

# Five Things...

## (we wish we had known)

British Gas Connected Homes

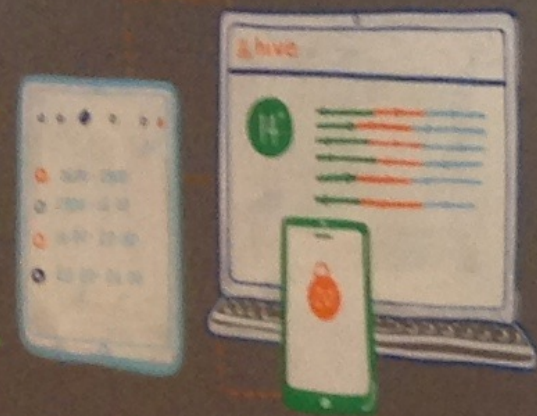
Josep Casals - Lead Data Engineer  
Jim Anning - Head of Data & Analytics



hive

Waking Up  
EVERY HOME  
in Britain by  
2020

Heating

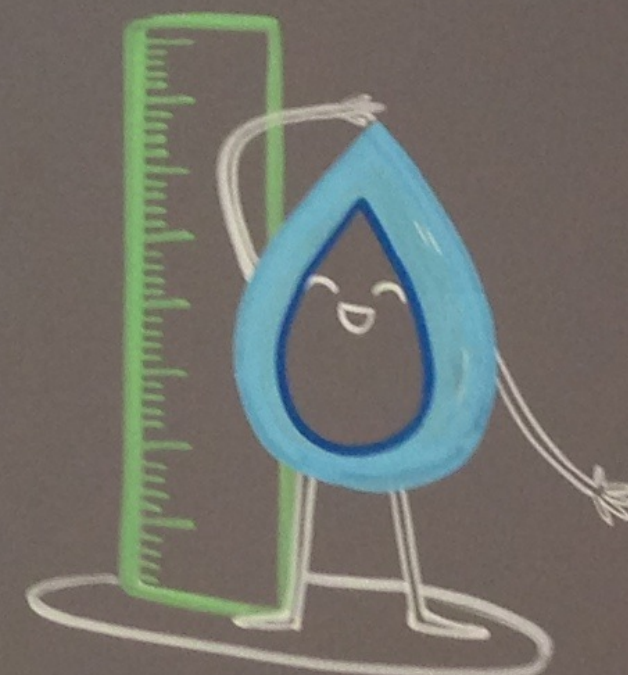


Wri  
RUN CLUB  
WHEN? WEDNESDAY  
12.30pm  
WHERE? REGENT'S PARK  
AND BACK (5K)  
ASK NICK IF YOU WANT TO



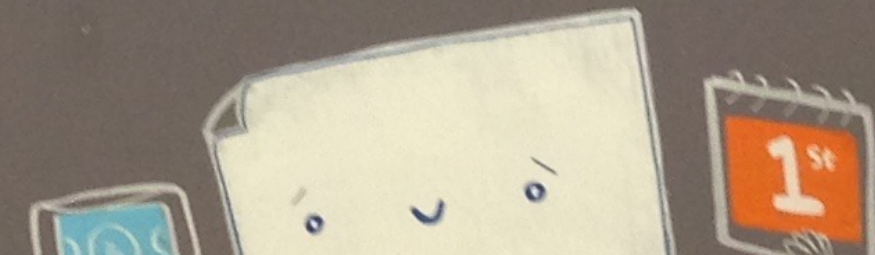
## JOINING THE DOTS

MY SMART METER  
LETS ME SEE  
HOW MUCH I'M  
€ Spending €  
ON MY ENERGY...



SUPPORTED BY MY  
SMART ENERGY REPORT  
WHICH HELPS ME UNDERSTAND  
WHAT AFFECTS my  
energy  
usage...

... I THEN GET A  
bill  
BASED ON ACCURATE  
&



Gresse  
Street

Rathbone Place

Lakeside West

British



The background is a large, grey exhibition wall covered in colorful, hand-drawn illustrations and text. On the left, there's a cloud with the 'hive' logo and the text 'Waking Up Every Home in Britain by 2020'. Below this, there are drawings of a laptop and a smartphone displaying energy-related data. In the center, a red SMEG refrigerator is visible. To its right, a whiteboard lists 'RUN CLUB' details. Further right, there are clouds labeled 'Rathbone Place', 'Lakeside West', and 'British'. A large, stylized blue house with a smiling face is in the center. At the bottom, there are drawings of a calendar showing '1st' and a document titled 'I THEN GET A bill BASED ON ACCURATE'.

**Hive : Control Your Heating from your Phone**

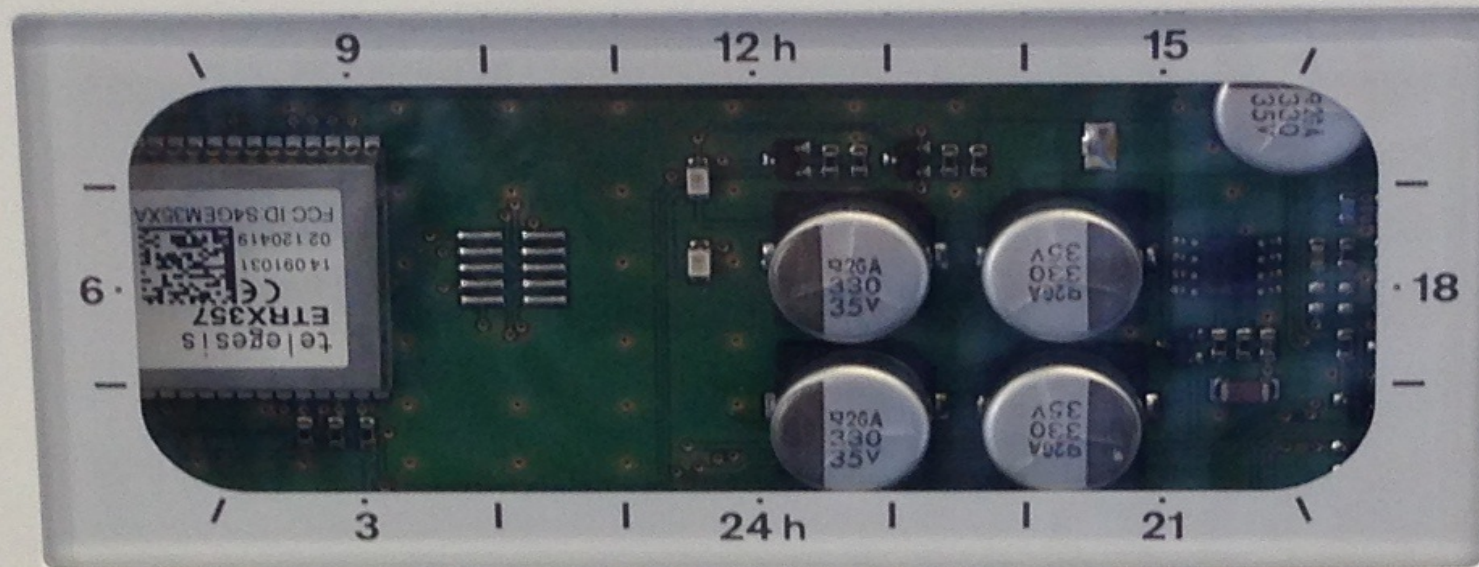
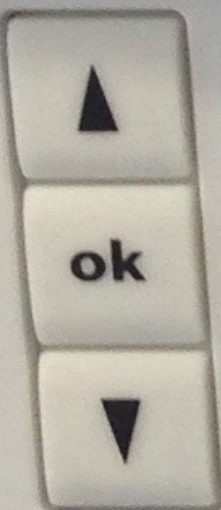
**Connected Boiler: Proactive Maintenance**

**MyEnergy: Understand your Energy Usage**





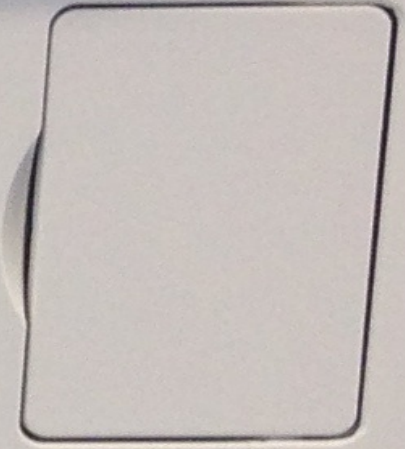
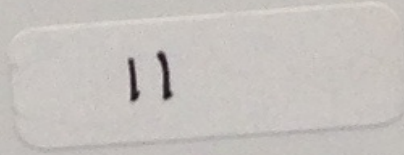




advance



advance





**600K - Smart Meter**

**3.8M Monthly**

**Future - 10 seconds**









my energy

[My usage](#)

[Energy saving advice](#)

