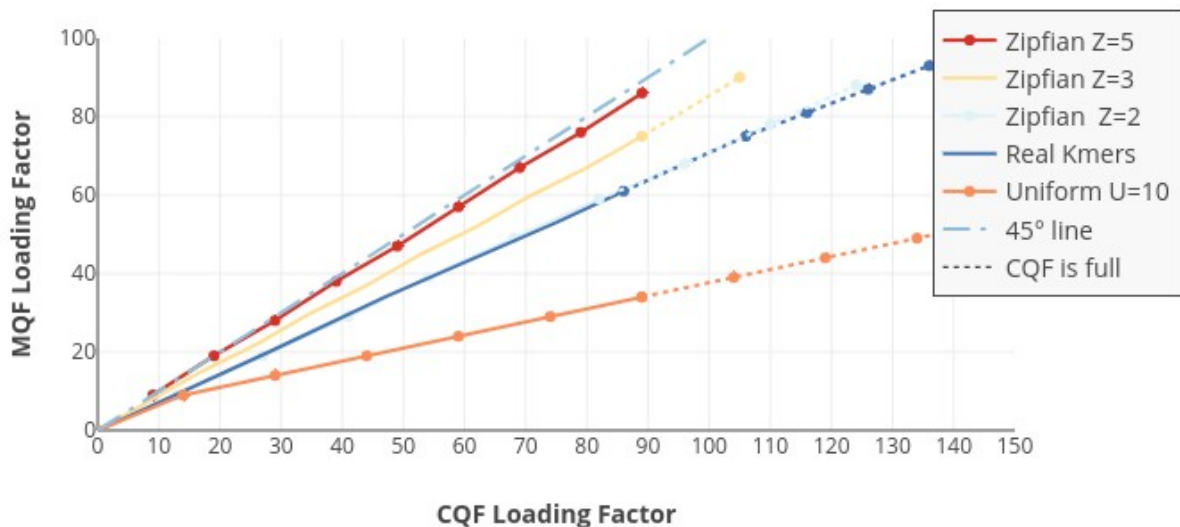
85

86

87

88

89

90

91

## 92   Results

93   **MQF has a lower load factor than CQF**

94    The load factor is defined as the actual space utilized divided by the total space

95    assigned for the data structure, and is an important measure of data structure performance. To

96    compare load factors between the CQF and MQF data structures, instances of both structures

97    were created using the same number of slots ($2^{27}$). Chunks of items from five datasets with

98    different distributions of item frequencies were inserted iteratively to both data-structures while

99    recording the load factor after the insertion of each chunk. The experiments stopped when

100   MQF's load factor reached 90%. MQF had lower loading factors for all tested datasets but the

101   difference was minimal for the dataset with the highest Zipfian distribution (Z=5). The lower the

102   tested Zipfian distribution the lower the loading factor of MQF (Figure 1). A lower loading factor

103   enabled MQF to accommodate > 30% of the CQF capacity from a dataset of real k-mers and to

104   exceed the double CQF capacity with uniform distribution (Figure 1 and supplementary table 1).

105



106

107   **Figure 1: MQF has a lower load factor compared to CQF**. Chunks of items, from different

108   distributions of item's frequencies, were inserted iteratively to matching CQF and MQF
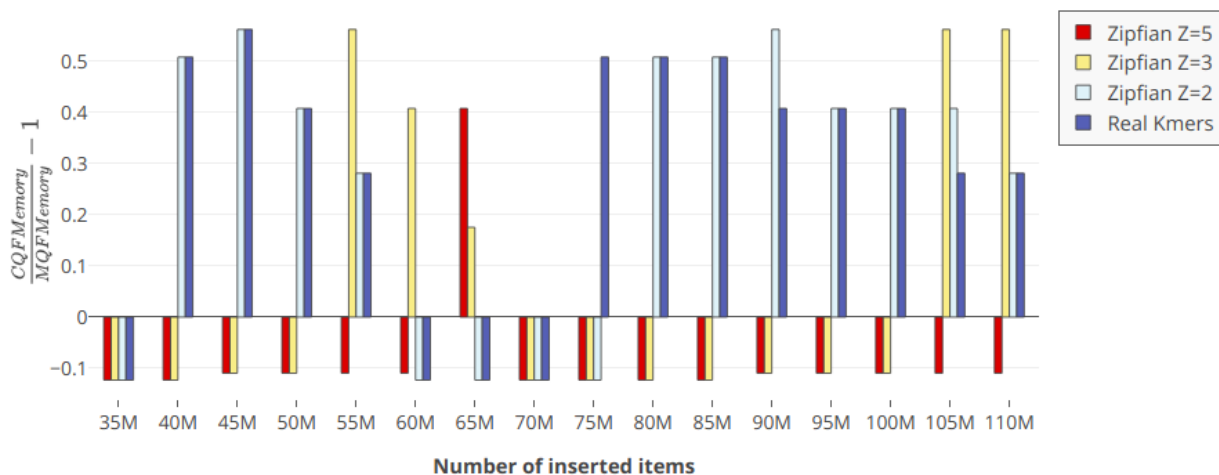
109   structures. MQF had lower loading factors for all tested datasets with better performance with

110   more uniform distributions (The further from the 45° line the better the MQF).

111

112   **MQF is usually more memory efficient than CQF**

113       Progressively increasing numbers of items were sampled from the real and Zipfian-

114   simulated datasets.  The smallest CQF and MQF to store the same number of items from each

115   dataset were created. To do that, the q parameter of CQF versus the q and $F_{size}$ parameters of

116   MQF were calculated empirically. MQF was more memory efficient for real k-mers and Zipfian-

117   simulated distributions with low coefficients in 75% of the cases (Figure 2). The tuning of the

118   $F_{size}$ enabled MQF to grow in size gradually compared to CQF which has to double in size to fit

119   the minimal increase in items beyond the capacity of a given q value (Supplementary Figure 1).
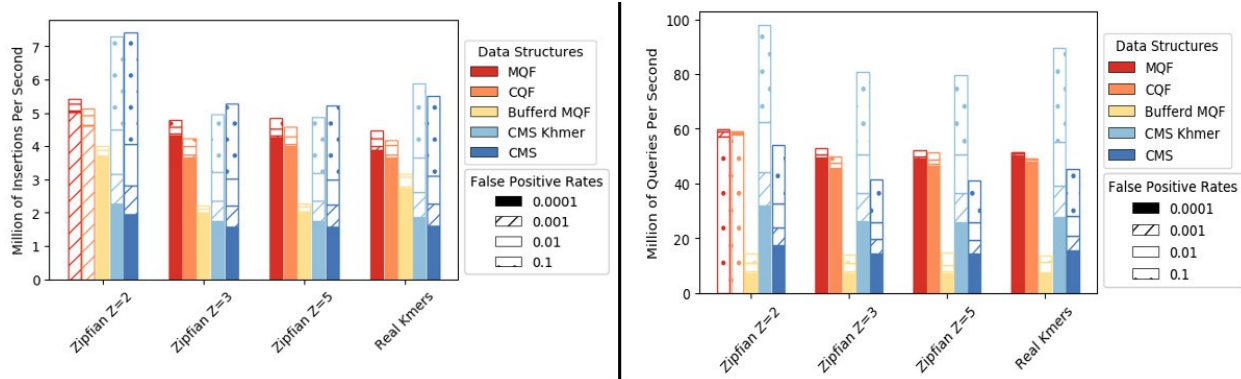
120



121

122   **Figure 2: Memory consumption comparison between CQF and MQF.** The graph compares

123   the memory consumption of the smallest CQF and MQF that fits different datasets. The bigger

124   the value on the y-axes, the more memory the MQF saved.

