

# Una Implementación Eficiente en C++ y Análisis de Rendimiento de Filtros de Bloom

Jorge Carlos Burgos Rivero, Adiane Cueto Portuondo, Yusmelys González Saavedra, Elizabeth Molina Mena,  
Cinthya Beatriz Verde Días,

**Resumen**—La prueba de pertenencia en grandes conjuntos de datos es un problema fundamental en ciencias de la computación. Las estructuras deterministas tradicionales como las tablas hash ofrecen respuestas exactas, pero a un costo de memoria significativo, especialmente para flujos de datos masivos. Las estructuras de datos probabilistas ofrecen una alternativa al intercambiar una pequeña tasa de error controlable por ganancias sustanciales en eficiencia de espacio y tiempo. Este artículo presenta un análisis detallado del Filtro de Bloom (Bloom Filter), una prominente estructura de datos probabilista. Cubrimos sus fundamentos teóricos, bases matemáticas y aplicaciones del estado del arte. Posteriormente, detallamos el diseño e implementación de un Filtro de Bloom genérico y de alto rendimiento en C++. La implementación utiliza un esquema de doble hash con el algoritmo MurmurHash2 y un arreglo de bits eficiente en memoria. Un benchmark comprehensivo compara nuestra implementación contra una `std::unordered_set` estándar para inserciones, búsquedas positivas y búsquedas negativas. Los resultados demuestran que el Filtro de Bloom proporciona un rendimiento superior y una huella de memoria drásticamente menor, con una tasa de falsos positivos observada empíricamente que se alinea estrechamente con el objetivo teórico. Este trabajo valida al Filtro de Bloom como una solución óptima para aplicaciones donde la exactitud no es un requisito estricto.

**Index Terms**—Filtro de Bloom, Estructura de Datos Probabilista, Hashing, C++, Análisis de Rendimiento, Falsos Positivos.

## I. INTRODUCCIÓN

En la era del big data, determinar eficientemente si un elemento es miembro de un conjunto grande es una operación crítica en numerosos dominios, incluyendo bases de datos, seguridad de redes y sistemas distribuidos. Las estructuras de datos convencionales como los árboles binarios balanceados o las tablas hash proporcionan respuestas deterministas. Sin embargo, su consumo de memoria crece linealmente con el número de elementos almacenados, lo cual puede volverse prohibitivo al tratar con miles de millones de ítems.

Las estructuras de datos probabilistas ofrecen una alternativa convincente. Están diseñadas para responder consultas de forma aproximada, utilizando significativamente menos espacio que sus contrapartes deterministas. El Filtro de Bloom, introducido por Burton H. Bloom en 1970 [1], es una de las estructuras probabilistas más influyentes y utilizadas. Representa un conjunto usando un arreglo de bits compacto y múltiples funciones hash. El intercambio fundamental es su disposición a aceptar falsos positivos —indicar que un elemento está en el

conjunto cuando no lo está— a cambio de una garantía de no tener falsos negativos y ganancias sustanciales en eficiencia.

El objetivo de este artículo es triple. Primero, proporcionar una visión general de los Filtros de Bloom, desde su base conceptual hasta las fórmulas matemáticas que gobiernan su rendimiento. Segundo, presentar el diseño de una implementación robusta y eficiente en C++, detallando elecciones clave como la estrategia de hashing y las técnicas de manipulación de bits. Tercero, realizar un benchmark riguroso contra una tabla hash estándar (`std::unordered_set`) para cuantificar los beneficios prácticos en términos de velocidad y uso de memoria. Nuestros resultados confirman que para aplicaciones tolerantes a una pequeña tasa de error predecible, el Filtro de Bloom es una herramienta excepcionalmente poderosa.

El resto de este artículo está estructurado de la siguiente manera. La Sección II revisa el estado del arte y trabajos relacionados. La Sección III profundiza en los conceptos clave y los fundamentos matemáticos. La Sección IV describe el diseño e implementación de nuestro Filtro de Bloom en C++. La Sección V presenta los resultados experimentales y el análisis. Finalmente, la Sección VI concluye el artículo y discute posibles trabajos futuros.

