# BIG DATA AND ANALYTICS

Daniel Gillblad, dgi@sics.se

Swedish Institute of Computer Science

SWEDISH ICT

SICS

# OUR RESEARCH: SURVIVING THE DATA FLOOD

# SOME QUESTIONS AROUND BIG DATA ANALYTICS

1.  What is Big Data?

2.  What is different this time?

3.  Why is it important?

4.  What does the future look like?

# BIG DATA – 3V, AND MORE?

SWEDISH ICT   SICS

# BIG DATA IS PARALLELIZATION



Read on 10 000 machines:
10 000 times faster

# BIG DATA IS SCALABILITY
## ON COMMODITY HARDWARE
### LOTS OF COMMODITY HARDWARE…
### …BIG DATA IS FAULT TOLERANCE

# BIG DATA IS PLATFORMS AND MIDDLEWARE

# FIRST, THERE WAS MAP REDUCE

- Distribute data over commodity hardware (HDFS etc.) in data center

- Map and distribute computation to computers storing the data

- Choose fix, batch-oriented form of calculation: Map -> Reduce

# MAP REDUCE, EXAMPLE

- Word count: "read -> map -> reduce -> write"

```
...
public static class Map extends MapReduceBase implements Mapper<LongWritable, Text, Text, IntWritable> {
    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();

    public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
      String line = value.toString();
      StringTokenizer tokenizer = new StringTokenizer(line);
      while (tokenizer.hasMoreTokens()) {
        word.set(tokenizer.nextToken());
        output.collect(word, one);
      }
    }
}

public static class Reduce extends MapReduceBase implements Reducer<Text, IntWritable, Text, IntWritable> {
    public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text, IntWritable> output, Reporter reporter)
    throws IOException {
      int sum = 0;
      while (values.hasNext()) {
        sum += values.next().get();
      }
      output.collect(key, new IntWritable(sum));
    }
}
...
```
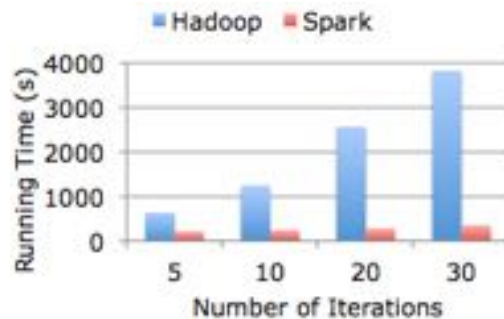
SWEDISH
ICT

SICS

# NEXT-GEN MAP REDUCE

- Handle real-time, interactive and iterative tasks: Necessary for machine learning etc.

- Resilient distributed data sets – in-memory caching etc.

1. Transformations: Map, filter, join…

2. Actions: Count, collect, save…

- Faster, easier, more powerful

```
file = spark.textFile("hdfs://...")

file.flatMap(line  => line.split(" "))
    .map(word => (word, 1))
    .reduceByKey(_  + _)
```

