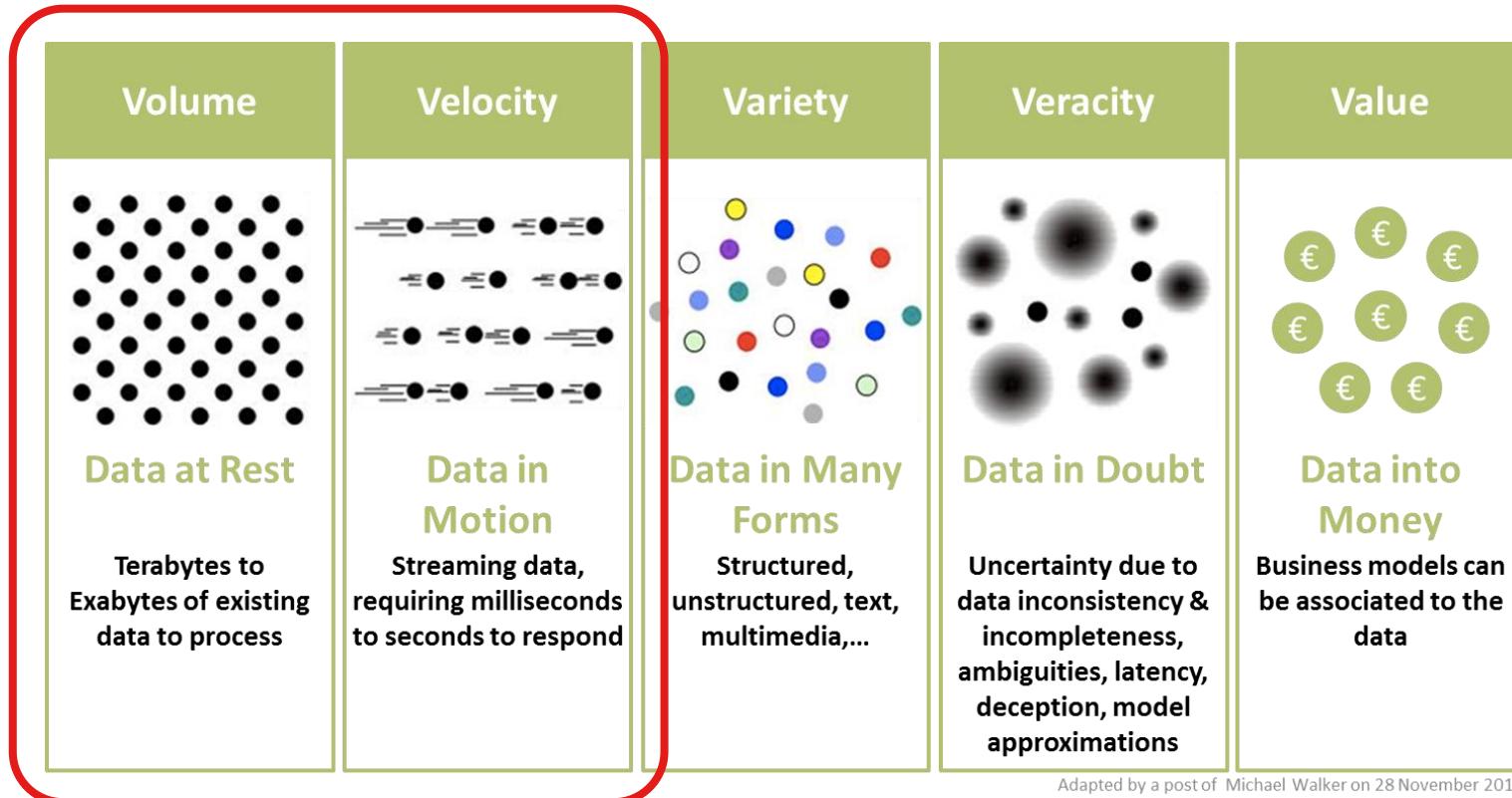


# DATA DRIVEN BUSINESS MODELS

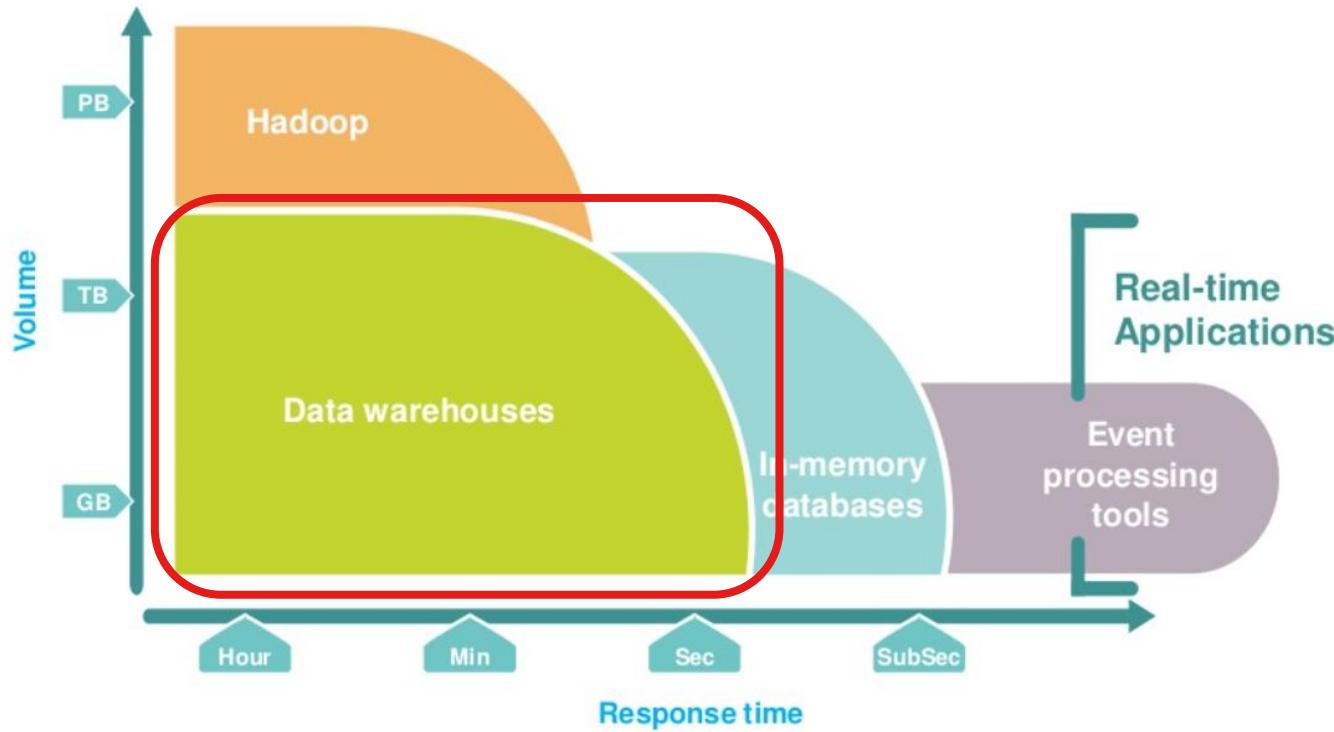
Big Data Technology, Data Science

# Store and Process

# The 5 V's of Big Data



Adapted by a post of Michael Walker on 28 November 2012



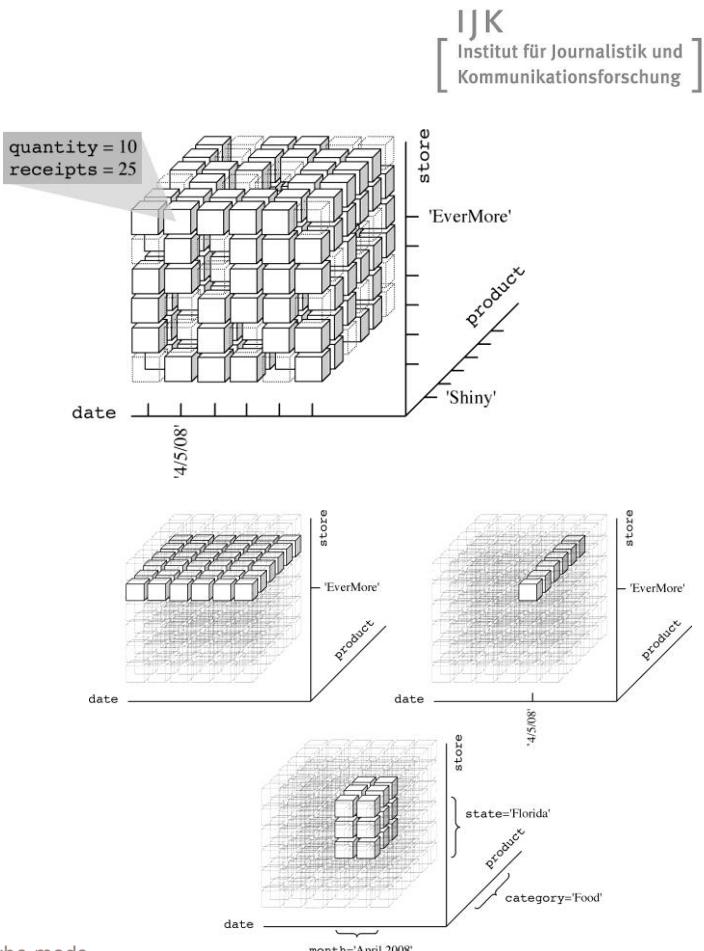
# OLTP – Transactional Workload

- Most business maintain a variety of databases to keep track of what happens in the organization.
- These databases are optimized for consistency and speed on a „transactional level“. It's fast and safe to update individual records about stock items, accounts receivable etc.
- They are not very good to create summaries and reports because this requires to scan all records.

09321		J.D. Edwards & Company Transaction Journal			
Company	Tax ID	Document No.	Description	Debit Amount	Credit Amount
PV	4155 00000	100.2060	Furniture & Office Equipment	2,487.61	
		100.2060	Furniture & Office Equipment	14,978.95	17,466.56
PV	4156 00000	100.4110	Accounts Payable-Trade	3,458.91	
		100.4110	Furniture & Office Equipment		3,458.91
PV	4216 00000	122.8155	Training Expenses	500.00	
		100.4110	Accounts Payable-Trade		500.00
PV	4252 00000	210.8360	Telephone Expense	1,465.61	
		100.4110	Accounts Payable-Trade		1,465.61
PV	4253 00000	400.8360	Telephone Expense	846.61	
		100.4110	Accounts Payable-Trade		846.61
PV	4254 00000	600.8360	Telephone Expense	1,006.74	
		100.4110	Accounts Payable-Trade		1,006.74
PV	4267 00000	529.6140	Tools Expense	1,867.00	
		529.6120	Prime Cost of Goods	8,855.00	
		529.6130	Scrap	4,155.00	
		529.6140	Freight	6,118.00	
		100.4110	Accounts Payable-Trade		18,625.00
PV	4277 00000	100.2060	Furniture & Office Equipment	15,967.62	
		100.2060	Furniture & Office Equipment	31,786.43	
		100.2060	Furniture & Office Equipment	15,689.23	
PV	4279 00000	529.8975	Water	755.00	
		529.8975	Water	615.00	
		529.8977	Sanitation	485.00	
		100.4110	Accounts Payable-Trade		66,443.28
PV	4297 00000	100.8360	Telephone Expense	1,245.95	
		100.4110	Accounts Payable-Trade		1,245.95
PV	4299 00000	210.8360	Telephone Expense	1,245.95	
		100.4110	Accounts Payable-Trade		1,245.95
PV	4300 00000	90.8350	Purchase Part Variance	500.55	
		100.4110	Accounts Payable-Trade		500.55
PV	4301 00000	600.8350	Rent Expense	1,801.00	
		100.4110	Accounts Payable-Trade		1,801.00
PV	4303 00000	90.8350	Rent Expense	2,200.00	
		100.4110	Accounts Payable-Trade		2,200.00
		90.8350	Rent Expense	200.15	
PV	4381 00000	100.4110	Accounts Payable-Trade		200.15
		90.8175	Uniforms	5,581.93	
PV	4623 00000	100.4110	Accounts Payable-Trade		5,581.93
		90.8175	Uniforms	11,428.84	
PV	4914 00000	100.4110	Accounts Payable-Trade		11,428.84
		90.8175	Uniforms	8,909.24	
PV	5447 00000	100.4110	Accounts Payable-Trade		8,909.24
		90.8320	Insurance - General Liability	5,560.50	
		90.8355	Maintenance & Repair	5,560.50	
		100.4110	Accounts Payable-Trade		11,121.00
PV	8324 00100	90.8700	Miscellaneous Expenses	1,500.00	
		100.4110	Accounts Payable-Trade		1,500.00
PV	8329 00100	90.8665	Entertainment	825.00	
		100.4110	Accounts Payable-Trade		825.00
PV	8330 00100	90.8740	Travel, Meals & Lodging		
		100.4110	Accounts Payable-Trade		

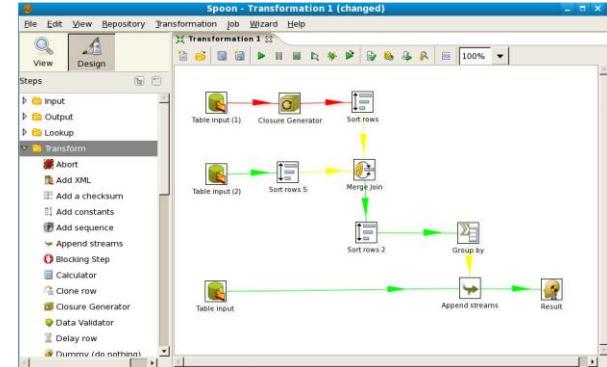
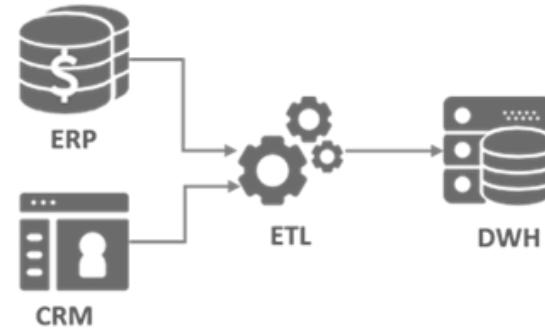
# OLAP – Analytical Workload

- Instead of scanning through all records for each analytical query specialized databases create „cubes“ that hold pre-aggregated sums, averages etc. for all possible combinations of dimensions.
- In addition some dimensions might be hierarchical which allows to aggregate on several hierarchy levels (e.g. day, month, year)

