

Bloom Filter

Jambura Anna, Pürstinger Kathrin,
Schnappauf Franziska, Thiele Coco

26. Februar 2026

Zusammenfassung

passend auszufüllen

6 1 Grundlagen und Motivation

7 1.1 Das Membership-Problem

8 Seien ein beliebiges Element x und eine Menge S gegeben. Das Membership-Problem ist
9 eine Bezeichnung für die Fragestellung: „Ist das Element x Teil der Menge S ?“ Diese Frage
10 tritt in vielen verschiedenen Bereichen und Anwendungen auf. Einige Beispiele dafür
11 sind Datenbankenabfragen und URL-Caching in Web-Browsern. Klassische Ansätze, wie
12 Listen, Hashtabellen oder Suchbäume liefern eine exakte Antwort auf die Frage, jedoch
13 benötigen sie alle entsprechend viel Zeit und Speicherplatz. [1]

14 1.2 Lösungsansatz - Bloomfilter

15 Bloomfilter wurden 1970 von Burton H. Bloom entwickelt, um den hohen Ressourcen-
16 bedarf zu umgehen. Sie sind probabilistische Datenstrukturen, das bedeutet sie arbeiten
17 mit Wahrscheinlichkeiten anstatt absoluter Sicherheit.

18 Dabei erlauben sie False-Positives in einem begrenzten Ausmaß. Ein Filter kann also
19 fälschlicherweise melden, das Element x sei Teil der Menge S , auch wenn dies nicht der
20 Fall ist. Umgekehrt sind False-Negatives jedoch ausgeschlossen. Wenn x tatsächlich ein
21 Element von S ist, wird das der Filter immer korrekt erkennen. Mit anderen Worten:
22 Ein vorhandenes Element wird nie als „nicht vorhanden“ gemeldet. [2]

23 1.3 Trade-off

24 Bloomfilter balancieren drei zentrale Faktoren. Neben der Reject-Time (Zeit zur Ab-
25 lehnung von Nicht-Mitgliedern) und dem benötigten Speicherplatz, die auch in konven-
26 tionellen Hashing-Methoden berücksichtigt werden müssen, wird hier auch die erlaubte
27 Fehlerrate betrachtet. Der zentrale Trade-off ist dabei zwischen dem akzeptablen Anteil
28 an False-Positives und der Speichereffizienz. Dieser ist bei der Implementierung eines
29 Bloomfilters individuell konfigurierbar.

30 Durch die kontrollierte Fehlerwahrscheinlichkeit wird der Speicherbedarf bedeutend re-
31 duziert, da er nicht von der Länge der Daten abhängt, sondern immer gleich viele Bits
32 pro Element beträgt. Je niedriger die Fehlerrate gewählt ist, desto mehr Bits pro Element
33 werden benötigt. Bloomfilter sind besonders hilfreich, wenn die Mehrheit der Anfragen
34 nicht-existente Elemente betrifft – hier liefern sie schnell ein definitives „Nein“ auf die
35 Membership-Frage. [1]

2 Funktionsweise und Mathematische Grundlagen

2.1 Aufbau

Der Bloomfilter besteht auf einem m -stelligen Bitarray, welches initial mit Nullen befüllt wird. Weiters werden k unabhängige Hashfunktionen definiert. Diese verwendet man um die Elemente der gewünschten Menge zu hashen. Abhängig von ihrem Hashwert werden die Elemente dann an der entsprechenden Position im Array eingefügt. Um also jedes Element erfolgreich einzufügen, muss die Hashfunktion $\text{mod } m$ angewandt werden. Somit erreicht man die Indizes 0 bis $m - 1$. [2]

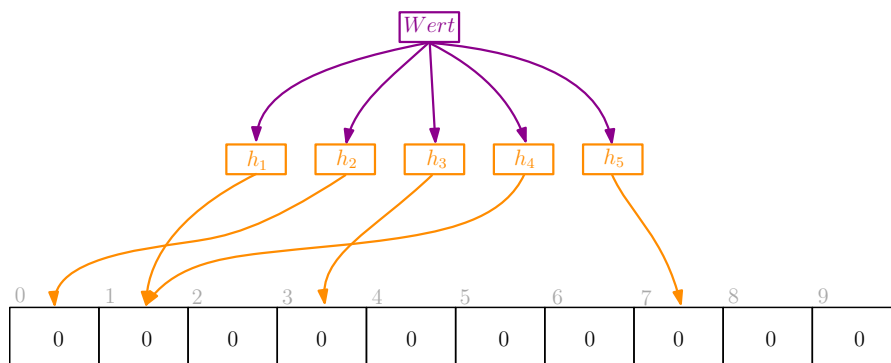


Abbildung 1: Visualisierung eines Bloomfilters

Da die Hashfunktionen keinem Sicherheitsstandard entsprechen, müssen keine kryptographischen Eigenschaften gelten. Kryptographische Eigenschaften bedeutet, minimale Eingabeänderungen müssen zu einer maximalen Änderung des Hashwerts führen. Die Eingabe darf nicht mittels der Hashfunktion wiederhergestellt werden können und zwei Eingaben haben fast unmöglich den selben Hashwert. Für Bloomfilter verwendet man schnelle und einfache Hashfunktionen, da die Effizienz im Vordergrund steht.

