# CSC4524: Streaming algorithms

**Bloom filter**

# Mini-project 1

- **Be able to know if a given Wikipedia page was visited the 01-08-2016**
  - Input: domain code + page title
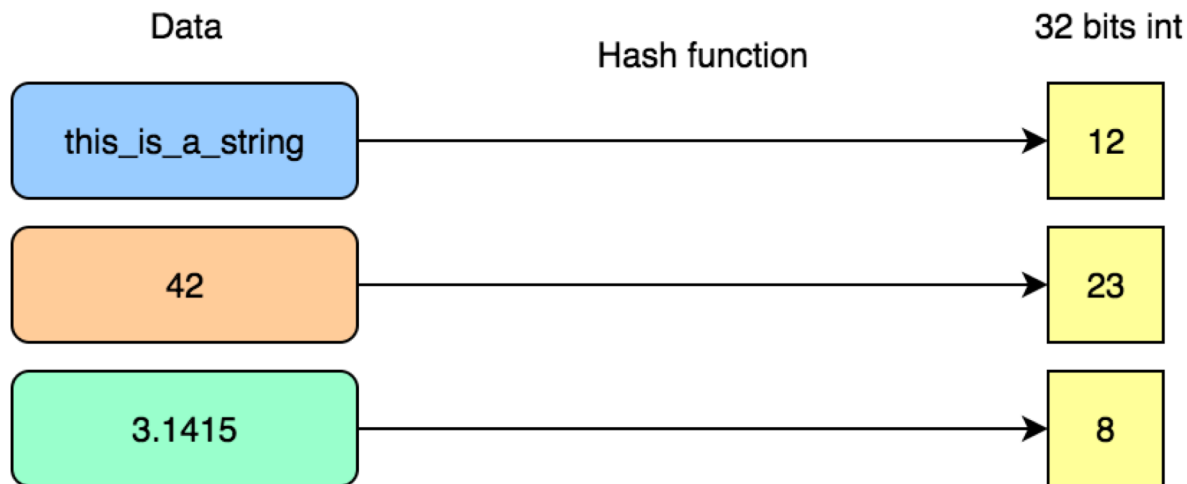  - Output: Yes | No
  - Work due to 17-11- 2019 23:59

# Bloom filter

# Hash function

- A **hash function** is any function that can be used to map data (string, integer, float …) of arbitrary size to fixed-size values.
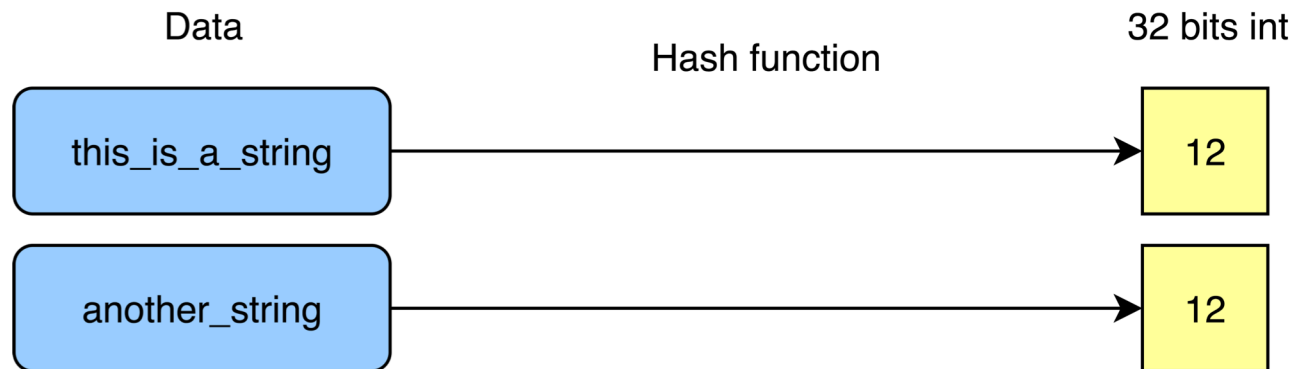
Example:

# Hash function

- **Uniformity**

A good hash function should map the expected inputs as evenly as possible over its output range. That is, every hash value in the output range should be generated with roughly the same probability. This is an important criteria to reduce collisions.
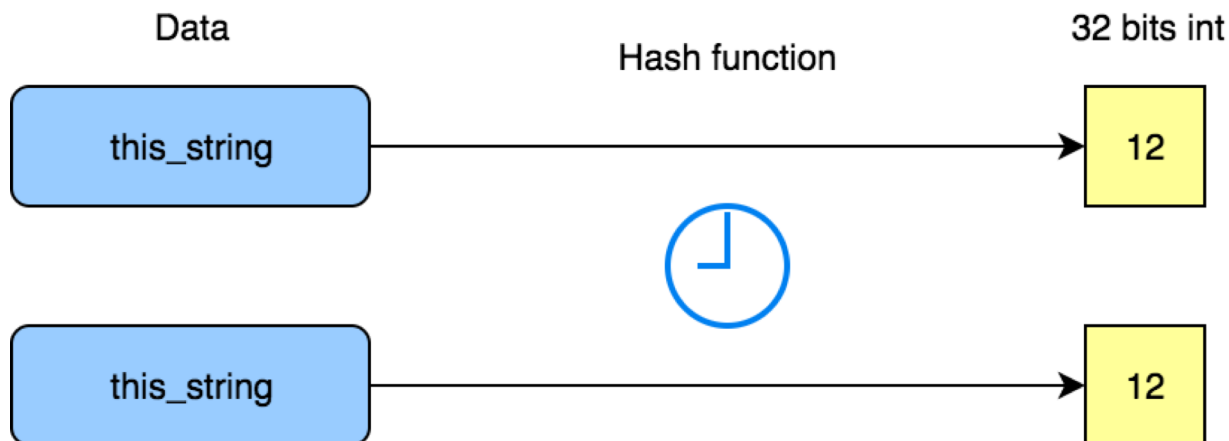
- **Collisions**

Data | Hash function | 32 bits int

this_is_a_string → 12

another_string → 12

Institut Mines-Télécom

TELECOM SudParis

# Hash function

■ **Deterministic**

A hash procedure must be **deterministic**, meaning that for a given input value it must always generate the same hash value.