## II. ESTADO DEL ARTE

Desde su creación, el Filtro de Bloom ha sido extensamente estudiado y adaptado. El artículo original de Bloom [1] sentó las bases, introduciendo el intercambio espacio/tiempo. La investigación posterior se centró en formalizar la probabilidad de falsos positivos y derivar parámetros óptimos. Las fórmulas clave para calcular el tamaño óptimo del arreglo de bits ( $m$ ) y el número de funciones hash ( $k$ ) basados en el número de elementos ( $n$ ) y la tasa de falsos positivos deseada ( $p$ ) son ahora conocimiento estándar en libros de texto [2].

Se han propuesto varias variantes del Filtro de Bloom clásico para abordar sus limitaciones, principalmente su incapacidad para soportar la eliminación de elementos. El *Filtro de Conteo* (Counting Bloom Filter) [3] reemplaza cada bit con un pequeño contador, permitiendo decrementos. Esto viene a costa de un aumento de memoria y una mayor probabilidad de desbordamiento del contador. Innovaciones más recientes incluyen el *Filtro Escalable* (Scalable Bloom Filter) [4], que adapta dinámicamente su tamaño para acomodar un número creciente de elementos, y el *Filtro Cociente* (Quotient Filter) [5], que ofrece una mejor localidad de caché y soporta eliminación con menos sobrecarga que un Filtro de Conteo.

Las aplicaciones de los Filtros de Bloom son vastas y diversas. En bases de datos, se utilizan para prevenir búsquedas

en disco innecesarias para claves no existentes en sistemas como Google Bigtable y Apache Cassandra [6]. En redes, se emplean para tareas como compartir caché, enrutamiento y detección de patrones de tráfico malicioso [7]. También son fundamentales en tecnología blockchain para simplificar la verificación de pagos y en bioinformática para acelerar la alineación de secuencias [8].

### III. CONCEPTOS CLAVE Y FUNDAMENTOS MATEMÁTICOS

Un Filtro de Bloom consta de dos componentes principales:

1. Un arreglo de bits de  $m$  bits, todos inicializados en 0.
2.  $k$  funciones hash independientes,  $h_1, h_2, \dots, h_k$ , cada una mapeando un elemento de entrada a un entero en el rango  $[0, m - 1]$ .

#### III-A. Operaciones Fundamentales

- **Add(item):** Para insertar un elemento, es procesado por todas las  $k$  funciones hash, obteniendo  $k$  índices. El bit en cada uno de estos índices en el arreglo se establece en 1.
- **Contains(item):** Para consultar un elemento, se calculan los mismos  $k$  índices. Si el bit en cualquiera de estos índices es 0, el elemento definitivamente no está en el conjunto. Si todos los  $k$  bits son 1, el elemento *probablemente* está en el conjunto. Un falso positivo ocurre cuando todos los  $k$  bits correspondientes han sido establecidos en 1 por otros elementos previamente insertados.

#### III-B. Probabilidad de Falsos Positivos

La probabilidad de un falso positivo es la métrica más crítica. Sea  $n$  el número de ítems insertados. Después de insertar un ítem, la probabilidad de que un bit específico siga siendo 0 es  $(1 - 1/m)^k$ . Después de insertar  $n$  ítems, esta probabilidad se convierte en  $(1 - 1/m)^{k \cdot n}$ . Para  $m$  grande, esto puede aproximarse usando la función exponencial ( $e^x$ ):

$$P(\text{bit es } 0) \approx e^{-k \cdot n / m} \quad (1)$$

Por lo tanto, la probabilidad de que un bit sea 1 es  $1 - e^{-k \cdot n / m}$ . Un falso positivo ocurre cuando todos los  $k$  bits para un elemento no existente son 1. Así, la tasa de falsos positivos  $p$  es:

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

#### III-C. Parámetros Óptimos

Dada una capacidad deseada  $n$  y una tasa de falsos positivos  $p$ , podemos resolver la Ecuación (2) para el  $m$  y  $k$  óptimos. Minimizando  $p$  con respecto a  $k$ , encontramos el número óptimo de funciones hash:

$$k = \frac{m}{n} \ln(2) \quad (3)$$

Sustituyendo este  $k$  óptimo de la Ecuación (3) de nuevo en la Ecuación (2) y resolviendo para  $m$ , obtenemos el tamaño óptimo del arreglo de bits:

$$m = -\frac{n \ln(p)}{(\ln(2))^2} \quad (4)$$

Estas fórmulas permiten que un Filtro de Bloom sea configurado con precisión para cumplir con requisitos específicos de rendimiento y precisión.

