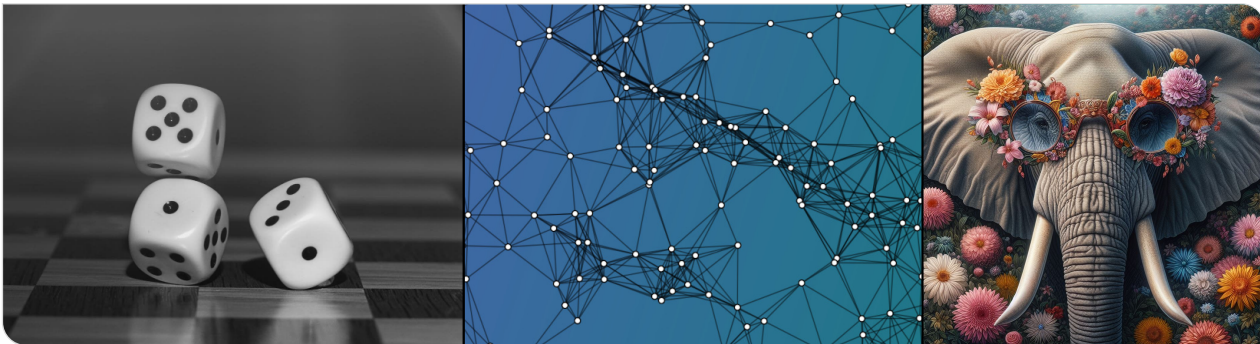# Probability and Computing – Bounded Differences and Bloom Filters

Stefan Walzer | WS 2024/2025

# Content

Method of Bounded Differences ●○○○

What is a Filter or AMQ? ○○

The Bloom Filter Data Structure ○○○

Analysis of Bloom Filters ○○○○○○○○○

Conclusion ○○

**2/20**     WS 2024/2025     Stefan Walzer: Bloom Filters     ITI, Algorithm Engineering

# **More Concentration Bounds**

Hoeffding: Let $X = X_1 + \cdots + X_n$ where $a_i \leq X_i \leq b_i$, and $X_1, \ldots, X_n$ independent

$\Pr[X - \mathbb{E}[X] \geq t] \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^{n}(b_i - a_i)^2}\right).$ // proof combines Chernoff with a different lemma

Method of Bounded Differences    What is a Filter or AMQ?    The Bloom Filter Data Structure    Analysis of Bloom Filters    Conclusion

**3/20**    WS 2024/2025    Stefan Walzer: Bloom Filters      ITI, Algorithm Engineering

# More Concentration Bounds

Hoeffding: Let $X = X_1 + \cdots + X_n$ where $a_i \leq X_i \leq b_i$, and $X_1, \ldots, X_n$ independent

$\Pr[X - \mathbb{E}[X] \geq t] \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^{n}(b_i - a_i)^2}\right)$. // proof combines Chernoff with a different lemma

Special Case: $X_1, \ldots, X_n$ are Bernoulli random variables

$\Pr[X - \mathbb{E}[X] \geq t] \leq \exp(-2t^2/n)$. // incomparable to the Chernoff bound we saw: $\Pr[X \geq (1 + \delta)\mathbb{E}[X]] \leq \left(\frac{e^\delta}{(1+\delta)^{1+\delta}}\right)^\mu$

Method of Bounded Differences    What is a Filter or AMQ?    The Bloom Filter Data Structure    Analysis of Bloom Filters    Conclusion

**3/20**    WS 2024/2025    Stefan Walzer: Bloom Filters    ITI, Algorithm Engineering

# More Concentration Bounds

Hoeffding: Let $X = X_1 + \cdots + X_n$ where $a_i \leq X_i \leq b_i$, and $X_1, \ldots, X_n$ independent

$\Pr[X - \mathbb{E}[X] \geq t] \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^{n}(b_i - a_i)^2}\right)$. // proof combines Chernoff with a different lemma

Special Case: $X_1, \ldots, X_n$ are Bernoulli random variables

$\Pr[X - \mathbb{E}[X] \geq t] \leq \exp(-2t^2/n)$. // incomparable to the Chernoff bound we saw: $\Pr[X \geq (1 + \delta)\mathbb{E}[X]] \leq \left(\frac{e^\delta}{(1+\delta)^{1+\delta}}\right)^\mu$

Generalisation: McDiarmid / Method of Bounded Differences

- Assume $X_1, \ldots, X_n$ are independent and $X = f(X_1, \ldots, X_n)$.
- Assume for each $i$, a change in $X_i$ changes $f(X_1, \ldots, X_n)$ by at most $c_i$.

$\Pr[X - \mathbb{E}[X] \geq t] \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^{n} c_i^2}\right)$ // same bound holds for $\Pr[\mathbb{E}[X] - X \geq t]$.

Method of Bounded Differences
○●○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

3/20    WS 2024/2025    Stefan Walzer: Bloom Filters    ITI, Algorithm Engineering

## More Concentration Bounds

Hoeffding: Let $X = X_1 + \cdots + X_n$ where $a_i \leq X_i \leq b_i$, and $X_1, \ldots, X_n$ independent

$\Pr[X - \mathbb{E}[X] \geq t] \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^{n}(b_i-a_i)^2}\right)$. // proof combines Chernoff with a different lemma

Special Case: $X_1, \ldots, X_n$ are Bernoulli random variables

$\Pr[X - \mathbb{E}[X] \geq t] \leq \exp(-2t^2/n)$. // incomparable to the Chernoff bound we saw: $\Pr[X \geq (1+\delta)\mathbb{E}[X]] \leq \left(\frac{e^\delta}{(1+\delta)^{1+\delta}}\right)^\mu$

Generalisation: McDiarmid / Method of Bounded Differences

- Assume $X_1, \ldots, X_n$ are independent and $X = f(X_1, \ldots, X_n)$.
- Assume for each $i$, a change in $X_i$ changes $f(X_1, \ldots, X_n)$ by at most $c_i$.

