



# Architecting a Data Platform For Enterprise Use

May 2018

Mark Madsen  
Todd Walter

# The business complaints about data and analytics

You really should fire IT.



**IT root causes**

**IT proximate causes**

**What is said in disputes**

1980s-era methods

Inappropriate technology

Data hygiene fetishes

Vendor lock

1990s-era procurement

IT skills deficit

Dysfunctional OLTP portfolio

Lack of agility

People & vendor cost basis

Client-server infrastructure

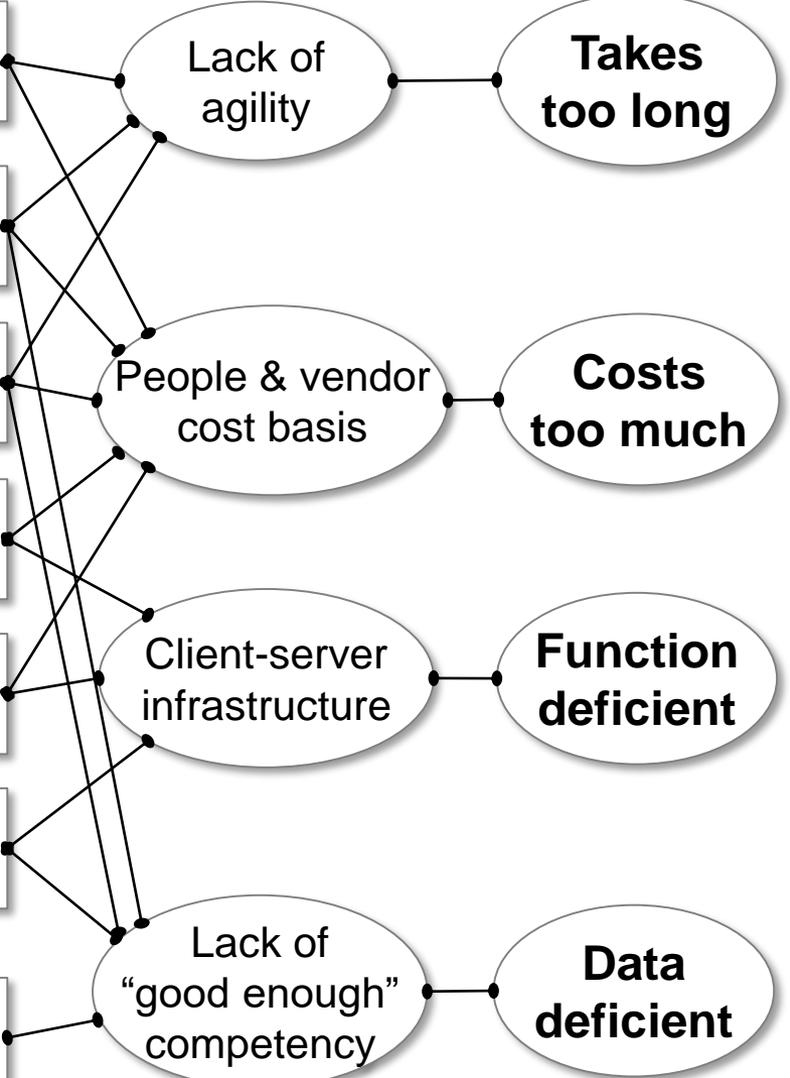
Lack of "good enough" competency

**Takes too long**

**Costs too much**

**Function deficient**

**Data deficient**



**IT root causes**

**IT proximate causes**

**What is said in disputes**

**How business responds**

**How IT responds**

1980s-era methods

Inappropriate technology

Data hygiene fetishes

Vendor lock

1990s-era procurement

IT skills deficit

Dysfunctional OLTP portfolio

Lack of agility

People & vendor cost basis

Client-server infrastructure

Lack of "good enough" competency

**Takes too long**

**Costs too much**

**Function deficient**

**Data deficient**

SaaS / Cloud BI

Consultants

Shadow BI

Self-service BI

Workarounds & spreadmarts

Loss of control

Loss of visibility

Loss of knowledge

Security risks

Hidden costs

Bad data risks

**IT root causes**

**IT proximate causes**

**What is said in disputes**

**How business responds**

**How IT responds**

1980s-era methods

Inappropriate technology

Data hygiene fetishes

Vendor lock

1990s-era procurement

IT skills deficit

Dysfunctional OLTP portfolio

Lack of agility

Cost basis too basic

Client-server infrastructure

Lack of "good enough" competency

**Takes too long**

Too much

**Function deficient**

**Data deficient**

SaaS / Cloud BI

Shadow BI

Self-service BI

Workarounds & spreadmarts

Loss of control

Loss of visibility

Loss of knowledge

Security risks

Hidden costs

Bad data risks

**This is not a technology problem – it is an architecture problem.**

# What do we hear in the field?

"I want to do analytics, beyond what I can do with BI tools"

"I need an analytics strategy"

"I need an analytics roadmap"

"I have a specific analytics project of type..."

"We have a data lake. What should we do with it?"

"Technology Y replaces technology X"

"I want to modernize my DW" aka "I want to speed up my DW"

"Don't need schema, curation, ETL, governance ... anymore"

*Good judgement is the result of experience.*

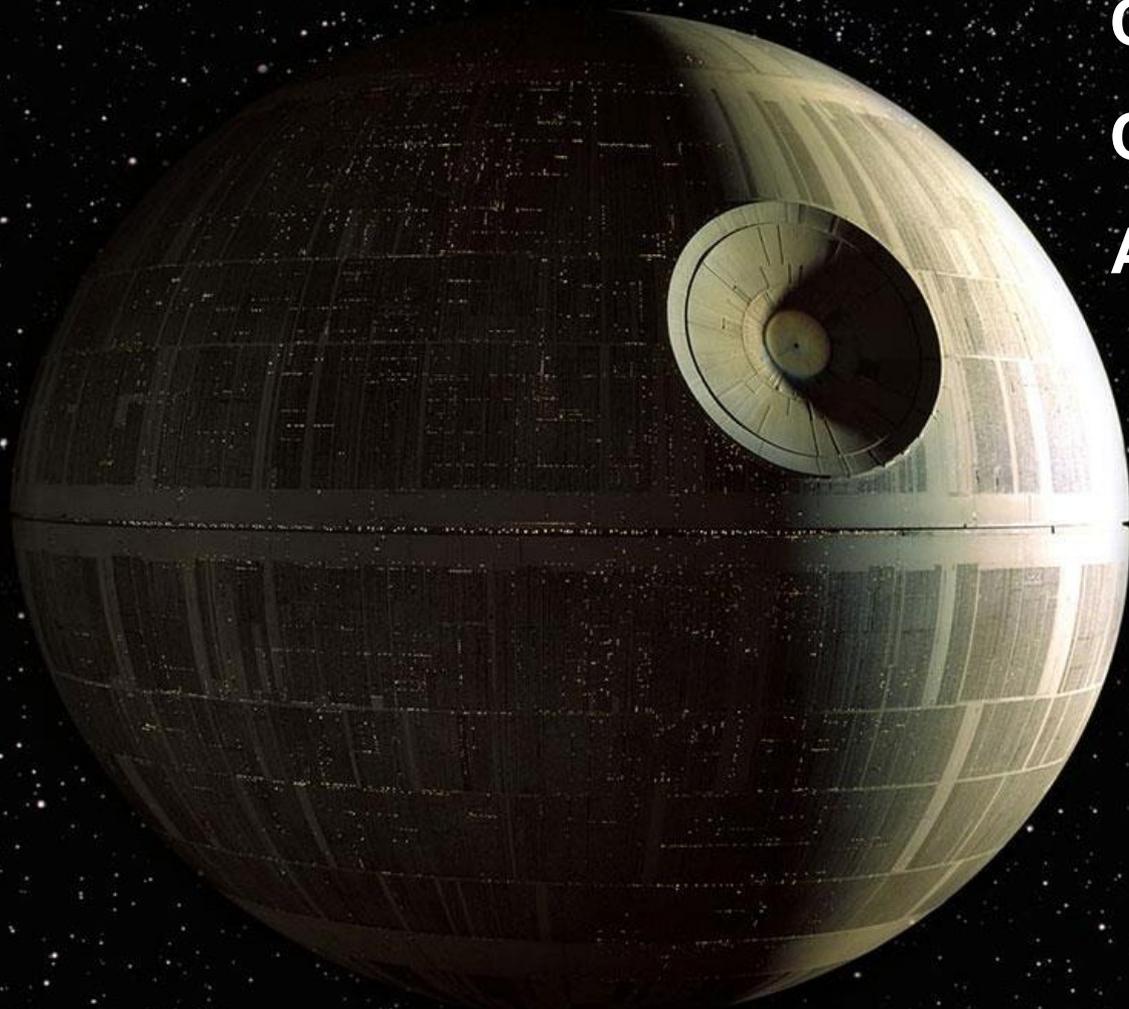
*Experience is the result of bad judgement.*

*—Fred Brooks*

**“Build a data lake and put all the data there first”**



# DW: Centralize, that solves all problems!



**Creates bottlenecks**  
**Causes scale problems**  
**Availability?**

# The data lake solution: no central authority

A satellite in space is shown looking at a large, bright, glowing data lake. The data lake is composed of many small, bright particles, creating a dense, glowing cloud. The satellite is a small, green and white object in the lower right corner. A speech bubble points to the satellite with the text "wtf, it was fully operational!".

wtf, it was fully operational!

# The data lake solution?



There's a problem: as the lake is envisioned, it is still a centralized data architecture, but this time there is no single global model. Instead it's files and not modeled. It can be operational while under construction.

*It's still a death star.*

# Eventually we run into the same problems



Seriously, wtf?  
It was agile  
and operational

Rising complexity and scale break centralized models

**“Standardize on one tech, it’s simpler for everyone”**

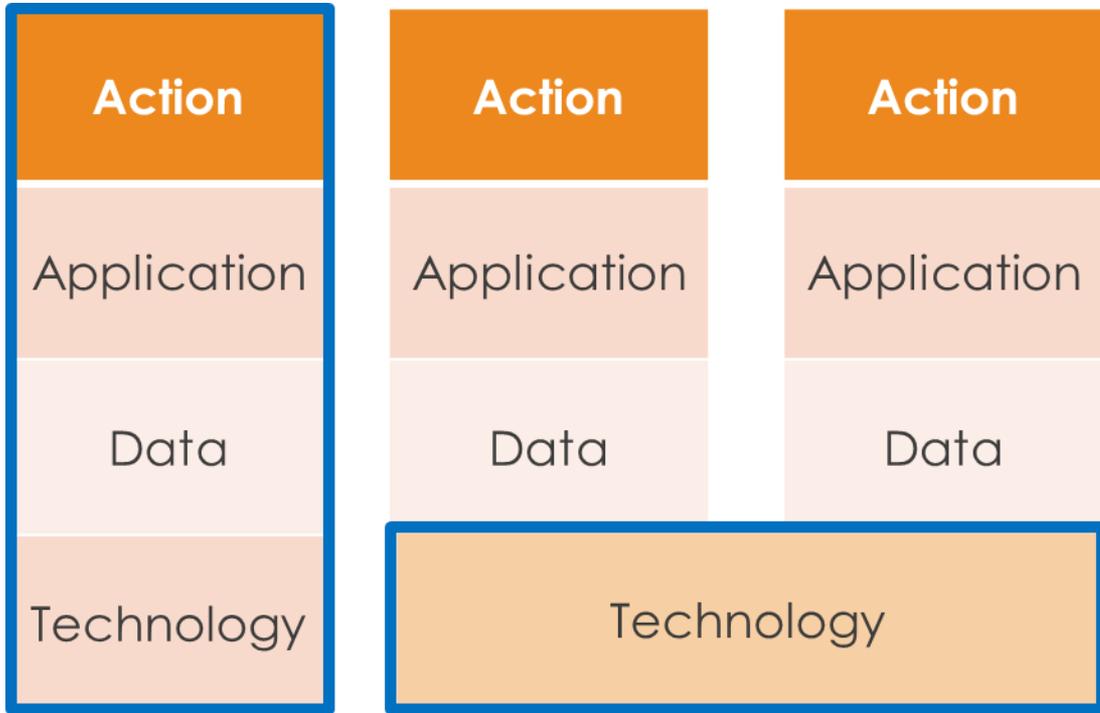


***It's called stack think. Pick your vendor. Cede all architecture.***

# Agile your way into the Analytics Shantytown

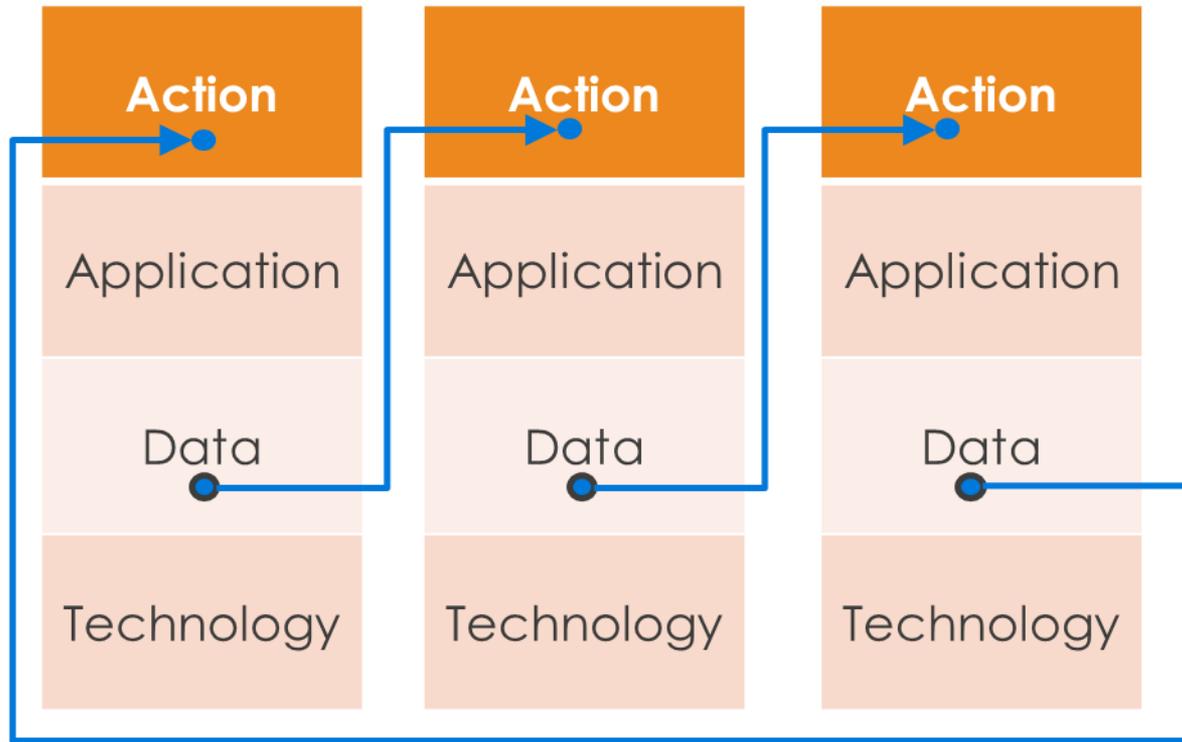


# “Infrastructure one PoC at a time”



- Use case driven
- Leverage (any) new technology
- Re-use of the technology stack
- Data is captive to the application or to the technology stack
- Developers tend to think “function first”
- Analytics people also think “function first”, but generally require integrated data

# A key aspect of operational analytics is “end-to-end”



- Most data applications are organized around a process, not a task or function.
- Real-time analytics systems pass their data over the network, but the dependencies can cross application boundaries, as well as requiring persistence.
- Developers are being forced to think of data outside the local context.

E2E does not allow you to take a random walk through technologies and BYOT because the complexity of integration leads to a mountain of technical debt.

**“Start with the platform. The rest will follow.”**





**Big data promise: What you see**



**Big data reality: What users see**

**Persisting data  
is not the end  
of the line.**

**If you stop  
here you win  
the battle and  
lose the war**



# We don't have an analytics problem, just like we didn't have a BI problem

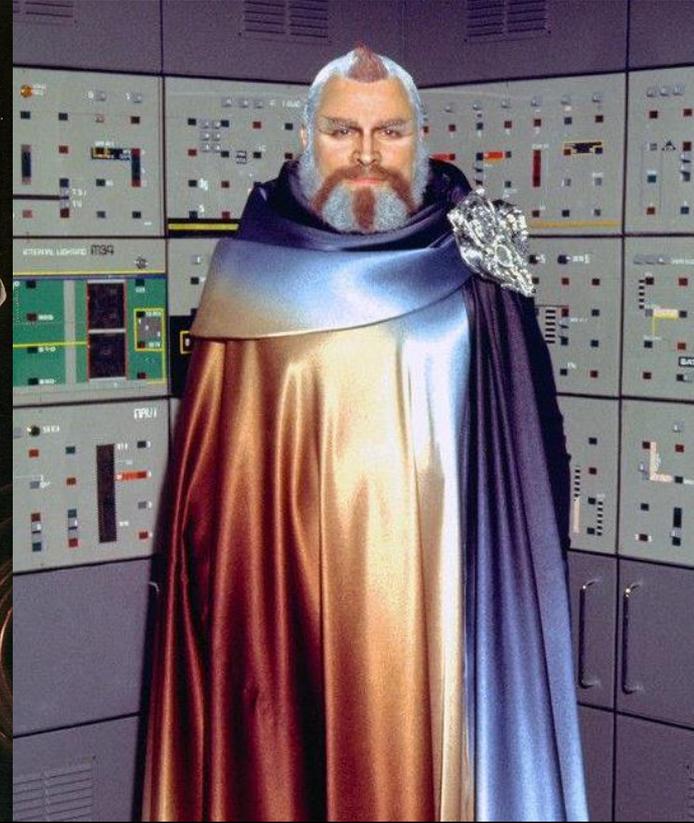
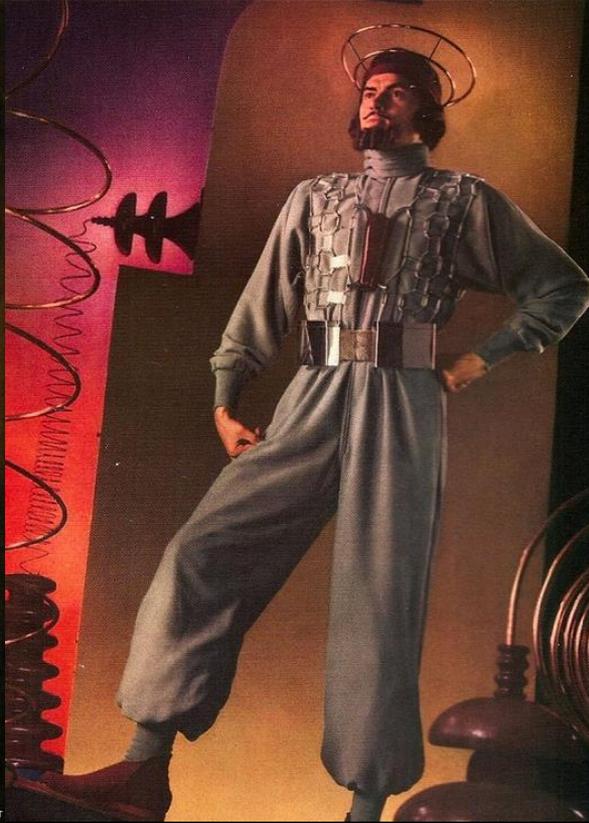
The origin of analytics as “business intelligence” was stated well in 1958:

“...the ability to apprehend the interrelationships of presented facts in such a way as to guide action towards a desired goal. ” ~ *H. P. Luhn*

Our goal is analytics as a **capability**, not a technology

“A Business Intelligence System”, <http://altaplana.com/ibmrd0204H.pdf>

# Three constituencies



Stakeholder  
aka the recipient

Analyst  
aka the data scientist

Builder  
aka the engineer

# Starting points for analytics strategy



Many organizations choose to start with the analysts. Create a data science team. Turn them loose to find a problem.



Many more start with builders: technology solutions looking for problems, e.g. 65% of the IT driven Hadoop and Spark projects over the last five years.

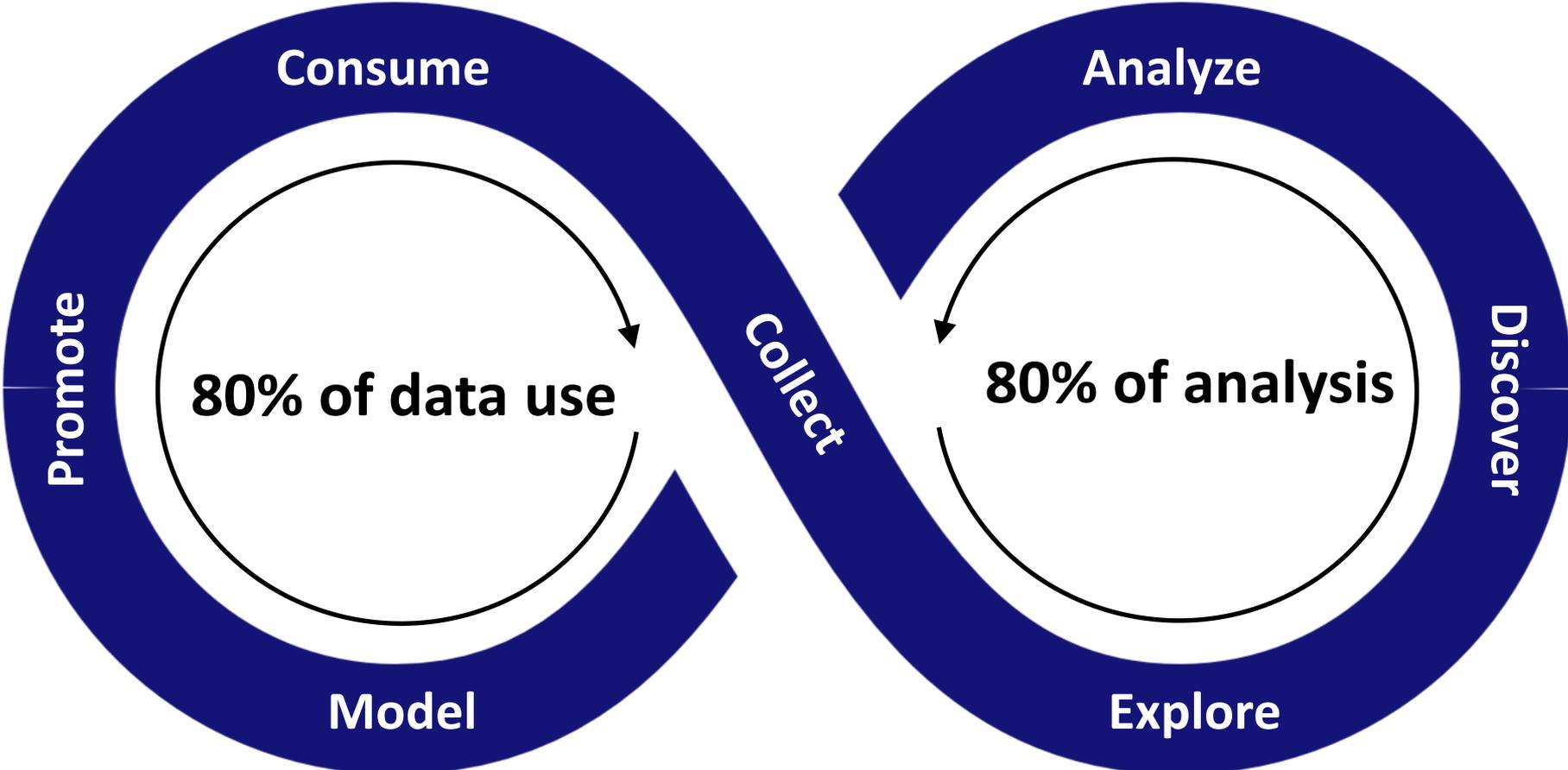


The right place to start? Stakeholders. The goal to achieve, the problem to solve.

Each constituency has their own set of problems to deal with

# **NATURE OF THE PROBLEM FROM THE STAKEHOLDER'S PERSPECTIVE**

# BI and Analysis: Repeatability and Discoverability



## Focus on repeatability

Application cycle time  
90% of users

## Focus on discoverability

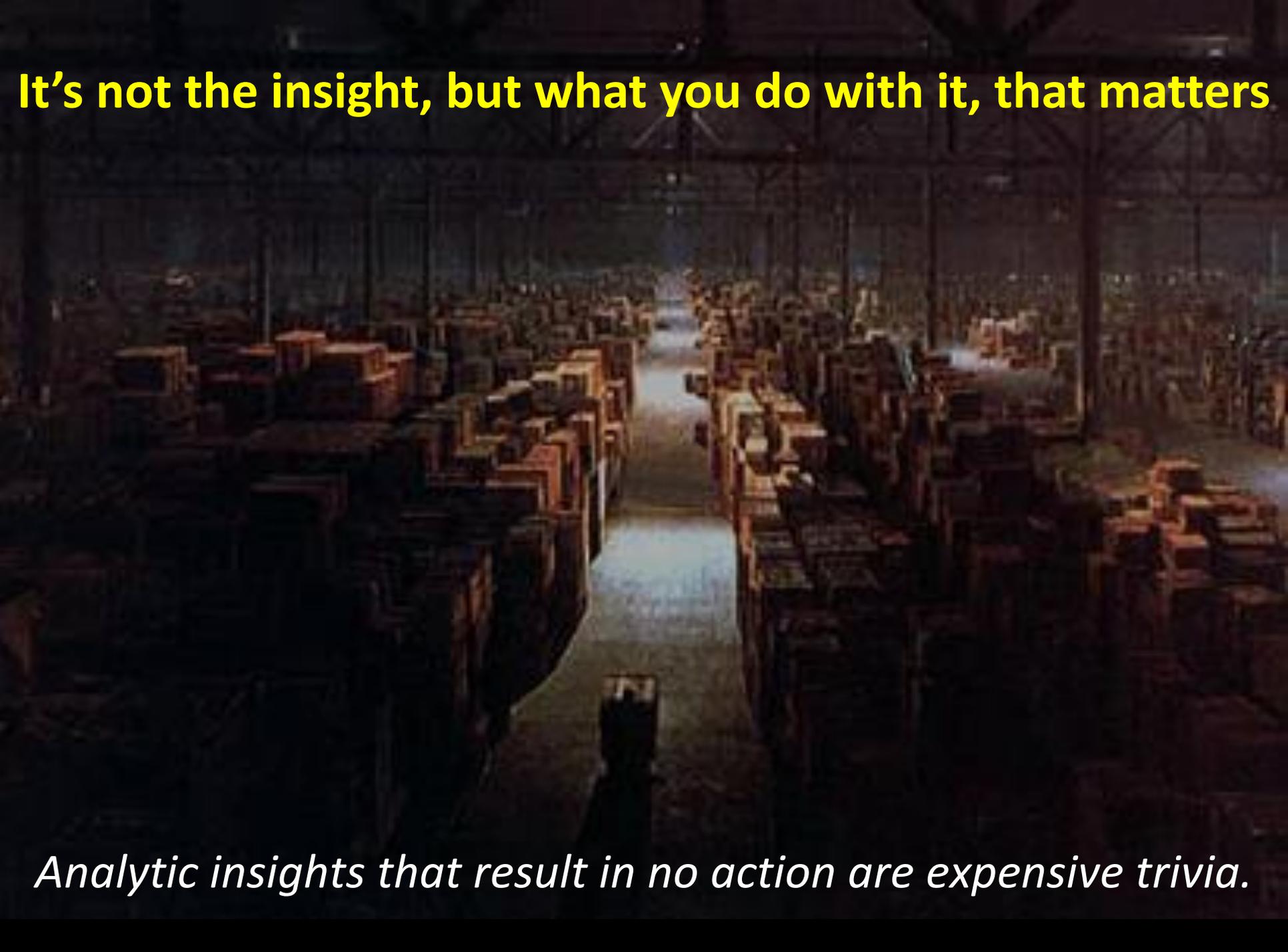
Analyst cycle time  
9% analysts, 1% data scientists

# The myth that still drives analytics – analytic gold

All we need is a fat pipe and pans working in parallel...



