# Bloom filter variants for multiple sets: a comparative assessment

Luca Calderoni[a,*], Dario Maio[a], Paolo Palmieri[b]

[a]*Department of Computer Science and Engineering, University of Bologna, Italy*
[b]*Department of Computer Science, University College Cork, Ireland*

## Abstract

In this paper we compare two probabilistic data structures for association queries derived from the well-known Bloom filter: the shifting Bloom filter (ShBF), and the spatial Bloom filter (SBF). With respect to the original data structure, both variants add the ability to store multiple subsets in the same filter, using different strategies. We analyse the performance of the two data structures with respect to false positive probability, and the inter-set error probability (the probability for an element in the set of being recognised as belonging to the wrong subset). As part of our analysis, we extended the functionality of the shifting Bloom filter, optimising the filter for any non-trivial number of subsets. We propose a new generalised ShBF definition with applications outside of our specific domain, and present new probability formulas. Results of the comparison show that the ShBF provides better space efficiency, but at a significantly higher computational cost than the SBF.

*Keywords:* probabilistic data structures, spatial bloom filter, shifting bloom filter, association queries

## 1. Introduction

The term *big data* informally refers to data sets that are so large and complex that cannot be processed efficiently using traditional data structures and algorithms. The recent introduction of inexpensive information-gathering devices has fuelled the growth of pervasive sensing, which led to an increase in the size and number of data sets to be computed. Where traditional, deterministic data structures are inadequate to process big data, probabilistic data structures have been widely adopted by computer scientists as a suitable alternative [1]. While standard data structures are designed to store elements and answer deterministically to queries, probabilistic data structures introduce an error probability. This drawback, however, is balanced by higher space efficiency and lower computational burden, which are crucial in big data applications [2].

Most probabilistic data structures, and specifically those relying on Bloom filters, use hash functions to randomize and compactly represent a set of items. The possibility of collisions introduces the potential for errors in the stored information, however the error probability is generally known and can be maintained under an arbitrary threshold by tuning the data structure parameters. Compared to errorless approaches, probabilistic data structures use significantly less memory, and have constant query time. They also usually support union and intersection operations and can therefore be easily parallelised.

Bloom filters (BF) are arguably the most prominent data structure belonging to this category. They were first introduced by Bloom [3] in 1970, and were extensively used during the last decades. Bloom filters were designed to support membership queries, i.e., to determine whether a given element is a member of a given set or not. The Bloom filter always determines positively if an element is in the set, while elements outside the set are generally determined negatively, but are subject to a (bounded) false positive error probability [4].

Several modifications of Bloom filters have been proposed over the years [5], with the purpose of adding additional functionalities. In this paper we compare two of these variants: the shifting Bloom filter (ShBF), and the spatial Bloom filter (SBF). Both variants allow the filter to store multiple sets[1] rather than a single one, and enable membership queries over these sets (called *association queries*). ShBF and SBF use different strategies to achieve this: the ShBF uses additional hash functions to determine a shift in the positions within the filter for each set, while the SBF writes a different index value for each set.

### 1.1. Related Works

Bloom filters and the derived data structures have been used in a number of database systems applications over the years [6]. Due to their ability to efficiently know whether or not a key belongs to a set of keys, they are widely adopted for first level indexes [2] and for query caching [7]. They were extensively tested over random access memories, traditional hard drives and, more recently on solid state drives. Bloom filters were also applied to semi-structured data queries [8], range queries and so forth.

---

*Corresponding author

*Email addresses:* `luca.calderoni@unibo.it` (Luca Calderoni), `dario.maio@unibo.it` (Dario Maio), `p.palmieri@cs.ucc.ie` (Paolo Palmieri)

[1]We note here that the relevant literature often mixes the definition of multi-set (as in multiple sets) and multiset. In this paper, we use the term multiset following the mathematical definition of a set which permits the coexistence of multiple instances of any element. In order to remove the ambiguity with multi-set, we only use the full expression "multiple sets" for datasets which include elements of different sets.

Concerning association queries, following the classification proposed in [5], a number of Bloom filter variants proposed in the literature allow for multiple sets to be stored in the same filter. However, most data structures implement this functionality using a naive strategy: that is, using multiple instances of the data structure, one for each set.

Many such structures were proposed in the context of network routing strategies, where Bloom filters are widely used to address packet filtering and forwarding, and the application scenario often requires the use of multiple sets [9, 10]. For instance, the *longest prefix match* solution proposed in [11] is one of the first strategies using the naive approach of one separate data structure for each set. In general, multiple BF strategies are viable when there is some a-priori knowledge regarding the number of sets to be mapped (and potentially the approximate number of items per set). This condition is often required for the design of hardware routing boards, where Bloom filters are parallelised and may thus be queried simultaneously. Several proposed solutions follow this principle: the IPv6 lookup scheme discussed in [12] or the packet forwarding strategy presented in [13] are widely cited examples. The difference Bloom filter (DBF) is similar in this regard as it requires additional data structures linked to the main Bloom filter [14]. However, as the number of sets increases this approach becomes intuitively inefficient, as it increases the overhead linearly to the number of sets. It is also unsuitable when the number of sets cannot be predicted, or changes over time. We therefore exclude such Bloom filter variants from the comparison we present in this paper, to focus instead on structures that allow for multiple sets by design. A number of variants of the Bloom filter allowing multiple sets exist in the literature: shifting Bloom filters [15], combinatorial Bloom filters [16, 17], spatial Bloom filters [18, 19], the kBF [20], and the Bloomier filter [21, 22]. However, only a few are designed specifically for storing multiple sets. The kBF is instead a Bloom filter designed for key-value storage, with applications on approximate state machines; while Bloomier filters are a data structure for static support lookup tables. While both data structures can implement multiple sets, we do not include them in the comparison for reasons explained below. In the case of the kBF, the focus of the data structure is to store a value associated to each element of the originating set (following the key-value paradigm, where the key is the element and the value is the information stored in the data structure). The value is stored in a 32-bits space, divided into 29 bits used for encoding and a 3-bit counter. While we can use the value for storing a set identifier (and therefore group the elements into subsets) this is not the core functionality of the data structure, and, as it implies writing a 32-bit value for each element, is relatively inefficient. Similarly, the core functionality of Bloomier filters is not association queries on multiple sets: instead, Bloomier filters are designed for computing arbitrary functions defined only on the originating set. It is technically possible to implement a Bloomier filter defined on a membership function returning different values for elements of different subsets (and indeed this is used as an example by the authors). However, Bloomier filters can be implemented on arbitrary functions, which is a more complex problem and therefore

introduce a significant overhead. In particular, Bloomier filters allow dynamic updates to the function, and also use two separate tables within the data structure, increasing memory usage. Filters that are designed for the purpose of storing multiple sets normally adopt either of two different strategies: making the filter non-binary (to allow for the set information to be stored in a filter position) as in the case of the spatial Bloom filter; or using an increased number of hash functions, where different hashes are used for different sets. The latter strategy is used by both the shifting Bloom filters and the combinatorial Bloom filters. For the reasons detailed above, we restrict the comparative assessment presented in this paper to the data structures that are designed specifically for membership queries over multiple sets, commonly referred to as association queries. As our objective is to compare the two main strategies to do so (detailed above), we select the spatial Bloom filter to represent the non-binary filter approach, and the shifting Bloom filter, which is the most recent structure following the multiple hashes approach.

