



Metadata Patterns Decide Who Lives and Who Dies

Delivered by **TERADATA**

Powered by **AB INITIO**



Objectives & Priorities

1. Accurately analyse data lineage and usage/queries end to end across ALL PLATFORMS



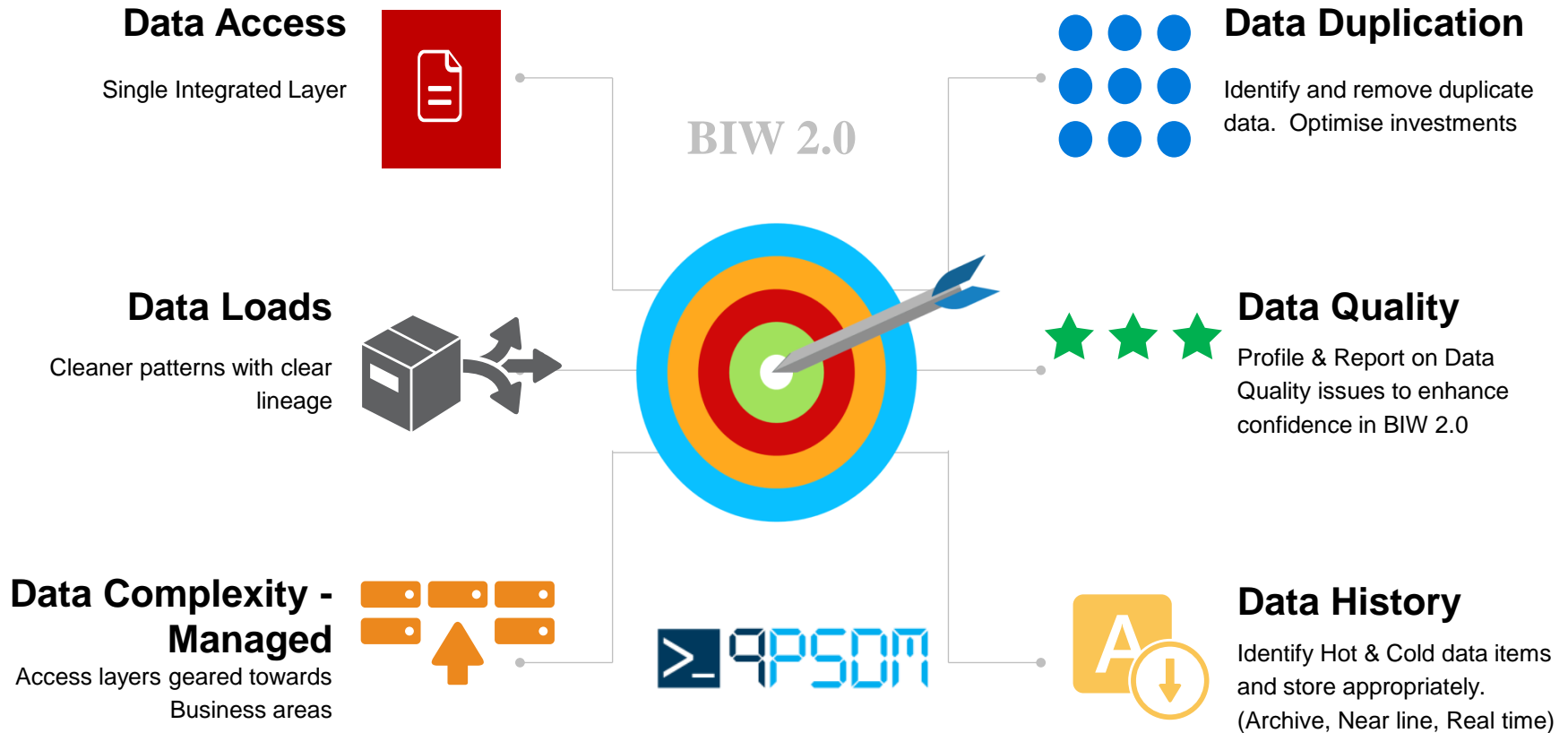
2. Use insights from (1) to prioritise
 - a. Which DBs are addressed first
 - b. Which new platform the DB goes onto
3. Automate the production of Info Packs to support discussions with DB business owners – persuade them to agree new platform decision



Two Parallel Programmes: Objectives

New Estate Build

Current Estate Cleanup



Challenges

- At Barclays, we have >2 Petabyte estate in our BI area. Like the ocean, the data evolves and changes daily

Usage patterns change – sometimes predictably, sometimes not



- With over 50 different business units using the environments – and a recent reorganisation – the need to distill and categorise quickly and succinctly was critical
- But what do you do with an ocean of data ?

Question: How do you drink a 2 Petabyte ocean?

Answer: Segment by segment



- How to get it into the buckets and quantify quickly so that we could make decisions ?
- Pattern based analysis - how does data flow and usage change over time, and how does that result in carving data up into different buckets
- We used Teradata MDE* for lineage, usage analysis and segmentation
- Barclays in-house knowledge gave context, then we agreed policy decisions for each segment
- Finally we began the process of moving to a new generation of BI at Barclays



*powered by Ab Initio



So what did we do . . .

