



The diagram illustrates the Bloom filter process for key-value storage. It is organized into three columns: **QUERY**, **FILTER**, and **STORAGE**.

- QUERY 1:** "Do you have 'key1'?" The **FILTER** response is "No". The **STORAGE** response is "No".
- QUERY 2:** "Do you have 'key2'?" The **FILTER** response is "Yes". The **STORAGE** response is "Yes". A green arrow labeled "necessary disk access" points from the filter to storage.
- QUERY 3:** "Do you have 'key3'?" The **FILTER** response is "Yes" (labeled "False Positive"). The **STORAGE** response is "No". A red arrow labeled "unnecessary disk access" points from the filter to storage.

At the bottom, a text box states: "Bloom filter used to speed up answers in a key-value storage".

implementation, however, the Bloom filter shines because its  $k$  lookups are independent and can be parallelized.

To understand its space efficiency, it is instructive to compare the general Bloom filter with its special case when  $k = 1$ . If  $k = 1$ , then in order to keep the false positive rate sufficiently low, a small fraction of bits should be set, which means the array must be very large and contain long runs of zeros. The [information content](#) of the array relative to its size is low. The generalized Bloom filter ( $k$  greater than 1) allows many more bits to be set while still maintaining a low false positive rate; if the parameters ( $k$  and  $m$ ) are chosen well, about half of the bits will be set,<sup>[2]</sup> and these will be apparently random, minimizing redundancy and maximizing information content.

## Probability of false positives [\[ edit \]](#)

Assume that a [hash function](#) selects each array position with equal probability. If  $m$  is the number of bits in the array, the probability that a certain bit is not set to 1 by a certain hash function during the insertion of an element is

$$1 - \frac{1}{m}.$$

If  $k$  is the number of hash functions, the probability that the bit is not set to 1 by any of the hash functions is

$$\left(1 - \frac{1}{m}\right)^k.$$

If we have inserted  $n$  elements, the probability that a certain bit is still 0 is

$$\left(1 - \frac{1}{m}\right)^{kn};$$

the probability that it is 1 is therefore

$$1 - \left(1 - \frac{1}{m}\right)^{kn}.$$

Now test membership of an element that is not in the set. Each of the  $k$  array positions computed by the hash functions is 1 with a probability as above. The probability of all of them being 1, which would cause the [algorithm](#) to erroneously claim that the element is in the set, is often given as

$$\left(1 - \left[1 - \frac{1}{m}\right]^{kn}\right)^k \approx \left(1 - e^{-kn/m}\right)^k.$$

This is not strictly correct as it assumes independence for the probabilities of each bit being set. However, assuming it is a close approximation we have that the probability of false positives decreases as  $m$  (the number of bits in the array) increases, and increases as  $n$  (the number of inserted elements) increases.

An alternative analysis arriving at the same approximation without the assumption of independence is given by Mitzenmacher and Upfal.<sup>[3]</sup> After all  $n$  items have been added to the Bloom filter, let  $q$  be the fraction of the  $m$  bits that are set to 0. (That is, the number of bits still set to 0 is  $qm$ .) Then, when testing membership of an element not in the set, for the array position given by any of the  $k$  hash functions, the probability that the bit is found set to 1 is  $1 - q$ . So the probability that all  $k$  hash functions find their bit set to 1 is  $(1 - q)^k$ . Further, the expected value of  $q$  is the probability that a given array position is left untouched by each of the  $k$  hash functions for each of the  $n$  items, which is (as above)

$$E[q] = \left(1 - \frac{1}{m}\right)^{kn}.$$

It is possible to prove, without the independence assumption, that  $q$  is very strongly concentrated around its expected value. In particular, from the [Azuma–Hoeffding inequality](#), they prove that<sup>[4]</sup>

$$\Pr(|q - E[q]| \geq \frac{\lambda}{m}) \leq 2 \exp(-2\lambda^2/m)$$

Because of this, we can say that the exact probability of false positives is

$$\sum_t \Pr(q = t)(1 - t)^k \approx (1 - E[q])^k = \left(1 - \left[1 - \frac{1}{m}\right]^{kn}\right)^k \approx \left(1 - e^{-kn/m}\right)^k$$

as before.