# SPARK, REAL-WORLD EXAMPLE

- Incomprehensible example:
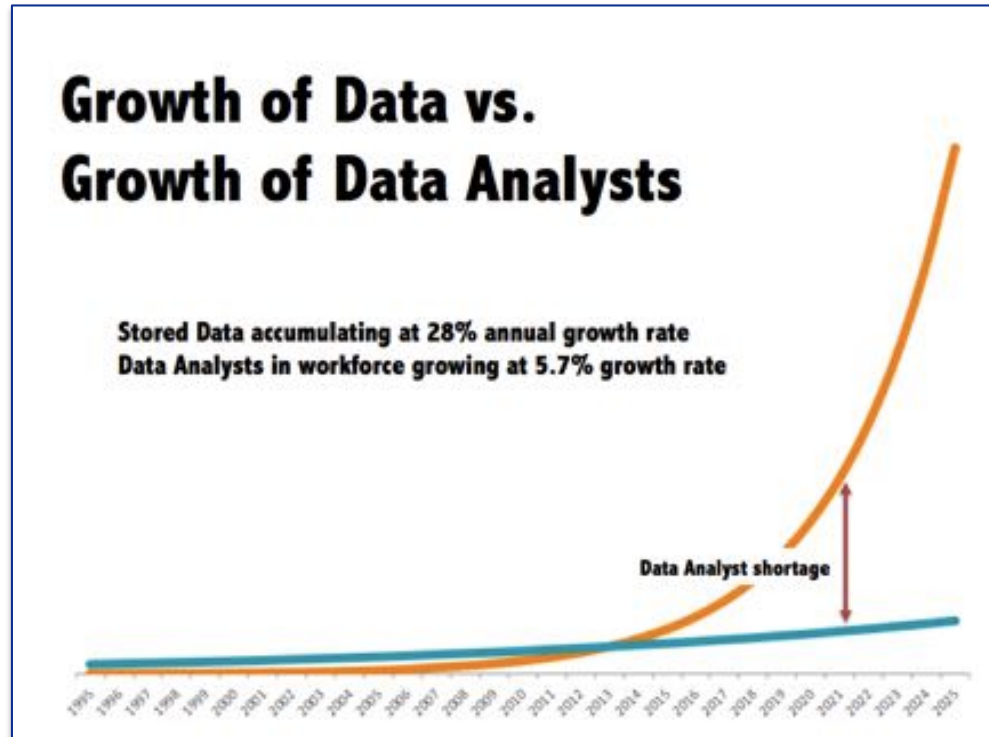
```
logs.map(log => (log.data, log)).groupByKey()
        .filter(_._2.size >= minimumValue)
        .map(xLogsPair => new Aggregate(Pair(xLogsPair._1,
          logsExpanded(xLogsPair._2))))
```

- Defines both *what/how* and hints on *how to parallelize*

SWEDISH ICT    SICS

# IS THAT IT?

- Still not that easy to write – special kind of developers
- Fantastic productivity in the right hands
- Developer support limited
- Parallelization is difficult
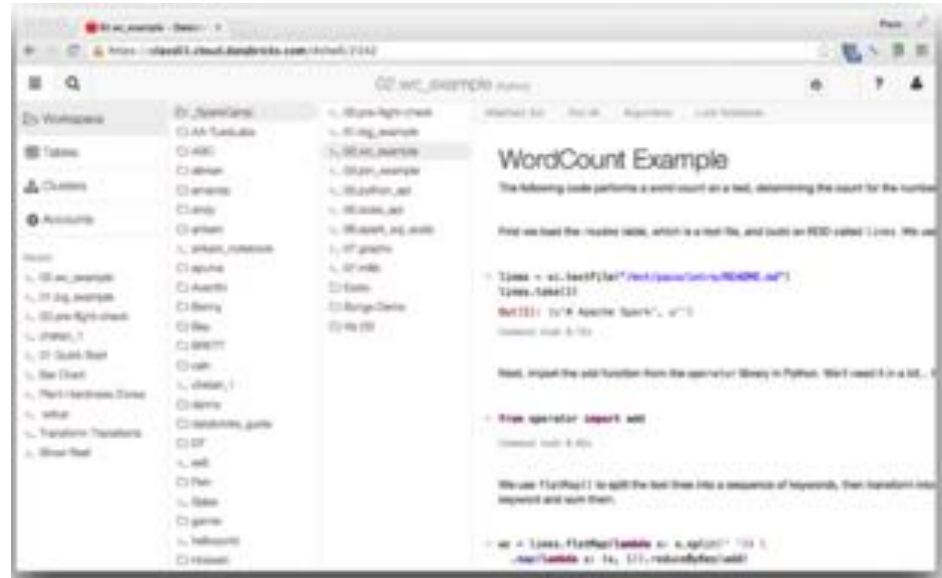- Pretty difficult to write advanced analytics

SWEDISH ICT SICS

# WHO USES THIS?



**Growth of Data vs.
Growth of Data Analysts**

Stored Data accumulating at 28% annual growth rate
Data Analysts in workforce growing at 5.7% growth rate