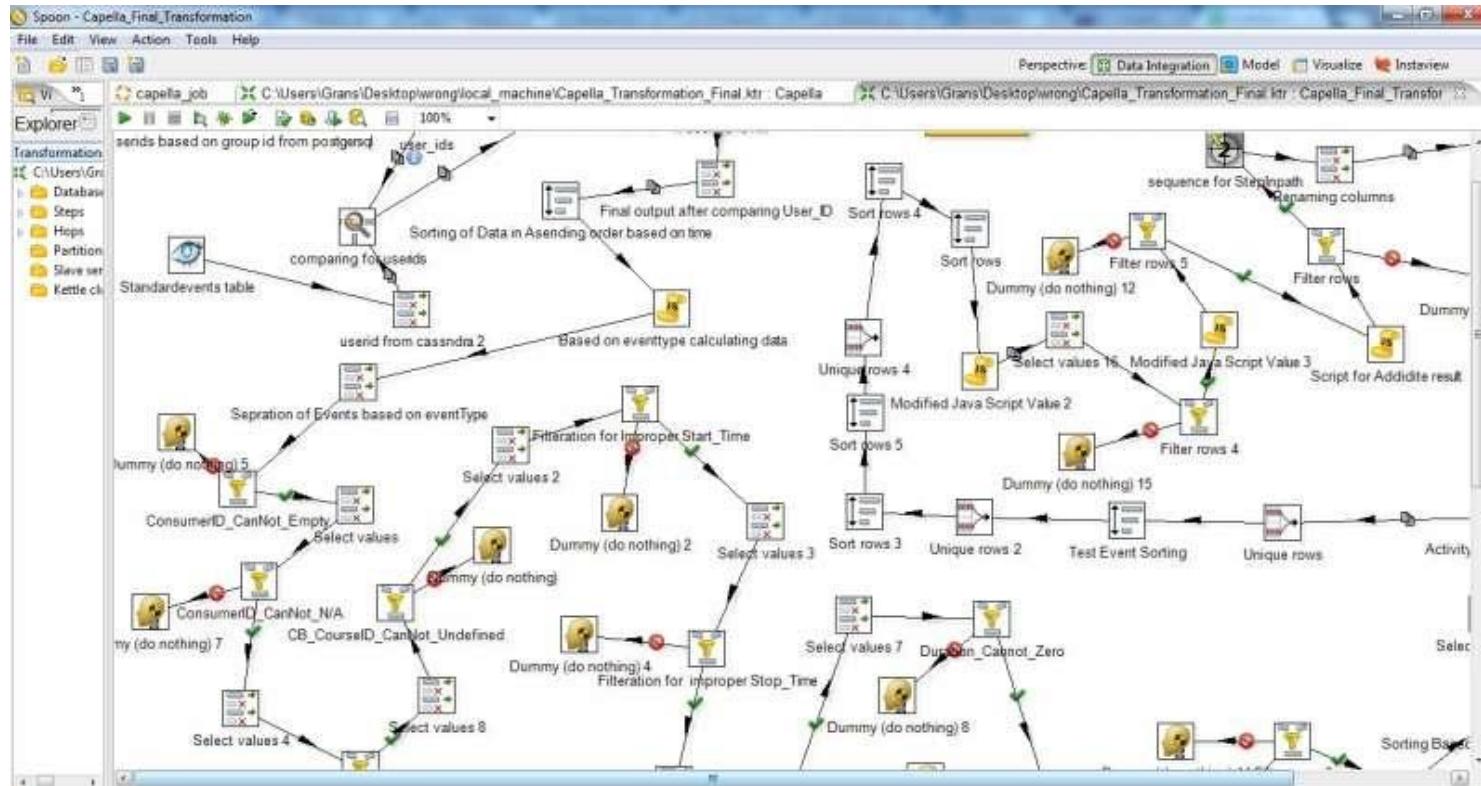


# ETL – Extract, Transform, Load

- For defining and periodically updating OLAP cubes there are specialized tools available that help to
- connect and find relevant data in operational OLTP systems (Extract)
- clean, augment and aggregate data across dimensions and levels (Transform)
- Store aggregated results into the data warehouse (Load)
- The ETL pipelines tend to become complex very quickly and are notoriously hard to maintain



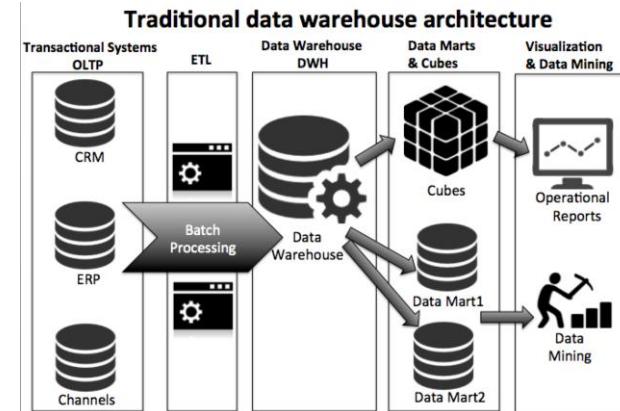
# ETL – Extract, Transform, Load



<https://www.truelancer.com/portfolio/pentaho-kettle-complex-transformation-example-293>

## What's wrong with that?

- Data Warehouses do contain aggregate data only. Lack of transactional details.
- Aggregates and Data Prep is „fixed“. Long turn around time for changes.
- Data often is stale. Update frequency often limited to 1-2 days or worse.
- Integration of new data sources can be a tedious process. Some (unstructured) data types cannot be integrated at all
- ETL and DWH are often owned by a team which becomes a bottleneck quickly
- Specialized hardware and commercial tools are very costly (typically starting at >\$500k)
- **Data Warehouses cannot keep up with the changing requirements**



Increased number and variety of data sources that generate large quantities of data.



Realization that data is “too valuable” to delete.

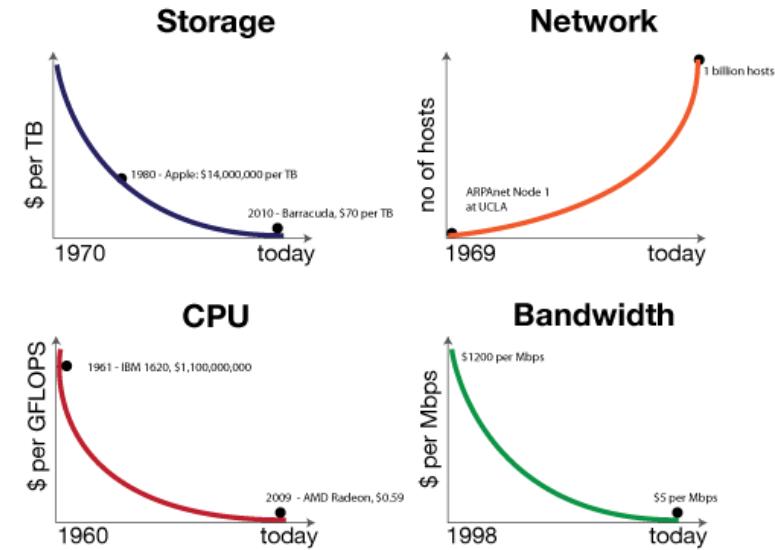


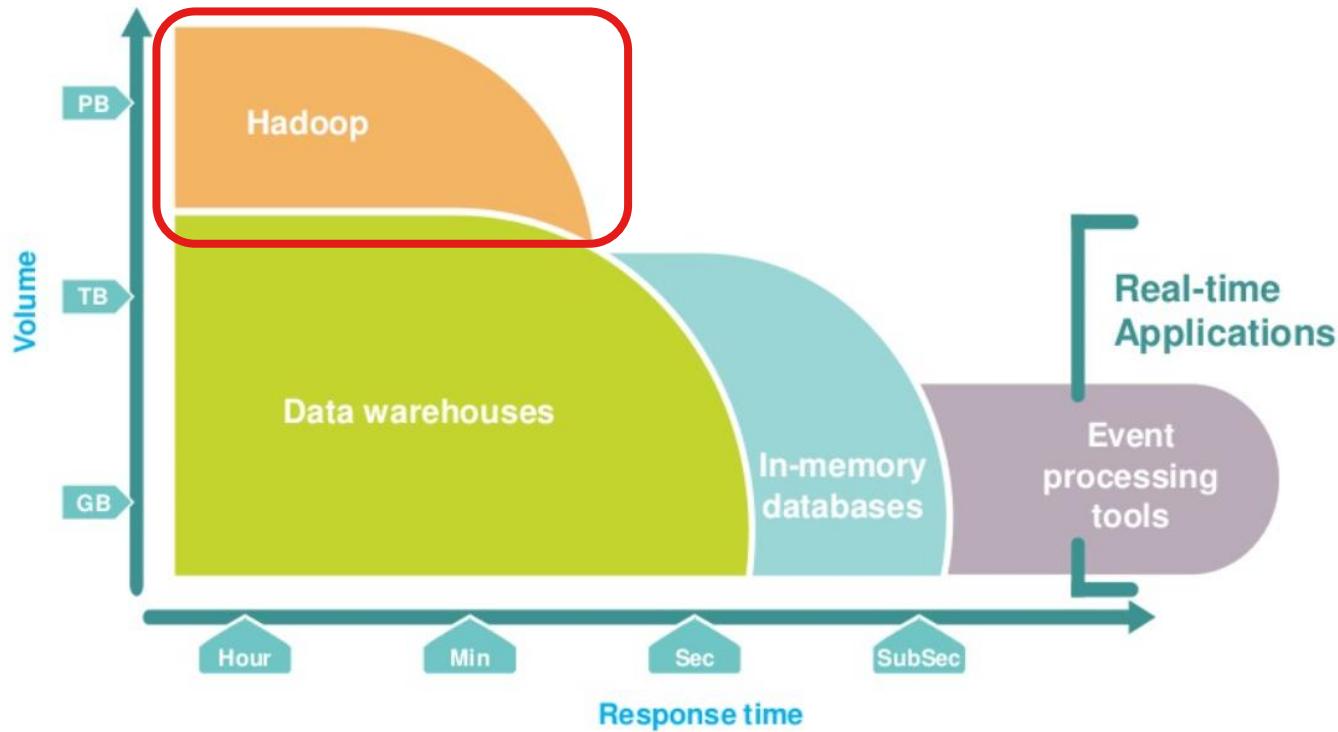
Dramatic decline in the cost of hardware, especially storage.



# Technology allows to Scale Out

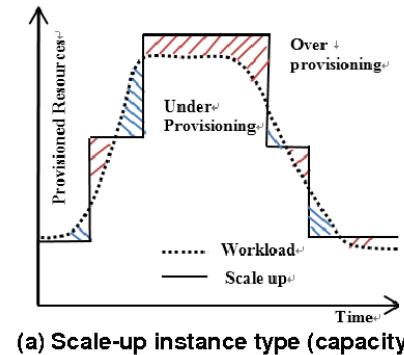
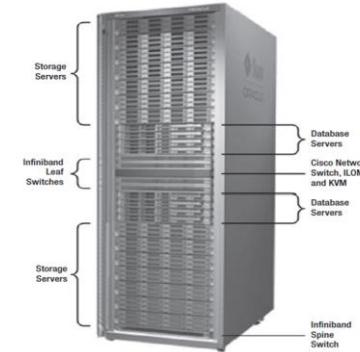
- Historically storage was very expensive.  
You ought to be very clever about how  
to minimize storage requirements.
- Traditional OLAP systems therefore use  
ETL to aggregate early and transfer  
data to where the logic is.
- Using commodity hardware, cost of  
storage and network bandwidth is  
dramatically decreasing
- Keep all the data and run analytics  
directly on the raw data.



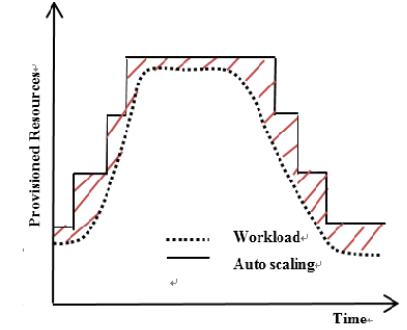


# Scale-Up vs Scale-Out

- Traditional systems have been able to scale up. You have to buy a bigger machine to prepare for heavier workloads.
- Economically very challenging to plan for. What is the right timing? What do you do with the old machine?
- The decline of network cost allows for a much more „elastic“ approach by combining many small machines on-demand.



(a) Scale-up instance type (capacity)



(b) Scale-out in instance quantity

# Add Software to the mix.

- If you can access all your data in its raw form you have maximum flexibility.
- Because it's a lot of data you need to distribute the load across many machines
- This happens to be very similar to what search engines do.
- Google popularized the idea of „distributed storage“ and „distributed processing“ by publishing papers about how they approach the search problem

**The Google File System**

Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung  
*Google*

**ABSTRACT**  
We have designed and implemented the Google File System, a scalable distributed file system for large distributed data-intensive applications. It provides fault tolerance while running on inexpensive commodity hardware, and it delivers high aggregate performance to a large number of clients.  
While sharing many of the same goals as previous distributed file systems, our design has been driven by observations of our application workloads and technological environment, both current and anticipated, that reflect a marked

**1. INTRODUCTION**  
We have designed and implemented the Google File System (GFS) to meet the rapidly growing demands of Google's data processing needs. GFS shares many of the same goals as previous distributed file systems such as performance, scalability, reliability, and availability. However, its design has been driven by key observations of our application workloads and technological environment, both current and anticipated, that reflect a marked departure from some earlier file system design assumptions. We have reexamined tradit-

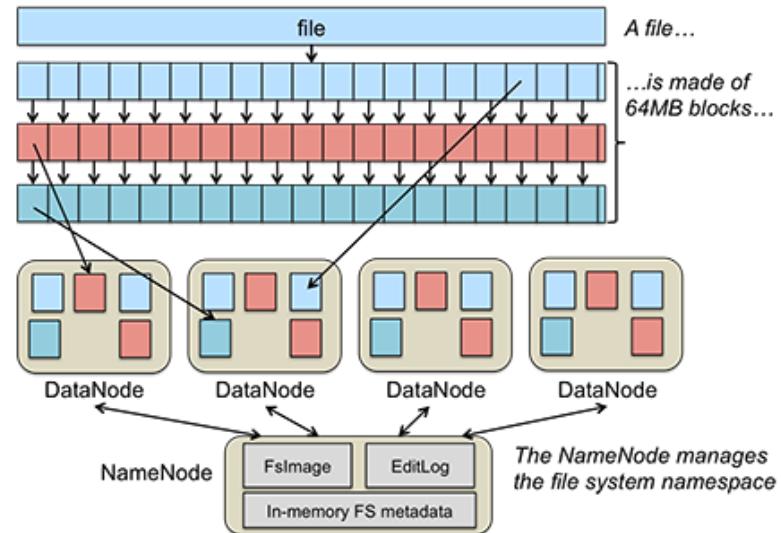
**MapReduce: Simplified Data Processing on Large Clusters**

Jeffrey Dean and Sanjay Ghemawat  
*jeff@google.com, sanjay@google.com*  
*Google, Inc.*

**Abstract**  
MapReduce is a programming model and an associated implementation for processing and generating large data sets. Users specify a *map* function that processes a key/value pair to generate a set of intermediate key/value pairs, and a *reduce* function that merges all intermediate values associated with the same intermediate key. Many real world tasks are amenable in this model, as shown given day, etc. Most such computations are conceptually straightforward. However, the input data is usually large and the computations have to be distributed across hundreds or thousands of machines in order to finish in a reasonable amount of time. The issues of how to parallelize the computation, distribute the data, and handle failures conspire to obscure the original simple computation with large amounts of complex code to deal with

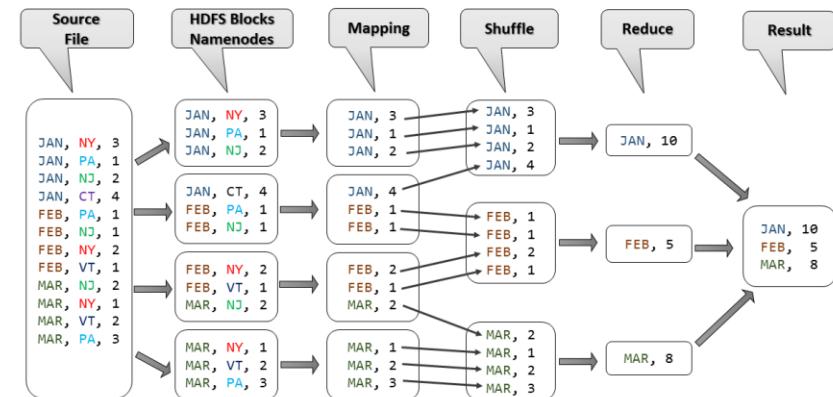
# Hadoop implements an Open Source Version of GFS

- Yahoo Engineers read the Google Papers and start building their own version called Hadoop
- Data is stored in flat files which are split up and distributed across many (small) machines. Blocks are stored redundantly because failure of single machines is inevitable.
- Adding / Removing machines allows to scale up and down.