125

126    **MQF is faster than CQF and low-FPR CMS** The in-memory and buffered MQFs were

127    evaluated for speed of insertion and query in comparison to three in-memory counting

128    structures: CQF, the original CMS (12), and khmer's CMS (13). To test the effect of FPR on the

129    performance, the experiment was repeated for 4 different FPRs (0.1, 0.01, 0.001, 0.0001). All

130    tested structures were constructed to have approximately the same memory space except for

131    buffered MQF which used only one-third of this memory for buffering while the full-size filter is

132    on the disk. MQF is guaranteed to hold the same number of items as a CQF having the same

133    number of slots. The number of slots in CQF was chosen so that the load factor was more than

134    85% and the MQFs were created with an equal number of slots. Items were sampled for

135    insertion from the real and Zipfian-simulated datasets. After finishing the insertion, to assess the

136    query rate, 5M items from the same distribution as the insertion datasets were queried. Half of

137    the query items didn't exist in the insertion datasets.

138    MQF has a faster insertion and query rates compared to CQF with minimal, if any, effect

139    of the FPR on either structure. The performance of CMS is better with higher FPR and Khmer's

140    implementation of CMS doubles the query rate of the original one. However, MQF is always

141    faster than both CMS unless the FPR is more than 0.01 (Figure 3).

142

143



144    **Figure 3: Performance comparison of four data-structures:** MQF, CQF, buffered MQF,

145  Khmer implementation of CMS, and Original implementation of CMS: Insertions rate (left panel)

146  and query rate (right panel).

147

148  **MQF outperforms CMS in real-world problems**

149      Khmer is a software package deploying a new implementation of CMS for k-mer

150  counting, error trimming and digital normalization (13). To test MQF in real-life applications, we

151  assessed the performance of the Khmer software package using CMS (13) versus our new

152  implementation using MQF (https://github.com/dib-lab/khmer/tree/MQFIntegration2). A real RNA

153  seq dataset with 51 million reads from the Genome in a Bottle project (14) was used for error

154  trimming and digital normalization; two real-world applications that involve both k-mer insertions

155  and queries. An exact MQF was used to create a benchmark for the approximate data

156  structures. It took 5Gb RAM to create the data structure and 45 and 43 minutes to perform

157  trimming and digital normalization respectively. The optimal memory for MQF and the optimal

158  number of hash functions for CMS were calculated to achieve the specified false-positive rates.

159  The CMS was constructed with the same size as the corresponding MQFs. The CMS and MQF

160  versions of Khmer were compared regarding the speed and accuracy (Table 1).

161

162

163

164

165

166

167

168

169

170

171

| FPR | Memory in GB | Error Trimming | | | | Digital normalization | | | | Error Bound in CMS | Hash func. In CMS |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Time in Min. | | Missed reads with Errors | | Time in Min. | | Reads kept by Error | | | |
| | | MQF | CMS | MQF | CMS | MQF | CMS | MQF | CMS | | |
| $10^{-1}$ | 1.8 | 42 | 39 | 11011 | 445817* | 39 | 37 | 3253 | 31143* | 13,11 | 3 |
| $10^{-2}$ | 2.6 | 43 | 48 | 1304 | 404354 | 41 | 45 | 416 | 24987 | 14 | 5 |
| $10^{-3}$ | 3.4 | 44 | 61 | 130 | 311464 | 42 | 54 | 58 | 21000 | 15 | 7 |
| $10^{-4}$ | 4.5 | 44 | 75 | 3 | 292746 | 42 | 68 | 4 | 18449 | 16 | 10 |
| Exact | 5 | 45 | - | 0 | - | 43 | - | 0 | - | - | - |

172 **Table 1: Khmer performance in error trimming and digital normalization using MQF and**

173 **CMS.** *Percentages of wrong decisions made by CMS at FPR = 0.1 in error trimming and digital

174 normalization are 0.8% and 0.13% of the total number of decisions versus 0.02% and 0.01%

175 made by MQF.

176 **MQF is faster than MPHF**

177 MPHF is constructed by default to fit the input k-mers while MQF would have different

178 load factor that might affect its performance. To address this question, four growing subsets of

179 real k-mers were inserted into MQFs of size 255 MB to achieve 60%, 70%, 80%, and 90% load

180 factors. MPHFs were constructed with sizes ranging from 15 to 22 MB to fit the four datasets. All

181 data structures were queried with 35M existing k-mers and the query times were reported. The

182 MQFs were ~10 folds faster than the MPHFs. The query time of the MQF was invariable over

183 the different load factors (Supplementary Figure 2).

# Discussion

184

185 MQF is a new variant of counting quotient filters with novel counting and labeling

186 systems. The new counting system increases memory efficiency as well as the speed of

187 insertions and queries for a wide range of data distributions. The labeling system provides a

188    flexible framework for labeling the member items while maintaining good data locality and a

189    concise memory representation.

190

191        MQF is built on the foundation of CQF. MQF has the same ability to behave as an exact

192    or approximate membership query data structure while tracking the count of its members. The

193    insertion/query algorithm developed for CQF enables this family of compact hashtables to

194    perform fast under high load factor (up to 95%) (4). CQFs are designed to work best for data

195    from high Zipfian distributions. However, previous k-mer spectral analysis of RNAseq datasets

196    showed substantial deviations from a Zipfian distribution in thousands of samples(6). Such

197    variations in distribution are expected given the variety of biosamples, the broad spectrum of

198    sequencing techniques, and different approaches to data preprocessing.

199        MQF implements a new counting system that allows the data structure to work efficiently

200    with a broader range of data distributions. The counting system adopts a simple encoding

201    scheme that uses a fixed small space alone or with a variable number of the filter's slots to

202    record the count of member items (Figure 4). Items with small counts utilize the small fixed-size

203    counters. Therefore, slots, used to be consumed by CQF as counters for these items, are freed

204    to accommodate more items in the filter. The MQF's load factor grows slower than CQF with all

205    distributions except the extreme Zipfian case (Z=5) where the load factor is almost the same

206    (Figure 1). This is why the memory requirement for MQFs is usually smaller compared to CQFs

207    under most distributions despite the extra space taken by the fixed counters (Figure 2). The size

208    of the fixed-size counter is constant independent from the slot size, therefore the memory

209    requirement for this counter will be trivial with big slots for smaller FPRs and almost negligible in

210    the exact mode. However, this fixed-size counter comes with an additional advantage for MQF.

211    Tuning the size of the fixed-size counter enables the filter to accommodate more items with a

212    slightly larger slot size. This allows the memory requirement for MQF to grow gradually instead

213    of the obligatory size doubling seen in CQF (Figure 2 and Supplementary Figure 1).

214    Moreover, the new counting scheme in MQF is simplified compared to that of the CQF.

215    MQF defines the required memory for any item based solely on its count. Therefore, an

216    accurate estimation of the required memory for any dataset can be done extremely quickly by

217    an approximate estimation of data distribution(15)(16). This is unlike CQF which needs to add a

218    safety margin to account for the special slots used by the counter encoding technique since it is

219    impossible to estimate the number of these slots.