# Hash function

■ **Applications**

- Hash tables
- Bloom filter
- Pseudonymization
- A/B test engine

# Bloom filter

- A **Bloom filter** is a space-efficient probabilistic data structure used to test whether an element is a member of a set. False positive matches are possible, but false negatives are not. In other words, a query returns either "possibly in set" or "definitely not in set". The more elements that are added to the set, the larger the probability of false positives.
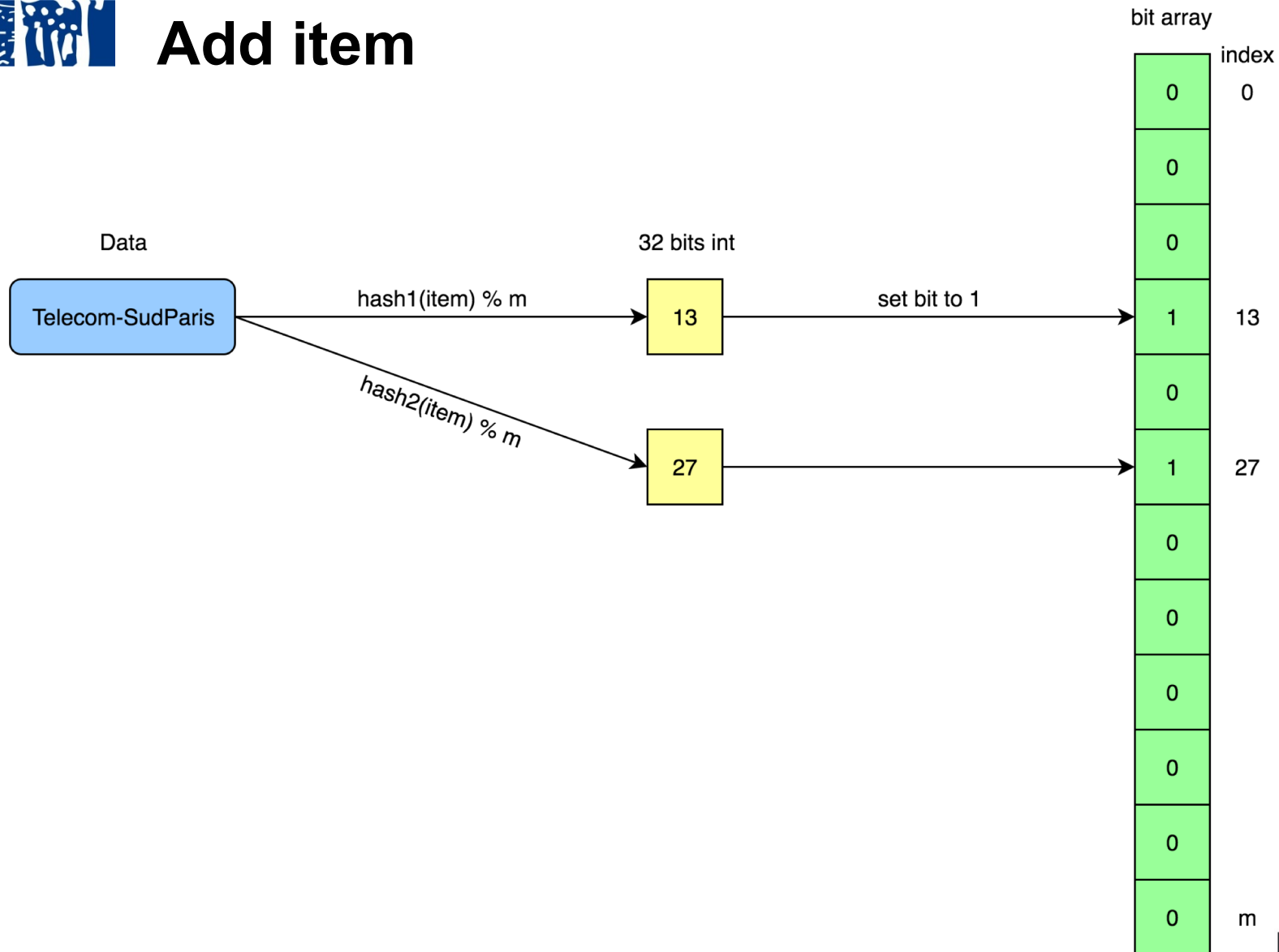
# Bloom filter

- Structure: **bit array** of *m* bits, all initialized to 0

- *k* different **hash functions** defined, each of which maps or hashes some set element to one of the *m* array positions. Example: f(item) = hash(item) % m

