Lambda architectures
Probabilistic Data Structures

λ vs ϰ

# Big Data Definition

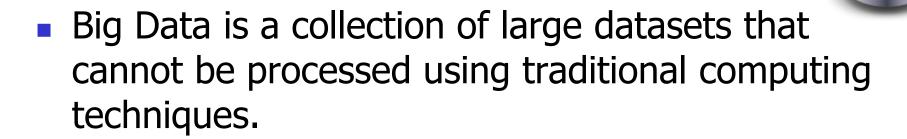- Big Data is a collection of large datasets that cannot be processed using traditional computing techniques.
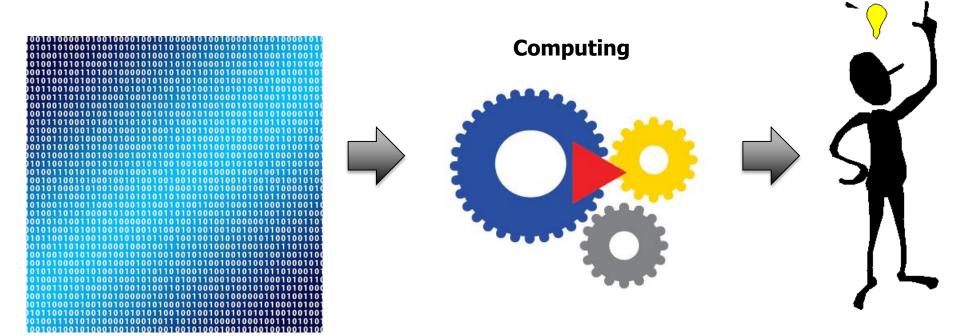
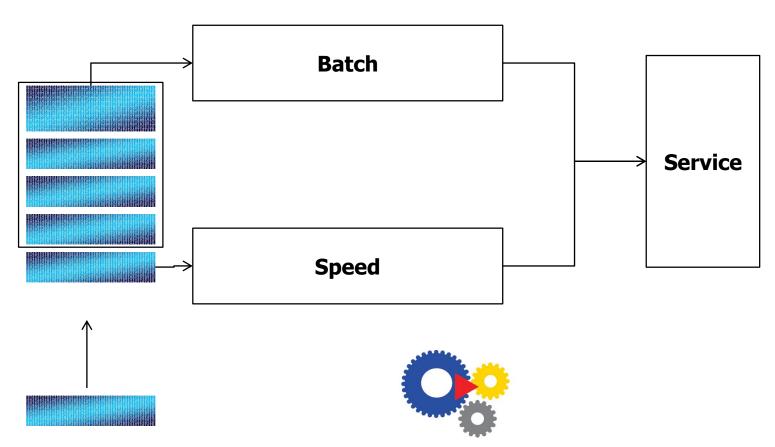- Any data that can challenge our current technology

# Big Data

> **Big data is about the Insight that we want to extract from information.**

**Data**

**Computing**

# Lambda Architecture
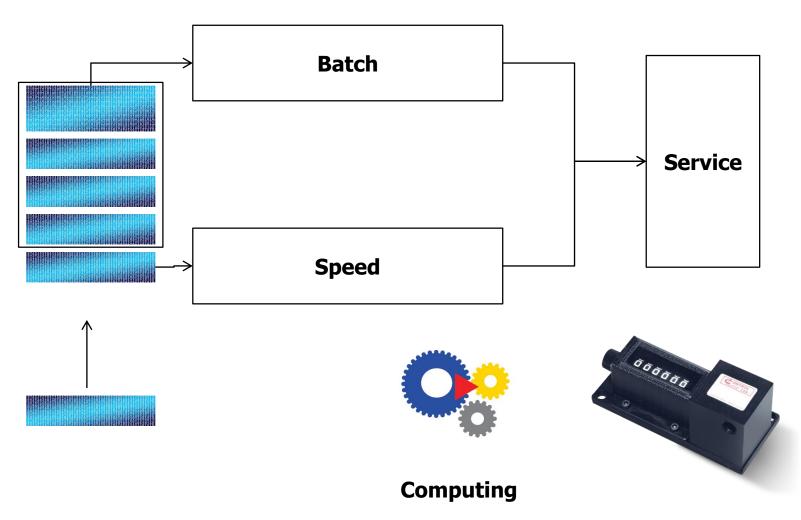
**Batch**

**Speed**

**Service**

**Computing**

# Lambda Architecture

**Hadoop MapReduce**

**Batch**

**Spark**

**Speed**

**Service**

**Hbase**

**Computing**

# Lambda Architecture

**Batch**

**Speed**

**Service**

**Computing**

# Lambda architecture

- **1) Speed/event layer :**
  - Stream processing architecture.
  - The data are appended to the master data periodically over a time window.
  - Real time views
- **2) Batch layer :**
  - Manages the master dataset:
    - Immutable, append-only raw data and pre-computing batch views.
  - Batch views
- **3) Serving layer:**
  - Merges batch and real-time views.

