

METHODOLOGY ARTICLE

Open Access



MQF and buffered MQF: quotient filters for efficient storage of k-mers with their counts and metadata

Moustafa Shokrof¹, C. Titus Brown² and Tamer A. Mansour^{2,3*} 

*Correspondence:

drtamermansour@gmail.com

² Department of Population Health and Reproduction, School of Veterinary Medicine, University of California, Davis, CA, USA
Full list of author information is available at the end of the article

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 ~tenfold faster than BBhash, with no need to re-analyze the original data for further insertions or deletions.

Conclusions: 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

Online algorithms effectively support streaming analysis of large data sets, which is important for analyzing data sets with large volume and high velocity [1]. Approximate data structures are commonly used in online algorithms to provide better average space and time efficiency [2]. For example, the Bloom filter supports approximate set membership queries with a predefined false positive rate (FPR) [3]. The count-min sketch (CMS) is similar to Bloom filters and can be used to count items with a tunable rate of overestimation. However, there are a number of problems with Bloom filters and the CMS—in particular, they do not support data locality.

The counting quotient filter (CQF) is a more efficient data structure that serves similar purposes with better efficiency for skewed distributions and much better data locality



© The Author(s) 2021. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

[4]. The CQF is a recent variant of quotient filters that tracks the count of its items using a variable size counter. As a compact hashtable, CQF can perform in either probabilistic or exact modes and supports deletes, merges, and resizing.

Analysis of k-mers in biological sequencing data sets is an ongoing challenge [5]. K-mers in raw sequencing data often have a Zipfian distribution, and the CQF was built to minimize memory requirements for counting such items [4]. However, this advantage deteriorates in applications that require frequent random access to the data structure, and where the k-mer count distribution may change in response to different sampling approaches, library preparation and/or sequencing technologies. For example, k-mer frequency across 1000 thousands of RNAseq experiments shows different patterns of abundant k-mers [6].

Data structures like CMS [7] and CQF [4] also do not natively support associating k-mers with multiple values, which can be useful for coloring in De Bruijn graphs as well as other features [8–10]. Classical hash tables are designed to associate their keys with a generic data type but they are expensive memory-wise [11]. Tools like jellyfish [12] and CHTKC [13] use lock free hash tables in their k-mer counting algorithms. Minimal Perfect Hash Functions (MPHFs) can provide a more compact solution by mapping each k-mer into a unique integer. These integers can then be used as indices for the k-mers to label them in other data structures [14]. An implementation capable of handling large scale datasets with fast performance requires ~ 3 bits per element [15]. However, such a concise representation comes with a high false-positive rate on queries for non-existent items. Moreover, unlike hashtables, MPHf does not support insertions or deletions thus any change in the k-mer set would require rehashing of the original dataset.

In this paper, we introduce the mixed-counters quotient filter (MQF), a modified version of the CQF with a new encoding scheme and labeling system supporting high data locality. We further show how Buffered MQF can be used to scale MQF to solid-state disks. We compare between MQF and the CQF, CMS, and MPHf data structures regarding memory efficiency, speed performance, and applicability to specific data analysis challenges. We further do a direct comparison of the CMS to MQF in the khmer software package for sequencing data analysis, to showcase the benefits of MQF is in real world applications.

Results

MQF has a lower load factor than CQF

The load factor is defined as the actual space utilized divided by the total space assigned for the data structure, and is an important measure of data structure performance. To compare load factors between the CQF and MQF data structures, instances of both structures were created using the same number of slots (2^{27}). Chunks of items from thirteen datasets with different distributions of item frequencies were inserted iteratively to be counted in both data-structures while recording the load factor after the insertion of each chunk. After the insertion of each chunk, both data structures were checked to confirm having the exact same hashes. The experiments stopped when MQF's load factor reached 90%. MQF had lower loading factors for all tested datasets but the difference was minimal for the dataset with the highest Zipfian distribution ($Z=5$). The lower the tested Zipfian distribution the lower the loading factor of MQF.

A lower loading factor enabled MQF to exceed the double CQF capacity with uniform distribution. Analysis of real k-mers from 8 sequencing datasets showed ~40% increase on average in the capacity of MQF in comparison to the matching CQF (Fig. 1, Additional file 1: Figure S1 and Table S1).

MQF is usually more memory efficient than CQF

Progressively increasing numbers of items were sampled from the real and Zipfian-simulated datasets. The smallest CQF and MQF to store the same number of items from each dataset were created. To do that, the q parameter of CQF versus the q and F_{size} parameters of MQF were calculated empirically. Starting with a small q, CQF and MQF were tested to count the items. The q was incremented to find the smallest value that enables the structure to hold the items. MQF was more memory efficient for most real k-mers and Zipfian-simulated distributions (Fig. 2, Additional file 1: Figure S2). The tuning of the F_{size} enabled MQF to grow in size gradually compared to CQF which has to double in size to fit the minimal increase in items beyond the capacity of a given q value (Additional file 1: Figure S3). Therefore, comparing the optimum parameters used to generate both data structures (Additional file 2: Table S2) shows that MQF is always more memory efficient if $F_{size} > 1$, while CQF could be slightly more memory efficient if $F_{size} = 1$.