**It's not the insight, but what you do with it, that matters**

A large industrial warehouse filled with rows of stacked wooden pallets, viewed from a central aisle looking down the length of the building. The lighting is dim, with a bright light source at the far end of the aisle, creating a strong perspective and highlighting the repetitive nature of the stacks.

*Analytic insights that result in no action are expensive trivia.*

# Applying analytics is not an analytics problem



Applying analytics is not in the analyst's control.

It's not in the engineer's control.

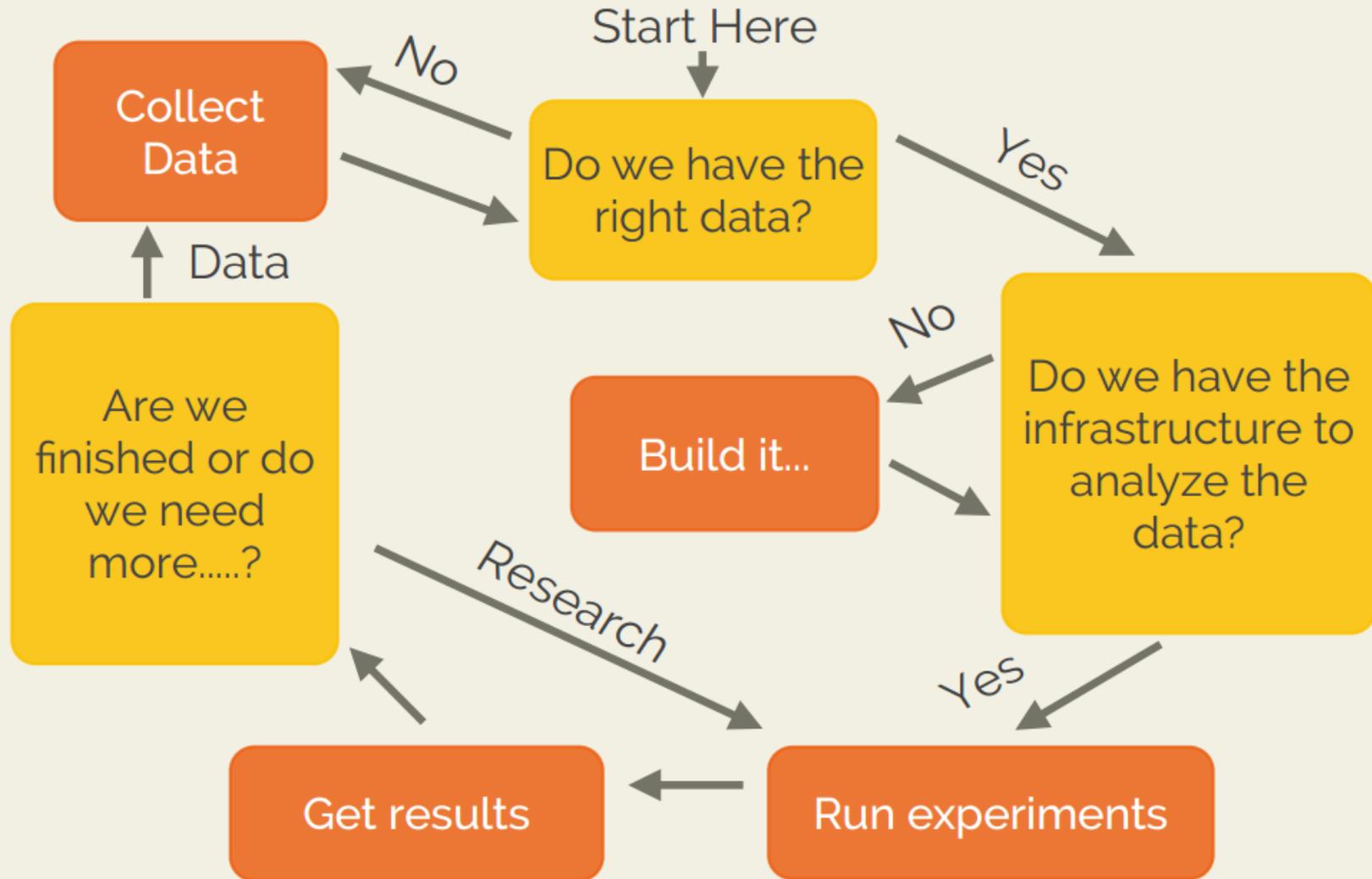
It's in the control of the people involved in the process.

**Failures are often in execution, not in analytics development.**

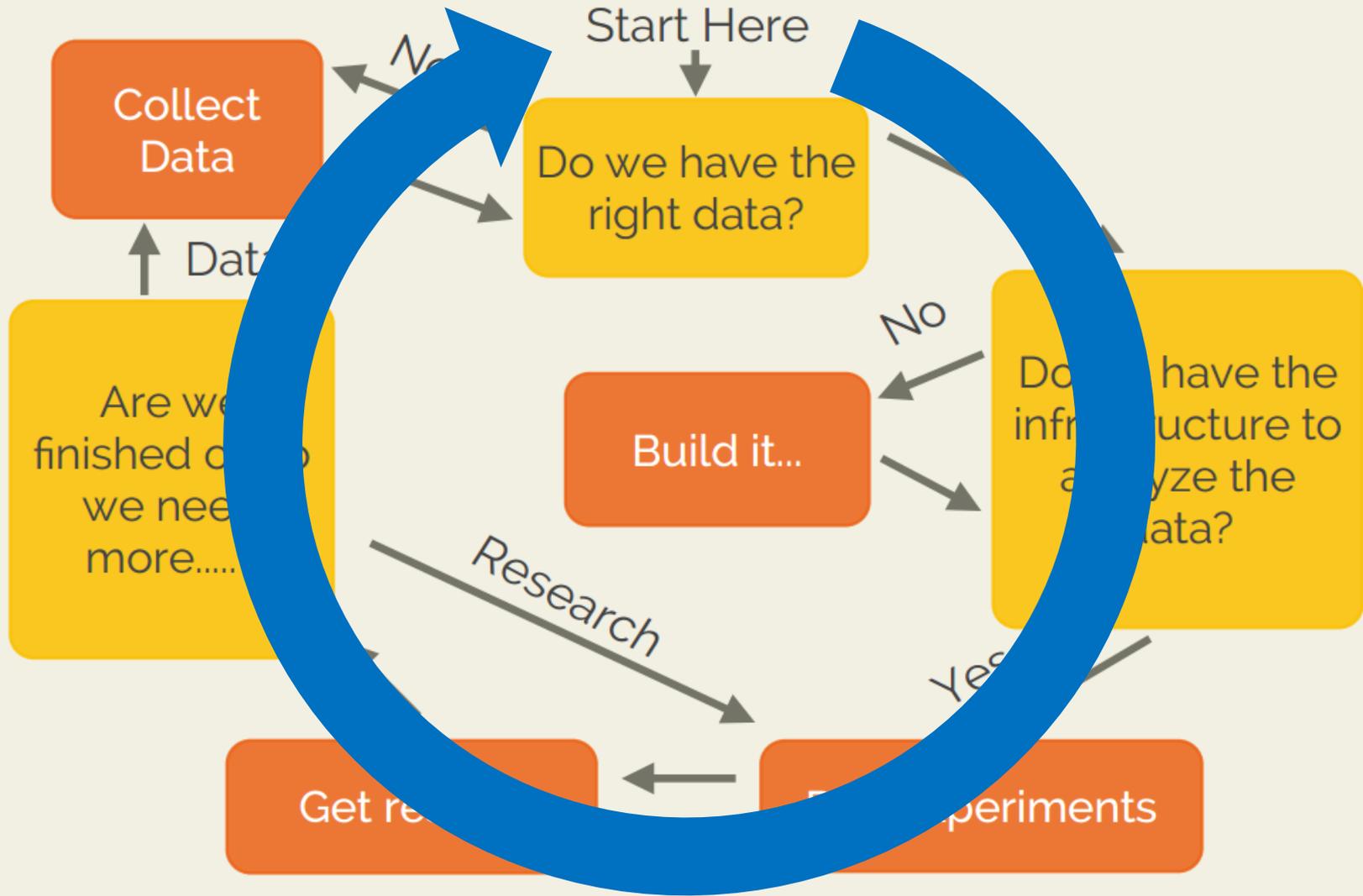
*For example, we saw unexpectedly poor performance in a number of geographies. Was it the new analytics we tried? Was it a data problem? No, it was a simple compliance problem.*

# **NATURE OF THE PROBLEM FROM THE ANALYST'S PERSPECTIVE**

# The analytics process at a high level

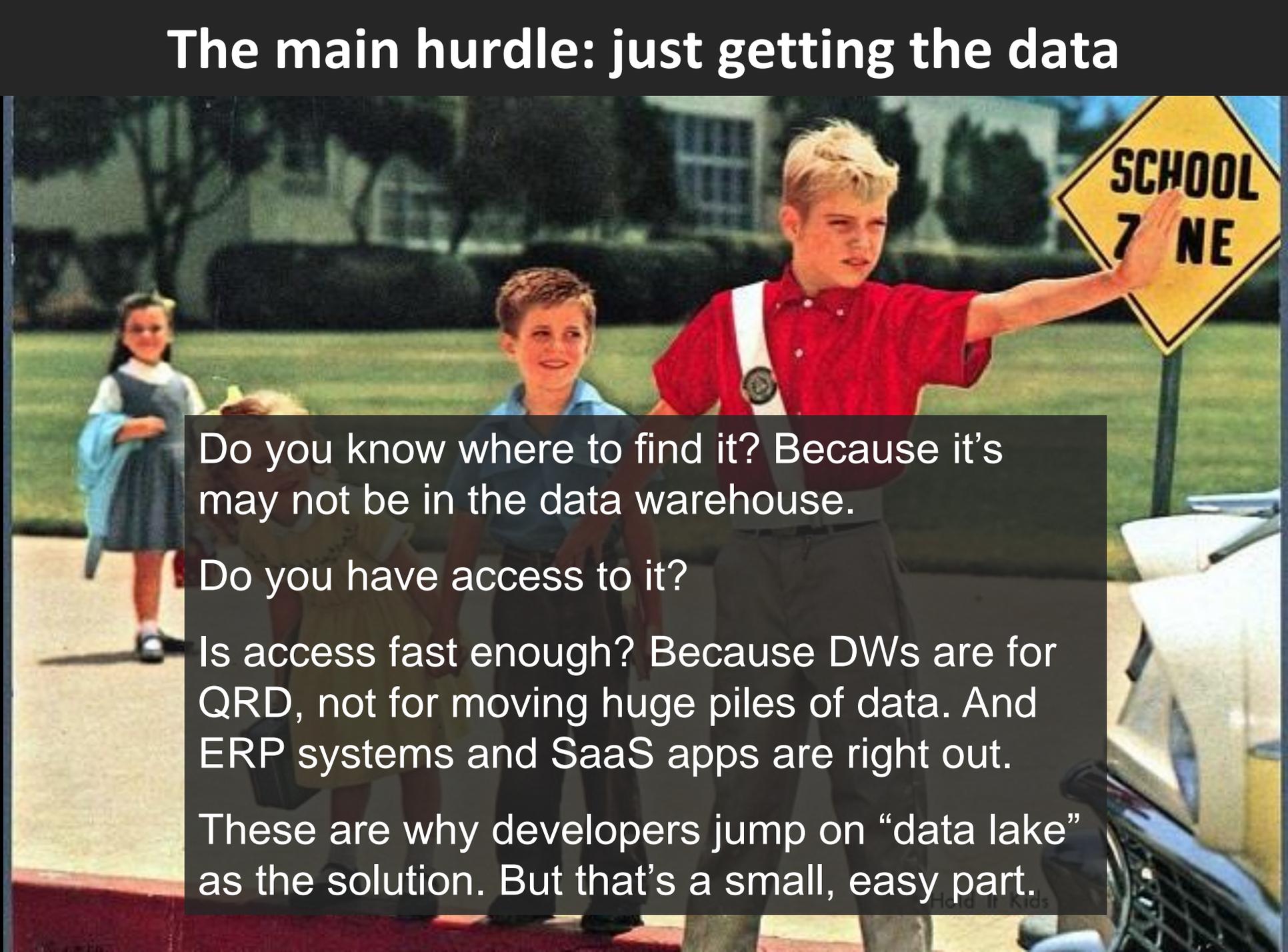


# Repeat for each new problem



The nature of analytics problems is researching the unknown rather than accessing the known.

# The main hurdle: just getting the data



Do you know where to find it? Because it's may not be in the data warehouse.

Do you have access to it?

Is access fast enough? Because DWs are for QRD, not for moving huge piles of data. And ERP systems and SaaS apps are right out.

These are why developers jump on “data lake” as the solution. But that's a small, easy part.



## Do you have the right data?

Many machine learning techniques require labeled (known good) training data:

***Supervised learning:*** a person has to define the correct output for some portion of the data. Data is divided into training sets used for model building and test sets for validating the results.

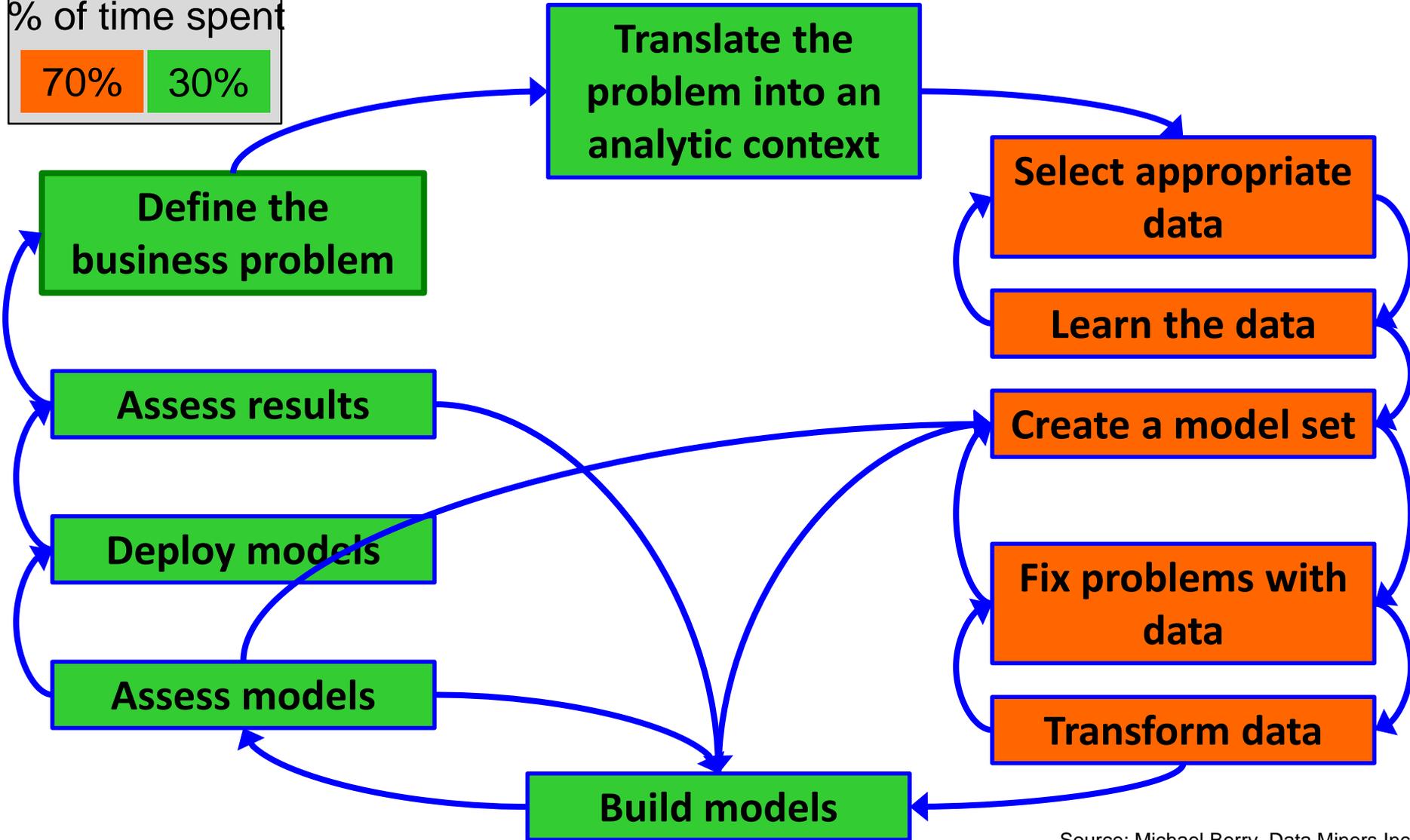
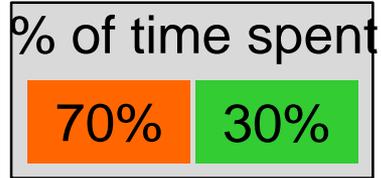
- *What is spam and what isn't?*
- *What does a fraudulent transaction look like*

**Do you have enough of the right data?**



*ML needs a lot, you may be disappointed in your efforts*

# Where do analysts spend their time? *mostly data work*



Source: Michael Berry, Data Miners Inc.

The analyst's workspace in BI is relatively spare



**The analyst's workspace needs to be more like a kitchen than like BI vending machines**



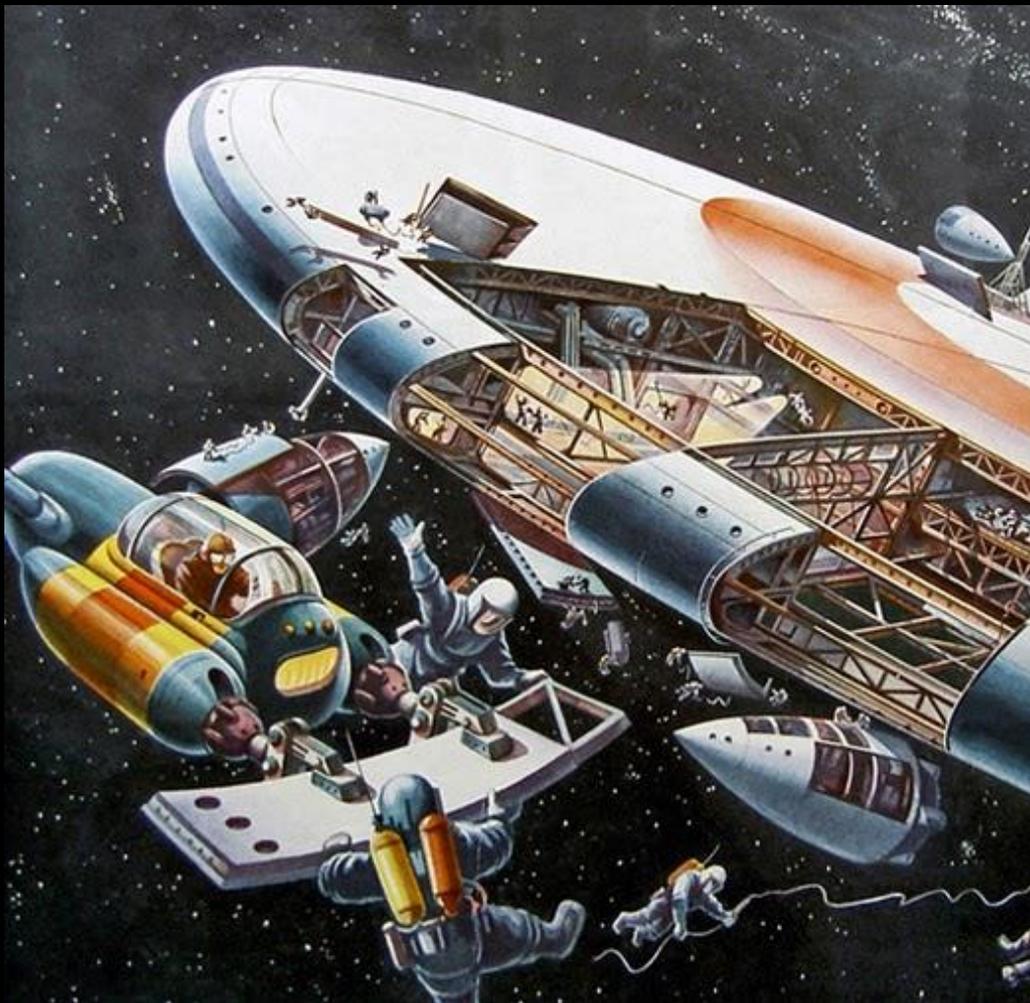
**Analytics work is a production workload, but...**



**Schema control, space limitations, resource limitations, production lockdown**

# **NATURE OF THE PROBLEM FROM THE BUILDER'S PERSPECTIVE**

# IT and Ops people want to know “what to build?”



Giant data platform?



Self service tools?

# Analytics has different processes and workloads

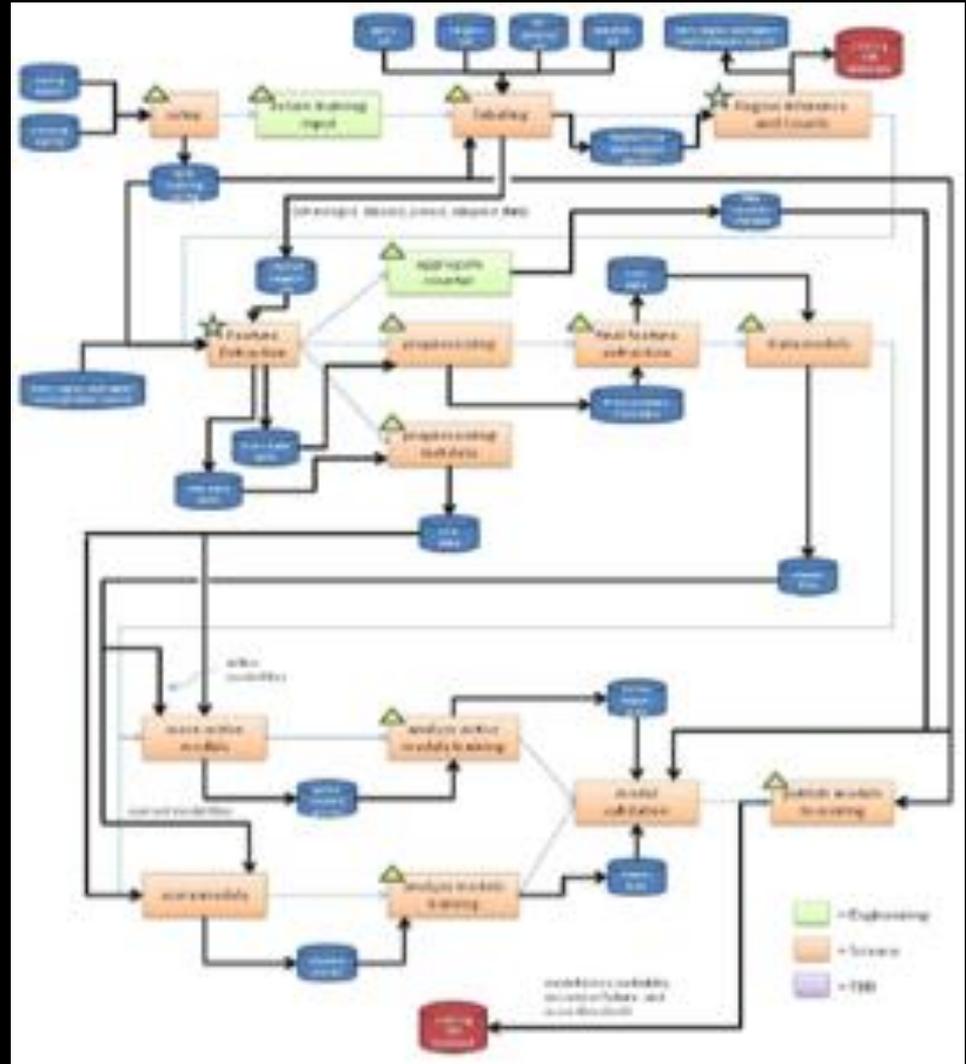
None of this analytics work is the same as what IT considered “analysis” to be, which is usually equated with BI or ad-hoc query.