*1.2. Contribution*

In this paper we present a comparative assessment of two probabilistic data structures for association queries, the shifting Bloom filter and the spatial Bloom filter. In the case of the shifting Bloom filter, we also provide a new, generalised definition that allows an unlimited number of sets to be stored (contrary to the original definition that focuses on two sets only), and we provide updated error probability formulas to reflect this change. The comparative assessment is based on an experimental analysis of the behaviour of the two data structures over large test datasets. Extensive tests have been performed using comparable implementations of both primitives. The empirical results validate the theoretical analysis of the error probabilities, space efficiency and computational costs.

## 2. Shifting Bloom Filter

The shifting Bloom filter [15] is a variant of Bloom filters which supports membership queries as well as association queries and multiplicity queries. However, this data structure cannot provide all the aforementioned features at once. The user needs to select which kind of queries the filter should answer to, and construct it according to procedures that are specific to each functionality. The interrogation procedures are also feature-specific. As the comparison proposed in this work focuses on association queries, we limit our discussion to ShBF filters set up for this specific type of queries. In the following we refer to the ShBF built for association queries simply as ShBF. The ShBF as proposed in the original work by Yang et al. [15] focuses only on two sets, and provides a mechanism for identifying elements that are present in both sets (therefore the two sets are not disjoint, and their intersection can be $\neq \emptyset$). However, the limitation of being able to store only two sets is a direct consequence of allowing the two originating sets to be non-disjoint, and requires the use of an auxiliary data structure thus increasing memory usage. We believe that this limitation significantly reduces the potential application scenarios for

the ShBF, and therefore in this paper we extend the ShBF data structure to allow for an unlimited number of originating sets. However, in doing so we remove the possibility to store sets with intersecting elements, and we allow for disjoint sets only. We feel the benefits of our proposed modification far outweigh the downsides, and we note that intersections can potentially be re-introduced using a naive strategy (defining an intersection as a separate set). This modification also allows for a more meaningful comparison of the ShBF with the SBF data structure, with can also store an unlimited number of disjoint sets.

The shifting Bloom filter allows multiple sets to be stored by shifting the position at which the values are stored in the filter: that is, an *offset* is added to the output of each hash function, and the offset specifies which set the element belongs to. The offset is not static, but is determined by an additional hash function, so that each set uses a different hash function to determine the offset. In practice, this means that two elements of the same set will have different offsets, but these offsets are both calculated using the same hash function, which is set-specific. For memory-efficiency reasons, Yang et al. propose in the original paper to limit the offset space (that is, the output size of this additional hash function), so that the position pointed by the hash function plus the offset will be in the same memory word as the position pointed by the hash function only (without the offset). The authors do not specify a particular size for the memory word. While this approach can reduce memory reads when only two sets are stored, in the proposed generalisation to unlimited number of sets a small word size can significantly increase the number of elements for which the set cannot be determined with accuracy.

In Figure 1 we provide the results of an experiment where multiple word lengths are used over the same test dataset. For this experiment, we store 255 sets each containing 256 elements, for a total of 65280 elements in a filter of length $2^{23}$ bits. As the two figures clearly depict, a small word value ($2^{10}$) significantly increases the number of errors returned by the filter (a formal definition of such errors is provided in the following). The best results are obtained for word sizes of $2^{16}$ or above. Considering that the original goal of reducing the memory reads cannot be achieved with such a high word size, we propose here to discard the limitation on the offset space entirely. As shown by the experiment, limiting the offset space to any meaningful word size would negatively affect the correctness of the association queries result. Moreover, removal of the word requirement reduces the computational burden by eliminating a modulo operation for each offset calculation.

In the following, we provide a formal definition of the proposed generalised ShBF, as described above. In particular, the definition we propose differs from the original ShBF as proposed by Yang et al. [15] by introducing an unlimited number of disjoint sets, and removing the offset space limitation.

**Definition 1.** *Given the originating sets* $\Delta_1, \Delta_2, \ldots, \Delta_s$ *to be represented in the filter, let* $\bar{S}$ *be the union set* $\bar{S} = \bigcup_{\Delta_i \in S} \Delta_i$ *and* $S$ *be the set of sets such that* $S = \{\Delta_1, \ldots, \Delta_s\}$. *Let* $H = \{h_1, \ldots, h_k\}$ *be a set of* $k$ *hash functions such that each* $h_i \in H : \{0, 1\}^* \rightarrow \{1, \ldots, m\}$, *that is, each hash function takes binary strings as input and outputs a number uniformly chosen in* $\{1, \ldots, m\}$. *Let also* $H^{\curlyvee} = \left\{h_1^{\curlyvee}, \ldots, h_{s-1}^{\curlyvee}\right\}$ *be a set of* $s - 1$ *hash functions such that each* $h_i^{\curlyvee} \in H^{\curlyvee} : \{0, 1\}^* \rightarrow \{1, \ldots, m\}$, *where* $m$ *is a positive integer.*

*We define a* shifting Bloom filter $B^{\curlyvee}(S)$ *over* $S$ *as the set*

$$B^{\curlyvee}(S) = \bigcup_{\delta \in \bar{S}, h \in H} h(\delta) + o(\delta) \ , \tag{1}$$

*where the offset* $o(\delta)$ *is intended as follows:*

$$o(\delta) = \begin{cases} 0 & \text{if } \delta \in \Delta_1 \\ h_i^{\curlyvee}(\delta) & \text{if } \delta \in \Delta_{i+1}, i \neq 0 \end{cases} \ . \tag{2}$$

A shifting Bloom filter $B^{\curlyvee}(S)$ can be represented as a binary vector $b^{\curlyvee}$ composed of $2m$ bits (or *cells*), where the $i$-th bit

$$b^{\curlyvee}[i] = \begin{cases} 1 & \text{if } i \in B^{\curlyvee}(S) \\ 0 & \text{if } i \notin B^{\curlyvee}(S) \end{cases} \ . \tag{3}$$

In the following, when referring to a shifting Bloom filter, we refer to its vector representation $b^{\curlyvee}$.