220

221    Regarding the speed of insertions and queries, MQF is slightly faster than CQF (Figure

222    3). This could be explained partially by the lower load factor of MQF and partially by the

223    simplicity of the coding/decoding scheme of its counting system. Both MQF and CQF are faster

224    than CMS unless the target FPR is really high (e.g. FPR > 0.1) (Figure 3). CMS controls its FPR

225    by increasing the number of its hash tables requiring more time for insertions and queries to

226    happen. In comparison, quotient filters use always one function but with more hash-bits to

227    control the FPR, with a minimal effect on the insertion/query performance (Figure 3). With high

228    FPR (e.g. FPR = 0.1), CMS uses fewer hash functions and is better performing than MQF. A

229    quotient filter or CMS with a FPR = $\delta$ should have the same probability of item collisions.

230    However, the quotient filter will be more accurate because CMS has another type of error with a

231    probability (1-$\delta$), which incorrectly increases the count of its items. This error is a "bounded

232    error" with a threshold that inversely correlates with the width of the CMS(12). In another sense,

233    some applications might deploy CMS with a smaller table's width to be more memory efficient

234    than MQF if the application can tolerate a high bounded error.

235    Buffered MQF can trade some of the speed of insertions and queries for significant

236    memory reduction by storing data on disk. The buffered structure was developed to make use of

237    the optimized sequential read and write on SSD. The buffered structure processes most of the

238    insertion operations using the bufferMQF that resides on memory, thereby limiting the number

239    of access requests to the MQF stored on the SSD hard drive. Sequential disk access happens

240    when the bufferMQF needs to be merged to the disk. This approach is very efficient for

241    insertions but not for random queries which require more frequent SSD data access. In k-mer

242    analysis of huge raw datasets, buffered MQF can be used initially to filter out the low abundant

243    k-mers (i.e. likely erroneous k-mers), then an in-memory MQF holding the filtered list of k-mers

244    could be used for subsequent application requiring frequent random queries. This allows

245    multistage analyses where a first pass eliminates likely errors (17–20).

246

247    CMS is commonly used for online or streaming applications as long as their high error

248    rate can be tolerated (21). MQF has a better memory footprint in the approximate mode for

249    lower error rates and thus can compute with CMS for online applications. A major advantage of

250    quotient filters compared to CMS is the dynamic resizing ability in response to the growing input

251    dataset (4). The buffered version of MQF can be very useful when the required memory is still

252    bigger than the available RAM. We should, however, notice that online applications on MQF

253    cannot make use of the memory optimization that could be achieved with an initial estimation of

254    the filter parameters. A new version of the Khmer software that replaces CMS with MQF proves

255    the new data structure more efficient in real-life applications. The MQF version is faster than the

256    one with CMS unless the target FPR is high. Also, MQF is always more accurate than CMS

257    although both structures have the same FPR. This behavior of CMS is due to the high error

258    bound of its counts.

259

260    MQF comes with a novel labeling system that supports associating each k-mer with

261    multiple values. There are two types of labels: Internal labels adjacent to each item to achieve

262    the best cache locality but has a fixed size and thus practically useful when a small size label is

263    needed. The second labeling system is to label the k-mers with one or more labels stored in

264    external arrays while using the k-mer order in the MQF as an index. External labeling is very

265    memory efficient mimicking the idea of the minimal perfect hash function (MPHF) (10,11).

266   MPHF undoubtedly has the least memory requirement of all the associative data structures (11).

267   However, MQF has better performance in both the construction and query phases. For

268   construction, both structures require initial k-mer counting. MQF needs just an extra O(N)

269   operation to update the block labels where N is the number of its unique k-mers. MPHF has to

270   read then rehash the list of unique k-mers possibly more than once which makes it slower than

271   MQF. For query operations, MQF is 10x faster regardless of the load factor of MQF

272   (supplementary figure 2).

273         Furthermore, MQF offers more functionality and has fewer limitations than MPHF. MQF

274   is capable of labeling a subset of its items which saves significant space for many applications.

275   For example, k-mer analysis applications may want to only label the frequent k-mers, as an

276   intermediate solution between pruning all the infrequent k-mers and labeling all the k-mers.

277   Moreover, MQF allows online insertions and deletions of items as well as merging of multiple

278   labeled MQFs (See the methods) while MPHF - which doesn't store the items - needs to be

279   rebuilt over the whole dataset, which requires reading and rehashing the datasets. Furthermore,

280   MQF can be exact, while MPHF has false positives when queried with novel items that don't

281   belong to the indexed dataset.

## 282   Conclusions

283         MQF is a new counting quotient filter with a simplified encoding scheme and an efficient

284   labeling system. MQF adapts well to a wide range of k-mer datasets to be more memory and

285   time-efficient than its predecessor in many situations. A buffered version of MQF has a fast

286   insertion algorithm while storing most of the structure on external memory. MQF combines a

287   fast access labeling system with MPHF-like associative functionality. MQF performance,

288   features, and extensibility make it a good fit for many online algorithms of sequence analysis.
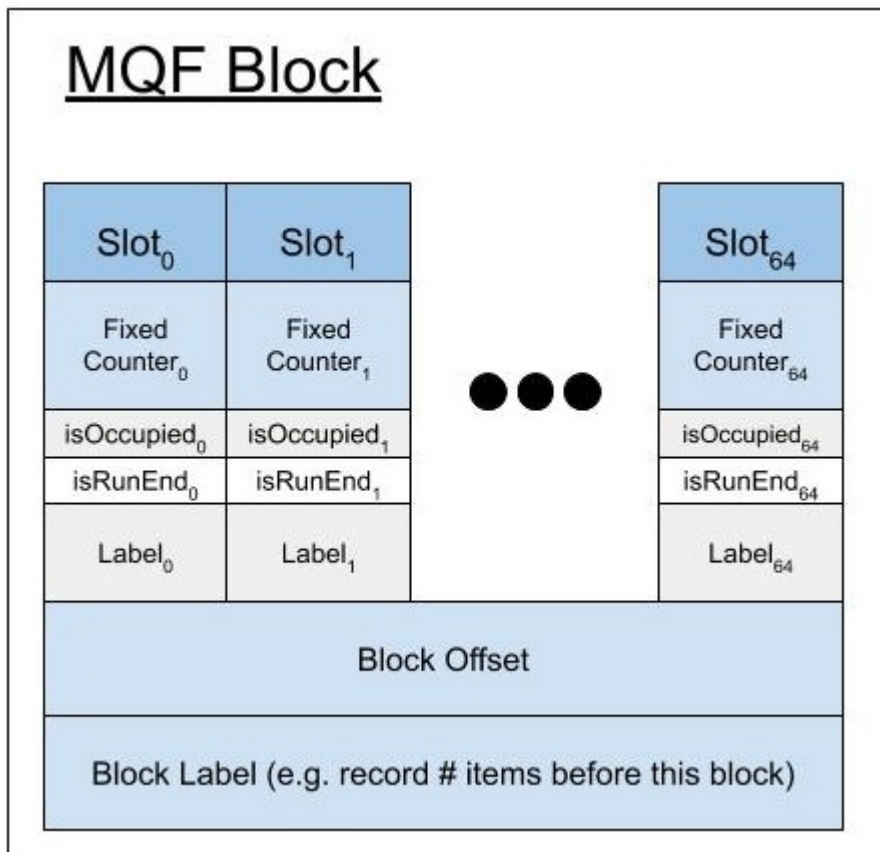
289

# Methods

290

291

292   **MQF Data structure:** MQF has a similar structure to CQF with a different scheme of metadata

293   that enables different counting and labeling systems (Figure 4). Like the CQF, the MQF requires

294   2 parameters, $r$ and $q$, and creates an array of $2^q$slots; each slot has $r$-bits. In Q MQF, $Q_i$  is the

295   slot at position i where $i =$  $1 \dots 2^q$. The MQF maintains the block design of CQF where each

296   block has 64 slots with their metadata and one extra byte of metadata called *Offset* to enhance

297   the query of items(4). Both MQF and CQF have two metadata bits to accompany each slot:

298   $isRunEnd_i$ and $isOccupied_i$. In the MQF, each slot $i$ has extra metadata, a fixed-size counter

299   with a value ($F_i$) and a configurable size (F$_{size}$).  There are also two optional fixed-size parts of

300   metadata allocated to allow different styles of labeling. Every slot has specific labeling ($ST_i$) with

301   a configurable size ($ST_{size}$>=0), and every block ( $j$ ) has an optional space of a configurable size

302   designed to store the number of items in the previous blocks.