## Optimal number of hash functions [\[ edit \]](#)

For a given  $m$  and  $n$ , the value of  $k$  (the number of hash functions) that minimizes the false positive probability is

$$k = \frac{m}{n} \ln 2,$$

which gives

$$2^{-k} \approx 0.6185^{m/n}.$$

The required number of bits  $m$ , given  $n$  (the number of inserted elements) and a desired false positive probability  $p$  (and assuming the optimal value of  $k$  is used) can be computed by substituting the optimal value of  $k$  in the probability expression above:

$$p = \left(1 - e^{-(m/n \ln 2)n/m}\right)^{(m/n \ln 2)}$$

which can be simplified to:

$$\ln p = -\frac{m}{n} (\ln 2)^2.$$

This results in:

$$m = -\frac{n \ln p}{(\ln 2)^2}.$$

This means that for a given false positive probability  $p$ , the length of a Bloom filter  $m$  is proportionate to the number of elements being filtered  $n$ .<sup>[5]</sup> While the above formula is asymptotic (i.e. applicable as  $m, n \rightarrow \infty$ ), the agreement with finite values of  $m, n$  is also quite good; the false positive probability for a finite bloom filter with  $m$  bits,  $n$  elements, and  $k$  hash functions is at most

$$\left(1 - e^{-k(n+0.5)/(m-1)}\right)^k.$$

So we can use the asymptotic formula if we pay a penalty for at most half an extra element and at most one fewer bit.<sup>[6]</sup>

## Approximating the number of items in a Bloom filter [\[ edit \]](#)

[Swamidass & Baldi \(2007\)](#) showed that the number of items in a Bloom filter can be approximated with the following formula,

$$n^* = -\frac{m \ln \left[1 - \frac{X}{m}\right]}{k}$$

where  $n^*$  is an estimate of the number of items in the filter,  $m$  is the length (size) of the filter,  $k$  is the number of hash functions, and  $X$  is the number of bits set to one.

## The union and intersection of sets [\[ edit \]](#)

Bloom filters are a way of compactly representing a set of items. It is common to try to compute the size of the intersection or union between two sets. Bloom filters can be used to approximate the size of the intersection and union of two sets. [Swamidass & Baldi \(2007\)](#) showed that for two bloom filters of length  $m$ , their counts, respectively can be estimated as

$$n(A^*) = -m \ln [1 - n(A)/m] / k$$

and

$$n(B^*) = -m \ln [1 - n(B)/m] / k.$$

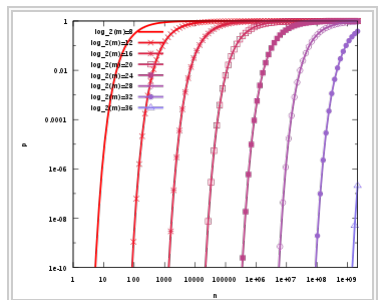
The size of their union can be estimated as

$$n(A^* \cup B^*) = -m \ln [1 - n(A \cup B)/m] / k$$

where  $n(A \cup B)$  is the number of bits set to one in either of the two bloom filters. Finally, the intersection can be estimated as

$$n(A^* \cap B^*) = n(A^*) + n(B^*) - n(A^* \cup B^*).$$

system. Values are stored on a disk which has slow access times. Bloom filter decisions are much faster. However some unnecessary disk accesses are made when the filter reports a positive (in order to weed out the false positives). Overall answer speed is better with the Bloom filter than without the Bloom filter. Use of a Bloom filter for this purpose, however, does increase memory usage.



The false positive probability  $p$  as a function of number of elements  $n$  in the filter and the filter size  $m$ . An optimal number of hash functions  $k = (m/n) \ln 2$  has been assumed.

using the three formulas together.

## Interesting properties [ [edit](#) ]

- Unlike a standard [hash table](#), a Bloom filter of a fixed size can represent a set with an arbitrarily large number of elements; adding an element never fails due to the data structure "filling up." However, the false positive rate increases steadily as elements are added until all bits in the filter are set to 1, at which point *all* queries yield a positive result.
- [Union](#) and [intersection](#) of Bloom filters with the same size and set of hash functions can be implemented with [bitwise](#) OR and AND operations, respectively. The union operation on Bloom filters is lossless in the sense that the resulting Bloom filter is the same as the Bloom filter created from scratch using the union of the two sets. The intersect operation satisfies a weaker property: the false positive probability in the resulting Bloom filter is at most the false-positive probability in one of the constituent Bloom filters, but may be larger than the false positive probability in the Bloom filter created from scratch using the intersection of the two sets.
- Some kinds of [superimposed code](#) can be seen as a Bloom filter implemented with physical [edge-notched cards](#). An example is [Zatocoding](#), invented by [Calvin Mooers](#) in 1947, in which the set of categories associated with a piece of information is represented by notches on a card, with a random pattern of four notches for each category.