The ShBF is built as follows. Initially all bits are set to 0. Then, for each $\Delta_i \in S$, for each element $\delta \in \Delta_i$ and for each $h \in H$ we calculate $h(\delta) + o(\delta) = i$, and set the corresponding $i$-th bit of $b^{\curlyvee}$ to 1. Thus, $2m$ bits are needed in order to store $b^{\curlyvee}$. As discussed in Section 4, we implemented a circular ShBF version composed of $m$ bits.

The verification procedure is formalized as follows. Let us suppose to test an element $\delta \in \mathcal{E}$ against the filter (where $\mathcal{E}$ represents a generic domain), in order to learn which originating set it belongs to (or whether it does not belong to any originating set). First of all we need to compute $k$ hash digests, i.e., $\forall h \in H$, we compute $h(\delta)$. Then, as each originating set is combined with a specific hash offset, we need to check separately whether the element belongs to each set. We check whether $\delta \in \Delta_i$ if
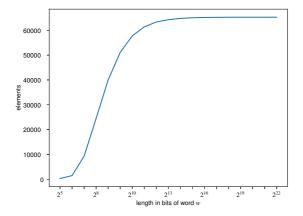
$$\forall h \in H, b^{\curlyvee}[h(\delta) + o(\delta)] = 1 \ . \tag{4}$$

During the verification process of the set $\Delta_i$, should we find an index $i$ such that $b^{\curlyvee}[i] = 0$, we can conclude $\delta \notin \Delta_i$ and we proceed to another set $\Delta_j \in S$. Conversely, when the condition (4) is fulfilled, we add $\Delta_i$ to the set of positives matches $\Gamma$. At the end of the verification procedure three are the possible scenarios:
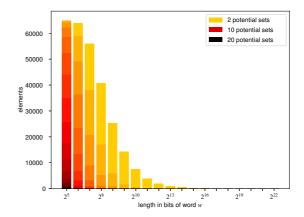
- $\Gamma = \emptyset$. Then we may assert $\delta \notin \bar{S}$.

- $|\Gamma| = 1, \Delta_i \in \Gamma$. Hence $\delta \in \Delta_i$.

- $|\Gamma| > 1$. Hence $\delta$ belongs to one among the sets included in $\Gamma$.

Both the second and the third case are subject to a false positive probability, as ShBF preserves the probabilistic nature of the classic Bloom filter. Concerning the latter case, please also note that it is not possible to assign a different weight to those sets included in $\Gamma$. Consequently, the element $\delta$ has the same probability to belong to each set, resulting in an inter-set error.

A schematization of the insertion and verification procedures of an ShBF is proposed in Figure 2.

3

(a) Correctly recognised elements.



(b) Elements returning multiple potential sets.

Figure 1: The word size.

## 2.1. False positives

As we redefined the ShBF data structure it is important to investigate its behaviour in terms of false positives.

**Definition 2.** *Given a filter $b^\curlyvee$, a false positive event occurs when the verification procedure performed on an element $\delta \notin \bar{S}$ terminates with a non-empty set of positives matches, i.e. $\Gamma \neq \emptyset$.*

As we may see, a false positive event reported by a ShBF may pertain several originating sets, due to the verification procedure which performs a separate lookup for each originating set. Consequently, we define a false positive event on a specific set as follows:

**Definition 3.** *Given a filter $b^\curlyvee$, a false positive event on a specific originating set $\Delta_i$ occurs when the verification procedure performed on an element $\delta \notin \bar{S}$ terminates with the set $\Delta_i$ included in the set of positives matches, i.e. $\Delta_i \in \Gamma$.*

Since during the verification procedure each set is checked for membership through a separate filter lookup performed on the same data structure, each originating set has the same probability of being included in the set of positives matches $\Gamma$. Specifically, the set-specific false positive probability ($\text{FPP}_i^\curlyvee$) in a ShBF coincides with the overall false positives probability of a common BF, assuming both filters were filled with the same amount of items, i.e. $b^\curlyvee$ was filled with $n = |\bar{S}|$ items:

$$\text{FPP}_i^\curlyvee = \left(1 - \left(1 - \frac{1}{m}\right)^{kn}\right)^k . \tag{5}$$

Conversely, the overall false positives probability ($\text{FPP}^\curlyvee$) is increased due to the multiple filter lookups performed during the verification procedure. In order for the verification procedure to produce a false positive, it sufficient that one among the $s$ set-specific lookups produces a set-specific false positive. As the probability for this event to occur is known ($\text{FPP}_i^\curlyvee$), we may address this problem using the well known success-failure scheme (Bernoulli trials). Given an event which occurs with probability $p$, the probability to observe $b$ successes among $a$ trials is:

$$\binom{a}{b} p^b (1-p)^{a-b} . \tag{6}$$

The probability to report at least one set-specific false positive (i.e. at least one success) among $s$ trials is the probability of the certain event minus the probability to never report a set-specific false positive among $s$ trials. Hence

$$\text{FPP}^\curlyvee = 1 - \left( \binom{s}{0} \text{FPP}_i^{\curlyvee 0} (1 - \text{FPP}_i^\curlyvee)^{s-0} \right) ; \tag{7}$$

it follows that:

$$\text{FPP}^\curlyvee = 1 - \left(1 - \left(1 - \left(1 - \frac{1}{m}\right)^{kn}\right)^k\right)^s . \tag{8}$$

## 2.2. Inter-set errors

As discussed above, when the verification procedure produces a set of positives matches composed of more than one set, it is not possible to discern which among those sets is the one that effectively contains the element. Thus, supposing $|\Gamma| = u$, we have only $1/u$ chances to assign the element to the correct set.