2.2 Einfügen/Suchen

Einfügen

Eine Menge S wird nun wie folgt in einem Bloomfilter eingefügt:
Für jedes Element $x \in S$ werden die Hash Werte aller k Hash Funktionen berechnet. Nun wird an diesen Positionen im Array die 0 auf eine 1 gesetzt. Sollte an einer dieser Positionen bereits eine 1 stehen, wird dies ignoriert. Dieser Vorgang wird für alle n Elemente der Menge S wiederholt.

58 2.2.1 Beispiel Einfügen

59 Betrachte folgende Menge $S = \{2, 4, 9\}$ und einen Bloomfilter der Länge $m = 10$ mit
 60 $k = 3$ Hashfunktionen.

61 Als beispielhafte Hashfunktionen verwenden wir: $h_1(x) = x \bmod 10$ $h_2(x) = (2x+3) \bmod$
 62 10 und $h_3(x) = (3x + 7) \bmod 10$.

63 Nun berechnen wir die Hashwerte für jedes Element der Menge S :

- 64 • Für $x = 2$:

$$h_1(2) = 2 \bmod 10 = 2$$

$$h_2(2) = (2 \cdot 2 + 3) \bmod 10 = 7$$

$$h_3(2) = (3 \cdot 2 + 7) \bmod 10 = 3$$

- 65 • Für $x = 4$:

$$h_1(4) = 4 \bmod 10 = 4$$

$$h_2(4) = (2 \cdot 4 + 3) \bmod 10 = 1$$

$$h_3(4) = (3 \cdot 4 + 7) \bmod 10 = 9$$

- 66 • Für $x = 9$:

$$h_1(9) = 9 \bmod 10 = 9$$

$$h_2(9) = (2 \cdot 9 + 3) \bmod 10 = 1$$

$$h_3(9) = (3 \cdot 9 + 7) \bmod 10 = 4$$

67 Nun fügt man die Elemente in den Bloomfilter ein. Für das erste Element 2 werden die
 68 Positionen 2, 7 und 3 auf 1 gesetzt. Daraus resultiert der folgende Bloomfilter:

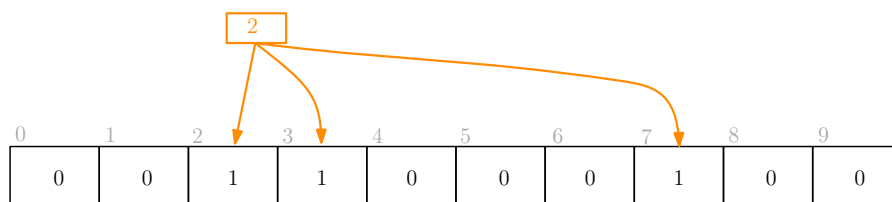


Abbildung 2: Bloomfilter nach Einfügen des Elements 2

69

70 Für das zweite Element 4 werden die Positionen 4, 1 und 9 auf 1 gesetzt.

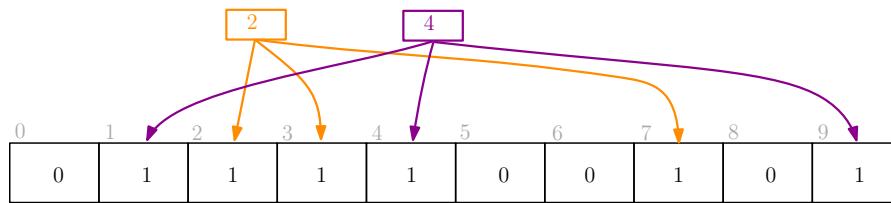


Abbildung 3: Bloomfilter nach Einfügen des Elements 4

71 Für das dritte Element 9 werden die Positionen 9, 1 und 4 auf 1 gesetzt. Da die Positionen
 72 1, 4 und 9 bereits auf 1 gesetzt wurden, ändert sich der Bloomfilter nicht weiter.