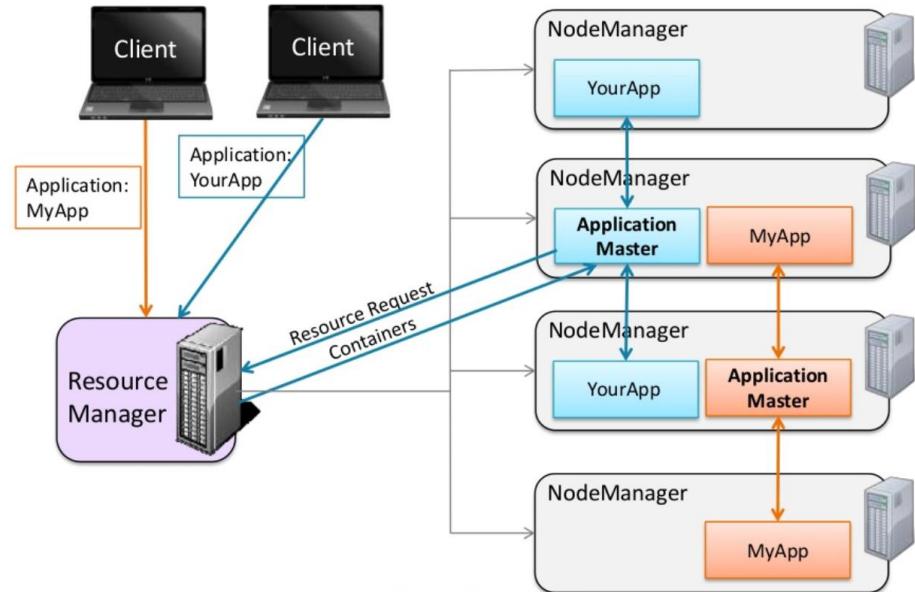
# Hadoop implements an Open Source Version of M/R

- The same Yahoo engineers also build a data processing engine into Hadoop utilizing the „Map Reduce“ approach.
- Very simple and pretty flexible
- Very easy to parallelize
- Leverages data locality of HDFS
- Hard to apply to complex questions



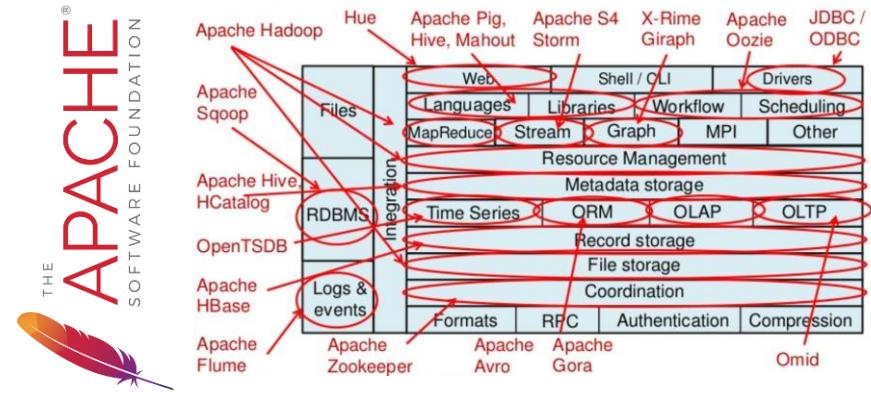
# Storage is close but separate from Processing

- Data has inertia – it's hard to move
- Big Data even more so
- If data sits in a distributed system it's easier to run your analytical application where the data is then moving data to the application.
- Hadoop supports this by leveraging info about „data locality“



# Apache Foundation vs Commercial Distributions

- Many relevant technology components for Big Data are Open Source. Core elements of the stack are maintained under the umbrella of the Apache Software Fundation.
- Enterprise users typically need a stable stack that is proven to works well together. Also they need a phone numer in case things get

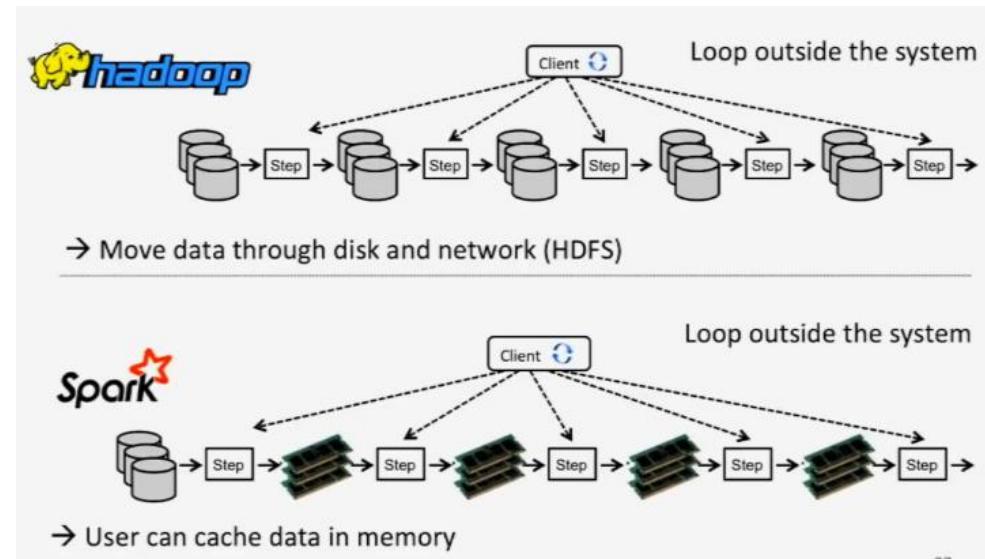


**cloudera**

**MAPR**

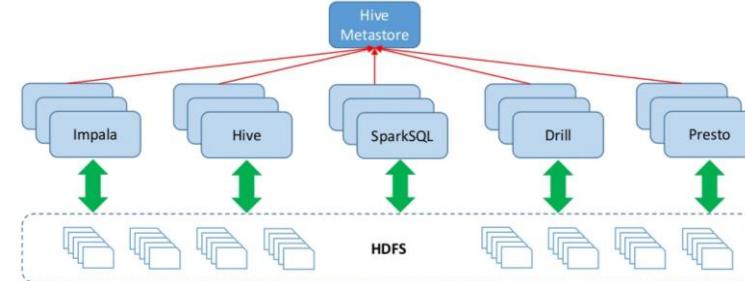
# Processing becomes more powerful when in-memory

- Map Reduce heavily relies on writing temporary results to disc.
- Spark leverages fast memory to cache temporary data between steps.
- Spark offers a richer API that allows for more complex steps than Map Reduce would enable.
- Multiple language bindings: JAVA, Scala, Python, R
- Spark can even translate SQL queries to its API calls.



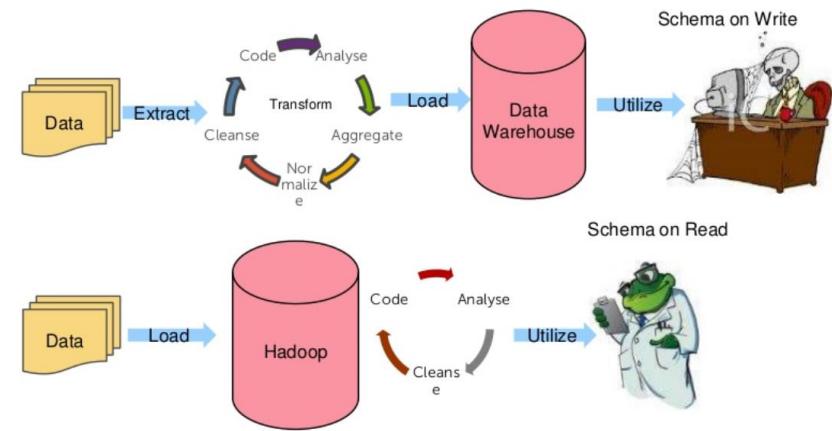
# Storage vs. Query Engines

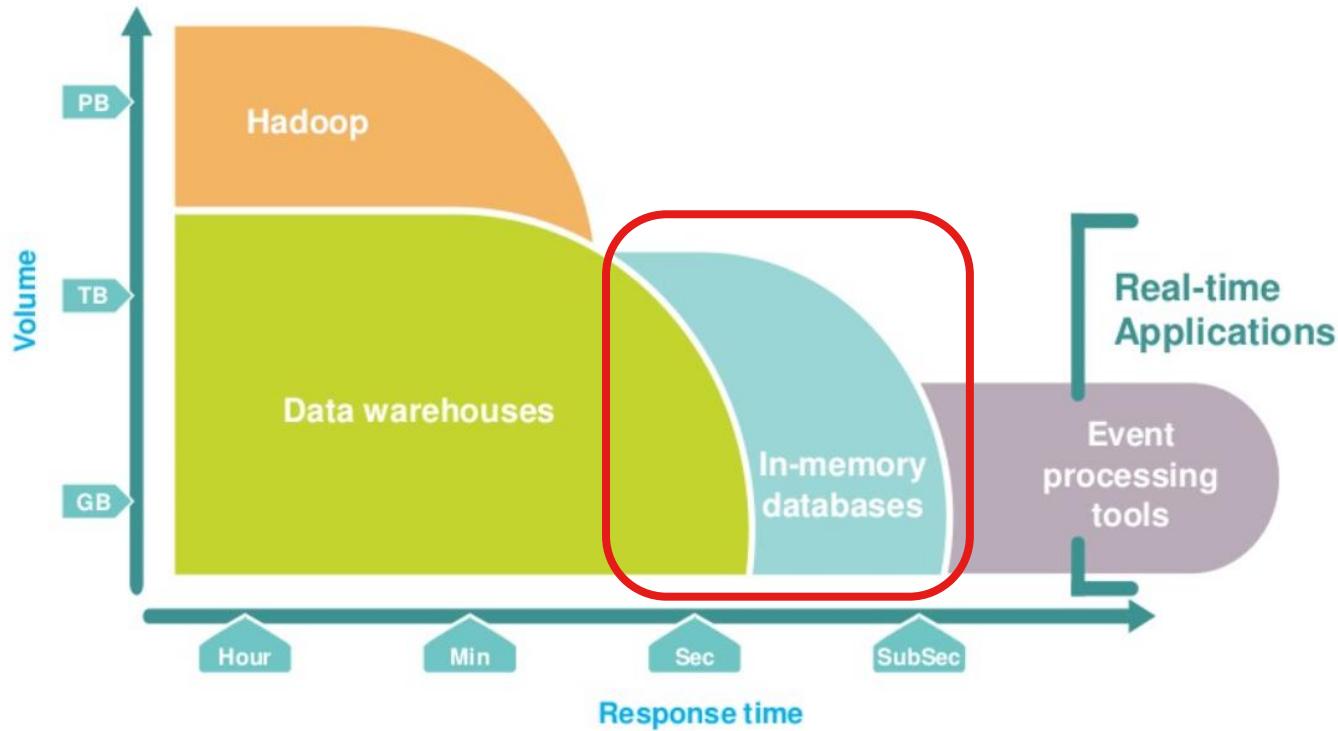
- Applications are either custom made using Python, JAVA, etc. or written in a high level scripting language like PIG or SQL
- SQL allows easy integration with BI tools which are often SQL-based
- There are competing engines optimized for specific use cases.
- The way you store your data greatly impacts the efficiency of your query.
- column-based vs. row-based



# „Schema on Read“

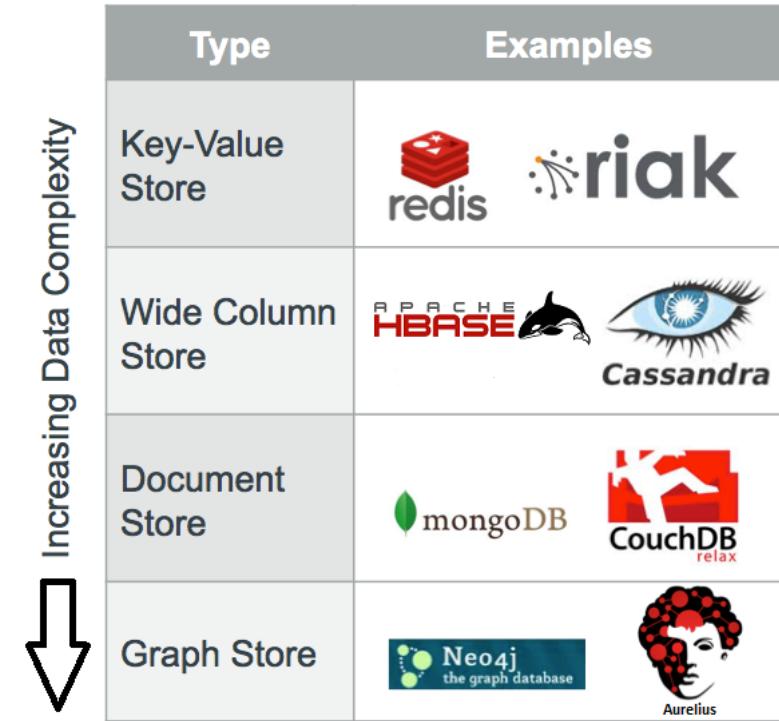
- Storing all the raw data prevents you from knowing all analytical use cases upfront.
- You cannot and want not to optimize for certain queries early in the process (=while writing the data).
- Instead you decide how to interpret the data once you know what the actual question is (=while reading the data)



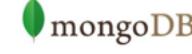
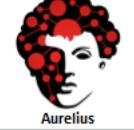


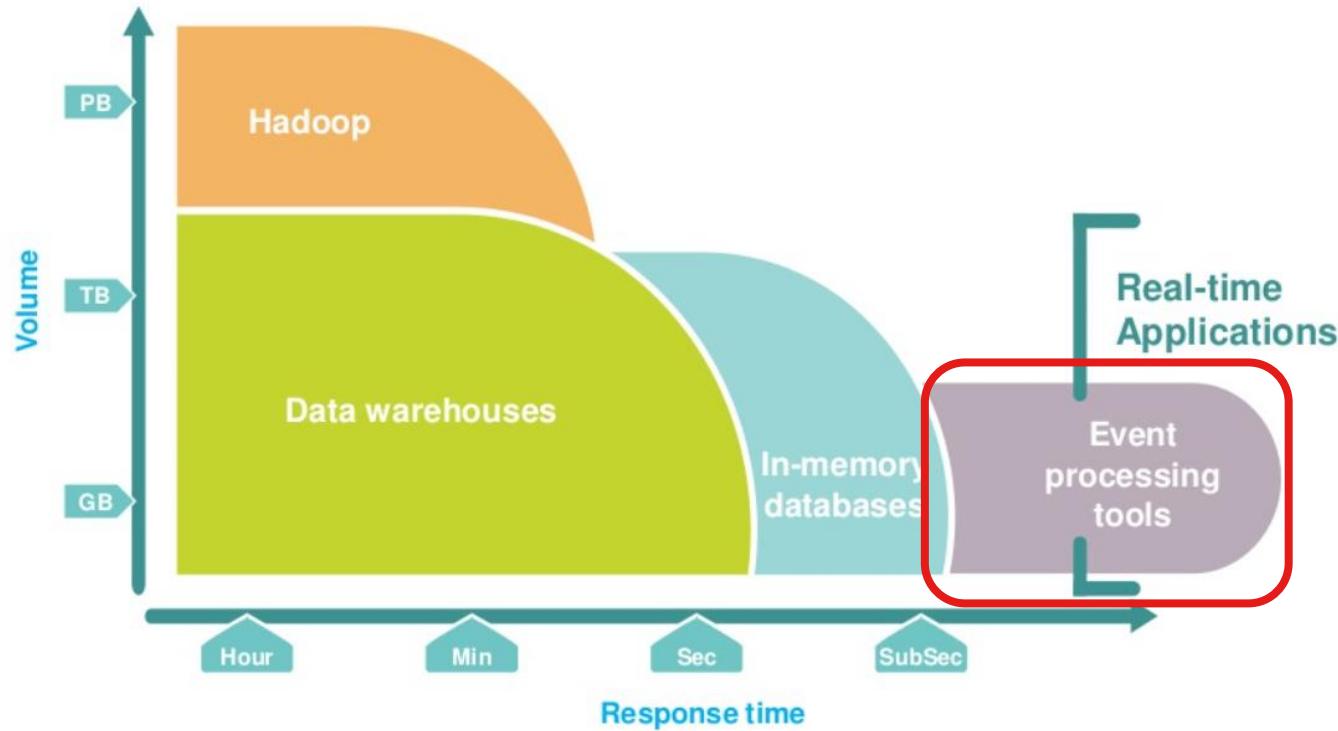
# NoSQL

- Sometimes storing data in a flat file is not enough. Especially if you need to optimize for certain access patterns.
- A lot of innovation happened in this area during the last 10 years.
- Some of these (e.g. HBase) do integrate with Hadoop for data storage and metadata management.
- Most of these engines are of limited use for analytical queries.