**Definition 4.** *Given a filter $b^\curlyvee$, an inter-set error event occurs when the verification procedure performed on an element $\delta \in \bar{S}$ terminates with a set of positives matches such that $|\Gamma| > 1$.*

As per the classic BF, the ShBF is not affected by false negatives. Thus, when we test the filter for an element $\delta \in \Delta_i$, the correct set ($\Delta_i$) is for sure included in $\Gamma$. However, the verification procedure is composed of $s$ ShBF-lookups, one for each set. During the remaining $s - 1$ lookups, should the procedure report a false-positive, the corresponding set would be wrongly added to the set of positives matches. Hence, the probability to
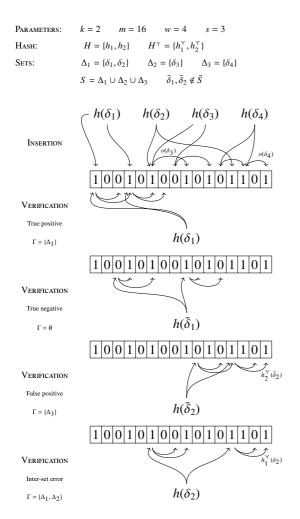
4

Figure 2: Insertion of three originating sets $\Delta_1, \Delta_2$ and $\Delta_3$ and verification of elements in a shifting Bloom filter of length $m = 16$, featuring two hash functions ($k = 2$). With respect to the verification process, it is possible to observe true positives, true negatives, false positives and inter-set errors. Each scenario is depicted in sequence.

observe an inter-set error (ISEP$^\vee$) may be derived as the probability of the certain event minus the probability to never report a false positive throughout $s - 1$ trials.

Following the same principle described in (7) we may derive:

$$\text{ISEP}^\vee = 1 - \left( \binom{s-1}{0} \text{FPP}_i^{\vee \, 0} (1 - \text{FPP}_i^\vee)^{s-1-0} \right) \; ; \tag{9}$$

hence

$$\text{ISEP}^\vee = 1 - \left( 1 - \left( 1 - \left( 1 - \frac{1}{m} \right)^{kn} \right)^k \right)^{s-1} \; . \tag{10}$$

As the amount of correct information produced by the verification procedure depends on the number of sets included in $\Gamma$, it is also useful to derive the probability ISEP$^\vee_{u_i}$ for it to produce a set of positives matches with a specific cardinality $i$. Following the aforementioned principle, it is straightforward that the probability to obtain $|\Gamma| = i$ is the probability to report exactly $i - 1$ false positives among $s - 1$ trials. Hence:

$$\text{ISEP}^\vee_{u_i} = \binom{s-1}{i-1} \text{FPP}_i^{\vee \, i-1} (1 - \text{FPP}_i^\vee)^{s-i} \; . \tag{11}$$

Finally, for ease of comparison with SBF, we discuss the probability for a specific set to be involved in a inter-set error event (ISEP$^\vee_i$).

**Definition 5.** *Given a filter $b^\vee$, an inter-set error event on a specific originating set $\Delta_i$ occurs when the verification procedure performed on an element $\delta \in \bar{S}$ terminates with a set of positives matches such that $\Delta_i \in \Gamma$, $|\Gamma| > 1$.*

Following the same principle discussed in Section 2.1 concerning false positive events on specific sets, it is straightforward to note that ISEP$^\vee_i$ coincides with FPP$^\vee_i$ (see (5) for reference).

## 3. Spatial Bloom Filter

The spatial Bloom filter was originally introduced in [18, 19] with the purpose of efficiently storing an arbitrary number of disjoint sets representing geographic areas. Although its first application was in the location privacy domain (as the name suggests), the data structure can store any type of element and thus represents an efficient solution for any kind of scenarios where association queries over an unlimited number of sets are required [23].

Similarly to the generalised ShBF we propose earlier in this paper, SBF natively supports multiple sets and element mapping is performed through a single data structure, without any need of auxiliary data structures or meta data.

The spatial Bloom filter (SBF) is defined as follows [19]:

**Definition 6.** *Given the originating sets $\Delta_1, \Delta_2, \ldots, \Delta_s$ to be represented in the filter, let $\bar{S}$ be the union set $\bar{S} = \bigcup_{\Delta_i \in S} \Delta_i$ and $S$ be the set of sets such that $S = \{\Delta_1, \ldots, \Delta_s\}$. Let $O$ be the strict total order over $S$ for which $\Delta_i < \Delta_j$ for $i < j$. Let also $H = \{h_1, \ldots, h_k\}$ be a set of $k$ hash functions such that each $h_i \in H : \{0, 1\}^* \rightarrow \{1, \ldots, m\}$, that is, each hash function in $H$ takes binary strings as input and outputs a random number uniformly chosen in $\{1, \ldots, m\}$. We define the* spatial Bloom filter *over $(S, O)$ as the set of couples*

$$B^\# (S, O) = \bigcup_{i \in I} \langle i, \max L_i \rangle \tag{12}$$

*where $I$ is the set of all values output by hash functions in $H$ for elements of $\bar{S}$*

$$I = \bigcup_{\delta \in \bar{S}, h \in H} h(\delta) \tag{13}$$

*and $L_i$ is the set of labels $l$ such that:*

$$L_i = \{l \mid \exists \delta \in \Delta_l, \exists h \in H : h(\delta) = i\} \; . \tag{14}$$

A spatial Bloom filter $B^\# (S, O)$ can be represented as a vector $b^\#$ composed of $m$ values (or *cells*), where the $i$-th value

$$b^\# [i] = \begin{cases} l & \text{if } \langle i, l \rangle \in B^\# (S, O) \\ 0 & \text{if } \langle i, l \rangle \notin B^\# (S, O) \end{cases} \quad \text{[19]}. \tag{15}$$
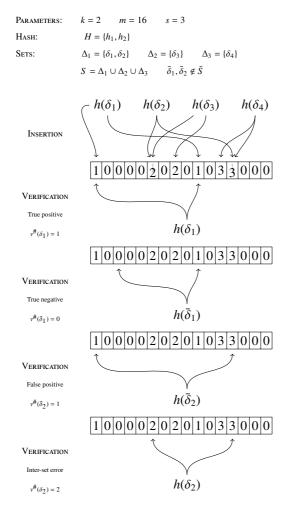
$h(\delta_1)$   $h(\delta_2)$   $h(\delta_3)$   $h(\delta_4)$