303

304         The MQF uses the same insertion/query algorithm of CQF (4). In brief, suppose item $I$,

305   repeated $c$ times, is to be inserted into Q. A hash function $H$ is applied to $I$ to generate a $p$-bit

306   fingerprint ($H(I)$). $H(I)$ value is split into two parts, a quotient and remainder. The quotient ($q_i$) is

307   the most significant $q$ bits while the remainder ($r_i$) is the remaining least significant r bits. The

308   filters store $r_i$ in a slot $Q_j$ where $j \geq q_i$. One or more slots can be used to store the count of the

309   same item. If the required slots for the item or its count are not free, all the consecutive

310   occupied slots starting from this position will be shifted to free the required space. All items

311   having the same $q$ are stored into consecutive slots and are called a run. Items in the run are

312   sorted by $r_i$, and *isRunEnd* of the last slot in the run is set to one. *isOccupied ($q_i$)* is set to one if

313    and only if there is a run for $q_i$. Therefore, there is one bit set to one in each *isOccupied* and

314    *isRunEnd* for each run. To query item *I*, a Rank and Select method is applied on the metadata

315    arrays to get the run start and end for $q_i$. Then all the items in the run are searched linearly for

316    the slot containing $r_i$. The subsequent one or more slots can be decoded to get the count of

317    item *I*. CQF uses a special encoding scheme to recognize these counting slots but MQF utilizes

318    the fixed-counter metadata element (see below).

319



320

321    **Figure 4: MQF block structure.** Each MQF block contains 64 slots with their metadata, a one-

322    byte block offset, and configurable size space to hold the number of items inserted in the filter

323    before the current block. The metadata of each slot consumes *r* bits, one bit for each

324    *isOccupied* and *isRunEnd* metadata, and configurable *f*-bits and *t*-bits for the fixed counter and

325    the slot-specific label respectively.

326

327   **_Counting scheme_**: MQF uses two types of counters for storing the values of the count ($c$): A

328   small fixed-size space ($F_i$) is slot specific and used to store the count of the item in its own slot

329   if this count is smaller than $F_{max}$, where $F_{max}$ is the maximum possible value for the fixed space.

330   A variable size space ($V_i$) is composed of one or more slots next to the item's slot and is used

331   to store larger values. For an item with high count $c$, the number of required slots for $V_i$ is

332   calculated as $\dfrac{\log_2\left(c - F_{max}\right) - F_{sizes}}{r}$ slots. The $F_i$ spaces of this item's slot and its $V_i$ slots are

333   used to mark the last slot for the item where all of them will be saturated to $F_{max}$ except the last

334   one (Figure 5). This counting scheme can be summarized into 2 rules:

335       _Rule 1:_ MQF requires $F_i < F_{max}$ if and only if _i_ is the index of the item's last slot.
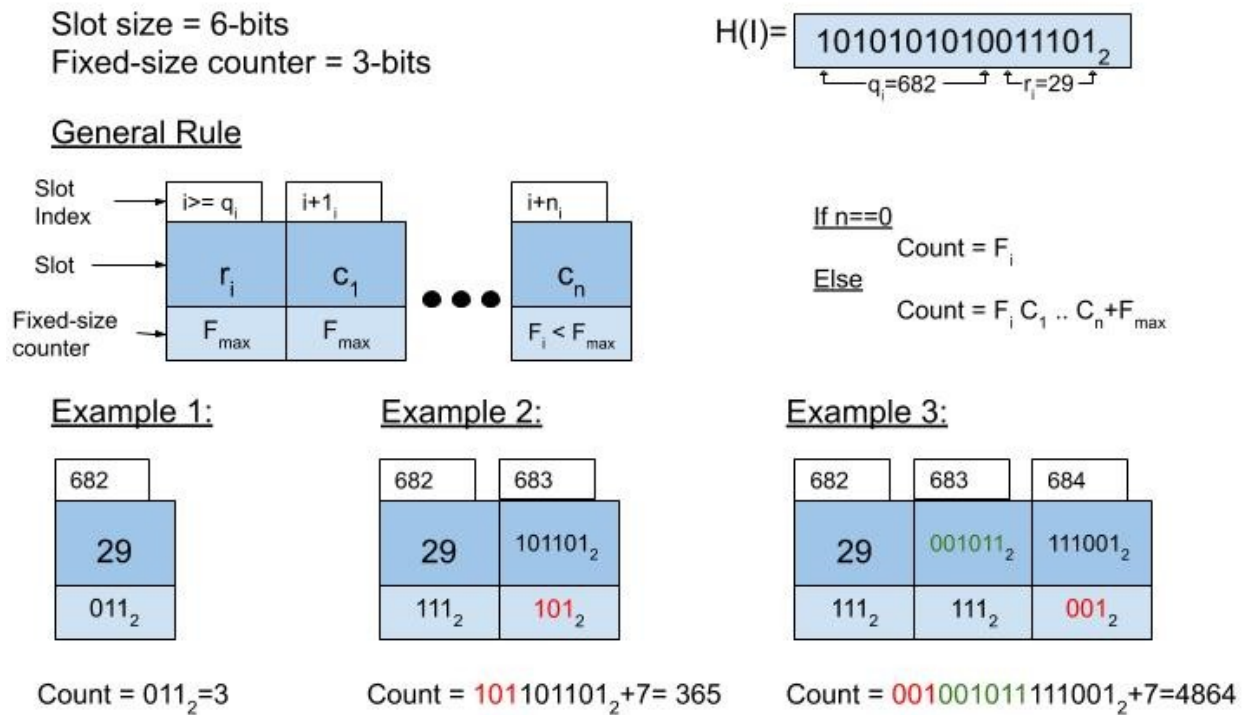
336       _Rule 2:_ If $c < F_{max}$, _c is stored in $F_i$ only. Otherwise, c - F$_{max}$ is stored in $V_i$_

337

338   In comparison to the CQF, the MQF does not use special slots to resolve ambiguities, which is

339   more memory efficient. The counter encoding algorithm is described in Supplementary Figure 3.

340

Slot size = 6-bits
Fixed-size counter = 3-bits

$H(I)=$ | $1010101010011101_2$ |
$\underset{q_i=682}{\longleftrightarrow}$ $\underset{r_i=29}{\longleftrightarrow}$

General Rule

| | Slot Index | $i>=q_i$ | $i+1_i$ | | $i+n_i$ |
|---|---|---|---|---|---|
| Slot | | $r_i$ | $C_1$ | $\bullet\bullet\bullet$ | $C_n$ |
| Fixed-size counter | | $F_{max}$ | $F_{max}$ | | $F_i < F_{max}$ |

If n==0
    Count = $F_i$
Else
    Count = $F_i\ C_1\ ..\ C_n+F_{max}$

Example 1:

| 682 |
|---|
| 29 |
| $011_2$ |

Count = $011_2$=3

Example 2:

| 682 | 683 |
|---|---|
| 29 | $101101_2$ |
| $111_2$ | $101_2$ |

Count = $101101101_2+7$= 365

Example 3:

| 682 | 683 | 684 |
|---|---|---|
| 29 | $001011_2$ | $111001_2$ |
| $111_2$ | $111_2$ | $001_2$ |

Count = $001001011111001_2+7$=4864

341

342 **Figure 5: MQF counters encoding scheme.** Items and their counts are stored in n slots and n

343 fixed counter as shown in the general rule. Each example stores the same item but different

344 count (count = 3, 365, or 4864).

