# PRACTICAL IOT DEVELOPMENT USING ORACLE BIG DATA AND ORACLE DV ...AND A WIFI KETTLE
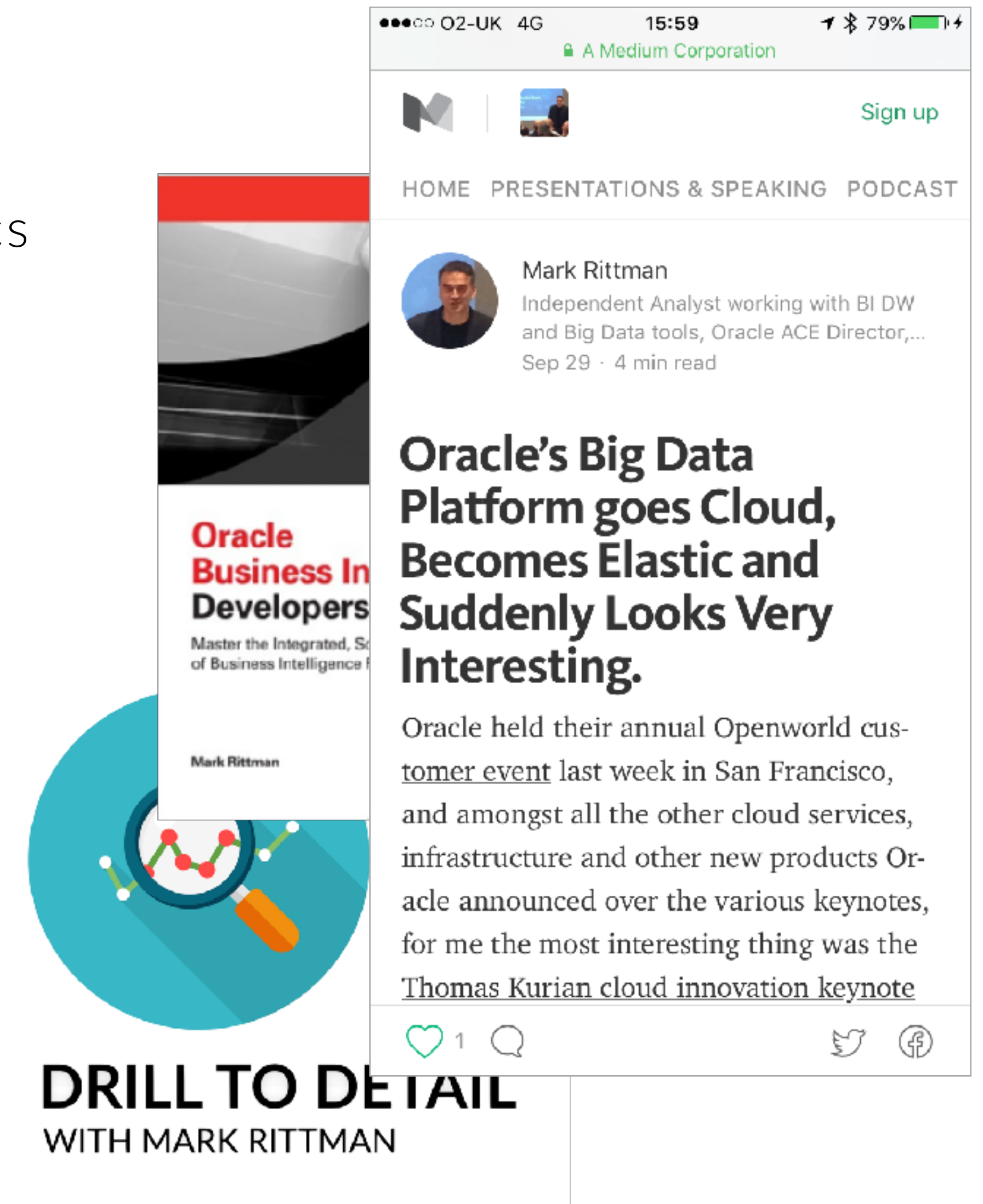
Mark Rittman, Oracle ACE Director & Independent Analyst
MJR Analytics ltd (http://www.mjr-analytics.com)