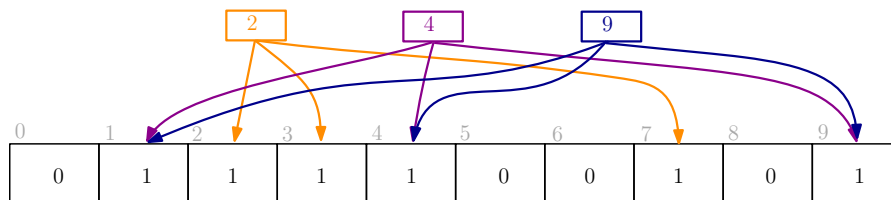


Abbildung 4: Bloomfilter nach Einfügen des Elements 9

73

74 Suchen

75 Um ein Element x in einem Bloomfilter zu suchen, werden dieselben Hashfunktionen wie
 76 beim Einfügen verwendet. Die Hashwerte werden berechnet und an den entsprechenden
 77 Positionen im Array geprüft. Wenn alle Positionen auf 1 gesetzt sind, so ist das Element
 78 wahrscheinlich in der Menge enthalten. Wenn mindestens eine Position auf 0 gesetzt ist,
 79 so ist das Element sicher nicht in der Menge enthalten. [2]

80 2.2.2 Beispiel Suchen

81 Betrachten wir den zuvor erstellten Bloomfilter und suchen nach dem Element 4. Be-
 82 rechnen wir die Hashwerte für 4:

$$h_1(4) = 4 \bmod 10 = 4$$

$$h_2(4) = (2 \cdot 4 + 3) \bmod 10 = 1$$

$$h_3(4) = (3 \cdot 4 + 7) \bmod 10 = 9$$

83 Nun prüfen wir die Positionen 4, 1 und 9 im Bloomfilter. Alle drei Positionen sind auf
 84 1 gesetzt, daher ist das Element 4 wahrscheinlich in der Menge enthalten.

85 Betrachten wir nun das Element 5 und berechnen die Hashwerte:

$$h_1(5) = 5 \bmod 10 = 5$$

$$h_2(5) = (2 \cdot 5 + 3) \bmod 10 = 3$$

$$h_3(5) = (3 \cdot 5 + 7) \bmod 10 = 2$$

86 Nun prüfen wir die Positionen 5, 3 und 2 im Bloomfilter. Die Position 2 ist auf 0 gesetzt,
87 daher ist das Element 5 sicher nicht in der Menge enthalten.

88 Ein wichtiger Aspekt des Bloomfilters ist, dass er fälschlicherweise angeben kann, dass
89 ein Element in der Menge enthalten ist, obwohl es tatsächlich nicht vorhanden ist. Dies
90 wird als *False Positive* bezeichnet. Wenn alle Positionen, die durch die Hashfunktionen
91 eines Elements angegeben werden, auf 1 gesetzt sind, obwohl das Element nicht in der
92 Menge enthalten ist, führt dies zu einem False Positive. Ein Beispiel hierfür wäre das
93 Element 12:

$$h_1(12) = 12 \bmod 10 = 2$$

$$h_2(12) = (2 \cdot 12 + 3) \bmod 10 = 7$$

$$h_3(12) = (3 \cdot 12 + 7) \bmod 10 = 3$$

94 Die Positionen 2, 7 und 3 sind alle auf 1 gesetzt, obwohl das Element 12 nicht in der
95 Menge enthalten ist. Daher würde der Bloomfilter fälschlicherweise angeben, dass 12 in
96 der Menge enthalten ist.

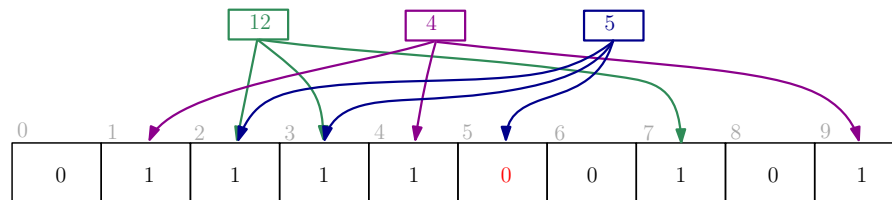


Abbildung 5: Suchen nach den Elementen 4, 5 und 12 im Bloomfilter

97 2.3 Formeln zur Evaluierung

98 2.3.1 False Positive Probability

99 Zur Erstellung des optimalen Bloomfilter ist es wichtig, die False Positive Probability
100 (FPP) zu berechnen. Diese gibt an, wie wahrscheinlich es ist, dass der Bloomfilter fälsch-
101 licherweise angibt, dass ein Element in der Menge enthalten ist, obwohl es tatsächlich
102 nicht vorhanden ist. Laut [2] entsteht die Formel zur Berechnung aus folgenden Kompo-
103 nenten:

104 Unter der Annahme, dass die Hashfunktionen unabhängig und gleichverteilt sind, ergibt
 105 sich die Wahrscheinlichkeit dass ein bestimmtes der m Bits nicht gesetzt ist durch:

$$1 - \frac{1}{m} \quad (1)$$

106 Weiters werden nun die k Hashfunktionen mitbetrachtet, immer noch für den Fall, dass
 107 ein bestimmtes Bit nicht gesetzt ist.