Data Analyst shortage

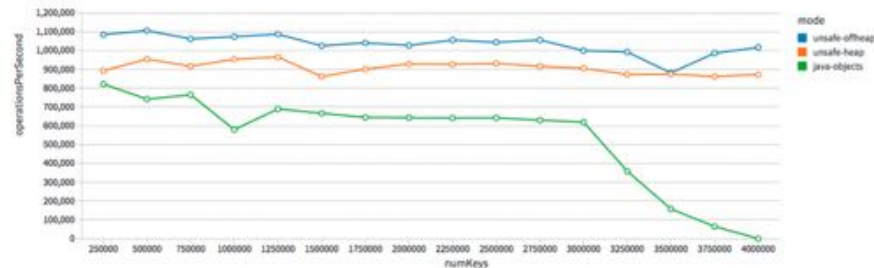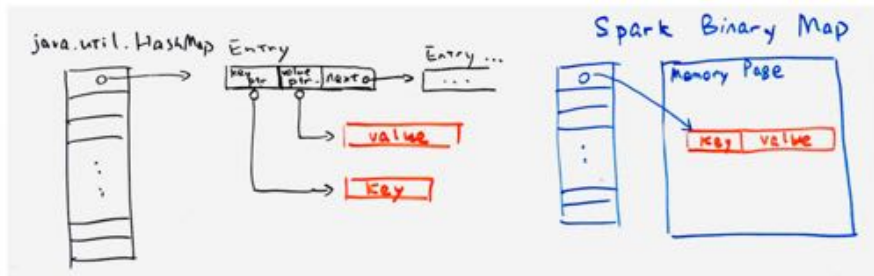http://www.delphianalytics.net

SWEDISH ICT    SICS

# SCALABILITY IS OFTEN NO LONGER THE MAIN ISSUE – ANALYTICS, EFFICIENCY AND PRODUCTIVITY ARE

# ACCESSIBLE BIG DATA TOOLS?

- Higher-level abstractions (RDDs -> Dataframes etc.)
- R and Python interfaces
- Libraries:
  - Spark MLib, GraphX
  - FlinkML, Gelly
  - …
- Higher level languages
- Imperative Big Data processing
- …

# RAW PERFORMANCE MATTERS



https://databricks.com/blog/2015/04/28/project-tungsten-bringing-spark-closer-to-bare-metal.html

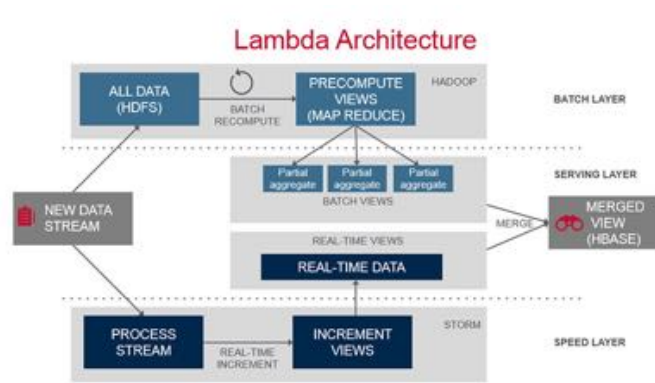http://bid2.berkeley.edu/bid-data-project/

SWEDISH ICT    SICS

# DATA AT REST…

# … TO DATA IN MOTION

# FROM LAMBDA ARCHITECTURE…



https://www.mapr.com/sites/default/files/otherpageimages/lambda-architecture-2-800.jpg