BIWA SUMMIT 2017, SAN FRANCISCO
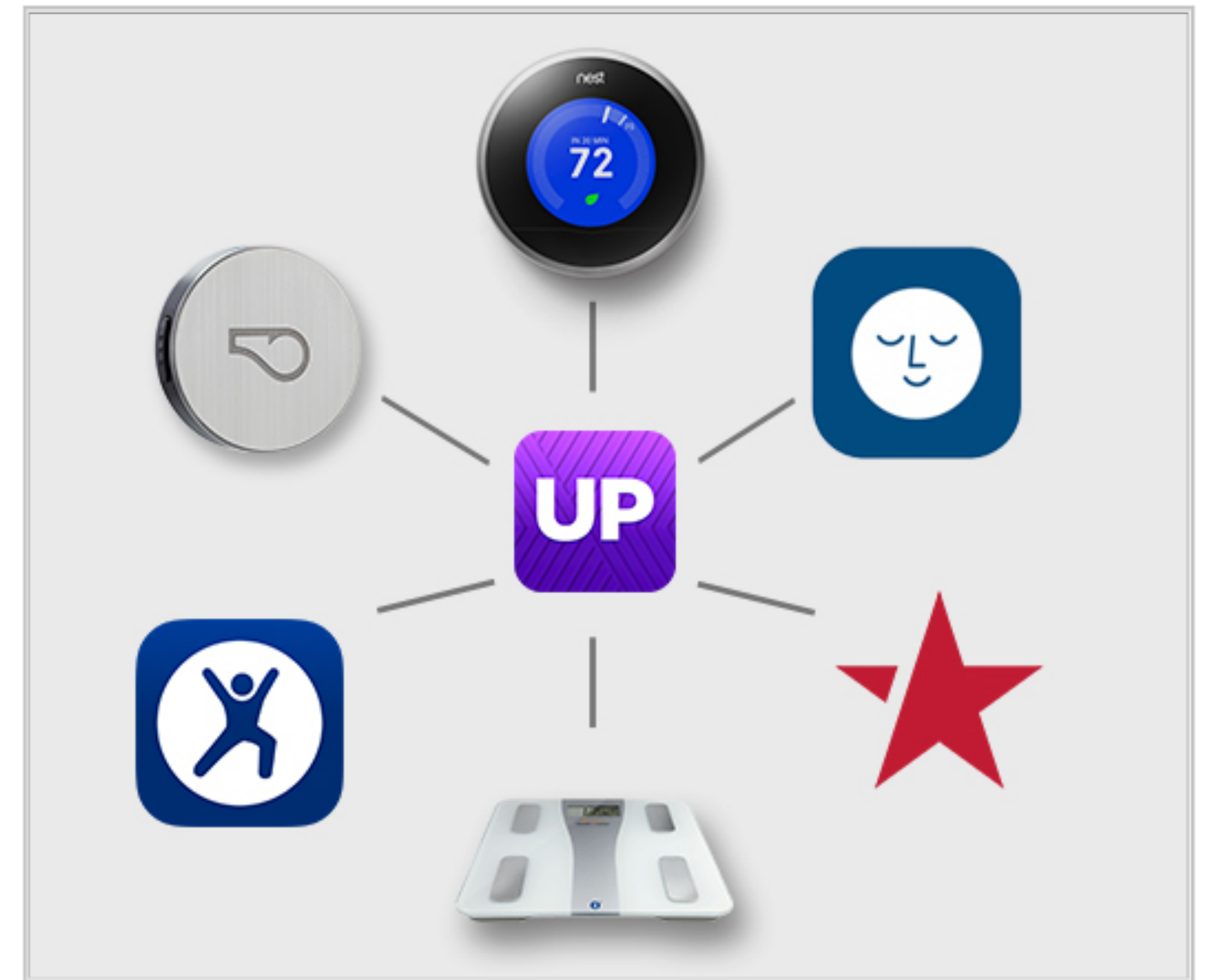
T : @markrittman

# About Mark Rittman

- Oracle ACE Director, Independent Analyst
- Company founder, Oracle ACE Director, product specialist
- Now working in product management around big data & analytics
- Regular columnist for Oracle Magazine, OTN
- Author of two books on Oracle BI & Engineered Systems
- 15+ Years in Oracle BI, DW, ETL + now Big Data
- Personal blog at medium.com/mark-rittman
- Podcast on iTunes and drilltodetail.com
- Contact me at mark@rittman.co.uk

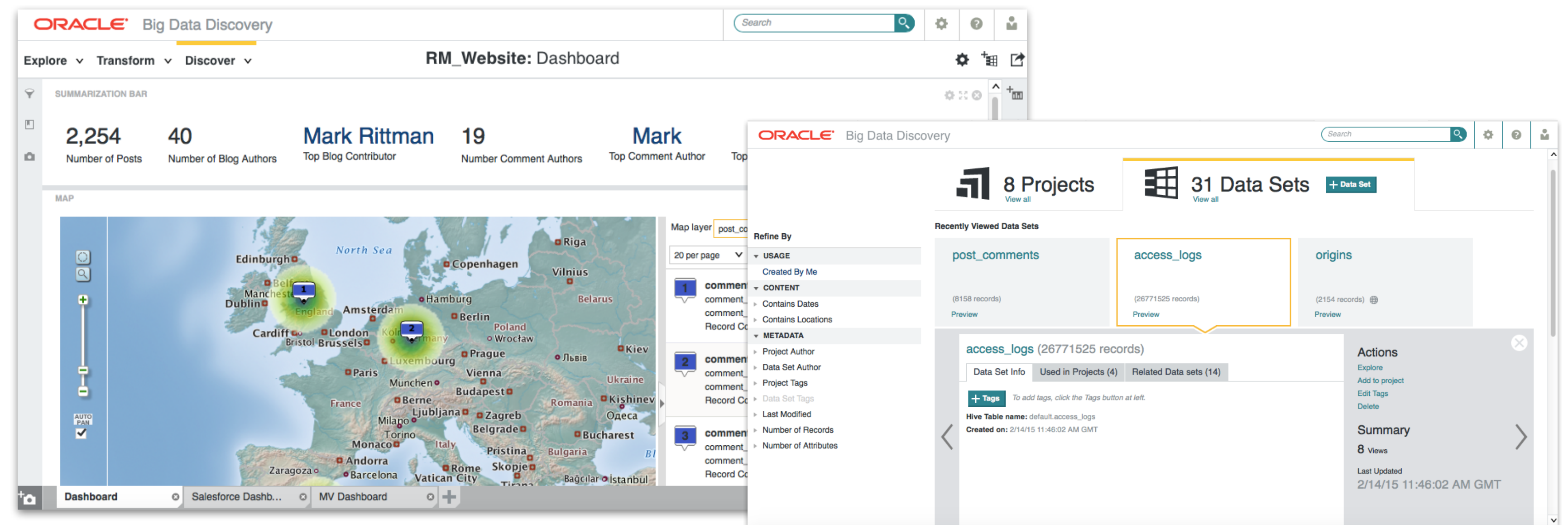# Wearables & Smart Devices - Our Data Ecosystem

- How many of you are using health bands, smartphone apps, other life logging services?
- It's likely fair proportion of you log workouts, steps and other activities daily
- Some of you may have Nest, Hue or other home smart devices
- All of these services capture and generate useful data
- What if we could capture, combine and mine this data for insights, correlations, trends and patterns?
  - And what if we used Oracle Big Data Discovery to bring the data together, and mine for those insights?

