

# Hashing 2: Saving Space with Bloom Filters

---



**Rasmus Resen Amossen**

SOLUTION ARCHITECT

[rasmus.resen.org](http://rasmus.resen.org)

# The Match Finder App

Hashing

Hash functions

Exact tables

Probabilistic filters

My Pictures



Spatial index trees

Disjoint-Set structures

Tries

Suffix trees

ing

Cooking

Fishing

Ingredients

King Kong

Kings

Cooking

Guitar

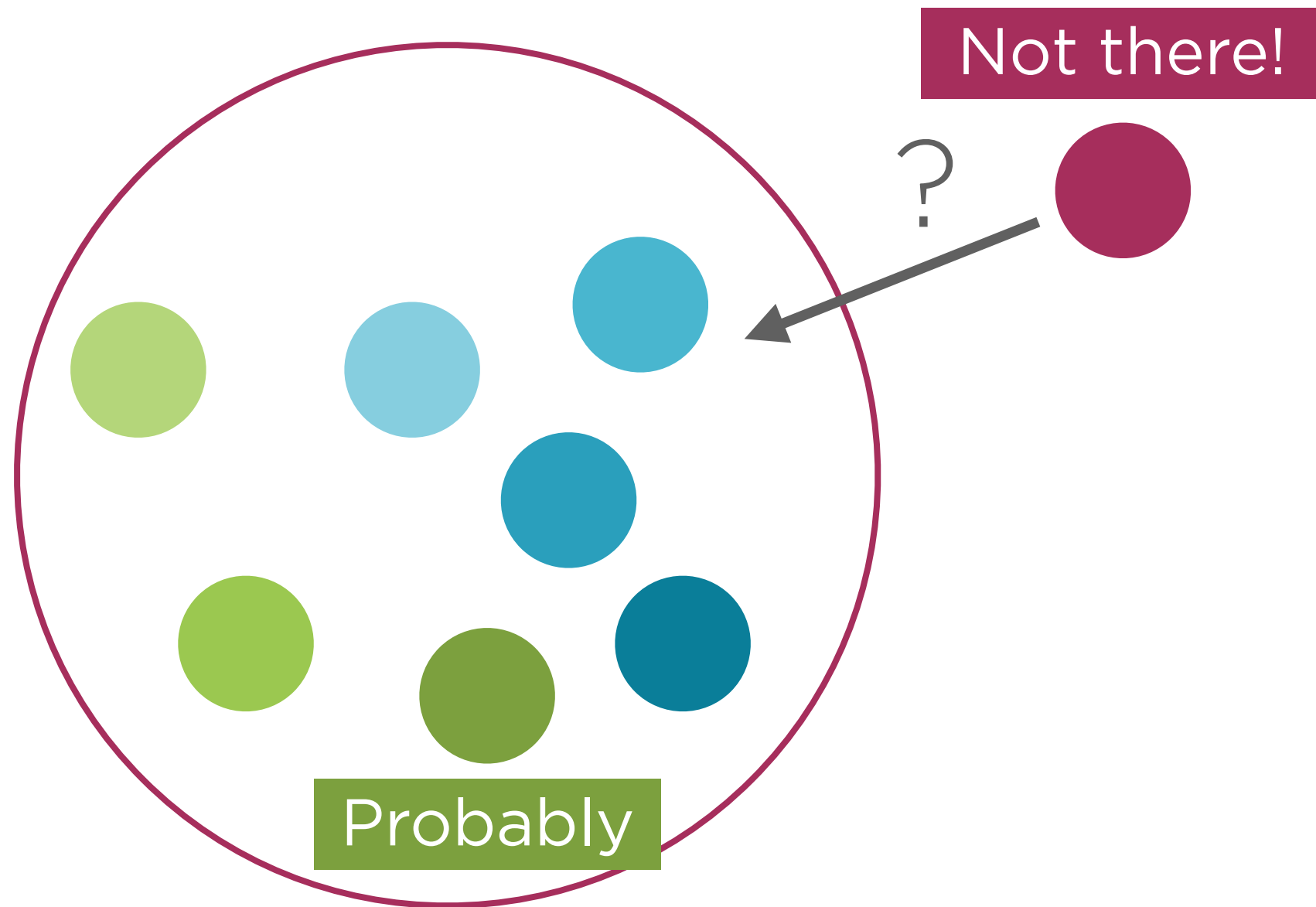
Fishing



Demo

**Wasting cache space**

# Bloom Filters



# Bloom Filters

16 bits (2 bytes)



0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

Use  $k$  independent hash functions:  $h_1(\dots)$ ,  $h_2(\dots)$ , ...,  $h_k(\dots)$

Set the bit specified by each hash function.