## Examples [ [edit](#) ]

- Google [BigTable](#), [Apache HBase](#) and [Apache Cassandra](#) use Bloom filters to reduce the disk lookups for non-existent rows or columns. Avoiding costly disk lookups considerably increases the performance of a database query operation.<sup>[7][8]</sup>
- The [Google Chrome](#) web browser used to use a Bloom filter to identify malicious URLs. Any URL was first checked against a local Bloom filter, and only if the Bloom filter returned a positive result was a full check of the URL performed (and the user warned, if that too returned a positive result).<sup>[9][10]</sup>
- The [Squid Web Proxy Cache](#) uses Bloom filters for [cache digests](#).<sup>[11]</sup>
- [Bitcoin](#) uses Bloom filters to speed up wallet synchronization.<sup>[12][13]</sup>
- The [Venti](#) archival storage system uses Bloom filters to detect previously stored data.<sup>[14]</sup>
- The [SPIN model checker](#) uses Bloom filters to track the reachable state space for large verification problems.<sup>[15]</sup>
- The [Cascading](#) analytics framework uses Bloom filters to speed up asymmetric joins, where one of the joined data sets is significantly larger than the other (often called Bloom join<sup>[16]</sup> in the database literature).<sup>[17]</sup>
- The [Exim](#) mail transfer agent (MTA) uses bloom filters in its rate-limit feature.<sup>[18]</sup>

## Alternatives [ [edit](#) ]

Classic Bloom filters use 1.44 log<sub>2</sub>⁡(1/ϵ) bits of space per inserted key, where ϵ is the false positive rate of the Bloom filter. However, the space that is strictly necessary for any data structure playing the same role as a Bloom filter is only log<sub>2</sub>⁡(1/ϵ) per key ([Pagh, Pagh & Rao 2005](#)). Hence Bloom filters use 44% more space than a hypothetical equivalent optimal data structure. The number of hash functions used to achieve a given false positive rate ϵ is proportional to log(1/ϵ) which is not optimal as it has been proved that an optimal data structure would need only a constant number of hash functions independent of the false positive rate.

[Stern & Dill \(1996\)](#) describe a probabilistic structure based on [hash tables](#), [hash compaction](#), which [Dillinger & Manolios \(2004b\)](#) identify as significantly more accurate than a Bloom filter when each is configured optimally. Dillinger and Manolios, however, point out that the reasonable accuracy of any given Bloom filter over a wide range of numbers of additions makes it attractive for probabilistic enumeration of state spaces of unknown size. Hash compaction is, therefore, attractive when the number of additions can be predicted accurately; however, despite being very fast in software, hash compaction is poorly suited for hardware because of worst-case linear access time.

[Putze, Sanders & Singler \(2007\)](#) have studied some variants of Bloom filters that are either faster or use less space than classic Bloom filters. The basic idea of the fast variant is to locate the k hash values associated with each key into one or two blocks having the same size as processor's memory cache blocks (usually 64 bytes). This will presumably improve performance by reducing the number of potential memory [cache misses](#). The proposed variants have however the drawback of using about 32% more space than classic Bloom filters.

The space efficient variant relies on using a single hash function that generates for each key a value in the range [0, *n*/ϵ] where ϵ is the requested false positive rate. The sequence of values is then sorted and compressed using [Golomb coding](#) (or some other compression technique) to occupy a space close to *n* log<sub>2</sub>⁡(1/ϵ) bits. To query the Bloom filter for a given key, it will suffice to check if its corresponding value is stored in the Bloom filter. Decompressing the whole Bloom filter for each query would make this variant totally unusable. To overcome this problem the sequence of values is divided into small blocks of equal size that are compressed separately. At query time only half a block will need to be decompressed on average. Because of decompression overhead, this variant may be slower than classic Bloom filters but this may be compensated by the fact that a single hash function need to be computed.

Another alternative to classic Bloom filter is the one based on space efficient variants of [cuckoo hashing](#). In this case once the hash table is constructed, the keys stored in the hash table are replaced with short signatures of the keys. Those signatures are strings of bits computed using a hash function applied on the keys.

## Extensions and applications [ [edit](#) ]

### Counting filters [ [edit](#) ]

Counting filters provide a way to implement a *delete* operation on a Bloom filter without recreating the filter afresh. In a counting filter the array positions (buckets) are extended from being a single bit to being an n-bit counter. In fact, regular Bloom filters can be considered as counting filters with a bucket size of one bit. Counting filters were introduced by [Fan et al. \(1998\)](#).