$\Pr[X - \mathbb{E}[X] \geq t] \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^{n} c_i^2}\right)$ // same bound holds for $\Pr[\mathbb{E}[X] - X \geq t]$.

Proof uses Martingales... not here.

Method of Bounded Differences    What is a Filter or AMQ?    The Bloom Filter Data Structure    Analysis of Bloom Filters    Conclusion
○●○○                              ○○                          ○○○                                 ○○○○○○○○○                    ○○

**3/20**    WS 2024/2025    Stefan Walzer: Bloom Filters                                                            ITI, Algorithm Engineering

# A Simple Use Case for the Method of Bounded Differences

## Setting

- fixed graph $G = (V, E)$ and $s, t \in V$
- independent edge lengths $X_1, \ldots, X_m \sim \mathcal{U}([0, 1])$
- let $P$ be the shortest $s - t$ path. // unique with probability 1

## McDiarmid / Bounded Differences

- $X_1, \ldots, X_n$ independent and $X = f(X_1, \ldots, X_n)$.
- changing $X_i$ changes $f(X_1, \ldots, X_n)$ by at most $c_i$.

$\Rightarrow \Pr[X - \mathbb{E}[X] \geq t] \leq \exp(-\frac{2t^2}{\sum_{i=1}^{n} c_i^2})$.

Method of Bounded Differences
○○●○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**4/20**    WS 2024/2025    Stefan Walzer: Bloom Filters    ITI, Algorithm Engineering

# A Simple Use Case for the Method of Bounded Differences

## Setting

- fixed graph $G = (V, E)$ and $s, t \in V$
- independent edge lengths $X_1, \ldots, X_m \sim \mathcal{U}([0, 1])$
- let $P$ be the shortest $s - t$ path. // unique with probability 1

## McDiarmid / Bounded Differences

- $X_1, \ldots, X_n$ independent and $X = f(X_1, \ldots, X_n)$.
- changing $X_i$ changes $f(X_1, \ldots, X_n)$ by at most $c_i$.

$\Rightarrow \Pr[X - \mathbb{E}[X] \geq t] \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^{n} c_i^2}\right).$

## Example where McDiarmid might be useful

Let $L :=$ length of $P$.

- $L = f(X_1, \ldots, X_m)$ for independent $X_1, \ldots, X_m$. ✓
- Changing $X_i$ changes $L$ by at most 1. $\rightsquigarrow c_i = 1$. ✓

$\Pr[L - \mathbb{E}[L] \geq t] \leq \exp\left(-\frac{2t^2}{m}\right)$ // might be useful

Method of Bounded Differences
○○●○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**4/20**   WS 2024/2025   Stefan Walzer: Bloom Filters                    ITI, Algorithm Engineering

# A Simple Use Case for the Method of Bounded Differences

## Setting

- fixed graph $G = (V, E)$ and $s, t \in V$
- independent edge lengths $X_1, \ldots, X_m \sim \mathcal{U}([0, 1])$
- let $P$ be the shortest $s - t$ path. // unique with probability 1

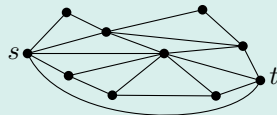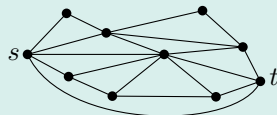## McDiarmid / Bounded Differences

- $X_1, \ldots, X_n$ independent and $X = f(X_1, \ldots, X_n)$.
- changing $X_i$ changes $f(X_1, \ldots, X_n)$ by at most $c_i$.

$\Rightarrow \Pr[X - \mathbb{E}[X] \geq t] \leq \exp(-\frac{2t^2}{\sum_{i=1}^n c_i^2})$.

## Example where McDiarmid seems useless

Let $H :=$ number of edges in $P$. // "hops"

- $H = f(X_1, \ldots, X_m)$ for independent $X_1, \ldots, X_m$. ✓
- Changing $X_i$ might change $H$ by $n - 2$. $\rightsquigarrow c_i = O(n)$.

$\Pr[H - \mathbb{E}[H] \geq t] \leq \exp(-\frac{2t^2}{\mathcal{O}(mn^2)})$ // too weak

Method of Bounded Differences
○○●○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**4/20**    WS 2024/2025    Stefan Walzer: Bloom Filters                                    ITI, Algorithm Engineering

# Content

Method of Bounded Differences     What is a Filter or AMQ?     The Bloom Filter Data Structure     Analysis of Bloom Filters     Conclusion

**5/20**    WS 2024/2025    Stefan Walzer: Bloom Filters     ITI, Algorithm Engineering

# Filter = Approximate Membership Query Data Structure

## Setting

- universe $D$ of possible keys
- a set $S \subseteq D$ of $n = |S|$
- a false positive probability $\varepsilon$

Want: Data structure representing $S$.

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
●○

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**6/20**   WS 2024/2025    Stefan Walzer: Bloom Filters                                                ITI, Algorithm Engineering

# Filter = Approximate Membership Query Data Structure

## Setting

- universe $D$ of possible keys
- a set $S \subseteq D$ of $n = |S|$
- a false positive probability $\varepsilon$

Want: Data structure representing $S$.

## Space Requirement

- want $\mathcal{O}(n \log(1/\varepsilon))$ bits
- *much* smaller than $\mathcal{O}(n \log |D|)$ bits needed for hash table

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
●○

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**6/20**   WS 2024/2025   Stefan Walzer: Bloom Filters   ITI, Algorithm Engineering

# Filter = Approximate Membership Query Data Structure

## Setting

- universe $D$ of possible keys
- a set $S \subseteq D$ of $n = |S|$
- a false positive probability $\varepsilon$

Want: Data structure representing $S$.