The diagram illustrates the increasing complexity of data storage types from top-left to bottom-right. A vertical arrow labeled "Increasing Data Complexity" points downwards between the first two columns. The table has four columns: Type, Examples, and two intermediate columns for Wide Column and Document stores.

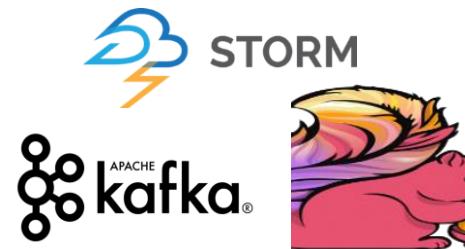
Type	Examples		
Key-Value Store			
Wide Column Store			
Document Store			
Graph Store			



# Stream and Realtime

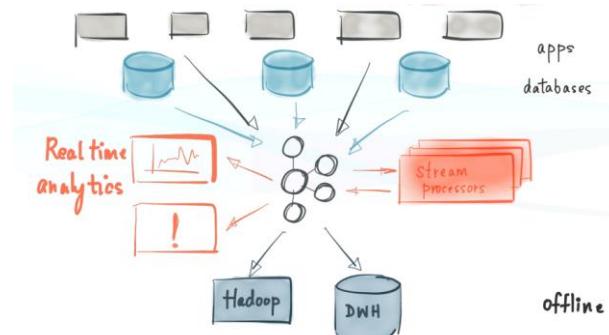
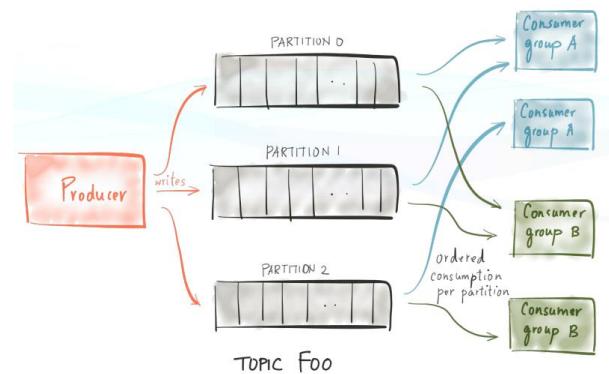
- The Hadoop ecosystem has been optimized for high throughput and for that sacrificed latency.
- For workloads that do require very low latency there are data systems that apply analytical functions while data is streaming in. Often these system use the notion of a message based Publish / Subscribe Queue.

- Hard problems:
- Exactly once delivery.
- Order of Delivery.



<https://blog.parse.ly/post/3886/pykafka-now/>

© GfK 2017 | IJK Seminar “Data Driven Business Models”



# Analyse and Visualize

# Simplistic Hadoop User Interfaces

- Apache Hue
- HDFS File Browser
- SQL Query Workbench
- Design and Schedule workflows
- Manage Processing Job

The screenshot shows the Apache Hue interface. On the left, there's a sidebar with a 'Navigator' section containing a tree view of tables and databases. The main area has tabs for 'Hive Editor' and 'Query Editor'. The 'Query Editor' tab is active, displaying a sample query for 'Salary growth' and its results. Below the query results is a bar chart titled 'salary' with 'X-Axis: description' and 'Y-Axis: salary'.

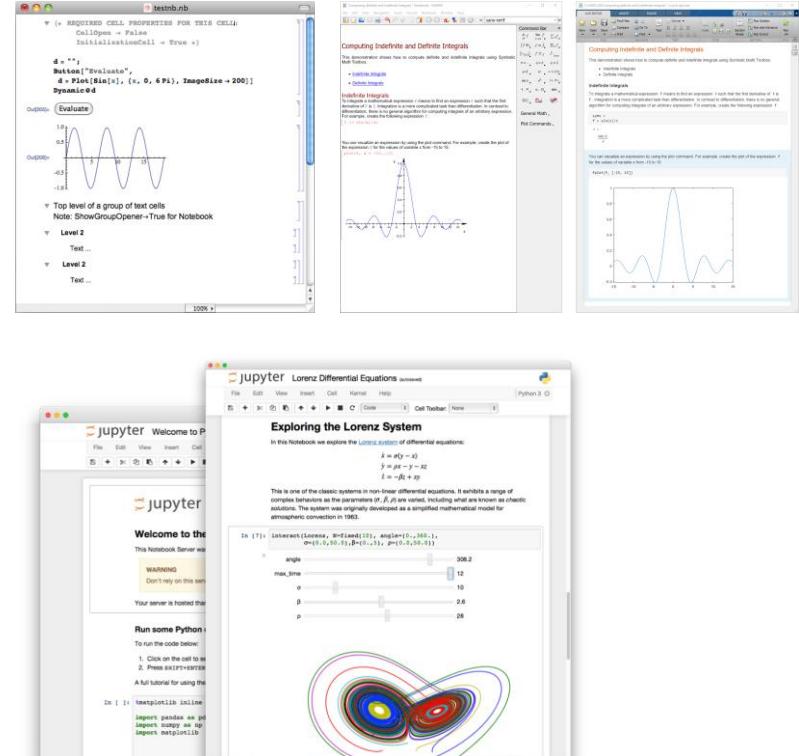
- Apache Zeppelin
- Interactive Data Science Notebooks
- Data Visualizations and Forms
- Supports multiple interpreters: R, Python, SQL

The screenshot shows the Apache Zeppelin interface. At the top, there's a navigation bar with tabs for 'Notebook' and 'Interpreter'. The main area is titled 'Zeppelin Tutorial' and contains a code cell with Scala code for reading a CSV file and performing some data processing. Below the code are three separate charts: a grouped bar chart, a stacked bar chart, and a line chart, all visualizing data from the processed bank data.

# Notebooks are pretty popular too

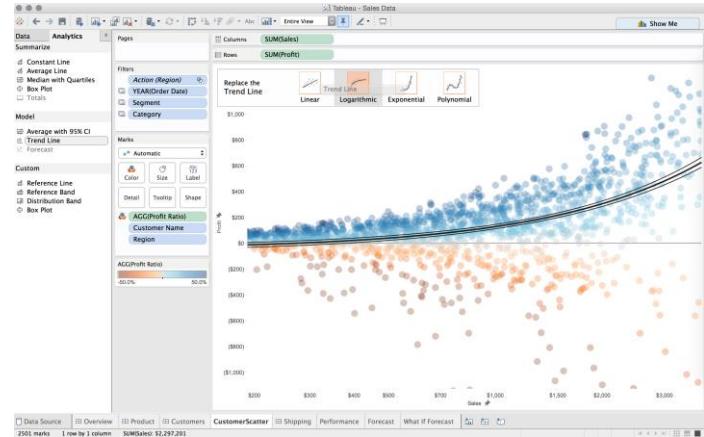
- Notebooks mix text, graphics and code
- Originates from MATLAB and Mathematica
- Browser-based IPython started ca. 2001
- Renamed to Apache Jupyter.

- Interactive execution of >40 languages
- Output of code execution inline
- Notebooks are stored as JSON files which makes it easy to version them
- Great tool for collaboration



# BI Tooling vs Data Science Tooling

- Self-Service Visual Data Exploration and Dashboard Generation
- Department Level users and subject matter experts without engineering or database background.



 Power BI

 REVOLUTION  
ANALYTICS

 SPSS

 Qlik

 + tableau®

 platfora®

 sas®

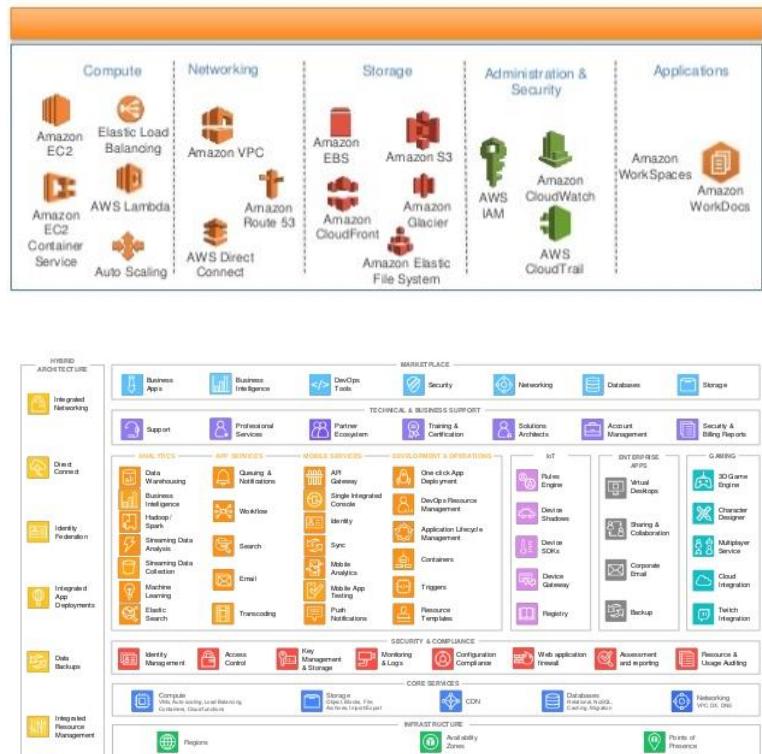
 rapidminer

 python™

# Data in the Clouds

# Cloud

- Public Cloud providers offer unlimited scale for global storage and compute
  - Typically storage is independent from compute so you can scale both independently
  - In top of the infrastructure there is a plethora of managed services, among them databases and big data clusters



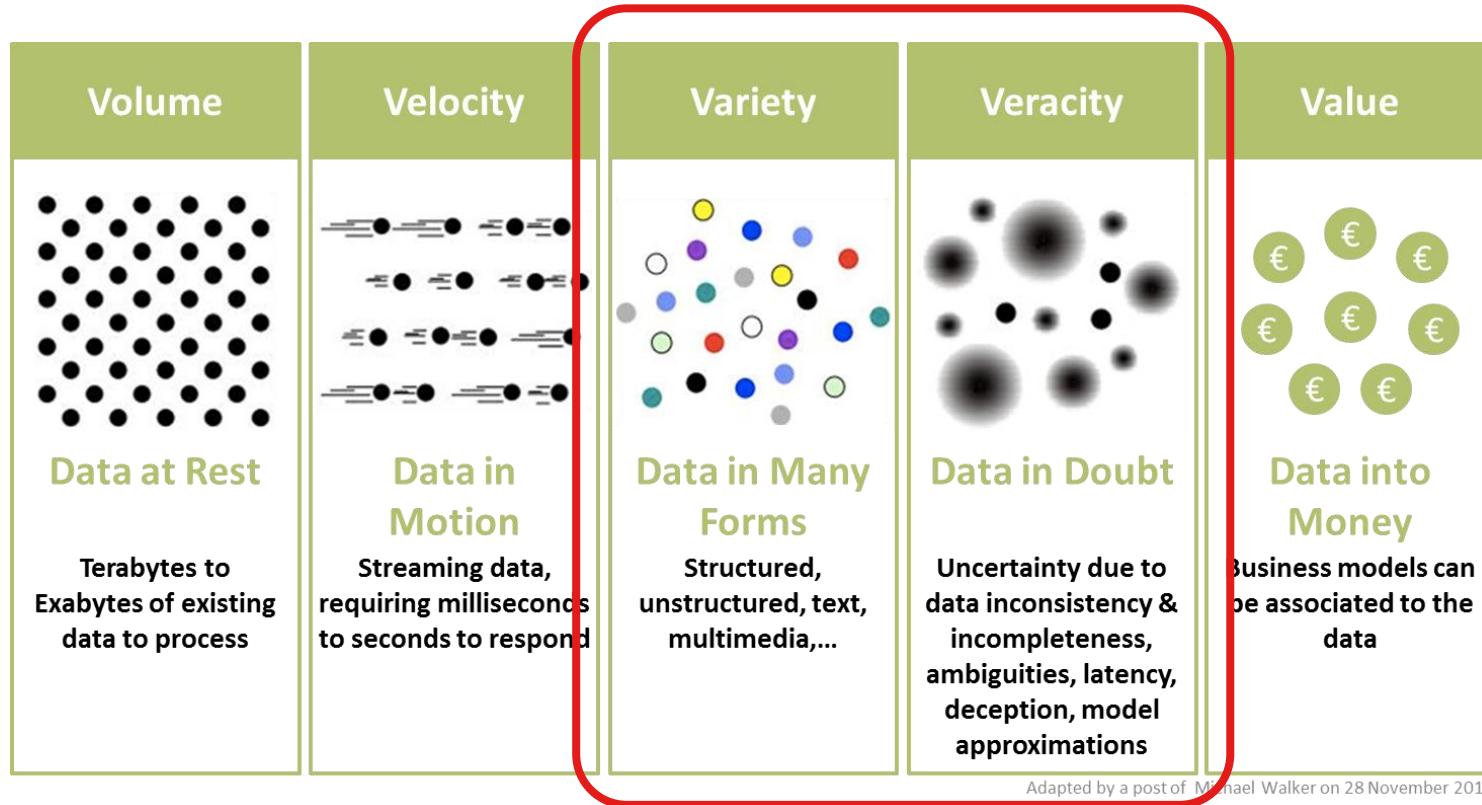
# Cloud Managed Data Services



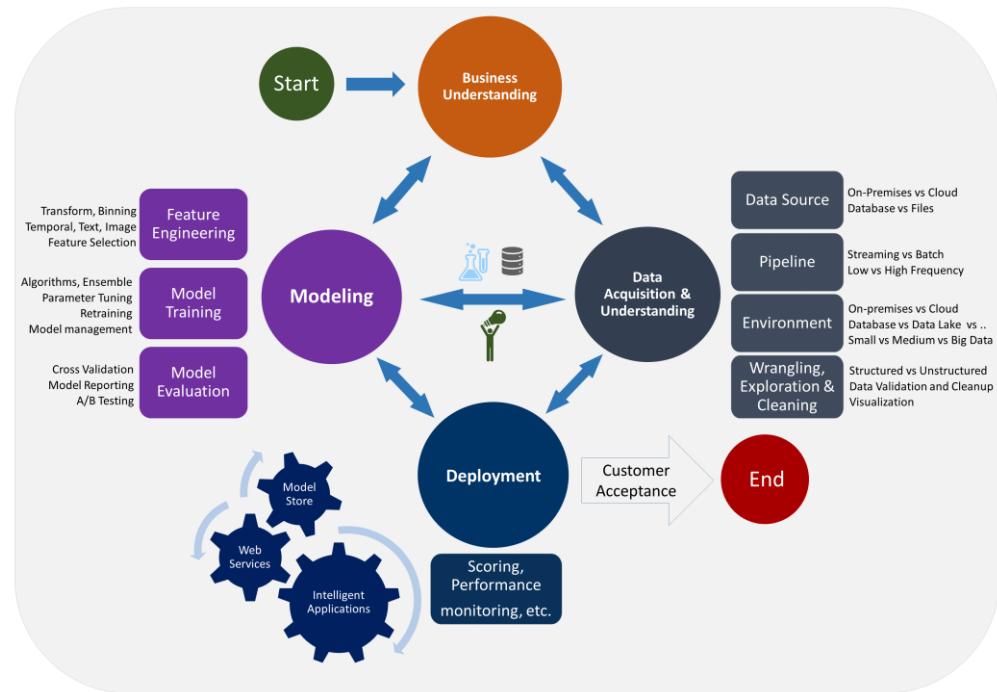
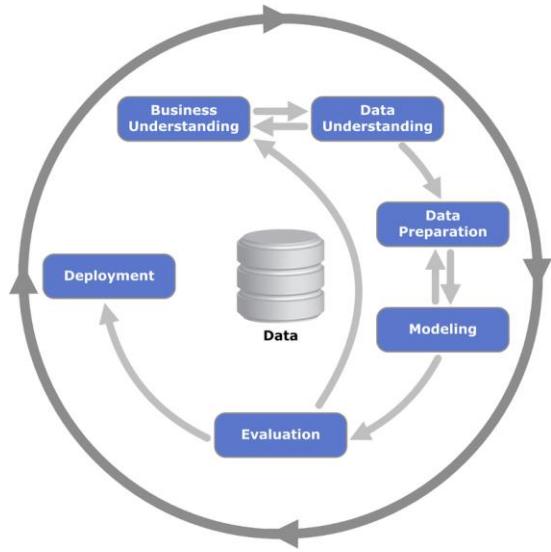
<b>Compute</b>	<i>n/a</i>	EC2	Compute	GCE
<b>Storage</b>	<i>HDFS</i>	S3	Blob	GCS
<b>Database</b>	<i>HBase</i>	DynamoDB	CosmosDB	BigQuery
<b>Processing</b>	<i>Hadoop, Spark</i>	EMR, Glue	HDInsight	DataProc
<b>Stream</b>	<i>Kafka</i>	Kinesis	Stream Analytics	DataFlow
<b>Notebook</b>	<i>Zeppelin</i>	Athena	Notebooks	DataLab
<b>AI</b>	<i>MLib</i>	AWS AI	Azure ML	TensorFlow