# FOR THE PAST SIX MONTHS, I DID JUST THAT

# Oracle Big Data Discovery - What Is It?

- A visual front-end to the Hadoop data reservoir, providing end-user access to datasets
- Data sampled and loaded from Hadoop (Hive) into NoSQL Dgraph engine for fast analysis
- Catalog, profile, analyse and combine schema-on-read datasets across the Hadoop cluster
- Visualize and search datasets to gain insights, potentially load in summary form into DW

# Key Features in Oracle Big Data Discovery 1.1.x
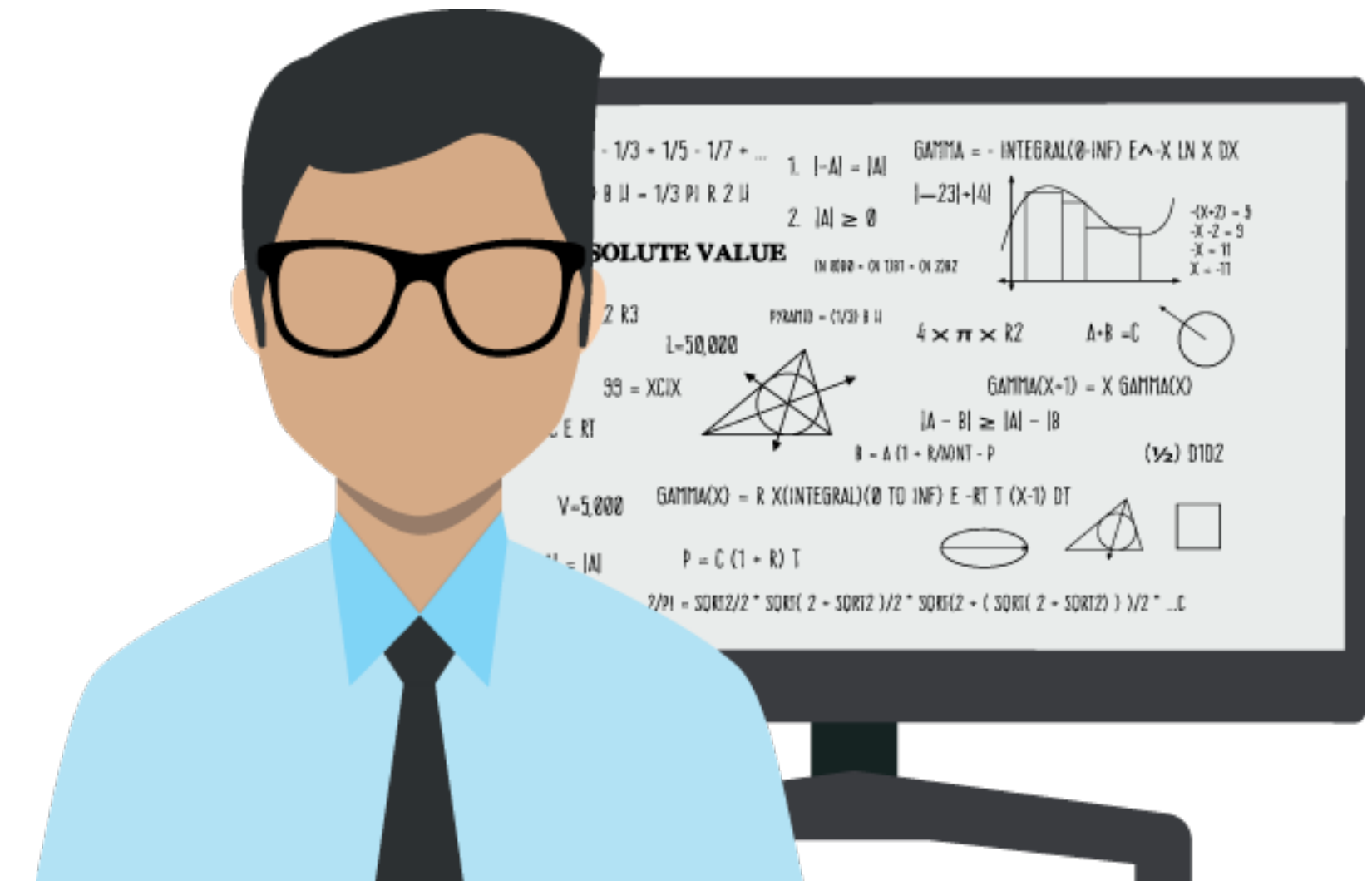
VISUALISING AND
TRANSFORMING DATA

- Provide a visual catalog and search function across data in the data reservoir
- Profile and understand data, relationships, data quality issues
- Apply simple changes, transformations to data
- Add enrichment to incoming data including sentiment, geo-location

COMMUNICATING
AND BUNDLING

- Visualize datasets using rich chart types
- Join datasets at visualisation level
- Add data from JDBC + file sources
- Prepare more structured Hadoop datasets for use with other tools