# …TO UNIFIED ARCHITECTURE

SWEDISH ICT    SICS
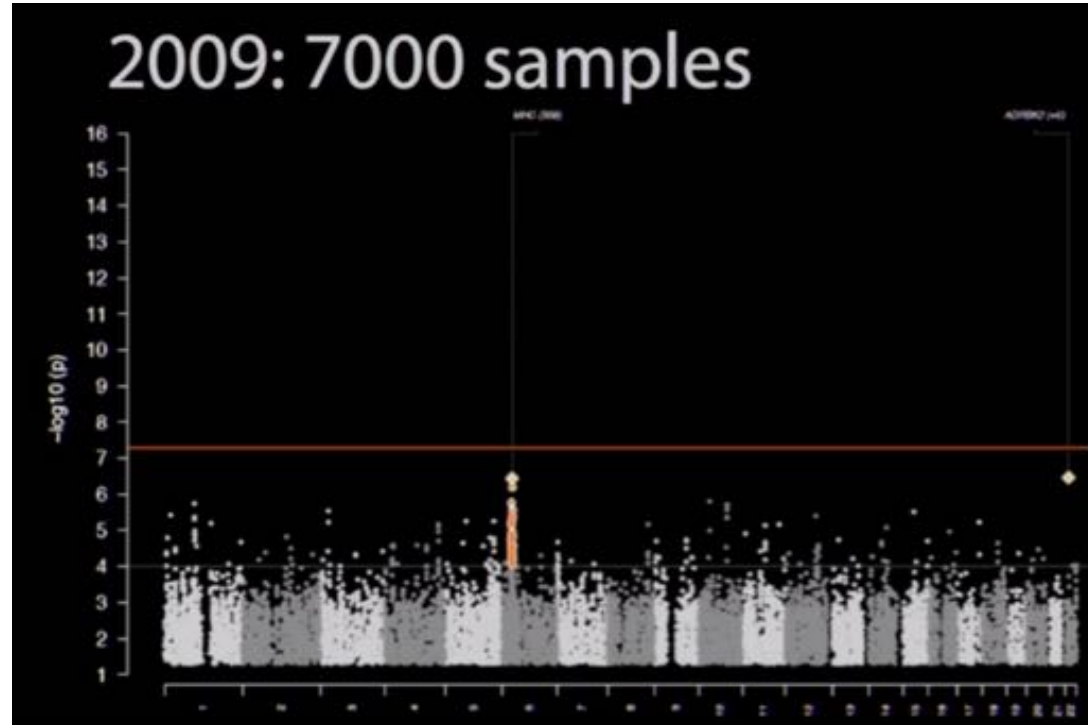
# WHAT IS DIFFERENT THIS TIME?

- Haven't we seen this before? AI, Data Mining etc.?
- Things are a bit different this time around:
    1. We have the computational power and middleware (cloud, Big Data middleware)
    2. We have the data (Human generated, IoT)
    3. We have the algorithms
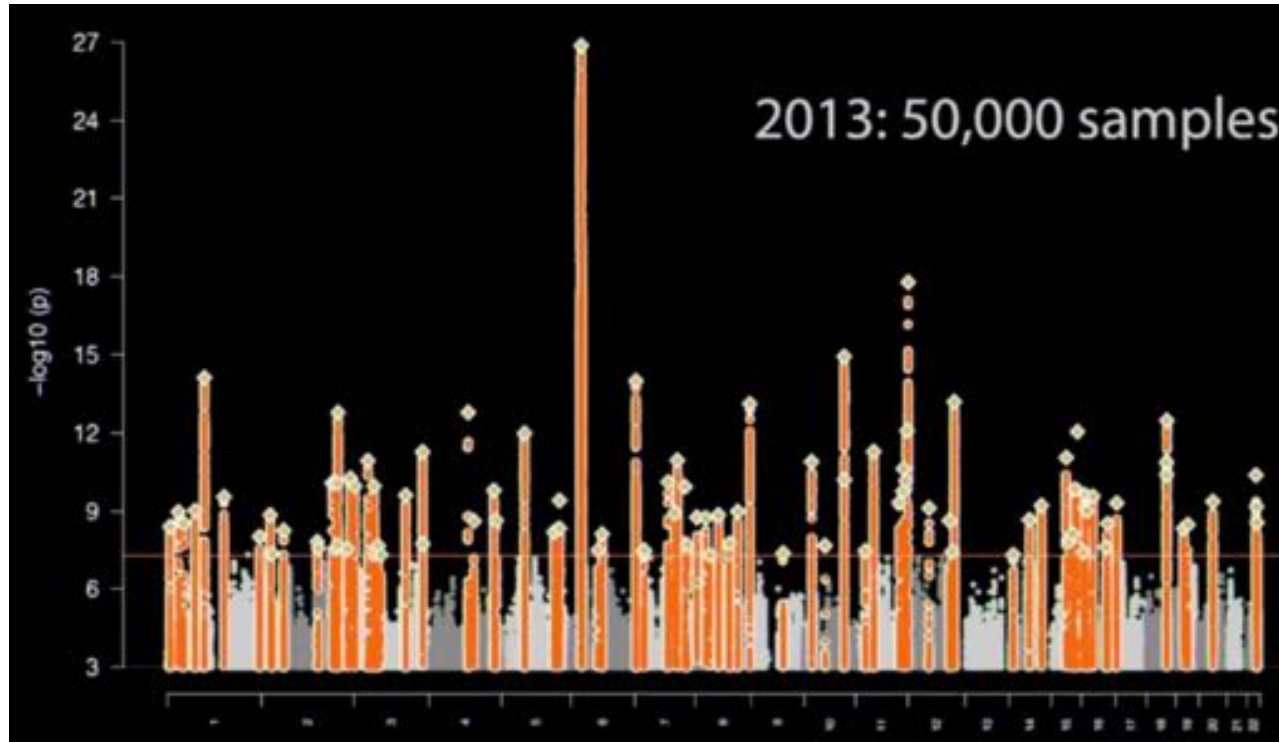
# THE UNREASONABLE EFFECTIVENESS OF DATA