# Describe and Govern

# The 5 V's of Big Data



## Data Science Lifecycle



[https://en.wikipedia.org/wiki/Cross-industry\\_standard\\_process\\_for\\_data\\_mining](https://en.wikipedia.org/wiki/Cross-industry_standard_process_for_data_mining)

<https://docs.microsoft.com/en-us/azure/machine-learning/team-data-science-process/overview>

# Example 1: GfK Data Activation

Working with DMP partners as platforms allows GfK to address new business

# Marketing Services Initiative

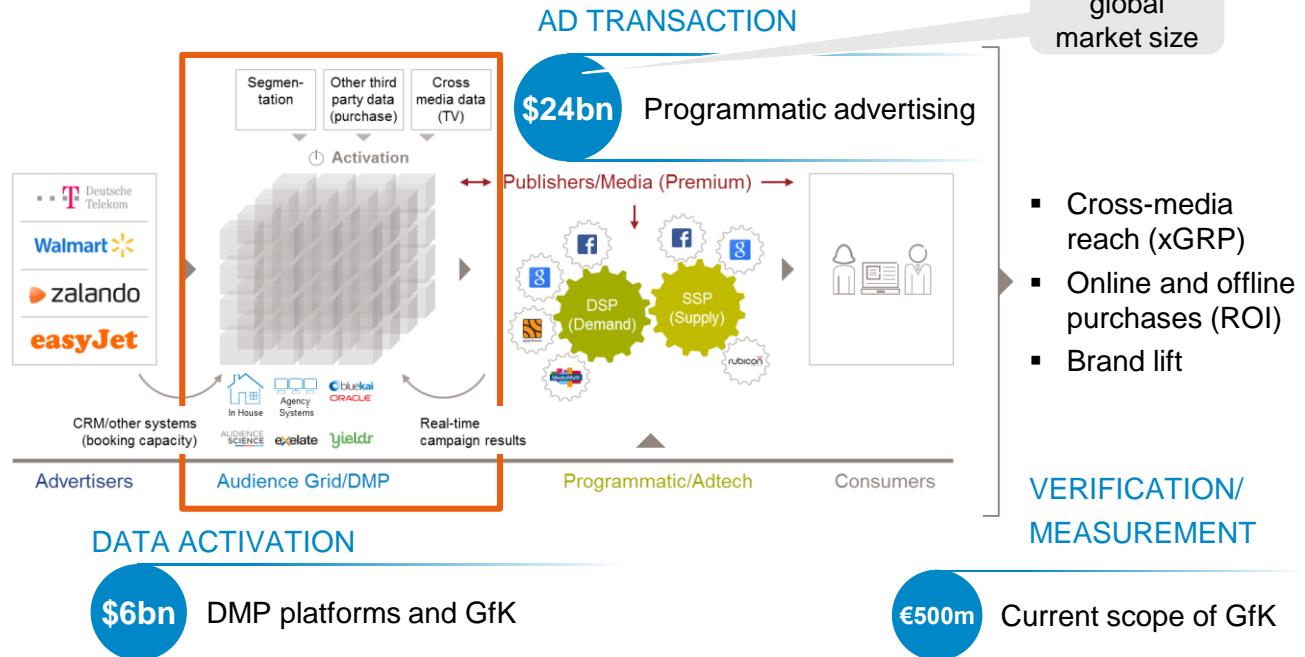
Working with DMPs opens up four different growth opportunities

- ## 1. Standard Audience Segments

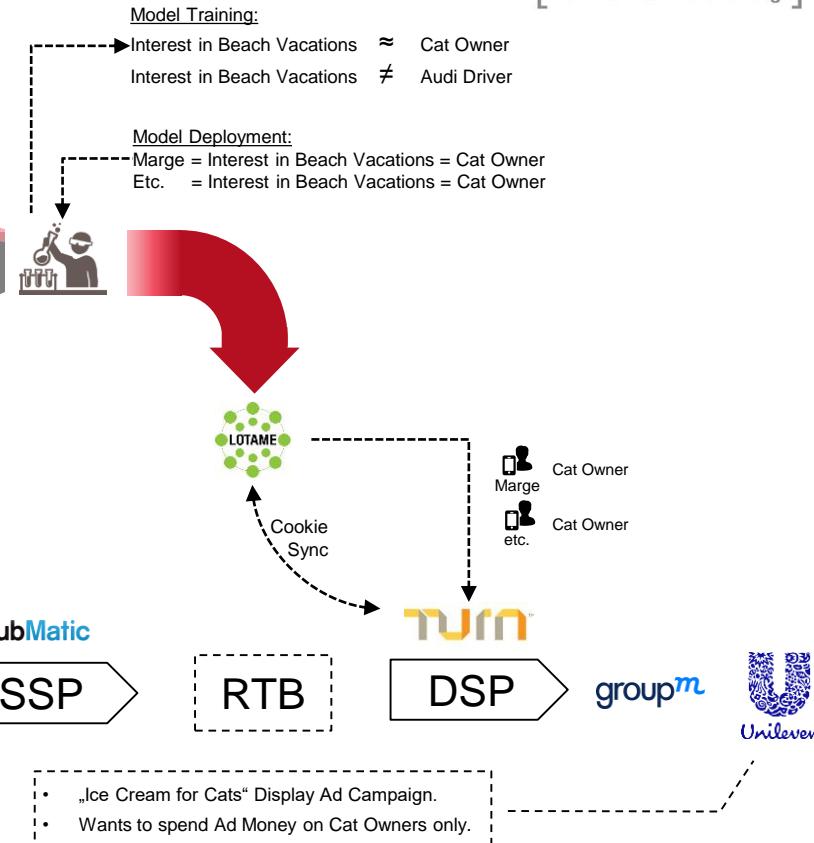
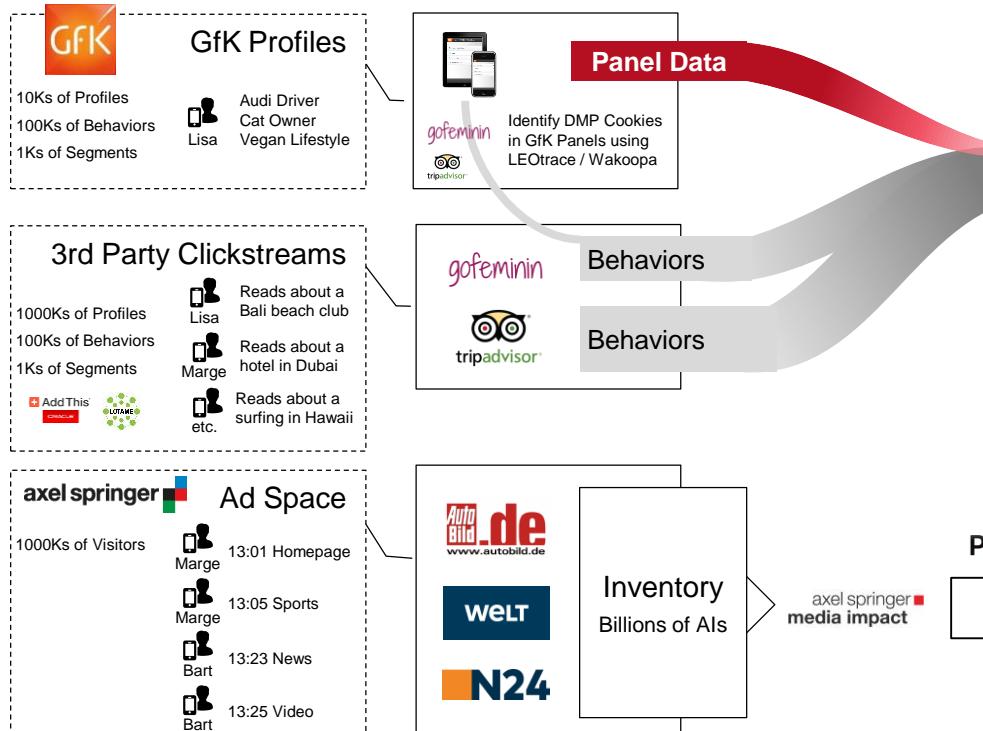
- ## 2. Custom Audience Segments

- ### 3. Personalization, Creative Ad Optimization

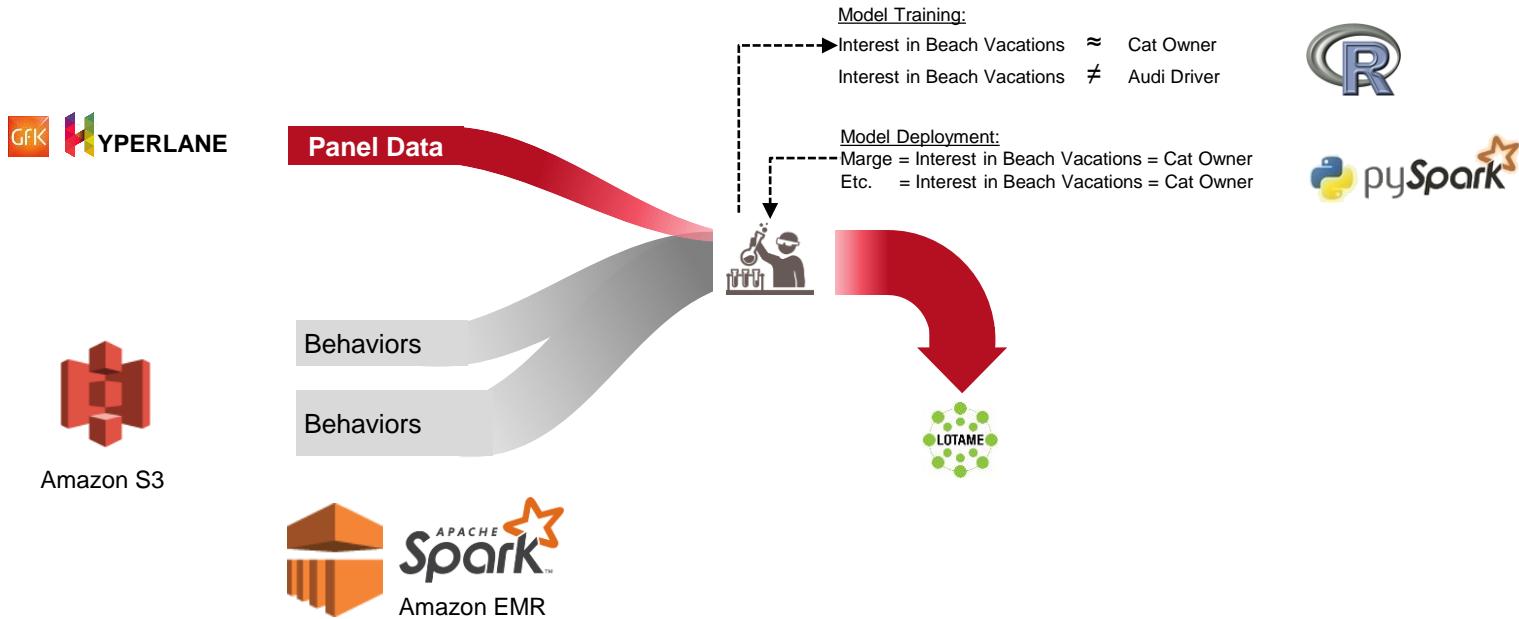
- ## 4. Measurement and Attribution



# Who wants to have Segments?



# Who wants to have Segments?



# Modelling Lotame Data

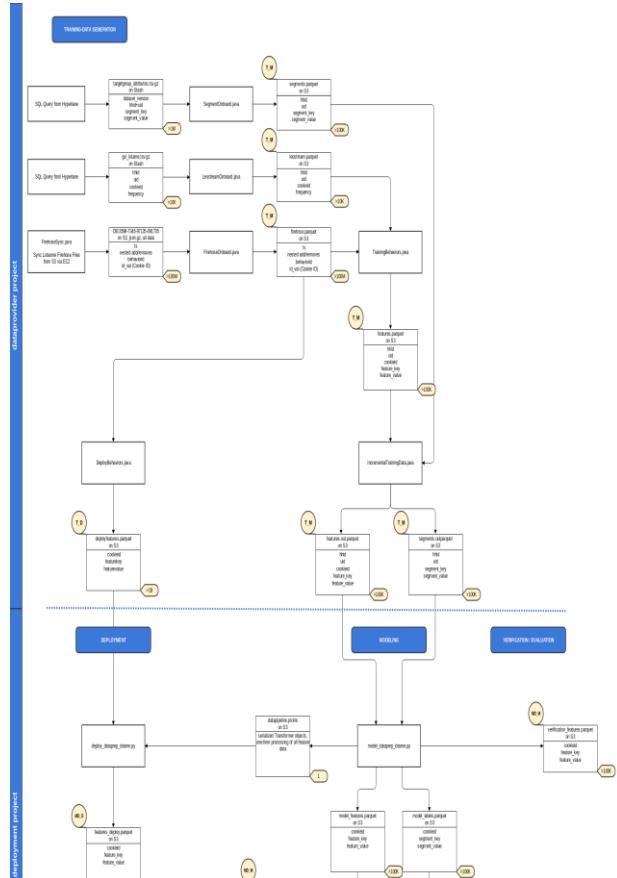
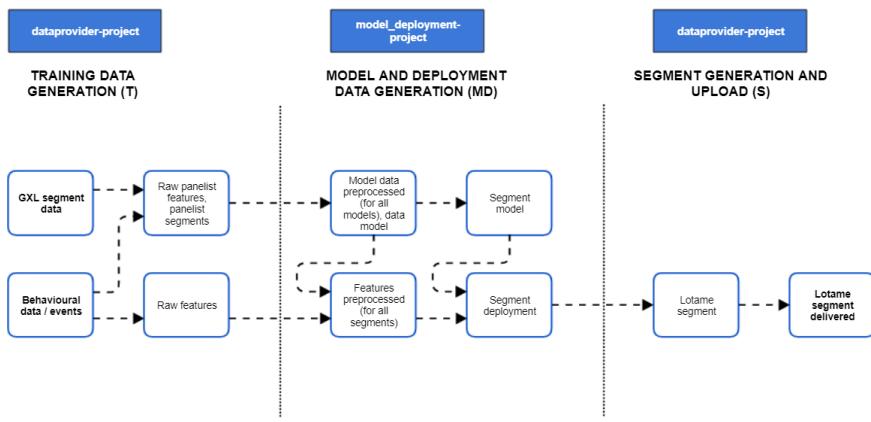
## Preprocessing:

- Aggregating hierarchy nodes (optional)
- Remove duplicated columns
- Delete sparse Cookies / Behaviours
- Add Geomarketing Data