INSERTION

| 1 | 0 | 0 | 0 | 0 | 2 | 0 | 2 | 0 | 1 | 0 | 3 | 3 | 0 | 0 | 0 |

VERIFICATION

True positive

$v^{\#}(\delta_1) = 1$

$h(\delta_1)$

| 1 | 0 | 0 | 0 | 0 | 2 | 0 | 2 | 0 | 1 | 0 | 3 | 3 | 0 | 0 | 0 |

VERIFICATION

True negative

$v^{\#}(\bar{\delta}_1) = 0$

$h(\bar{\delta}_1)$

| 1 | 0 | 0 | 0 | 0 | 2 | 0 | 2 | 0 | 1 | 0 | 3 | 3 | 0 | 0 | 0 |

VERIFICATION

False positive

$v^{\#}(\bar{\delta}_2) = 1$

$h(\bar{\delta}_2)$

| 1 | 0 | 0 | 0 | 0 | 2 | 0 | 2 | 0 | 1 | 0 | 3 | 3 | 0 | 0 | 0 |

VERIFICATION

Inter-set error

$v^{\#}(\delta_2) = 2$

$h(\delta_2)$

Figure 3: Insertion of three originating sets $\Delta_1, \Delta_2$ and $\Delta_3$ and verification of elements in a spatial Bloom filter of length $m = 16$, featuring two hash functions ($k = 2$). With respect to the verification process, it is possible to observe true positives, true negatives, false positives and inter-set errors. Each scenario is depicted in sequence.

Throughout this paper, we refer to spatial Bloom filters using the vector representation $b^{\#}$.

The construction procedure of a SBF may be summarized as follows [19]. The first step consists in setting each value in $b^{\#}$ to 0. Then we compute $h(\delta) = i$ for each item $\delta \in \Delta_1$ and for each $h \in H$. The corresponding cell of the filter ($b^{\#}[i]$) is set to 1 (as 1 represents the label of $\Delta_1$). The same procedure is performed for elements included in $\Delta_2$ (writing the value 2), then for those included in $\Delta_3$ and so forth. The construction procedure terminates when the last set $\Delta_s$ has been processed. It is important to point out that, as the construction procedure follows the strict total order defined over $S$, sets with lower labels are more likely to be overwritten than sets with higher label values, in case of collision.

Concerning collisions, we note that they occur when a specific hash function produces the same digest for two different elements (this is indeed what we commonly refer to as a hash collision), but also when two separate hash functions produce the same digest for a single element or for two distinct elements. When one among the aforementioned conditions is met, the

same filter cell is written twice. We refer to this phenomenon as *cell overwrite*, or simply as *collision*, indifferently. An in-depth analysis concerning collisions and how they affect the probabilistic properties of this data structure is beyond the scope of this work and may be found in [24].

During the filter construction, the filter crosses several *states* [24]:

**Definition 7.** *Let us consider the spatial Bloom filter $b^{\#}$. Given $i \in L$, we say the filter is in* state $i$ *(and we refer to its vector representation as $b^{\#}_{\triangleleft i}$) if and only if all the elements of set $\Delta_i$ have been inserted into the filter.*

Therefore, state 0 ($b^{\#}_{\triangleleft 0}$) represents the empty filter (with all of its cells set to zero). At the end of the construction process, the filter is in state $s$ ($b^{\#}_{\triangleleft s}$).

In order to know whether or not an element belongs to one of the originating sets, we need to perform a single filter lookup. Given $\delta \in \mathcal{E}$ (where $\mathcal{E}$ represents a generic domain), we compute $k$ hash digest, i.e., $\forall h \in H$, we compute $h(\delta)$. We check whether $\delta \in \Delta_i$ if

$$\exists h \in H : b^{\#}[h(\delta)] = i \quad \text{and} \quad \forall h \in H, b^{\#}[h(\delta)] \geq i \ . \quad (16)$$

Should one or more cells contain the value 0, we can conclude $\delta \notin \bar{S}$, i.e., it does not belong to any of the originating sets.

The verification procedure which classifies an element $\delta \in \mathcal{E}$ (where $\mathcal{E}$ represents a generic domain) as belonging to one of the originating sets may be formalized using a functional notation as follows:

$$\begin{aligned} v^{\#} : \quad & \mathcal{E} \to L_0 \\ & \delta \mapsto v^{\#}(\delta) \end{aligned} \quad (17)$$

where $L_0 = \{0, \ldots, s\}$. This function processes a generic element of the domain $\mathcal{E}$ and outputs a value included in $\{0, \ldots, s\}$. This integer indicates the set to which the element is supposed to belong to (if the output is $> 0$), or indicates that the element does not belong to any of the originating sets (when the output is 0). As thoroughly investigated in [24], false negatives are not possible, while positive matches are subject to false-positives and inter-set errors. An exemplification of insertion and verification procedures concerning a SBF is proposed in Figure 3.

### 3.1. False positives

While an in-depth investigation concerning SBF false positives falls outside the scope of this work, it is meaningful to recall some outcome for ease of comparison with the ShBF. Specifically, the overall false positive probability in a SBF coincides with the one of a common BF, assuming both filters were filled with the same amount of items, i.e. $|\bar{S}| = n$. In this case, the probability to observe a false positive when we query a SBF for membership of an element which does not belong to any of the originating sets is

$$\text{FPP}^{\#} = \left(1 - \left(1 - \frac{1}{m}\right)^{kn}\right)^k \ . \quad (18)$$

As widely discussed in [24], this probability may be divided in a set-specific one. Specifically, the probability to report a false positive on a given set $\Delta_i$ is:

$$\text{FPP}_i^{\#} = \left(1 - \left(1 - \frac{1}{m}\right)^{k \sum_{j=i}^{s} n_j}\right)^k - \sum_{j=i+1}^{s} \text{FPP}_j^{\#} \ . \qquad (19)$$

### 3.2. Inter-set errors

As stated above, an in-depth investigation concerning SBF inter-set errors falls outside the scope of this work. However, we report the main outcomes discussed in [24] for ease of comparison with the ShBF. The probability to observe a set-specific inter-set error is:

$$\text{ISEP}_i^{\#} = \left(1 - \left(1 - \frac{1}{m}\right)^{k n_i^{\text{FILL}}}\right)^k \ , \qquad (20)$$

where $n_i^{\text{FILL}}$ represents the number of elements left for insertion after the filter $b^{\#}$ reaches the state $i$ (see [24] for details).