. . . And how did we start?



Phase 1: Pilot



EXPECTATION: Take slice of ocean, ask Teradata MDE team to determine patterns and recommend segments for future action i.e. what lives (goes to BIW2.0 or Hadoop), or Queryable archive, or dies (quarantined then deleted)

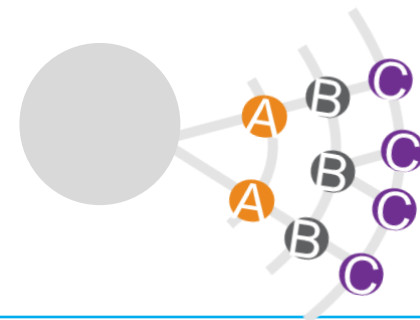


Pilot area chosen because well understood & Simple. **Only 5 Databases** ...



REALITY: MDE produced end to end data lineage to show how data fragmented across 74 Databases not 5. Some on Teradata, most on MS SQL Server, then data fed back to Teradata

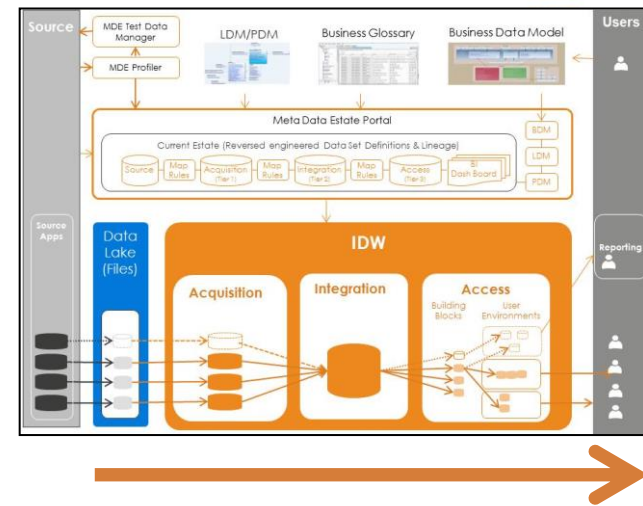
Reality: 74 Databases fragmented across different platforms



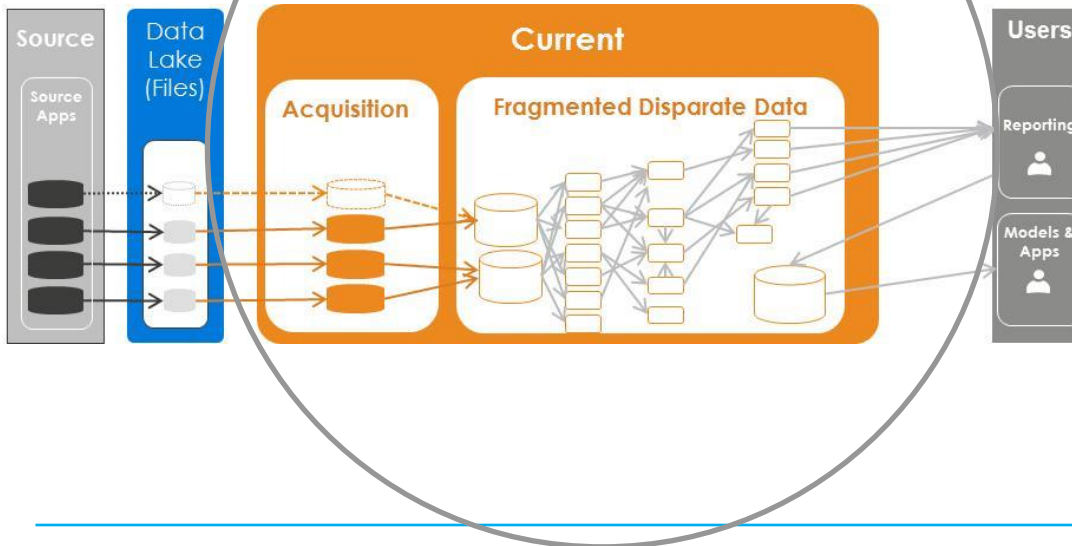
Systematic Approach:

Step A: Left To Right Ingest

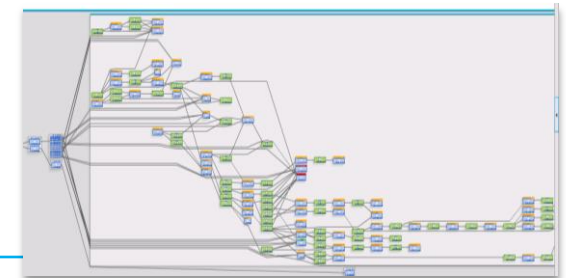
- Ingest metadata of current estate
- Derive Lineage - reveal complexity
- Identify breaks in Lineage
- Iterate