my energy



January 2015 **£367.17**



Sun

Mon

Tue

Wed

Thu

Fri

Sat

 View graph

28

29

30

31

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

27

28

29

30

31

Highest day

**£15.15**

Sat 17th

Average day

**£11.84**

Lowest day

**£8.54**

Sat 10th

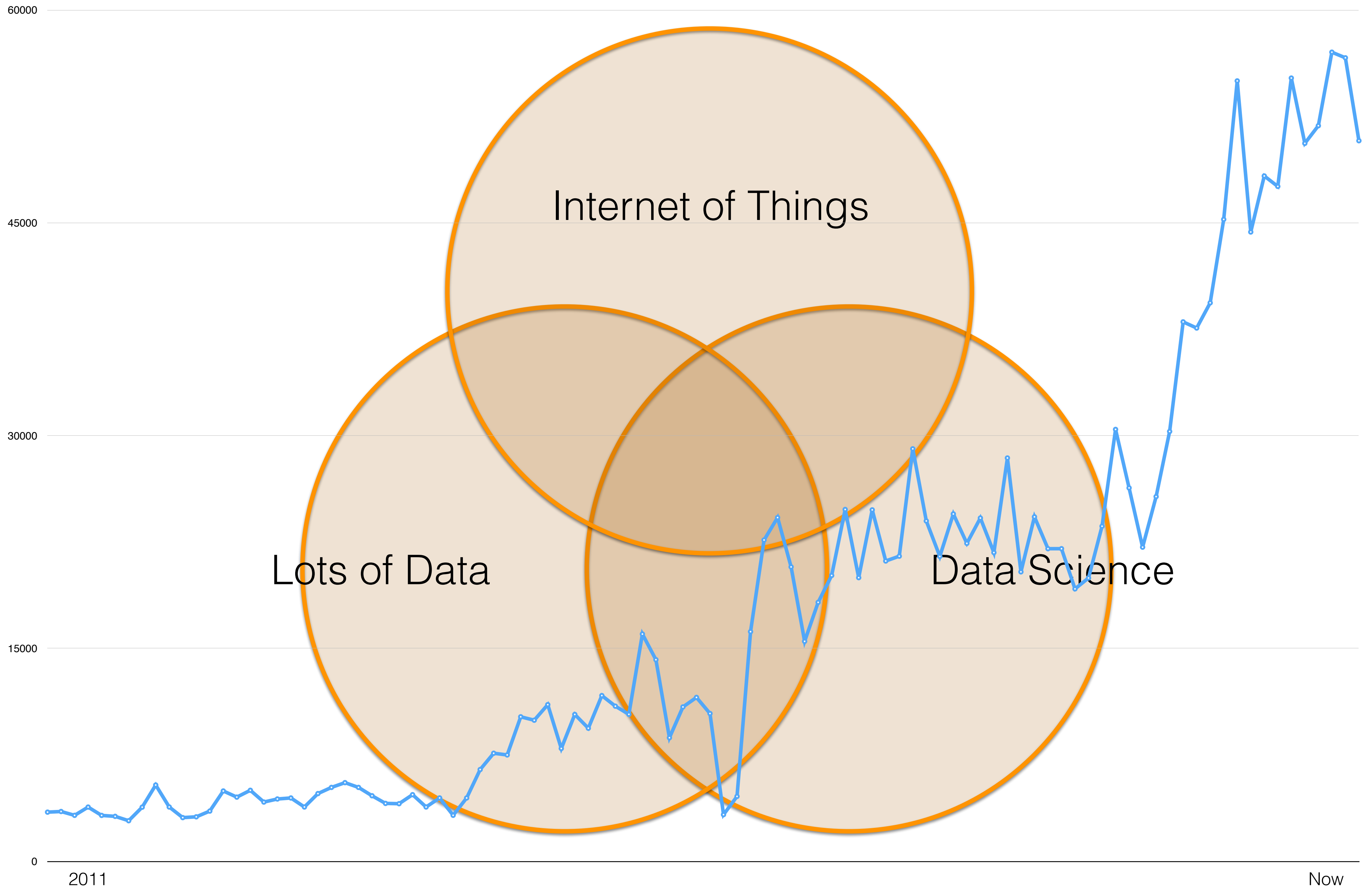
Send us a message

Give feedback or ask for help

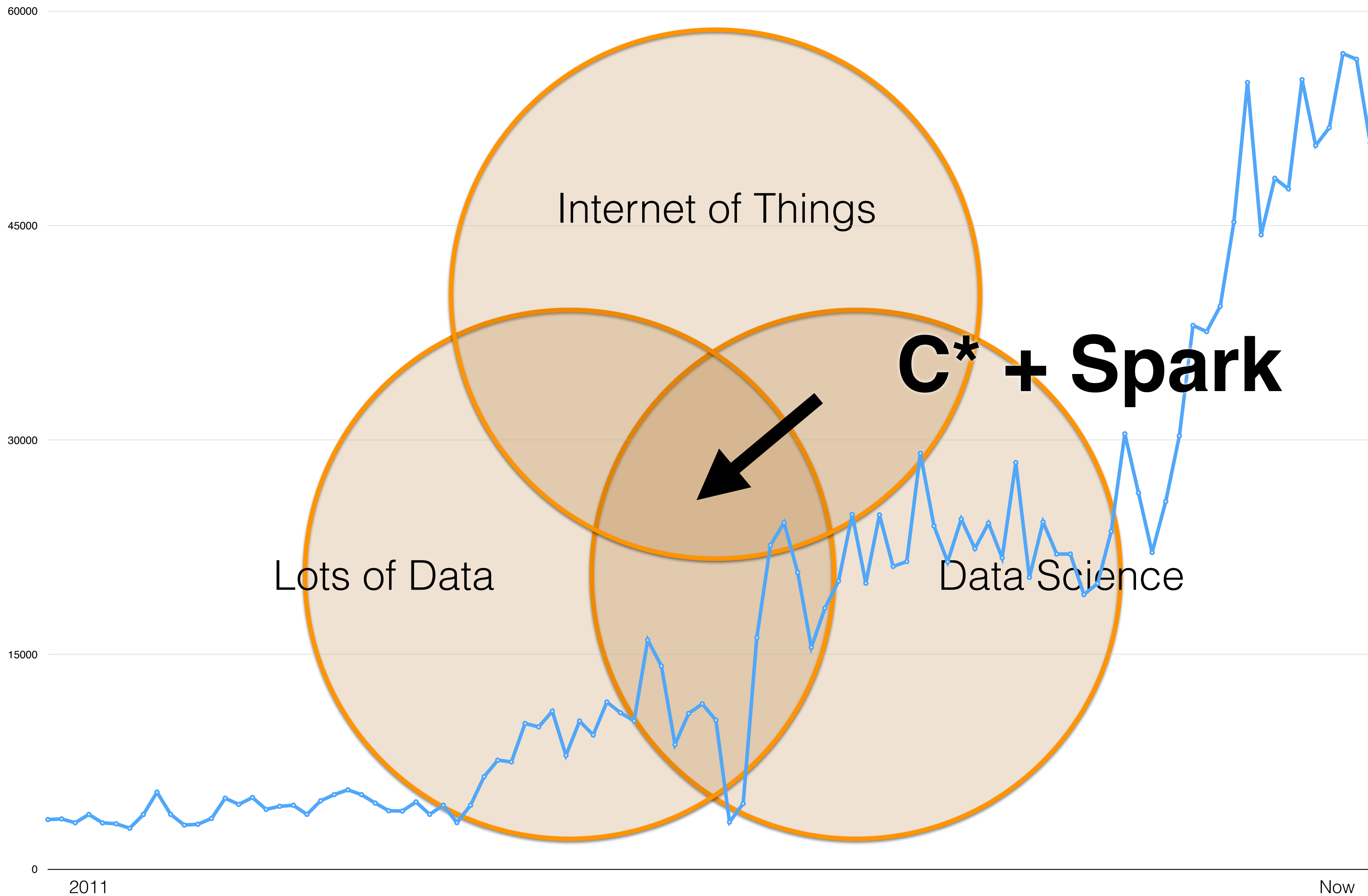
☐ Include a screenshot of this page

Next









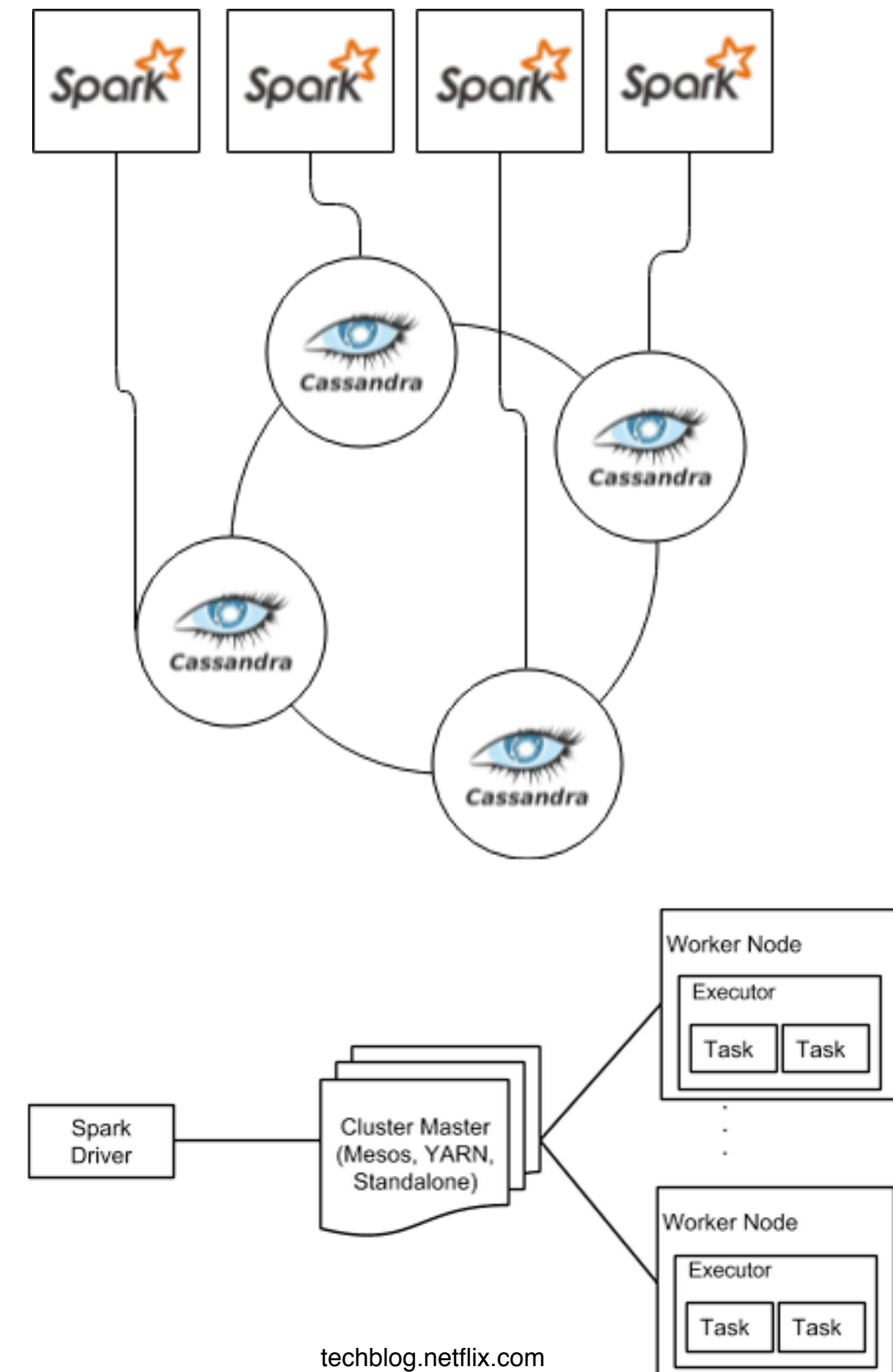


Lesson 1 : Not to race  
against bicycles



# Spark is for parallel execution

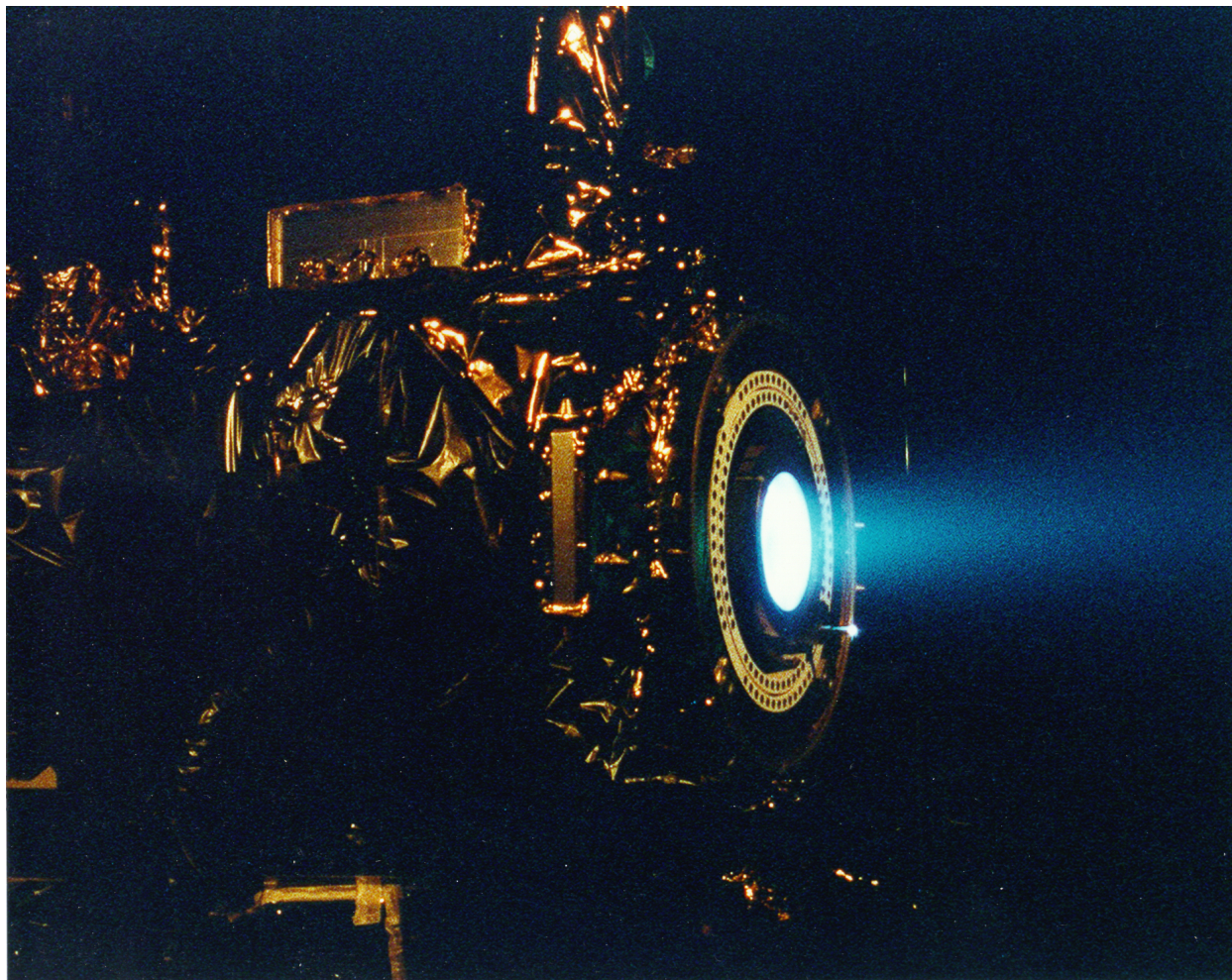