$$\left(1 - \frac{1}{m}\right)^k \quad (2)$$

108 Nun werden die n Elemente der Menge S betrachtet, welche in den Bloomfilter eingefügt
 109 werden.

$$\left(1 - \frac{1}{m}\right)^{kn} \quad (3)$$

110 Die Wahrscheinlichkeit, dass ein bestimmtes Bit auf 1 gesetzt ist, ergibt sich aus der
 111 Gegenwahrscheinlichkeit:

$$1 - \left(1 - \frac{1}{m}\right)^{kn} \quad (4)$$

112 Da es bei Bloomfiltern um Membership-Tests geht, muss die Wahrscheinlichkeit berech-
 113 net werden, dass alle k Positionen eines Elements auf 1 gesetzt sind, obwohl das Element
 114 nicht in der Menge enthalten ist.

$$\left(1 - \left(1 - \frac{1}{m}\right)^{kn}\right)^k \quad (5)$$

115 Aus der Formel lässt sich schließen, dass je **größer** m gewählt wird, desto **kleiner** wird
 116 die False Positive Probability. Je **größer** n gewählt wird, desto **größer** wird die False
 117 Positive Probability.

118 Da die False Positive Probability so klein wie möglich gehalten werden soll, ist auch
 119 die Wahl der Anzahl Hashfunktionen von großer Bedeutung. Setzt man die Formel für
 120 die False Positive Probability gleich 0 und löst sie nach k auf, erhält man die optimale
 121 Anzahl an Hashfunktionen:

$$k_{opt} = \frac{m}{n} \ln 2 \approx \frac{9m}{13n} \quad (6)$$

122 3 Pseudocode und Implementierung

123 Ein Bloom-Filter lässt sich mit drei grundlegenden Operationen beschreiben: Initialisie-
124 rung, Einfügen und Abfragen.

125 3.1 Initialisierung

126 Bei der Initialisierung werden alle m Bits im Array auf 0 gesetzt und k Hash-Funktionen
127 festgelegt. Je kleiner die gewünschte Fehlerrate ist, desto größer muss m gewählt wer-
128 den.

Algorithm 1 Initialisierung eines Bloom-Filters

- 1: Erzeuge Bit-Array $B[0 \dots m - 1]$ und setze alle Bits auf 0
 - 2: Definiere Hash-Funktionen h_1, h_2, \dots, h_k
-

129 3.2 Einfügen

130 Beim Einfügen wird für jedes Element x eine Schleife genau k -mal ausgeführt. In jeder
131 Iteration wird mithilfe der jeweiligen Hash-Funktion h_i ein Index berechnet und das
132 entsprechende Bit im Bitarray auf 1 gesetzt. Der Modulo-Operator stellt sicher, dass der
133 berechnete Index immer innerhalb des gültigen Bereichs von 0 bis $m - 1$ liegt. [2]

Algorithm 2 Einfügen eines Elements x

- 1: **for** $i = 1$ **to** k **do**
 - 2: $index \leftarrow h_i(x) \bmod m$
 - 3: $B[index] \leftarrow 1$
 - 4: **end for**
-

134 3.3 Abfragen

135 Für eine Abfrage werden dieselben k Hash-Werte berechnet und die entsprechenden
136 Positionen im Array überprüft. Existiert mindestens ein Bit mit dem Wert 0, kann man
137 mit absoluter Sicherheit sagen, dass das Element nicht enthalten ist – es gibt keine
138 False Negatives. Sind hingegen alle k Bits gleich 1, gilt das Element als wahrscheinlich
139 enthalten. Diese probabilistische Aussage ist das zentrale Merkmal des Bloom-Filters:
140 Es sind False Positives möglich.

Algorithm 3 Abfrage eines Elements x

```
1: for  $i = 1$  to  $k$  do
2:    $index \leftarrow h_i(x) \bmod m$ 
3:   if  $B[index] = 0$  then
4:     return FALSE
5:   end if
6: end for
7: return TRUE
```

4 Komplexitätsanalyse

4.1 Zeitkomplexität

Sowohl das Einfügen als auch das Abfragen eines Elements haben eine Zeitkomplexität von $\mathcal{O}(k)$, wobei k die Anzahl der Hash-Funktionen bezeichnet. Entscheidend ist dabei, dass diese Zeit *unabhängig* von der Anzahl n der bereits im Filter gespeicherten Elemente ist. Der Grund dafür liegt in der Struktur des Filters: Es werden keine Elemente explizit gespeichert, sondern lediglich Bits in einem Array der Größe m gesetzt oder gelesen. Egal ob sich 1.000 oder 100 Millionen Elemente im Filter befinden – die Abfragezeit bleibt konstant [1].