345

346 # Parameter Estimation

347     For offline counting applications, the MQF parameters (q, r, $F_{size}$) can be even more

348 optimized for each dataset to create the most memory-efficient filter that has enough slots to fit

349 all unique items and their counts. The $q$ parameter defines the number of slots (N) in MQF

350 where $q=\log_2(N)$. The required numbers of slots for items and their count can be estimated from

351 the cardinality of the target dataset, as with CQF. The $r$ parameter is calculated from the

352 equation $r = p\text{-}q$ where p is the total number of hash-bits used to represent each item. In the

353 exact mode, $p$ equals the exact output of a reversible hash function. In the inexact mode, $p$ is
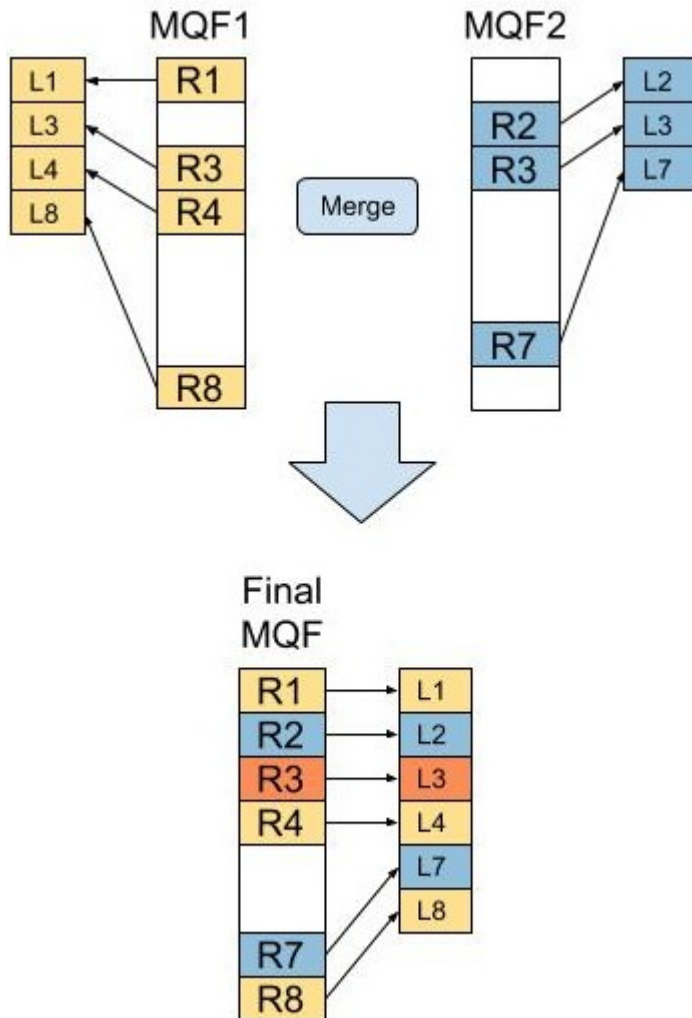
354    controlled by the target FPR ($\delta$) according to the equation $p = \log_2 \dfrac{N}{\delta}$ described before (4). The

355    $F_{size}$ parameter defines the size of the fixed-size counter. This is critical because if a given MQF

356    has too few slots for items in a dataset, the bigger MQF would have to double the number of

357    slots causing a big jump in the memory requirement. To avoid that jump, MQF can use larger

358    fixed size counters to decrease the number of slots required in counting on the expense of a

359    slight increase in the slot size.

360

## Labeling System:

362      MQF can map each item to its count as well as other values, which we call "labels".

363    Labels in MQF have two different systems. An internal labeling system stores the associated

364    value for every key in the data structure, like a hash table. This label has a fixed size defined at

365    the initialization of the MQF and is practically useful when a small size label is needed (e.g. one

366    or two bits). The second labeling system labels the block. We use this label to store the number

367    of items inserted in the MQF before each block. This enables labeling the items of the filter by

368    separate arrays matching the order of the items in the filter, a behavior that can act as a minimal

369    perfect hash function (11). The naive way to compute the items' order is to find the item in the

370    MQF and iterate backward until the beginning of the filter to count the number of the preceding

371    items, which is an O(N) operation. The MQF stores the number of items that exist before each

372    block; therefore, the MQF iterates only to the beginning of each block, which is an O(1)

373    operation. The number of previous items for each block is computed after the MQF is

374    constructed. Any additional insertions or deletions of items would only require re-calculation of

375    the block label values with no need to re-analyze the original data. Moreover, labeled MQFs can

376    be updated by merging multiple labeled MQFs and their external labeling arrays. External label

377    arrays need to be merged after merging the labeled MQFs. To do so, the new items' order is

378    recomputed in the final MQF. Then, labels in the input external arrays can be copied into a new

379    external array according to the new item order. Such a function has to consider resolving the

380    conflicts of items happening in multiple-input MQF and labeled by different external labels

381    (Figure 6).

382



383

384    **Figure 6: Merging MQFs with external labels.** $R_i$ is the remaining part of item i, and $T_i$ is the

385    external label of the item. Merging the input MQF produces a final MQF with a new order of its

386    member items. All labels in the input external arrays are copied into a new external array

387    according to this new order of the items. However, the implementation of the merge function has

388    to resolve the conflict of R3 labels which exist in both input structures with two labels.

## Buffered MQF

The Buffered MQF is composed of two MQF structures: a big structure stored on SSD called onDiskMQF, and an insertion buffer stored in the main memory called bufferMQF. OnDiskMQF uses stxxl vectors(22) because of the performance of their asynchronous IO. The bufferMQF is used to limit the number of accesses on the OnDiskMQF and change the access pattern to the on-disk structure from random to sequential. As shown in the insertion algorithm in Figure 7, all the insertions are done first on bufferMQF; when it is full, the items are copied from bufferMQF to OnDiskMQF, and bufferMQF is cleared. The copy operation edits the onDiskMQF in a serial pattern which is preferred while working on SSD because many edits will be grouped together in one read/write operation. Figure 8 shows the query algorithm. The queried items are inserted first to temporary MQF and sequential access is done to query the items from the OnDiskMQF. The final count is the sum of the bufferMQF and the ondiskMQF.