# New Features In Oracle Big Data Discovery 1.2

### IMPORTING AND TIDYING DATA

- **Aggregation**
- **Materialised Joins**
- Better Pan and Zoom
- Speed and Scale

### METADATA AND DEVELOPER PRODUCTIVITY

- Metadata Curation
- Attribute-level Search from Catalog
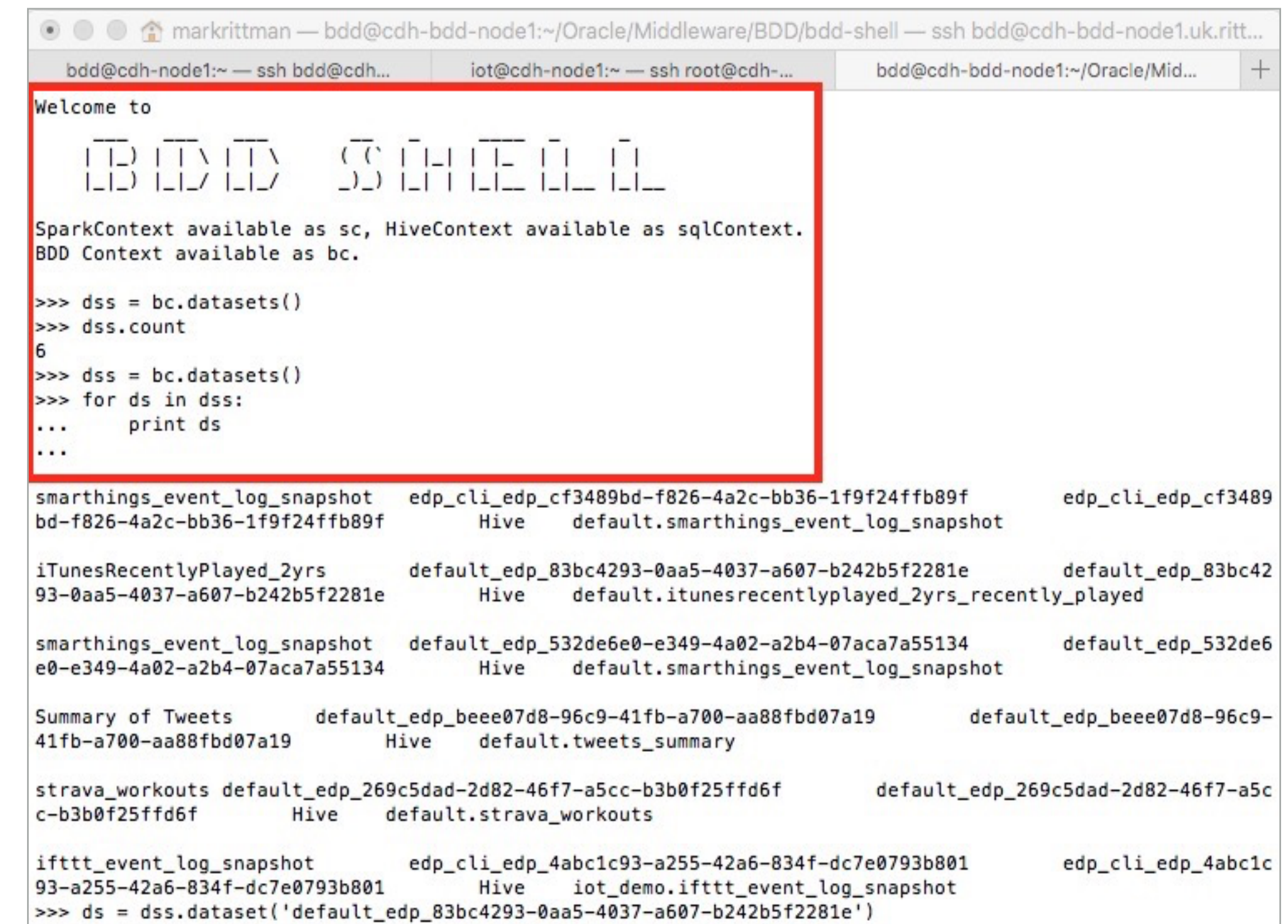- Activity Hub
- **Python Interface to BDD Datasets**

### COMMUNICATING AND BUNDLING

- Streamlined UI
- Faster Data Indexing
- Activity Hub
- Sunburst Visualization

# BDD Shell - pySpark Command-Line

- Interactive tool designed to work with BDD without using Studio's front-end
- Exposes all BDD concepts
  (views, datasets, data sources etc)
- Supports Apache Spark
- HiveContext and SQLContext exposed
- BDD Shell SDK for easy access to BDD features, functionality
- Access to third-party libraries such as Pandas, Spark ML, numPy
- Use with web-based notebook such as iPython, Jupyter, Zeppelin

# Oracle Big Data Discovery as the Data Scientists' Toolkit

VISUALISING AND
TRANSFORMING DATA

MODELING AND INFERRING

COMMUNICATING
AND BUNDLING

IMPORTING AND
TIDYING DATA

# Using Wearables To Enhance & Improve Workouts

- Over the past year or so, I've getting into cycling and generally trying to keep fit and lose weight
- Also using these activities as data sources for this project
  - Using Wahoo Elemnt + Strava for workout recording
  - Withings Wifi scales for weight + body fat measurement
  - Jawbone UP3 for steps, sleep, resting heart rate
  - All the time, collecting data and storing it in Hadoop

# HOME AUTOMATION

**Shades**
Open the shades and let the morning light in.

**Electric Kettle**
Start the electric kettle so the water is ready for your tea.

**Thermostat**
Turn up the thermostat before you get out of bed.

11

# Home Automation and Smart 'IoT' Devices

- Another personal project has been home Automation, IoT and the "smart home"

- Started with Nest thermostat and Philips Hue lights

- Extended the Nest system to include Nest Protect and Nest Cam

- Used Apple HomeKit, HomeBridge, Apple TV for Siri voice control

- Added Samsung Smart Things hub for Z-wave, Zigbee compatibility

"Hey Siri turn on the Office lights"
tap to edit

OK, the lights are turned on.