- Makes sense when we have jobs that can't run on a single machine
- The Spark master needs to distribute the job to workers
- If the job shuffles all data to one single node, parallelism is lost
- For small tasks, many times a simple script is better





# Things that look like a Spark / C\* cluster

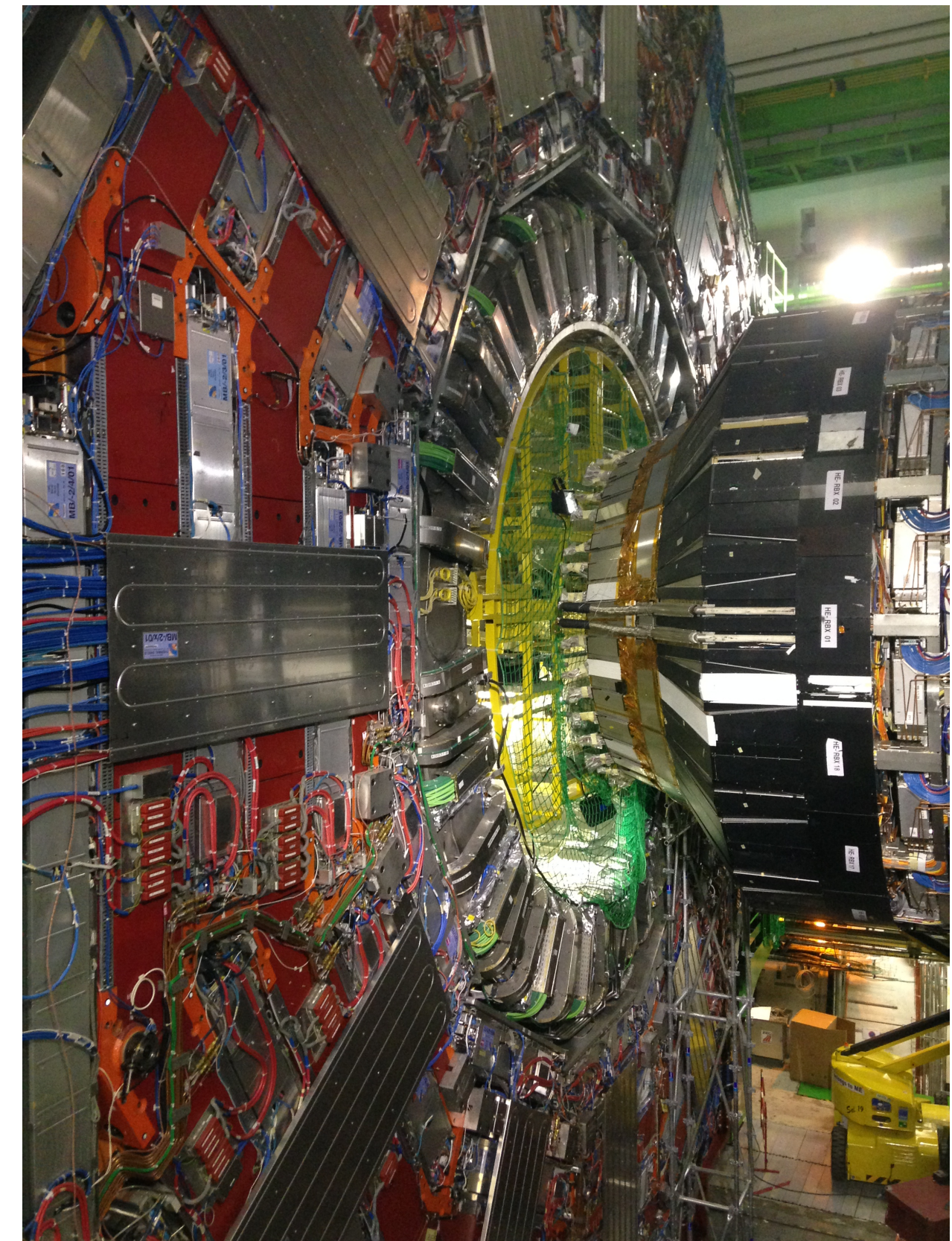
A Ion Thrust Engine



- It starts slow but in the long run goes very fast

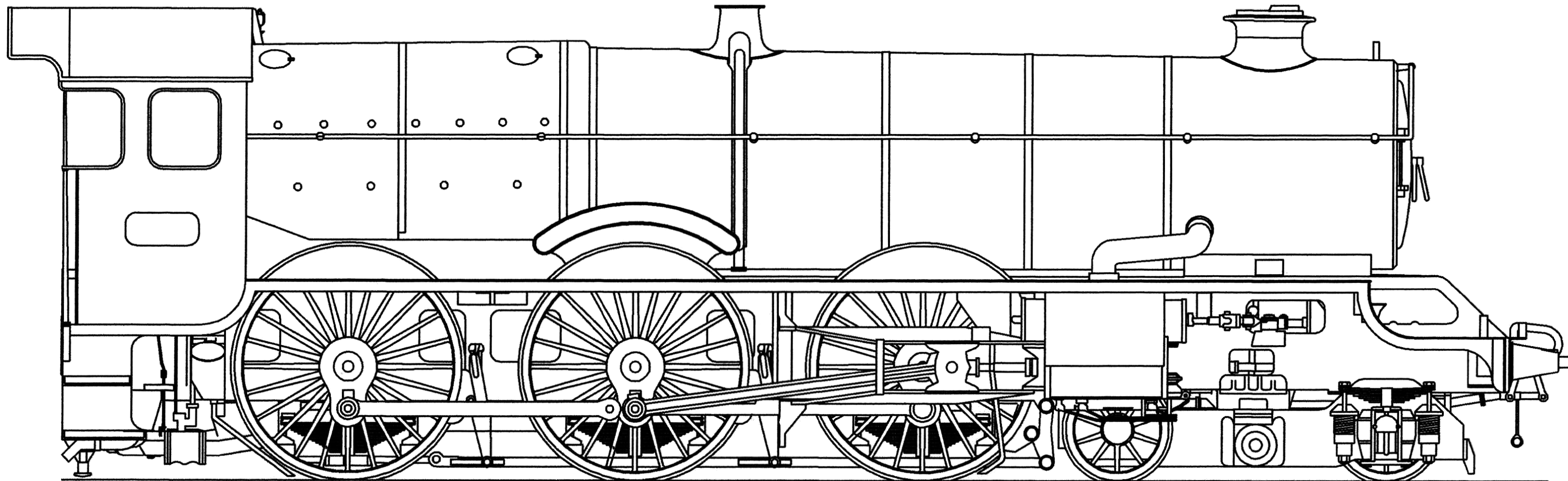
- It can achieve big energies
- It takes a lot of fine tuning

A Large Hadron Collider





# Who wins?



It depends on how far you go...