Key	$h_1$	$h_2$	$h_3$
“MyFirstLongKeyValue”	14	6	5
“MySecondLargeKeyValue”	2	10	6



Filter is  
never full

Only increased  
rate of false  
positives

## Lookup



1 0 0 1 1 0 0 0 1 0 0 0 1 0  
2 3 4 5 6 7 8 9 10 11 12 13 14 15

All  $k$  bits set: It's probably there.

Not all  $k$  bits set: Definitely not there.



Key	$h_1$	$h_2$	$h_3$	
"MyFirstLongKeyValue"	14	6	5	Probably
"MySecondLargeKeyValue"	2	10	6	Probably
"MyNonExistingKeyValue"	14	2	12	Not there!
"MyFourthLargeKeyValue"			14	Probably



False positive



# Deletion

0 0 1 0 0 1 1 0  
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

Not  
supported

...in this scheme  
(check “*counting bloom filters*”)

All  $k$  bits set: It's probably there.

Not all  $k$  bits set: Definitely not there.

	Key	$h_1$	$h_2$	$h_3$	
✓	“MyFirstLongKeyValue”	14	6	5	Probably
✓	“MySecondLargeKeyValue”	2	10	6	Probably
⊖	“MyNonExistingKeyValue”	14	2	12	Not there!
⊖	“MyFourthLargeKeyValue”			14	Probably

False positive

# Sizing

Targets: Minimize filter size in number of bits ( $m$ ) Deduce  
Minimize risk of false positives ( $p$ ) Given

Other knobs: Number of hash functions ( $k$ ) Deduce  
Number of elements ( $n$ ) Given

$$p = 10\%$$

$$m = 4.79 \cdot n$$

$$k = 4$$

$$p = 1\%$$

$$m = 9.59 \cdot n$$

$$k = 7$$

$$p = 0.1\%$$

$$m = 14.38 \cdot n$$

$$k = 10$$

If  $n = 1\text{M}$  then...

$$m = 4.8\text{M} \text{ (0.6 MB)}$$

$$m = 9.6\text{M} \text{ (1.2 MB)}$$

$$m = 14.4\text{M} \text{ (1.8 MB)}$$

If  $n = 100\text{M}$  then...

$$m = 480\text{M} \text{ (56 MB)}$$

$$m = 959\text{M} \text{ (117 MB)}$$

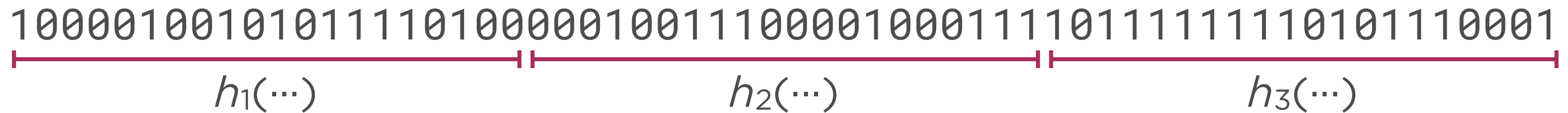
$$m = 1.4\text{G} \text{ (175 MB)}$$



# Trick #1

How do we get that many independent hash functions?

Share the bits!



Bloom filter is  $m$  bits wide.

Requires  $\log_2(m)$  bits in hash value to address.

# Sizing

Targets: Minimize filter size in number of bits ( $m$ ) **Deduce**  
Minimize risk of false positives ( $p$ ) **Given**

Other knobs: Functions ( $k$ ) **Deduce**  
Items ( $n$ ) **Given**

Requires a 20 bit hash value

Requires a 30 bit hash value

$$p = 10\%$$

$$m = 4.79 \cdot n$$

$$k = 4$$

If  $n = 1\text{M}$  then...

$$m = 4.8\text{M} \text{ (0.6 MB)}$$

If  $n = 100\text{M}$  then...

$$m = 480\text{M} \text{ (56 MB)}$$

$$p = 1\%$$

$$m = 9.6 \cdot n$$

$$k = 10$$

$$m = 9.6\text{M} \text{ (1.2 MB)}$$

$$m = 959\text{M} \text{ (117 MB)}$$

$$p = 0.1\%$$

$$m = 14.38 \cdot n$$

$$k = 10$$

$$m = 14.4\text{M} \text{ (1.8 MB)}$$

$$m = 1.4\text{G} \text{ (175 MB)}$$

# Trick #2

How do we get that many independent hash functions?