4.2 Speicherkomplexität

Die Speicherkomplexität beträgt $\mathcal{O}(m)$, wobei m die Größe des Bit-Arrays ist. Im Gegensatz zu klassischen Datenstrukturen hängt dieser Speicherbedarf *nicht* von der Größe der gespeicherten Elemente ab, sondern nur von zwei Faktoren: der Anzahl der zu speichernden Elemente n und der akzeptierten False-Positive-Rate ε [1].

Als praktische Faustregel gilt: Bei einer Fehlerrate von etwa 1 % benötigt ein Bloom-Filter weniger als 10 Bits pro Element. Das ist bemerkenswert effizient – unabhängig davon, ob es sich bei den Elementen um kurze Zeichenketten oder lange URLs handelt [3].

4.3 Vergleich mit anderen Datenstrukturen

Tabelle 1 stellt die Komplexitätseigenschaften des Bloom-Filters denen einer Hash-Tabelle mit verketteten Listen sowie eines balancierten Baums gegenüber.

Die **Hash-Tabelle mit verketteten Listen** erreicht im Durchschnitt $\mathcal{O}(1)$ für Einfüge- und Suchoperationen, kann im schlechtesten Fall jedoch auf $\mathcal{O}(n)$ anwachsen. Da jedes Element explizit gespeichert wird, beträgt der Speicherbedarf $\mathcal{O}(n)$ – typischerweise 64 Bits oder mehr pro Element (Nutzdaten plus Pointer). Der wesentliche Vorteil liegt

Eigenschaft	Bloom-Filter	Hash-Tabelle mit Chaining	Balancierter Baum
Zeitkomplexität	$\mathcal{O}(k)$	$\mathcal{O}(1)$, worst $\mathcal{O}(n)$	$\mathcal{O}(\log n)$
Speicherkomplexität	$\mathcal{O}(m)$	$\mathcal{O}(n)$	$\mathcal{O}(n)$
Genauigkeit	Probabilistisch	Exakt	Exakt

Tabelle 1: Vergleich ausgewählter Datenstrukturen

in der Exaktheit: Es gibt keine False Positives, und Elemente können jederzeit wieder abgerufen werden [4].

Der **balancierte Baum** hat eine Zeitkomplexität von $\mathcal{O}(\log n)$ für Suche und Einfügen. Die Laufzeit steigt mit wachsender Elementanzahl langsam an, da bei jedem Schritt etwa die Hälfte der verbleibenden Elemente verworfen wird. Der Speicherbedarf ist ebenfalls $\mathcal{O}(n)$. Der Vorteil liegt in der Möglichkeit, Elemente geordnet zu speichern, was zusätzliche Operationen wie Bereichsabfragen erlaubt [4].

4.4 Speichereffizienz in der Praxis

Um die Speicherersparnis greifbar zu machen, betrachten wir ein konkretes Beispiel: Für 100 Millionen URLs benötigt ein Bloom-Filter bei einer Fehlerrate von 1 % rund 120 Megabyte. Eine Hash-Tabelle mit denselben Einträgen würde hingegen über ein Gigabyte beanspruchen. Das ist nicht nur ein quantitativer, sondern oft ein qualitativer Unterschied – nämlich der zwischen einem System, das auf einem Endgerät lauffähig ist, und einem, das einen dedizierten Server erfordert.

Dieser enorme Vorteil hat allerdings seinen Preis: Ein Bloom-Filter beantwortet ausschließlich die Frage „*Ist das Element möglicherweise in der Menge?*“ Er kann weder Elemente aufzählen noch löschen, noch gibt er die Elemente selbst zurück [2].

183 5 Probleme von Bloom-Filtern und Lösungen

184 5.1 Das Löschen von Elementen

185 Der klassische Bloom-Filter besitzt unter anderem die Einschränkung, dass er das Lö-
186 schen von Elementen nicht unterstützt. Möchte man ein Element entfernen, liegt es
187 zunächst nahe, die entsprechenden Bits im Bit-Array wieder auf 0 zu setzen. Genau
188 hier entsteht jedoch ein fundamentales Problem. Mehrere Elemente können auf dieselbe
189 Position im Bit-Array hashen. Wird ein Bit zurückgesetzt, entfernt man daher nicht nur
190 das gewünschte Element, sondern gleichzeitig auch alle anderen Elemente, die an die-
191 ser Position gespeichert wurden. Das eigentliche Element ist zwar entfernt, aber andere
192 Elemente gelten nun ebenfalls als nicht mehr vorhanden, obwohl sie eigentlich noch im
193 Filter sein sollten.