The insert operation is extended to *increment* the value of the buckets and the lookup operation checks that each of the required buckets is non-zero. The delete operation, obviously, then consists of decrementing the value of each of the respective buckets.

[Arithmetic overflow](#) of the buckets is a problem and the buckets should be sufficiently large to make this case rare. If it does occur then the increment and decrement operations must leave the bucket set to the maximum possible value in order to retain the properties of a Bloom filter.

The size of counters is usually 3 or 4 bits. Hence counting Bloom filters use 3 to 4 times more space than static Bloom filters. In theory, an optimal data structure equivalent to a counting Bloom filter should not use more space than a static Bloom filter.

Another issue with counting filters is limited scalability. Because the counting Bloom filter table cannot be expanded, the maximal number of keys to be stored simultaneously in the filter must be known in advance. Once the designed capacity of the table is exceeded, the false positive rate will grow rapidly as more keys are inserted.

[Bonomi et al. \(2006\)](#) introduced a data structure based on d-left hashing that is functionally equivalent but uses approximately half as much space as counting Bloom filters. The scalability issue does not occur in this data structure. Once the designed capacity is exceeded, the keys could be reinserted in a new hash table of double size.

The space efficient variant by [Putze, Sanders & Singler \(2007\)](#) could also be used to implement counting filters by supporting insertions and deletions.

[Rottenstreich, Kanizo & Keslassy \(2012\)](#) introduced a new general method based on variable increments that significantly improves the false positive probability of counting Bloom filters and their variants, while still supporting deletions. Unlike counting Bloom filters, at each element insertion, the hashed counters are incremented by a hashed variable increment instead of a unit increment. To query an element, the exact values of the counters are considered and not just their positiveness. If a sum represented by a counter value cannot be composed of the corresponding variable increment for the queried element, a negative answer can be returned to the query.

### Decentralized aggregation [ [edit](#) ]

Bloom filters can be organized in distributed [data structures](#) to perform fully decentralized computations of [aggregate functions](#). Decentralized aggregation makes collective measurements locally available in every node of a distributed network without involving a centralized computational entity for this purpose.<sup>[19][20]</sup>

### Data synchronization [ [edit](#) ]

Bloom filters can be used for approximate [data synchronization](#) as in [Byers et al. \(2004\)](#). Counting Bloom filters can be used to approximate the number of differences between two sets and this approach is described in [Agarwal & Trachtenberg \(2006\)](#).

### Bloomier filters [ [edit](#) ]

[Chazelle et al. \(2004\)](#) designed a generalization of Bloom filters that could associate a value with each element that had been inserted, implementing an [associative array](#). Like Bloom filters, these structures achieve a small space overhead by accepting a small probability of false positives. In the case of "Bloomier filters", a *false positive* is defined as returning a result when the key is not in the map. The map will never return the wrong value for a key that *is* in the map.

### Compact approximators [ [edit](#) ]

[Boldi & Vigna \(2005\)](#) proposed a [lattice](#)-based generalization of Bloom filters. A **compact approximator** associates to each key an element of a lattice (the standard Bloom filters being the case of the Boolean two-element lattice). Instead of a bit array, they have an array of lattice elements. When adding a new association between a key and an element of the lattice, they compute the maximum of the current contents of the *k* array locations associated to the key with the lattice element. When reading the value associated to a key, they compute the minimum of the values found in the *k* locations associated to the key. The resulting value approximates from above the original value.

### Stable Bloom filters [ [edit](#) ]

[Deng & Rafiei \(2006\)](#) proposed Stable Bloom filters as a variant of Bloom filters for streaming data. The idea is that since there is no way to store the entire history of a stream (which can be infinite),

Stable Bloom filters continuously evict stale information to make room for more recent elements. Since stale information is evicted, the Stable Bloom filter introduces false negatives, which do not appear in traditional bloom filters. The authors show that a tight upper bound of false positive rates is guaranteed, and the method is superior to standard bloom filters in terms of false positive rates and time efficiency when a small space and an acceptable false positive rate are given.

#### Scalable Bloom filters [\[ edit \]](#)

[Almeida et al. \(2007\)](#) proposed a variant of Bloom filters that can adapt dynamically to the number of elements stored, while assuring a minimum false positive probability. The technique is based on sequences of standard bloom filters with increasing capacity and tighter false positive probabilities, so as to ensure that a maximum false positive probability can be set beforehand, regardless of the number of elements to be inserted.