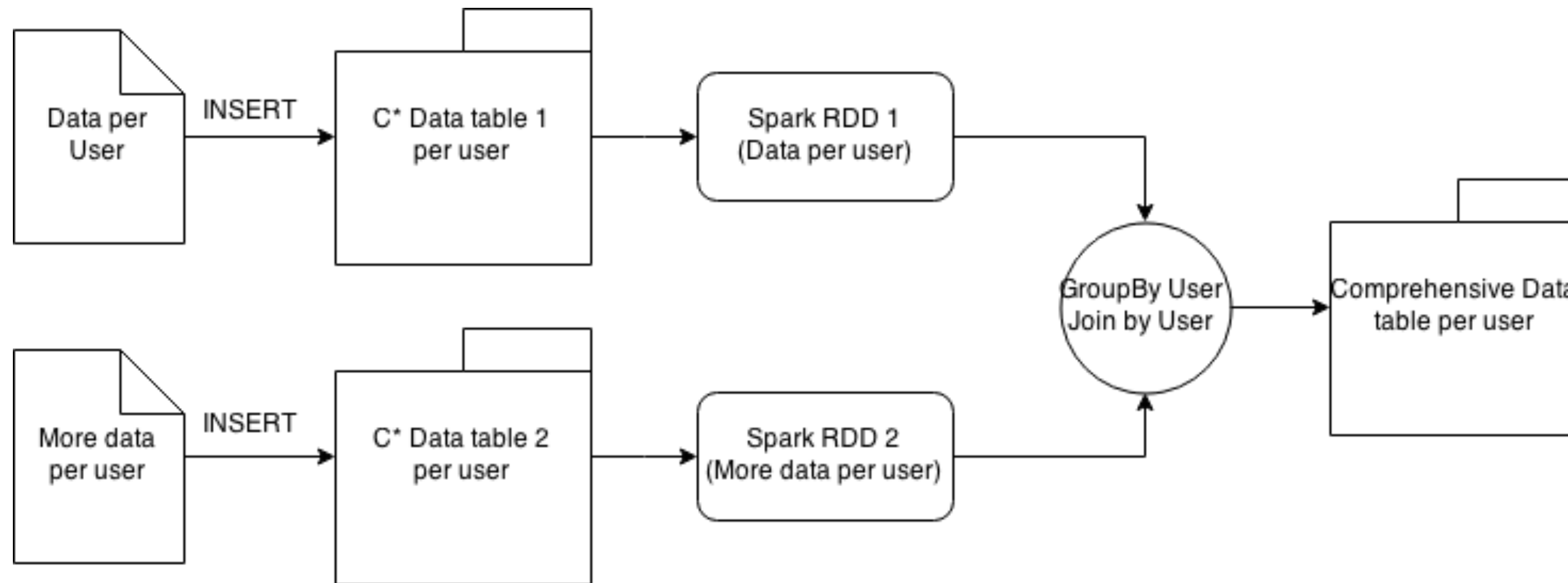


Lesson 2 : Not to use Spark  
too much



# Joining data from multiple sources

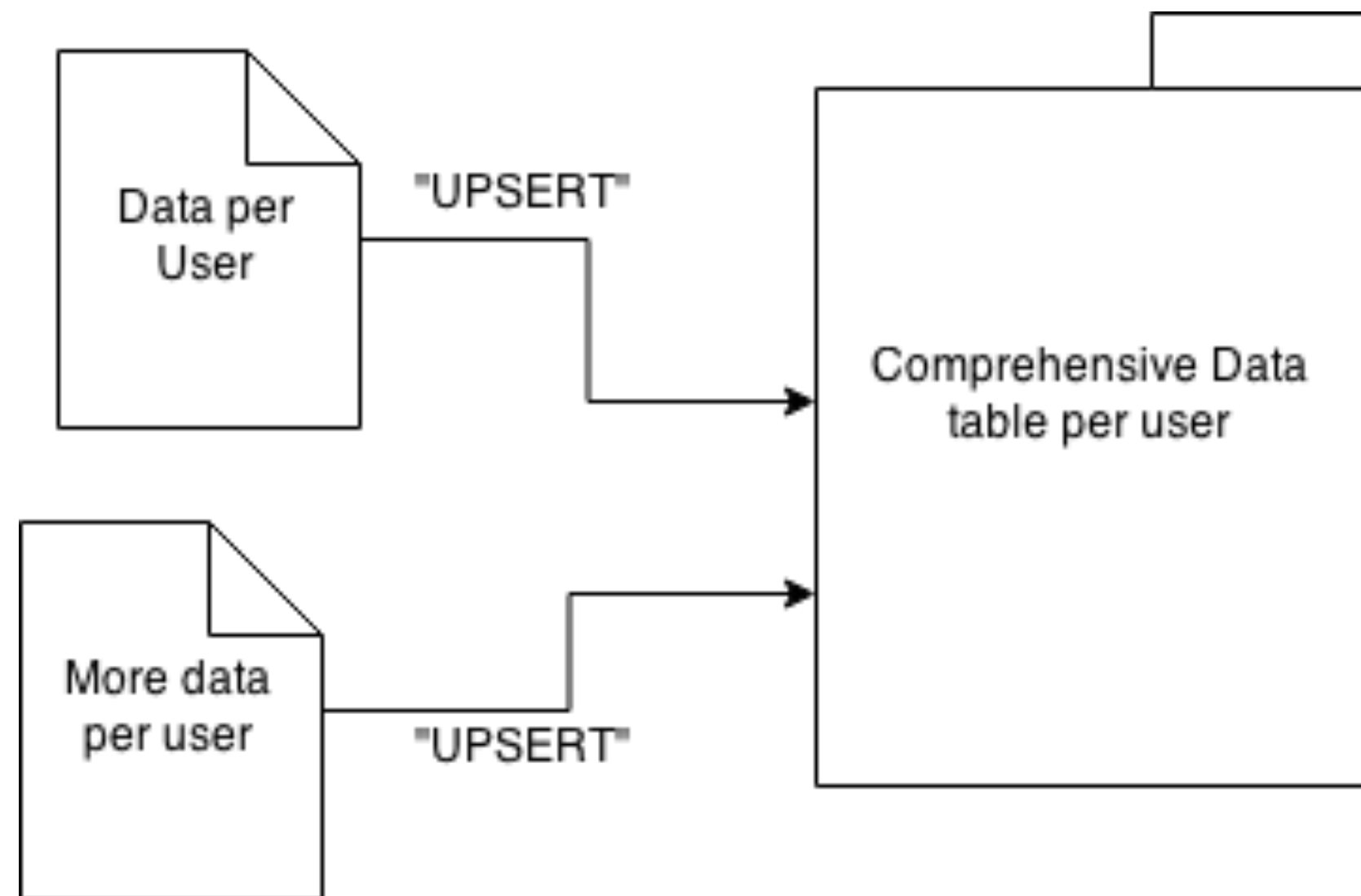
Think twice when you do that





# Upserting data from multiple sources

Do that if possible





# Upserting data from multiple sources

```
CREATE KEYSPACE meter_data
WITH REPLICATION = { 'class' : 'SimpleStrategy', 'replication_factor' : 1 };

CREATE TABLE meter_data.consumption (
  meter_id text,
  elec_read_kwh bigint,
  gas_read_m3 bigint,
  read_date timestamp,