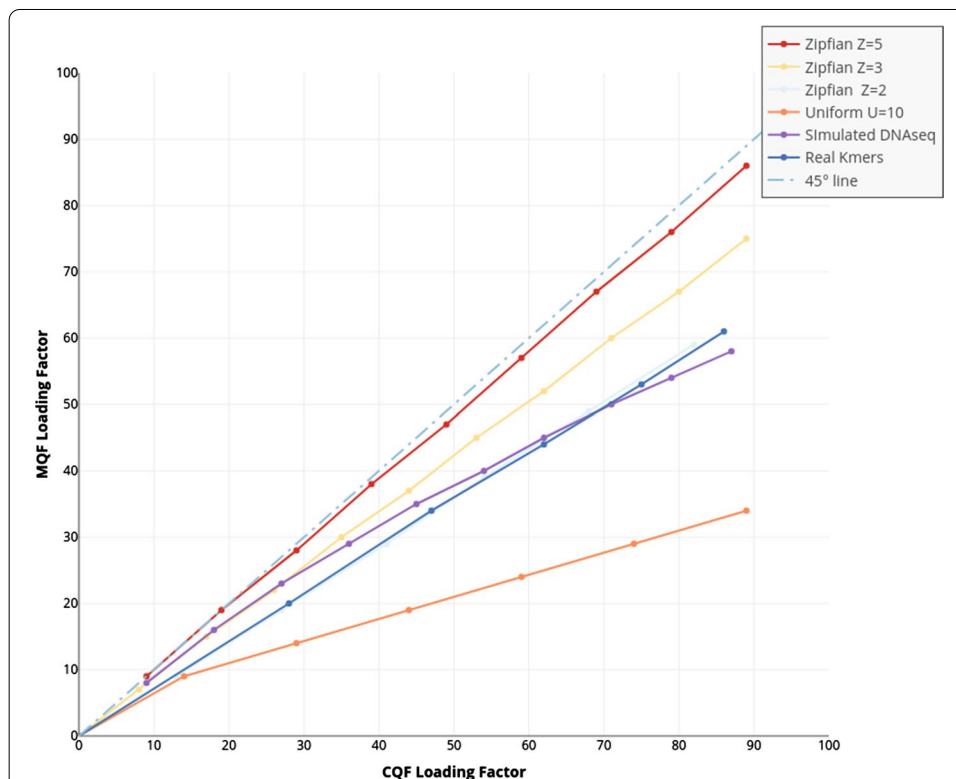
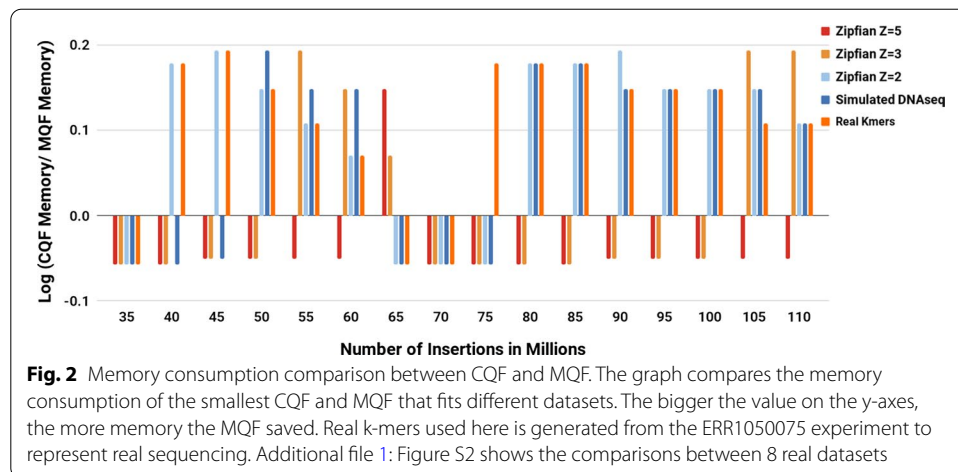


Fig. 1 MQF has a lower load factor compared to CQF. Chunks of items, from different frequency distributions, were inserted iteratively to matching CQF and MQF structures. MQF had lower loading factors for all tested datasets with better performance with more uniform distributions (The further from the 45° line the better the MQF). Real k-mers used here is generated from the ERR1050075 experiment to represent real sequencing. Additional file 1: Figure S1 shows the comparisons between 8 real datasets



This is more likely to happen if sequencing datasets are highly erroneous or very polymorphic (e.g. mitochondrial genome sequencing in SRR12989394). Also, if sequencing has a very low coverage causing most of the k-mers to appear once (e.g. small subsets of whole genome sequencing data in ERR992657). MQF comes with a utility file (<https://github.com/dib-lab/MQF/blob/master/include/utills.h>) to allow predicting the optimum MQF parameters.

MQF is faster than CQF and low-FPR CMS

The in-memory and buffered MQFs were evaluated for speed of insertion and query in comparison to three in-memory counting structures: CQF, the original CMS [7], and khmer's CMS [16]. To test the effect of FPR on the performance, the experiment was repeated for 4 different FPRs (0.1, 0.01, 0.001, 0.0001). All tested structures were constructed to have approximately the same memory space except for buffered MQF which used only one-third of this memory for buffering while the full-size filter is on the disk. MQF is guaranteed to hold the same number of items as a CQF having the same number of slots. The number of slots in CQF was chosen so that the load factor was more than 85% and the MQFs were created with an equal number of slots. Items were sampled to be counted from the real and Zipfian-simulated datasets. After finishing the insertion, to assess the query rate, 5 M items from the same distribution as the insertion datasets were queried. Half of the query items did not exist in the insertion datasets.