Ad-hoc analysis  $\neq$

Exploratory data analysis  $\neq$

Batch analytics  $\neq$

Real-time analytics



*A real analytics production workflow*

Hatch, CIKM '11

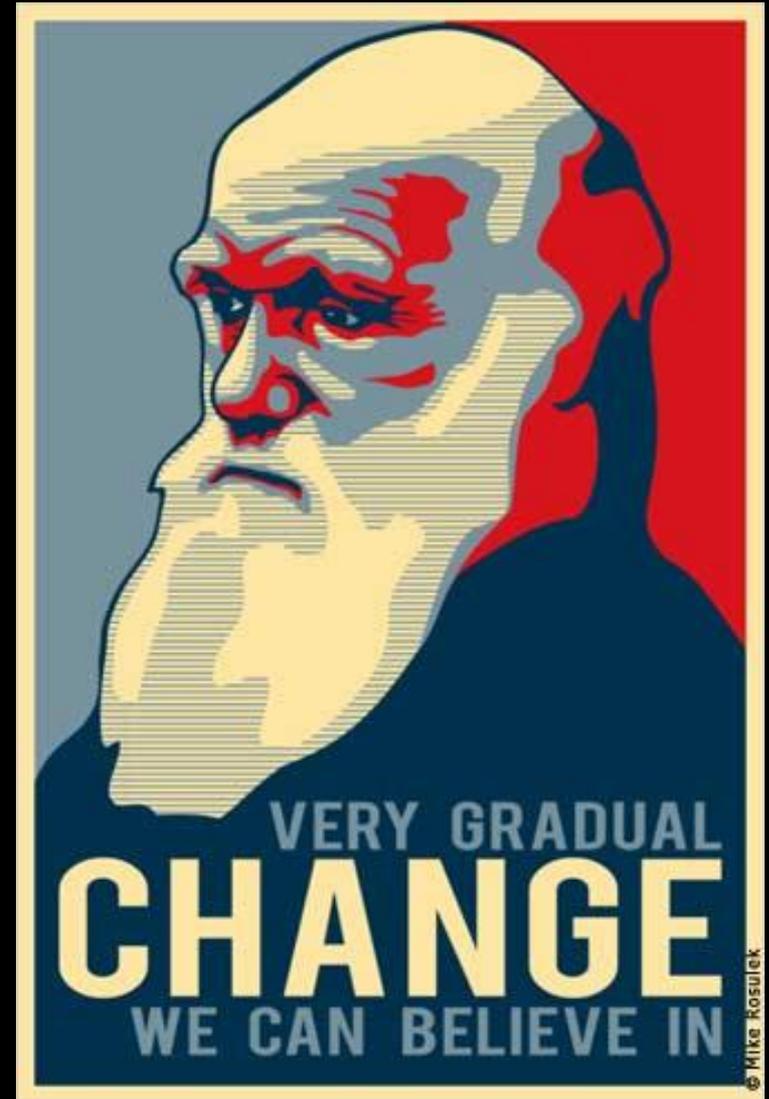
# The world changes, do the models?

In BI you maintain ETL and schemas, in ML you maintain models.

“Model decay” happens as the assumptions around which a model is built change, e.g. spam techniques change.

When you adjust the model you need to know it is normal again

- Better save the data used to build the model
- Better save the model
- Baseline and measurements



There are requirements from all constituents. You need to put them together to have a complete picture of what's needed.

**THREE PERSPECTIVES, ONE SOLUTION?**

# The missing stakeholder

There is another stakeholder: analytics management - the CAO, CDO, VP of analytics, aka “your boss” if you’re a data scientist.

This is the perspective and problems of the person responsible for oversight of the team and efforts is across the organization and across multiple projects



# Job #1 - Repeatability





# Job #3 - Reproducibility



**You need a system of record for analytics**



# There is an extensive list of requirements to support

Primary requirements needed by constituents	S	D	E
Data catalog and ability to search it for datasets		X	X
Self-service access to curated data		X	
Self-service access to uncurated (unknown, new) data		X	X
Temporary storage for working with data		X	
Data integration, cleaning, transformation, preparation tools and environment		X	X
Persistent storage for source data used by production models		X	X
Persistent storage for training, testing, production data used by models		X	X
Storage and management of models		X	X
Deployment, monitoring, decommissioning models			X
Lineage, traceability of changes made for data used by models		X	X
Lineage, traceability for model changes	X	X	X
Managing baseline data / metrics for comparing model performance	X	X	X
Managing ongoing data / metrics for tracking ongoing model performance	X	X	X

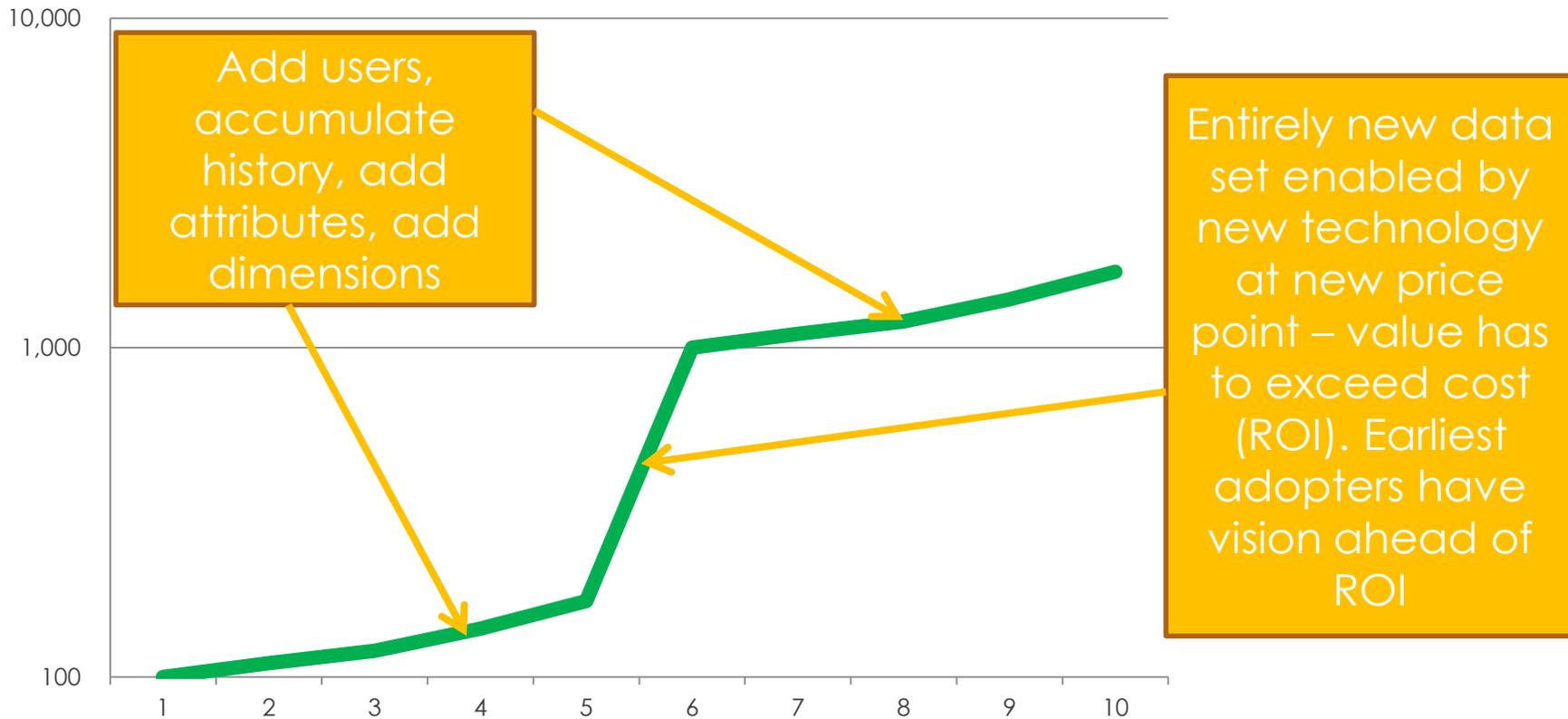
S = stakeholder, user, D = data scientist, analyst, E = engineer, developer



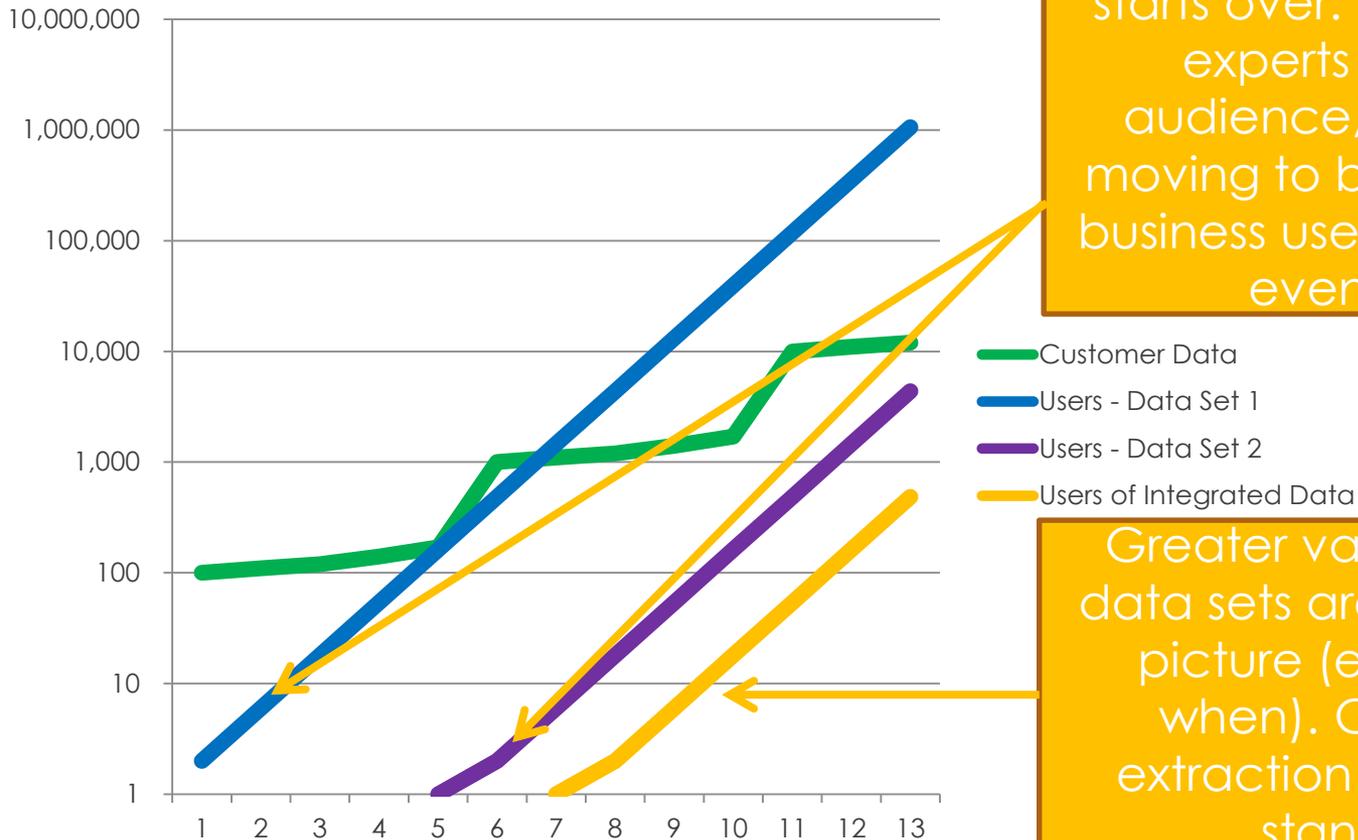
**How did we get to this point  
with BI & analytics?**

There's a difference  
between having no past  
and actively rejecting it.

## Customer Data Plateaus



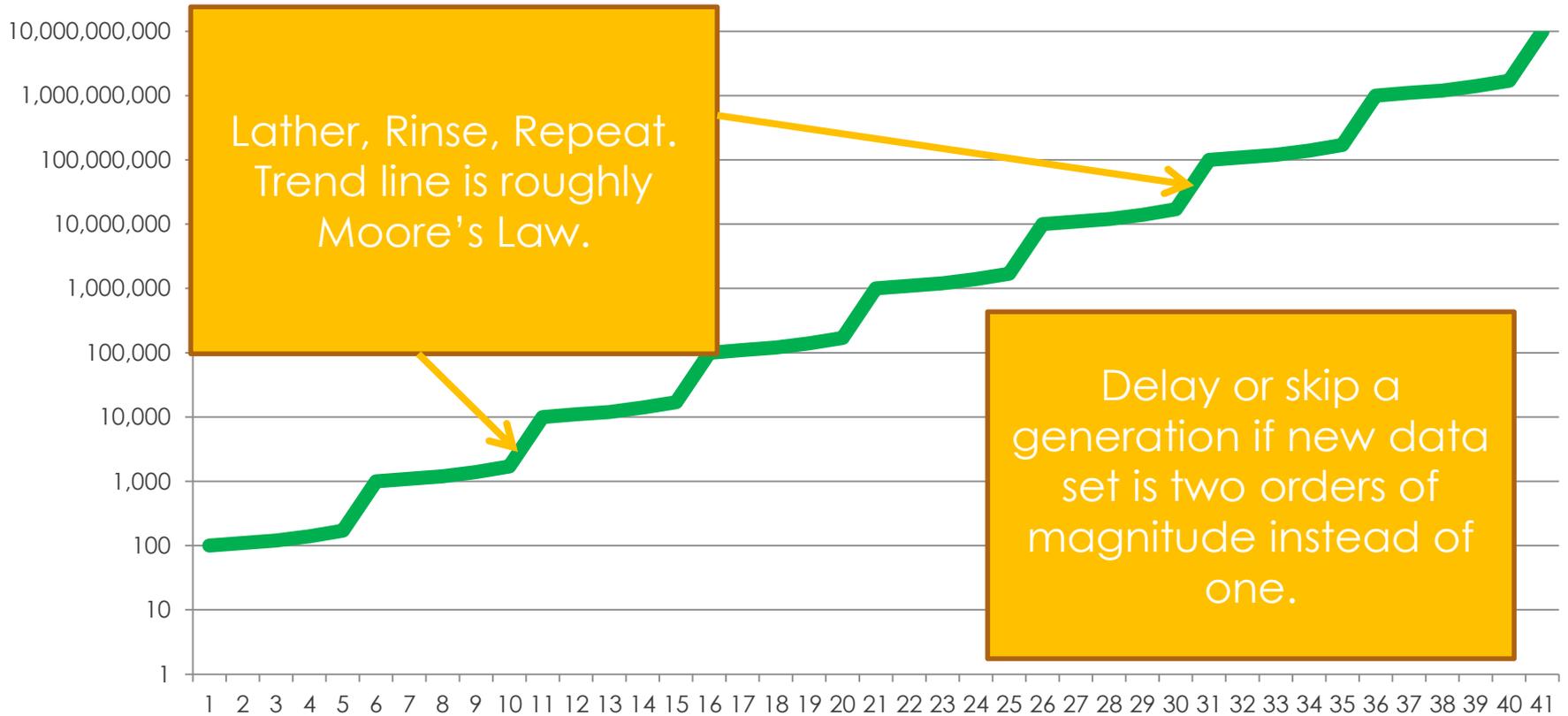
## User Adoption of New Data



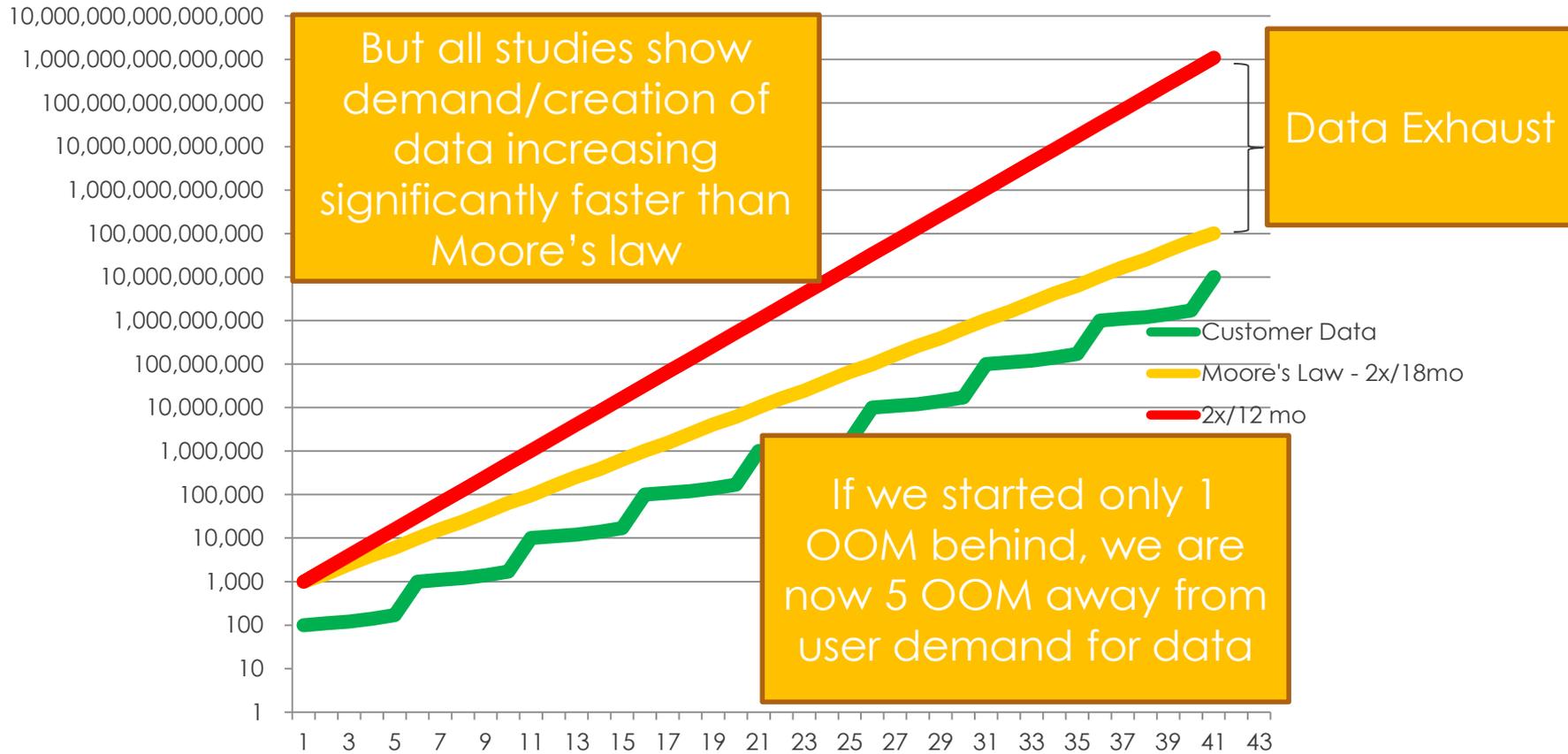
User adoption of new data sets starts over. Very small number of experts growing to wider audience, sophisticated users moving to business analysts, then business users, B2B customers and even to consumers

Greater value is derived when data sets are linked – see bigger picture (eg who buys what, when). Comes after initial extraction of easy value from standalone data

## Data Size Plateaus over Time



## Customer Data Size vs. 2x per 12mo



# Retail Plateaus

- <Store, Item, Week>
- <Store, Item, Day>
  - Simple aggregations
- Market Basket
  - Affinity
  - Link to person, demographics, HR
- Inventory by SKU by store
  - Temporal, time series, forecasting
  - Link to product, marketing, market basket



2B records total  
for 9 quarters



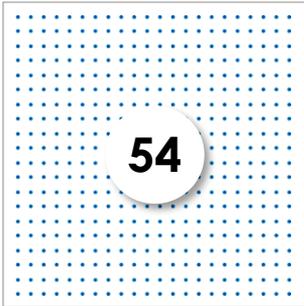
2B records per  
day, keep 9  
quarters

# Retail Plateaus

30B records per day

- Web Logs and traffic
  - Behavioral patterns – eg path linked to person, offers, other channels
  - Operations of the web site
- Supply chain sensors – sampled at major event
  - Activity Based Costing
  - Link to customer, product, HR, planning
- Social Media
  - Text analysis, Filtering, languages
  - Link to customer, sales, other channel interactions
- Supply chain sensors – sampled at minutes or seconds
  - Telematics
  - Real time, Event detection, trending, static and dynamic rules
  - Link to HR, thresholds, forecasts, routing, planning

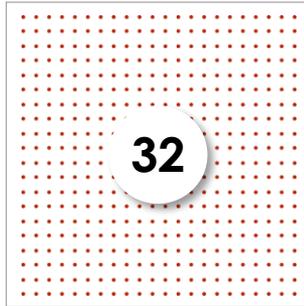
How many batteries are in inventory by plant?



## OPERATIONS

Inventory  
Returns  
Manufacturing  
Supply Chain

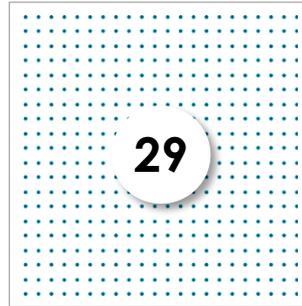
What is the trend of warranty costs?



## FINANCE

Revenue  
Expenses  
Customers

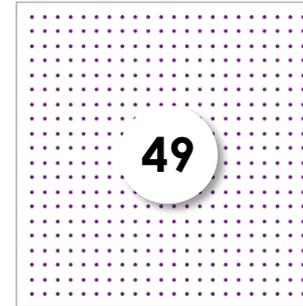
How many people made a warranty claim last week?



## CUSTOMER CARE

Customer  
Products  
Orders  
Case History

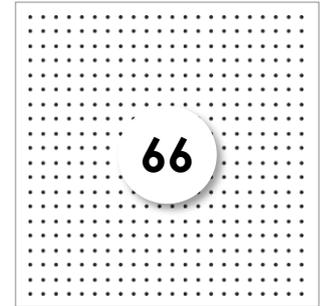
How many sales have been made quarter to date?



## SALES

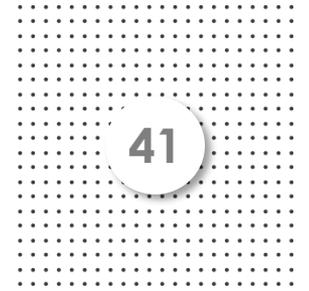
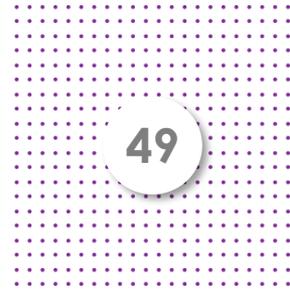
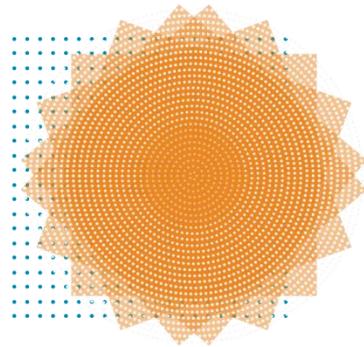
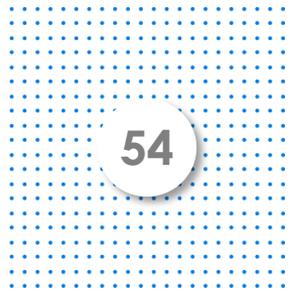
Orders  
Customers  
Products

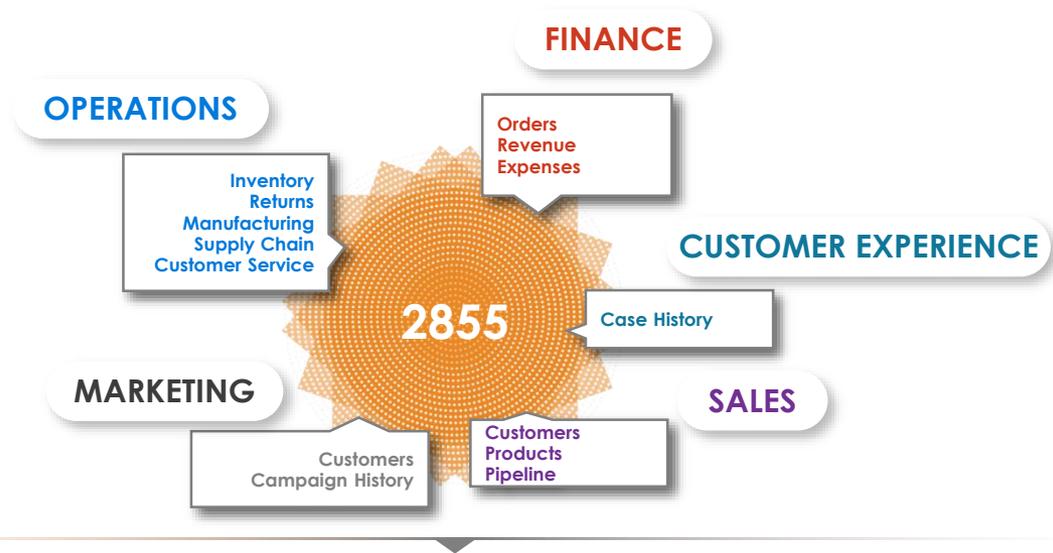
Which customers should get a communication on extended warranties?



## MARKETING

Customers  
Orders  
Campaign History





Given the rise in **warranty costs**, isolate the problem to be a **plant**, then to a **battery lot**. Communicate with **affected customers**, who have not made a warranty claim on batteries, through **Marketing** and **Customer Service** channels to recall cars with affected batteries.

# Manufacturing: Data Overlap Analysis

New Business Improvement Opportunities through Data Leverage

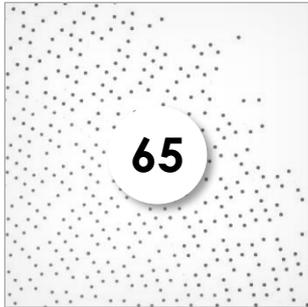
	<b>Sales Force Profitability Analysis</b>	<b>Transportation Planning</b>	<b>Production Planning</b>	<b>Vendor Managed Inventory</b>	<b>Global Pricing Rationalization</b>	<b>Fulfillment (Perfect Order)</b>	<b>Manufacturing Quality Optimization</b>	<b>Preventative Maintenance Analysis</b>	<b>Warranty Claims Analysis</b>	<b>Quality Life Cycle Improvement</b>
<b>Sales Force Profitability Analysis</b>	<b>100%</b>	80%	66%	24%	41%	66%	0%	24%	0%	24%
<b>Transportation Planning</b>	11%	<b>100%</b>	<b>87%</b>	28%	64%	56%	22%	34%	13%	45%
<b>Production Planning</b>	6%	57%	<b>100%</b>	19%	<b>83%</b>	35%	17%	28%	9%	40%
<b>Vendor Managed Inventory</b>	12%	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>78%</b>	<b>100%</b>	61%	<b>100%</b>	39%	<b>100%</b>
<b>Global Pricing Rationalization</b>	4%	50%	<b>100%</b>	17%	<b>100%</b>	27%	15%	29%	11%	41%
<b>Fulfillment (Perfect Order)</b>	16%	<b>94%</b>	<b>90%</b>	48%	58%	<b>100%</b>	35%	53%	19%	56%
<b>Manufacturing Quality Optimization</b>	0%	<b>76%</b>	<b>88%</b>	60%	66%	<b>71%</b>	<b>100%</b>	<b>79%</b>	43%	<b>73%</b>
<b>Preventative Maintenance Analysis</b>	7%	<b>75%</b>	<b>93%</b>	61%	<b>80%</b>	68%	50%	<b>100%</b>	27%	<b>82%</b>
<b>Warranty Claims Analysis</b>	0%	<b>94%</b>	<b>100%</b>	<b>83%</b>	<b>100%</b>	<b>83%</b>	<b>94%</b>	<b>93%</b>	<b>100%</b>	<b>98%</b>
<b>Quality Life Cycle Improvement</b>	4%	57%	<b>75%</b>	35%	64%	41%	26%	47%	16%	<b>100%</b>

Then

If

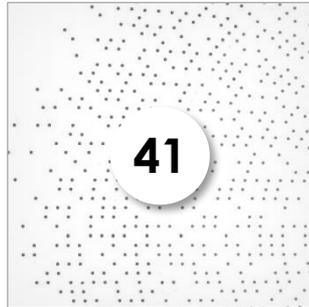


How many visitors did we have to our hybrid cars microsite yesterday?