### IV. DISEÑO E IMPLEMENTACIÓN

Nuestra implementación, `SimpleBloomFilter`, es una clase plantilla en C++ diseñada para el rendimiento y la genericidad.

#### IV-A. Estructura de la Clase

La clase es una plantilla (`template<typename T>`) para soportar cualquier tipo de dato que pueda ser serializado en un arreglo de bytes. Las variables miembro clave son:

- `size_t m_size`: El tamaño calculado  $m$  del arreglo de bits.
- `size_t m_hash_count`: El número calculado de funciones hash  $k$ .
- `std::vector<uint64_t> m_bits`: El arreglo de bits en sí. Elegimos `std::vector<uint64_t>` sobre `std::vector<bool>` para mayor eficiencia. Las operaciones en palabras de 64 bits son típicamente más rápidas y tienen mejor localidad de caché. El número total de bloques de 64 bits es  $(m\_size + 63)/64$ .

El constructor toma `capacity (n)` y `error_rate (p)` como entrada y calcula `m_size` y `m_hash_count` usando las Ecuaciones (4) y (3), respectivamente.

#### IV-B. Estrategia de Hashing

Implementar  $k$  funciones hash distintas puede ser engorroso e ineficiente. Empleamos la técnica de "doble hash"[9], que genera  $k$  valores hash a partir de solo dos funciones hash independientes,  $h_1$  y  $h_2$ :

$$h_i(x) = h_1(x) + i \cdot h_2(x) \pmod{m} \quad (5)$$

Seleccionamos **MurmurHash2**, una función hash no criptográfica conocida por su excelente distribución y alto rendimiento. En nuestros métodos `add` y `contains`, generamos dos hashes base:

1.  $h_1 = \text{MurmurHash2}(\text{data}, 0)$
2.  $h_2 =$

Usar  $h_1$  como semilla para  $h_2$  asegura su independencia. Los  $k$  índices subsecuentes se calculan luego usando la fórmula anterior.

#### IV-C. Manipulación de Bits

Para manipular bits individuales dentro del `std::vector<uint64_t>`, implementamos dos métodos auxiliares privados:

- `set_bit(size_t index)`: Este método calcula el bloque de 64 bits ( $\text{index} / 64$ ) y la posición del

bit dentro del bloque ( $\text{index} \% 64$ ). Luego usa una máscara de bits y el operador OR a nivel de bits ( $|=$ ) para establecer el bit en 1.

- **get\_bit(size\_t index):** Este método calcula de manera similar el bloque y la posición, pero usa el operador AND a nivel de bits ( $&$ ) para verificar si el bit es 1, devolviendo un resultado booleano.

## V. RESULTADOS Y ANÁLISIS

Para evaluar nuestra implementación, realizamos un benchmark comparándola contra `std::unordered_set`. Las pruebas se realizaron en un sistema con una CPU Intel Core i7-10750H, 16GB de RAM, compilado con GCC 11.2 usando la bandera de optimización `-O3`.

### V-A. Configuración Experimental

- **Conjunto de Datos:**  $n = 500,000$  cadenas de texto únicas.
- **Tasa de Falsos Positivos Objetivo:**  $p = 0,01$  (1%).
- **Operaciones:**
  1. Inserción de todos los  $n$  elementos.
  2. Búsquedas positivas para todos los  $n$  elementos.
  3. Búsquedas negativas para un conjunto separado de  $n$  elementos.
- **Métricas:** Tiempo de ejecución (en milisegundos) y consumo de memoria.

### V-B. Huella de Memoria

Usando las fórmulas, los parámetros óptimos para  $n = 500,000$  y  $p = 0,01$  son  $m \approx 4,799,259$  bits y  $k \approx 7$ . Esto se traduce en un uso de memoria de aproximadamente 600 KB para el Filtro de Bloom. Una `std::unordered_set<std::string>` que almacene las mismas 500,000 cadenas requeriría significativamente más memoria (ej. asumiendo una longitud de cadena promedio de 15 bytes más sobrecarga, excedería fácilmente los 15-20 MB), destacando la eficiencia espacial del Filtro de Bloom.

### V-C. Resultados de Rendimiento

Los tiempos de ejecución se resumen en la Tabla I.