- In a wide array of academic fields, the ability to effectively process data is superseding other more classical modes of research.

*"More data trumps better algorithms"**

*"The Unreasonable Effectiveness of Data" [Halevey et al 09]

SWEDISH ICT    SICS

# EXAMPLE: POPULATION SCALE GENOMICS



[Image source: Patterson, Fighting the Big C with the Big D, 2014]

# EXAMPLE: POPULATION SCALE GENOMICS



[Image source: Patterson, Fighting the Big C with the Big D, 2014]

# BIG DATA – A BIG MISTAKE?

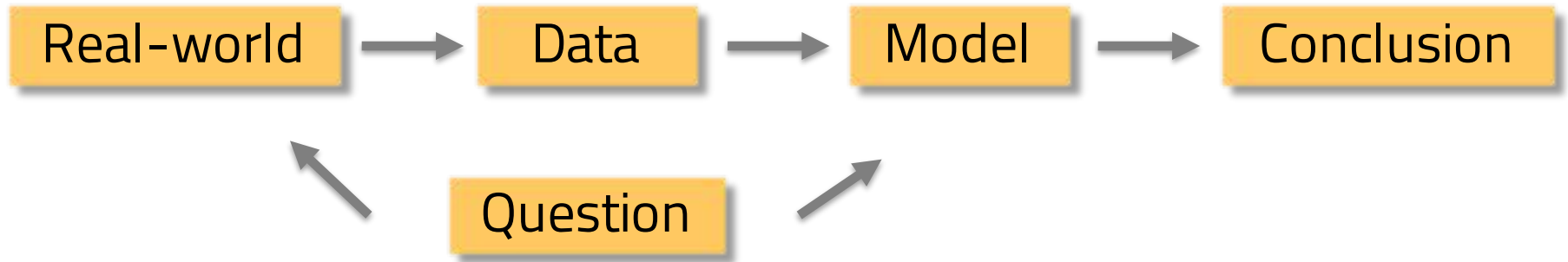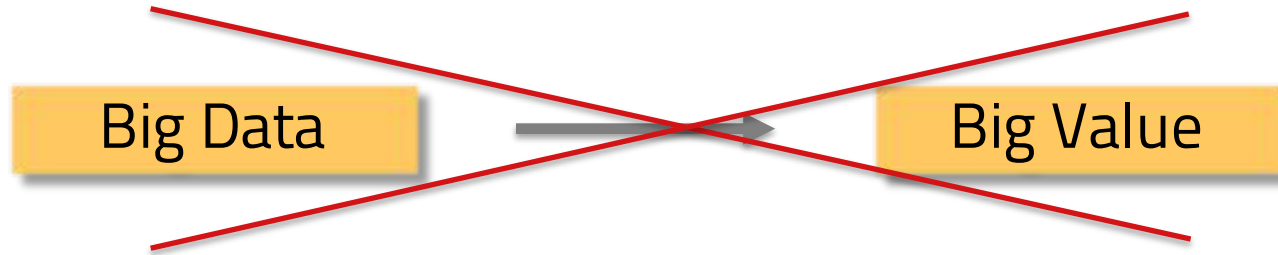- "Four oversimplified articles of faith:"
    1. Data analysis produces uncannily accurate results
    2. Every single data point can be captured, making old statistical sampling techniques obsolete
    3. It is passé to fret about what causes what, because statistical correlation tells us what we need to know
    4. "The End of Theory": With enough data, the numbers speak for themselves
- Practitioners do not necessarily believe this, and neither should you