# Add item

Data

Telecom-SudParis

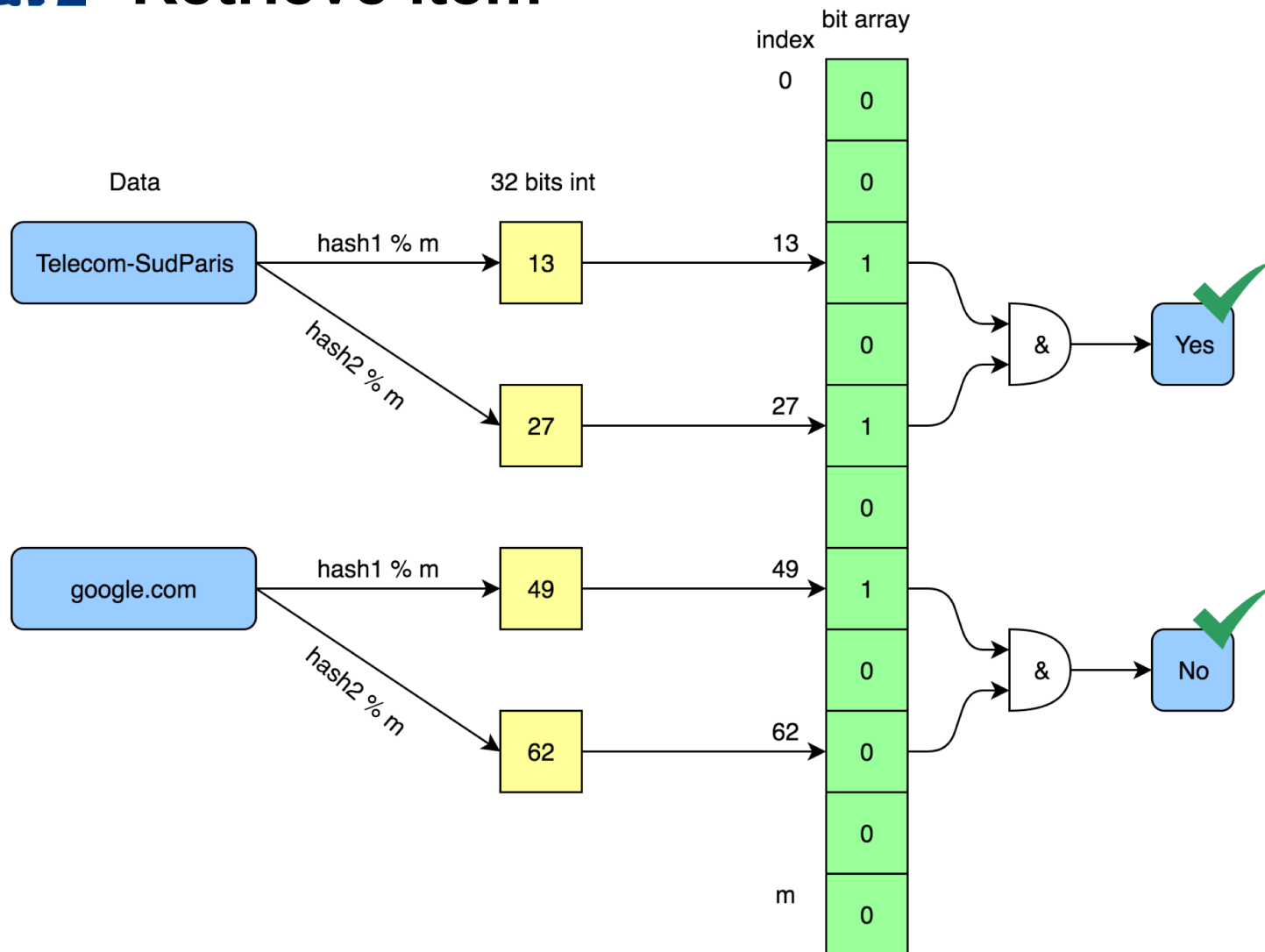32 bits int

bit array

index

hash1(item) % m → 13 → set bit to 1 → 1 (index 13)

hash2(item) % m → 27 → 1 (index 27)

| | index |
|---|---|
| 0 | 0 |
| 0 | |
| 0 | |
| 1 | 13 |
| 0 | |
| 1 | 27 |
| 0 | |
| 0 | |
| 0 | |
| 0 | |
| 0 | |
| 0 | m |