Nest Protect (X2), Thermostat, Cam

Apple Homekit, Apple TV, Siri

Homebridge Homekit / Smarthings Connector

Philips Hue Lighting

Samsung Smart Things Hub (Z-Wave, Zigbee)

Withings Smart Scales

Airplay Speakers

Door, Motion, Moisture, Presence Sensors

sign in    become a supporter    subscribe    search    jobs    dating    more ▾    UK edition ▾

**theguardian**
website of the year

UK   world   politics   sport   football   opinion   culture   business   lifestyle   fashion   environment   tech   travel    ☰ browse all sections

home › tech

**Smart homes**

# English man spends 11 hours trying to make cup of tea with Wi-Fi kettle

Data specialist Mark Rittman spent an entire day attempting to set up his new appliance so that it would boil on command

Bonnie Malkin

🐦 @bonniemalkin

Wednesday 12 October 2016 02.29 BST

Shares
**2,747**

Comments
**591**

🔖 Save for later



Mark Rittman set about trying to make a cup of tea at 9am but night had fallen by the time his new Wi-Fi

# And The Third Hobby : Land All That Data Into Hadoop

- Data extracted or transported to target platform using LogStash, CSV file batch loads
- Landed into HDFS as JSON documents, then exposed as Hive tables using Storage Handler
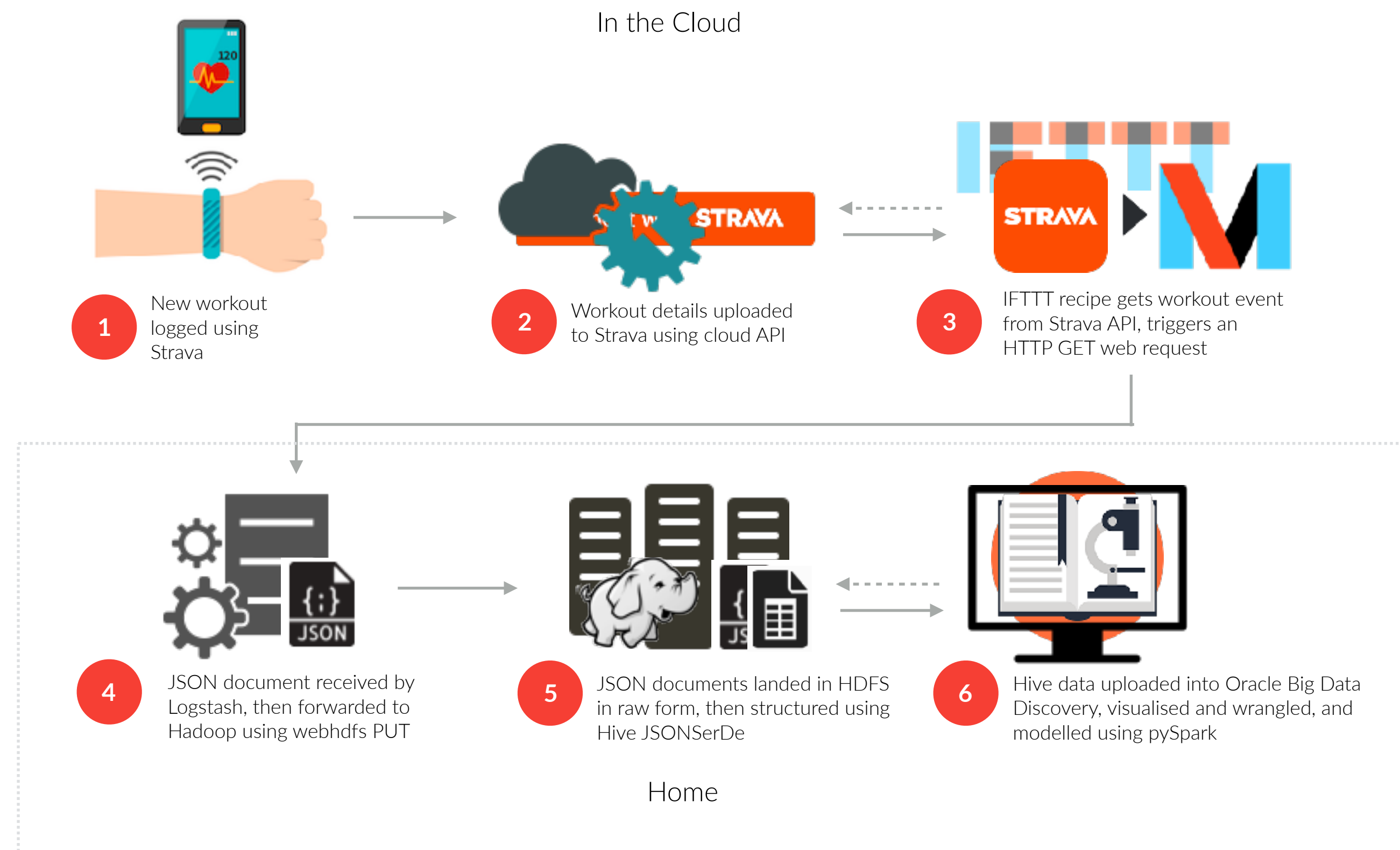- Cataloged, visualised and analysed using Oracle Big Data Discovery + Python ML