## Space Requirement

- want $\mathcal{O}(n \log(1/\varepsilon))$ bits
- *much* smaller than $\mathcal{O}(n \log |D|)$ bits needed for hash table

## Operations

- **insert** elements to $S$ and **delete** elements from $S$ (optional)
- **query**: given $x \in D$ answer "is $x \in S$?" *approximately*:

$$\textbf{query}(x) = \textcolor{blue}{\textbf{YES}} \text{ for } x \in S$$

$$\Pr[\textbf{query}(x) = \textcolor{red}{\textbf{NO}}] \geq 1 - \varepsilon \text{ for } x \notin S$$

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
●○

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**6/20**   WS 2024/2025   Stefan Walzer: Bloom Filters   ITI, Algorithm Engineering

# Filter = Approximate Membership Query Data Structure

**a set $S$:**

Anja Blancani
Peter Sanders
Florian Kurpicz
Hape Lehmann
Thomas Worsch

## Setting

- universe $D$ of possible keys
- a set $S \subseteq D$ of $n = |S|$
- a false positive probability $\varepsilon$

Want: Data structure representing $S$.

## Space Requirement

- want $\mathcal{O}(n \log(1/\varepsilon))$ bits
- *much* smaller than $\mathcal{O}(n \log |D|)$ bits needed for hash table

## Operations

- **insert** elements to $S$ and **delete** elements from $S$ (optional)
- **query**: given $x \in D$ answer "is $x \in S$?" *approximately*:

$$\text{query}(x) = \textbf{YES} \text{ for } x \in S$$
$$\Pr[\textbf{query}(x) = \textbf{NO}] \geq 1 - \varepsilon \text{ for } x \notin S$$

**a filter for $S$:**

Anj
Pete
Flo
Hap
Tho

query(Peter Sanders) = **YES**

query(Petra Mutzel) = **YES**

query(Donald Knuth) = **NO**

The **YES** answers are **unreliable**.
The **NO** answers are **reliable**.

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
●○

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**6/20**   WS 2024/2025   Stefan Walzer: Bloom Filters

ITI, Algorithm Engineering

# Simplified Use Case for Filters

CPU

get data for:
"Thomas Bläsius"

Bernhard Beckert
Klemens Böhm
Anne Koziolek
André Platzer
Alexander Weigl

**NO**

Anja Blancani
Peter Sanders
Florian Kurpicz
Hape Lehmann
Thomas Worsch

**NO**

Michael Beigl
Thomas Bläsius
Carsten Dachsbacher
Pascal Friederich
Hannes Hartenstein

**YES**

10 000 000 ns

10 000 000 ns

10 000 000 ns

Method of Bounded Differences    What is a Filter or AMQ?    The Bloom Filter Data Structure    Analysis of Bloom Filters    Conclusion

**7/20**    WS 2024/2025    Stefan Walzer: Bloom Filters    ITI, Algorithm Engineering

# Simplified Use Case for Filters

CPU · · · (hard disk)

get data for:
"Thomas Bläsius"

50 ns → Ber / Kle / A / An / Ale   **NO** ········

Bernhard Beckert
Klemens Böhm
Anne Koziolek
André Platzer
Alexander Weigl

50 ns → An / Pe / Flo / Hap / Th   **YES**  10 000 000 ns →

Anja Blancani
Peter Sanders
Florian Kurpicz
Hape Lehmann
Thomas Worsch    **NO**

50 ns → Mic / Th / Car / Pá / Ha   **YES**  10 000 000 ns →

Michael Beigl
Thomas Bläsius
Carsten Dachsbacher
Pascal Friederich
Hannes Hartenstein    **YES**

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○●

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○○○○○○○

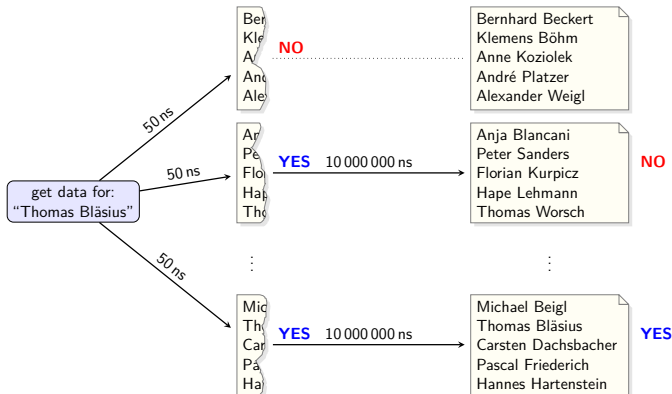Conclusion
○○

# Simplified Use Case for Filters



## General Idea

If the reliable **NO** answers are frequent, a filter access can replace a (costly) access to a reliable data structure.

Method of Bounded Differences ○○○○  
What is a Filter or AMQ? ○●  
The Bloom Filter Data Structure ○○○  
Analysis of Bloom Filters ○○○○○○○○○  
Conclusion ○○

**7/20**   WS 2024/2025   Stefan Walzer: Bloom Filters   ITI, Algorithm Engineering
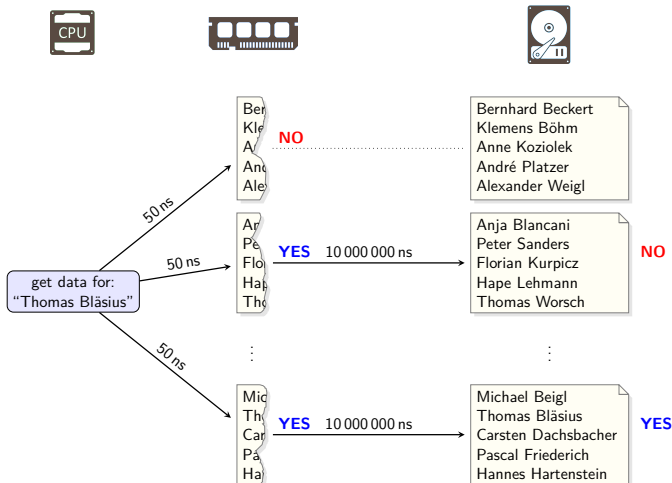
# Simplified Use Case for Filters



## General Idea

If the reliable **NO** answers are frequent, a filter access can replace a (costly) access to a reliable data structure.