Assuming $|\Delta_i| = n_i$, the overall inter-set error probability ($\text{ISEP}^{\#}$) may be derived as the weighted sum of each set-specific probability:

$$\text{ISEP}^{\#} = \frac{1}{n} \sum_{i=1}^{s} n_i \text{ISEP}_i^{\#} \ . \qquad (21)$$

## 4. Comparison

We compare the performances of the shifting Bloom filters and the spatial Bloom filters over the two most meaningful filter characteristics: the false positive probability and the inter-set error rate. Because of the different strategies used by the data structures to store multiple sets, these probabilities need to be evaluated with respect to two different parameters: the filter length in bits $l$ (which expresses the memory usage of the corresponding filter), and the number of cells $m$ (where each cell is a unique position within the filter). In the case of the ShBF, these two parameters coincide and $l = m$, as each cell can store either a 0 or a 1, following the binary nature of the original Bloom filter. It is important to point out that, by definition, ShBF is subject to an offset function which may cause indexed cells to exceed $m$. For ease of comparison, we opted for a circular implementation of the ShBF: during both construction and verification phase, cell indexing is subject to a $m$ modulus, that is, given an element $\delta$ to be inserted or tested against the filter, each hash function is evaluated in combination with the proper offset as $(h(\delta) + o(\delta))\%m$. The SBF allows instead multi-bit cells, where the data structure stores the index of the relevant set. The length of the filter in bits is therefore the number of cells times the size in bits of each cell, which is determined as the smallest power of 2 which is equal or greater than the number of sets to be stored plus one. So, if we need to store 255 sets, $l = 8 \cdot m$, as $2^8 = 256$.

We performed extensive experimental tests over implementations of both the shifting Bloom filter and the spatial Bloom filter using the test datasets described in Table 1. Both datasets have a total of 65280 elements distributed over 255 sets. In the uniform dataset, each set contains exactly 256 elements, while in the random dataset the elements are distributed randomly between sets. In order to test the false positive probability, we use an additional dataset of 500000 "non-elements" that are not part of either test datasets, and therefore should ideally be detected as not belonging to any of the originating sets.

Results of the false positive probability experiment are depicted in Figure 4. The figure plots the observed ratio of false positives over the uniform dataset. On the left hand side, the ratio is calculated over a varying filter legth in bits $l$, while on the right hand side this is done over the number of cells $m$. Given that $l = m$ for the ShBF, as described above, the experimentally observed probability remains the same, while the ratio changes for the SBF. From the experiment, we can see that the SBF performs better if we consider the number of cells $m$, resulting in a negligible number of false positives for $m \geq 2^{20}$. However, if we consider the memory space used by the filter and therefore the filter length in bits, ShBF can provide optimal results starting at $l = 2^{21}$ while the SBF at $l = 2^{23}$.

While the experiment depicted in Figure 4 was carried out keeping the number of originating sets fixed, it is important to point out that, concerning false positive probability, the SBF outperforms the ShBF as the number $s$ of originating sets increases. This condition is evident from (8) and (18), and is outlined in Figure 5. In particular, for $m = 20$ the false positive probability for the ShBF increases from close to 0 for a single set to over 0.1 for 250 sets (Figure 5a). In the case of a larger filter with $m = 23$, the false positive probability for the ShBF ranges from close to 0 to $1.4 \times 10^9$. The false positive probability of the SBF is instead independent of the number of sets, as expected.

Figure 6 depicts the behaviour of the ShBF and SBF over both the uniform and random datasets with respect to inter-set errors. The bar on the left hand side for both graphs shows results for the SBF over a number of cells $m = 2^{20}$ and filter length in bits $l = 2^{23}$. The ShBF is tested over the same number of cells $m = 2^{20}$ as the SBF in the bar at the centre; and over the same filter length in bits with $l = 2^{23}$ in the bar on the right hand side. Since some outcomes concerning errors are not easily visible, results of the experiment are also provided in Table 2.

Given that the computational cost (in terms of number of hash computations) increases significantly as the number of sets increases for the ShBF (as discussed in the following section, and Table 3), we performed the tests over large but compact datasets of 65280 elements and 255 sets. Larger datasets would

Table 1: Test datasets used for the experimental comparison of the two datastructures.

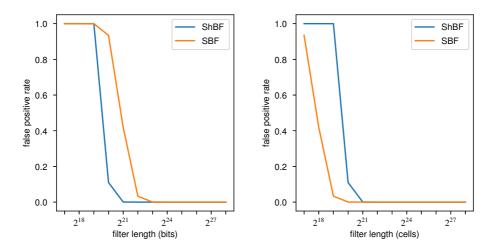| dataset | sets ($s$) | elem. per set | elements ($n$) |
|---|---|---|---|
| uniform | 255 | 256 | 65280 |
| random | 255 | [209, 298] | 65280 |
| non-elements | - | - | 500000 |

Figure 4: False positive probability for the two data structures over a uniformly distributed dataset, calculated over the same filter length in bits (left) and the same number of cells (right).
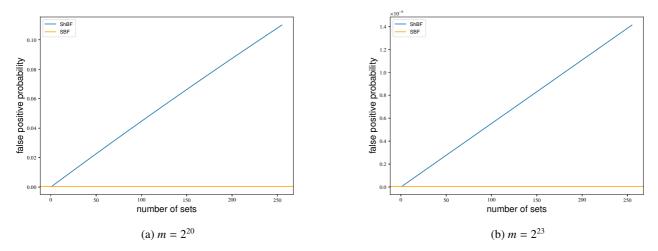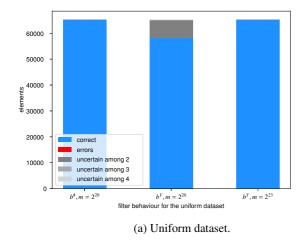


(a) $m = 2^{20}$

(b) $m = 2^{23}$

Figure 5: A plot of the false positive probability as expressed in (8) and (18) as the number of originating sets $s$ increases. The considered filter length is $m = 2^{20}$ (a) and $m = 2^{23}$ (b).

