# Feature Engineering for Texts
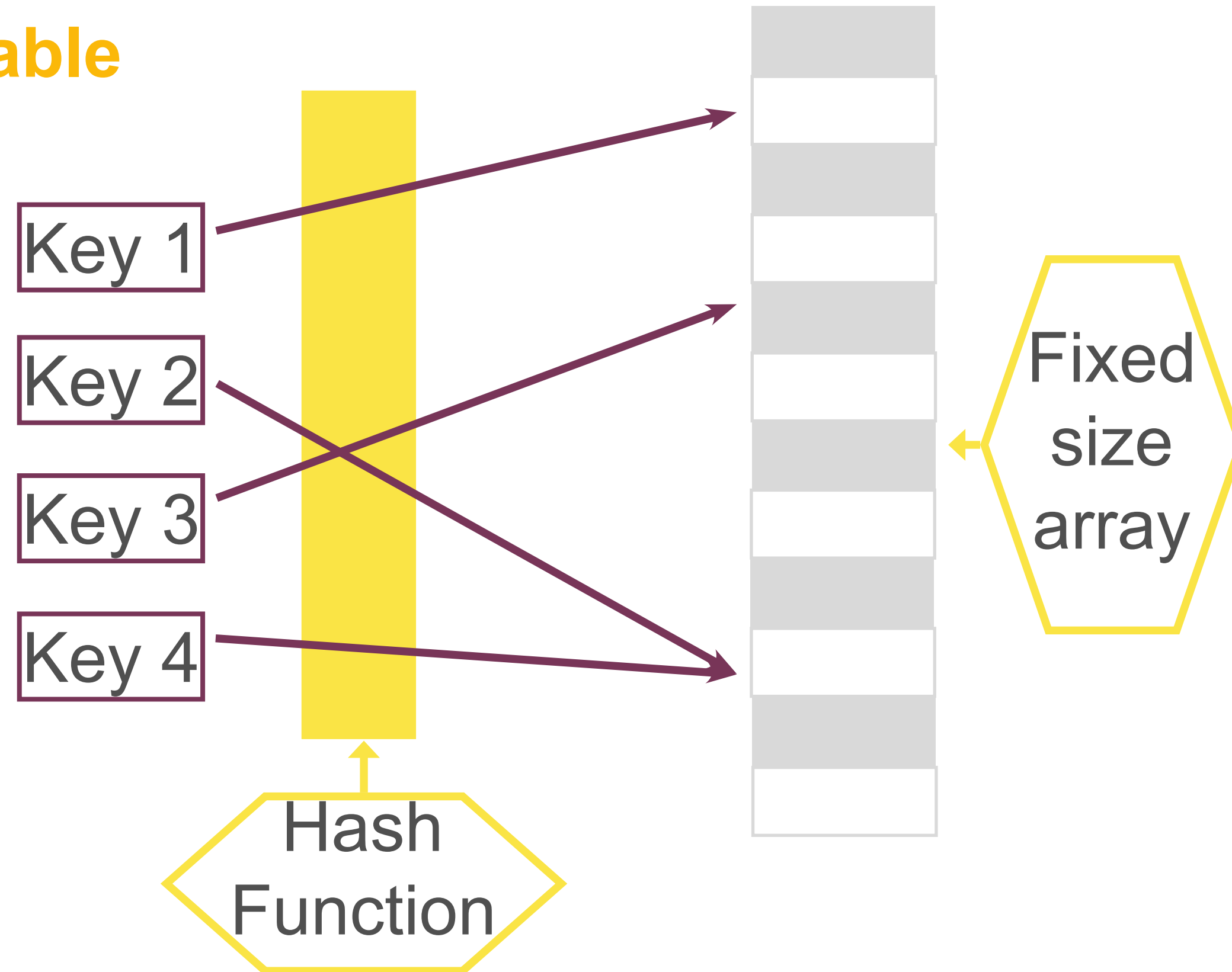## Part 2

# Hashing trick

## Hash Table

**Hash function  f(s):**

f("this") = 2

f("apple") = 10

f("fruit") = 5

…

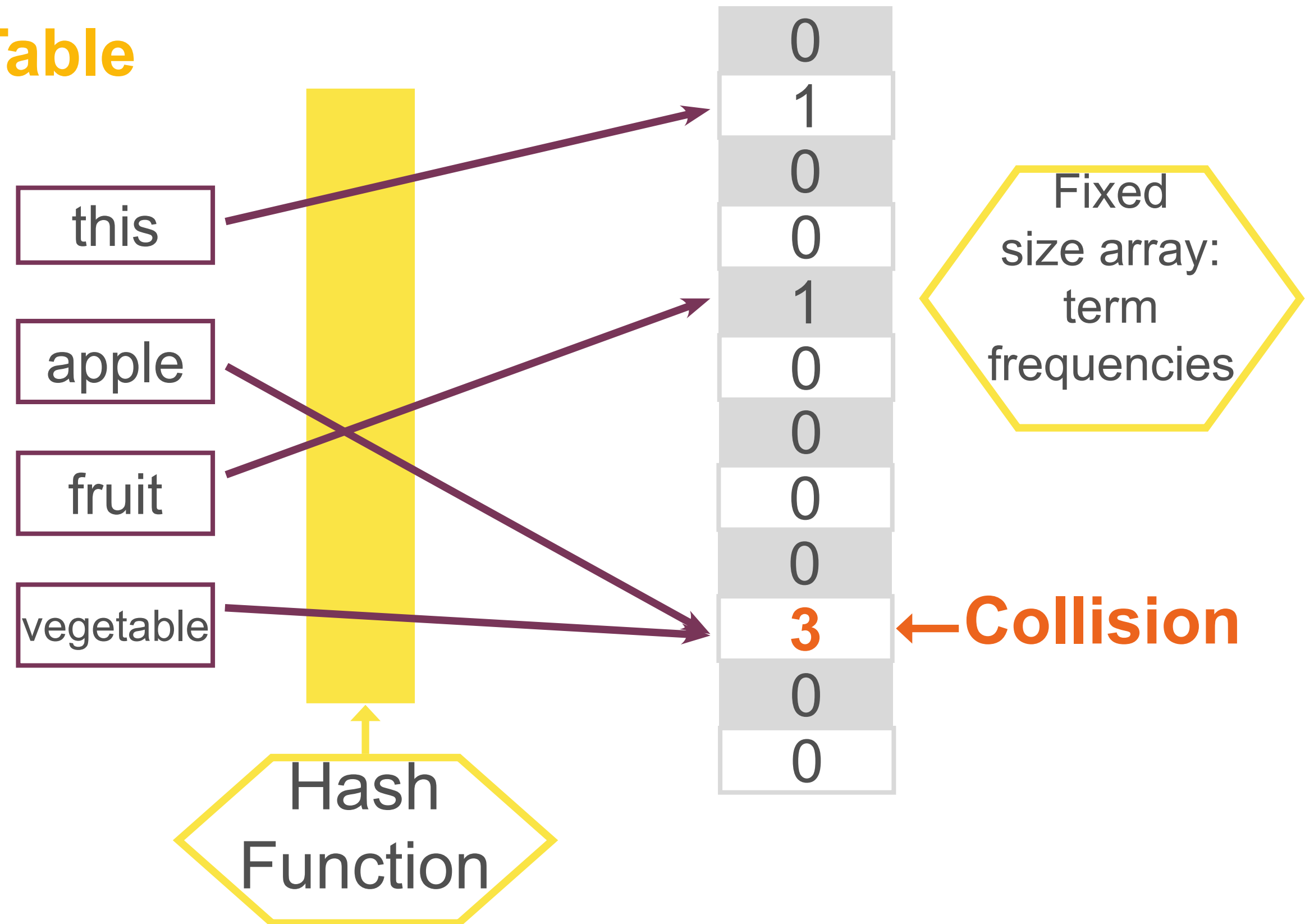**Common hash functions:**

- MurMur3 hash
- Jenkins hash
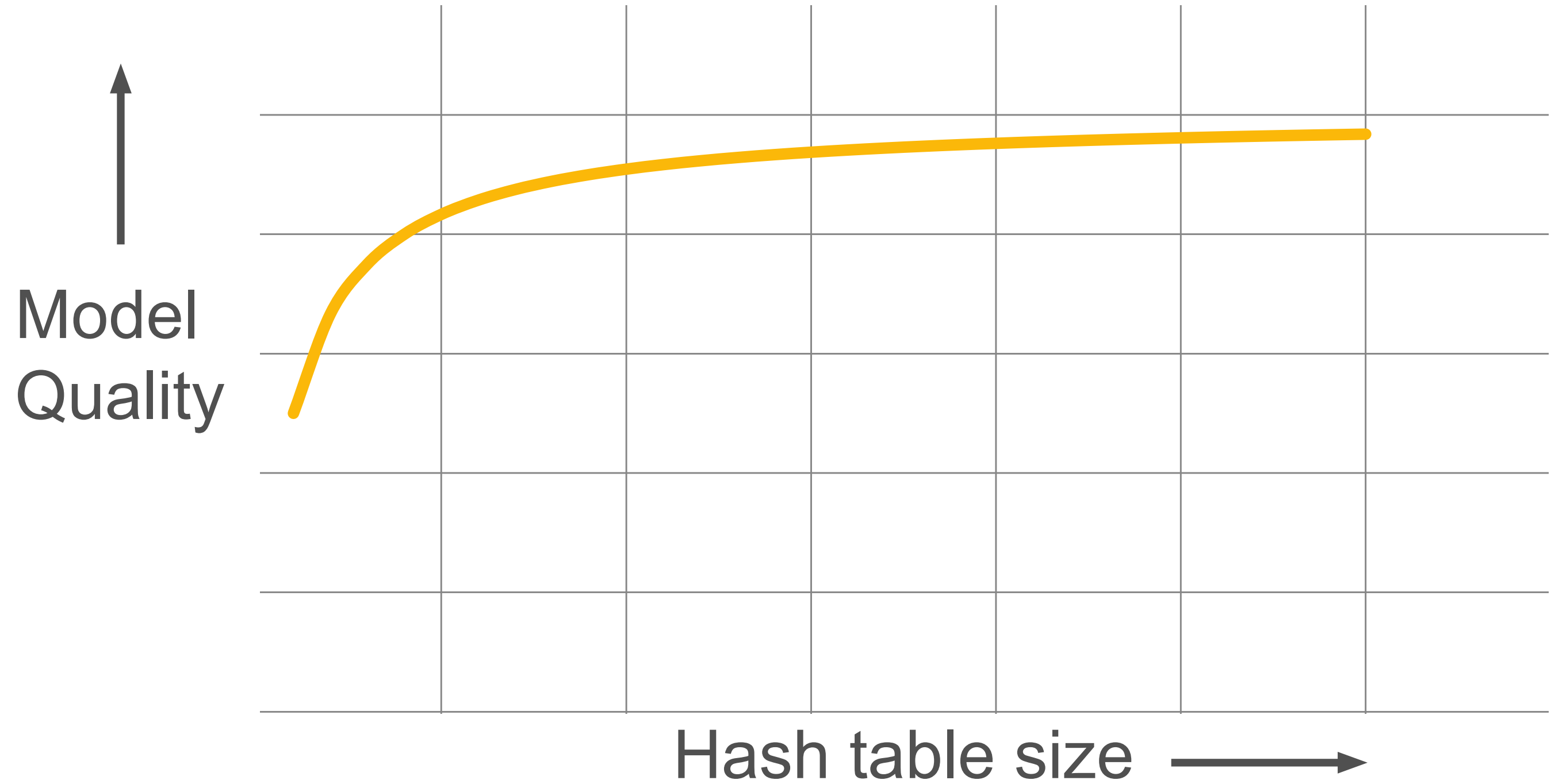- CityHash
- md5 hash

# Hashing trick

## Hash Table

| |
|---|
| 0 |
| 1 |
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 3 |
| 0 |
| 0 |

this

apple

fruit

vegetable

Hash Function

Fixed size array: term frequencies

← **Collision**

How to find an index of a word in the feature vector?

$$[0, \mathbf{0.01}, 0, \ldots, 0, \mathbf{0.2}, 0, \ldots, \mathbf{0.05}, 0, \ldots, 0]$$

this

apple     vegetable

fruit

Dictionary: **lookup** in dictionary
Hashing trick: calculate the **hash function** of the word

Hash table size vs. Model quality

# Dictionary vs. Hashing trick

| Dictionary | Hashing trick |
|---|---|
| No collisions | **Collisions** |
| **Need to store dictionary for learning and in production** | No dictionary<br>Calculations on-the-fly |
| **Slow if dictionary is large – O(log \|D\|)** | Fast – O(1) |
| **Feature vector size = unique words and n-grams count**<br>**Variable memory footprint** | Feature vector size = size of the hash table (fixed)<br>Fixed memory footprint |