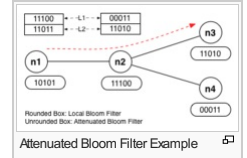
#### Layered Bloom filters [\[ edit \]](#)

A layered bloom filter consists of multiple bloom filter layers. Layered bloom filters allow keeping track of how many times an item was added to the bloom filter by checking how many layers contain the item. With a layered bloom filter a check operation will normally return the deepest layer number the item was found in.<sup>[[21](#)]</sup><sup>[[22](#)]</sup>

#### Attenuated Bloom filters [\[ edit \]](#)

An attenuated bloom filter of depth D can be viewed as an array of D normal bloom filters. In the context of service discovery in a network, each node stores regular and attenuated bloom filters locally. The regular or local bloom filter indicates which services are offered by the node itself. The attenuated filter of level i indicates which services can be found on nodes that are i-hops away from the current node. The i-th value is constructed by taking a union of local bloom filters for nodes i-hops away from the node.<sup>[[23](#)]</sup>

Let's take a small network shown on the graph below as an example. Say we are searching for a service A whose id hashes to bits 0,1, and 3 (pattern 11010). Let n1 node be the starting point. First, we check whether service A is offered by n1 by checking its local filter. Since the patterns don't match, we check the attenuated bloom filter in order to determine which node should be the next hop. We see that n2 doesn't offer service A but lies on the path to nodes that do. Hence, we move to n2 and repeat the same procedure. We quickly find that n3 offers the service, and hence the destination is located.<sup>[[24](#)]</sup>



By using attenuated Bloom filters consisting of multiple layers, services at more than one hop distance can be discovered while avoiding saturation of the Bloom filter by attenuating (shifting out) bits set by sources further away.<sup>[[23](#)]</sup>

#### Chemical structure searching [\[ edit \]](#)

Bloom filters are often used to search large chemical structure databases (see [chemical similarity](#)). In the simplest case, the elements added to the filter (called a fingerprint in this field) are just the atomic numbers present in the molecule, or a hash based on the atomic number of each atom and the number and type of its bonds. This case is too simple to be useful. More advanced filters also encode atom counts, larger substructure features like carboxyl groups, and graph properties like the number of rings. In hash-based fingerprints, a hash function based on atom and bond properties is used to turn a subgraph into a [PRNG](#) seed, and the first output values used to set bits in the Bloom filter.

Molecular fingerprints arose as a way to screen out obvious rejects in [molecular subgraph searches](#). They are more often used in calculating molecular similarity, computed as the [Tanimoto](#) similarity between their respective fingerprints. The Daylight and Indigo fingerprints, use a variable length (which are restricted to be powers of two) that adapts to the number of items to ensure the final filter is not oversaturated. A length for the filter is chosen to keep the saturation of the filter approximately at a fixed value (for example 30%).

Strictly speaking, the MACCS keys and CACTVS keys are not Bloom filters, rather they are deterministic bit-arrays. Similarly, they are often incorrectly considered to be molecular fingerprints, but they are actually "structural keys". They do not use hash functions, but use a dictionary to map specific substructures to specific indices in the filter. In contrast with Bloom filters, this is a one-to-one mapping that does not use hash functions at all.