Combine two hash functions,  $h_1(\dots)$  and  $h_2(\dots)$ .

$$h_i(\dots) = h_1(\dots) + i \cdot h_2(\dots)$$

Example:

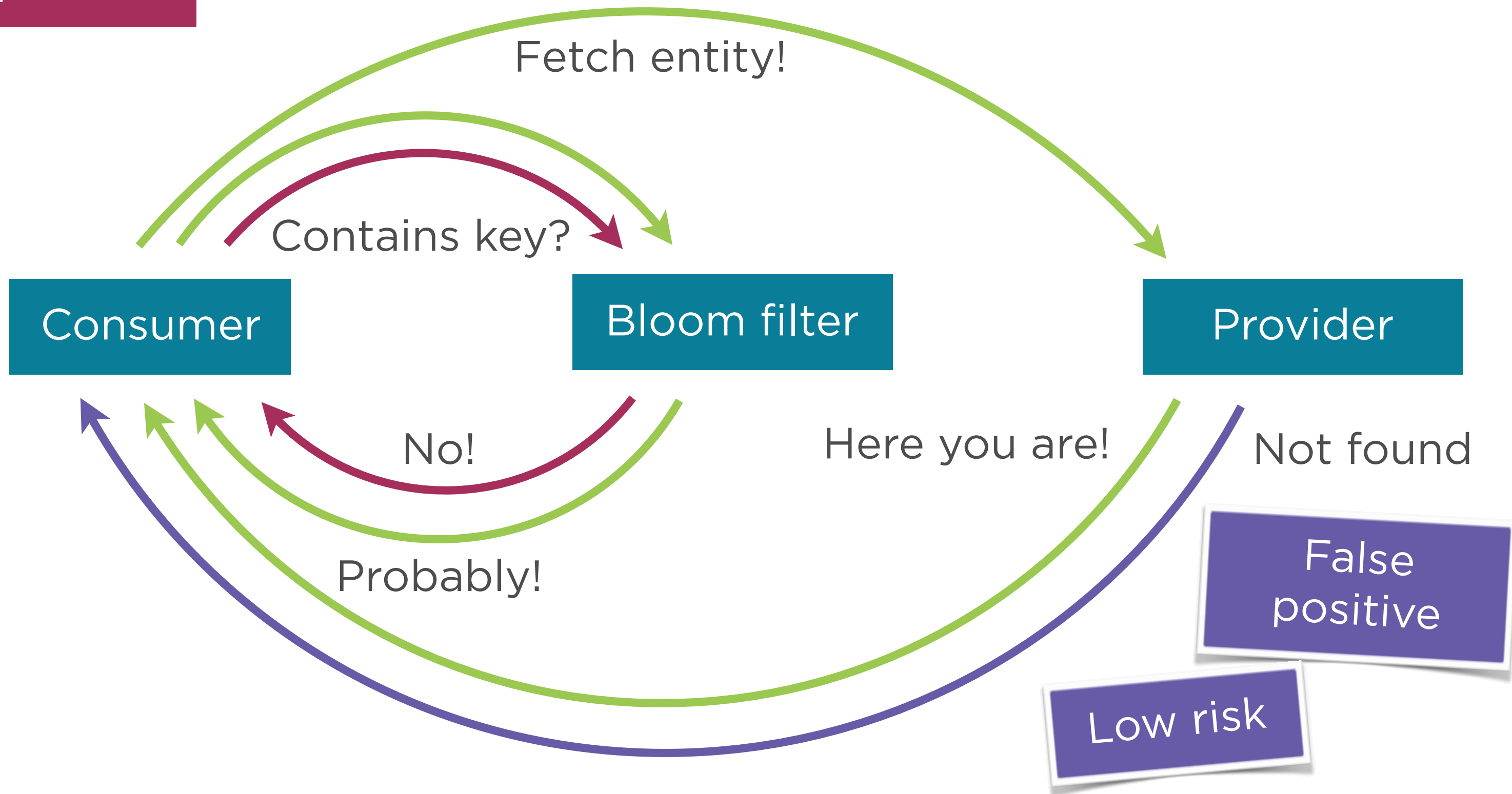
$$h_4(\dots) = h_1(\dots) + 4 \cdot h_2(\dots)$$

$$h_5(\dots) = h_1(\dots) + 5 \cdot h_2(\dots)$$

$$h_6(\dots) = h_1(\dots) + 6 \cdot h_2(\dots)$$

# How to Use It

No need to  
ask provider



# Applications

Cache saver



First-responder in  
NoSQL databases



Identify malicious  
URLs



Hy-phen-a-tion

Wikipedia: “Bloom filter”

Bloom filter by example: [tinyurl.com/bloomByExample](http://tinyurl.com/bloomByExample)

In-depth details: [tinyurl.com/cuckoobloom](http://tinyurl.com/cuckoobloom)

# Demo

**Saving the cache**



# Lessons Learned

111001001010100111001001001010110100011110110011110101011001

Bits are set using  
multiple hash  
functions

Lookups: Check all  
bits specified by  
hash functions

False positives

No false negatives!

Size and number of hash  
functions depends on false  
positive rate