194 Eine Lösung für dieses Problem ist der Counting Bloom Filter. [5] Das Grundprinzip
195 beim Einfügen bleibt dabei gleich wie beim klassischen Bloom-Filter. Der Unterschied
196 besteht darin, dass man an jeder Position nicht nur ein einzelnes Bit speichert, sondern
197 einen kleinen Zähler. Dieser Zähler wird beim Einfügen eines Elements um 1 erhöht und
198 beim Löschen wieder um 1 verringert. Typischerweise sind diese Zähler 4 Bit groß und
199 können somit Werte von 0 bis 15 speichern. Dadurch wird es möglich, Elemente sicher
200 zu löschen, ohne andere Einträge unbeabsichtigt zu beeinflussen.

201 Allerdings hat der Counting Bloom Filter auch Nachteile. Der Speicherverbrauch ist
202 deutlich höher, da statt eines einzelnen Bits nun 4 Bits pro Position benötigt werden.
203 Bei gleicher Genauigkeit benötigt ein Counting Bloom Filter somit ungefähr das Drei-
204 bis Vierfache an Speicher im Vergleich zum klassischen Bloom-Filter. Außerdem können
205 die Zähler überlaufen. Wenn mehr als 15 Elemente auf dieselbe Position hashen, reichen
206 4 Bits nicht mehr aus. Man könnte größere Zähler verwenden, allerdings würde das den
207 Speicherbedarf weiter erhöhen. Der Counting Bloom Filter eignet sich daher besonders
208 dann, wenn häufig gelöscht werden muss - man bezahlt diese Möglichkeit jedoch mit
209 einem deutlich höheren Speicherverbrauch. An Verbesserungen wird zwar gearbeitet,
210 doch eine genauere Betrachtung würde den Rahmen dieser Arbeit sprengen.[6]

211 5.2 Größenplanung

212 Ein weiteres grundlegendes Problem klassischer Bloom-Filter ist die Größenplanung. In
213 der Regel muss man vorher festlegen, wie groß der Filter sein soll. Ist er zu klein dimen-
214 sioniert, steigt die Fehlerwahrscheinlichkeit stark an. Die Bits werden sehr schnell gesetzt
215 und die False-Positive-Rate nimmt deutlich zu. Ist der Filter hingegen zu groß gewählt,
216 wird Speicherplatz verschwende, da möglicherweise Kapazitäten reserviert werden, die
217 nie vollständig genutzt werden.

218 Der Scalable Bloom Filter bietet hier eine Lösung durch dynamisches Wachstum. [5] Er
219 besteht aus mehreren klassischen Bloom Filtern, die nacheinander erstellt werden. Sobald

ein Filter eine bestimmte Auslastung erreicht, wird ein neuer, größerer Filter mit einer strengeren Fehlerrate hinzugefügt. Auf diese Weise bleibt die Gesamtfehlerwahrscheinlichkeit über alle Filter hinweg kontrollierbar. Selbst wenn mehrere Filter hinzukommen, bleibt die kombinierte Fehlerrate in akzeptablen Grenzen. Der große Vorteil ist, dass der Filter beliebig wachsen kann, ohne komplett neu aufgebaut werden zu müssen. Ein Nachteil ist jedoch, dass Abfragen mit jedem zusätzlichen Filter etwas langsamer werden, da mehrere Filter überprüft werden müssen. Neben Counting- und Scalable-Varianten gibt es noch viele weitere spezielle Varianten von Bloom-Filtern, die jedoch den Rahmen dieser Arbeit überschreiten würden.[7]

6 Cuckoo Filter

Eine alternative Datenstruktur stellt der Cuckoo Filter dar. Hierbei handelt es sich nicht mehrh wirklich um einen Bloom-Filter, dennoch verfolgt er dasselbe Ziel: speichereffiziente Mengenabfragen bei geringen Fehlerraten. Der Cuckoo Filter basiert nicht auf einem Bit-Array, sondern auf einer Hash-Tabelle mit kleinen Fächern, sogenannten Buckets. In diesen Buckets werden Fingerabdrücke, sogenannte Fingerprints, gespeichert. Das sind kurze, eindeutige Kennungen der Elemente mit nur wenigen Bits Länge.

Beim Einfügen eines Elements wird mithilfe einer Hash-Funktion berechnet, in welches Fach es gehört. Jedes Element besitzt dabei genau zwei mögliche Buckets, in denen es abgelegt werden kann. Ist in einem dieser Buckets noch Platz vorhanden, wird der Fingerprint dort gespeichert. Sind jedoch beide Buckets belegt, greift das sogenannte Cuckoo-Prinzip. Hier verdrängt das neue Element einen bestehenden Eintrag aus einem der beiden Buckets. Das verdrängte Element muss sich anschließend einen neuen Platz in seinem alternativen Bucket suchen. Dieser Prozess kann sich fortsetzen, bis schließlich alle Elemente einen Platz gefunden haben.