MQF has slightly yet persistent faster insertion and query rates compared to CQF with minimal, if any, effect of the FPR on either structure. The performance of CMS is better with higher FPR and Khmer's implementation of CMS doubles the query rate of the original one. However, MQF is always faster than both CMS unless the FPR is more than 0.01 (Fig. 3, Additional file 1: Figures S4, S5).

MQF outperforms CMS in real-world problems

Khmer is a software package deploying a new implementation of CMS for k-mer counting, error trimming and digital normalization [16]. To test MQF in real-life applications, we assessed the performance of the Khmer software package using CMS versus our new implementation using MQF (<https://github.com/dib-lab/khmer/tree/MQFIntegration2>

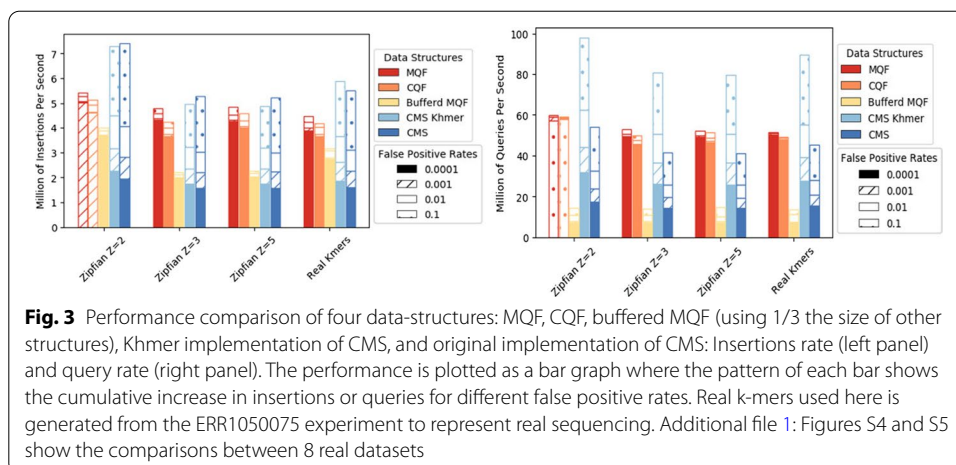


Table 1 Khmer performance in error trimming and digital normalization using MQF and CMS

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	11,011	445,817 ^a	39	37	3253	31,143 ^a	13,11	3
10 ⁻²	2.6	43	48	1304	404,354	41	45	416	24,987	14	5
10 ⁻³	3.4	44	61	130	311,464	42	54	58	21,000	15	7
10 ⁻⁴	4.5	44	75	3	292,746	42	68	4	18,449	16	10
Exact	5	45	-	0	-	43	-	0	-	-	-

^a Percentages of wrong decisions made by CMS at FPR=0.1 in error trimming and digital normalization are 0.8% and 0.13% of the total number of decisions versus 0.02% and 0.01% made by MQF

). A real RNA seq dataset with 51 million reads from the Genome in a Bottle project [17] was used for error trimming and digital normalization; two real-world applications that involve both k-mer insertions and queries. An exact MQF was used to create a benchmark for the approximate data structures. It took 5 Gb RAM to create the data structure and 45 and 43 min to perform trimming and digital normalization respectively. The optimal memory for MQF and the optimal number of hash functions for CMS were calculated to achieve the specified false-positive rates. The CMS was constructed with the same size as the corresponding MQFs. The CMS and MQF versions of Khmer were compared regarding the speed and accuracy (Table 1).

MQF is faster than MPHF

MPHF is constructed by default to fit the input k-mers while MQF would have different load factor that might affect its performance. To address this question, four growing subsets of real k-mers were inserted into MQFs of size 255 MB to achieve 60%, 70%, 80%, and 90% load factors. For labeling, the order of the 1st k-mer in each block was stored in the external labeling space (See the methods section). MPHFs were constructed with sizes ranging from 15 to 22 MB to fit the four datasets. All data structures were queried

with 35 M existing k-mers and the query times were reported. The MQFs were ~ 10 folds faster than the MPHFs. The query time of the MQF was invariable over the different load factors (Additional file 1: Figure S6).