From *Big Data: are we making a big mistake?*, Tim Harford, Financial Times March 28 2014

SWEDISH ICT    SICS

# BIG DATA TO BIG VALUE?

Big Data → Big Value

SWEDISH ICT    SICS

# BIG DATA TO BIG VALUE?

Big Data → Big Value

Real-world → Data → Model → Conclusion

Question

SWEDISH ICT    SICS

# BDA: WORKFLOW NOT FUNDAMENTALLY DIFFERENT

Acquisition

Cleaning

Representation

| Case-based | Statistical | Logical | Kernel-based |

Validation

Deployment

SWEDISH ICT    SICS

# EVERYTHING IS ABOUT MACHINE LEARNING ANYWAY…

*A Few Useful Things to Know About Machine Learning*
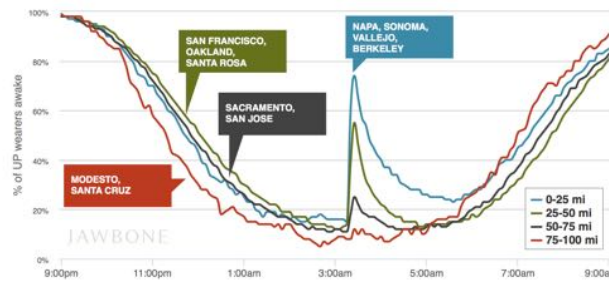Pedro Domingos
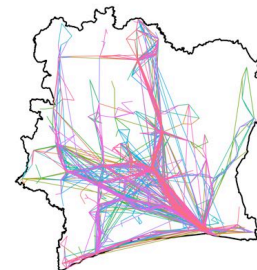CACM 55:10
http://dl.acm.org/citation.cfm?id=2347755

# LET'S START WITH THE OBVIOUS



Lindh, J et al: PLoS ONE 7(11), 2012.

It is a key component of digital services    It will help us understand the world better

# …BUT IT IS REALLY, REALLY DIFFICULT



Sample bias, data veracity



Model complexity

…+ user and analyst bias…

SWEDISH ICT    SICS

# SCALABLE, PRIVACY-PRESERVING MOBILITY ANALYTICS FROM NETWORK DATA

Urban planning

Traffic management

Network management

Crisis management

Consumer applications

A3, A4, A5, *C1*, *D1*, D2, *B4*, B1
F2, F3, F6, *G1*, *C5*, *D3*, *E2*, E1, *B5*, B2
F7, F5, F6, *G1*, G2, *D4*, E2, E3, *B6*

SWEDISH ICT    SICS

# AN UNSOLVABLE PROBLEM?

- Moving datasets can be difficult
    - Large and sensitive (business value, integrity issues, etc.)
    - Can allow for re-identification in anonymized data
- Combining several datasets outside their collection entity is key to many (public) applications, *but...*

| | | |
|---|---|---|
| ***Centralized datasets:*** | *Technically feasible* | *Politically hard* |
| ***Federated datasets:*** | *Technically difficult* | *Politically solvable* |

SWEDISH ICT

SICS

# BIG DATA: KEY COMPONENT IN AUTOMATION

# MASSIVE LEARNING SYSTEMS

SWEDISH ICT    SICS

# LARGER AND LARGER REPRESENTATIONS



**Convolution**
**Pooling**
**Softmax**
**Other**

# UNIFICATION OF MACHINE LEARNING?



*Alternating Direction Method of Multipliers*
S. Boyd, N. Parikh, et al.
http://stanford.edu/~boyd/papers/admm_distr_stats.html

# MOVING TO PROBABILISTIC APPROXIMATIONS



$10^7$ elements
$10^6$ distinct values
domain of 32-bit integers

**40 MB** Raw Data

**0.6 MB** Membership Query with 4% error – Bloom Filter

Exact Membership Query, Cardinality Estimation – Sorted IDs or Hash Table

**4 MB**

**48 KB** Frequences of top-100 most frequent elements with 4% error – Count-Min Sketch

**14 KB** Top-100 most frequent elements with 4% error – Stream-Summary

**7 MB**

$10^6$ pairs
{
  32-bit value,
  24-bit counter
}

**2 KB** Cardinality Estimation with 4% error – Loglog Counter

**125 KB** Cardinality Estimation with 4% error – Linear Counter

Exact Frequency Estimation, Range Query – Sorted Table or Hash Map

https://highlyscalable.wordpress.com/2012/05/01/probabilistic-structures-web-analytics-data-mining/