#### See also [\[ edit \]](#)

- Count–min sketch
- Feature hashing
- Min-Hash
- Quotient filter
- Skip list



#### Notes [\[ edit \]](#)

- ↑ Mitzenmacher & Upfal (2005).
- ↑ Blustein & El-Maazawi (2002), pp. 21–22
- ↑ Mitzenmacher & Upfal (2005), pp. 109–111, 308.
- ↑ Mitzenmacher & Upfal (2005), p. 308.
- ↑ Starobinski, Trachtenberg & Agarwal (2003)
- ↑ Goel & Gupta (2010).
- ↑ (Chang et al. 2006).
- ↑ "11.6. Schema Design" ⓘ. *apache.org*.
- ↑ Yakunin, Alex (2010-03-25). "Alex Yakunin's blog: Nice Bloom filter application" ⓘ. *Blog.alexyakunin.com*. Retrieved 2014-05-31.
- ↑ "Issue 10896048: Transition safe browsing from bloom filter to prefix set. - Code Review" ⓘ. *Chromiumcodereview.appspot.com*. Retrieved 2014-07-03.
- ↑ Wessels, Duane (January 2004), "10.7 Cache Digests", *Squid: The Definitive Guide* (1st ed.), O'Reilly Media, p. 172, ISBN 0-596-00162-2, "Cache Digests are based on a technique first published by Pei Cao, called Summary Cache. The fundamental idea is to use a Bloom filter to represent the cache contents."
- ↑ *Bitcoin 0.8.0* ⓘ
- ↑ "The Bitcoin Foundation - Supporting the development of Bitcoin" ⓘ. *bitcoinfoundation.org*.
- ↑ "Plan 9 /sys/man/8/venti" ⓘ. *Plan9.bell-labs.com*. Retrieved 2014-05-31.
- ↑ *http://spinroot.com/* ⓘ
- ↑ Mullin (1990)
- ↑ "BloomJoin: BloomFilter + CoGroup | LiveRamp Blog" ⓘ. *Blog.liveramp.com*. Retrieved 2014-05-31.
- ↑ "Exim source code" ⓘ. *github*. Retrieved 2014-03-03.
- ↑ Poumaras, Warnier & Brazier (2013)
- ↑ "DIAS - Dynamic Intelligent Aggregation Service" ⓘ.
- ↑ "pmylund/go-bloom" ⓘ. *Github.com*. Retrieved 2014-06-13.
- ↑ First Name Middle Name Last Name (2010-08-22). "IEEE Xplore Abstract - A multi-layer bloom filter for duplicated URL detection" ⓘ. *leexplore.ieee.org*. doi:10.1109/ICACTE.2010.5578947 ⓘ. Retrieved 2014-06-13.
- ↑ *a b* Koucheryav et al. (2009)

#### References [\[ edit \]](#)

- Koucheryav, Y.; Giambene, G.; Staehle, D.; Barcelo-Arroyo, F.; Braun, T.; Siris, V. (2009), "Traffic and QoS Management in Wireless Multimedia Networks", *COST 290 Final Report* (USA): 111
- Kubiatowicz, J.; Bindel, D.; Czerwinski, Y.; Geels, S.; Eaton, D.; Gummadi, R.; Rhea, S.; Weatherspoon, H. et al. (2000), "Oceanstore: An architecture for global-scale persistent storage" ⓘ (PDF), *ACM SIGPLAN Notices* (USA): 190–201
- Agarwal, Sachin; Trachtenberg, Ari (2006), "Approximating the number of differences between remote sets" ⓘ (PDF), *IEEE Information Theory Workshop* (Punta del Este, Uruguay): 217, doi:10.1109/ITW.2006.1633815 ⓘ, ISBN 1-4244-0035-X
- Ahmadi, Mahmood; Wong, Stephan (2007), "A Cache Architecture for Counting Bloom Filters", *15th international Conference on Networks (ICON-2007)* ⓘ, p. 218, doi:10.1109/ICON.2007.4444089 ⓘ, ISBN 978-1-4244-1229-7
- Almeida, Paulo; Baquero, Carlos; Pregaica, Nuno; Hutchison, David (2007), "Scalable Bloom Filters" ⓘ (PDF), *Information Processing Letters* **101** (6): 255–261, doi:10.1016/j.ipl.2006.10.007 ⓘ
- Byers, John W.; Considine, Jeffrey; Mitzenmacher, Michael; Rost, Stanislav (2004), "Informed content delivery across adaptive overlay networks", *IEEE/ACM Transactions on Networking* **12** (5): 767, doi:10.1109/TNET.2004.836103 ⓘ
- Bloom, Burton H. (1970), "Space/Time Trade-offs in Hash Coding with Allowable Errors" ⓘ, *Communications of the ACM* **13** (7): 422–426, doi:10.1145/362686.362692 ⓘ
- Boldi, Paolo; Vigna, Sebastiano (2005), "Mutable strings in Java: design, implementation and lightweight text-search algorithms", *Science of Computer Programming* **54** (1): 3–23, doi:10.1016/j.scico.2004.05.003 ⓘ
- Bonomi, Flavio; Mitzenmacher, Michael; Panigrahy, Rina; Singh, Sushil; Varghese, George (2006), "An Improved Construction for Counting Bloom Filters", *Algorithms – ESA 2006, 14th Annual European Symposium* ⓘ (PDF), *Lecture Notes in Computer Science* **4168**, pp. 684–695, doi:10.1007/11841036\_61 ⓘ, ISBN 978-3-540-38875-3
- Broder, Andrei; Mitzenmacher, Michael (2005), "Network Applications of Bloom Filters: A Survey" ⓘ (PDF), *Internet Mathematics* **1** (4): 485–509, doi:10.1080/15427951.2004.10129096 ⓘ
- Chang, Fay; Dean, Jeffrey; Ghemawat, Sanjay; Hsieh, Wilson; Wallach, Deborah; Burrows, Mike; Chandra, Tushar; Fikes, Andrew; Gruber, Robert (2006), "Bigtable: A Distributed Storage System for Structured Data", *Seventh Symposium on Operating System Design and Implementation* ⓘ
- Charles, Denis; Chellapilla, Kumar (2008), "Bloomier Filters: A second look", *The Computing Research Repository (CoRR)*, arXv:0807.0928 ⓘ
- Chazelle, Bernard; Kilian, Joe; Rubinfeld, Ronitt; Tal, Ayellet (2004), "The Bloomier filter: an efficient data structure for static support lookup tables", *Proceedings of the Fifteenth Annual ACM-SIAM Symposium on Discrete Algorithms* ⓘ (PDF), pp. 30–39
- Cohen, Saar; Matias, Yossi (2003), "Spectral Bloom Filters", *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data* ⓘ (PDF), pp. 241–252, doi:10.1145/872757.872787 ⓘ, ISBN 1-55860-426-2 ⓘ, direct link