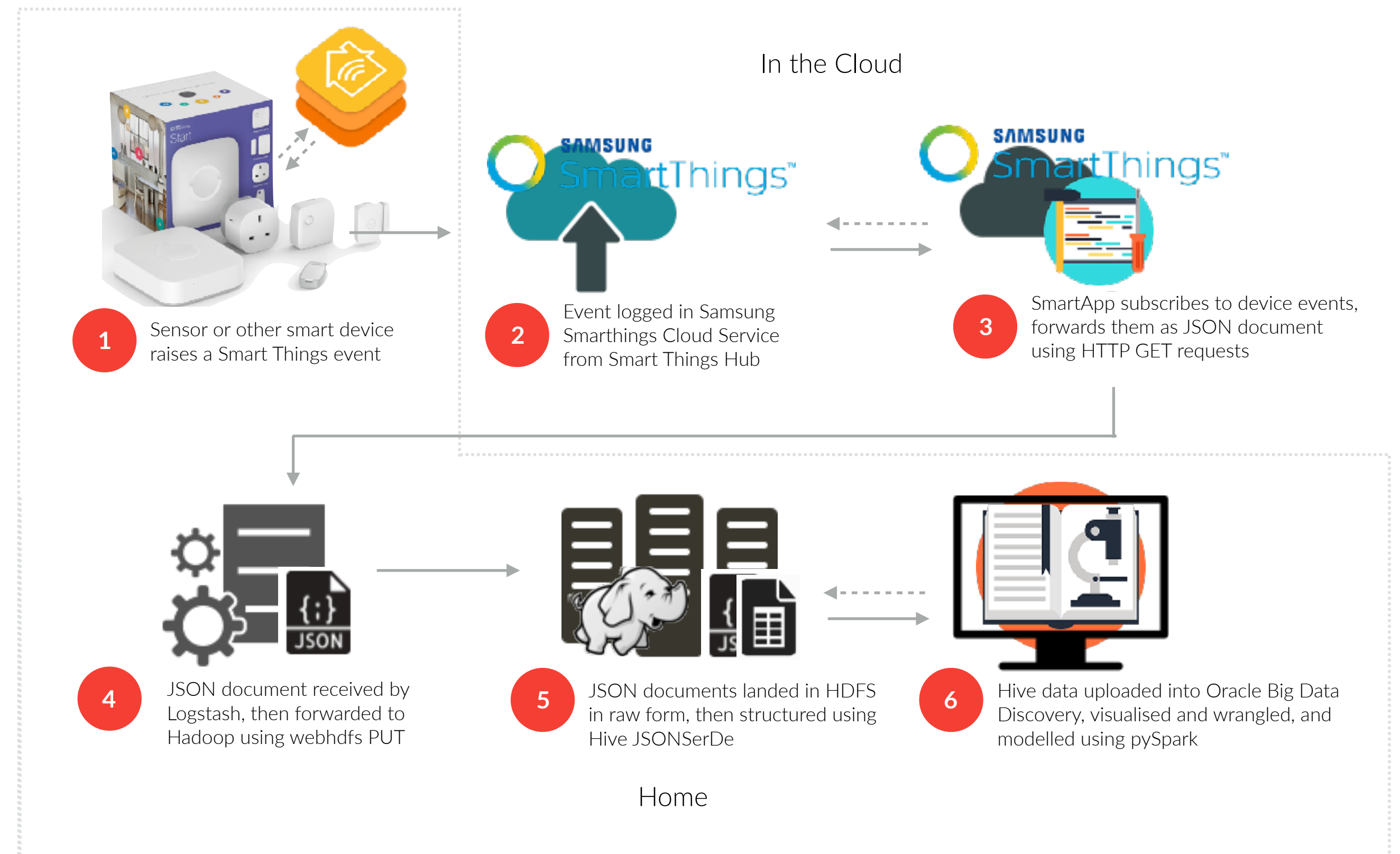# Real-Time Logging of IoT + Wearable Activity Data



- Gmail
- Withings Scales
- Strava
- Jawbone UP
- Weather
- Youtube
- IOS Photos

- Twitter
- RescueTime
- Pocket
- Instagram
- Google Calendar
- Facebook

(real-time)

Apple Homekit, Apple TV, Siri

Nest Protect (X2), Thermostat, Cam

Homebridge Homekit / Smarthings Connector

IFTTT Maker Channel JSON via HTTP POST

(real-time)

6-Node CDH5.8 Hadoop Cluster in garage, + Oracle Big Data Discovery 1.2.0 on VMWare ESXi 4-node cluster

Philips Hue Lighting

Samsung Smart Things Hub (Z-Wave, Zigbee)

SmartThings

(real-time)

(real-time)

Withings Smart Scales

Airplay Speakers

Door, Motion, Moisture, Presence Sensors

LogStash

18

# Landing Wearables Data In Real-Time

- Uses IFTTT cloud workflow service to subscribe to events on wearables' APIs
- Triggers HTTP GET request via IFTTT Maker Channel to Logstash running at home
- Event data sent as JSON documents, loaded into HDFS via webhdfs protocol
- Structured in Hadoop using Hive JSONSerDe
- Then loaded hourly into DGraph using Big Data Discovery dataprocessing CLI
  - Event data automatically enriched, and can be joined to smart home data for analysis