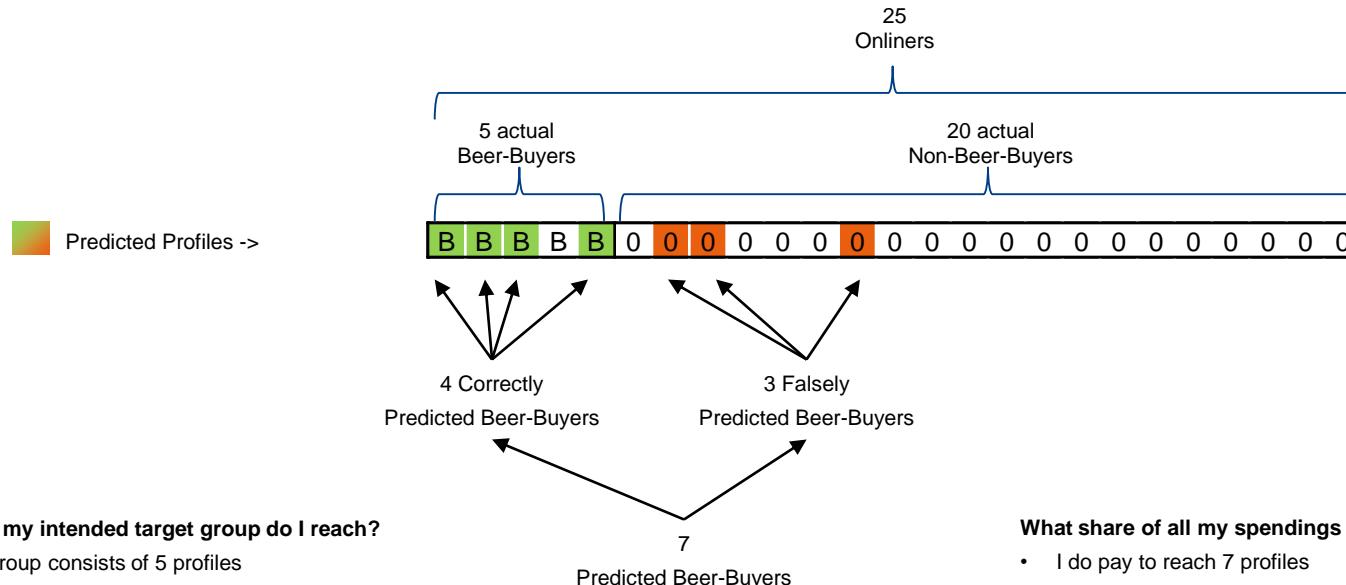


- If a cookie has multiple uid's we solve the multi-user problem the following way (If the behavior differs)
  - Assign the behavior randomly (with observed frequencies)
- Principal Component Analysis
- Aggregation by Correlation
- Sober (developed within GfK)

# Pipelines quickly become complex



# How do we define our Quality?



**What share of my intended target group do I reach?**

- My target group consists of 5 profiles
  - I do reach 4 out of these 5 profiles
  - this means I cover 80% of my target group

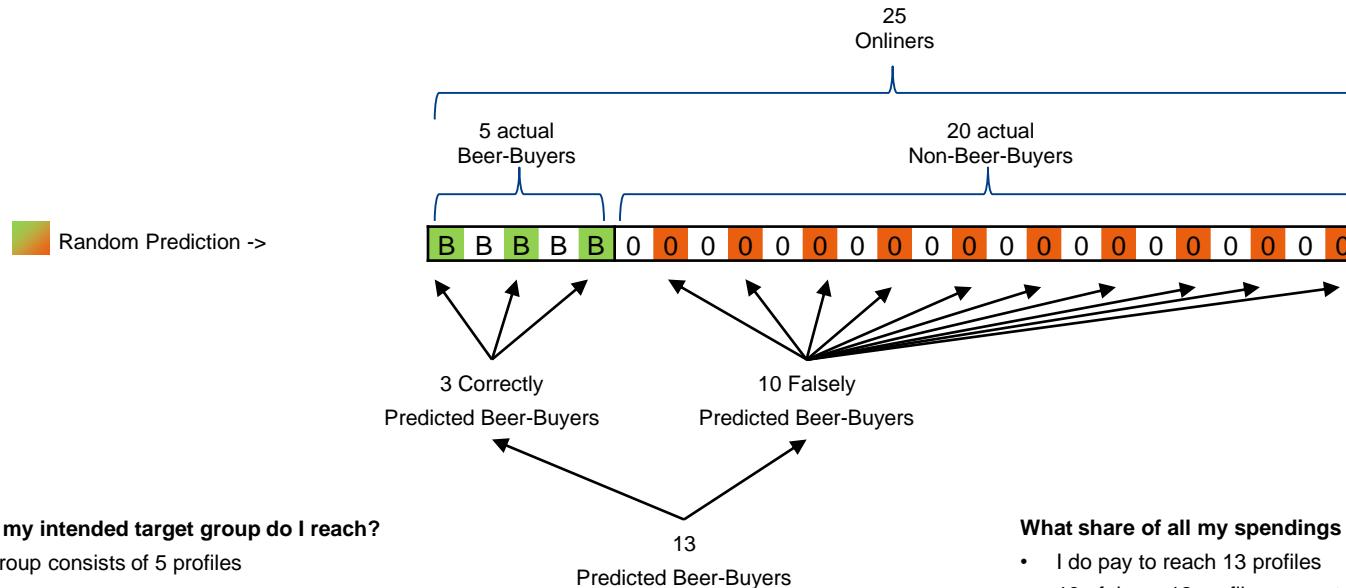
$$Hit\ Rate = \frac{TP}{P} = \frac{True\ Positive}{Actual\ Target\ Group} = \frac{4}{5} = 80\%$$

## What share of all my spendings is wasted?

- I do pay to reach 7 profiles
  - 3 of these 7 profiles are not in my target group
  - this means 43% of my spendings are wasted

$$Waste = \frac{FP}{TP+FP} = \frac{\text{False Positive}}{\text{All Predicted}} = \frac{3}{7} = 43\%$$

# What is good quality? Better than the random baseline:



## What share of my intended target group do I reach?

- My target group consists of 5 profiles
- I do reach 3 out of these 5 profiles
- this means I cover 60% of my target group

## What share of all my spendings is wasted?

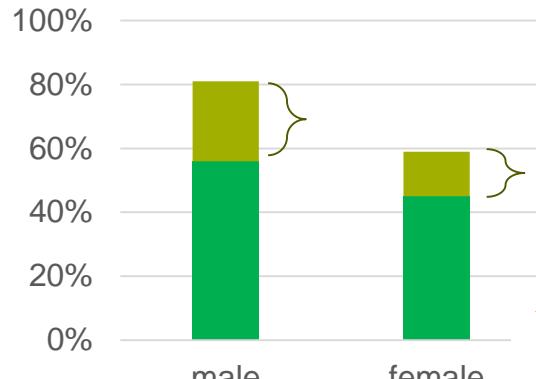
- I do pay to reach 13 profiles
- 10 of these 13 profiles are not in my target group
- this means 77% of my spendings are wasted

$$\text{Hit Rate} = \frac{TP}{P} = \frac{\text{True Positive}}{\text{Actual Target Group}} = \frac{3}{5} = 60\%$$

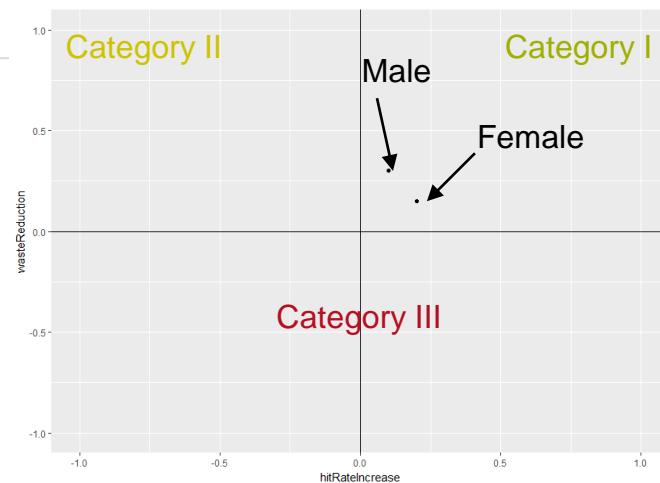
$$\text{Waste} = \frac{FP}{TP+FP} = \frac{\text{False Positive}}{\text{All Predicted}} = \frac{10}{13} = 77\%$$

# How do we visualize these relations in a concise way?

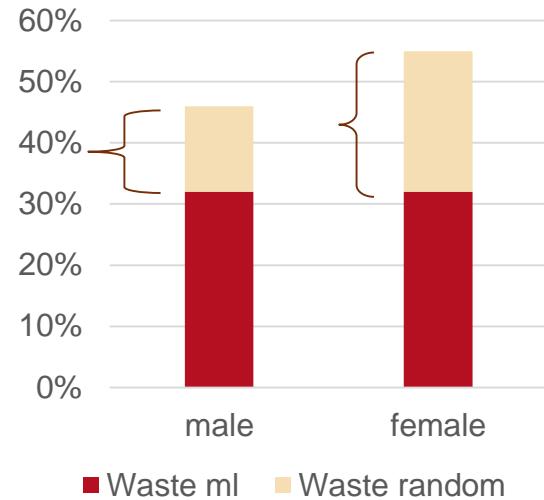
## Hit Rate



- Hit-Rate increase
- Male: 25%
  - Female: 14%



## Waste

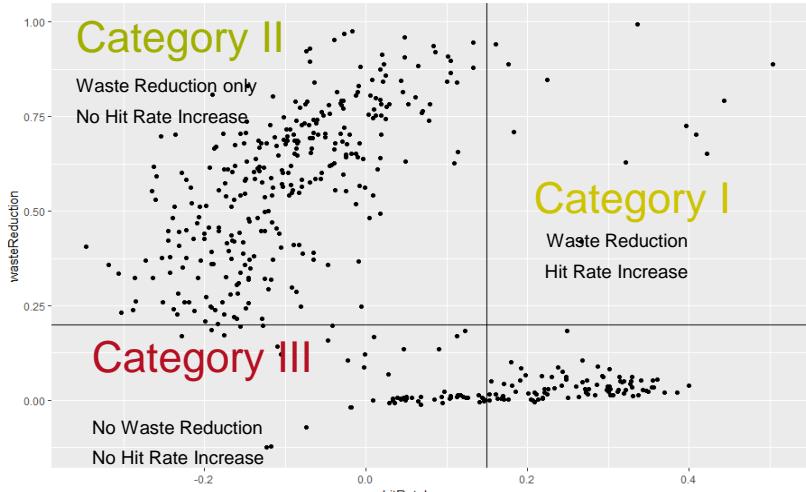


- Waste Reduction
- Male: 14%
  - Female: 23%

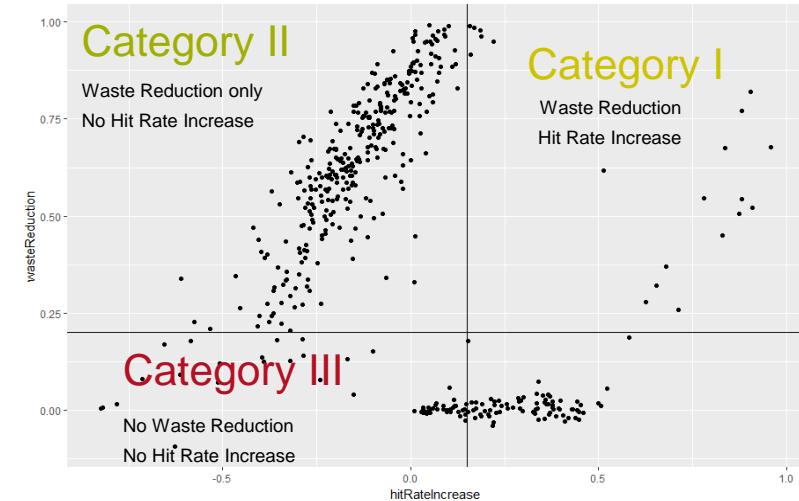
All numbers are illustrative only

# Overview

## Raw Data Input A



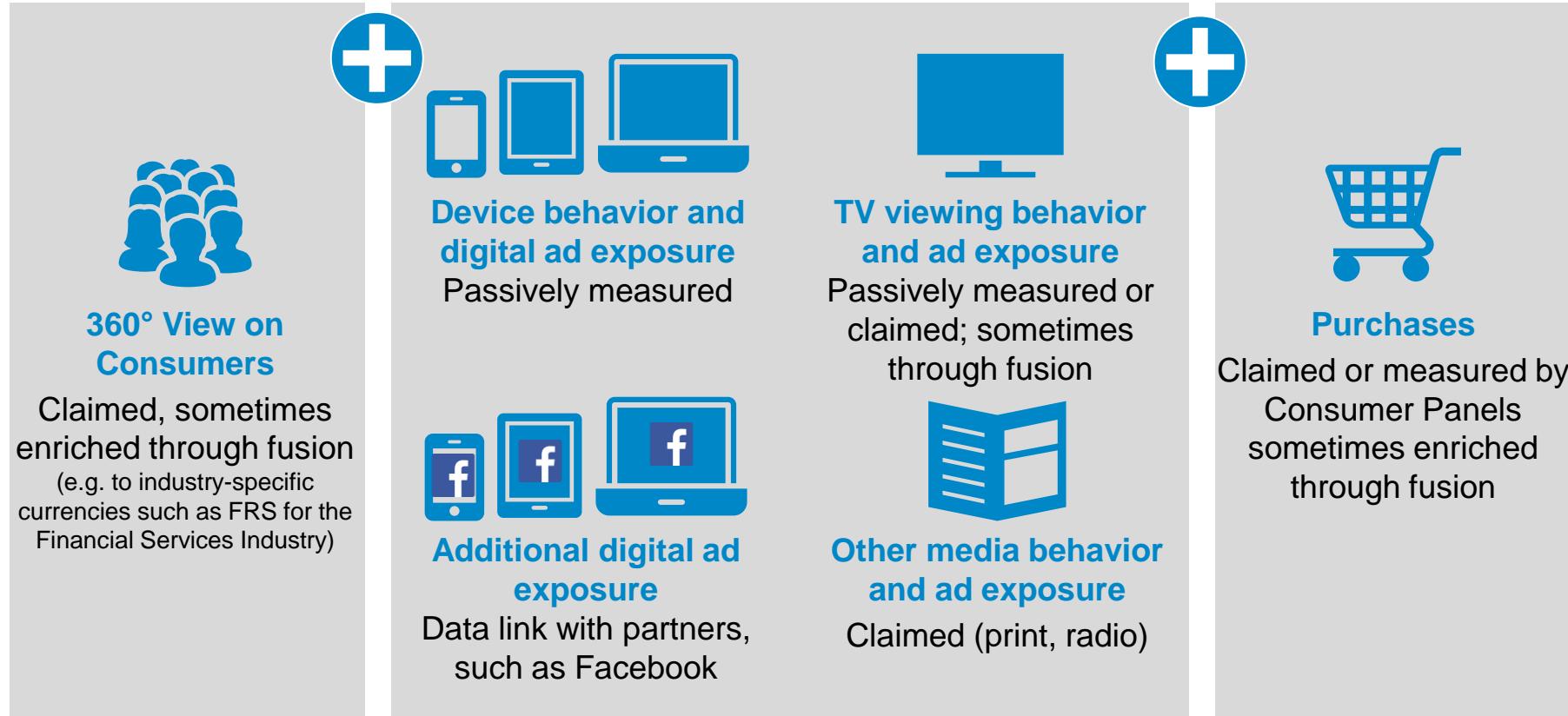
## Raw Data Input B



► Data and modelling quality is better with Input A.