ISBN 1581136347unsub mmy

- Deng, Fan; Rafiei, Davood (2006), "Approximately Detecting Duplicates for Streaming Data using Stable Bloom Filters", *Proceedings of the ACM SIGMOD Conference*  (PDF), pp. 25–36
- Dharmapuriar, Sarang; Song, Haoyu; Turner, Jonathan; Lockwood, John (2006), "Fast packet classification using Bloom filters", *Proceedings of the 2006 ACM/IEEE Symposium on Architecture for Networking and Communications Systems*  (PDF), pp. 61–70, doi:10.1145/1185347.1185356 , ISBN 1595935800
- Dietzfelbinger, Martin; Pagh, Rasmus (2008), "Succinct Data Structures for Retrieval and Approximate Membership", *The Computing Research Repository (CoRR)*, arXiv:0803.3693
- Swamidass, S. Joshua; Baldi, Pierre (2007), "Mathematical correction for fingerprint similarity measures to improve chemical retrieval", *Journal of chemical information and modeling* (ACS Publications) **47** (3): 952–964, doi:10.1021/ci600526a , PMID 17444629
- Dillinger, Peter C.; Manolios, Panagiotis (2004a), "Fast and Accurate Bitstate Verification for SPIN", *Proceedings of the 11th International Spin Workshop on Model Checking Software* , Springer-Verlag, Lecture Notes in Computer Science 2989
- Dillinger, Peter C.; Manolios, Panagiotis (2004b), "Bloom Filters in Probabilistic Verification", *Proceedings of the 5th International Conference on Formal Methods in Computer-Aided Design* , Springer-Verlag, Lecture Notes in Computer Science 3312
- Donnet, Benoit; Baynat, Bruno; Friedman, Timur (2006), "Retouched Bloom Filters: Allowing Networked Applications to Flexibly Trade Off False Positives Against False Negatives", *CoNEXT 06 – 2nd Conference on Future Networking Technologies*
- Eppstein, David; Goodrich, Michael T. (2007), "Space-efficient straggler identification in round-trip data streams via Newton's identities and invertible Bloom filters", *Algorithms and Data Structures, 10th International Workshop, WADS 2007*, Springer-Verlag, Lecture Notes in Computer Science 4619, pp. 637–648, arXiv:0704.3313
- Fan, Li; Cao, Pei; Almeida, Jussara; Broder, Andrei (2000), "Summary Cache: A Scalable Wide-Area Web Cache Sharing Protocol", *IEEE/ACM Transactions on Networking* **8** (3): 281–293, doi:10.1109/90.851975 . A preliminary version appeared at SIGCOMM '98.
- Goel, Ashish; Gupta, Pankaj (2010), "Small subset queries and bloom filters using ternary associative memories, with applications", *ACM Sigmetrics 2010* **38**: 143, doi:10.1145/1811099.1811056
- Kirsch, Adam; Mitzenmacher, Michael (2006), "Less Hashing, Same Performance: Building a Better Bloom Filter", in Azar, Yossi; Erlebach, Thomas, *Algorithms – ESA 2006, 14th Annual European Symposium*  (PDF), Lecture Notes in Computer Science **4168**, Springer-Verlag, Lecture Notes in Computer Science 4168, pp. 456–467, doi:10.1007/11841036 , ISBN 978-3-540-38875-3
- Mortensen, Christian Worm; Pagh, Rasmus; Patrascu, Mihai (2005), "On dynamic range reporting in one dimension", *Proceedings of the Thirty-seventh Annual ACM Symposium on Theory of Computing*, pp. 104–111, doi:10.1145/1060590.1060606 , ISBN 1581139608
- Pagh, Anna; Pagh, Rasmus; Rao, S. Srinivasa (2005), "An optimal Bloom filter replacement", *Proceedings of the Sixteenth Annual ACM-SIAM Symposium on Discrete Algorithms*  (PDF), pp. 823–829
- Porat, Ely (2008), "An Optimal Bloom Filter Replacement Based on Matrix Solving", *The Computing Research Repository (CoRR)*, arXiv:0804.1845