## Further Example

Filter stores set of malicious URLs.

- Most accessed URLs will be non-malicious.
- Only false positives and true positives have to access reliable data structure (e.g. web-service).

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○●

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**7/20**  WS 2024/2025  Stefan Walzer: Bloom Filters  ITI, Algorithm Engineering

# Content

Method of Bounded Differences    What is a Filter or AMQ?    The Bloom Filter Data Structure    Analysis of Bloom Filters    Conclusion

**8/20**   WS 2024/2025   Stefan Walzer: Bloom Filters      ITI, Algorithm Engineering

# Reminder: SUHA

## Simple Uniform Hashing Assumption (SUHA)

- We have access to $h \sim \mathcal{U}(R^D)$ for any $R$ and $D$.
- $h$ takes $\mathcal{O}(1)$ time to evaluate.
- $h$ takes no space to store.

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○●○

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**9/20**   WS 2024/2025   Stefan Walzer: Bloom Filters                                                                ITI, Algorithm Engineering

# The Bloom Filter Data Structure

## Parameters

| | |
|---|---|
| $m$ | length of a bit array $A[1..m]$ that we use |
| $k \in \mathcal{O}(1)$ | number of hash functions $h_1, \ldots, h_k \sim \mathcal{U}([m]^D)$ |
| $n$ | number of keys in $S \subseteq D$ (dynamic) |
| $\alpha \in \mathcal{O}(1)$ | load $n/m$ (dynamic) |

each cell stores a bit

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

1  2  3  $\cdots$ $\cdots$ $m$

### insert($x$):

**for** $i \in [k]$ **do**
$\quad A[h_i(x)] = 1$

### query($x$):

**for** $i \in [k]$ **do**
$\quad$ **if** $A[h_i(x)] = 0$ **then**
$\quad\quad$ **return NO**

**return YES**

### delete($x$):

$\langle$ not supported $\rangle$

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
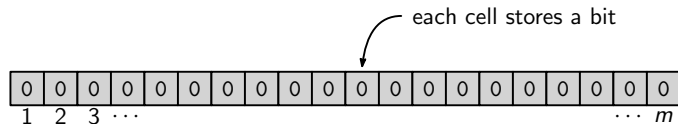○○●

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**10/20**  WS 2024/2025  Stefan Walzer: Bloom Filters  ITI, Algorithm Engineering

# The Bloom Filter Data Structure

## Parameters

| | |
|---|---|
| $m$ | length of a bit array $A[1..m]$ that we use |
| $k \in \mathcal{O}(1)$ | number of hash functions $h_1, \ldots, h_k \sim \mathcal{U}([m]^D)$ |
| $n$ | number of keys in $S \subseteq D$ (dynamic) |
| $\alpha \in \mathcal{O}(1)$ | load $n/m$ (dynamic) |

$(h_1(x_1), h_2(x_1), h_3(x_1)) = (6, 9, 14) \in [m]^3$

$k = 3$



$$1 \quad 2 \quad 3 \quad \cdots \qquad\qquad\qquad\qquad \cdots \quad m$$

### insert($x$):

**for** $i \in [k]$ **do**
$\quad A[h_i(x)] = 1$

### query($x$):

**for** $i \in [k]$ **do**
$\quad$ **if** $A[h_i(x)] = 0$ **then**
$\quad\quad$ **return NO**

**return YES**

### delete($x$):

⟨ not supported ⟩

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○●

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**10/20**  WS 2024/2025  Stefan Walzer: Bloom Filters

ITI, Algorithm Engineering

# The Bloom Filter Data Structure

## Parameters

| | |
|---|---|
| $m$ | length of a bit array $A[1..m]$ that we use |
| $k \in \mathcal{O}(1)$ | number of hash functions $h_1, \ldots, h_k \sim \mathcal{U}([m]^D)$ |
| $n$ | number of keys in $S \subseteq D$ (dynamic) |
| $\alpha \in \mathcal{O}(1)$ | load $n/m$ (dynamic) |



$S = \{\ x_1,\ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \}$

| 0 | 0 | 0 | 0 | 0 | 0 | **1** | 0 | **1** | 0 | 0 | 0 | 0 | **1** | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

1 2 3 $\cdots$ $\cdots$ $m$

### insert($x$):

**for** $i \in [k]$ **do**
$\quad A[h_i(x)] = 1$

### query($x$):

**for** $i \in [k]$ **do**
$\quad$ **if** $A[h_i(x)] = 0$ **then**
$\quad\quad$ **return NO**

**return YES**

### delete($x$):

⟨ not supported ⟩

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○●

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**10/20**   WS 2024/2025   Stefan Walzer: Bloom Filters

ITI, Algorithm Engineering

# The Bloom Filter Data Structure

## Parameters

| | |
|---|---|
| $m$ | length of a bit array $A[1..m]$ that we use |
| $k \in \mathcal{O}(1)$ | number of hash functions $h_1, \ldots, h_k \sim \mathcal{U}([m]^D)$ |
| $n$ | number of keys in $S \subseteq D$ (dynamic) |
| $\alpha \in \mathcal{O}(1)$ | load $n/m$ (dynamic) |



$S = \{ \quad x_1, x_2, \qquad \}$

| 0 | 0 | **1** | 0 | **1** | 0 | **1** | 0 | **1** | 0 | 0 | **1** | 0 | **1** | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | $\cdots$ | | | | | | | | | | | | | | $\cdots$ | $m$ | |

### insert($x$):

**for** $i \in [k]$ **do**
$\quad A[h_i(x)] = 1$

### query($x$):

**for** $i \in [k]$ **do**
$\quad$ **if** $A[h_i(x)] = 0$ **then**
$\quad\quad$ **return NO**

**return YES**

### delete($x$):

$\langle$ not supported $\rangle$

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○●

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**10/20**  WS 2024/2025   Stefan Walzer: Bloom Filters