Discussion

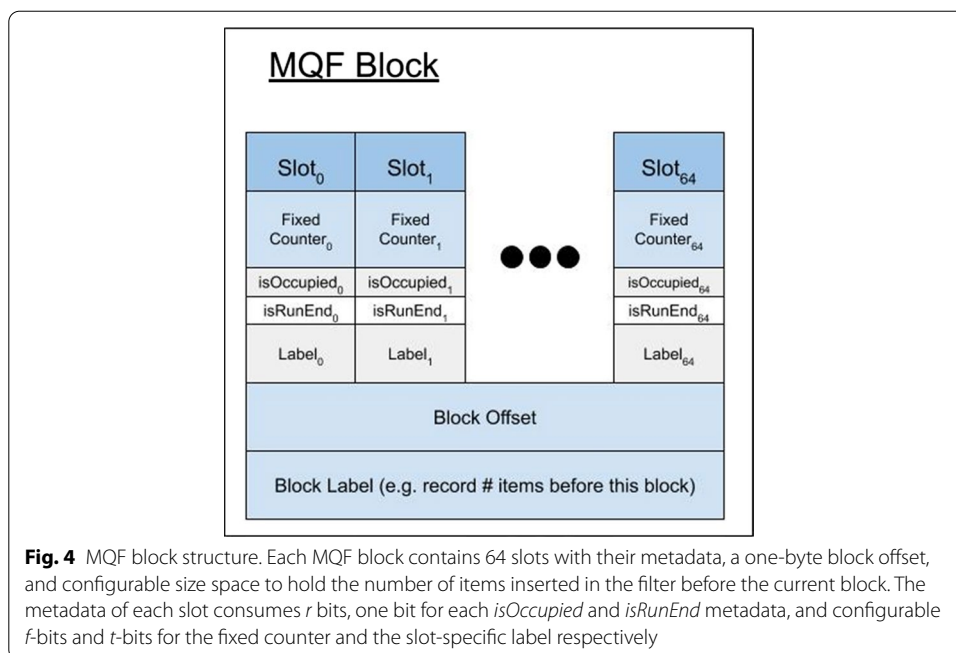
MQF is a new variant of counting quotient filters with novel counting and labeling systems. The new counting system increases memory efficiency as well as the speed of insertions and queries for a wide range of data distributions. The labeling system provides a flexible framework for labeling the member items while maintaining good data locality and a concise memory representation.

MQF is built on the foundation of CQF. MQF has the same ability to behave as an exact or approximate membership query data structure while tracking the count of its members. The insertion/query algorithm developed for CQF enables this family of compact hashtables to perform fast under high load factor (up to 95%) [4]. CQFs are designed to work best for data from high Zipfian distributions. However, previous k-mer spectral analysis of RNAseq datasets showed substantial deviations from a Zipfian distribution in thousands of samples [6]. Such variations in distribution are expected given the variety of biosamples, the broad spectrum of sequencing techniques, and different approaches to data preprocessing.

MQF implements a new counting system that allows the data structure to work efficiently with a broader range of data distributions. The counting system adopts a simple encoding scheme that uses a fixed small space alone or with a variable number of the filter's slots to record the count of member items (Fig. 4). Items with small counts utilize the small fixed-size counters. Therefore, slots, used to be consumed by CQF as counters for these items, are freed to accommodate more items in the filter. The MQF's load factor grows slower than CQF with all distributions except the extreme Zipfian case ($Z=5$) where the load factor is almost the same (Fig. 1). This is why the memory requirement for MQFs is usually smaller compared to CQFs under most distributions despite the extra space taken by the fixed counters (Fig. 2). The size of the fixed-size counter is constant independent of the slot size, therefore the memory requirement for this counter will be trivial with big slots for smaller FPRs and almost negligible in the exact mode. However, this fixed-size counter comes with an additional advantage for MQF. Tuning the size of the fixed-size counter enables the filter to accommodate more items with a slightly larger slot size. This allows the memory requirement for MQF to grow gradually instead of the obligatory size doubling seen in CQF (Fig. 2 and Additional file 1: Figures S2, S3).

Moreover, the new counting scheme in MQF is simplified compared to that of the CQF. MQF defines the required memory for any item based solely on its count (Additional file 1: Figures S7, S8). Therefore, an accurate estimation of the required memory for any dataset can be done extremely quickly by an approximate estimation of data distribution [18, 19]. This is unlike CQF which needs to add a safety margin to account for the special slots used by the counter encoding technique since it is impossible to estimate the number of these slots.

Regarding the speed of insertions and queries, MQF is slightly faster than CQF (Fig. 3). This could be explained partially by the lower load factor of MQF and



partially by the simplicity of the coding/decoding scheme of its counting system. Both MQF and CQF are faster than CMS unless the target FPR is really high (e.g. $FPR > 0.1$) (Fig. 3). CMS controls its FPR by increasing the number of its hash tables requiring more time for insertions and queries to happen. In comparison, quotient filters use always one function but with more hash-bits to control the FPR, with a minimal effect on the insertion/query performance (Fig. 3). With high FPR (e.g. $FPR = 0.1$), CMS uses fewer hash functions and is better performing than MQF. A quotient filter or CMS with a $FPR = \delta$ should have the same probability of item collisions. However, the quotient filter will be more accurate because CMS has another type of error with a probability $(1 - \delta)$, which incorrectly increases the count of its items. This error is a “bounded error” with a threshold that inversely correlates with the width of the CMS [7]. In another sense, some applications might deploy CMS with a smaller table’s width to be more memory efficient than MQF if the application can tolerate a high bounded error.

Buffered MQF can trade some of the speed of insertions and queries for significant memory reduction by storing data on disk. The buffered structure was developed to make use of the optimized sequential read and write on SSD. The buffered structure processes most of the insertion operations using the bufferMQF that resides in memory, thereby limiting the number of access requests to the MQF stored on the SSD hard drive. Sequential disk access happens when the bufferMQF needs to be merged to the disk. This approach is very efficient for insertions but not for random queries which require more frequent SSD data access. Therefore, bufferMQF is best used for the applications where insertions represent the performance bottleneck e.g. k-mer counting applications [5] and sequencing dataset indexing [20]. Moreover, in k-mer analysis of huge raw datasets, buffered MQF can be used initially to filter out the low abundant k-mers (i.e. likely erroneous k-mers), then an in-memory MQF holding the filtered list of k-mers could be used for subsequent application requiring frequent

random queries. This allows multistage analyses where a first pass eliminates likely errors to minimize the memory requirements of computationally demanding applications like in the case of widely used graph-based algorithms [21, 22].