Table 2: Results of the experiment also depicted in Figure 6. $c$ is the number of elements whose set was correctly identified. For the SBF, the number of elements recognised as belonging to a wrong set is $e$. For the ShBF, $u_i$ is the number of elements that were identified as belonging to one out of $i$ sets. The resulting entropy $ent$ is also provided for both data structures, calculated as described in Section 4.

| filter | $m$ | $c$ | $e$ | $u_2$ | $u_3$ | $u_4$ | $u_5$ | $ent$ |
|--------|-----|-----|-----|-------|-------|-------|-------|-------|
| UNIFORM DATASET | | | | | | | | |
| SBF  | 20 | 65276 | 4 | -    | -   | -  | -  | 0.99994 |
| ShBF | 20 | 58174 | - | 6739 | 352 | 15 | 0  | 0.94462 |
| ShBF | 23 | 65276 | - | 4    | 0   | 0  | 0  | 0.99997 |
| RANDOM DATASET | | | | | | | | |
| SBF  | 20 | 65277 | 3 | -    | -   | -  | -  | 0.99995 |
| ShBF | 20 | 58282 | - | 6600 | 379 | 18 | 1  | 0.94536 |
| ShBF | 23 | 65278 | - | 2    | 0   | 0  | 0  | 0.99998 |

not provide significantly different results from the error probability perspective, given the filter parameters $m$ and $k$ are adjusted accordingly, but would incur in a largely increased computation time. Again, we notice here that the ShBF outperforms the SBF in terms of space efficiency, but not in terms of cell numbers. It is important to note here the different behaviour of the two data structures with respect to the occurrence of an inter-set error. In the case of the ShBF, an inter-set error happens when the originating set cannot be determined exactly, because more than one potential set is returned by the filter over a query. However, the correct set is always included in the returned list. In the case of the SBF, instead, and inter-set error returns a single set, which is not the correct originating set. In order to compare the different behaviours, we introduce an *entropy* metric here (modeled around the entropy concept first proposed by Shannon) which assesses the amount of informa-

(a) Uniform dataset.
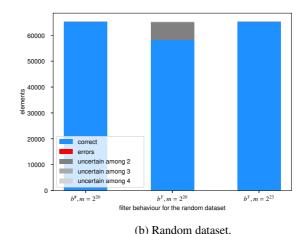


(b) Random dataset.

Figure 6: Filter behaviour of the two data structures over two datasets, (a) with a uniform distribution of elements among the 255 sets, (b) with a random distribution. The SBF is tested over a number of cells $m = 2^{20}$ and filter length in bits $l = 2^{23}$ (left). The ShBF is tested over the same number of cells $m = l = 2^{20}$ (centre); and over the same filter length in bits with $m = l = 2^{23}$ (right). For a ShBF, the number of cells and the filter length in bits coincide.

Table 3: Computational cost (expressed as number of hash computations required) for each lookup and query.

| filter | lookups/ query | hashes/ query | cells read / query [min, max] |
|--------|---------------|---------------|-------------------------------|
| ShBF   | $s$           | $k + s - 1$   | $[s, (s \cdot k)]$            |
| SBF    | $1$           | $k$           | $[1, k]$                      |

tion returned by a query. For the purpose of this work, we consider a correct result as having entropy 1, an incorrect result (possible only for the SBF) as having entropy 0, and a doubtful result (when multiple sets are returned by querying an ShBF) as $1/u$ where $u$ is the number of sets returned. Average entropy results (*ent*) for the experiment are provided in Table 2, where we call an event where $u = i$ as $u_i$. The results indicate that the ShBF has a marginally higher entropy for the same filter length $l$ for both distributions. Finally, we note here that while inter-set errors are distributed uniformly in the case of the ShBF, the SBF has a higher occurrence of inter-set errors for sets that have low index (that is, those that are entered first into the filter during construction). Therefore, the filter can be tuned to have different error probabilities for different sets, which may be an advantage in certain application scenarios.

### 4.1. Computational cost

In analysing the error probability of the two data structures, we referred to the memory usage of the filters. However, the operation of both data structures also implies a computational cost, which is different between the two constructions. In particular, in the following we analyse the computational cost associated with a single query to a filter. In general, cryptographic operations are the most computationally expensive, and therefore we analyse here the number of hash computations required for a query. Results are shown in Table 3.
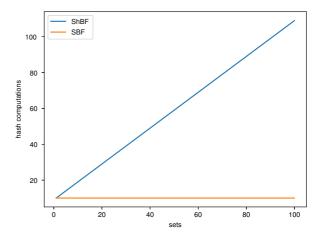


Figure 7: Number of hash computations required for a single query over the ShBF (linear) and SBF (constant) over the number of sets.

A shifting Bloom filter requires $(k + s - 1)$ hash digests to be computed for each query (where $k$ is the number of hash functions chosen at filter construction, and $s$ is the number of sets), and each query also involves $s$ lookups. Moreover, a single query involves reading a minimum of $s$ cells (if a 0 is returned in the first position of all set offsets) and a maximum of $(s \cdot k)$ (in the worst case scenario for which a 1 is returned for each cell read). A spatial Bloom filter, instead, requires a constant number of hash computations $k$ and a single lookup for each query, with a minimum of 1 and a maximum of $k$ cells read. We can therefore conclude that the computational cost of a ShBF is significantly higher per query than the cost of an equivalent SBF, as evident in Figure 7.

## 5. Conclusions

In this paper we compared the shifting Bloom filters and the spatial Bloom filters, two probabilistic data structures designed to support association queries. For the former data structure, we also provide a novel generalised definition allowing an unlimited amount of originating sets, and we discuss the resulting probabilistic model in detail. We implemented and tested both data structures over several datasets. Results show that the two data structures provide different benefits, and an adoption choice should depend on the application context. In particular, the probabilistic model of the spatial Bloom filter shows that it outperforms the shifting Bloom filter concerning both false positives and inter-set errors when the number of cells is considered. However, when effective memory consumption is considered, the shifting Bloom filter can achieve similar error rates to the spatial Bloom filter using less memory. With regards to computational costs, SBF requires a constant number of hash calculations, while ShBF requires a number of hash digest computations that increases linearly with the number of sets.

## 6. Acknowledgements

## References

[1] S. Tarkoma, C. E. Rothenberg, E. Lagerspetz, Theory and practice of bloom filters for distributed systems, IEEE Communications Surveys and Tutorials 14 (1) (2012) 131–155. doi:10.1109/SURV.2011.031611.00024.

[2] F. Chang, J. Dean, S. Ghemawat, W. C. Hsieh, D. A. Wallach, M. Burrows, T. Chandra, A. Fikes, R. E. Gruber, Bigtable: A distributed storage system for structured data, ACM Trans. Comput. Syst. 26 (2) (2008) 4:1–4:26. doi:10.1145/1365815.1365816.

[3] B. H. Bloom, Space/time trade-offs in hash coding with allowable errors, Communications of the ACM 13 (7) (1970) 422–426. doi:10.1145/362686.362692.

[4] F. Grandi, On the analysis of bloom filters, Inf. Process. Lett. 129 (2018) 35–39. doi:10.1016/j.ipl.2017.09.004.

[5] L. Luo, D. Guo, R. T. B. Ma, O. Rottenstreich, X. Luo, Optimizing bloom filter: Challenges, solutions, and comparisons, CoRR abs/1804.04777. arXiv:1804.04777.