ITI, Algorithm Engineering

# The Bloom Filter Data Structure



insert($x$):

**for** $i \in [k]$ **do**
$\quad A[h_i(x)] = 1$

query($x$):

**for** $i \in [k]$ **do**
$\quad$ **if** $A[h_i(x)] = 0$ **then**
$\quad\quad$ **return NO**

**return YES**

delete($x$):

$\langle$ not supported $\rangle$

## Parameters

| | |
|---|---|
| $m$ | length of a bit array $A[1..m]$ that we use |
| $k \in \mathcal{O}(1)$ | number of hash functions $h_1, \ldots, h_k \sim \mathcal{U}([m]^D)$ |
| $n$ | number of keys in $S \subseteq D$ (dynamic) |
| $\alpha \in \mathcal{O}(1)$ | load $n/m$ (dynamic) |

$S = \{ \quad x_1, x_2, x_3, \quad \quad \}$

Method of Bounded Differences
OOOO

What is a Filter or AMQ?
OO

The Bloom Filter Data Structure
OO●

Analysis of Bloom Filters
OOOOOOOOO

Conclusion
OO

**10/20**   WS 2024/2025    Stefan Walzer: Bloom Filters

ITI, Algorithm Engineering

# The Bloom Filter Data Structure

## Parameters

| | |
|---|---|
| $m$ | length of a bit array $A[1..m]$ that we use |
| $k \in \mathcal{O}(1)$ | number of hash functions $h_1, \ldots, h_k \sim \mathcal{U}([m]^D)$ |
| $n$ | number of keys in $S \subseteq D$ (dynamic) |
| $\alpha \in \mathcal{O}(1)$ | load $n/m$ (dynamic) |



$S = \{\ x_1, x_2, x_3, x_4, x_5\ \}$

insert($x$):

**for** $i \in [k]$ **do**
$\quad A[h_i(x)] = 1$

query($x$):

**for** $i \in [k]$ **do**
$\quad$ **if** $A[h_i(x)] = 0$ **then**
$\quad\quad$ **return NO**

**return YES**

delete($x$):

$\langle$ not supported $\rangle$

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○●

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**10/20**　WS 2024/2025　　Stefan Walzer: Bloom Filters　　　　　　　　　　　　　　　　　ITI, Algorithm Engineering

# The Bloom Filter Data Structure

## Parameters

| | |
|---|---|
| $m$ | length of a bit array $A[1..m]$ that we use |
| $k \in \mathcal{O}(1)$ | number of hash functions $h_1, \ldots, h_k \sim \mathcal{U}([m]^D)$ |
| $n$ | number of keys in $S \subseteq D$ (dynamic) |
| $\alpha \in \mathcal{O}(1)$ | load $n/m$ (dynamic) |

| 0 | 0 | **1** | 0 | **1** | 0 | **1** | 0 | **1** | **1** | 0 | **1** | 0 | **1** | **1** | 0 | **1** | 0 | **1** | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

1   2   3 $\cdots$                            $\cdots$   $m$

## insert($x$):

**for** $i \in [k]$ **do**
   | $A[h_i(x)] = 1$

## query($x$):

**for** $i \in [k]$ **do**
   | **if** $A[h_i(x)] = 0$ **then**
      | **return NO**

**return YES**

## delete($x$):

⟨ not supported ⟩

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○●

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**10/20**   WS 2024/2025   Stefan Walzer: Bloom Filters            ITI, Algorithm Engineering

# The Bloom Filter Data Structure

## Parameters

| | |
|---|---|
| $m$ | length of a bit array $A[1..m]$ that we use |
| $k \in \mathcal{O}(1)$ | number of hash functions $h_1, \ldots, h_k \sim \mathcal{U}([m]^D)$ |
| $n$ | number of keys in $S \subseteq D$ (dynamic) |
| $\alpha \in \mathcal{O}(1)$ | load $n/m$ (dynamic) |



insert(x):

**for** $i \in [k]$ **do**
$\quad A[h_i(x)] = 1$

query(x):

**for** $i \in [k]$ **do**
$\quad$ **if** $A[h_i(x)] = 0$ **then**
$\quad\quad$ **return NO**

**return YES**

delete(x):

⟨ not supported ⟩

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○●

Analysis of Bloom Filters
○○○○○○○○○

Conclusion
○○

**10/20**  WS 2024/2025   Stefan Walzer: Bloom Filters                                      ITI, Algorithm Engineering

# Content

Method of Bounded Differences   What is a Filter or AMQ?   The Bloom Filter Data Structure   Analysis of Bloom Filters   Conclusion

**11/20**   WS 2024/2025   Stefan Walzer: Bloom Filters   ITI, Algorithm Engineering

# **Preparation**

### Exercise: Some approximations of $e$

$$\forall n \in \mathbb{N} : (1 + \tfrac{1}{n})^n \leq e \leq (1 + \tfrac{1}{n})^{n+1}$$
$$\text{and} \quad (1 - \tfrac{1}{n})^n \leq e^{-1} \leq (1 - \tfrac{1}{n})^{n-1}.$$

### Corollaries

$$\forall n \in \mathbb{N} : (1 + \tfrac{1}{n})^n = e - \mathcal{O}(1/n)$$
$$\text{and} \quad (1 - \tfrac{1}{n})^n = e^{-1} - \mathcal{O}(1/n).$$

Method of Bounded Differences     What is a Filter or AMQ?     The Bloom Filter Data Structure     **Analysis of Bloom Filters**     Conclusion
○○○○         ○○         ○○○         ○●○○○○○○○         ○○

**12/20**    WS 2024/2025    Stefan Walzer: Bloom Filters          ITI, Algorithm Engineering
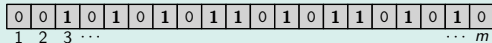
# Bloom Filter Analysis (i)