# Add item: collision



bit array

Data

32 bits int

**Telecom-SudParis** → hash1(item) % m → 13 → set bit to 1 → 1 (index 13)

hash2(item) % m → 27

**Veepee.com** → hash1(item) % m → 27

hash2(item) % m → 49

Bit array values (top to bottom): 0 (index 0), 0, 0, 1 (index 13), 0, 1 (index 27), 0, 1 (index 49), 0, 0, 0, 0 (m)
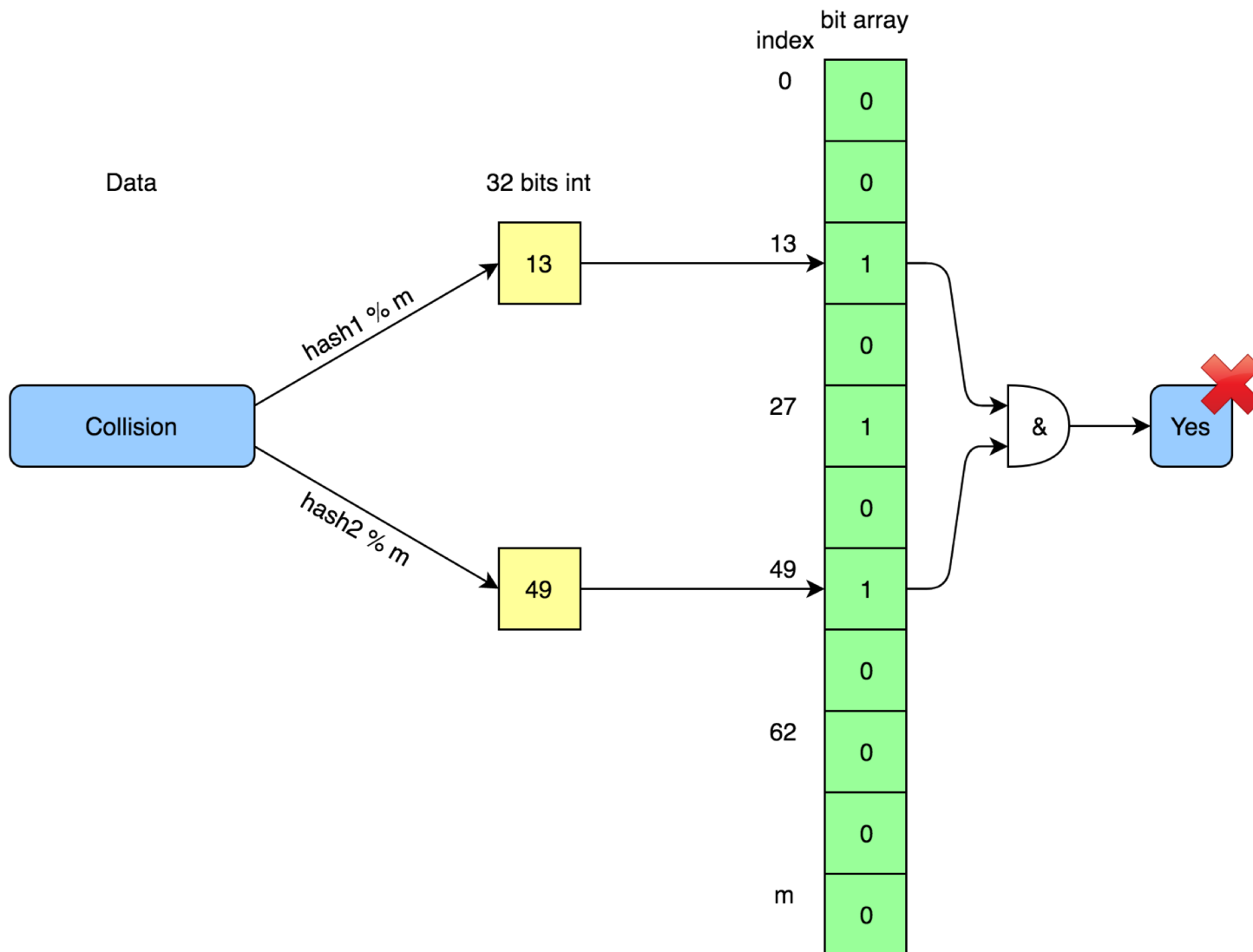
# Retrieve item

# Retrieve item: collision

# Error rate

■ **False negatives rate**          **0**


■ **False positives rate**

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

m: bit array size
k: number of hash functions
n: input cardinality

TELECOM
SudParis

# Complexity

- **Time**

$$O(n)$$

- **Memory**

$$O(m)$$

TELECOM
SudParis

# Advantages

- 1 unique pass through the data

- Fast processing

- Drastic dimension reduction

# Applications

- Cache filtering (solve the "one-hit-wonders » issue)

- Malicious sites listing

- Data base item existence

TELECOM
SudParis

# Mini-project 2

- **What is the number of unique (domain name, page title) couples visited the 01-08-2016 ?**
  - Output: Input cardinality
  - Work due to 25-11- 2019 23:59