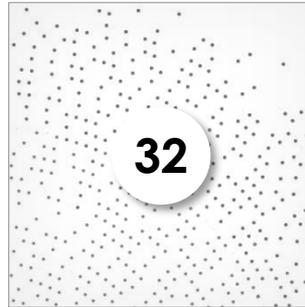
**CLICKSTREAM**

What are the temperature readings for batteries by Manufacturer?



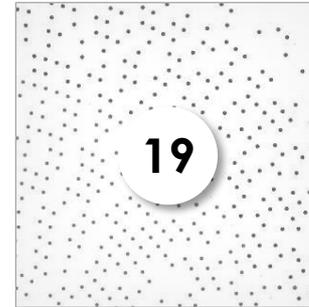
**PRODUCT SENSOR**

What is the sentiment towards line of hybrid vehicles?



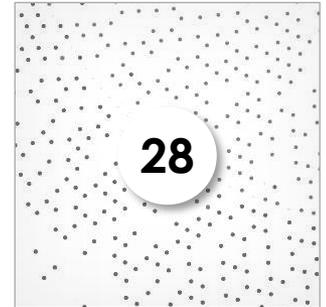
**SOCIAL MEDIA**

Which customers likely expressed anger with customer care?



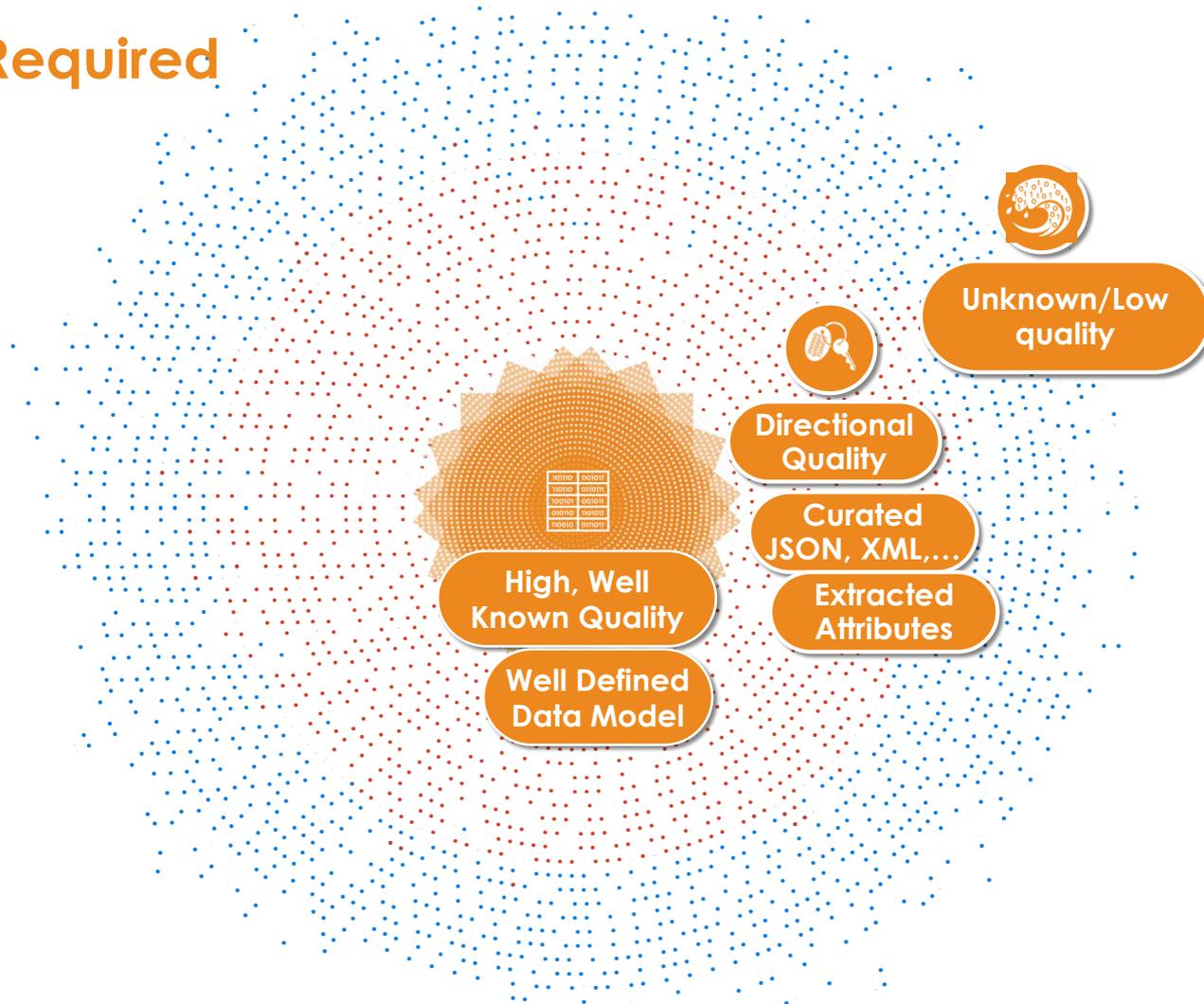
**CUSTOMER CARE AUDIO RECORDINGS**

Which ad creative generated the most clicks?