SWEDISH ICT          SICS

www.sics.se

# ANALYTICS EVERYWHERE

# ANALYTICS AS A SERVICE

# WHERE ARE WE GOING?

① Advanced analytics is moving towards large-scale Machine Learning

② Computation and storage platforms need to adapt and develop

⬇

*LearningMachines@SICS*

*Analytics and system development on real data and use cases*

SWEDISH ICT · SICS

# THE DATA DRIVEN SYSTEMS STACK, SICS



*Analytics*

| Diagnosis / classification | Anomaly detection | Clustering |

*Algorithms and frameworks*

| Graph algorithms | Scalable / On-line | Structural / Deep Learning |

*Data processing engines*

| Flink | Spark | Storm |

*Compute and resource management*

| Mesos | Yarn |

*Storage and streams*

| HDFS | HBase | Kafka |

*Hops* PaaS (http://hops.io)

*Networking*

| ICN, SDN | Autonomous RAN |

*Data collection*

| SicsthSense | SDN Monitoring | Text, Social Media |

○ *SICS Contributions*

SWEDISH ICT   SICS

www.sics

# DATA DRIVEN SYSTEMS, APPLICATION EXAMPLES



|  | Network analytics | Vehicles & Transport | eHealth |
|---|---|---|---|
| *Analytics* | Diagnosis / classification | Anomaly detection | Clustering |
| *Algorithms and frameworks* | Graph algorithms | Scalable / On-line | Structural / Deep Learning |
| *Data processing engines* | Flink | Spark | Storm |
| *Compute and resource management* | Mesos | | Yarn |
| *Storage and streams* | HDFS | HBase | Kafka |
| | *Hops* PaaS (http://hops.io) | | |
| *Networking* | ICN, SDN | | Autonomous RAN |
| *Data collection* | SicsthSense | SDN Monitoring | Text, Social Media |

○ *SICS Contributions*

SWEDISH ICT    SICS

www.sics.se

# EXPERIENCE FROM APPLICATIONS

- Big Data can be *made* small – consider the complete application

- Big Data can become small very quickly

- Beware of sample bias

- Distribute models, not data

- Distributed solutions, Statistical Machine Learning, and Bayesian statistics can help

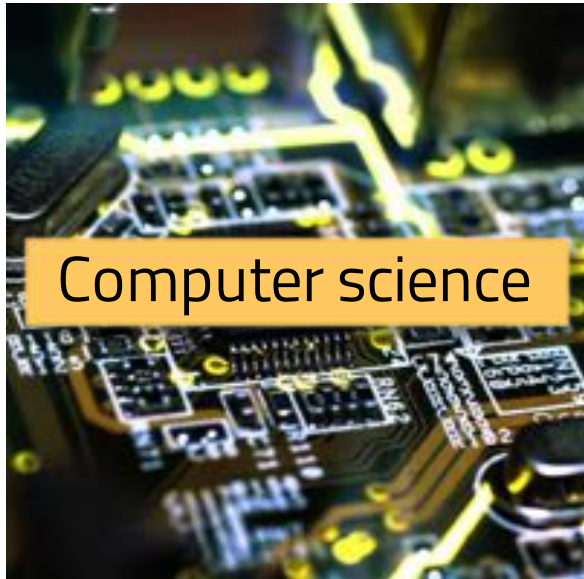SWEDISH ICT

SICS

# DISCUSSING APPLICATIONS: IT IS NOT OBVIOUS WHAT IS DIFFICULT

# WHAT ABOUT INTEGRITY?

- Not all Big Data data is integrity sensitive!
    - Media, measurements, science, …
- How data is (or potentially is) used is everything
- Surveillance and data driven services are very different
- It all comes down to trust
    - Transparency is key
    - Laws and regulations? How do we manage sales, sharing?

SWEDISH ICT  SICS

# THE END OF COMPUTER SCIENCE



Computer science → Data science

SWEDISH ICT | SICS

# …BUT JUST REMEMBER…

*We're not there yet – things are moving fast!*

# WWW.SICS.SE