  PRIMARY KEY ((meter_id),read_date));

INSERT INTO meter_data.consumption (meter_id, elec_read_kwh, read_date) VALUES ('10293847856', 2567, '2015-04-20');

INSERT INTO meter_data.consumption (meter_id, gas_read_m3, read_date) VALUES ('10293847856', 18363, '2015-04-20');

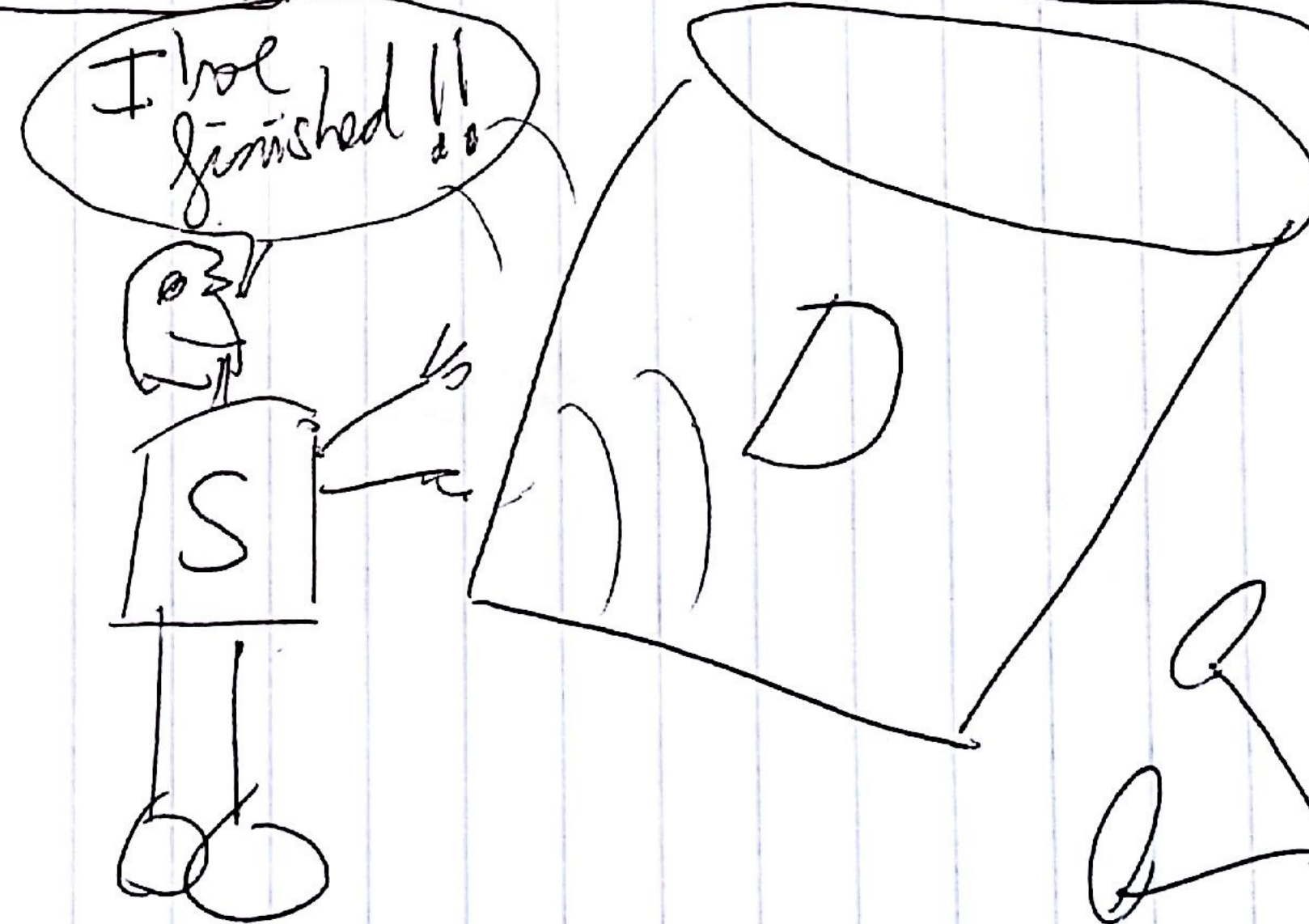
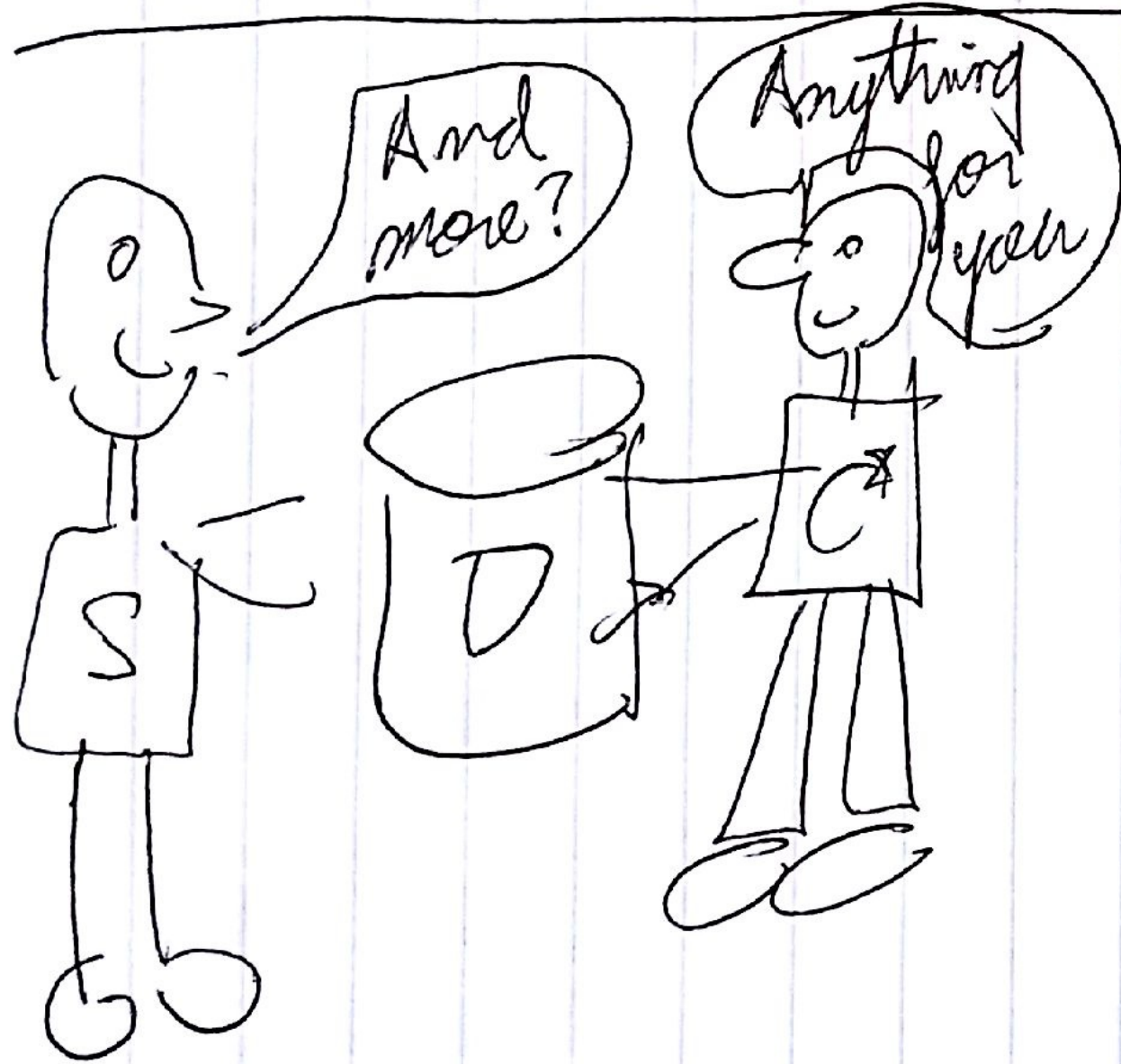
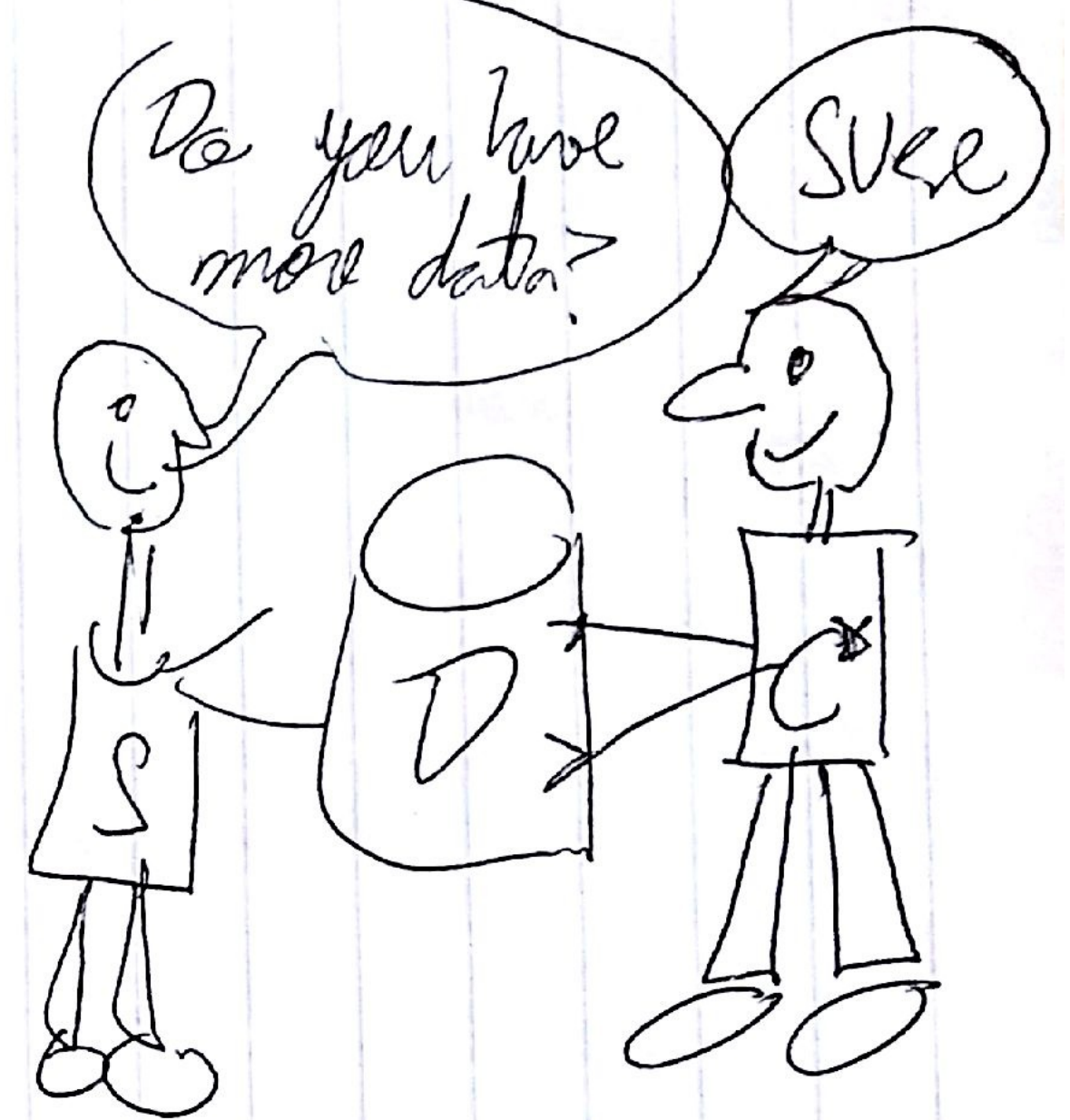
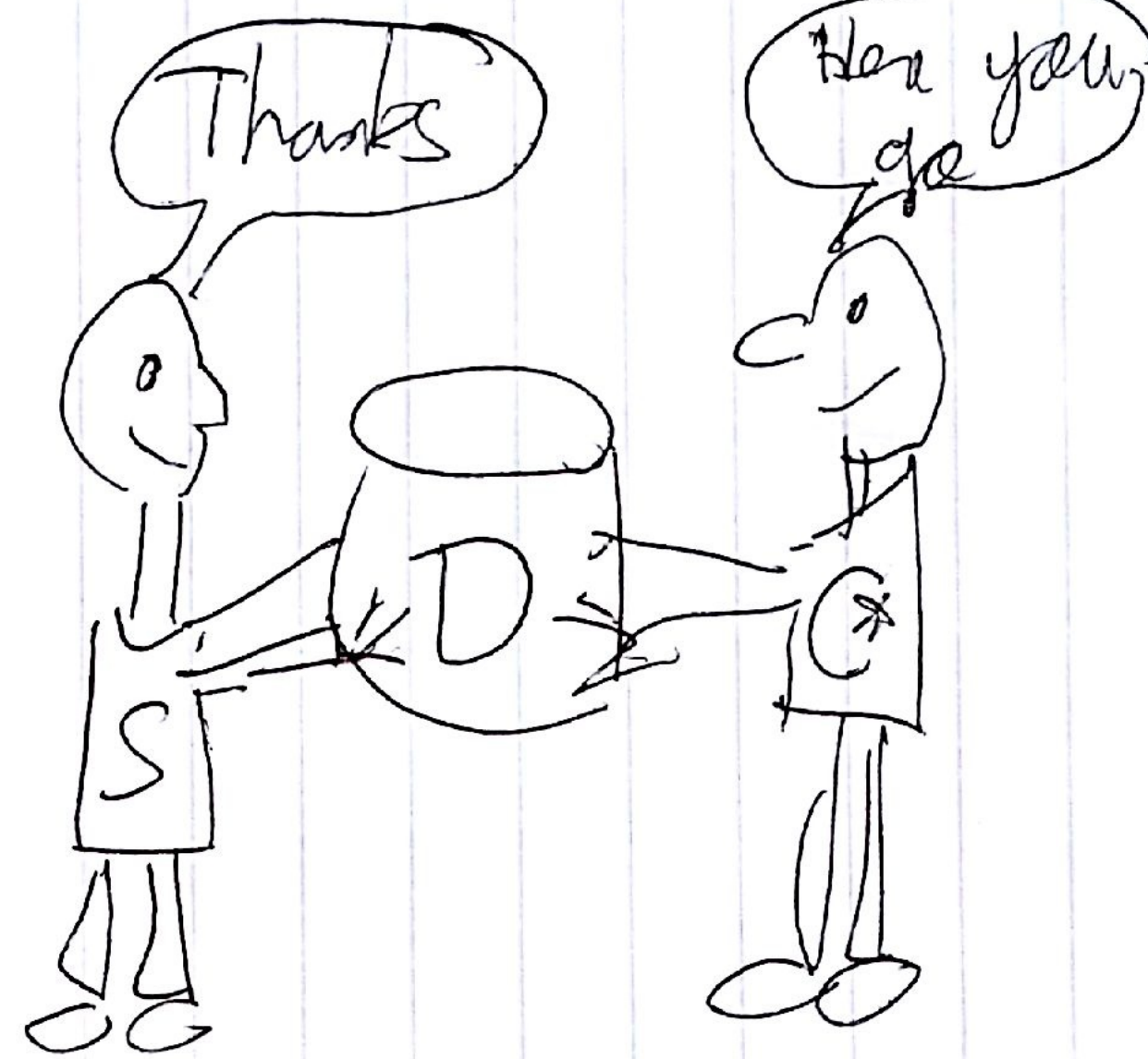
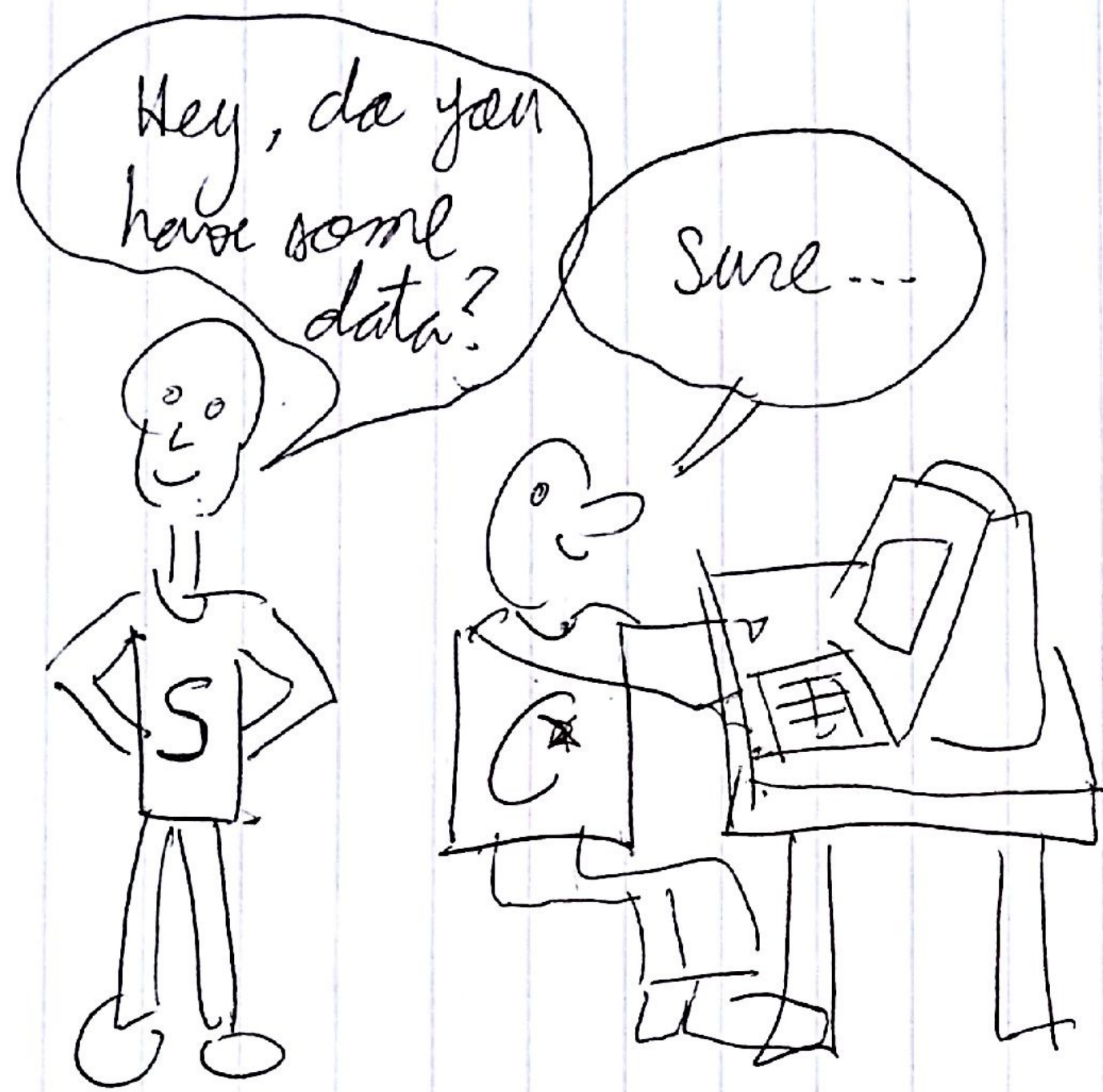
SELECT * FROM meter_data.consumption ;
```