- Putze, F.; Sanders, P.; Singler, J. (2007), "Cache-, Hash- and Space-Efficient Bloom Filters", in Demetrescu, Camil, *Experimental Algorithms, 6th International Workshop, WEA 2007*  (PDF), Lecture Notes in Computer Science **4525**, Springer-Verlag, Lecture Notes in Computer Science 4525, pp. 108–121, doi:10.1007/978-3-540-72845-0 , ISBN 978-3-540-72844-3
- Sethumadhavan, Simha; Desikan, Rajagopalan; Burger, Doug; Moore, Charles R.; Keckler, Stephen W. (2003), "Scalable hardware memory disambiguation for high ILP processors", *36th Annual IEEE/ACM International Symposium on Microarchitecture, 2003. MICRO-36*  (PDF), pp. 399–410, doi:10.1109/MICRO.2003.1253244 , ISBN 0-7695-2043-X
- Shanmugasundaram, Kulesh; Brönnimann, Hené; Memon, Nasir (2004), "Payload attribution via hierarchical Bloom filters", *Proceedings of the 11th ACM Conference on Computer and Communications Security*, pp. 31–41, doi:10.1145/1030083.1030089 , ISBN 1581139616
- Starobinski, David; Trachtenberg, Ari; Agarwal, Sachin (2003), "Efficient PDA Synchronization", *IEEE Transactions on Mobile Computing* **2** (1): 40, doi:10.1109/TMC.2003.1195150
- Stern, Ulrich; Dill, David L. (1996), "A New Scheme for Memory-Efficient Probabilistic Verification", *Proceedings of Formal Description Techniques for Distributed Systems and Communication Protocols, and Protocol Specification, Testing, and Verification: IFIP TC6/WG6.1 Joint International Conference*, Chapman & Hall, IFIP Conference Proceedings, pp. 333–348, CiteSeerX 10.1.1.47.4101
- Haghighat, Mohammad Hashem; Tavakoli, Mehdi; Kharrazi, Mehdi (2013), "Payload Attribution via Character Dependent Multi-Bloom Filters", *Transaction on Information Forensics and Security*, *IEEE* **99** (5): 705, doi:10.1109/TIFS.2013.2252341
- Mitzenmacher, Michael; Upfal, Eli (2005), *Probability and computing: Randomized algorithms and probabilistic analysis* , Cambridge University Press, pp. 107–112, ISBN 9780521835404
- Mullin, James K. (1990), "Optimal semijoins for distributed database systems", *Software Engineering, IEEE Transactions on* **16** (5): 558–560, doi:10.1109/32.52778
- Rottenstreich, Ori; Kanizo, Yossi; Keslassy, Isaac (2012), "The Variable-Increment Counting Bloom Filter", *31st Annual IEEE International Conference on Computer Communications, 2012. Infocom 2012*  (PDF), pp. 1880–1888, doi:10.1109/INFCOM.2012.6195563 , ISBN 978-1-4673-0773-4

## External links

[ edit ]

- Why Bloom filters work the way they do (Michael Nielsen, 2012)
- Bloom Filters — A Tutorial, Analysis, and Survey (Blustein & El-Maazawi, 2002)  at Dalhousie University
- Table of false-positive rates for different configurations  from a University of Wisconsin–Madison website
- Interactive Processing demonstration  from ashcan.org
- "More Optimal Bloom Filters," Ely Porat (Nov/2007) Google TechTalk video  on YouTube
- "Using Bloom Filters"  Detailed Bloom Filter explanation using Perl
- "A Garden Variety of Bloom Filters"  - Explanation and Analysis of Bloom filter variants
- "Bloom filters, fast and simple"  - Explanation and example implementation in Python

Categories: Hashing | Probabilistic data structures | Lossy compression algorithms

This page was last modified on 2 September 2015, at 03:37.

Text is available under the Creative Commons Attribution-ShareAlike License; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy. Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit organization.

Privacy policy About Wikipedia Disclaimers Contact Wikipedia Developers Mobile view