CMS is commonly used for online or streaming applications as long as their high error rate can be tolerated [23]. MQF has a better memory footprint in the approximate mode for lower error rates and thus can compute with CMS for online applications. A major advantage of quotient filters compared to CMS is the dynamic resizing ability in response to the growing input dataset [4]. The buffered version of MQF can be very useful when the required memory is still bigger than the available RAM. We should, however, note that online applications on MQF cannot make use of the memory optimization that could be achieved with an initial estimation of the filter parameters. A new version of the Khmer software that replaces CMS with MQF proves that the new data structure is more efficient in real-life applications. The MQF version is faster than the one with CMS unless the target FPR is high. Also, MQF is always more accurate than CMS although both structures have the same FPR. This behavior of CMS is due to the high error bound of its counts.

Unlike CQF, MQF is designed to be a more comprehensive associative data structure. MQF comes with a novel labeling system that supports associating each k-mer with multiple values to avoid redundant duplication of k-mers' keys in separate data structures. There are two types of labels: Internal labels adjacent to each item to achieve the best cache locality by storing labels next to the k-mer. However, it has a fixed size and thus practically useful when a small size label is needed. The second labeling system is to label the k-mers with one or more labels stored in external arrays while using the k-mer order in the MQF as an index. External labeling is very memory efficient mimicking the idea of the minimal perfect hash function (MPHF) [14, 15]. MPHF undoubtedly has the lowest memory requirement of all the associative data structures [15]. However, MQF has better performance in both the construction and query phases. For construction, both structures require initial k-mer counting. MQF needs just an extra $O(N)$ operation to update the block labels where N is the number of its unique k-mers. MPHF has to read then rehash the list of unique k-mers possibly more than once which makes it slower than MQF. For query operations, MQF is $10 \times$ faster regardless of the load factor of MQF (Additional file 1: Figure S6).

Furthermore, MQF offers more functionality and has fewer limitations than MPHF. MQF is capable of labeling a subset of its items which saves significant space for many applications. For example, k-mer analysis applications may want to only label the frequent k-mers, as an intermediate solution between pruning all the infrequent k-mers and labeling all the k-mers. Moreover, MQF allows online insertions and deletions of items as well as merging of multiple labeled MQFs (See the methods) while MPHF—which doesn't store the items—needs to be rebuilt over the whole dataset, which requires reading and rehashing the datasets. Furthermore, MQF can be exact, while MPHF has false positives when queried with novel items that don't belong to the indexed dataset.

Conclusions

MQF is a new counting quotient filter with a simplified encoding scheme and an efficient labeling system. MQF adapts well to a wide range of k-mer datasets to be more memory and time-efficient than its predecessor in many situations. A buffered version of MQF has a fast insertion algorithm while storing most of the structure on external memory.

MQF combines a fast access labeling system with MPHf-like associative functionality. MQF performance, features, and extensibility make it a good fit for many online algorithms of sequence analysis.

Methods

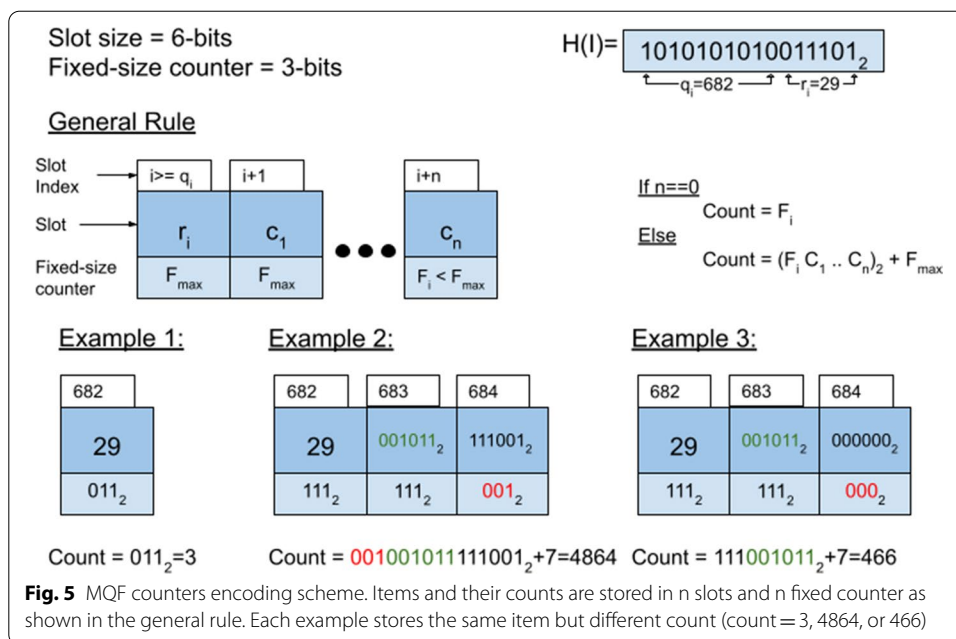
MQF data structure

MQF has a similar structure to CQF with a different scheme of metadata that enables different counting and labeling systems (Fig. 4). Like the CQF, the MQF requires 2 parameters, r and q , and creates an array of 2^q slots; each slot has r -bits. In MQF, Q_i is the slot at position i where $i = 1 \dots 2^q$. The MQF maintains the block design of CQF where each block has 64 slots with their metadata and one extra byte of metadata called *Offset* to enhance the query of items [4]. Both MQF and CQF have two metadata bits to accompany each slot: *isRunEnd_i* and *isOccupied_i*. In the MQF, each slot i has extra metadata, a fixed-size counter with a value (F_i) and a configurable size (F_{size}). There are also two optional fixed-size parts of metadata allocated to allow different styles of labeling. Every slot has specific labeling (ST_i) with a configurable size ($ST_{size} \geq 0$), and every block (j) has an optional space of a configurable size designed to store the number of items in the previous blocks.