meter_id	read_date	elec_read_kwh	gas_read_m3
10293847856	2015-04-20 00:00:00+0100	2567	18363



Lesson 3 : Spark is stronger  
than Cassandra







# Spark Properties & Cassandra-specific properties tuning

## Write properties

You can set the following properties in SparkConf to fine tune the saving process.

### **spark.cassandra.output.batch.size.bytes**

Default = auto. Number of bytes per single batch. The default, auto, means the connector adjusts the number of bytes based on the amount of data.

### **spark.cassandra.output.consistency.level**

Default = LOCAL\_ONE. Consistency level to use when writing.

### **spark.cassandra.output.concurrent.writes**

Default = 5. Maximum number of batches executed in parallel by a single Spark task.

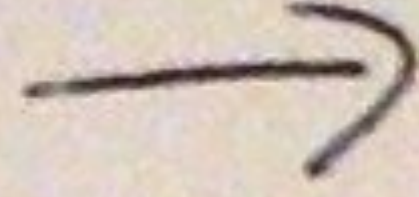
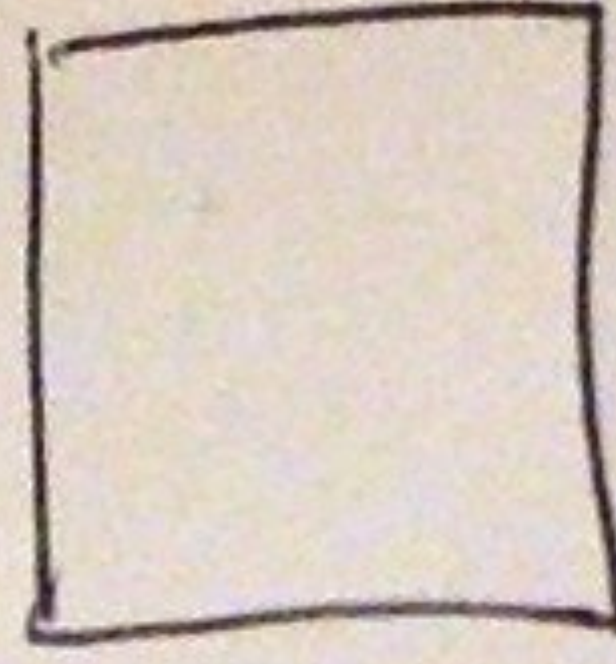
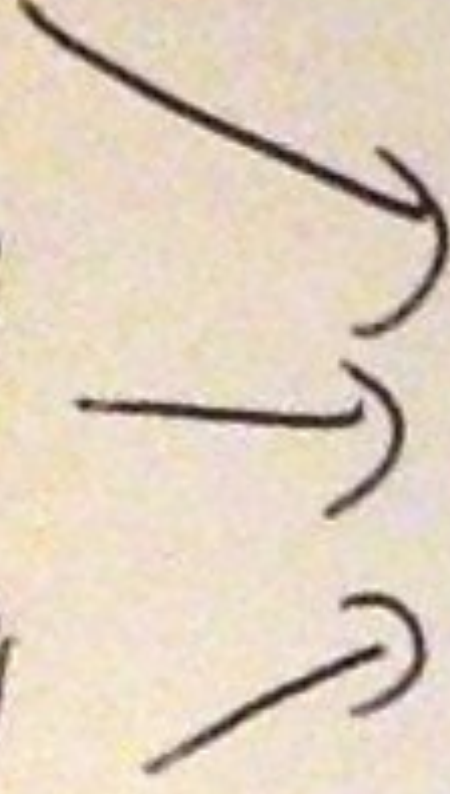
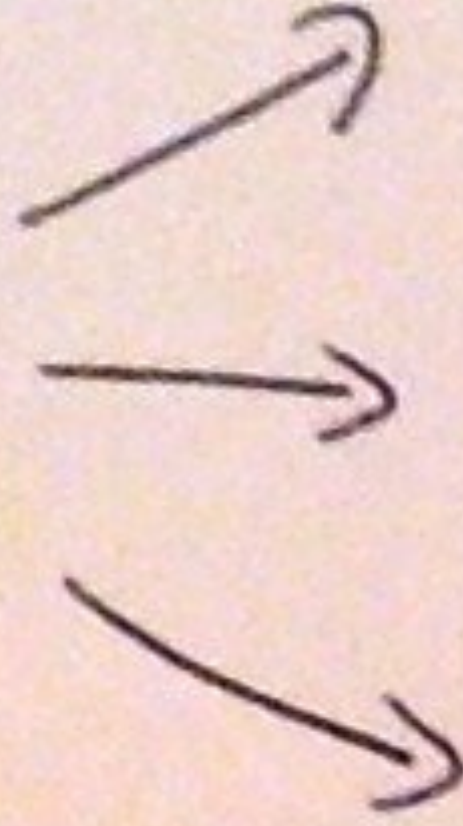
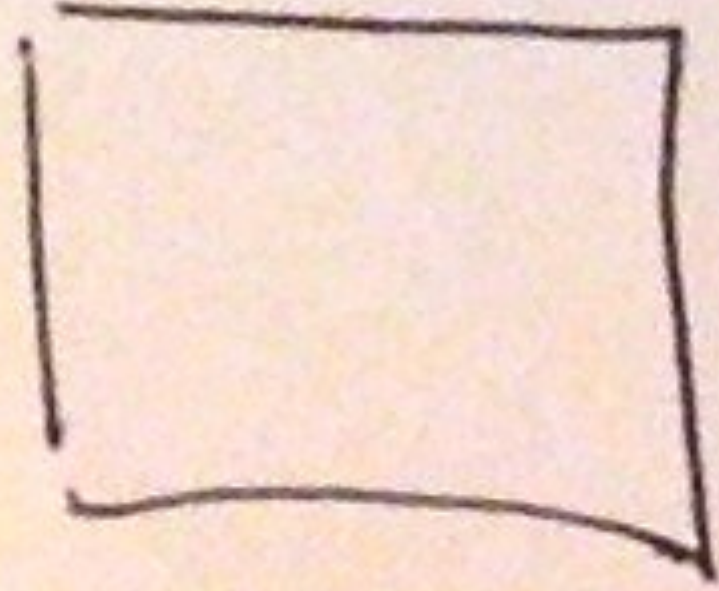
### **spark.cassandra.output.batch.size.rows**

Default = 64K. The maximum total size of the batch in bytes.