In the Cloud

**1** New workout logged using Strava

**2** Workout details uploaded to Strava using cloud API

**3** IFTTT recipe gets workout event from Strava API, triggers an HTTP GET web request

**4** JSON document received by Logstash, then forwarded to Hadoop using webhdfs PUT

**5** JSON documents landed in HDFS in raw form, then structured using Hive JSONSerDe

**6** Hive data uploaded into Oracle Big Data Discovery, visualised and wrangled, and modelled using pySpark

Home

# Landing Smart Home Data In Real-Time

- All smart device events and sensor readings are routed through Samsung Smart Things hub
  - Including Apple HomeKit devices, through custom integration
- Event data uploads to Smart Things cloud service + storage
- Custom Groovy SmartApp subscribes to device events, transmits JSON documents to Logstash using HTTP GET requests
- Then process flow the same as with wearables and social media / comms data

In the Cloud

**1** Sensor or other smart device raises a Smart Things event

**2** Event logged in Samsung Smarthings Cloud Service from Smart Things Hub

**3** SmartApp subscribes to device events, forwards them as JSON document using HTTP GET requests

**4** JSON document received by Logstash, then forwarded to Hadoop using webhdfs PUT

**5** JSON documents landed in HDFS in raw form, then structured using Hive JSONSerDe

**6** Hive data uploaded into Oracle Big Data Discovery, visualised and wrangled, and modelled using pySpark

Home

# Initial Focus Area : What Drives Weight Gain/Loss?

- This combined dataset can potentially be used to answer some interesting questions
- For example … "which of my daily activities or behaviours has most influence on my weight?"
  - Is it amount of exercise? amount of sleep? What I eat? How much work I'm doing in evenings?
- Objective is to work out which variable has the most influence on % weight change wk/wk
  - Will require tidying/reformatting of data feeds to standardise dates, bin and transform data
  - Dealing with nulls where workouts, weight readings were missed on certain days
  - Aggregating and joining different datasets
  - Build linear regression model to identify
    most influential variable

MODELING AND INFERRING

# Perform Exploratory Analysis On Data

- Understand the "spread" of data using histograms
- Use box-plot charts to identify outliers and range of "usual" values
- Sort attributes by strongest correlation to a target attribute

# Transform ("Wrangle") Data To Standardise & Tidy

• Initial row-wise preparation and transformation of data using Groovy transformations

# Dealing With Missing Recordings In The Data

- Very typical with self-recorded healthcare and workout data
- Most machine-learning algorithms expect every attribute to have a value per row
- Self-recorded data is typically sporadically recorded, lots of gaps in data
- Need to decide what to do with columns of poorly populate values

# Joining Wearables Data With Comms + Smart Devices

- Previous versions of BDD allowed you to create joins for views
  - Used in visualisations, equivalent to a SQL view i.e. SELECT only
- BDD 1.2.x allows you to add new joined attributes to data view, i.e. materialise
- In this instance, use to bring in data on emails, and on geolocation

# Aggregate Data Up To The Week Level

- Only sensible option when looking at change in weight compared to prior period
  - Change compared to previous day too granular

# Use of BDD Shell, Python Pandas + Jupyter Notebook

- Now we have the data organised into weekly reading rows, we now switch to **Python Pandas**
- Use this Python statistics and data visualisation library to **calculate w/o/w weight change**, and **identify most influential variable** (i.e. action, activity type I've recorded)
  - Use BDD Shell to connect to BDD SDK from pySpark environment
  - Work with BDD datasets as Spark dataframes
  - Import and use Python Pandas and SparkML packages
  - Shape and transform dataframes further if needed
  - Use visualizations to understand correlations between variables
  - Create linear regression ML model to identify most influential variable

MODELING AND INFERRING

# Use BDD Shell API to Identify Main Dataset ID

```
In [1]:  execfile('ipython/00-bdd-shell-init.py')

In [31]: dss = bc.datasets()
         dss.count

Out[31]: 76

In [32]: for ds in dss:
             print('Name: %s\t

         Name: ifttt_comms_ema
         Name: ifttt_comms_twi
         e8c8240
         Name: ifttt_health_ev
         Name: Combined Health
         Name: Combined Health
         7b-df15d8580b29
         Name: Jawbone UP Mark
         Name: Combined Health
         70-ac05e7699f50
         Name: ifttt_comms_eve
         Name: smarthings_log_
         Name: Combined Aggreg
         eec6914
         Name: Combined Health
         Key:default_edp_4f930
         Name: Jawbone UP Mark
         Name: ifttt_health_ev
         Name: Heart Rate - Ma
         Name: ifttt_comms_fac
         73f6e60
```

```
In [34]: ds = dss.dataset('default_edp_07c07ea5-891e-40ae-b2ca-b6d85f68b9e1')

         import json
         print json.dumps(ds.properties(),indent=2,sort_keys=True)

         {
           "accessType": "public_default",
           "attributeCount": "14",
           "attributeDisplayNames": "workout_duration__mins__avg",
           "attributeKeys": "workout_duration__mins__avg",
           "attributeNotes": null,
           "attributeSemanticTypes": null,
           "attributeTags": null,
           "authorizedGroup": null,
           "authorizedReadGroup": null,
           "authorizedReadUser": "10098",
           "authorizedUser": "10098",
           "collectionIdToBeReplaced": null,
           "collectionKey": "default_edp_07c07ea5-891e-40ae-b2ca-b6d85f68b9e1",
           "creationTime": "2016-06-23T22:50:52.787Z",
           "curated": "false",
           "databaseKey": "default_edp_07c07ea5-891e-40ae-b2ca-b6d85f68b9e1",
           "dateTimePresent": "false",
```