# Lambda architecture

- **Offline or batch processing : insight**
    - No time constraint : no data structures
    - Data are stored on disk
    - Data Lakes

- **Streaming data : near real-time processing**
    - Collecting analytics
    - Monitoring
    - Alerting,
    - Updating search indexes
    - Caches.

# Data Stream

IOT ⟹

about | after | again | air | all | along | also | an | and | another | any | are | around

# Data Stream

IOT

around
are
any
another
and
an
also
along
all
air
again
after
about

# Data Stream



I O T

about | after | again | air | all | along | also | an | and | another | any | are | around | as | at | away | back | be | becaus | been | before | below | betwee | both | but | by | came | can | come | could | day | did | differen | do | does | don't | down | each | end | even | every | few | find | first | for | found | from | get | give | go | good | great | had | has | have | he | help | her | here | him | his | home | house | how | I | if | in | into | is | it | its | just | know

# Data Stream

I
O
T

about | after | again | air | all | along | also | an | and | another | any | are | around | as | at | away | back | be | becaus | been | before | below | betwee | both | but | by | came | can | come | could | day | did | differen | do | does | don't | down | each | end | even | every | few | find | first | for | found | from | get | give | go | good | great | had | has | have | he | help | her | here | him | his | home | house | how | I | if | in | into | is | it | its | just | know

# Data Stream

# Data Stream

# Data Stream



Who ?
How Many ?
What ?
Together ?
How many of them ?
Why ?

APPROXIMATE IPV4 ADDRESS SPACE USAGE BY YEAR

# Speed/event layer

- Online, real-time processing
- Collecting real-time analytics
- In stream processing
- Large (Big) data streams :
  - Impractical to store that data in memory
  - Trade off to get real-time insight

- Stream Processing with Probabilistic counting
  - In-stream counters

# Advantages Probabilistic Data Structures

- ## Use small amount of memory
  - ### Can be controled
- ## Can be easily parallelizable
  - ### Hashes are independent
- ## Have constant query time
  - ### Not even amortized constant like in dictionary

# Usage

- **Computation Of Advanced Metrics :**
    - Number Of Unique Visitor
    - Most Frequent Items
- **Web Crawlers:**
    - Process A Stream Of Urls And Documents Tp Produce Indexed Content.
- **Websites:**
    - Process A Stream Of Page Views And Update Various Counters And Gauges.
- **Email:**
    - Process A Stream Of Text And Produce A Filtered, Spam-free Inbox.

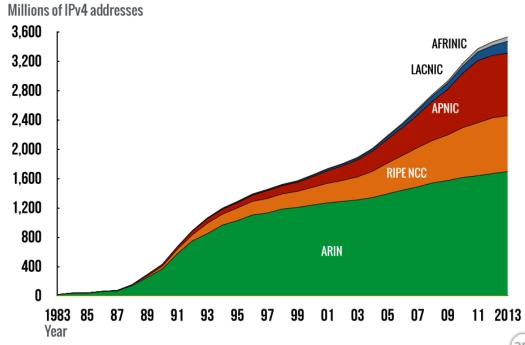# Exemple

- How to count the cardinality of IP addresses which access our website ?

- There are over four billion possible distinct IP V4 addresses.

## APPROXIMATE IPV4 ADDRESS SPACE USAGE BY YEAR

Millions of IPv4 addresses

AFRINIC
LACNIC
APNIC
RIPE NCC
ARIN

3,600
3,200
2,800
2,400
2,000
1,600
1,200
800
400
0

1983  85  87  89  91  93  95  97  99  01  03  05  07  09  11  2013
Year

ars

# Probabilistic data structures

- Probabilistic data structures can not provide definite answer

- Instead they provide a reasonable approximation of the answer and a way to approximate this estimation.

- Useful for big data and streaming application because they allow to decrease the amount of memory needed with comparison to data structures that give you exact answers.

# Membership

**Is this guy on board ?**

**How many of them are there ?**

# Presence evaluation

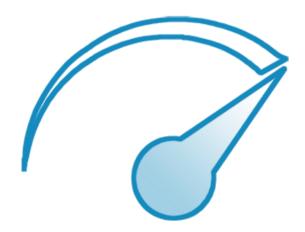- How to evaluate the presence of the element in a stream ?

# Frequency evaluation

- How to store a numerical value associated with each element ?

- How to evaluate the number of occurrences of the element in a stream ?

# Hash Functions

- Probabilistic data structures use hash functions to randomize items.

- They ignore collisions to keep the size constant

- But this is a reason why they can not give exact values.