Der Cuckoo Filter bringt sowohl Vorteile als auch Nachteile mit sich. Ein großer Vorteil ist, dass Elemente problemlos gelöscht werden können, da die Fingerprints direkt gespeichert sind und gezielt entfernt werden können. Außerdem sind Abfragen sehr schnell, da nur zwei Buckets geprüft werden müssen. Ein Nachteil zeigt sich bei sehr hoher Auslastung der Hash-Tabelle. In solchen Fällen kann die Verdrängungskette sehr lang werden, ohne dass ein freier Platz gefunden wird. Dann muss die gesamte Struktur vergrößert werden. Studien zeigen jedoch, dass Cuckoo Filter in vielen realen Anwendungen praktisch besser abschneiden als klassische Bloom-Filter.[8] Klassische Bloom-Filter sind dennoch besonders sinnvoll, wenn sehr große Datenmengen verarbeitet werden, der verfügbare Speicher knapp oder teuer ist, kleine Fehlerraten akzeptiert werden können und sie als Vorfilter von aufwendigen oder rechenintensiven Operationen eingesetzt werden. [4]

255 7 Anwendungsbeispiele

256 7.1 Web-Proxy-Caching

257 In verteilten Netzwerken arbeiten mehrere Proxy-Server zusammen und tauschen sich
258 untereinander aus. Bei einer Anfrage nach einer Webseite sucht ein Proxy zunächst im
259 eigenen Cache, ob er diese bereits gespeichert hat. Wenn das nicht der Fall ist, spricht
260 man von einem Cache-Miss und es wird geprüft, ob sich die Webseite im Cache eines
261 anderen Proxys befindet. Wird sie hier gefunden, wird die Anfrage an den entsprechenden
262 Proxy weitergeleitet, anstatt die Seite direkt aus dem Web zu laden.

263 Damit dieses System funktioniert, muss jeder Proxy über den Inhalt der Caches aller
264 anderen Proxies Bescheid wissen. Um den enormen Netzwerkverkehr, der beim wieder-
265 holten Austausch der kompletten URL-Listen entstehen würde, zu vermeiden, kommen
266 hier Bloomfilter zum Einsatz. Im Summary Cache Protokoll tauschen Proxies periodisch
267 Bloomfilter untereinander aus, die den Inhalt ihres Caches zusammenfassen. Wenn nun
268 ein Cache-Miss auftritt, werden die Bloomfilter jener anderen Proxies konsultiert, die
269 ein positives Ergebnis versprechen und die Anfrage wird entsprechend weitergeleitet.

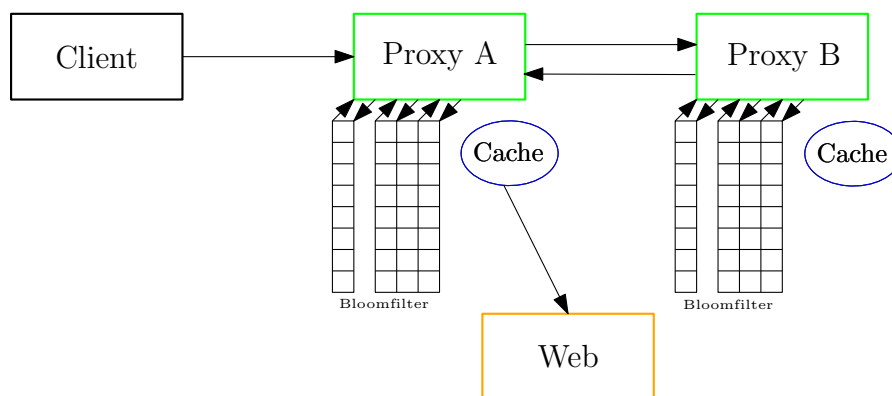


Abbildung 6: Web-Proxy-Caching mit Bloomfiltern

270 Hierbei können False-Positives auftreten, wobei es dann zu einer minimalen Verzögerung
271 kommt. Die massive Reduktion des Netzwerkverkehrs durch den Bloomfilter überwiegt
272 diesen Nachteil bei Weitem. Das Summary Cache Protokoll wird beispielsweise im Web-
273 Proxy-Cache „Squid“ eingesetzt. [4]

274 7.2 Google Bigtable

275 Bloomfilter werden oft in Datenbanksystemen verwendet, wobei Google Bigtable ein
276 bekanntes Beispiel hierfür ist. Bigtable speichert die Daten auf der Festplatte in Sorted-
277 String-Tables (SSTables). Wenn eine Leseoperation durchgeführt werden soll, müssen

278 potenziell mehrere dieser Tables durchsucht werden, bis die gewünschten Daten gefun-
279 den werden. Da jede Table auf der Festplatte liegt, verursacht jeder Zugriff auf eine
280 SSTable auch einen teuren Festplattenzugriff. Besonders problematisch im Bezug auf
281 die benötigten Ressourcen wird dies bei Abfragen nach nicht-existenten Daten.