# Lesson 4 : Mindset

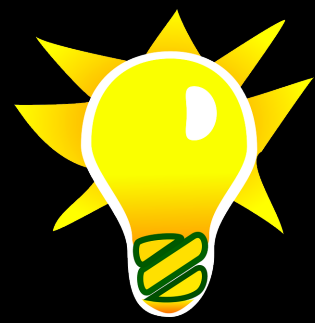




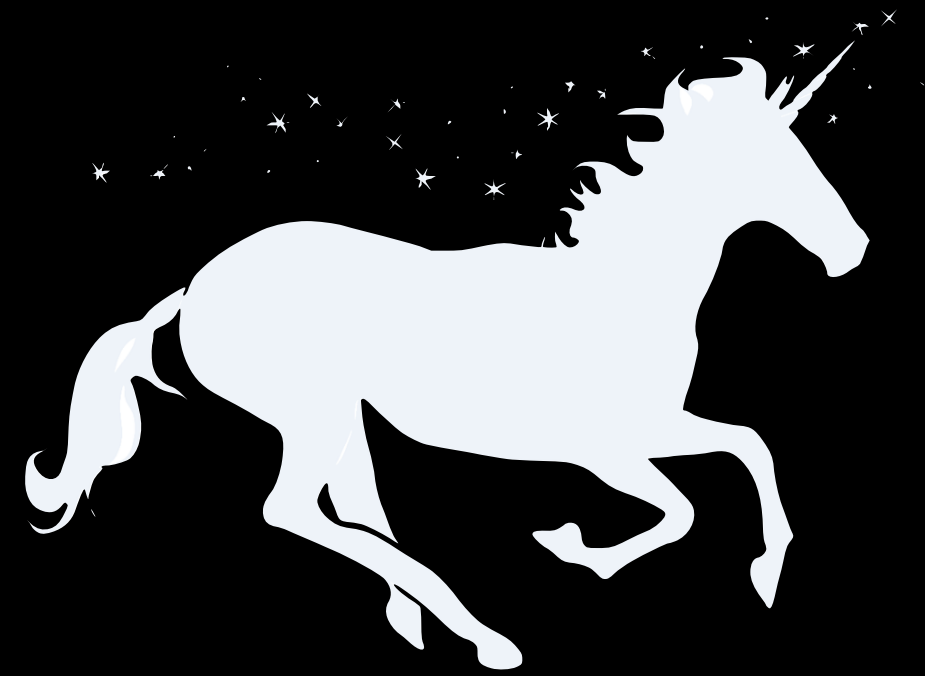


# Lesson 5 : Velocity

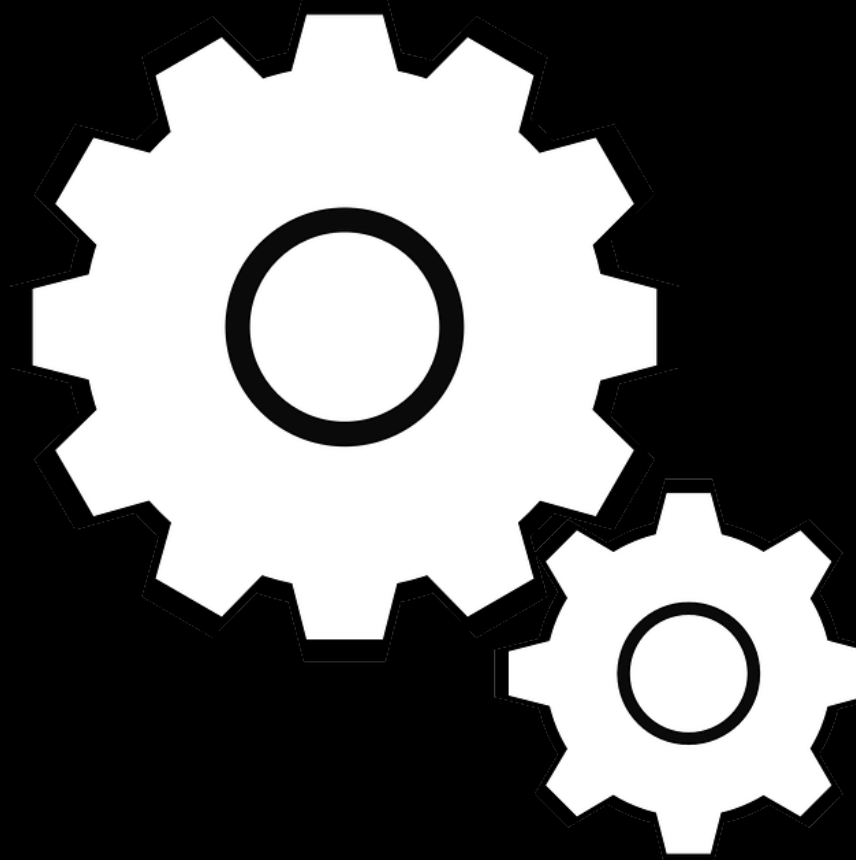




Idea



Data Science

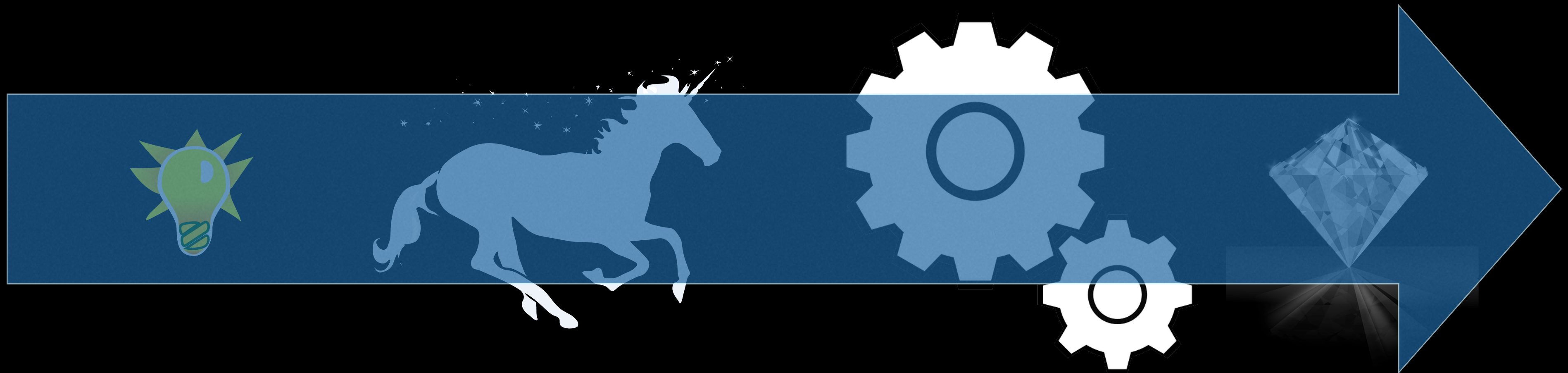


Data Engineering  
Data Operations



Value





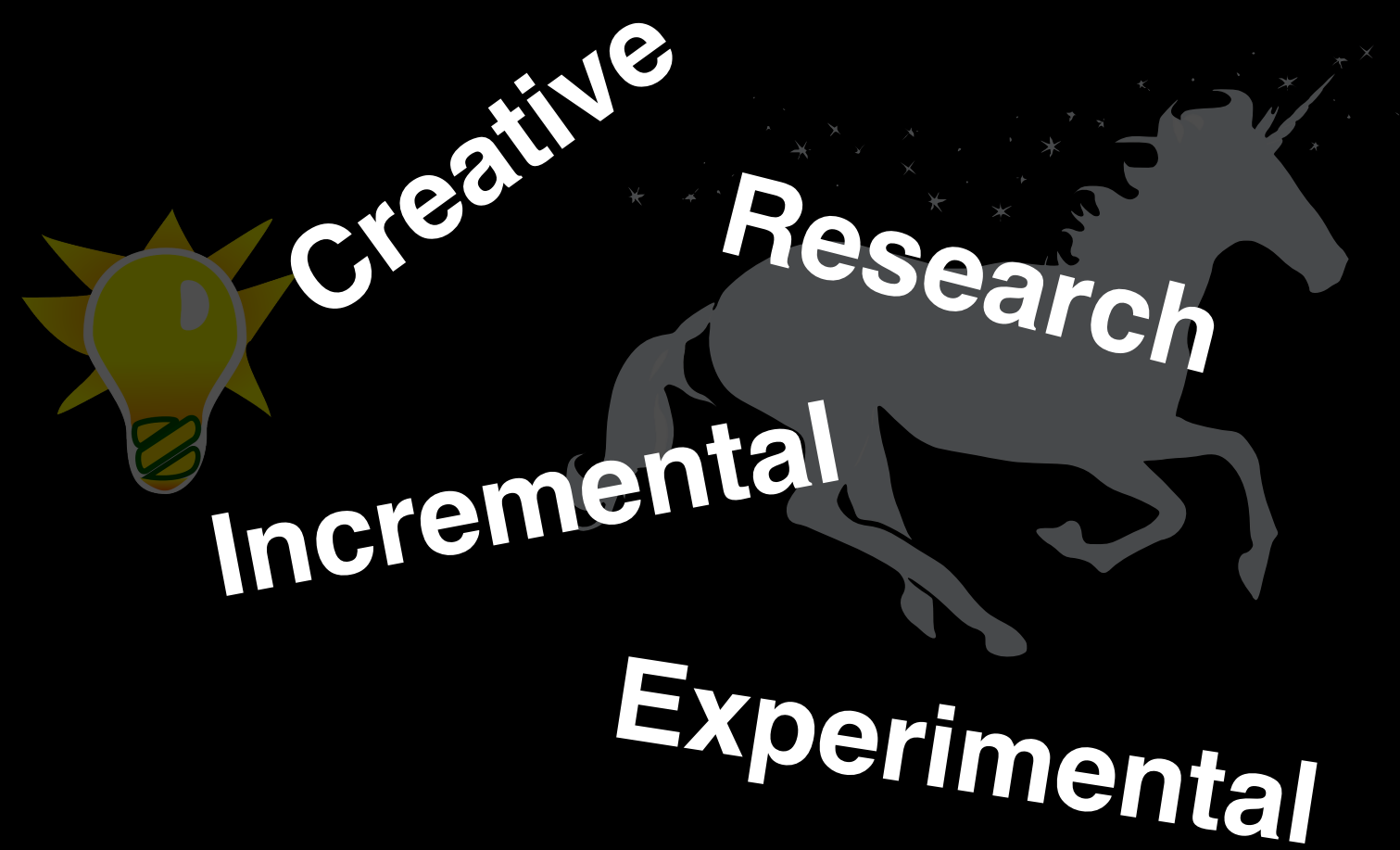
Idea

Data Science

Data Engineering  
Data Operations

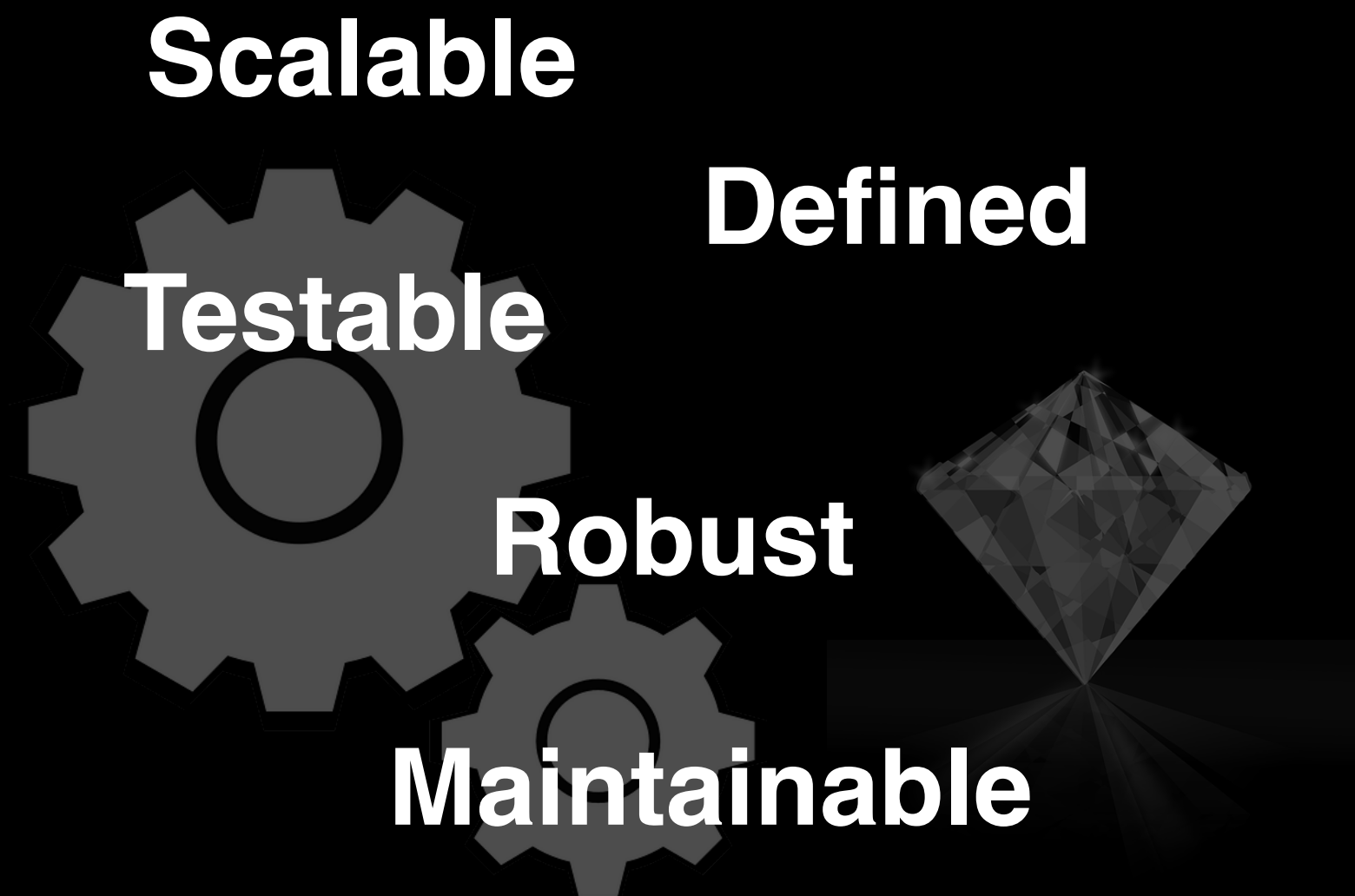
Value





Idea

Data Science



Data Engineering

Data Operations

Value



R

Python

Java

Scala

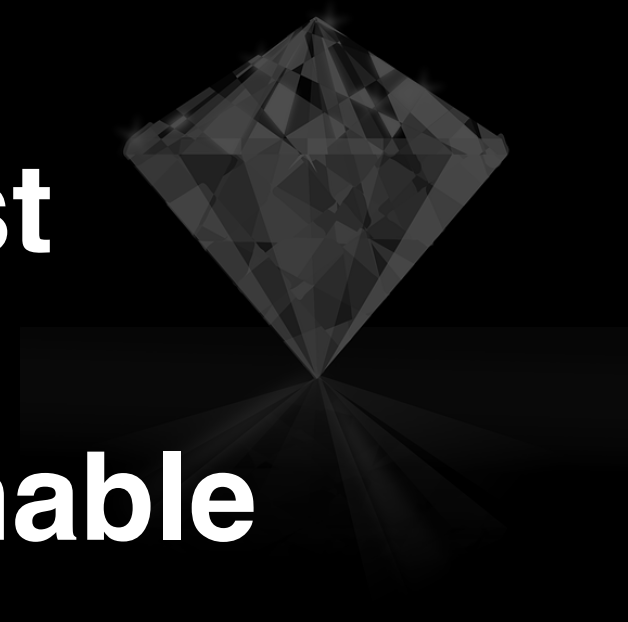
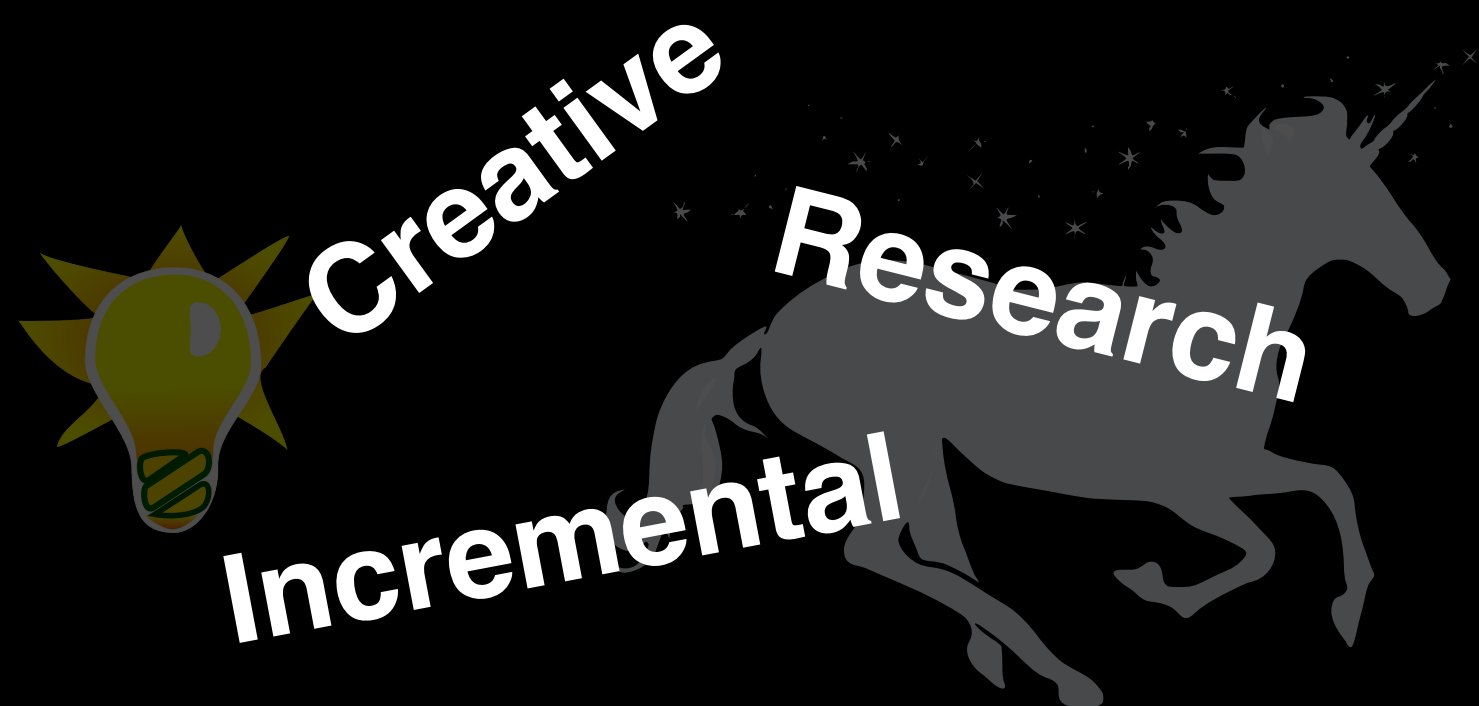
Scalable

Defined

Testable

Robust

Maintainable



Small Datasets

Idea

Data Science

Data Engineering

Value

Offline

#BigData

Data Operations

Clustered