# Hash Usage : Set Membership, Counting

SSN : Social Security Number

**N element Array**

999999999999999

**Is this guy on board ?**

| steve |
|-------|
|       |
|       |
|       |
| mike  |
|       |
|       |
| john  |
|       |
|       |
| bob   |
|       |
|       |
|       |
|       |
|       |
| max   |
|       |
|       |
| edward|

125675798988090

**KEY**

000000000000000

99

| steve |
|-------|
| bob   |
| john  |
|       |
| mike  |
| edward|

00

**hash**

**addresses**

**value**

# Hash Usage : Set Membership, Counting

SSN : Social Security Number

**Is this guy on board ?**

**N element Array**

**Probe O(1)**

| | |
|---|---|
| steve | 99 |
| | |
| bob | |
| john | |
| | |
| mike | |
| edward | 00 |

**addresses**

**hash**

**value**

| | |
|---|---|
| steve | 999999999999999 |
| | |
| | |
| | |
| mike | |
| | |
| | |
| john | 125675798988090 |
| | |
| | |
| bob | **KEY** |
| | |
| | |
| | |
| max | |
| | |
| | |
| edward | 000000000000000 |

# Hash Usage : Set Membership, Counting

SSN : Social Security Number

**How many of them are there ?**

**N element Array**

**Probe O(1)**

| | |
|---|---|
| steve | 99 |
| | |
| bob | |
| john | |
| | |
| mike | |
| edward | 00 |

**addresses**

**hash**

**value**

| steve | 999999999999999 |
|---|---|
| mike | |
| john | 125675798988090 |
| bob | **KEY** |
| max | |
| edward | 000000000000000 |

**Emmanuel fuchs Architectures des Systèmes de Bases de Données**

# Hash Usage : Set Membership, Counting

SSN : Social Security Number

**How many of them are there ?**

**N element Array**

steve

999999999999999

**Probe O(1)**

| | |
|---|---|
| steve | 99 |
| | |
| bob | |
| john | |
| | |
| mike | |
| edward | 00 |

mike

john

125675798988090

**addresses**

**hash**

bob

**KEY**

**value**

max

edward

000000000000000

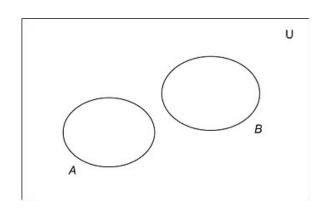# Hash Usage : Set Membership, Counting
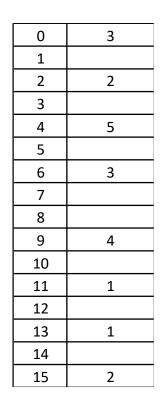
SSN : Social Security Number

**Is this guy on board ?**

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

# Hash Usage : Set Membership, Counting

SSN : Social Security Number

**How many of them are there ?**

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

SSN : Social Security Number

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

**M**

**What about memory size ?**

# Exemple with alphabet letters

## The set

| A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |

# Exemple with alphabet letters

**The set**

| A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |

**Bit Map : Probalistic Data Structure**

| 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 |

# Exemple with alphabet letters

**The set**

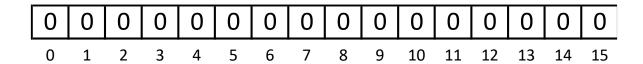| A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |

**Size = 27 * Short = 432 bytes**

**Probalistic Data Structure**

| 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 |

**Size = 2 bytes**

# Hash set

**Hash(A)**

**H(A)**

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

# Hash set

**Hash(B)**

**H(B)**

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

# Hash set

**Hash(Y)**

**H(Y)**

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

# Hash set

**Hash(E)= Hash(Y)**

H(E)   H(Y)

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

**False**

**-** **?** **+**

# Hash set

**Hash (K) = Hash (B)**

H(K)    H(B)    H(E)    H(Y)

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

**False**

**?**

-    +

# Fault Positive Prevention

- Two Hash functions : H1(), H2()
- Two Hash sets

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

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

# Fault Positive

**Collision**

H2(E)

H1(E)

H1(Y)

H2(Y)

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

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

# Fault Positive

**Collision**　　**If Y not In**

H2(E)

H1(E)

H1(Y) = collision with E

H2(Y)

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

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

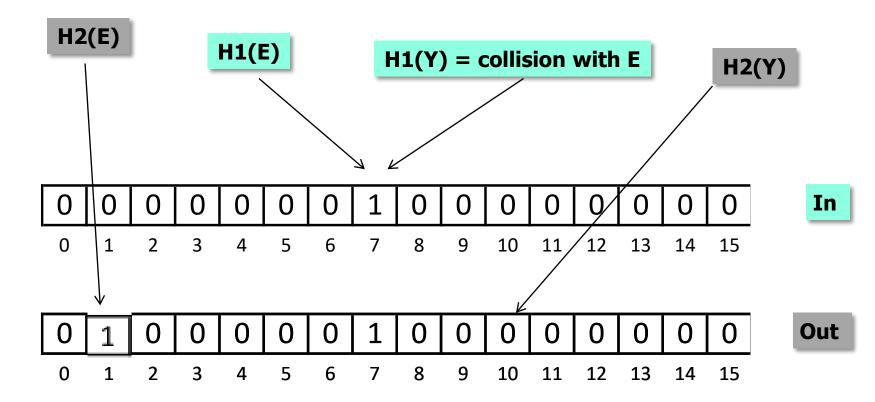**Y is definitiveley Out**

# Usual probabilistic data structures

- Bloom Filters
- Count-min Sketch
- Count-Mean-Min Sketch
- Hyperloglog