282 Kommen jetzt die Bloomfilter zum Einsatz, ändert sich dies drastisch. Für jede SSTable
283 wird ein Bloomfilter im Hauptspeicher gehalten, der Auskunft über deren Inhalt gibt.
284 Vor einem Festplattenzugriff wird also der Filter befragt, ob die gesuchten Daten in
285 der Table enthalten sind. Bei einem positiven Ergebnis wird der Zugriff durchgeführt,
286 ansonsten kann er eingespart werden. [9]

287 7.2.1 Beispiel Anfrage

288 Angenommen es wird eine Anfrage auf den Schlüssel X gestellt und auf der Festplatte
289 liegen drei SSTables. Ohne Verwendung von Bloomfiltern müssten alle drei Tables abge-
290 rufen und durchsucht werden, also drei Festplattenzugriffe durchgeführt werden. Unter
291 Einsatz von Bloomfiltern werden jedoch zuerst diese konsultiert. Die ersten beiden Fil-
292 ter könnten melden, dass Schlüssel X jeweils nicht in SSTable 1 bzw. SSTable 2 liegt,
293 sie können also beide übersprungen werden. Filter 3 sagt jetzt, dass sich X in Table 3
294 befinden könnte – dieser Zugriff wird durchgeführt. Demzufolge wurde nur ein Festplat-
295 tenzugriff durchgeführt, bis der gesuchte Schlüssel X gefunden wurde, das bedeutet eine
296 Ersparnis von zwei Zugriffen durch die Verwendung von Bloomfiltern.

297 7.3 Weitere Anwendungen

298 Heute kommen Bloomfilter in zahlreichen Systemen zum Einsatz. Google Chrome nutzt
299 sie beispielsweise für Safe-Browsing zur Malware-Erkennung. [10] Neben Google Bigta-
300 ble setzen auch weitere Datenbanksysteme, wie Apache Cassandra auf die Vorteile von
301 Bloomfiltern, um unnötige Festplattenzugriffe zu vermeiden. [11]

Literatur

- [1] Burton H. Bloom. Space/time trade-offs in hash coding with allowable errors. *Commun. ACM*, 13(7):422–426, 1970.
- [2] Sasu Tarkome, Christian Esteve Rothenberg, and Eemil Lagerspetz. Theory and practice of bloom filters for distributed systems. *IEEE Communications Surveys & Tutorials*, 14(1):131–155, 2012.
- [3] Li Fan, Pei Cao, Jussara Almeida, and Andrei Z Broder. Summary cache: A scalable wide-area web cache sharing protocol. *IEEE/ACM Transactions on Networking*, 8(3):281–293, 2000.
- [4] Andrei Broder and Michael Mitzenmacher. Network applications of bloom filters: A survey. *Internet Mathematics*, 1(4):485–509, 2004.
- [5] Paul A. Gagniuc, Ionel-Bujorel Păvăloiu, and Maria-Iuliana Dascălu. Bloom filters at fifty: From probabilistic foundations to modern engineering and applications. *Algorithms*, 18:1–15, 2025.
- [6] Flavio Bonomi, Michael Mitzenmacher, Rina Panigrahy, Sushil Singh, and George Varghese. An improved construction for counting bloom filters. In *Proceedings of the 14th conference on Annual European Symposium on Algorithms*, pages 684–695. Springer, 2006.
- [7] Paulo S. Almeida, Carlos Baquero, Nuno Preguiça, and David Hutchison. Scalable bloom filters. *Information Processing Letters*, 101(6):255–261, 2007.
- [8] Bin Fan, David G. Andersen, Michael Kaminsky, and Michael D. Mitzenmacher. Cuckoo filter: Practically better than bloom. In *Proceedings of the 10th ACM International Conference on Emerging Networking Experiments and Technologies*, pages 179–190. ACM, 2014.
- [9] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach, Mike Burrows, Tushar Chandra, Andrew Fikes, and Robert E. Gruber. Bigtable: A distributed storage system for structured data. *ACM Trans. Comput. Syst.*, 26(2), 2008.
- [10] Thomas Gerbet, Amrit Kumar, and Cédric Lauradoux. (un)safe browsing. Technical Report RR-8594, INRIA, 2010.
- [11] Avinash Lakshman and Prashant Malik. Cassandra: a decentralized structured storage system. *SIGOPS Oper. Syst. Rev.*, 44(2):35–40, 2010.