The MQF uses the same insertion/query algorithm of CQF [4]. In brief, suppose item I , repeated c times, is to be inserted into Q . A hash function H is applied to I to generate a p -bit fingerprint ($H(I)$). $H(I)$ value is split into two parts, a quotient and remainder. The quotient (q_i) is the most significant q bits while the remainder (r_i) is the remaining least significant r bits. The filters store r_i in a slot Q_j where j is determined by linear probing starting from q_i . One or more slots can be used to store the count of the same item. If the required slots for the item or its count are not free, all the consecutive occupied slots starting from this position will be shifted to free the required space. All items having the same q are stored into consecutive slots and are called a run. Items in the run are 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 and only if there is a run for q_i . Therefore, there is one bit set to one in each *isOccupied* and *isRunEnd* for each run. To query item I , a Rank and Select method is applied on the metadata arrays to get the run start and end for q_i . Then all the items in the run are searched linearly for the slot containing r_i . The subsequent one or more slots can be decoded to get the count of item I . CQF uses a special encoding scheme to recognize these counting slots but MQF utilizes the fixed-counter metadata element (see below).

Counting scheme

MQF uses variable-length integers to encode the k -mers count as shown in Fig. 5. It uses as many slots as required to encode the count while using the fixed-size space of the slots to mark the last slot. To do so, each fixed-size space in all slots before last is assigned its max value (F_{max}), while the fixed-size space of the last slot stores the most significant bits of the count (F_i) where ($F_i < F_{max}$). If the value of the most significant bits is equal to F_{max} , an extra slot with zero bits will be added (see example 3 in Fig. 5). For an item with count c , the fixed-size space of the item's slot is enough for counting until $c \geq F_{max}$ where the number of required slots for counting can be calculated as



$CEILING\left(\frac{abs(log_2(c-F_{max})-F_{size})}{r}\right) + FLOOR\left(\frac{c_s}{F_{max}}\right)$ where c_s represents the decimal value of the most significant bits. In comparison to the CQF, the MQF does not use special slots to resolve ambiguities, which is more memory efficient (Additional file 1: Figures S7 and S8). The counter encoding algorithm is described in Additional file 1: Figure S9.

Parameter estimation

For offline counting applications, the MQF parameters (q, r, F_{size}) can be even more optimized for each dataset to create the most memory-efficient filter that has enough slots to fit all unique items and their counts. The q parameter defines the number of slots (N) in MQF where $q = log_2(N)$. The required numbers of slots for items and their count can be estimated from the cardinality of the target dataset, as with CQF. The r parameter is calculated from the equation $r = p - q$ where p is the total number of hash-bits used to represent each item. In the exact mode, p equals the exact output of a reversible hash function. In the inexact mode, p is controlled by the target FPR (δ) according to the equation $p = log_2\left(\frac{N}{\delta}\right)$ described before [4]. The F_{size} parameter defines the size of the fixed-size counter. This is critical because if a given MQF has too few slots for items in a dataset, the bigger MQF would have to double the number of slots causing a big jump in the memory requirement. To avoid that jump, MQF can use larger fixed size counters to decrease the number of slots required in counting on the expense of a slight increase in the slot size. MQF comes with a utility file (<https://github.com/dib-lab/MQF/blob/master/include/utils.h>) to enable the calculation of the optimum parameters.

Labeling system

MQF can map each item to its count as well as other values, which we call “labels”. Labels in MQF have two different systems. An internal labeling system stores the associated