# Use Python PANDAS to Calculate % CHG W/w

# Identify Correlations Between Attributes

# Use Linear Regression on BDD Dataset via Python

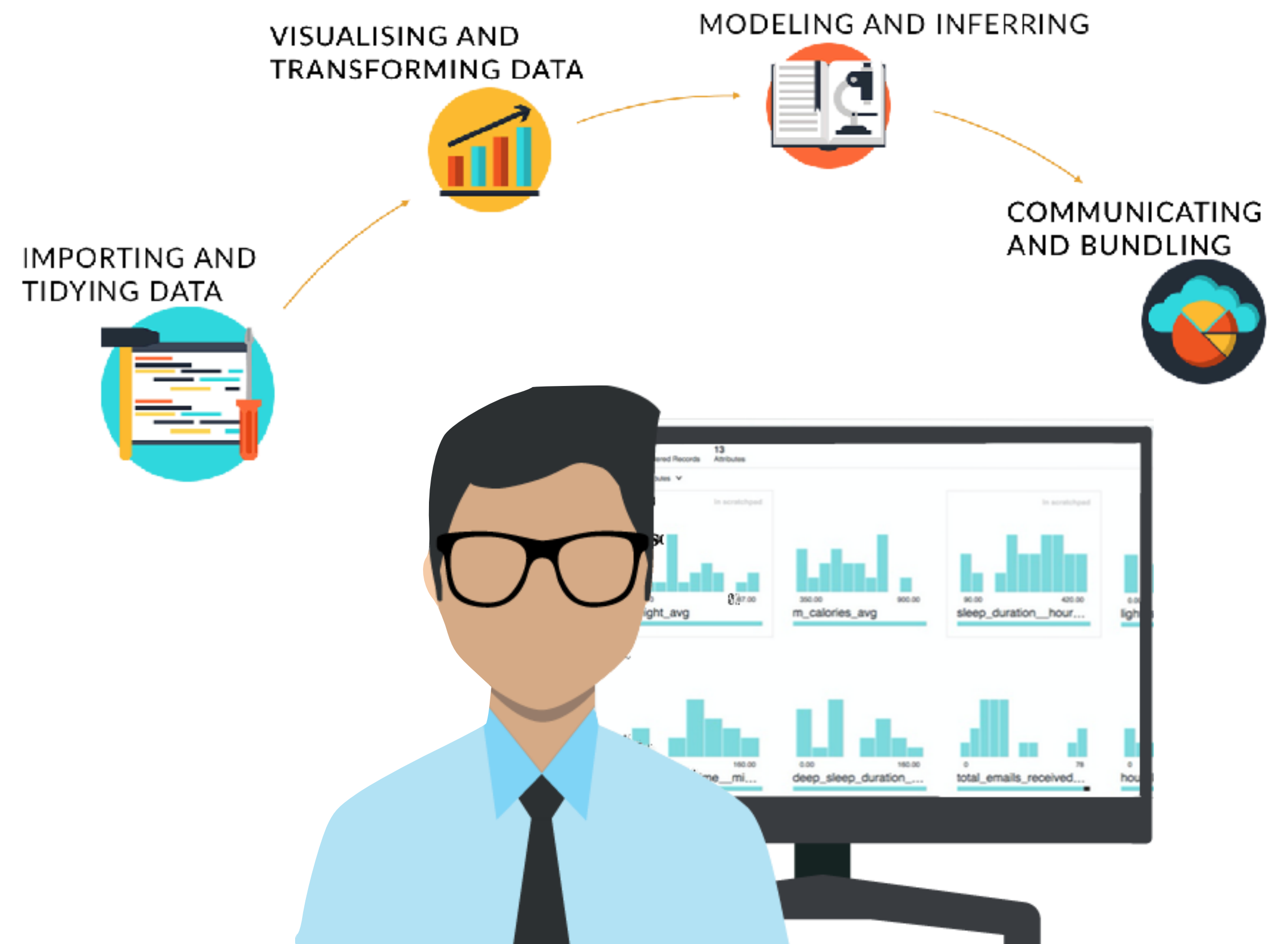- To answer the question - which metric is the most influential when it comes to weight change?

# The Answer? … Hours of Sleep Most Influential Activity

- Most influential variable/attribute in my weight / loss gain is hours of sleep per week
  - The more sleep i get, the more likely I am to exercise, walk somewhere, eat properly and lose weight
  - Weeks where less sleep recorded led to eating more carbs, driving rather than walking, weight gain
- Environment (internal, external) had less influence this time, but influential variables were:
  - Comms activity - emails sent late night, Facebook likes, Instagram photos - proxy for working/play
  - Heat/Temperature inside house - indicates warm/cold outside, driver of exercise activity
  - Geo-location - am I on holiday? At work that week?
  - Diet? Although fairly constant over perio



29

# How Did Oracle BDD Help With This Project?

- Visual, graphic way to understand shape, data distribution and outliers/completeness
- Simple user-driven graphical tools for data tidying and transformation
- Join and aggregate datasets to get to one row of data = set of weekly readings
- Enrich and bring in additional datasets to add comms and environment activity data
- Enable use of wide range of industry-standard stats and ML libraries on final dataset

# PRACTICAL IOT DEVELOPMENT USING ORACLE BIG DATA AND ORACLE DV …AND A WIFI KETTLE

Mark Rittman, Oracle ACE Director & Independent Analyst
MJR Analytics ltd (http://www.mjr-analytics.com)

BIWA SUMMIT 2017, SAN FRANCISCO

T : @markrittman