Single Machine

Realtime



R

Python



Creative

Research

Incremental

Experimental

Data Science

Offline

Small Datasets

Idea

Single Machine

Java

Scala

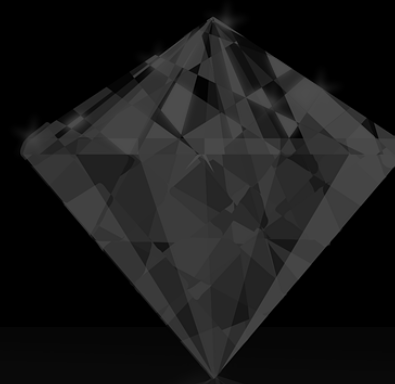
Scalable

Defined

Testable

Robust

Maintainable



Data Engineering

Value

Data Operations

#BigData

Clustered

Realtime



R

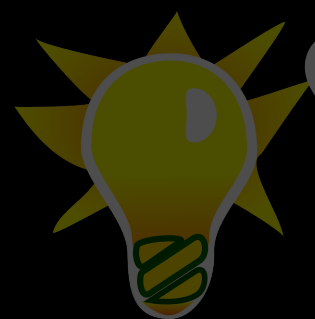
Python

Java

Scala

Scalable

Defined



Creative

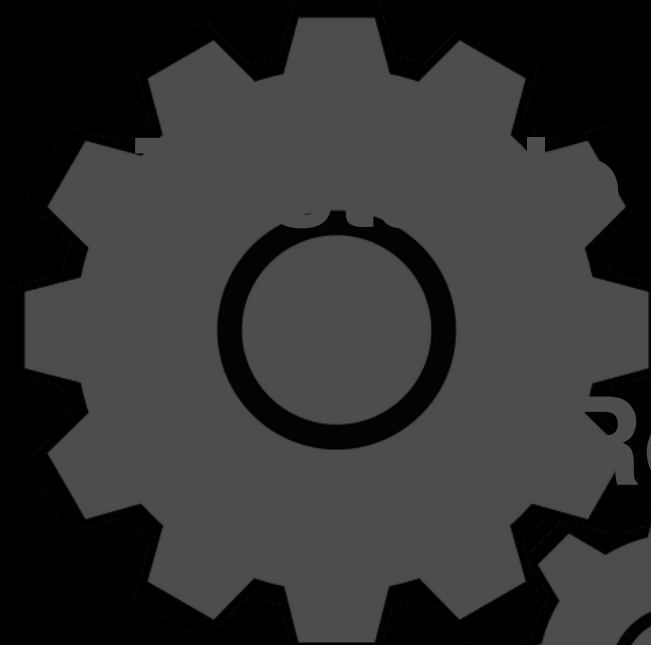
Research

Incremental

Experimental

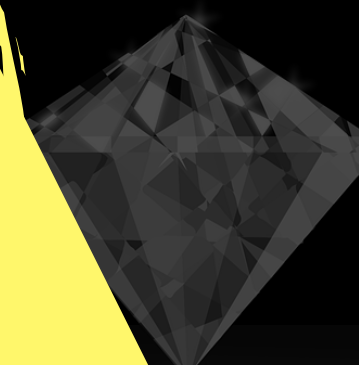
Data Science

Offline



Robust

Modifiable



Value

Data Engineering

Data Operations

Clustered

Small Datasets

Idea

Single Machine

Realtime



+



# Thankyou

@JimAnning : jim.anning@bgch.co.uk

@Jcasals :  josep.casals@bgch.co.uk