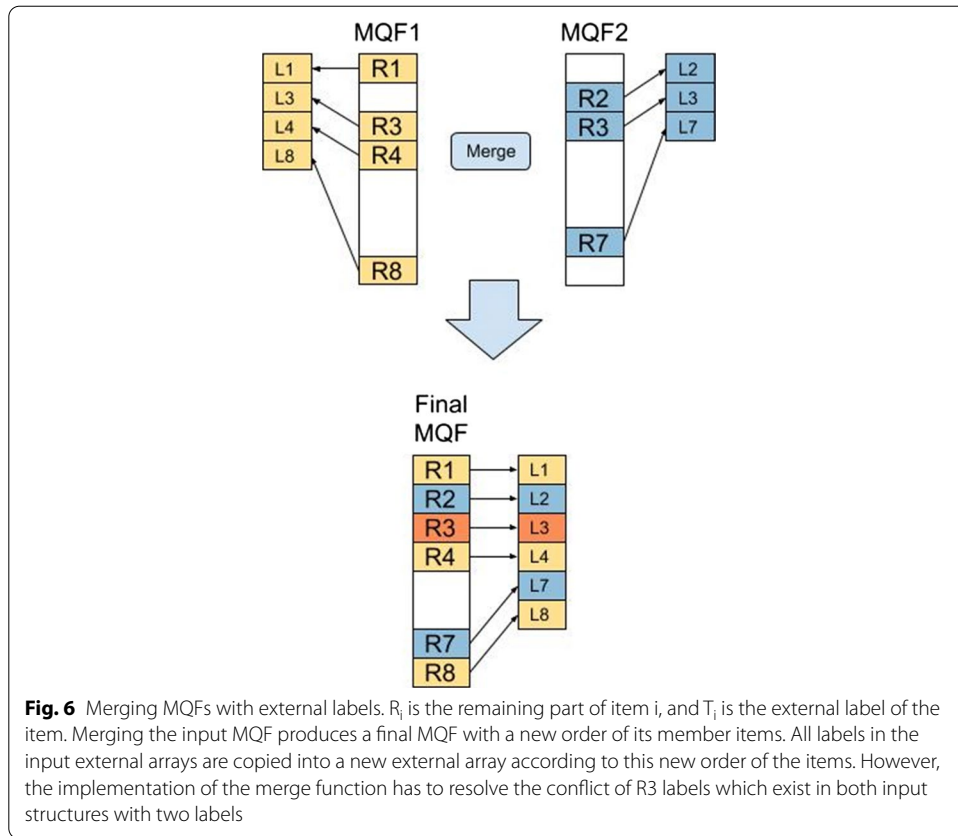
value for every key in the data structure, like a hash table. This label has a fixed size defined at the initialization of the MQF and is practically useful when a small size label is needed (e.g. one or two bits). The second labeling system labels the block. We use this label to store the number of items inserted in the MQF before each block. This enables labeling the items of the filter by separate arrays matching the order of the items in the filter, a behavior that can act as a minimal perfect hash function [15]. The naive way to compute the items' order is to find the item in the MQF and iterate backward until the beginning of the filter to count the number of the preceding items, which is an $O(N)$ operation. The MQF stores the number of items that exist before each block; therefore, the MQF iterates only to the beginning of each block, which is an $O(1)$ operation. The number of previous items for each block is computed after the MQF is constructed. Any additional insertions or deletions of items would only require re-calculation of the block label values with no need to re-analyze the original data. Moreover, labeled MQFs can be updated by merging multiple labeled MQFs and their external labeling arrays. External label arrays need to be merged after merging the labeled MQFs. To do so, the new items' order is recomputed in the final MQF. Then, labels in the input external arrays can be copied into a new external array according to the new item order. Such a function has to consider resolving the conflicts of items happening in multiple-input MQF and labeled by different external labels (Fig. 6).

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 [24] 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 Fig. 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.

Experimental setup of benchmarking

A total of 13 datasets (5 simulated and 8 real sequence data) were used in the experiments to cover a broad spectrum of fragment selection approaches and sequencing platforms. Three datasets called z2, z3, and z5 were simulated to follow Zipfian distribution using three different coefficients: 2, 3, and 5 respectively. The bigger the coefficient the more singletons in the dataset [25]. A fourth dataset was simulated from a uniform distribution with a frequency equal to 10. Another dataset was simulated to mimic single-end 100 bp DNA Illumina sequencing from human chromosome 20 using ART simulator [26]. The remaining 8 datasets represent k-mers sampled from real sequencing datasets (ERR1050075, SRR11551346, SRR12801265, SRR12924365, SRR12937177, ERR992657, SRR12873993, SRR12989394). Throughout the manuscript, each dataset is



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  $i \in \text{bufferMQF}$  do
5:       mqf_insert(bufferMQF, i)
6:     end for
7:     mqf_clear(bufferMQF)
8:   end if
9: end procedure

```

Fig. 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

referred to by its accession number. Additional file 1: Table S3 enlists all the datasets with their description. A k-mer size of 25 nucleotides was used in all experiments.

Experiments were conducted to compare the performance, memory, and accuracy of MQF with the state-of-the-art counting structures CQE, CMS, and MPHE. Unless stated otherwise, CQE and MQF used the same number of slots, and the same slot size while the fixed counter of MQF was set to two. The slot size was calculated to achieve the target FPR as described in the parameter estimation section (see Methods). To create comparable CMS, the number of the tables in the sketches was calculated using $\ln \frac{1}{\delta}$ as described before [7]. The table width was calculated by dividing the

Algorithm 3 Buffered MQF Query