## Example 2: GfK Crossmedia Link / GfK Hyperlane

# GXLL: A platform of single source cross-device digital measurement panels



# GfK's core global data assets

Survey based  
TV exposure      MS: MarketingScan

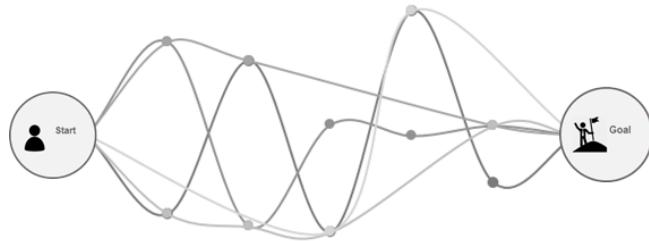


Data type	UK	USA	DE	NL	IT	FR	PL	RU	TU	INDO	BR	SP	ME	ARG	CO
TV exposure and behaviour				GXL	GXL	Sinottica	MS	GXL	GXL	GXL	GXL	GXL	Netquest	Netquest	Netquest
Desktop usage behaviour	GXL	GXL	GXL	GXL	Sinottica		GXL	GXL	GXL		GXL				
Mobile/tablet usage data	GXL	GXL	GXL	GXL	Sinottica		GXL	GXL	GXL	GXL	GXL	Netquest	Netquest	Netquest	Netquest
HH and individual demographic and attitudinal segments	GXL	MRI	GXL	GXL	Sinottica	MS	GXL	GXL	GXL	GXL	GXL	Netquest	Netquest	Netquest	Netquest
HH and individual purchase information				GXL	GXL	Sinottica	MS	GXL	GXL	GXL					
Purchasing power															
Point of sale/retail data	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

One size doesn't fit all. Journeys are different, allowing a large or narrow range of possibilities to influence on consumers in the purchase process

Depending on the category, the Journey may be active for a longer or shorter time period

### Highly Involved



### Autopilot



### Impulse



*Paul is looking for a new Coffee machine, and search internet to learn more about the different types that exist, look at the machines available in his usual hypermarket during his weekly shopping trip and asks friends & colleagues for tips....*

*Martha discovers she's out of beef flavoured bouillon cubes when preparing for dinner, and puts it on the shopping list for next shopping trip. She goes through the store leaflets to see if there's anything on promotion – one never knows – ....*

*It's September, but really hot still. Jack decided to walk to work. "40' walk would do me well" he was thinking. However, he didn't realise it was that hot until he saw the ad for Pellegrino in front of the Relay kiosk. "I definitely need a bottle of water" he said to himself and went into the kiosk...*

# Central Dashboards available for crossdevice website and mobile app usage accross target groups and categories...

## Available Features\* (I)

This screenshot shows the 'Ranking' tab of the Crossmedia Visualizer Beta. It displays a table of data for different brands and products. The columns include Brand, Product, Device, Unique Users, Net Reach, and Duration per User. The data shows, for example, Google having the highest net reach at 91.0% and duration at 02:02:34 h.

BRAND	PRODUCT	DEVICE	UNIQUE USERS	NET REACH	DURATION PER USER
1. Google	All Products (74)	Smartphone	48 469 706	91.0 %	02:02:34 h
2. Facebook	All Products (3)	Smartphone	36 660 867	68.8 %	06:36:27 h
		Tablet	27 590 323	51.8 %	04:47:13 h
		Desktop	19 325 850	36.3 %	04:40:20 h
		Smartphone + Tablet	4 045 591	7.6 %	01:53:32 h
3. YouTube	All Products (6)	Smartphone	35 584 943	63.0 %	02:22:46 h
4. Amazon	All Products (110)	Smartphone	32 434 180	60.9 %	00:49:26 h

This screenshot shows the 'Custom Analysis' tab. It lists 'Selected Objects' such as Opendo (brand), Expedia, Holidaycheck, TripAdvisor, Trivago, Expedia, and Opendo. For each object, it shows Product, Device, Unique Users, Net Reach, and Duration per User. The data indicates varying levels of engagement across different platforms.

SELECTED OBJECTS	PRODUCT	DEVICE	UNIQUE USERS	NET REACH	DURATION PER USER
Opendo (brand), Expedia	All Products	Smartphone	9 105 279	17.1 %	00:12:24 h
Holidaycheck		Smartphone	4 188 249	7.9 %	00:14:09 h
TripAdvisor		Smartphone	3 499 172	6.6 %	00:06:27 h
Trivago		Smartphone	2 460 613	4.6 %	00:03:20 h
Expedia		Smartphone	2 046 001	3.8 %	00:08:24 h
Opendo		Smartphone	1 180 244	2.2 %	00:04:45 h

This screenshot shows the 'Travel Dashboard'. It includes a chart titled 'DEVICE DISTRIBUTION - TRAVEL' comparing Smartphone, Tablet, and Desktop usage. Below the chart is a section for 'CUSTOM ANALYSIS - BASIC KPIs' showing a line graph of KPI Impressions over time from January 2014 to December 2014 for various travel brands.

DEVICE DISTRIBUTION - TRAVEL

CUSTOM ANALYSIS - BASIC KPIs

### Top ranking: websites & apps desktop + smartphone + tablet

- ranking of websites & apps in combination with target group filter
- various KPIs can be evaluated per devices or in total

### websites & apps: Category view & custom analysis

- helps to build up your own relevant universe
- enables flexible combinations & filter for brands and products as well as categories

### Dashboard starting point for analysis & navigation

- central dashboard for snapshot views and navigation
- outlines most important KPIs at a glance
- data will be automatically updated

Additional features will be implemented gradually (until end of 2015), back data available from January 2014, several features on demand.

...which is highly appropriate to continuously track and understand target segments' online and mobile behavior

## Available Features\* (II)

This screenshot shows the 'SELECT TARGET AUDIENCE' interface. It includes a sidebar with pre-defined buyer target groups like 'Chocolate - Chocolate', 'Sports', 'Sausage', etc. The main panel displays 'YOUR SELECTION' with gender: Male, age: 14-19, 20-24, 25-29, 30-34, 35-39, and 'BUYER PROFILES' for 'Travel - Destination Germany: buyer'. At the bottom, there's a summary table for 'Expedia' and 'Opodo'.

This screenshot shows the 'Crossmedia Visualizer Beta' interface under 'User Characteristics'. It displays a dashboard with three main sections: 'Basic KPIs', 'User Characteristics', and 'User Quality'. The 'User Characteristics' section shows audience share for gender (Male: 49.7%, Female: 50.3%) and age groups. Below this is a chart titled 'Frequency of use (usage days)' showing stacked bars for Holidaycheck, Expedia, and Oyota brands across different usage frequency categories.

This screenshot shows the 'Crossmedia Visualizer Beta' interface under 'User Quality'. It displays a dashboard with sections for 'Basic KPIs', 'User Characteristics', 'User Quality', and 'Overlaps'. The 'User Quality' section features a chart titled 'Frequency of use (usage days)' showing stacked bars for Holidaycheck, Expedia, and Oyota brands across different usage frequency categories.

### Target Group Filter

- pre defined buyer target groups & offline media target groups can be selected for all further analysis
- individual target groups on demand

### Demographics: audience profiles

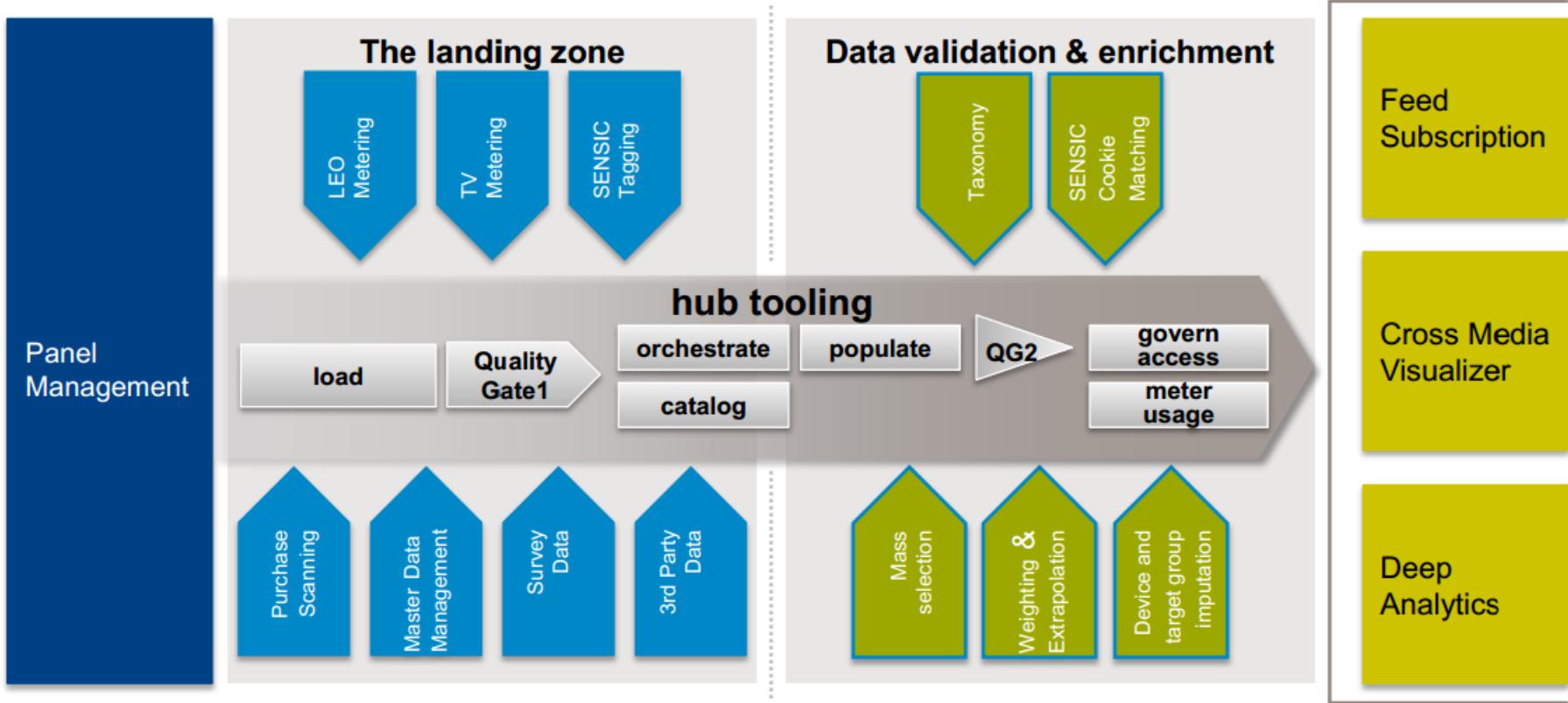
- who are the users of websites & apps
- analyze the data to focus your actions on your target audience

### User Quality

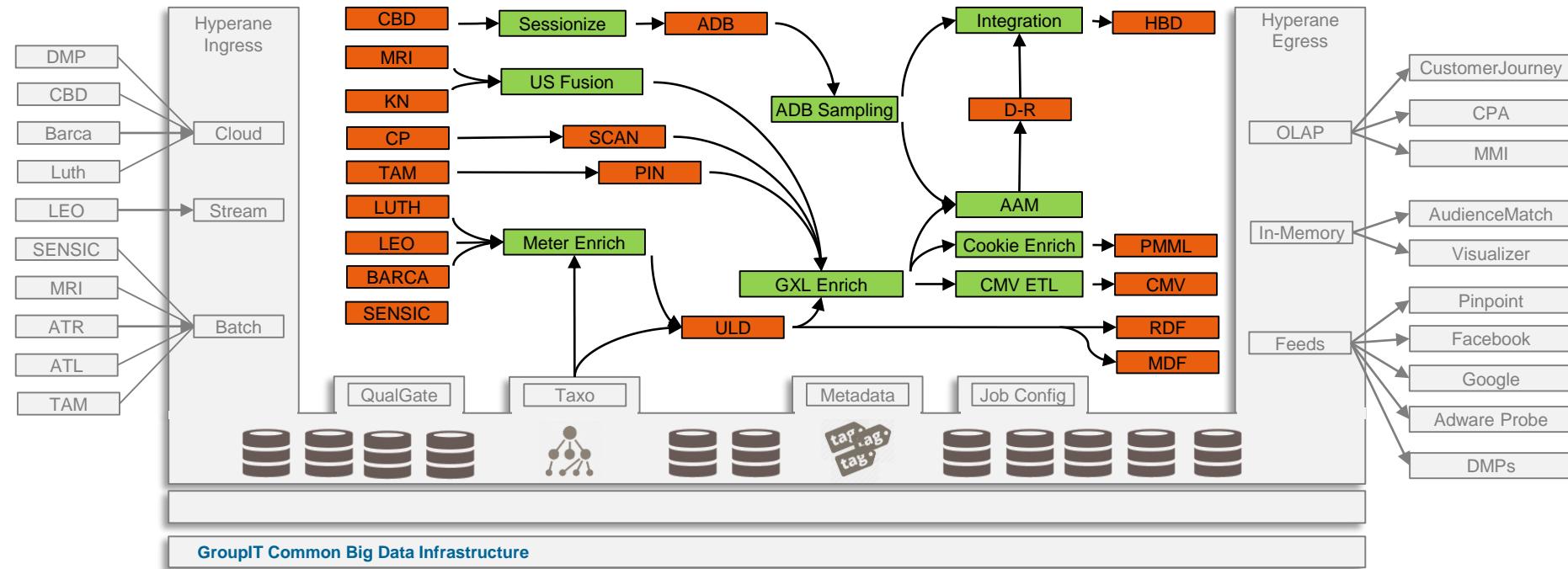
- explore the quality of visitors and their intensity of use
- visit frequency
- reach w/wo bouncers
- lost / recurring / new users
- number of usage days in a month

\* Additional features will be implemented gradually (until end of 2015), back data available from January 2014, several features on demand

# Hyperlane: GXL Business Application



# Hyperlane Apps: Data Assets and Job Flows



# Hyperlane: Broader Context

Sculley et al (2014):

## Undeclared Consumers:

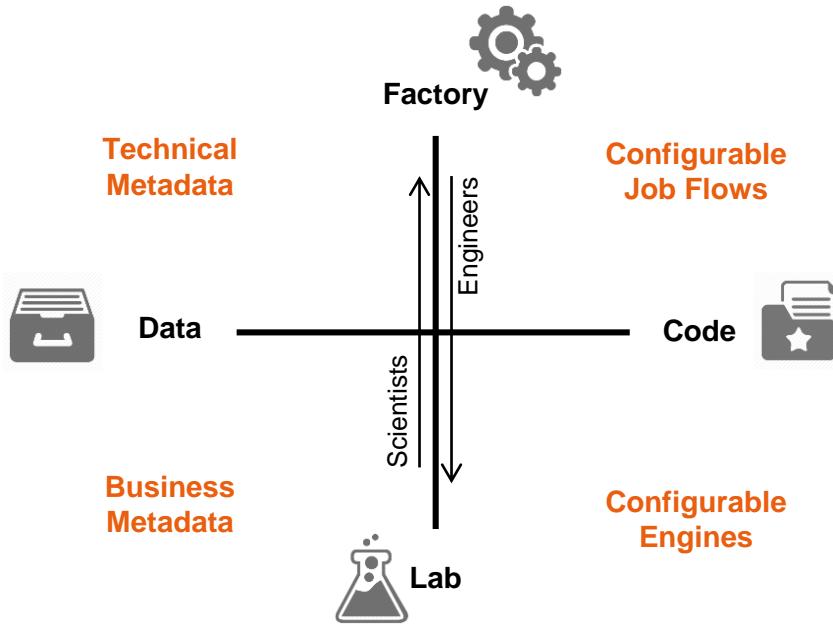
Consuming the output of a given prediction model as an input to another component of the system. Changes will very likely impact these other parts.

## Unstable Data Dependencies:

Some input signals are unstable, meaning that they qualitatively change behavior over time.

## Static Analysis of Data:

On teams with many engineers, or if there are multiple interacting teams, not everyone knows the status of every single feature.



Sculley et al (2014):

## Glue Code:

Using self-contained solutions often results in a glue code system design pattern, in which a massive amount of supporting code is written to get data into and out of general-purpose packages.

## Pipeline Jungles:

The system for preparing data in an ML-friendly format may become a jungle of scrapes, joins, and sampling steps, often with intermediate files output.

## Configuration Debt:

In a mature system which is being actively developed, the number of lines of configuration can far exceed the number of lines of the code.