Cuadro I  
COMPARACIÓN DE RENDIMIENTO (EN MILISEGUNDOS)

Operación	Filtro de Bloom	<code>std::unordered_set</code>
Inserción	45 ms	120 ms
Búsquedas Positivas	35 ms	50 ms
Búsquedas Negativas	40 ms	55 ms

### V-D. Análisis

Los resultados demuestran claramente las ventajas de rendimiento del Filtro de Bloom.

- **Inserción:** El Filtro de Bloom es más de 2.5 veces más rápido. La operación `add` implica algunos cálculos hash

y simples asignaciones de bits, lo cual es computacionalmente más barato que la lógica compleja de una inserción en tabla hash (manejo de colisiones, posible rehashing y asignación de memoria para cadenas).

- **Búsquedas:** Tanto las búsquedas positivas como las negativas son más rápidas con el Filtro de Bloom. La operación `contains` es similarmente ligera, requiriendo solo cálculos hash y accesos a memoria a un arreglo de bits compacto, lo que resulta en una mejor localidad de caché en comparación con el diseño de memoria más disperso de una tabla hash.

Para las búsquedas negativas, el Filtro de Bloom reportó **5,089 falsos positivos**. Esto resulta en una **tasa de falsos positivos real del 1.0178 %**, la cual es notablemente cercana a nuestro objetivo del 1.0 %. Esto valida empíricamente la precisión del modelo matemático descrito en la Sección III. La `std::unordered_set`, como se esperaba, tuvo cero falsos positivos.

## VI. CONCLUSIÓN

Este artículo ha presentado un estudio comprensivo del Filtro de Bloom, desde su base teórica hasta una implementación práctica y de alto rendimiento en C++. Nuestra implementación aprovecha características modernas de C++ y algoritmos eficientes como el doble hash con MurmurHash2. Los resultados del benchmark muestran inequívocamente que el Filtro de Bloom ofrece un intercambio superior, proporcionando un consumo de memoria drásticamente menor y tiempos de inserción/consulta más rápidos en comparación con una tabla hash estándar. La tasa de falsos positivos medida empíricamente coincide estrechamente con la predicción teórica, confirmando la fiabilidad de su modelo matemático subyacente.

Para aplicaciones donde una pequeña probabilidad de error predecible es aceptable —como caché, prefiltrado y enrutamiento de red— el Filtro de Bloom es una herramienta indispensable. El trabajo futuro podría implicar implementar y evaluar variantes como el Filtro de Conteo o el Filtro Cociente para evaluar sus intercambios, o integrar esta estructura en una aplicación real más grande.

## REFERENCIAS

- [1] B. H. Bloom, “Space/Time Trade-offs in Hash Coding with Allowable Errors,” *Communications of the ACM*, vol. 13, no. 7, pp. 422–426, Jul. 1970.
- [2] M. Mitzenmacher y E. Upfal, *Probability and Computing: Randomized Algorithms and Probabilistic Analysis*. Cambridge University Press, 2005.
- [3] L. Fan, P. Cao, J. Almeida, y A. Z. Broder, “Summary Cache: A Scalable Wide-Area Web Cache Sharing Protocol,” *IEEE/ACM Transactions on Networking*, vol. 8, no. 3, pp. 281–293, Jun. 2000.
- [4] P. S. Almeida, C. Baquero, N. Preguiça, y D. Hutchison, “Scalable Bloom Filters,” *Information Processing Letters*, vol. 101, no. 6, pp. 255–261, 2007.
- [5] M. A. Bender et al., “The Quotient Filter: Approximate Membership in the Logarithmic Space,” *SIAM Journal on Computing*, vol. 45, no. 1, pp. 254–294, 2016.
- [6] A. Lakshman y P. Malik, “Cassandra: A Decentralized Structured Storage System,” *ACM SIGOPS Operating Systems Review*, vol. 44, no. 2, pp. 35–40, Abr. 2010.
- [7] S. Geravand y M. Ahmadi, “A Survey on Bloom Filter-Based Structures for Network and Security Applications,” *Journal of Network and Computer Applications*, vol. 183, 2021.

- [8] M. C. Schatz, "Cloud-based Bioinformatics: Google Maps for Genomes," *Nature Biotechnology*, vol. 27, no. 12, pp. 1079–1080, Dic. 2009.
- [9] A. Kirsch y M. Mitzenmacher, "Less Hashing, Same Performance: Building a Better Bloom Filter," en *Algorithms - ESA 2006*, 2006, pp. 456–467.