---

**Algorithm 2** Buffered MQF Insertion

1: **procedure** INSERT($onDiskMQF, bufferMQF, item$)
2:     $mqf\_insert(bufferMQF, item)$
3:     **if** $mqf\_space(bufferMQF) > 90$ **then**
4:         **for all** $i \in bufferMQF$ **do**
5:             $mqf\_insert(bufferMQF, i)$
6:         **end for**
7:         $mqf\_clear(bufferMQF)$
8:     **end if**
9: **end procedure**

---

**Figure 7: Buffered MQF insertion algorithm.** Insertion Algorithm for inserting items in the Buffered MQF. It inserts the item in the in-memory data structure. The on-memory structure is merged into the on-disk structure when it is filled.

```
Algorithm 3 Buffered MQF Query
 1: procedure QUERY(onDiskMQF, bufferMQF, list_item)
 2:     for all i ∈ listItems do
 3:         mqf_insert(tmpMQF, i)
 4:     end for
 5:     for all i ∈ tmpMQF do
 6:         counts[i] ← mqf_query(onDiskMQF, i)
 7:         counts[i] ← counts[i] + mqf_query(bufferMQF, i)
 8:     end for
 9:     return counts
10: end procedure
```

407

408 **Figure 8: Buffered MQF query algorithm.** Query algorithm for retrieving counts for a list of

409 items in the Buffered MQF. First, insert all the items in the list into a temporary MQF. Second,

410 iterate over the list of items in the temporary MQF and query both the in-memory and on-disk

411 structures.

412

## Experimental Setup of Benchmarking

414 Five datasets were used in the experiments to cover most of the bioinformatics

415 applications. Three datasets called z2, z3, and z5 were simulated to follow Zipfian distribution

416 using three different coefficients: 2, 3, and 5 respectively. The bigger the coefficient the more

417 singletons in the dataset (23). A fourth dataset was simulated from a uniform distribution with a

418 frequency equal to 10. One more dataset, named k-mers, represented real k-mers generated in

419 the ERR1050075 RNA-seq experiment from humans(24). Experiments were conducted to

420 compare the performance, memory, and accuracy of MQF with the state-of-the-art counting

421 structures CQF, CMS, and MPHF.  Unless stated otherwise, CQF and MQF used the same

422 number of slots, and the same slot size while the fixed counter of MQF was set to two. The slot

423 size was calculated to achieve the target FPR as described in the parameter estimation section

424 (see Methods). To create comparable CMS, the number of the tables in the sketches was

425 calculated using $\ln\dfrac{1}{\delta}$ as described before (12). The table width was calculated by dividing the

426 MQF size by the number of tables. The MPHF was created using the default options in the

427 BBhash repo (https://github.com/rizkg/BBHash). An Amazon AWS t3.large machine with Ubuntu

428 Server 18.04 was used to run all the experiments. The instance had 2 VCPUS and 8GB RAM

429 with a 100GB provisioned IOPS SSD attached for storage. All codes used in the experiments

430 can be accessed through the MQF GitHub repository (https://github.com/dib-lab/2020-paper-

431 mqf-benchmarks).

432

# List of abbreviations

434 • MQF: mixed-counters quotient filter.

435 • CQF: counting quotient filter.

436 • FPR: false positive rate.

437 • CMS: count-min sketch.

438 • MPHF: Minimal Perfect Hash Functions

439

# Declarations

**Ethics approval and consent to participate**

442 Not applicable

**Consent for publication**

444 Not applicable

## Availability of data and materials

The datasets used in Benchmarking are available in the "2020-paper-mqf-benchmarks"
repository.
https://github.com/dib-lab/2020-paper-mqf-benchmarks

## Competing interests

The authors declare that they have no competing interests

## Funding

Not applicable

## Authors' contributions

TAM and MS developed theoretical formalism. MS carried out the implementation and
benchmarking. TAM conceived the original idea and supervised this work. All authors
contributed to the writing of the manuscript.

## Acknowledgements

Not applicable



# References

1.  Kolajo T, Daramola O, Adebiyi A. Big data stream analysis: a systematic literature review.
    Journal of Big Data. 2019 Jun 6;6(1):1–30.

2.  Matias Y, Vitter JS, Young NE. Approximate data structures with applications. In:
    Proceedings of the fifth annual ACM-SIAM symposium on Discrete algorithms. Society for
    Industrial and Applied Mathematics; 1994. p. 187–94.

3.  Bloom BH. Space/time trade-offs in hash coding with allowable errors. Commun ACM.
    1970 Jul 1;13(7):422–6.

4.  Pandey P, Bender MA, Johnson R, Patro R. A General-Purpose Counting Filter. In:
    Proceedings of the 2017 ACM International Conference on Management of Data - SIGMOD
    '17 [Internet]. 2017. Available from: http://dx.doi.org/10.1145/3035918.3035963

5.  Manekar SC, Sathe SR. A benchmark study of k-mer counting methods for high-throughput
    sequencing. Gigascience [Internet]. 2018 Dec 1;7(12). Available from:
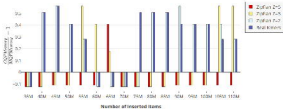    http://dx.doi.org/10.1093/gigascience/giy125

474  6.  Yu Y, Liu J, Liu X, Zhang Y, Magner E, Lehnert E, et al. SeqOthello: querying RNA-seq
475      experiments at scale. Genome Biol. 2018 Oct 19;19(1):167.

476  7.  Cormode G, Muthukrishnan S. An Improved Data Stream Summary: The Count-Min Sketch
477      and Its Applications. In: Farach-Colton M, editor. LATIN 2004: Theoretical Informatics.
478      Berlin, Heidelberg: Springer Berlin Heidelberg; 2004. p. 29–38. (Goos G, Hartmanis J, van
479      Leeuwen J, editors. Lecture Notes in Computer Science; vol. 2976).

480  8.  Muggli MD, Bowe A, Noyes NR, Morley PS, Belk KE, Raymond R, et al. Succinct colored
481      de Bruijn graphs. Bioinformatics. 2017 Oct 15;33(20):3181–7.

482  9.  Cormen TH, Leiserson CE, Rivest RL, Stein C. Introduction to Algorithms. MIT Press; 2009.
483      1292 p.

484  10.  Belazzougui D, Botelho FC, Dietzfelbinger M. Hash, Displace, and Compress. In: Fiat A,
485       Sanders P, editors. Algorithms - ESA 2009. Berlin, Heidelberg: Springer Berlin Heidelberg;
486       2009. p. 682–93. (Lecture Notes in Computer Science; vol. 5757).

487  11.  Limasset A, Rizk G, Chikhi R, Peterlongo P. Fast and scalable minimal perfect hashing for
488       massive key sets [Internet]. arXiv [cs.DS]. 2017. Available from:
489       http://arxiv.org/abs/1702.03154

490  12.  Cormode G, Muthukrishnan S. An improved data stream summary: the count-min sketch
491       and its applications. J Algorithm Comput Technol. 2005;55(1):58–75.

492  13.  Crusoe MR, Alameldin HF, Awad S, Boucher E, Caldwell A, Cartwright R, et al. The khmer
493       software package: enabling efficient nucleotide sequence analysis. F1000Res [Internet].
494       2015 Sep 25 [cited 2019 Jun 3];4. Available from: https://f1000research.com/articles/4-900/
495       v1/pdf

496  14.  Zook JM, Salit M. Genomes in a bottle: creating standard reference materials for genomic
497       variation - why, what and how? Genome Biol. 2011 Sep 19;12(1):P31.

498  15.  Flajolet P, Fusy É, Gandouet O, Meunier F. Hyperloglog: The analysis of a near-optimal
499       cardinality estimation algorithm. In: IN AOFA '07: PROCEEDINGS OF THE 2007
500       INTERNATIONAL CONFERENCE ON ANALYSIS OF ALGORITHMS [Internet]. 2007 [cited
501       2018 Nov 19]. Available from: http://citeseerx.ist.psu.edu/viewdoc/summary?
502       doi=10.1.1.76.4286

503  16.  Mohamadi H, Khan H, Birol I. ntCard: a streaming algorithm for cardinality estimation in
504       genomics data. Bioinformatics. 2017 May 1;33(9):1324–30.