## Lemma

Assume $S = \{x_1, \ldots, x_n\}$ is inserted into the Bloom filter. Let $(A_1, \ldots, A_m) \in \{0, 1\}^m$ be the random filter state and $Z := \sum_{i=1}^{m}(1 - A_i)$ the number of zeroes. Then

**i** $\mathbb{E}[\frac{Z}{m}] = (1 - \frac{1}{m})^{m\alpha k} = e^{-\alpha k} - o(1)$

**ii** For $y \notin S : \Pr[\mathbf{query}(y) = \mathbf{YES} \mid Z = z] = (1 - \frac{z}{m})^k$

| 0 | 0 | **1** | 0 | **1** | 0 | **1** | 0 | **1** | **1** | 0 | **1** | 0 | **1** | **1** | 0 | **1** | 0 | **1** | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

1 2 3 $\cdots$ $\cdots$ $m$

Method of Bounded Differences     What is a Filter or AMQ?     The Bloom Filter Data Structure     Analysis of Bloom Filters     Conclusion

**13/20**    WS 2024/2025    Stefan Walzer: Bloom Filters          ITI, Algorithm Engineering
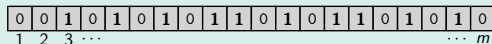
# Bloom Filter Analysis (i)

## Lemma

Assume $S = \{x_1, \ldots, x_n\}$ is inserted into the Bloom filter. Let $(A_1, \ldots, A_m) \in \{0, 1\}^m$ be the random filter state and $Z := \sum_{i=1}^{m}(1 - A_i)$ the number of zeroes. Then

**i** $\mathbb{E}[\frac{Z}{m}] = (1 - \frac{1}{m})^{m\alpha k} = e^{-\alpha k} - o(1)$

**ii** For $y \notin S$ : $\Pr[\textbf{query}(y) = \textbf{YES} \mid Z = z] = (1 - \frac{z}{m})^k$

| 0 | 0 | **1** | 0 | **1** | 0 | **1** | 0 | **1** | **1** | 0 | **1** | 0 | **1** | **1** | 0 | **1** | 0 | **1** | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

1  2  3 $\cdots$ $\qquad\qquad\qquad\qquad$ $\cdots$ $m$

## Proof of (i).

$$\mathbb{E}[\tfrac{Z}{m}] = \tfrac{1}{m}\mathbb{E}[\sum_{i=1}^{m}(1 - A_i)] = \tfrac{1}{m}\sum_{i=1}^{m}\Pr[A_i = 0] = \tfrac{1}{m}\sum_{i=1}^{m}\Pr[A_1 = 0] = \Pr[A_1 = 0]$$

$$= \Pr[\forall x \in S : \forall i \in [k] : h_i(x) \neq 1] \overset{\text{SUHA}}{=} \prod_{x\in S}\prod_{i\in[k]}\Pr[h_i(x) \neq 1] \overset{\text{SUHA}}{=} \prod_{x\in S}\prod_{i\in[k]}(1 - \tfrac{1}{m})$$

$$= (1 - \tfrac{1}{m})^{nk} = (1 - \tfrac{1}{m})^{m\alpha k} = (e^{-1} - o(1))^{\alpha k} = e^{-\alpha k} - o(1).$$

Method of Bounded Differences
0000

What is a Filter or AMQ?
00

The Bloom Filter Data Structure
000

Analysis of Bloom Filters
00●000000

Conclusion
00

**13/20**   WS 2024/2025   Stefan Walzer: Bloom Filters
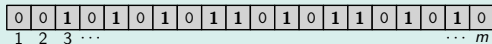
ITI, Algorithm Engineering

# Bloom Filter Analysis (i)

## Lemma

Assume $S = \{x_1, \dots, x_n\}$ is inserted into the Bloom filter. Let $(A_1, \dots, A_m) \in \{0, 1\}^m$ be the random filter state and $Z := \sum_{i=1}^{m}(1 - A_i)$ the number of zeroes. Then

**i** $\mathbb{E}[\frac{Z}{m}] = (1 - \frac{1}{m})^{m\alpha k} = e^{-\alpha k} - o(1)$

**ii** For $y \notin S : \Pr[\mathbf{query}(y) = \mathbf{YES} \mid Z = z] = (1 - \frac{z}{m})^k$

| 0 | 0 | **1** | 0 | **1** | 0 | **1** | 0 | **1** | **1** | 0 | **1** | 0 | **1** | **1** | 0 | **1** | 0 | **1** | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 $\cdots$ | | | | | | | | | | | | | | | | $\cdots$ | m |

## Proof of (ii).

$$\Pr[\mathbf{query}(y) = \mathbf{YES} \mid Z = z] = \Pr[\forall i \in [k] : A_{h_i(y)} = 1 \mid Z = z] \overset{\text{SUHA}}{=} \prod_{i \in [k]} \left( \frac{m - z}{m} \right) = (1 - \tfrac{z}{m})^k.$$

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○●○○○○○○

Conclusion
○○

**13/20** WS 2024/2025 Stefan Walzer: Bloom Filters

ITI, Algorithm Engineering

# How should a Bloom filter be configured?

*Approximate* false positive rate

From the previous Lemma we get for $y \notin S$:

$$\varepsilon = \Pr[\textbf{query}(y) = \textbf{YES}] \approx \Pr[\textbf{query}(y) = \textbf{YES} \mid Z = \mathbb{E}[Z]]$$

$$\overset{\text{ii}}{=} \left(1 - \frac{\mathbb{E}[Z]}{m}\right)^k \overset{\text{i}}{=} (1 - e^{-\alpha k} + o(1))^k \approx (1 - e^{-\alpha k})^k.$$