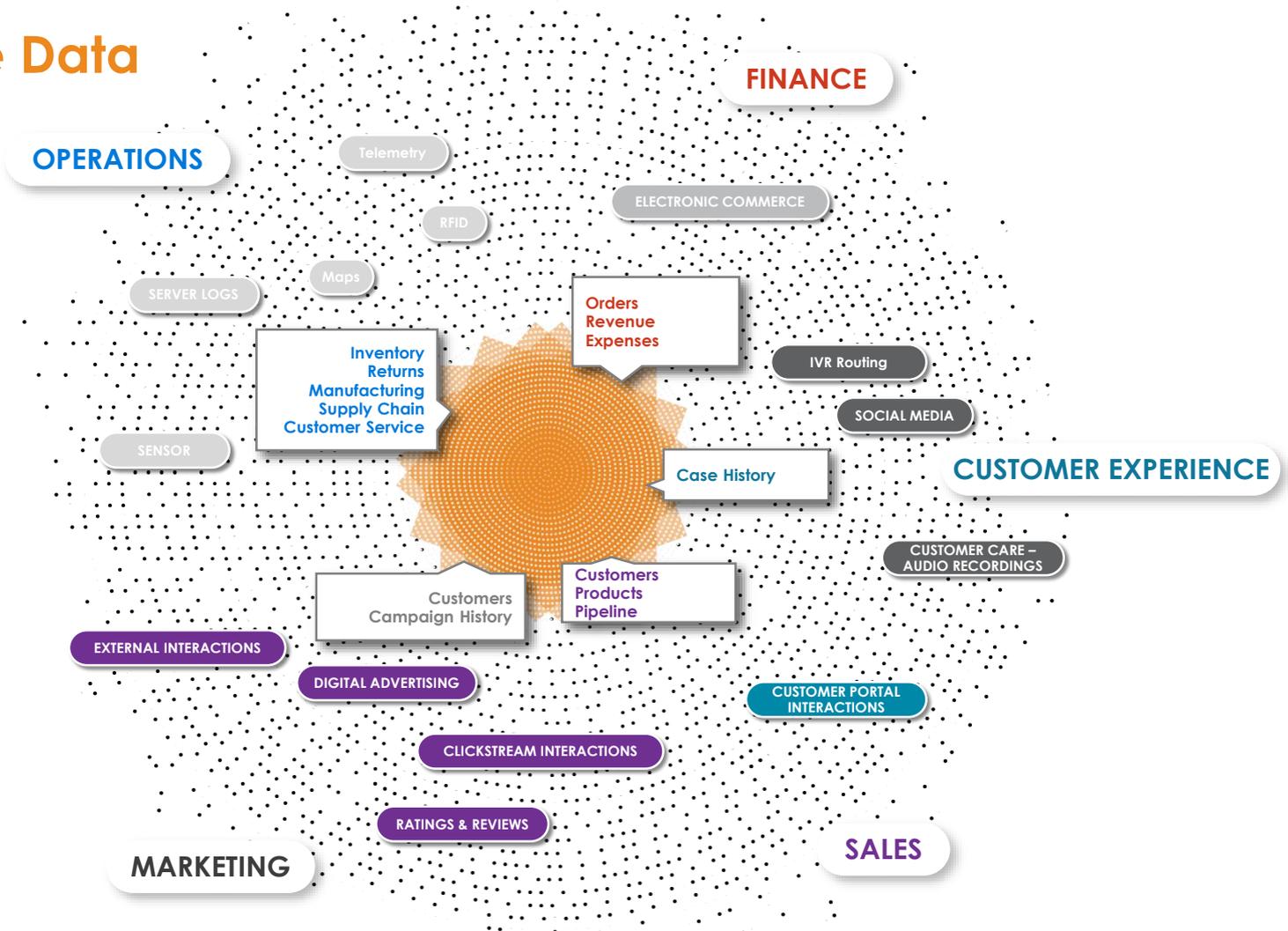


**DIGITAL ADVERTISING**

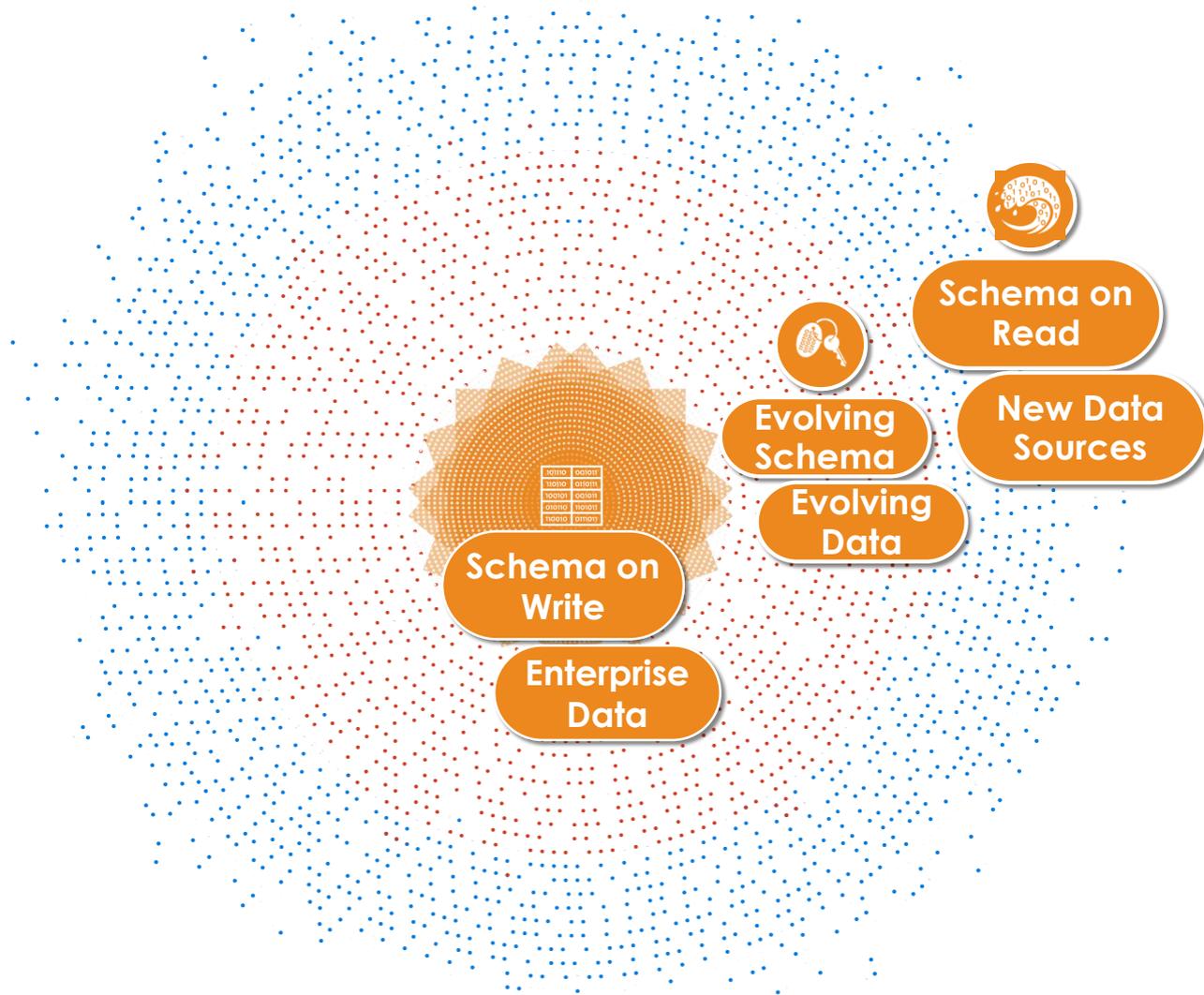
# Curation Required



# Enterprise Data

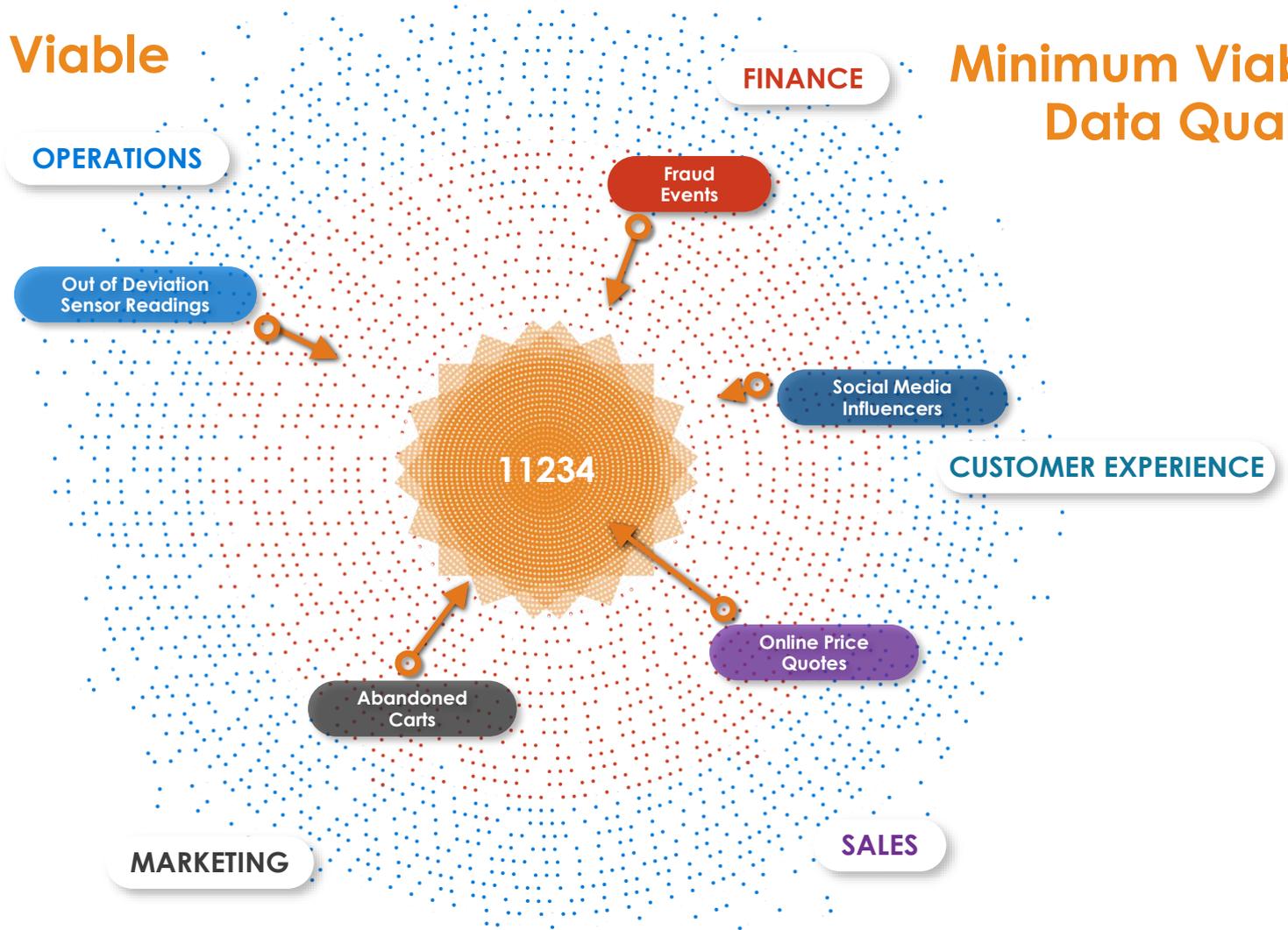


# Schema?

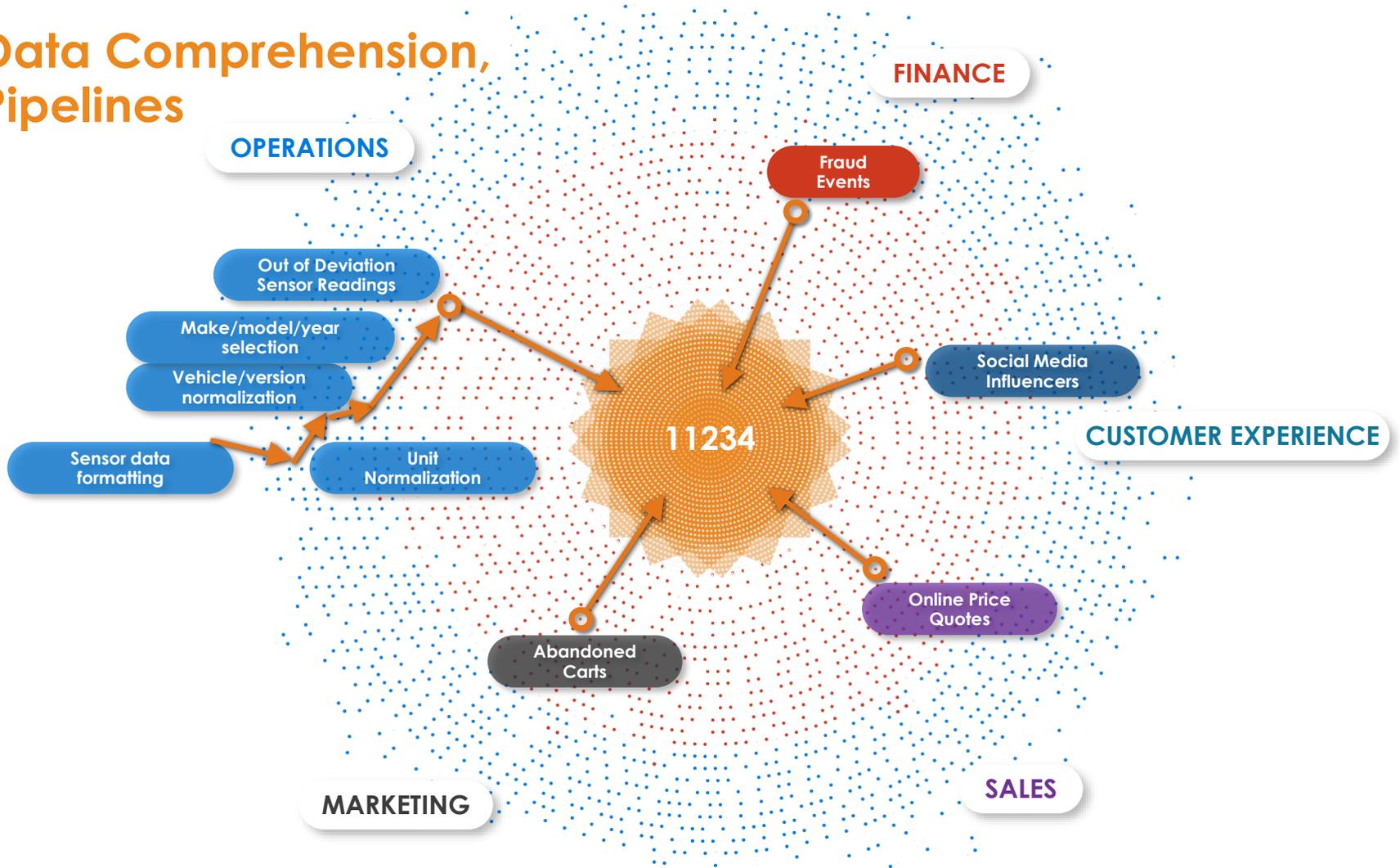


# Minimum Viable Curation

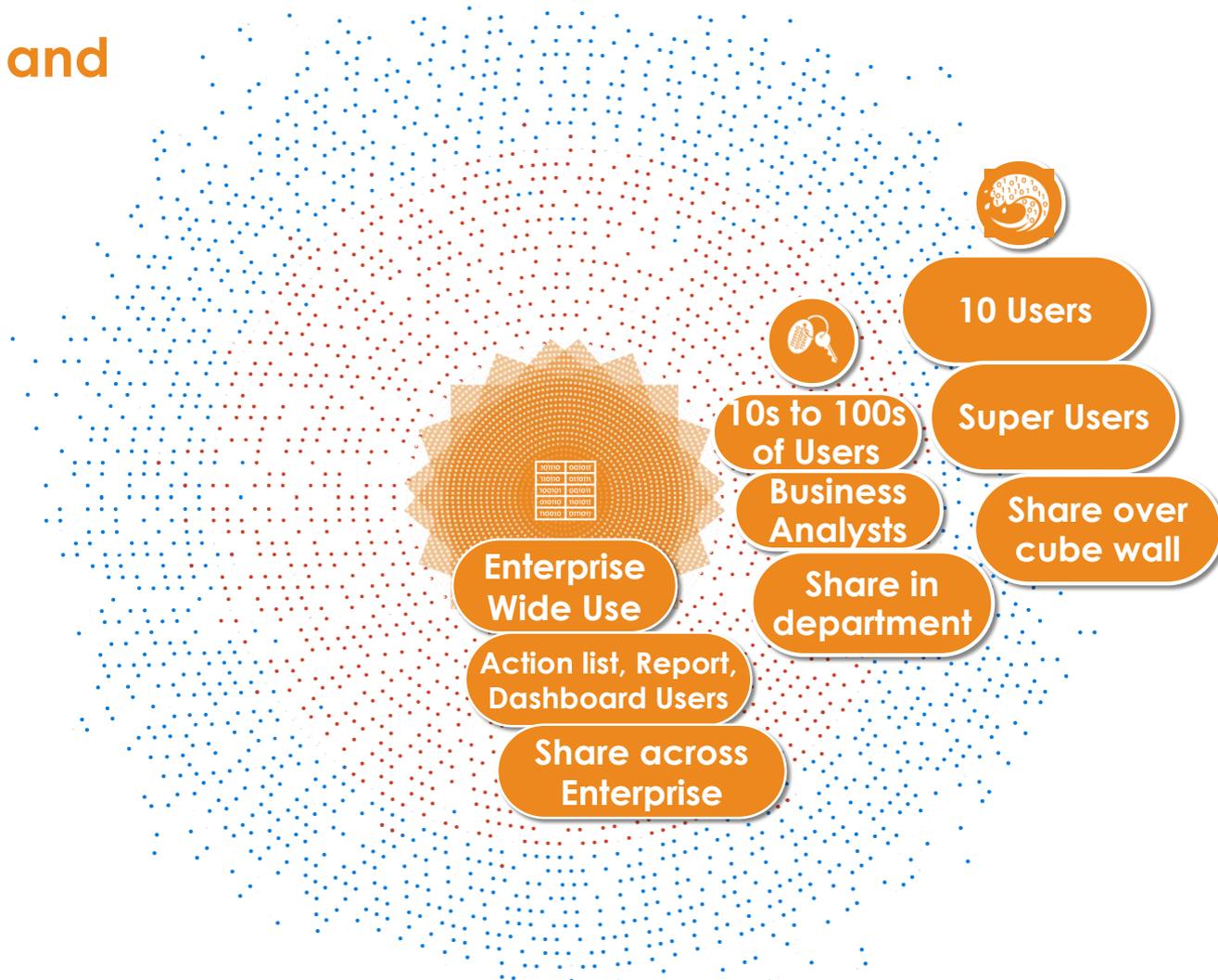
# Minimum Viable Data Quality



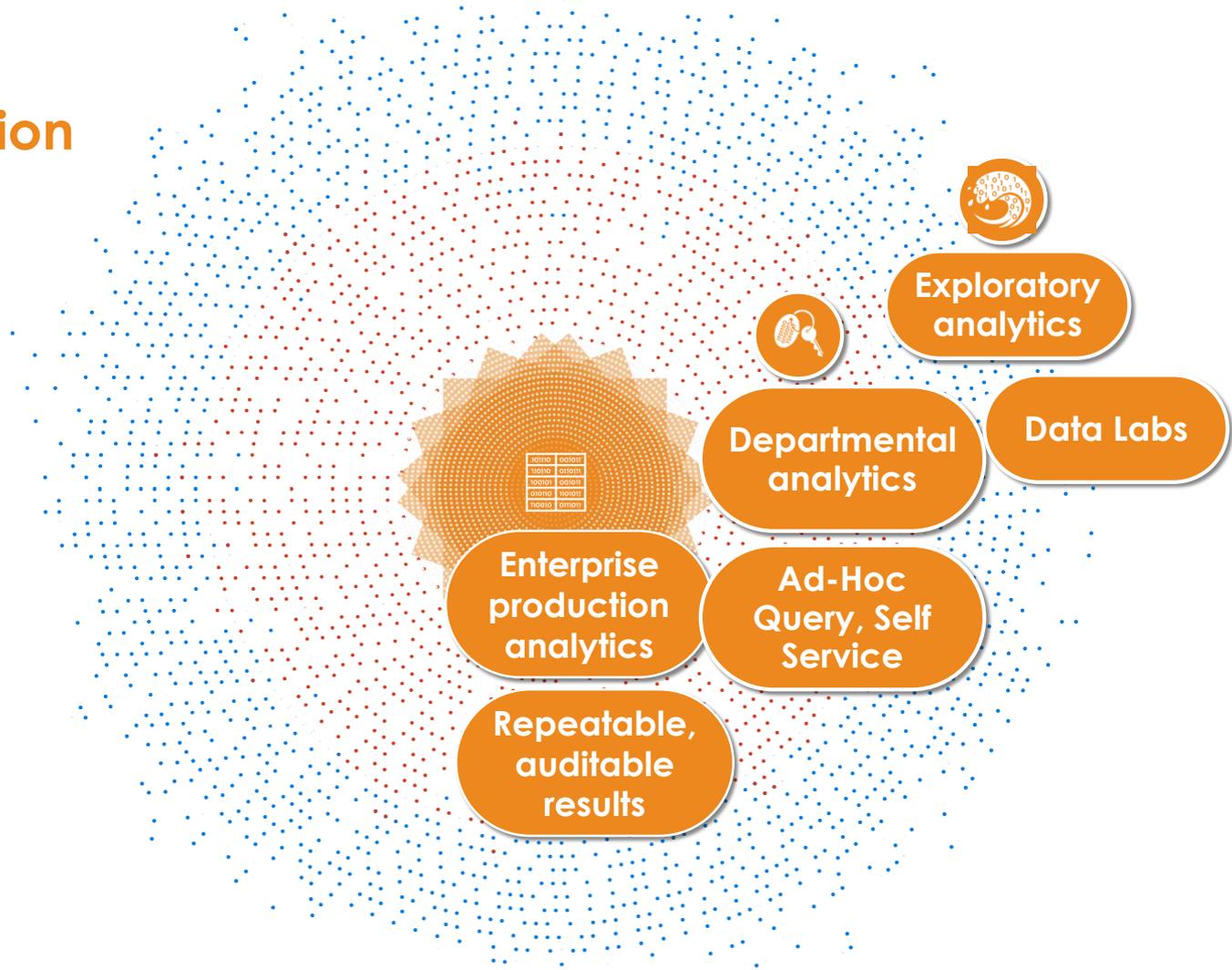
# Data Comprehension, Pipelines



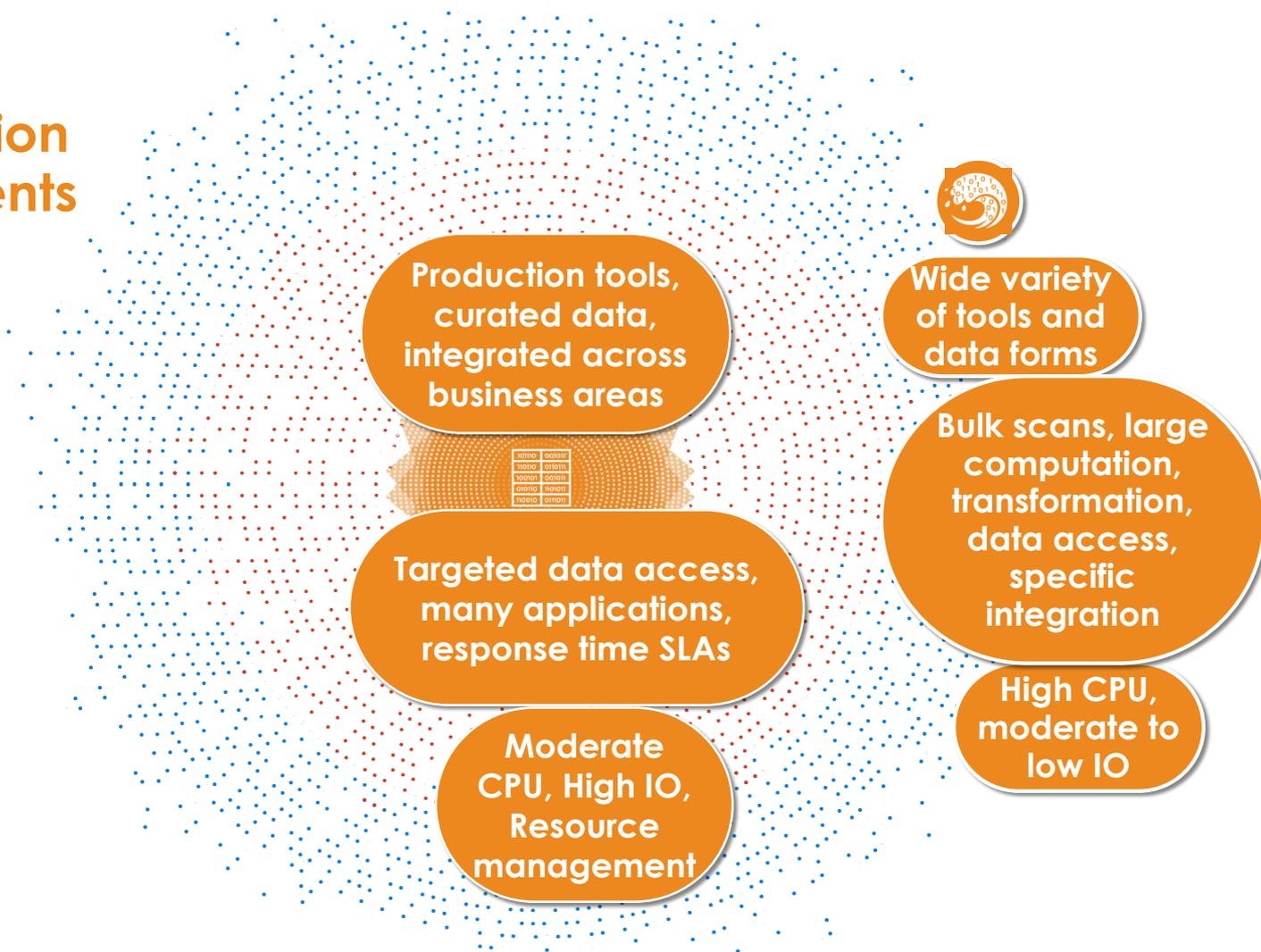
# User Base and Sharing



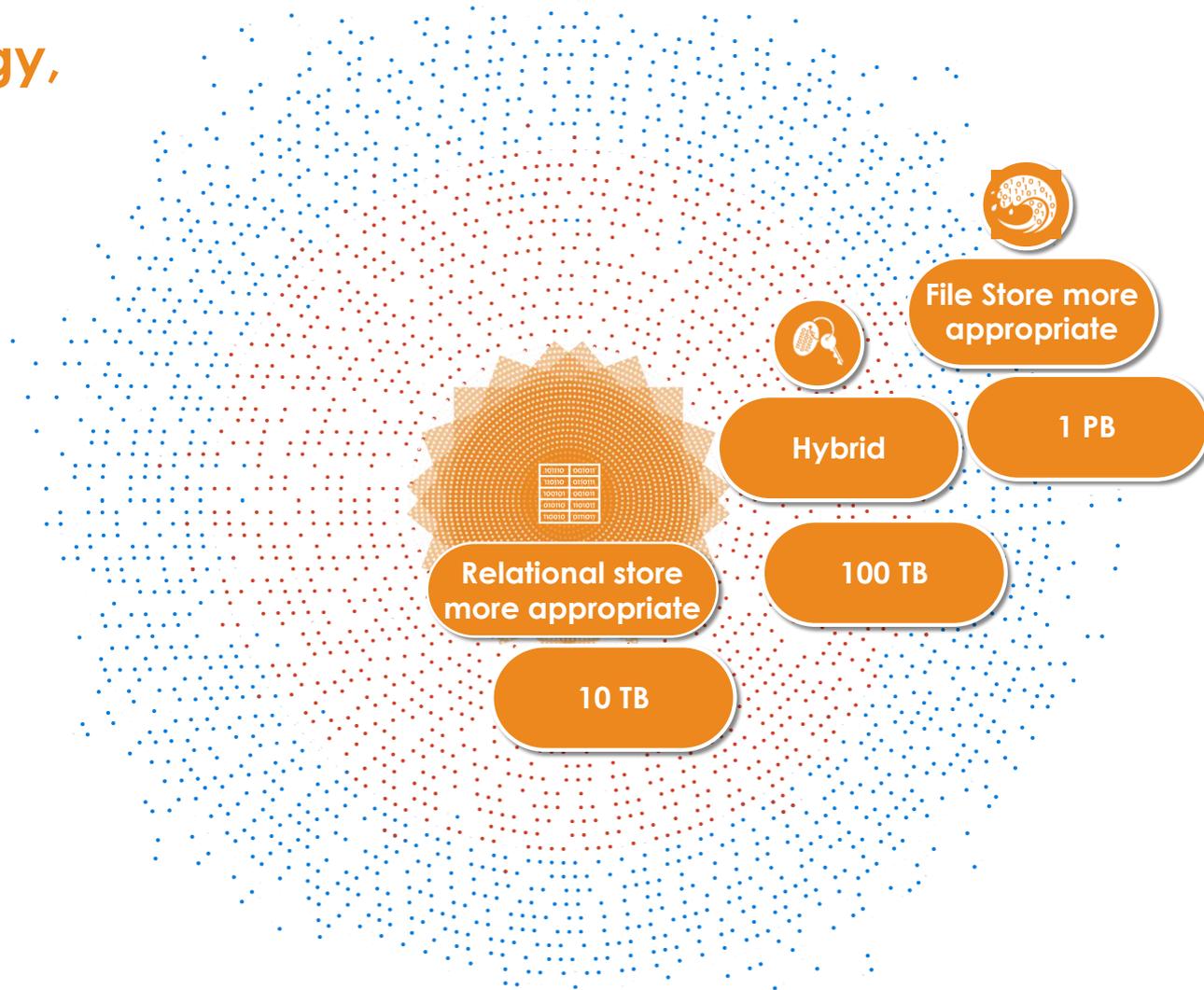
# Evolving Consumption



# Evolving Consumption Requirements



# Technology, Capacity



# Access Wide Variety of Data to Answer a Question

OPERATIONS

FINANCE

Given the rise in warranty costs, isolate the problem to be a plant and the specific lot.

Exclude 2/3<sup>rd</sup> of the batteries from the lot that are fine.

Communicate with affected customers, who have not made a warranty claim, through Marketing and Customer Service channels to recall cars with affected batteries.

SENSOR

MANUFACTURING

COSTS

CASE HISTORY

PRODUCTS

CUSTOMERS

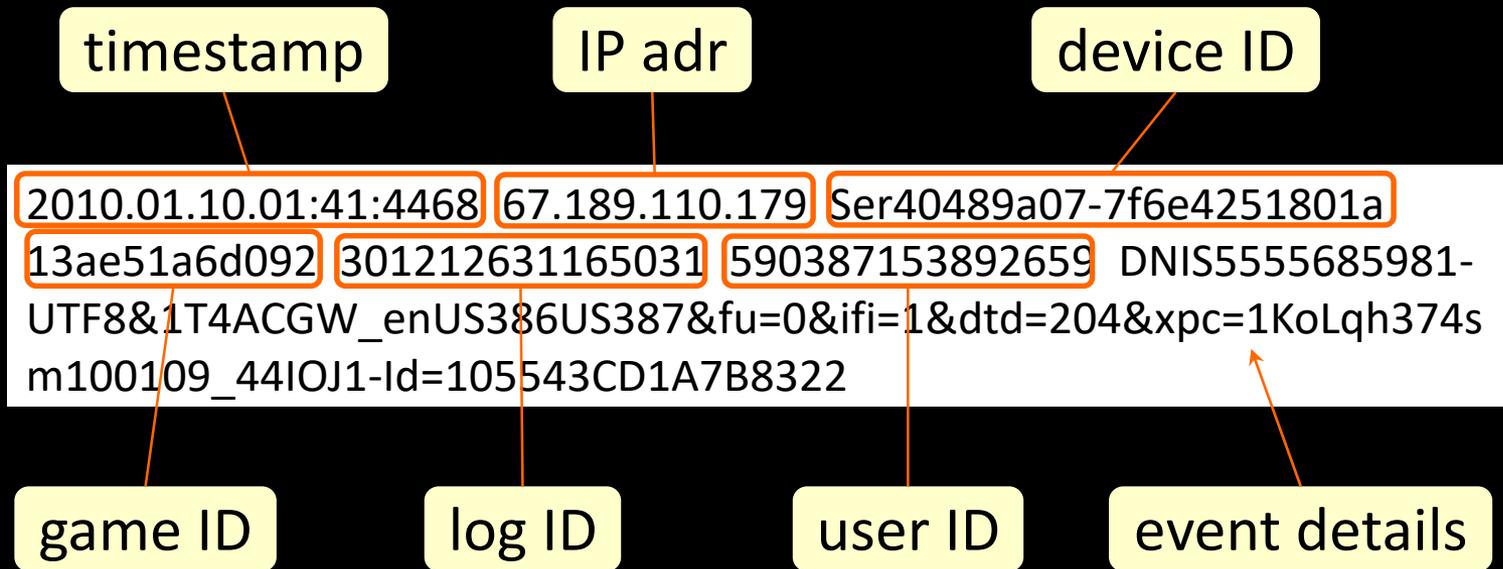
CAMPAIGN HISTORY

CUSTOMER EXPERIENCE

MARKETING

SALES

# An event contains mainly IDs...that reference other data



Log de-referencing and enrichment is difficult since you can't enforce integrity like you can in a DB.

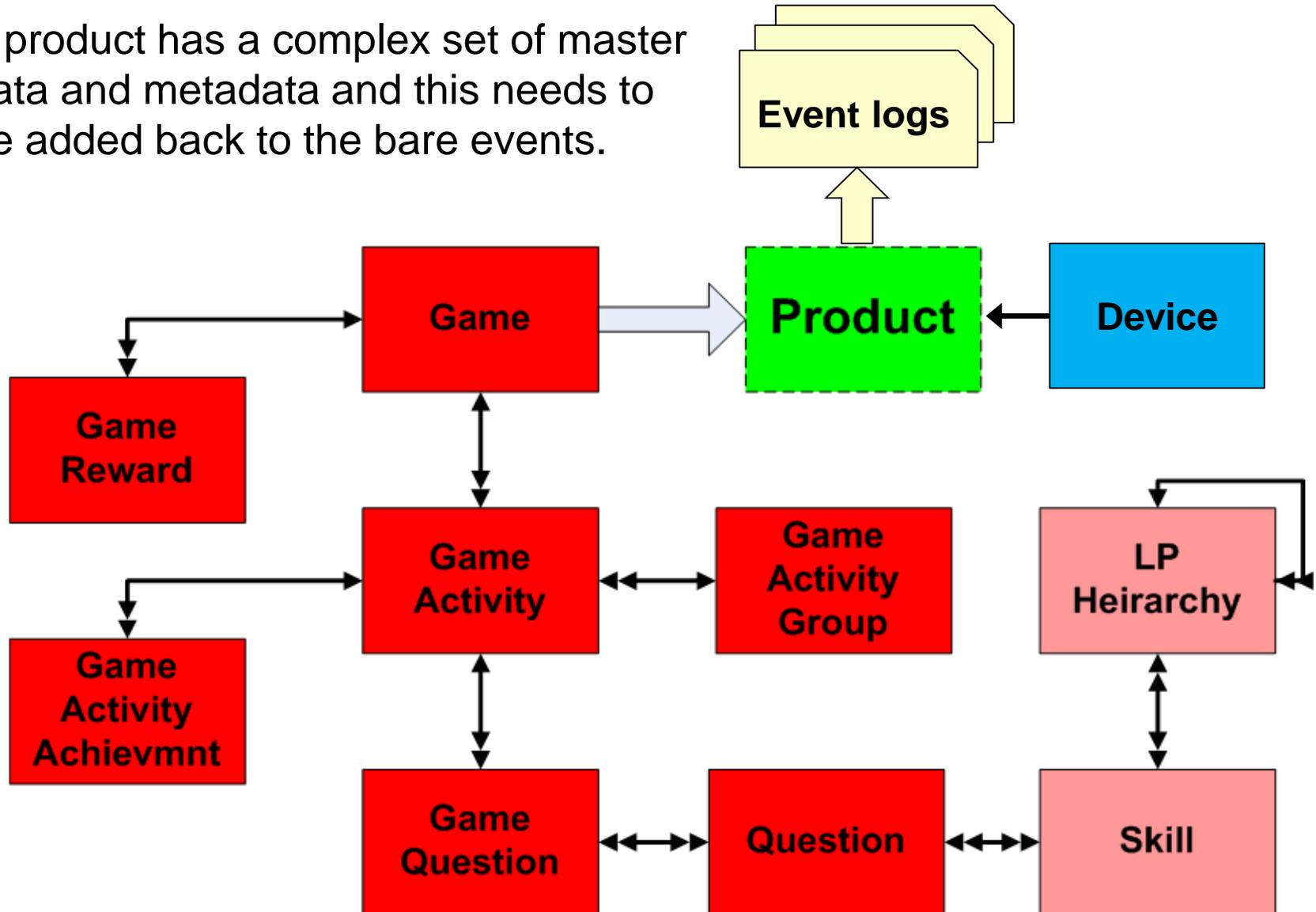
What's the glue that holds it together?

It's just keys to other data.

e.g. remember that device ID 0 problem?

# It's not just the log data, it's big data + small data

A product has a complex set of master data and metadata and this needs to be added back to the bare events.



# Where does the reference data come from?

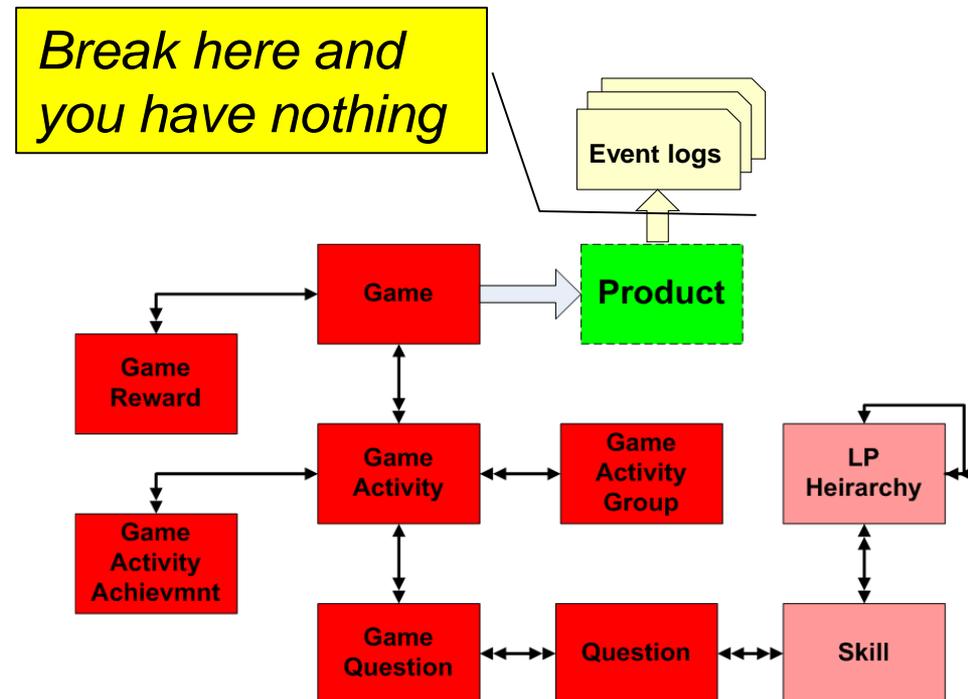
The keys come from somewhere. You don't just make up system-wide unique identifiers in your code.

The lack of local lookup data at event generation leads to development practices that lead to inconsistencies.

Problem we had: product identifiers that didn't match any known products

Had to fix by analyzing each log of bad-ID events.

Because developers used config files they copied from the PMS and put in the code



# MDM again

If you want to link datasets then you must manage the keys  
You need canonical forms for common data (in code too)

Event

2010.01.10 01:41:4468 67.189.110.179 Ser40489a07-7f6e4251801a  
13ae51a6d092 301217331137031 590387153892659 DNIS5555685981-  
UTF8&1T4ACGW\_enUS55555587&fu=0&ifi=1&dtd=204&xpc=1KoLqh374s  
m100109\_44IOJ1-Id=105543CD1A7B8322

date

IP adr

Click

2010.01.10. 14:26:2468 67.189.110.179 10098213 5046876319474403 MOZILLA/4.0  
(COMPATIBLE; TRIDENT/4.0; GTB6; .NET CLR 1.1.4322) https://w game ID ng.com/  
gifts/store/LogonForm?mmc=link-src-email\_m100109 http://www.google.com/search?  
sourceid=navclient&aq=0h&oq=Italian&ie=UTF8&pid=1T4ACGW\_13ae51a6d092&q=ita  
lian+rose&fu=0&ifi=1&dtd=204&xpc=1KoLqh374s

user ID

customer ID

Cust-  
user

UID	CID	Email	City	State	Country
590387153892659	10098213	barry.dylan@odin.com	Paris	Île-de-France	France

# What you think you have





**The solution to our problems isn't  
technology, it's architecture.**

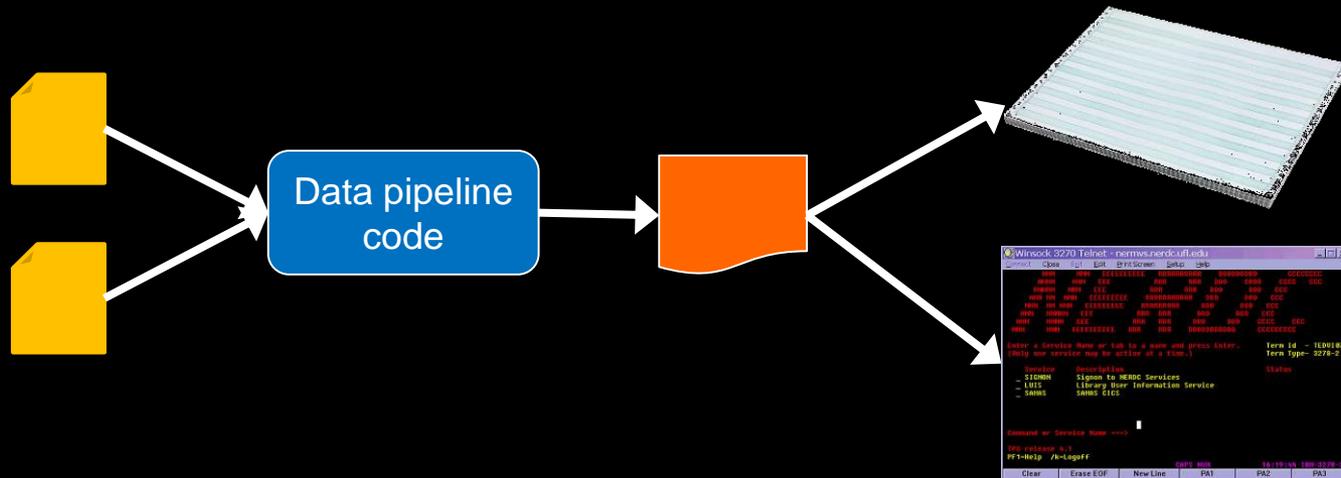


# History: This is how BI was done through the 80s

First there were files and reporting programs.

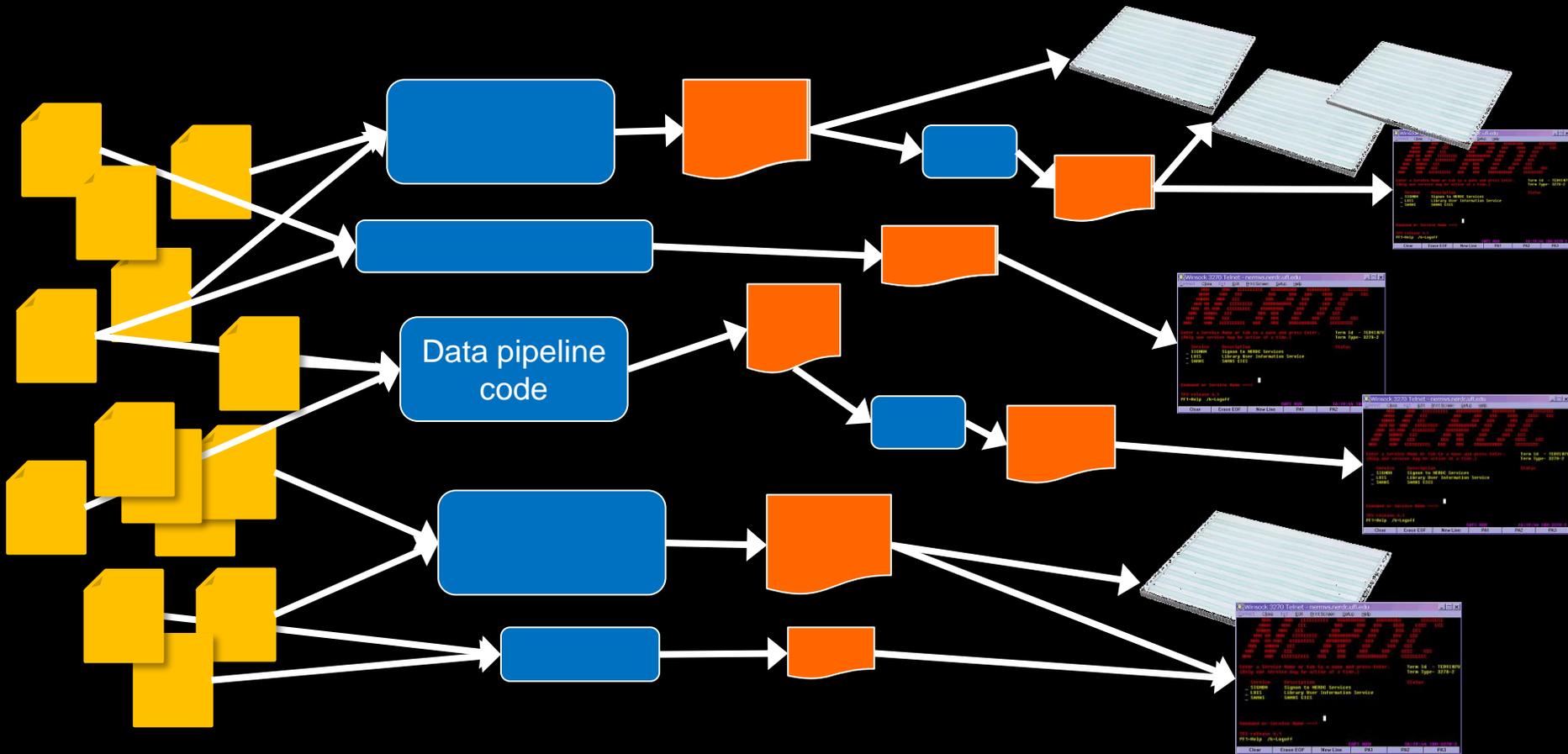
Application files feed through a data processing pipeline to generate an output file. The file is used by a report formatter for print/screen.