505  17.  Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, et al. Mash: fast
506       genome and metagenome distance estimation using MinHash. Genome Biol. 2016 Jun
507       20;17(1):132.

508  18.  Chikhi R, Limasset A, Medvedev P. Compacting de Bruijn graphs from sequencing data
509       quickly and in low memory. Bioinformatics. 2016 Jun 15;32(12):i201–8.

510  19.  Zhang Q, Awad S, Titus Brown C. Crossing the streams: a framework for streaming
511       analysis of short DNA sequencing reads [Internet]. PeerJ PrePrints; 2015 Mar [cited 2020
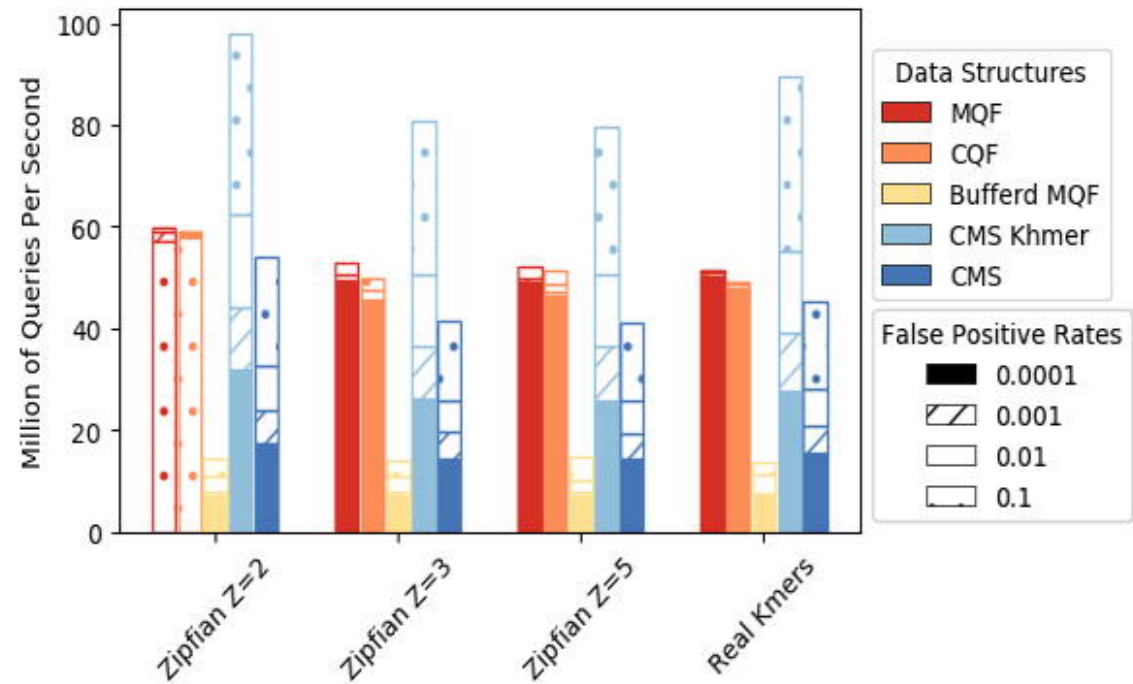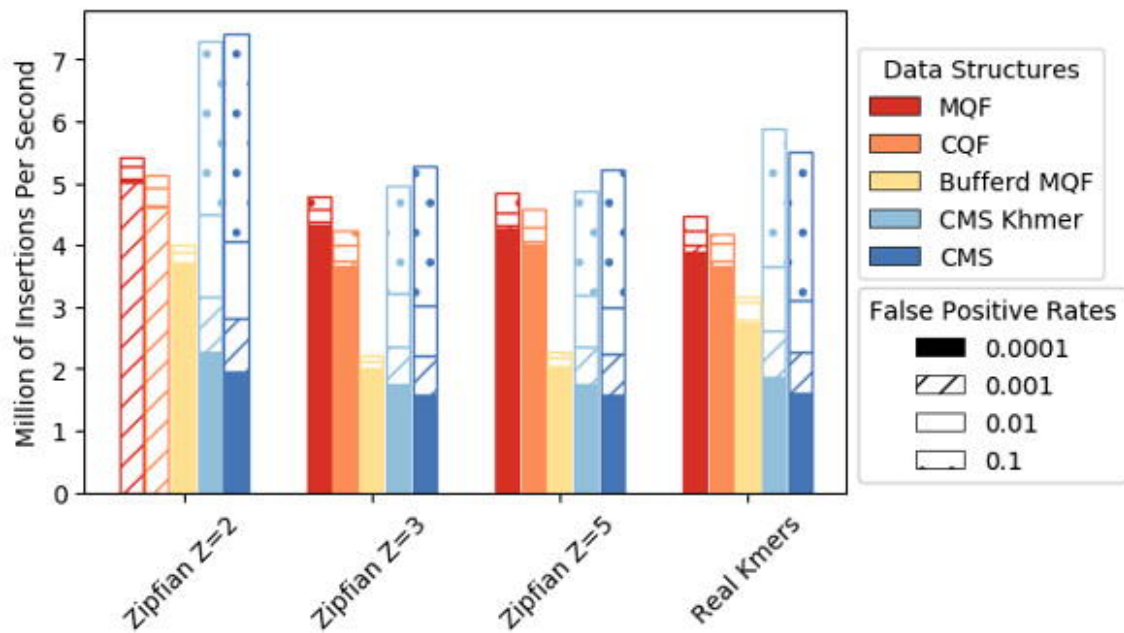512       Jul 23]. Report No.: e1100. Available from: https://peerj.com/preprints/890/

513   20. Titus Brown C, Howe A, Zhang Q, Pyrkosz AB, Brom TH. A Reference-Free Algorithm for
514       Computational Normalization of Shotgun Sequencing Data [Internet]. arXiv [q-bio.GN].
515       2012. Available from: http://arxiv.org/abs/1203.4802

516   21. Muthukrishnan S. Data Streams: Algorithms and Applications. TCS. 2005 Sep 27;1(2):117–
517       236.

518   22. Dementiev R, Kettner L, Sanders P. Stxxl: Standard Template Library for XXL Data Sets.
519       In: Algorithms – ESA 2005. Springer, Berlin, Heidelberg; 2005. p. 640–51. (Lecture Notes in
520       Computer Science).

521   23. Powers DMW. Applications and explanations of Zipf's law. In: Proceedings of the Joint
522       Conferences on New Methods in Language Processing and Computational Natural
523       Language Learning. Association for Computational Linguistics; 1998. p. 151–60.