# Hyperlane: Functional Scope

## Offline Business Process

Find suitable data science methods to deal with peculiarities of the data



**Data Science**

Feature Selection

Model Selection

Model Fitting

## Online Business Process

Continuously load and prepare data so it integrates well with science method



**Data Engineering**

Data Acquisition

Data Preparation

Data Enrichment

Feature Engineering

Model Deployment

Model Automation

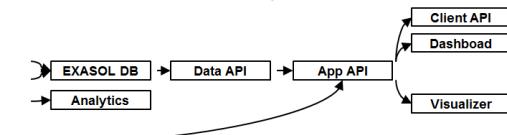
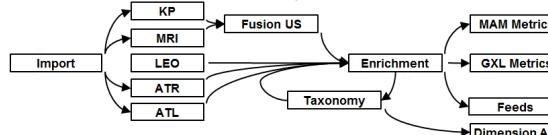
Model Monitoring

## Job Flow Automation

Collaboration between Scientists and Engineers to build scalable data apps that cover Online Business Processes. Configure, schedule and monitor for country setups based on templates.



**Job Flow Ops**



## Data Governance

Ingest cumulative and evolving Data Sets for which the quality needs to be monitored and documented. Business and technical meta data is tagged.



**Data Ops**



## Data Enrichment

Check quality of ingested data sets and approve for automation. Curate Taxonomy and maintain links to data.



**Coding Ops**

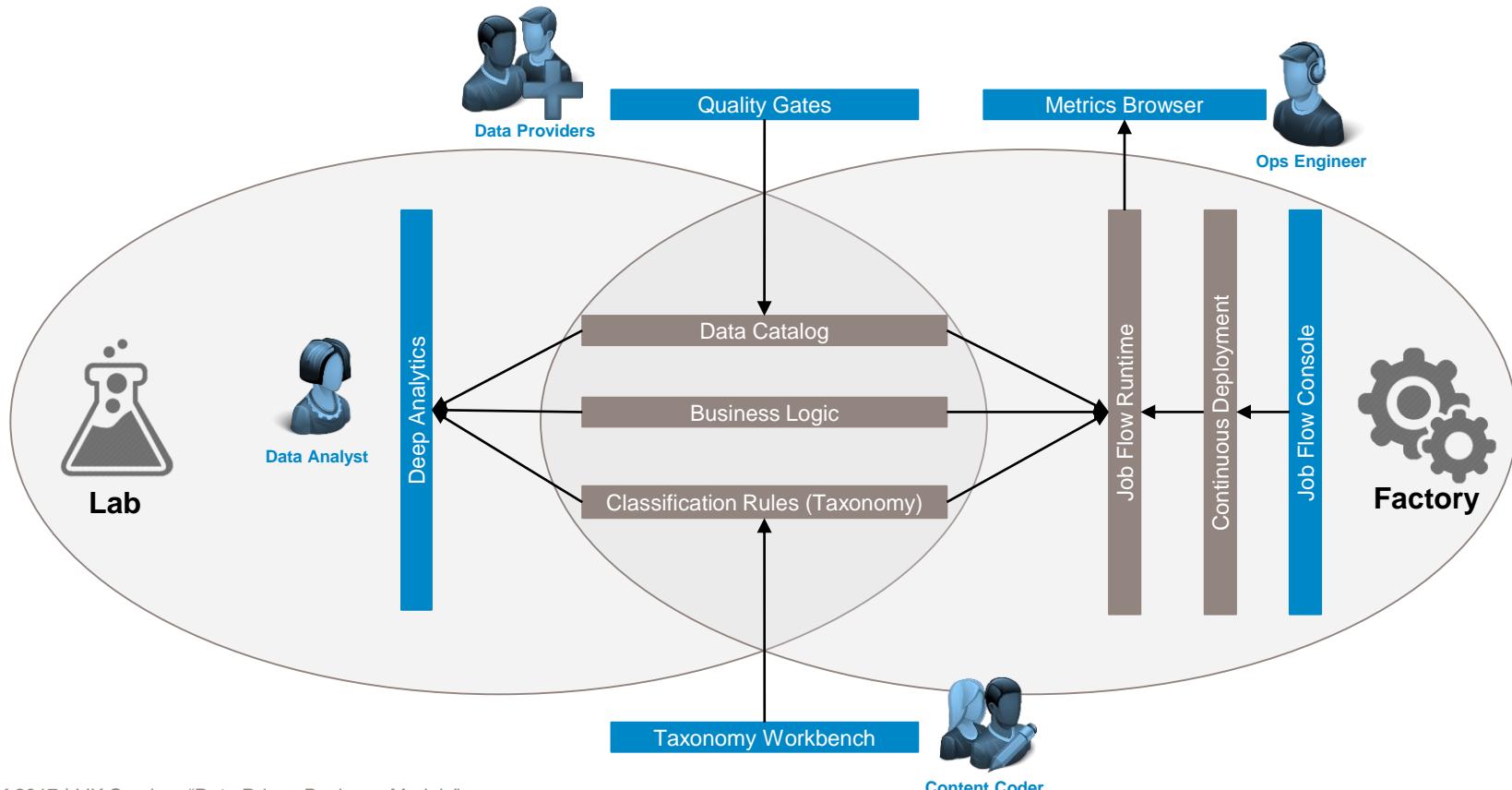


**Quality Gates**

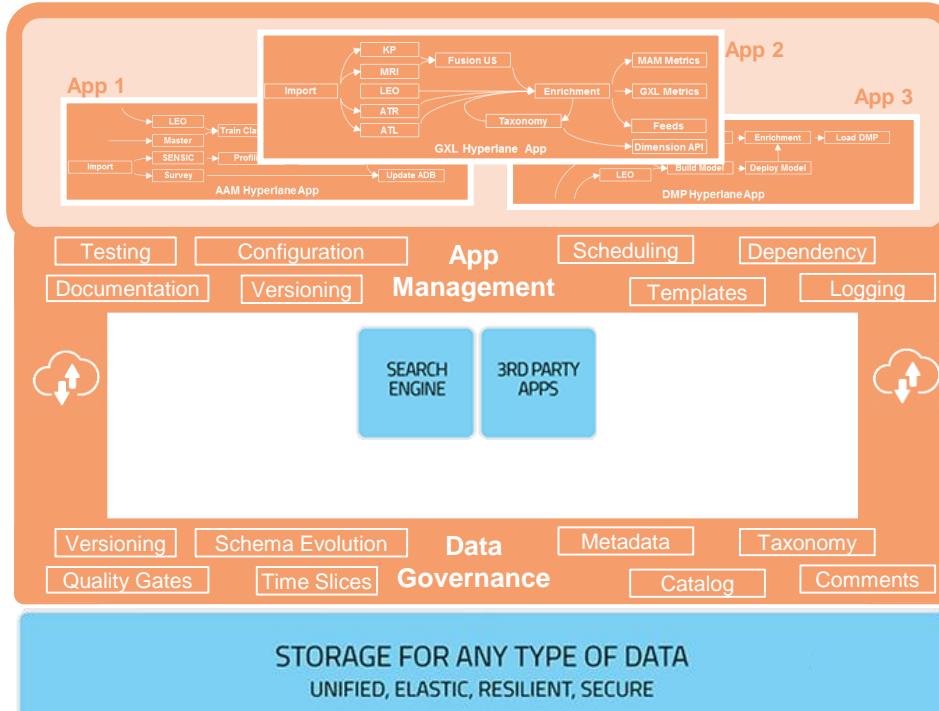


**Taxonomy**

# Hyperlane: Functional Scope



# Hyperlane: Platform vs. Apps



## Hyperlane hosts Data Apps

that combine Data Engineering and Data Science Models into automated Pipelines. Hyperlane Apps are configurable so they can be reused in multiple setups. Data Sets can be selectively fed into Apps based on their meta data. Tooling is available to manage configs, schedule execution and monitor progress.

## Hyperlane manages Data Sets

that can be annotated with technical and Business Meta Data to keep track of lineage and quality. The meta data powers tooling to search, navigate and query arbitrary data sets.

## Hyperlane productionizes Data-Science

capabilities by extending GfK's Common Big Data Technology Stack (Cloudera CDH5) into a platform for GXL and other projects with similar goals. It features an open architecture based on standard technologies.

# Hyperlane: Tooling

## Monitoring



## Analytics

## Taxonomy Workbench

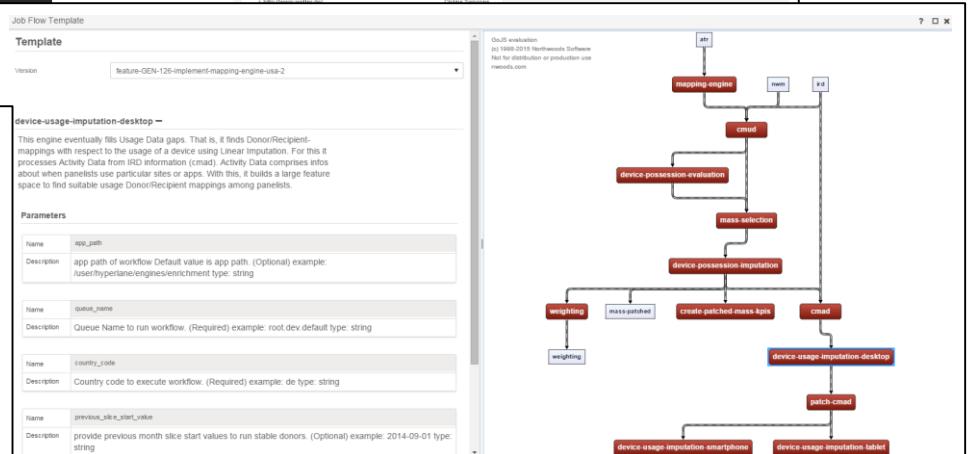
## Data Set Browser

The Data Set Browser interface displays a list of datasets on the left and their visual representations on the right. The datasets listed include:

- hypertane\_pi\_person\_weighting\_pim
- hypertane\_tr\_annotated\_int\_p1d
- hypertane\_tr\_annotated\_int\_with\_commonnames\_p1d
- hypertane\_tr\_atr
- hypertane\_tr\_id
- hypertane\_tr\_imw
- hypertane\_tr\_pasted\_mass\_p1m
- hypertane\_tr\_person\_exposure\_p1m
- hypertane\_tr\_person\_weighting\_p1m

The Taxonomy Manager interface for GMX Mail includes the following sections:

- Taxonomy Bundles:** Shows a list of bundles (Bundle 1, Bundle 2, Bundle 3, Bundle 4) and a search bar.
- Rule Set Timeboxes:** Displays timeboxes for GMX Webmail and GMX Webmail 2.
- GMX Webmail:** Shows classification rules for URLs and iOS/Android/UA patterns.
- Change Log:** A history of changes made to the taxonomy.



Job Flow Configuration Management

# Hyperlane Tooling: Taxonomy Workbench

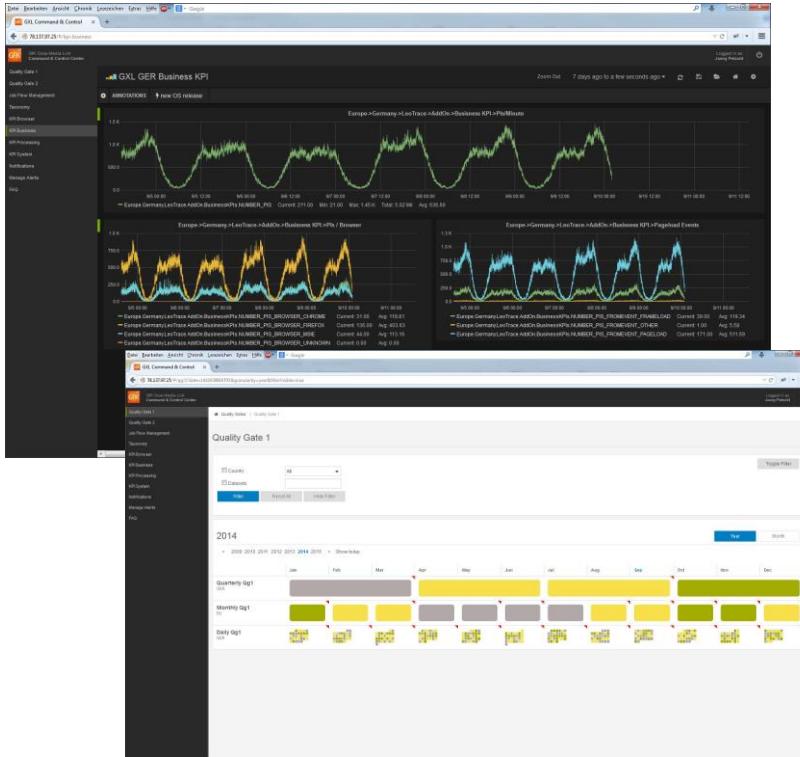
The screenshot shows the GMX Mail Taxonomy Manager interface. On the left, there's a sidebar with navigation links like 'All Bundles', 'Bundle 1', 'Bundle 2', 'Bundle 3', 'Bundle 4' (selected), 'Enter name...', 'Trunk', 'Company 1', 'Company 2', 'Brand 1', 'Brand 2', and 'Product 1'. The main area has two tabs: 'GMX Webmail' and 'GMX Webmail 2'. Each tab contains a 'Valid Time Manager' section with a 'Valid Time' dropdown set to 'Unbound - 8th March 14' and a 'Changed' date of '10th July 14'. Below this are 'Rule Set Timeboxes' and 'Classification' sections. The 'Classification' section lists categories: 'News / Information', 'Online Services', and 'Homecare / Cleaning'. Under 'Rule Set Timeboxes', there are sections for 'URL Patterns' (with four entries: http://www.wetter.de/, http://www.wetter.de/, http://www.wetter.de/, http://www.wetter.de/), 'iOS Patterns' (0), 'Android Patterns' (0), and 'UA Patterns' (0). To the right of the tabs is a 'Change Log' table:

Date	User	Action	Category
20. jan. 13	Thomas S.	New Pattern	GMX Webmail
20. jan. 13	Peter J.	Renamed X	GMX Webmail
20. jan. 13	Thomas S.	New Pattern	GMX Webmail
20. jan. 13	Peter J.	Renamed X	GMX Webmail
20. jan. 13	Thomas S.	New Pattern	GMX Webmail
20. jan. 13	Peter J. Thomas S.	Renamed X	GMX Webmail
20. jan. 13	Peter J.	Changed tests sda sd asdadasdasd asdsssdsssd a sda asd a sda sda sda sda	GMX Webmail

A tool that supports our Hub-based Coding Team to efficiently manage the A2C taxonomy.

- Specify URL Patterns and Mobile App IDs that comprise Products& Brands
- Map A2C Categories to Products
- Manage existing Pattern Sets as well as emerging URLs and App IDs
- Support of Custom Patterns for specific clients (Patching the Default Taxonomy)

# Hyperlane Tooling: Monitoring



A centralized Monitoring System that provides transparency about data collection and data production processes.

- Visualization of events as time series
- Browser to manage hundreds of metrics
- Compare and Overlay different metrics
- Calculate Delta to time periods in the past
- Define Thresholds to receive alerts
- Setup customized Dashboards

# Hyperlane Tooling: Deep Data Analytics

The screenshot shows the Hyperlane Query Editor interface. At the top, there are tabs for 'Hive Editor' and 'Query Editor'. Below the tabs, a 'Navigator' pane lists various database tables and their descriptions. The main area displays a SQL query titled 'Sample: Salary growth' with a timestamp 'Salary growth (selected) from 2007-08'. The query uses JOIN and ORDER BY clauses to calculate salary growth. Below the query is a 'Results' section containing a bar chart. The chart has 'description' on the X-axis and 'salary' on the Y-axis, ranging from 0 to 20000. The bars represent salary values for different job descriptions like 'Dentist, all dentists', 'Nurse', 'Physician/Physician's assistants', etc.

This screenshot of the Hyperlane Query Editor interface shows the same 'Sample: Salary growth' query and results as the previous one. However, the results section now displays a table instead of a chart. The table has columns for 'description', 'salary', and '\_c3'. It lists 23 rows of data corresponding to the job descriptions and their calculated salary growth values.

A web-based analytics console that allows analysts to answer advanced questions by browse, query, aggregate and visualize data.

- Interactively browse existing data sets
- Run queries using well-known SQL
- Maintain a repository of Query Templates
- Merge data sets and apply custom weights
- Replace parts of the existing analytics tool chain and provide an environment for collaboration.

# Wrap Up