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○●○○○○○

Conclusion
○○

**14/20**   WS 2024/2025   Stefan Walzer: Bloom Filters

ITI, Algorithm Engineering

# How should a Bloom filter be configured?
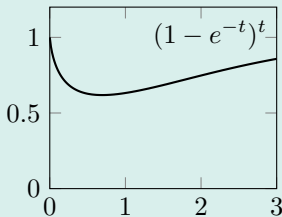
## *Approximate* false positive rate

From the previous Lemma we get for $y \notin S$:

$$\varepsilon = \Pr[\textbf{query}(y) = \textbf{YES}] \approx \Pr[\textbf{query}(y) = \textbf{YES} \mid Z = \mathbb{E}[Z]]$$

$$\overset{\text{ii}}{=} \left(1 - \frac{\mathbb{E}[Z]}{m}\right)^k \overset{\text{i}}{=} (1 - e^{-\alpha k} + o(1))^k \approx (1 - e^{-\alpha k})^k.$$

## Which $k$ minimises $\varepsilon$? (when $\alpha$ is fixed)

$$\underset{k \in \mathbb{N}}{\arg\min}(1 - e^{-\alpha k})^k$$

$$= \underset{k \in \mathbb{N}}{\arg\min}(1 - e^{-\alpha k})^{\alpha k}$$

$$\approx \frac{1}{\alpha} \underset{t \in \mathbb{R}_+}{\arg\min}(1 - e^{-t})^t$$

$$= \frac{1}{\alpha} \underset{t \in \mathbb{R}_+}{\arg\min}\, t \ln(1 - e^{-t})$$



- plot $(1 - e^{-t})^t \rightsquigarrow$ one global minimum.
- deriving $t \ln(1 - e^{-t})$ gives $\ln(1 - e^{-t}) + \frac{te^{-t}}{1 - e^{-t}}$
- $t = \ln(2)$ is root of the derivative.

$\hookrightarrow k = \ln(2)/\alpha$ is optimal for fixed $\alpha$.
$\hookrightarrow$ choose $\alpha$ and $k$ such that $\alpha k = \ln(2)$

Method of Bounded Differences     What is a Filter or AMQ?     The Bloom Filter Data Structure     **Analysis of Bloom Filters**     Conclusion

**14/20**    WS 2024/2025     Stefan Walzer: Bloom Filters       ITI, Algorithm Engineering

# **False Positive Probability and Space**

## Intuition for optimality of $\alpha k = \ln(2)$

- gives $\mathbb{E}[\frac{Z}{m}] \approx e^{-\alpha k} = \frac{1}{2}$
- maximises *entropy* of the filter bits

## Theorem

A Bloom filter with $k \in \mathbb{N}$ hash functions and load factor $\alpha = \ln(2)/k$ has
   **space requirement** $m = n/\alpha = \frac{kn}{\ln 2} \approx 1.44kn$ bits and
 **false positive probability** $\varepsilon = 2^{-k} + o(1)$.

- space requirement ✓
- false positive probability: need a concentration bound first.

Method of Bounded Differences    What is a Filter or AMQ?    The Bloom Filter Data Structure    Analysis of Bloom Filters    Conclusion
○○○○          ○○          ○○○         ○○○○●○○○○       ○○

**15/20**   WS 2024/2025    Stefan Walzer: Bloom Filters                   ITI, Algorithm Engineering

# **Concentration bound for** $Z$

### Lemma

i Pr$[Z \le \mathbb{E}[Z] - t \quad ] \le \exp(-\Theta(t^2/m))$ for any $t > 0$,

ii Pr$[Z \le \mathbb{E}[Z] - m^{2/3}] \le \exp(-\Theta(m^{1/3}))$ by setting $t = m^{2/3}$.

Method of Bounded Differences          What is a Filter or AMQ?          The Bloom Filter Data Structure          Analysis of Bloom Filters          Conclusion
○○○○                                   ○○                                ○○○                                     ○○○○○●○○○                          ○○

**16/20**   WS 2024/2025   Stefan Walzer: Bloom Filters                                                                                    ITI, Algorithm Engineering

# **Concentration bound for** $Z$

## Lemma

- **i** $\Pr[Z \leq \mathbb{E}[Z] - t \quad] \leq \exp(-\Theta(t^2/m))$ for any $t > 0$,
- **ii** $\Pr[Z \leq \mathbb{E}[Z] - m^{2/3}] \leq \exp(-\Theta(m^{1/3}))$ by setting $t = m^{2/3}$.

## McDiarmid / Bounded Differences

- $X_1, \ldots, X_n$ independent and $X = f(X_1, \ldots, X_n)$.
- changing $X_i$ changes $f(X_1, \ldots, X_n)$ by at most $c_i$.

$\Rightarrow \Pr[\mathbb{E}[X] - X \geq t] \leq \exp(-\frac{2t^2}{\sum_{i=1}^{n} c_i^2})$.

## Proof of (i) using the method of bounded differences.

- $Z$ is a function of $kn$ independent hash values // SUHA
- each hash value can change $Z$ by at most 1
- use method of bounded differences!

$$\Rightarrow \Pr[Z \leq \mathbb{E}[Z] - t] \leq \Pr[\mathbb{E}[Z] - Z \geq t] = \exp\left(\frac{-2t^2}{nk}\right) = \exp\left(\frac{-2t^2}{m\alpha k}\right) = \exp\left(\frac{-2t^2}{m\ln(2)}\right). \quad \square$$

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○○○●○○○

Conclusion
○○

**16/20**    WS 2024/2025    Stefan Walzer: Bloom Filters    ITI, Algorithm Engineering

# False Positive Probability of Bloom filters

## Proof of the Main Theorem on Bloom filters (false positive probability).

By choice of $k$ and $\alpha$ we have $\mathbb{E}[\frac{Z}{m}] = e^{-\alpha k} - o(1) = \frac{1}{2} - o(1)$.