524   24. Clarke L, Zheng-Bradley X, Smith R, Kulesha E, Xiao C, Toneva I, et al. The 1000
525       Genomes Project: data management and community access. Nat Methods. 2012 Apr
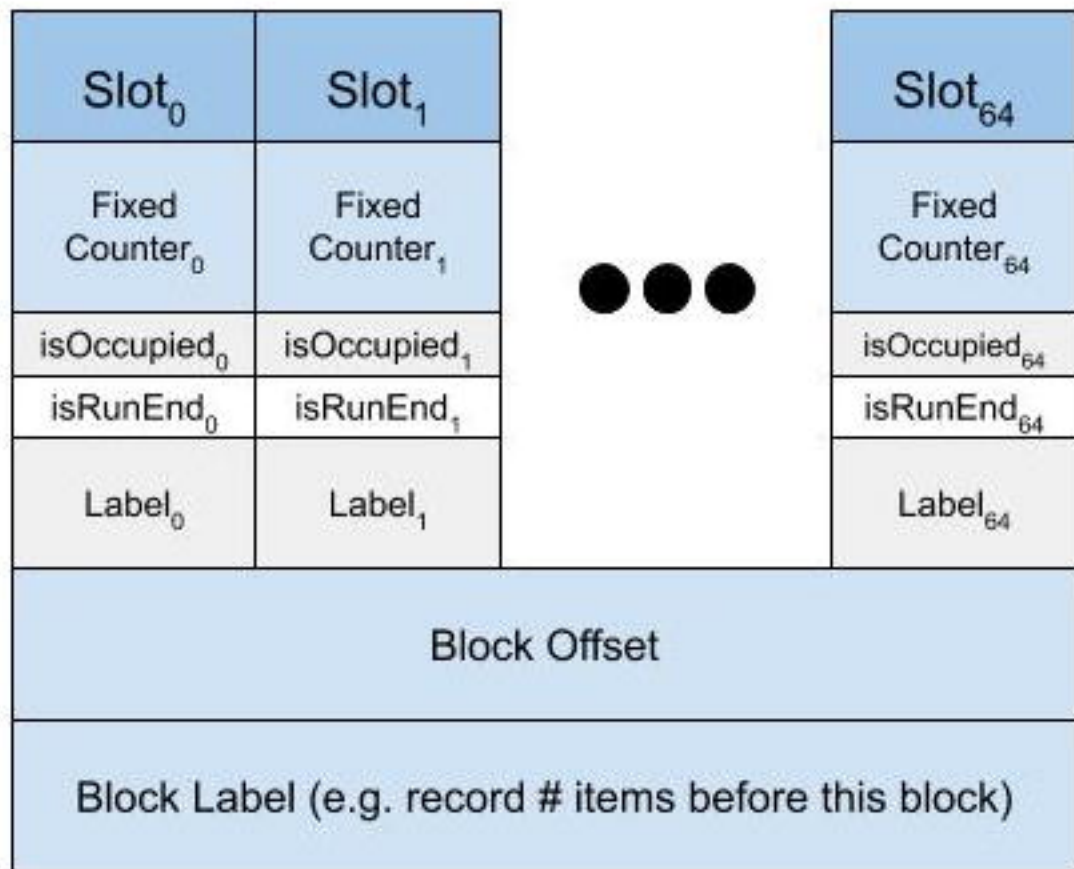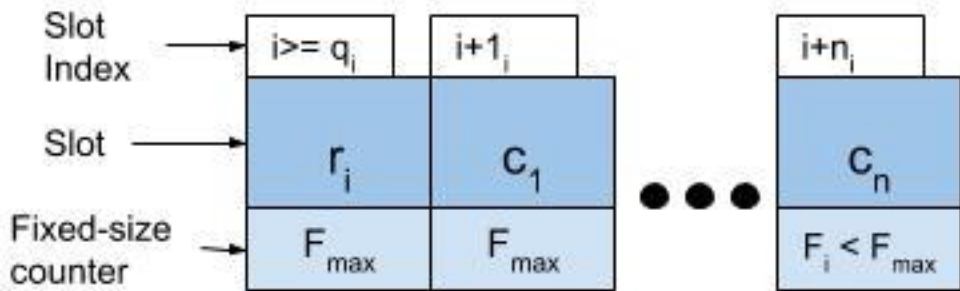526       27;9(5):459–62.

527

Number of inserted items

# MQF Block

| Slot$_0$ | Slot$_1$ | | Slot$_{64}$ |
|---|---|---|---|
| Fixed Counter$_0$ | Fixed Counter$_1$ | ● ● ● | Fixed Counter$_{64}$ |
| isOccupied$_0$ | isOccupied$_1$ | | isOccupied$_{64}$ |
| isRunEnd$_0$ | isRunEnd$_1$ | | isRunEnd$_{64}$ |
| Label$_0$ | Label$_1$ | | Label$_{64}$ |

Block Offset

Block Label (e.g. record # items before this block)

Slot size = 6-bits
Fixed-size counter = 3-bits

$H(I) = $ | $1010101010011101_2$ |

$\underbrace{\qquad}_{q_i=682}$ $\underbrace{}_{r_i=29}$

## General Rule

Slot Index → | $i >= q_i$ | $i+1_i$ | | $i+n_i$ |

Slot → | $r_i$ | $c_1$ | $\bullet\bullet\bullet$ | $c_n$ |

Fixed-size counter → | $F_{max}$ | $F_{max}$ | | $F_i < F_{max}$ |

If n==0

Count = $F_i$

Else

Count = $F_i\, C_1\, .. \,C_n + F_{max}$

## Example 1:

| 682 |
| 29 |
| $011_2$ |

Count = $011_2 = 3$

## Example 2:

| 682 | 683 |
| 29 | $101101_2$ |
| $111_2$ | $101_2$ |

Count = $101101101_2 + 7 = 365$

## Example 3:

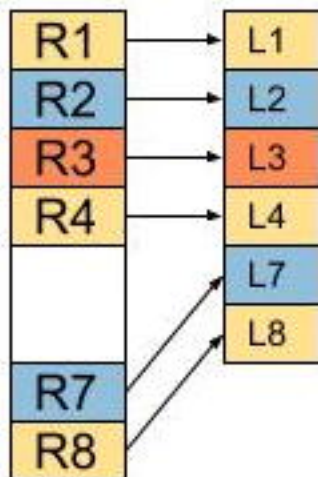| 682 | 683 | 684 |
| 29 | $001011_2$ | $111001_2$ |
| $111_2$ | $111_2$ | $001_2$ |

Count = $001001011111001_2 + 7 = 4864$

MQF1    MQF2

Merge

Final
MQF

**Algorithm 2** Buffered MQF Insertion

```
1: procedure INSERT(onDiskMQF, bufferMQF, item)
2:     mqf_insert(bufferMQF, item)
3:     if mqf_space(bufferMQF) > 90 then
4:         for all i ∈ bufferMQF do
5:             mqf_insert(bufferMQF, i)
6:         end for
7:         mqf_clear(bufferMQF)
8:     end if
9: end procedure
```
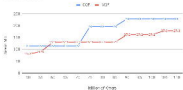
**Algorithm 3** Buffered MQF Query

```
1: procedure QUERY(onDiskMQF, bufferMQF, list_item)
2:    for all i ∈ listItems do
3:        mqf_insert(tmpMQF, i)
4:    end for
5:    for all i ∈ tmpMQF do
6:        counts[i] ← mqf_query(onDiskMQF, i)
7:        counts[i] ← counts[i] + mqf_query(bufferMQF, i)
8:    end for
9:    return counts
10: end procedure
```
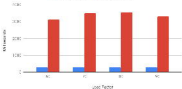
Size Comparison between MQF and CQF

Query time for 36M items

```
 1: procedure ENCODE(Q, r_i, count , start)
 2:     base ← 2^Q                                    ▷ Q.r is #bits in the slot
 3:     fcountMax ← 2^{Q.f} − 1                      ▷ Q.f is #bits in the fixed-size slot
 4:     stack ← ∅
 5:     while count < fountMax − 1 do
 6:         stack.push(count%base)
 7:         count ← count ≫ Q.r                       ▷ bit shift operation
 8:     end while
 9:     stack.push(r_i)
10:     i ← start
11:     while stack ≠ ∅ do
12:         Q_i.r ← stack.pop()                       ▷ Slot of index i
13:         Q_i.f ← fountMax − 1                       ▷ fixed-size counter of index i
14:         i ← i + 1
15:     end while
16: end procedure
```

**Algorithm 1** Counters Encoder