Every report is a program written by a developer.



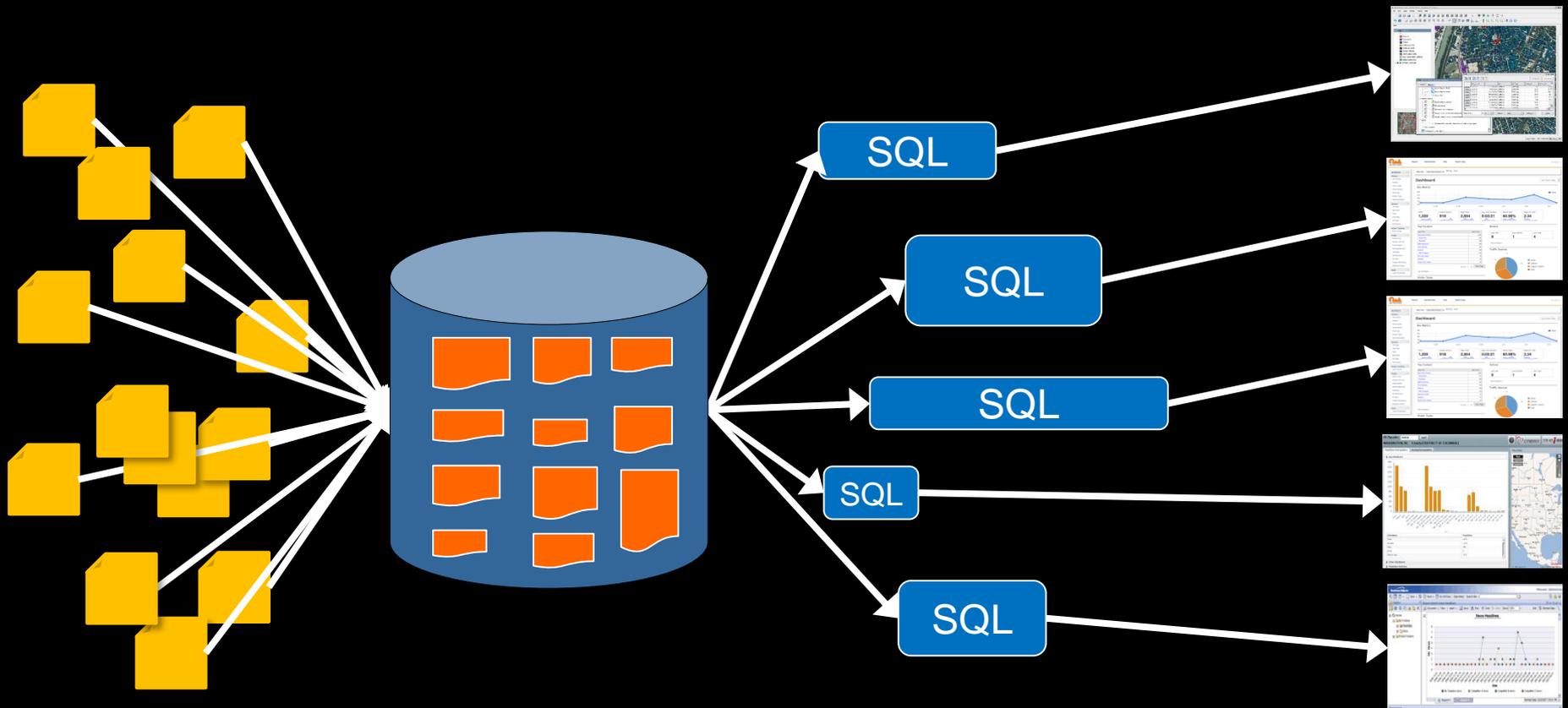
# History: This is how BI ended the 80s

The inevitable situation was...



# History: This is how we started the 90s

Collect data in a database. Queries replaced a *LOT* of application code because much was just joins. We learned about “dead code”





**One of the reasons the "Big Band Era" ended.**

# Pragmatism and Data

Lessons learned during the ad-hoc SQL era of the DW market:

When the technology is awkward for the users, the users will stop trying to use it.

Even "simple" schemas weren't enough for anyone other than analysts and their Brio...

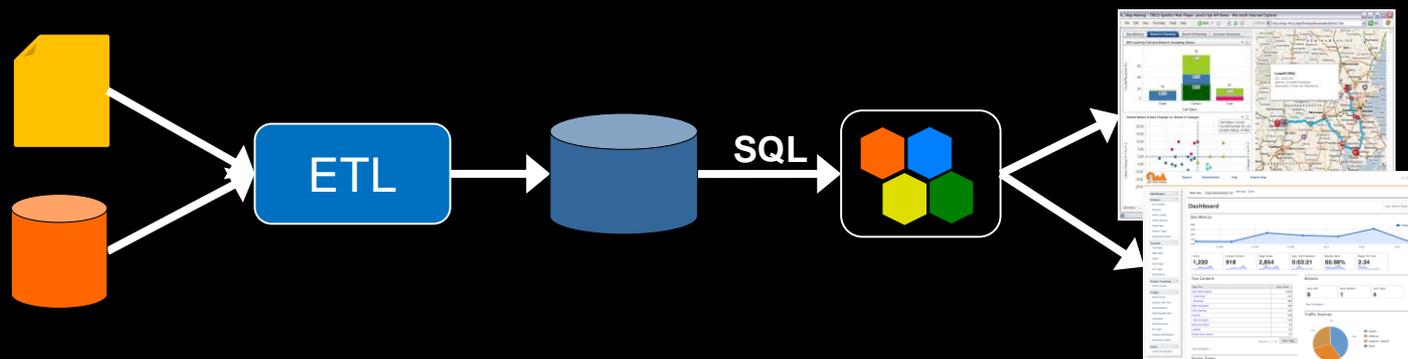
Led to the evolution of metadata-driven SQL-generating BI tools, ETL tools.

# BI evolved to hiding query generation for end users

With more regular schema models, in particular dimensional models that didn't contain cyclic join paths, it was possible to automate SQL generation via semantic mapping layers.

We developed data pipeline building tools (ETL).

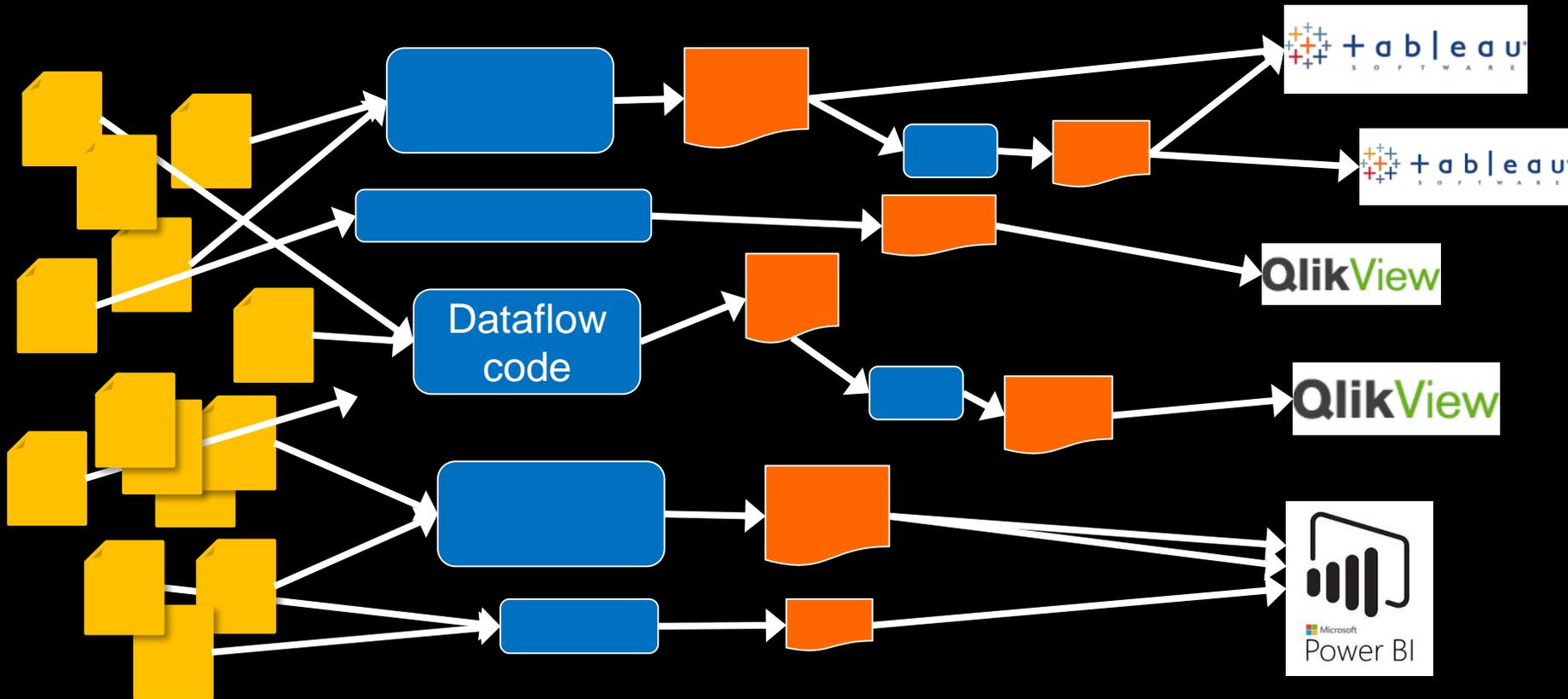
Query via business terms made BI usable by non-technical people.



*Life got much easier...for a while*

# Today's model: Lake + self-service, looks familiar...

The Lake with data pipelines to files or Hive tables is exactly the same pattern as the COBOL batch..



*We already know that people don't scale...*



We're so focused on the light switch that we're not talking about the light

# DATA ARCHITECTURE

# Decouple the Architecture

The core of the data warehouse isn't the database, it's the data architecture that the database and tools implement.

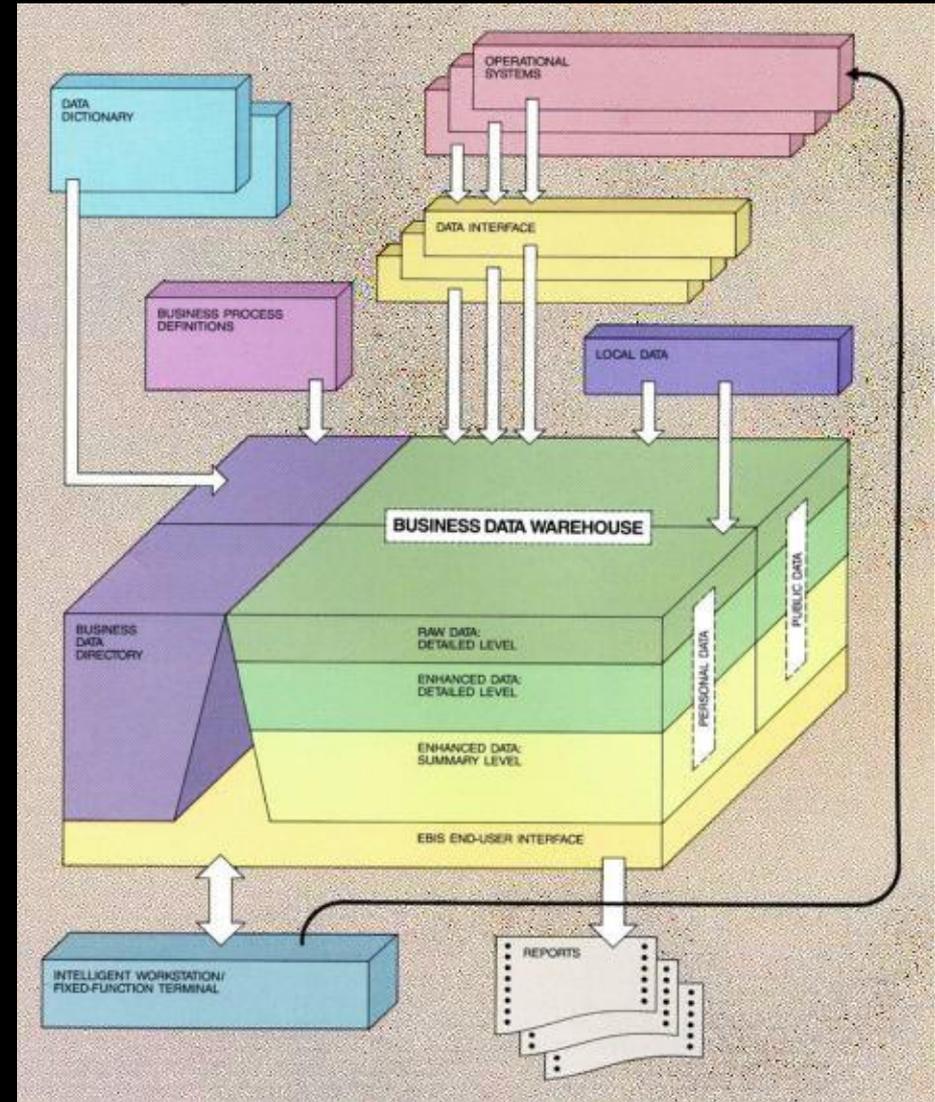
We need a data architecture that is not limiting:

- Deals with change more easily and at scale
- Does not enforce requirements and models up front
- Does not limit the format or structure of data
- Assumes the range of data latencies in and out, from streaming to one-time bulk
- Allows both reading and writing of data from outside

# The architecture from 1988 we SHOULD HAVE BEEN USING

The general concept of a separate architecture for BI has been around longer, but this paper by Devlin and Murphy is the first formal data warehouse architecture and definition published.

*“An architecture for a business and information system”, B. A. Devlin, P. T. Murphy, IBM Systems Journal, Vol.27, No. 1, (1988)*

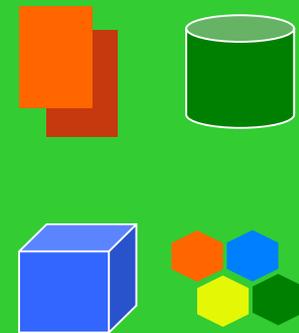


# The goal is to decouple: solve the application and infrastructure problems separately

Data access is already somewhat separate today. Make the separation of different access methods a formal part of the architecture. Don't force one model.

Storage

Data Access  
Deliver & Use



Platform Services

This separates BI from other uses of data, allowing each type of use to structure the data specific to its own requirements.

# The goal is to decouple: solve the application and infrastructure problems separately

Data Management  
Process & Integrate



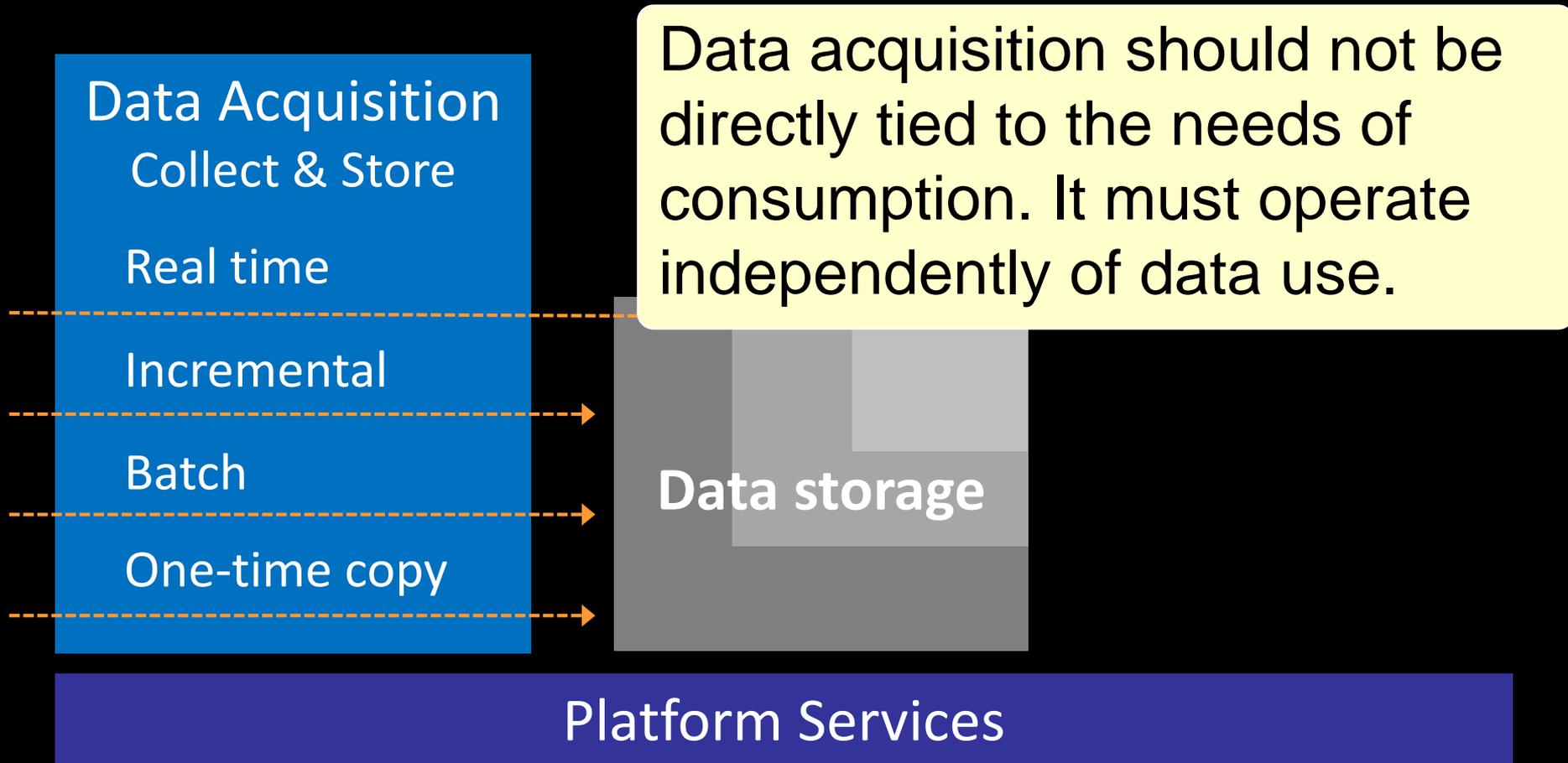
Data stor

Data management was blended with both data acquisition and structuring data for client tools. It should be an independent function.

Platform Ser

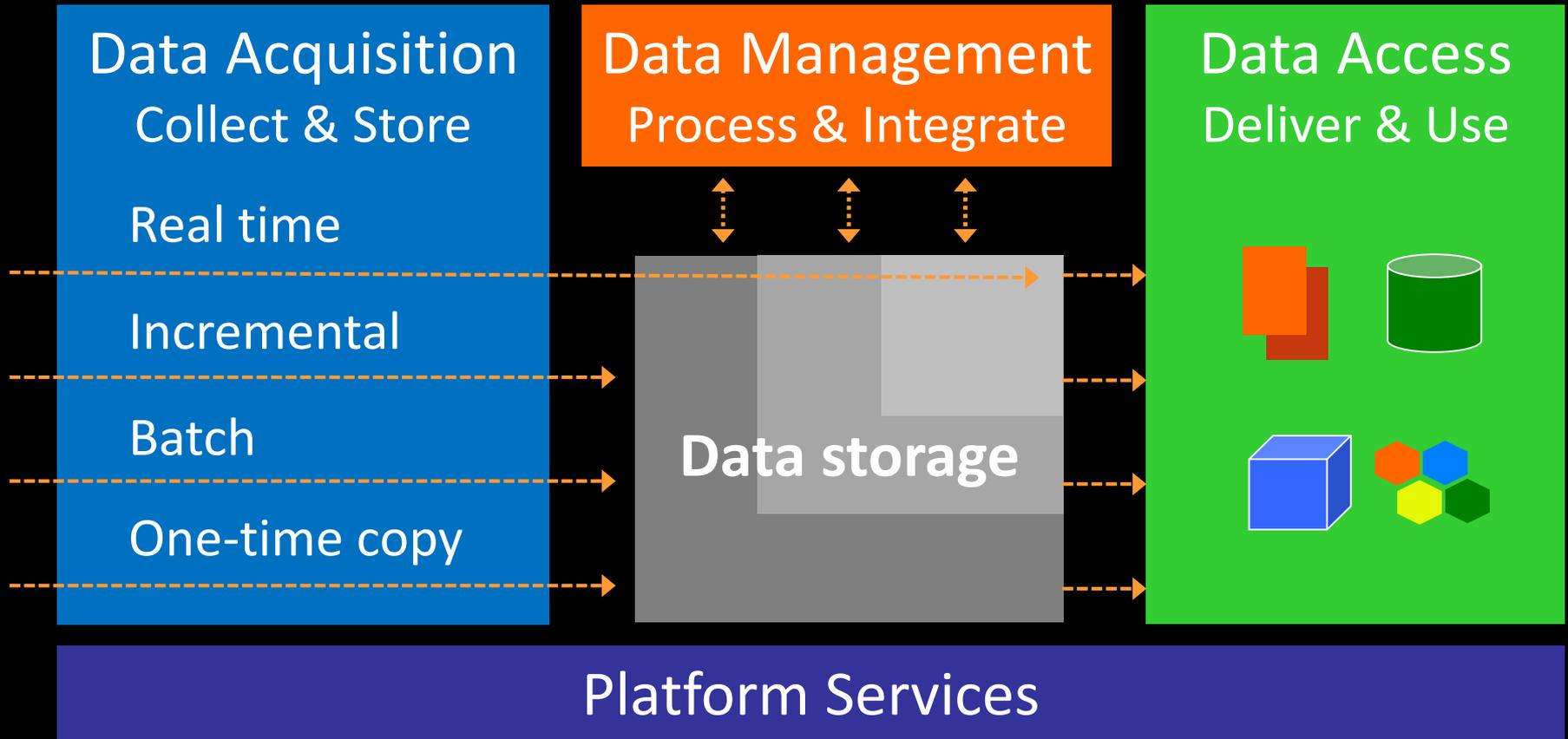
Data management should not be subject to the constraints of a single use

# The goal is to decouple: solve the application and infrastructure problems separately



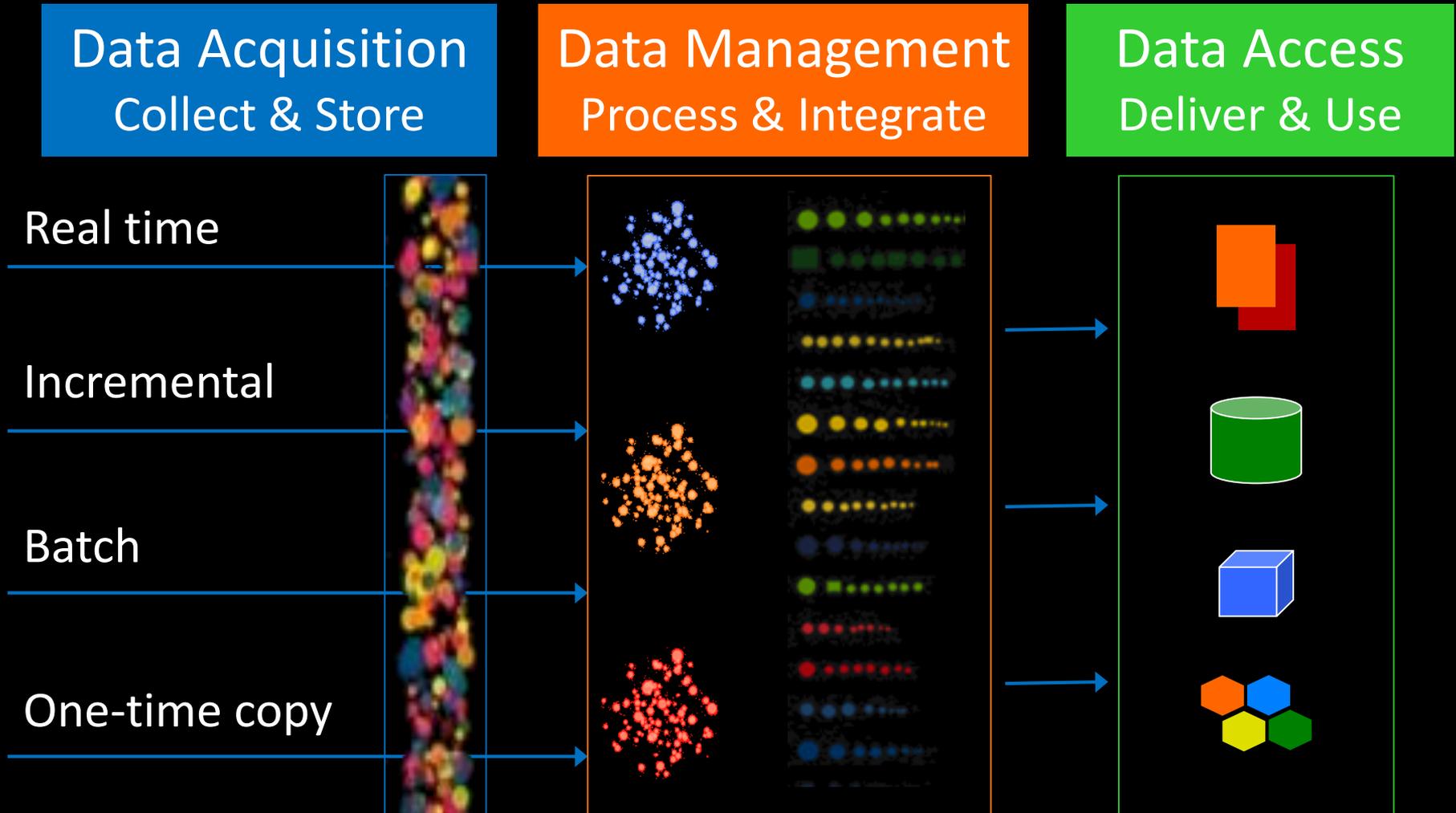
**Data arrives in many latencies, from real-time to one-time. Acquisition can't be limited by the management or consumption layers.**

# The full analytic environment subsumes all the functions of the data warehouse, and extends them



The platform has to do more than serve queries; it has to be read-write.