Let $y \notin S$ and $B = \lfloor \mathbb{E}[Z] - m^{2/3} \rfloor$.

$$\varepsilon = \Pr[\textbf{query}(y) = \textbf{YES}] \overset{\text{LTP}}{=} \sum_{z=1}^{m} \Pr[Z = z] \cdot \Pr[\textbf{query}(y) = \textbf{YES} \mid Z = z] = \sum_{z=1}^{m} \Pr[Z = z] \cdot (1 - \tfrac{z}{m})^k$$

$$\leq \sum_{z=1}^{B} \Pr[Z = z] + \sum_{z=B+1}^{m} \Pr[Z = z](1 - \tfrac{B+1}{m})^k \leq \Pr[Z \leq B] + (1 - \tfrac{B+1}{m})^k$$

$$\leq \Pr[Z \leq \mathbb{E}[Z] - m^{2/3}] + \big(1 - \tfrac{\mathbb{E}[Z] - m^{2/3}}{m}\big)^k \overset{\text{ii}}{\leq} \exp(-\Theta(m^{1/3})) + (1 - \tfrac{1}{2} + o(1))^k = 2^{-k} + o(1). \quad \square$$

Method of Bounded Differences
OOOO

What is a Filter or AMQ?
OO

The Bloom Filter Data Structure
OOO

Analysis of Bloom Filters
OOOOOO●OO

Conclusion
OO

**17/20**    WS 2024/2025    Stefan Walzer: Bloom Filters    ITI, Algorithm Engineering

# How to Configure Your Bloom Filter

## Theorem

A Bloom filter with $k \in \mathbb{N}$ hash functions and load factor $\alpha = \ln(2)/k$ has
      **space requirement** $m = n/\alpha = \frac{kn}{\ln 2} \approx 1.44kn$ bits and
  **false positive probability** $\varepsilon = 2^{-k} + o(1)$.

## How to determine $m$ and $k$ (the parameters you actually need)

1. $n$: determined by input

2. $\varepsilon$: choose a trade-off between space usage and false positive probability
   - If utility comes from negative answers "$x \notin S$, definitely" and running time is negligible, then:
     - want to maximise utility $-$ disutility, where: ($\propto$ means "proportional to")
     - utility $\propto$ negative answers $=$ queries $\cdot \Pr[x \notin S] \cdot (1 - \varepsilon)$
     - disutility $\propto$ space consumption $= 1.44 \log(1/\varepsilon)n$ bits of RAM or cache

3. compute $k = \lceil \log(1/\varepsilon) \rceil$ // effectively restricts $\varepsilon$ to powers of 2

4. compute $\alpha = \ln(2)/k$ and $m = \lceil n/\alpha \rceil$

Method of Bounded Differences      What is a Filter or AMQ?      The Bloom Filter Data Structure      Analysis of Bloom Filters      Conclusion
oooo        oo        ooo        oooooooo●o        oo

**18/20**    WS 2024/2025     Stefan Walzer: Bloom Filters                                   ITI, Algorithm Engineering

# Remarks

## Much, much more is known

- more functionality
  - $\hookrightarrow$ counting Bloom filters support deletions
- better query times
  - $\hookrightarrow$ blocked Bloom filters improve cache efficiency
- better space efficiency
  - $\hookrightarrow$ cuckoo filters use $n\log(1/\varepsilon) + \mathcal{O}(n)$ bits rather than $\approx 1.44n\log(1/\varepsilon)$ bits
  - $\hookrightarrow$ static filters (no insertions or deletions) use $n\log(1/\varepsilon) + o(n)$ bits.
- ...

Method of Bounded Differences
○○○○

What is a Filter or AMQ?
○○

The Bloom Filter Data Structure
○○○

Analysis of Bloom Filters
○○○○○○○○●

Conclusion
○○

**19/20**   WS 2024/2025   Stefan Walzer: Bloom Filters

ITI, Algorithm Engineering

# Conclusion

- **Approximate Membership Queries.**
  - Decide "is $x \in S$?" with *false positive probability* $\varepsilon$.
  - The Bloom filter is the most widespread AMQ.
- **Space Efficient.** $\approx 1.44 \log(1/\varepsilon)$ bits per element
  - often fit into cache or RAM when proper set data structure does not
- **Used to prevent costly accesses.**
  - Reliable on **NO** answers.
  - Useful if **NO** answers are frequent.

Method of Bounded Differences
0000

What is a Filter or AMQ?
00

The Bloom Filter Data Structure
000

Analysis of Bloom Filters
000000000

Conclusion
●○

WS 2024/2025    Stefan Walzer: Bloom Filters

ITI, Algorithm Engineering

# Anhang: Mögliche Prüfungsfragen I

- Approximate-Membership-Query Datenstrukturen im Allgemeinen
    - Welche Aufgabe hat eine AMQ Datenstruktur?
    - Was ist der Vorteil gegenüber einer exakten Datenstruktur?
    - Was wäre ein Anwendungsfall, in dem eine AMQ Datenstruktur nützlich ist?
- Bloomfilter
    - Wie ist ein Bloomfilter aufgebaut und welche Operationen unterstützt er?
    - Welche Parameter gibt es, und wie hängen diese zusammen?
    - Was hat unsere Analyse zur geschickten Wahl der Parameter zu sagen? Wie werden die übrigen Parameter gewählt? Welcher Speicherverbrauch ergibt sich?
    - Fragen zur Analyse
        - Welche Anzahl von Nullen bzw. Einsen erwarten wir?
        - Wie hängt die falsch-positiv Wahrscheinlichkeit mit der Anzahl Nullen bzw. Einsen zusammen?
        - Wir kann man argumentieren, dass die Anzahl Nullen bzw. Einsen im Bloomfilter nahe am Erwartungswert liegt?

Method of Bounded Differences        What is a Filter or AMQ?        The Bloom Filter Data Structure        Analysis of Bloom Filters        Conclusion
0000                                 00                             000                                   000000000                       0●

**21**/20    WS 2024/2025    Stefan Walzer: Bloom Filters                                                                    ITI, Algorithm Engineering