```

1: procedure QUERY(onDiskMQF, bufferMQF, list_item)
2:   for all  $i \in listItems$  do
3:     mqf_insert(tmpMQF,  $i$ )
4:   end for
5:   for all  $i \in tmpMQF$  do
6:      $counts[i] \leftarrow mqf\_query(onDiskMQF, i)$ 
7:      $counts[i] \leftarrow counts[i] + mqf\_query(bufferMQF, i)$ 
8:   end for
9:   return counts
10: end procedure

```

Fig. 8 Buffered MQF query algorithm. Query algorithm for retrieving counts for a list of items in the Buffered MQF. First, insert all the items in the list into a temporary MQF. Second, iterate over the list of items in the temporary MQF and query both the in-memory and on-disk structures

MQF size by the number of tables. The MPHf was created using the default options in the BBhash repo (<https://github.com/rizkg/BBHash>). An Amazon AWS t3.large machine with Ubuntu Server 18.04 was used to run all the experiments. The instance had 2 VCPUS and 8 GB RAM with a 100 GB provisioned IOPS SSD attached for storage. All codes used in the experiments can be accessed through the MQF GitHub repository (<https://github.com/dib-lab/2020-paper-mqf-benchmarks>).

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12859-021-03996-x>.

Additional file 1. Supplementary Figures, Supplementary Table 1 and Supplementary Table 3.

Additional file 2. Supplementary Table 2.

Abbreviations

MQF: Mixed-counters quotient filter; CQF: Counting quotient filter; FPR: False positive rate; CMS: Count-min sketch; MPHf: Minimal perfect hash functions.

Acknowledgements

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. All authors read and approved the final manuscript.

Funding

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>.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

¹ Department of Computer Science, University of California, Davis, CA, USA. ² Department of Population Health and Reproduction, School of Veterinary Medicine, University of California, Davis, CA, USA. ³ Department of Clinical Pathology, School of Medicine, University of Mansoura, Mansoura, Egypt.

Received: 15 September 2020 Accepted: 4 February 2021

Published online: 16 February 2021

References

1. Kolajo T, Daramola O, Adebiji A. Big data stream analysis: a systematic literature review. *J Big Data*. 2019;6:47.
2. Matias Y, Vitter JS, Young NE. Approximate data structures with applications. In: Proceedings of the fifth annual ACM-SIAM symposium on discrete algorithms. Arlington: Society for Industrial and Applied Mathematics; 1994. p. 187–194.
3. Bloom BH. Space/time trade-offs in hash coding with allowable errors. *Commun ACM*. 1970;13:422–6.
4. Pandey P, Bender MA, Johnson R, Patro R. A general-purpose counting filter: making every bit count. In: Proceedings of the 2017 ACM international conference on management of data. Chicago: Association for Computing Machinery; 2017. p. 775–87.
5. Manekar SC, Sathe SR. A benchmark study of k-mer counting methods for high-throughput sequencing. *Gigascience*. 2018;7:gij125.
6. Yu Y, et al. SeqOthello: querying RNA-seq experiments at scale. *Genome Biol*. 2018;19:167.
7. Cormode G, Muthukrishnan S. An improved data stream summary: the count-min sketch and its applications. *J Algorithms*. 2005;55:58–75.
8. Muggli MD, et al. Succinct colored de Bruijn graphs. *Bioinformatics*. 2017;33:3181–7.
9. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq quantification. *Nat Biotechnol*. 2016;34:525–7.
10. Iqbal Z, Caccamo M, Turner I, Flicek P, McVean G. De novo assembly and genotyping of variants using colored de Bruijn graphs. *Nat Genet*. 2012;44:226–32.
11. Cormen TH, Leiserson CE, Rivest RL, Stein CS. Introduction to algorithms. Cambridge: MIT Press; 2009.
12. Marçais G, Kingsford C. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics*. 2011;27:764–70.
13. Wang J, Chen S, Dong L, Wang G. CHTKC: a robust and efficient k-mer counting algorithm based on a lock-free chaining hash table. *Brief Bioinform*. 2020. <https://doi.org/10.1093/bib/bbaa063>.
14. Belazzougui D, Botelho FC, Dietzfelbinger M. Hash, displace, and compress. In: Fiat A, Sanders P, editors. Algorithms—ESA 2009. Berlin: Springer; 2009. p. 682–93.
15. Limasset A, Rizk G, Chikhi R, Peterlongo P. Fast and scalable minimal perfect hashing for massive key sets. In: 16th International Symposium on Experimental Algorithms. Vol. 11. London, United Kingdom; 2017; p. 1–11.
16. Crusoe MR, et al. The khmer software package: enabling efficient nucleotide sequence analysis. *F1000 Res*. 2015;4:900.
17. Zook JM, Salit M. Genomes in a bottle: creating standard reference materials for genomic variation—why, what and how? *Genome Biol*. 2011;12:P31–P31.
18. Flajolet P, Fusy R, Gandouet O, Meunier F. HyperLogLog: the analysis of a near-optimal cardinality estimation algorithm. In: Discrete mathematics & theoretical computer science; 2007. p. 137–56.
19. Mohamadi H, Khan H, Birol I. ntCard: a streaming algorithm for cardinality estimation in genomics data. *Bioinformatics*. 2017;33:1324–30.
20. Pandey P, et al. Mantis: a fast, small, and exact large-scale sequence-search index. *Cell Syst*. 2018;7:201–207 e4.
21. Schlick T. Adventures with RNA graphs. *Methods*. 2018;143:16–33.
22. Sohn JI, Nam JW. The present and future of de novo whole-genome assembly. *Brief Bioinform*. 2018;19:23–40.
23. Muthukrishnan S. Data streams: algorithms and applications. *Found Trends Theor Comput Sci*. 2003;1:117–236.
24. Dementiev R, Kettner L, Sanders P. STXXL: standard template library for XXL data sets. *Softw Prac Exp*. 2008;38:589–637.
25. Powers D. Applications and explanations of Zipf's Law. In: CoNLL; 1998.
26. Huang W, Li L, Myers JR, Marth GT. ART: a next-generation sequencing read simulator. *Bioinformatics*. 2012;28:593–4.

Publisher's Note

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

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

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

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