# The data architecture must align with system components because each of them addresses different data needs

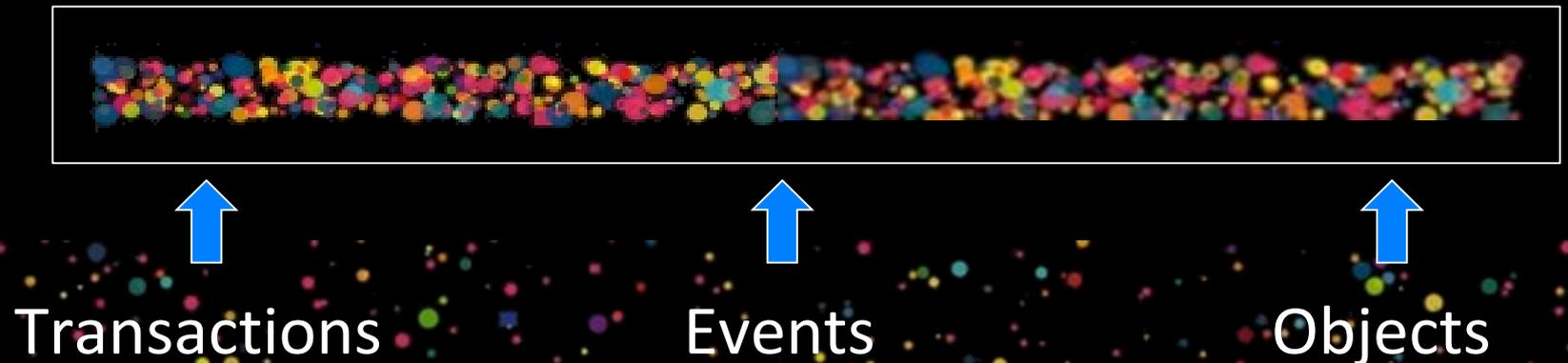


Separating concerns is part of the mechanism for change isolation

# Zone 1: Acquisition

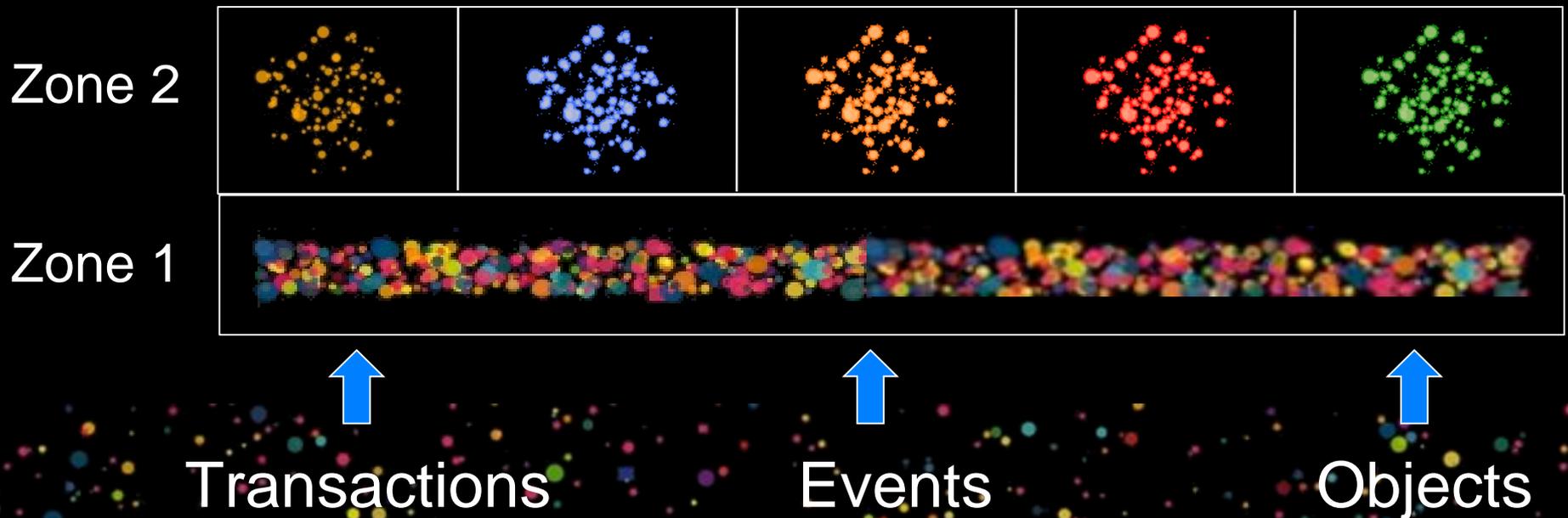
- Focus is only on collecting and tracking data
- Direct recording of data from sources
- Each dataset has as much schema as the source provides, could be explicit or implicit or none
- Foundation layer for all subsequent use

Zone 1



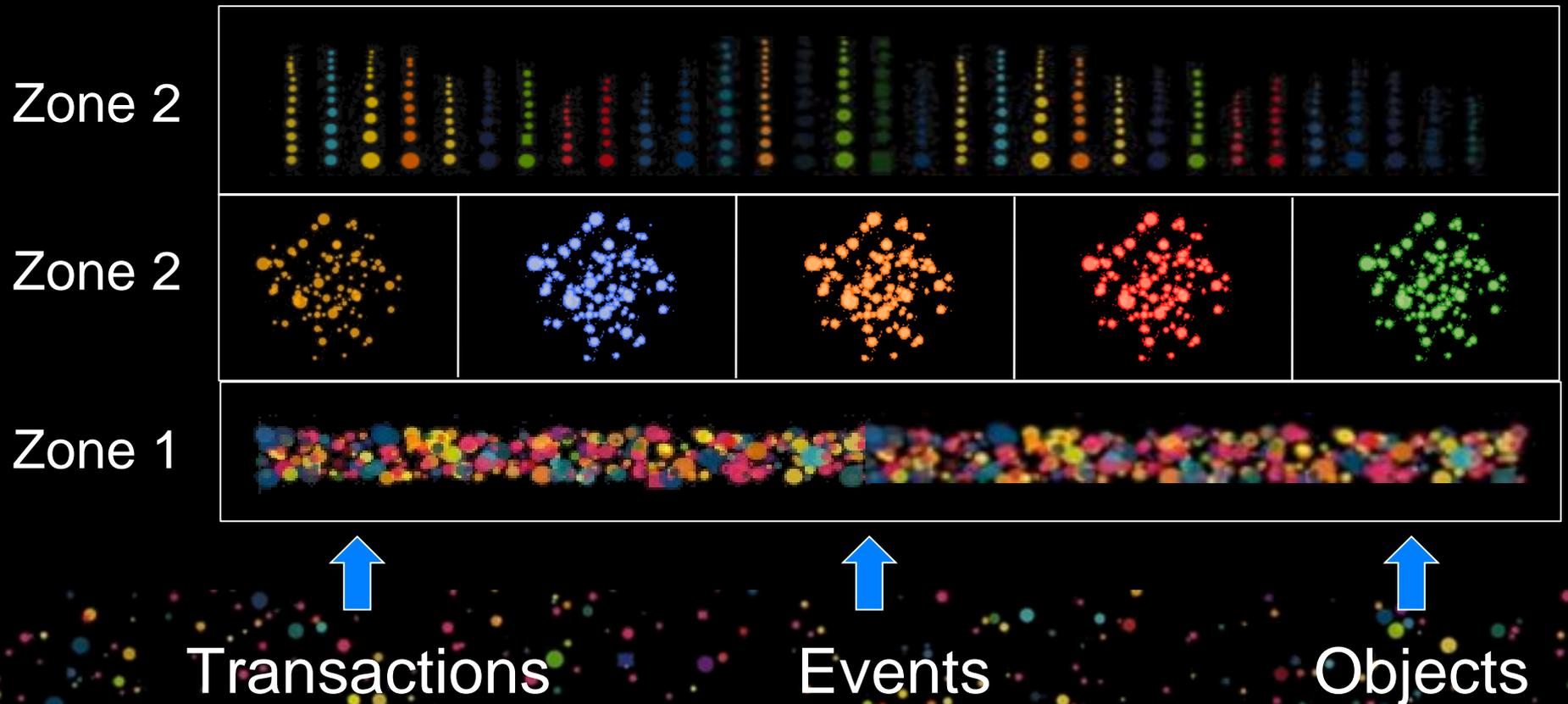
# Zone 2: Integration - Standardized

- Parsed and cataloged – all structure and form are made explicit so data can be accessed
- Namespace (e.g. keys) managed
- Common elements are standardized
- Datasets are profiled, indexed, available



# Zone 2: Integration - Enhanced

- Common shared KPIs, master datasets
- Extracted and derived data available, e.g. NEE output
- Data is linkable, labeled, possibly cleaned

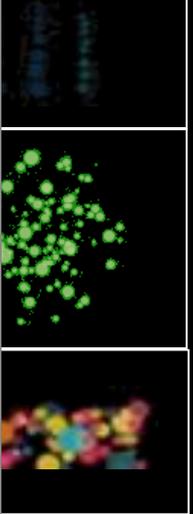


Zone 3

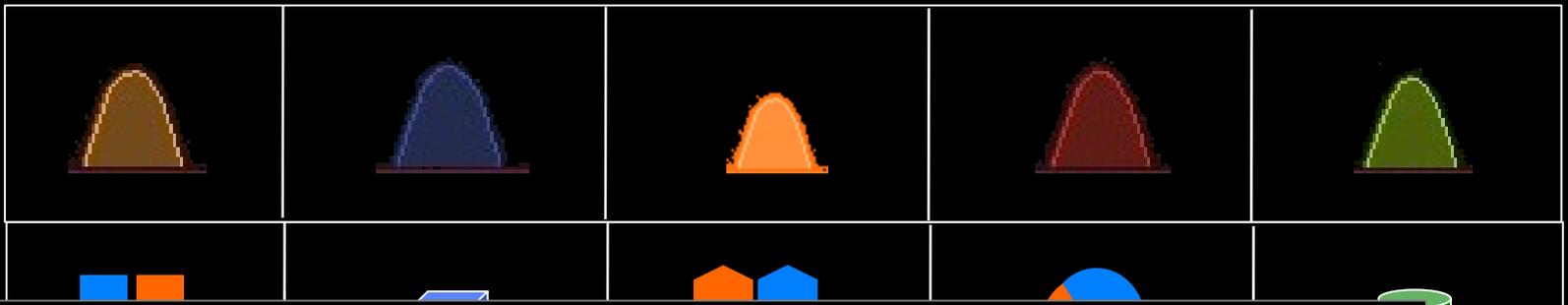


## Zone 3: Access

- Data structured to suit the workload
- Integrated data for a specific purpose
- Business rules applied, e.g. filters, controls, etc.
- Designed, explicit structures
- Generally repeatable use



Zone 4



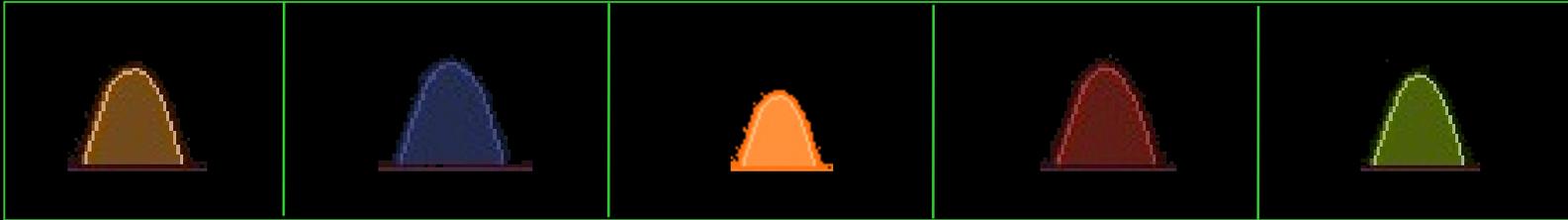
## Zone 4: Transient

- Under business / developer control
- The data for one-off projects of unknown value or repeatability, integrated from other layers
- Place for ephemeral analytics output
- The sandbox

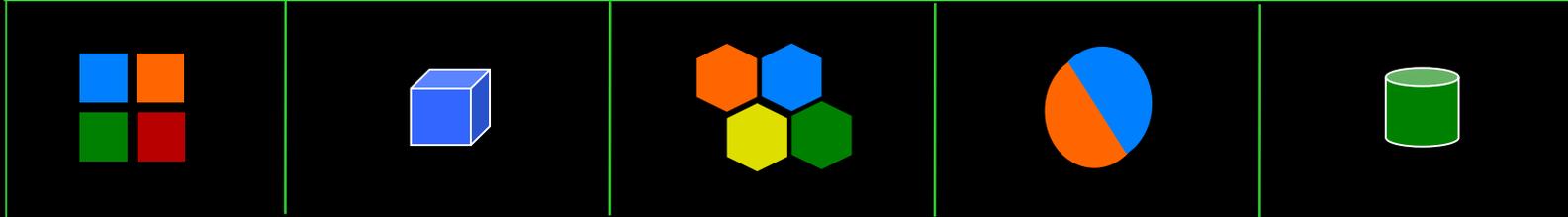


# Decoupled data architecture

Zone 4  
Transient



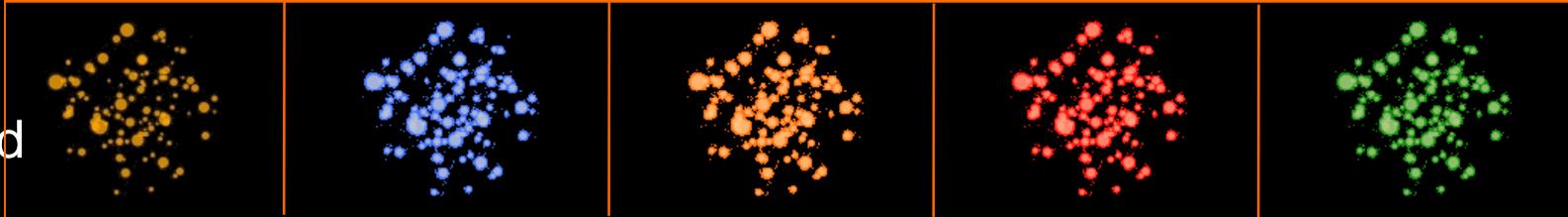
Zone 3  
Managed



Zone 2  
Enhanced



Zone 2  
Standardized



Zone 1  
Raw



# Food supply chain: an analogy for analytic data

Multiple contexts of use, differing quality levels



*You need to keep the original because just like baking, you can't unmake dough once it's mixed.*

# The design focus is different in each area



## Ingredients

Goal: available

User needs a recipe in order to make use of the data.



## Pre-mixed

Goal: discoverable and integrateable

User needs a menu to choose from the data available



## Meals

Goal: usable

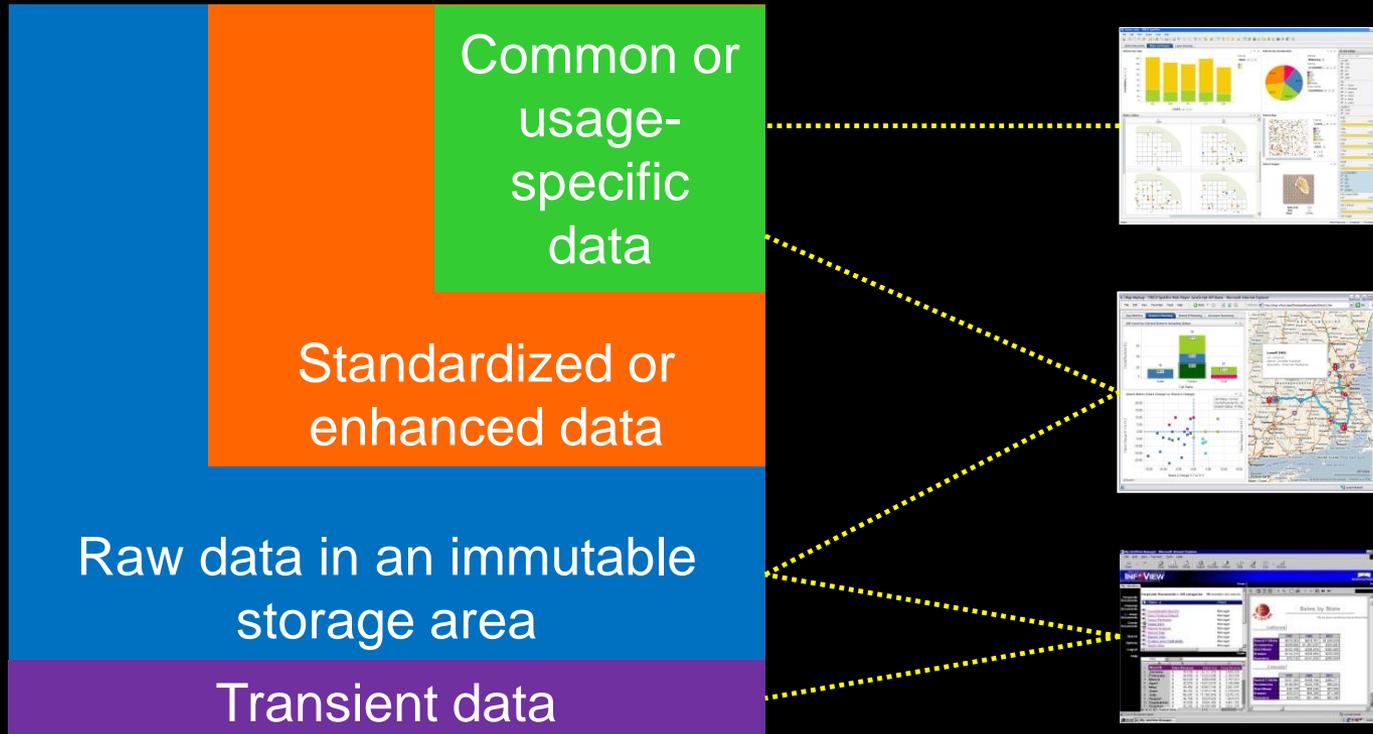
User needs utensils but is given a finished meal

# The data is in zones of management, *not* isolating layers

Relax control to enable self-service while avoiding a mess.

Do not constrain access to one zone or to a single tool.

Focus on visibility of data use, not control of data.

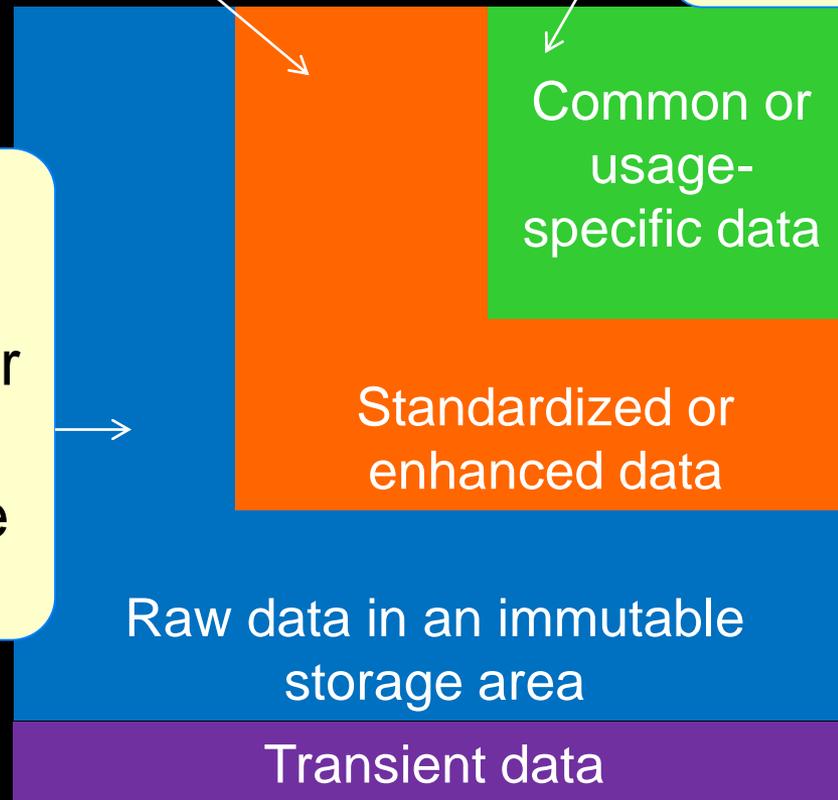


# This data architecture resolves rate of change problems

More effort applied to management, slower.

Optimized for specific uses / workloads. Generally the slowest change.

New data of unknown value, simple requests for new data can land here first, with little work by IT.

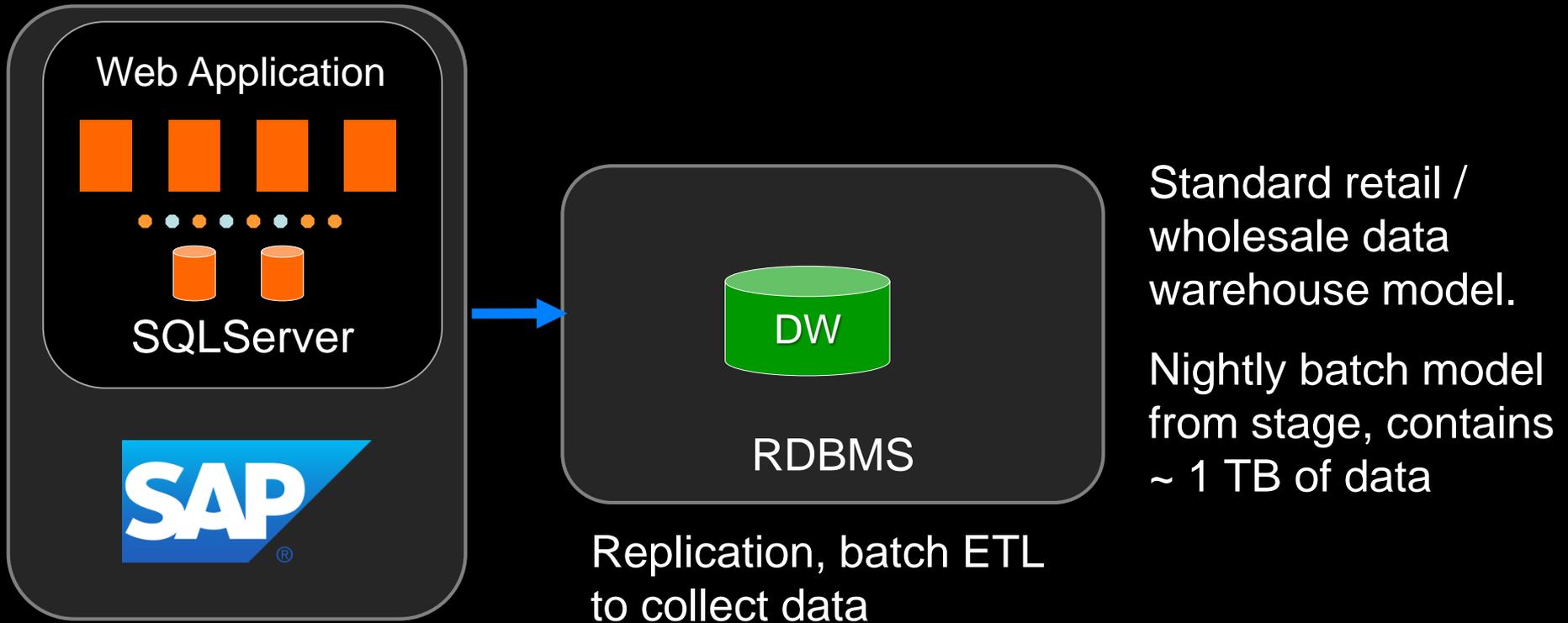


Not fast vs slow:  
*fast vs right*

Not flexibility vs control:  
*flexibility vs repeatability*

Agile for structure change  
vs agile for questions / use

# Example: data environment, mid-size retailer



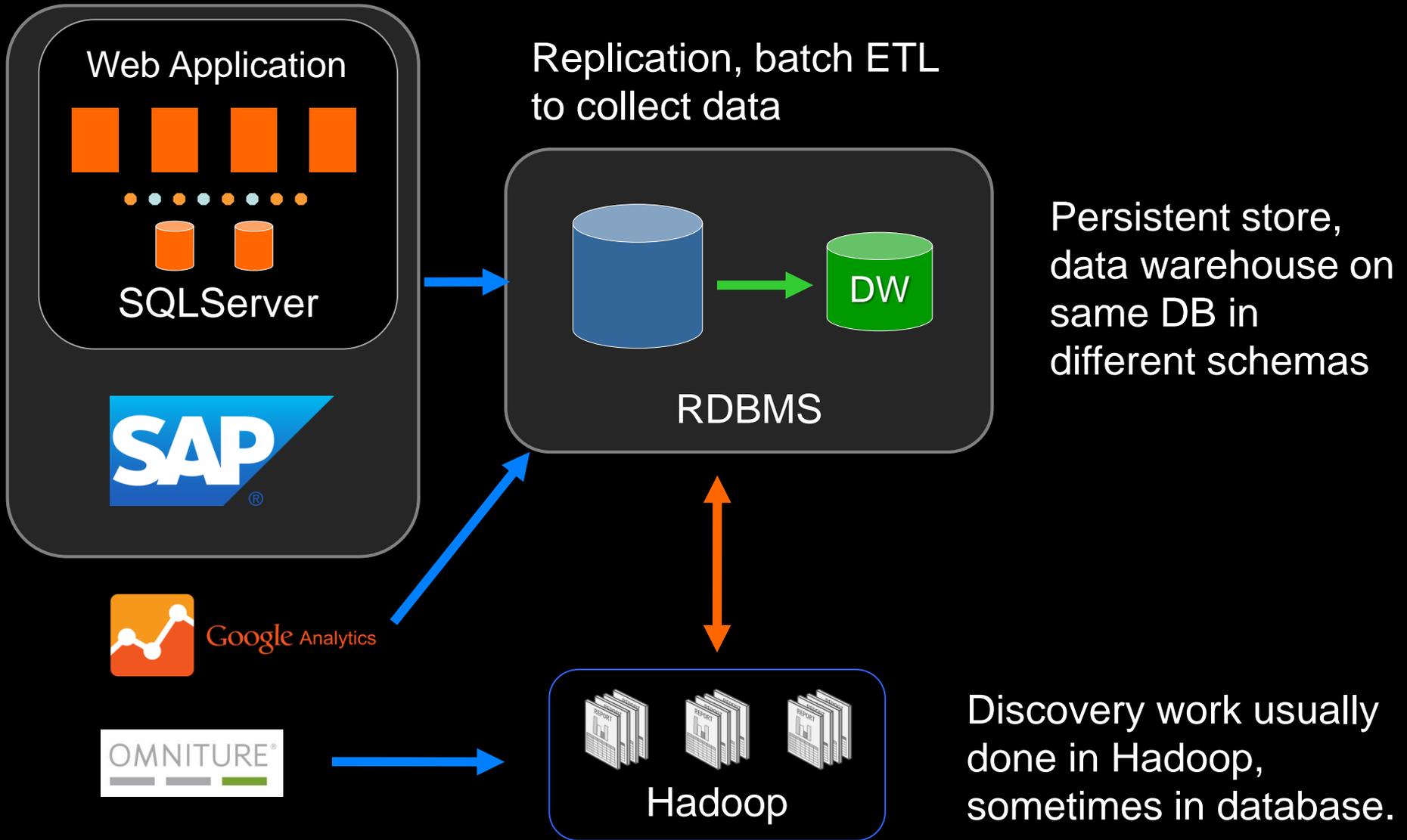
## New requirement:



Load all the online marketing and web activity for use in analytic models (not simple web analytics)

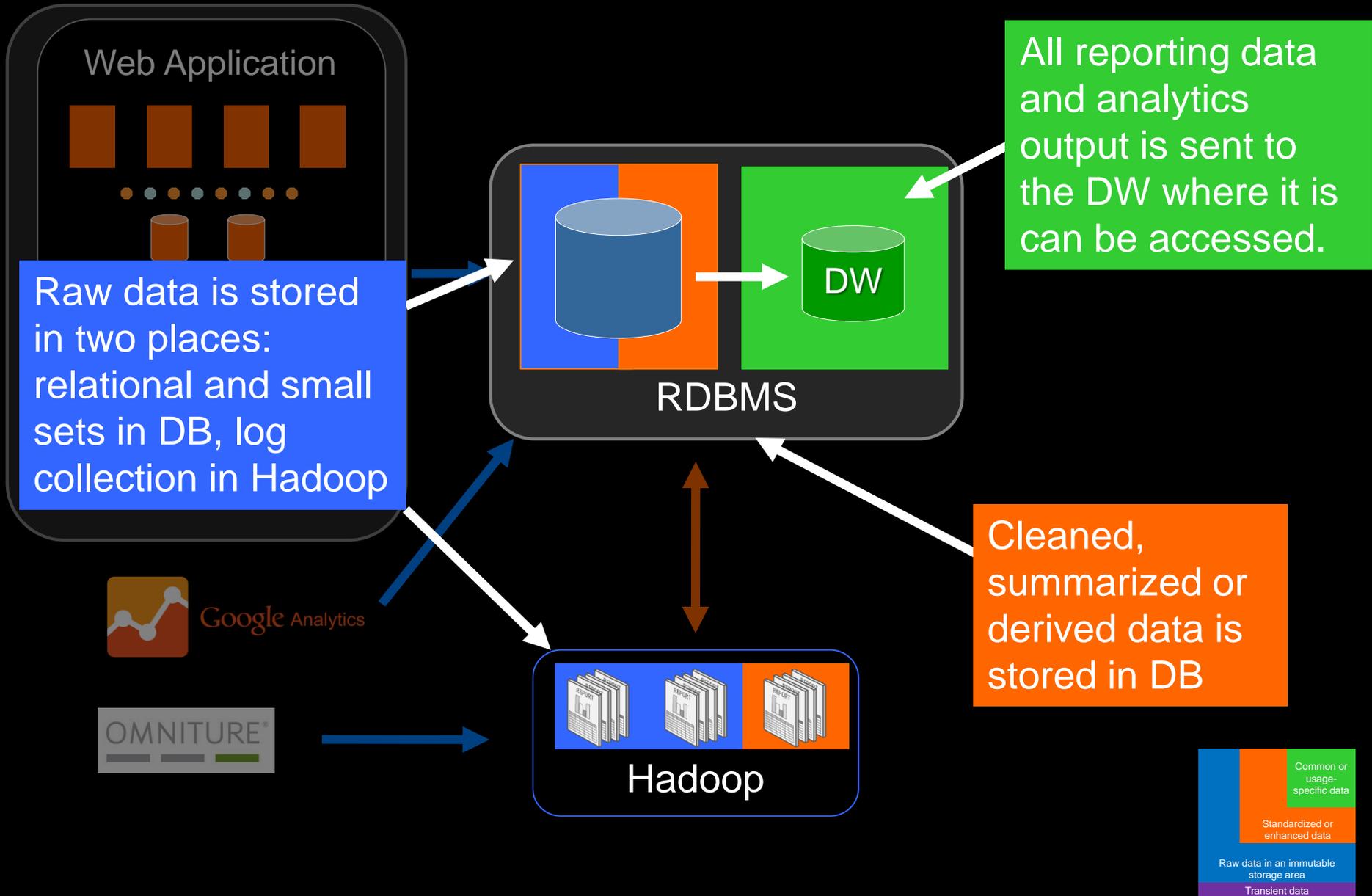
Will add ~10 TB of data, continuous or hourly loads

# Example: data environment, mid-size retailer

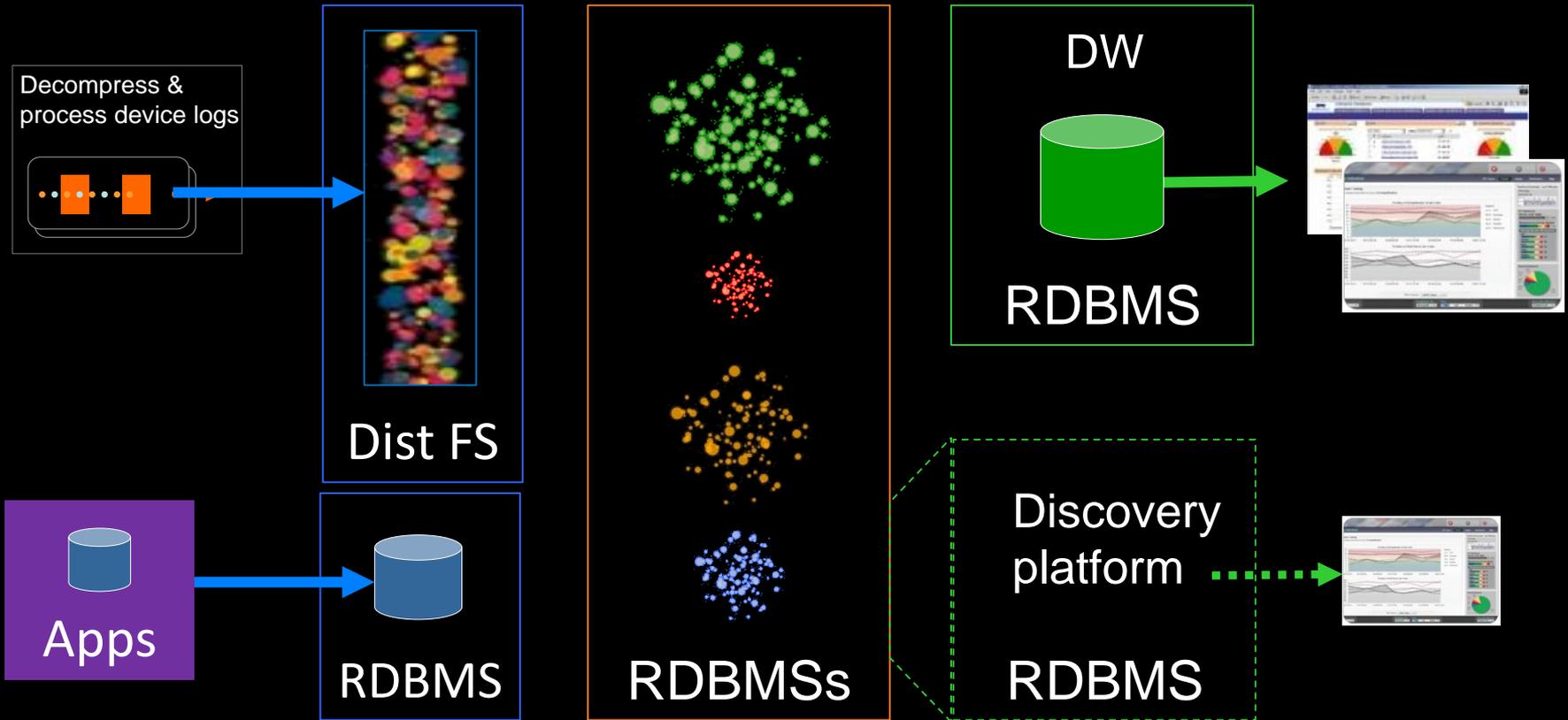


Log file fetch & load for clickstream, summaries sent to reporting env

# This data architecture uses the 3 zone pattern



# The concept of a zone is not a physical system. It's data architecture



The biggest decision is to separate all data collection from the data integration from consumption.

Physical system/technology overlays are separate, depend on the specific use cases and needs of the organization.

# What about the technology? Do I need an <X>?



# Blended Architectures Are a Requirement, Not an Option



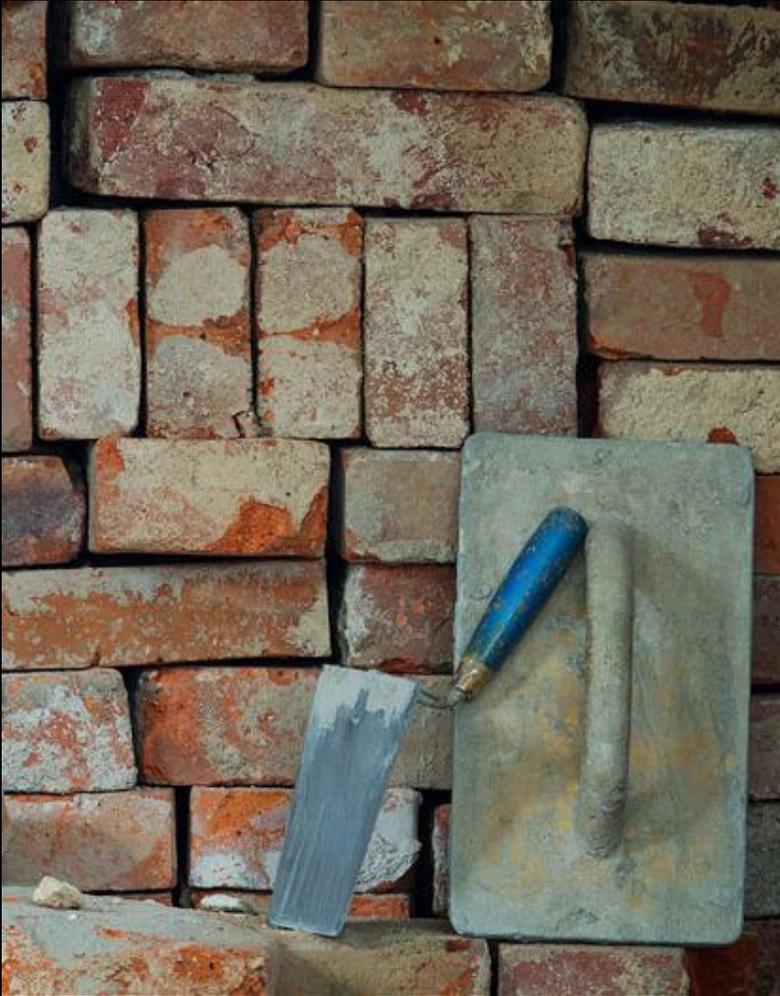
Data Warehouse + Data Lake

On Premise + Cloud

RDBMS + S3 + HDFS

Commercial + Open Source

# Bricks are not buildings



We don't think this

is equivalent to this

Architecture is not technology. It's not a product you can buy.

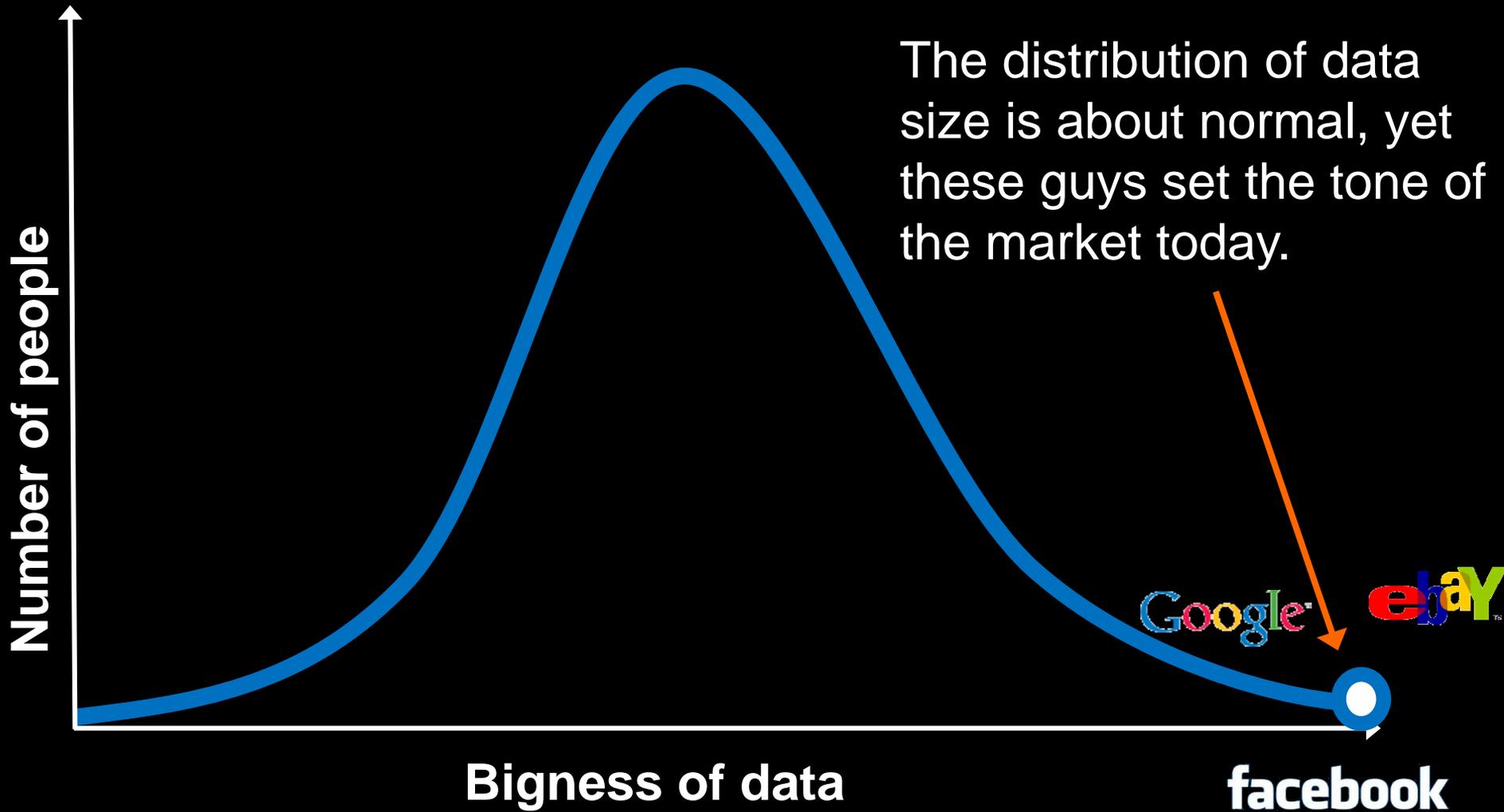
# New tools and tool architectures are required

One tool for all jobs doesn't work any more\*

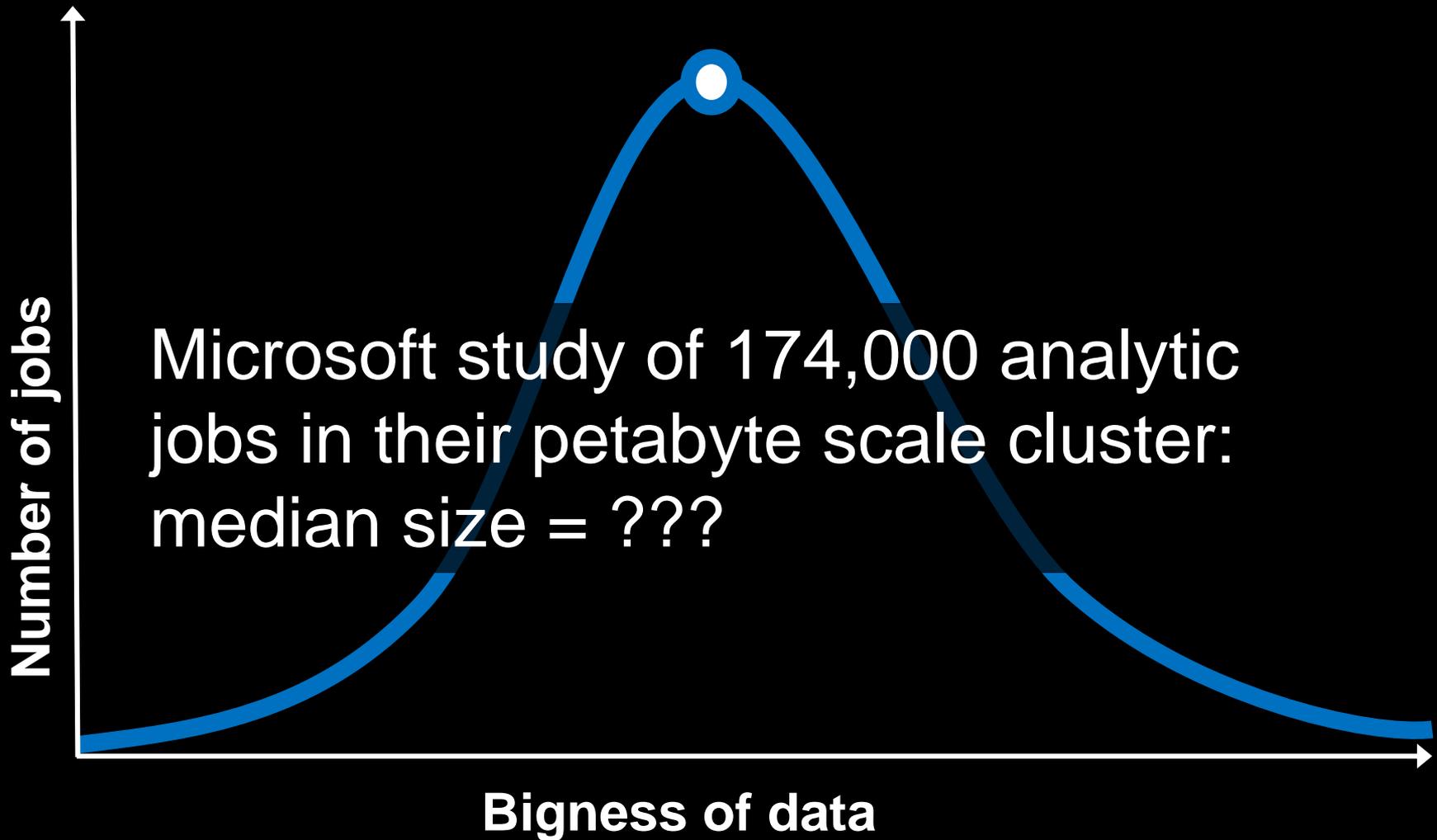


*\* it never did, we just had fewer jobs*

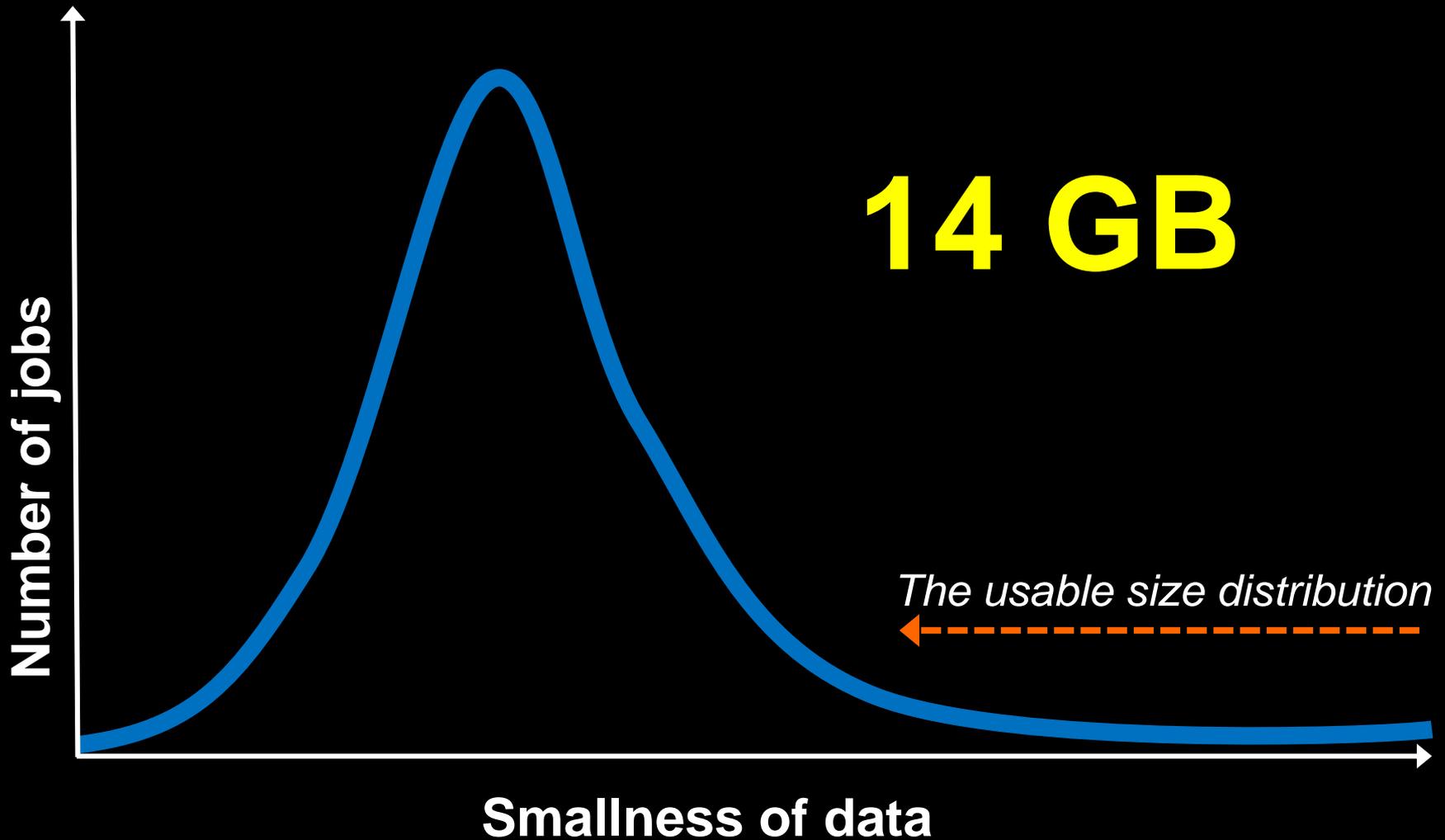
# Beware overengineering: do you really need special new technology infrastructure?



# Analytics: This is really *raw data under storage*



# Working data for analytics most often not big



**It's nice, but it'll never replace playing outside in the fresh air and getting plenty of exercise.**



## TANSTAAFL

When replacing the old with the new (or ignoring the new over the old) you always make tradeoffs, and usually you won't see them for a long time.

Technologies are not perfect replacements for one another. Often not better, only different.

# Conclusion

1. It's not about storing data, it's about using it
2. Use drives architecture. Understand the uses, what you are designing for, to drive decisions.
3. Put data at the center, not technology. Don't let the tech define what you can do or how you do it.
4. The death star is not the answer. The data model is not a flat earth. You are not building a monolith.
5. Know your history. Avoiding wheel reinvention saves time, money, careers.

## Manage your data (or it will manage you)

Data management is where developers are weakest.

Modern engineering practices are where data management is weakest.

You need to bridge these groups and practices in the organization if you want to do meaningful work with event stream data.



*In piena foresta indiana, un uomo aspetta il treno vicino alla linea ferroviaria. Improvvisamente un boa assale il malcapitato, stringendolo nelle proprie spire pericolose. Ma ecco una tigre slanciarsi a sua volta contro l'enorme rettile il quale avvolge, allora, anche la belva nella stretta mortale. Sul mostruoso groviglio sopraggiunge, trattanto, il treno. Il viaggino è spezzato sanguinosamente dalle ruote del convoglio. (Disegno di A. Bellucci)*



“Now is not the end.  
It is not even the  
beginning of the end.  
But it is, perhaps,  
the end of the beginning.”  
*Winston Churchill*

# About the Presenter

Mark Madsen is the global head of architecture at Teradata, Prior to that he was president of Third Nature, a research and consulting firm focused on analytics, data integration and data management. Mark is an award-winning author, architect and CTO whose work has been featured in numerous industry publications. Over the past ten years Mark received awards for his work from the American Productivity & Quality Center, TDWI, and the Smithsonian Institute. He is an international speaker, chairs several conferences, and is on the O'Reilly Strata program committee. For more information or to contact Mark, follow @markmadsen on Twitter or visit <http://ThirdNature.net>



# Todd Walter

## Chief Technologist - Teradata



- Chief Technologist for Teradata
- A pragmatic visionary, Walter helps business leaders, analysts and technologists better understand all of the astonishing possibilities of big data and analytics
- Works with organizations of all sizes and levels of experience at the leading edge of adopting big data, data warehouse and analytics technologies
- With Teradata for more than 30 years and served for more than 10 years as CTO of Teradata Labs, contributing significantly to Teradata's unique design features and functionality
- Holds more than a dozen Teradata patents and is a Teradata Fellow in recognition of his long record of technical innovation and contribution to the company