[6] L. Liu, M. T. Özsu (Eds.), Encyclopedia of Database Systems, Second Edition, Springer, 2018. doi:10.1007/978-1-4614-8265-9.

[7] L. Fan, P. Cao, J. M. Almeida, A. Z. Broder, Summary cache: a scalable wide-area web cache sharing protocol, IEEE/ACM Trans. Netw. 8 (3) (2000) 281–293. doi:10.1109/90.851975.

[8] W. Wang, H. Jiang, H. Lu, J. X. Yu, Bloom histogram: Path selectivity estimation for XML data with updates, in: M. A. Nascimento, M. T. Özsu, D. Kossmann, R. J. Miller, J. A. Blakeley, K. B. Schiefer (Eds.), (e)Proceedings of the Thirtieth International Conference on Very Large Data Bases, Toronto, Canada, August 31 - September 3 2004, Morgan Kaufmann, 2004, pp. 240–251.

[9] S. Geravand, M. Ahmadi, Bloom filter applications in network security: A state-of-the-art survey, Computer Networks 57 (18) (2013) 4047–4064. doi:10.1016/j.comnet.2013.09.003.

[10] A. Z. Broder, M. Mitzenmacher, Survey: Network applications of bloom filters: A survey, Internet Mathematics 1 (4) (2003) 485–509. doi:10.1080/15427951.2004.10129096.

[11] S. Dharmapurikar, P. Krishnamurthy, D. E. Taylor, Longest prefix matching using bloom filters, in: A. Feldmann, M. Zitterbart, J. Crowcroft, D. Wetherall (Eds.), Proceedings of the ACM SIGCOMM 2003 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communication, August 25-29, 2003, Karlsruhe, Germany, ACM, 2003, pp. 201–212. doi:10.1145/863955.863979.

[12] H. Song, F. Hao, M. S. Kodialam, T. V. Lakshman, Ipv6 lookups using distributed and load balanced bloom filters for 100gbps core router line cards, in: INFOCOM 2009. 28th IEEE International Conference on Computer Communications, Joint Conference of the IEEE Computer and Communications Societies, 19-25 April 2009, Rio de Janeiro, Brazil [25], pp. 2518–2526. doi:10.1109/INFCOM.2009.5062180.

[13] F. Chang, K. Li, W. Feng, Approximate caches for packet classification, in: Proceedings IEEE INFOCOM 2004, The 23rd Annual Joint Conference of the IEEE Computer and Communications Societies, Hong Kong, China, March 7-11, 2004, IEEE, 2004, pp. 2196–2207. doi:10.1109/INFCOM.2004.1354643.

[14] D. Yang, D. Tian, J. Gong, S. Gao, T. Yang, X. Li, Difference bloom filter: A probabilistic structure for multi-set membership query, in: IEEE International Conference on Communications, ICC 2017, Paris, France, May 21-25, 2017, IEEE, 2017, pp. 1–6. doi:10.1109/ICC.2017.7996678.

[15] T. Yang, A. X. Liu, M. Shahzad, Y. Zhong, Q. Fu, Z. Li, G. Xie, X. Li, a shifting bloom filter framework for set queries, Proceedings of the VLDB Endowment 9 (5) (2016) 408–419. doi:10.14778/2876473.2876476.

[16] F. Hao, M. S. Kodialam, T. V. Lakshman, H. Song, Fast multiset membership testing using combinatorial bloom filters, in: INFOCOM 2009. 28th IEEE International Conference on Computer Communications, Joint Conference of the IEEE Computer and Communications Societies, 19-25 April 2009, Rio de Janeiro, Brazil [25], pp. 513–521. doi:10.1109/INFCOM.2009.5061957.

[17] F. Hao, M. S. Kodialam, T. V. Lakshman, H. Song, Fast dynamic multiple-set membership testing using combinatorial bloom filters, IEEE/ACM Trans. Netw. 20 (1) (2012) 295–304. doi:10.1109/TNET.2011.2173351.

[18] P. Palmieri, L. Calderoni, D. Maio, Spatial bloom filters: Enabling privacy in location-aware applications, in: D. Lin, M. Yung, J. Zhou (Eds.), Information Security and Cryptology - 10th International Conference, Inscrypt 2014, Beijing, China, December 13-15, 2014, Revised Selected Papers, 2014, pp. 16–36. doi:10.1007/978-3-319-16745-9_2.

[19] L. Calderoni, P. Palmieri, D. Maio, Location privacy without mutual trust: The spatial bloom filter, Computer Communications 68 (2015) 4–16. doi:10.1016/j.comcom.2015.06.011.

[20] S. Xiong, Y. Yao, Q. Cao, T. He, kbf: A bloom filter for key-value storage with an application on approximate state machines, in: 2014 IEEE Conference on Computer Communications, INFOCOM 2014, Toronto, Canada, April 27 - May 2, 2014, IEEE, 2014, pp. 1150–1158. doi:10.1109/INFOCOM.2014.6848046.

[21] B. Chazelle, J. Kilian, R. Rubinfeld, A. Tal, The bloomier filter: an efficient data structure for static support lookup tables, in: J. I. Munro (Ed.), Proceedings of the Fifteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2004, New Orleans, Louisiana, USA, January 11-14, 2004, SIAM, 2004, pp. 30–39.

[22] D. X. Charles, K. Chellapilla, Bloomier filters: A second look, in: D. Halperin, K. Mehlhorn (Eds.), Algorithms - ESA 2008, 16th Annual European Symposium, Karlsruhe, Germany, September 15-17, 2008. Proceedings, Vol. 5193 of Lecture Notes in Computer Science, Springer, 2008, pp. 259–270. doi:10.1007/978-3-540-87744-8\_22.

[23] P. Palmieri, L. Calderoni, D. Maio, Private inter-network routing for wireless sensor networks and the internet of things, in: Proceedings of the Computing Frontiers Conference, CF'17, Siena, Italy, May 15-17, 2017, ACM, 2017, pp. 396–401. doi:10.1145/3075564.3079068.

[24] L. Calderoni, P. Palmieri, D. Maio, Probabilistic properties of the spatial bloom filters and their relevance to cryptographic protocols, IEEE Transactions on Information Forensics and Security 13 (7) (2018) 1710–1721. doi:10.1109/TIFS.2018.2799486.

[25] INFOCOM 2009. 28th IEEE International Conference on Computer Communications, Joint Conference of the IEEE Computer and Communications Societies, 19-25 April 2009, Rio de Janeiro, Brazil, IEEE, 2009.