1. INGEST SELECTED DATASETS/PROCESSING



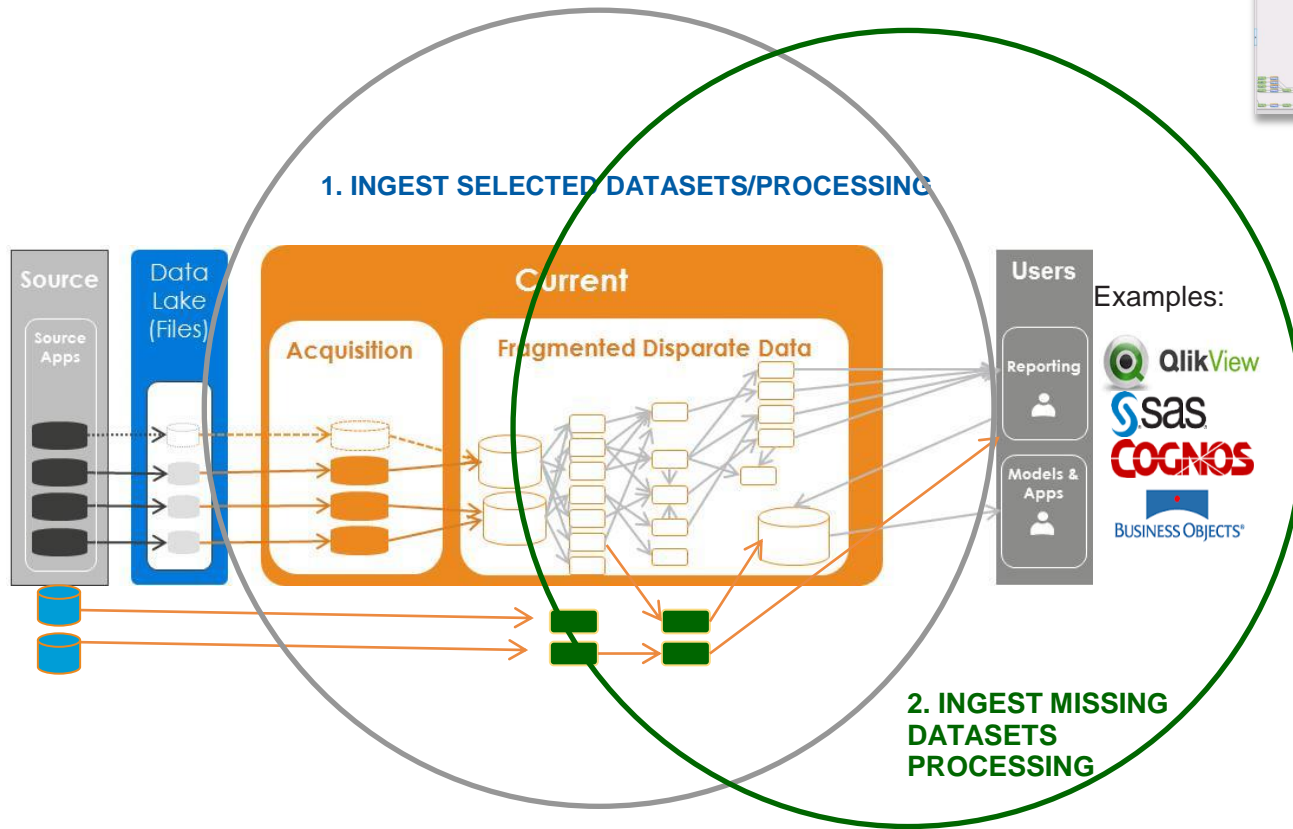
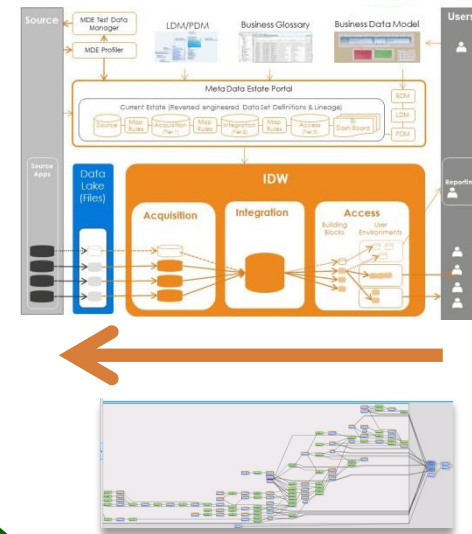
Example Lineage



Systematic Approach:

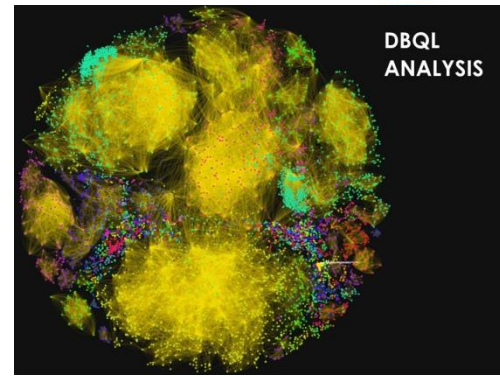
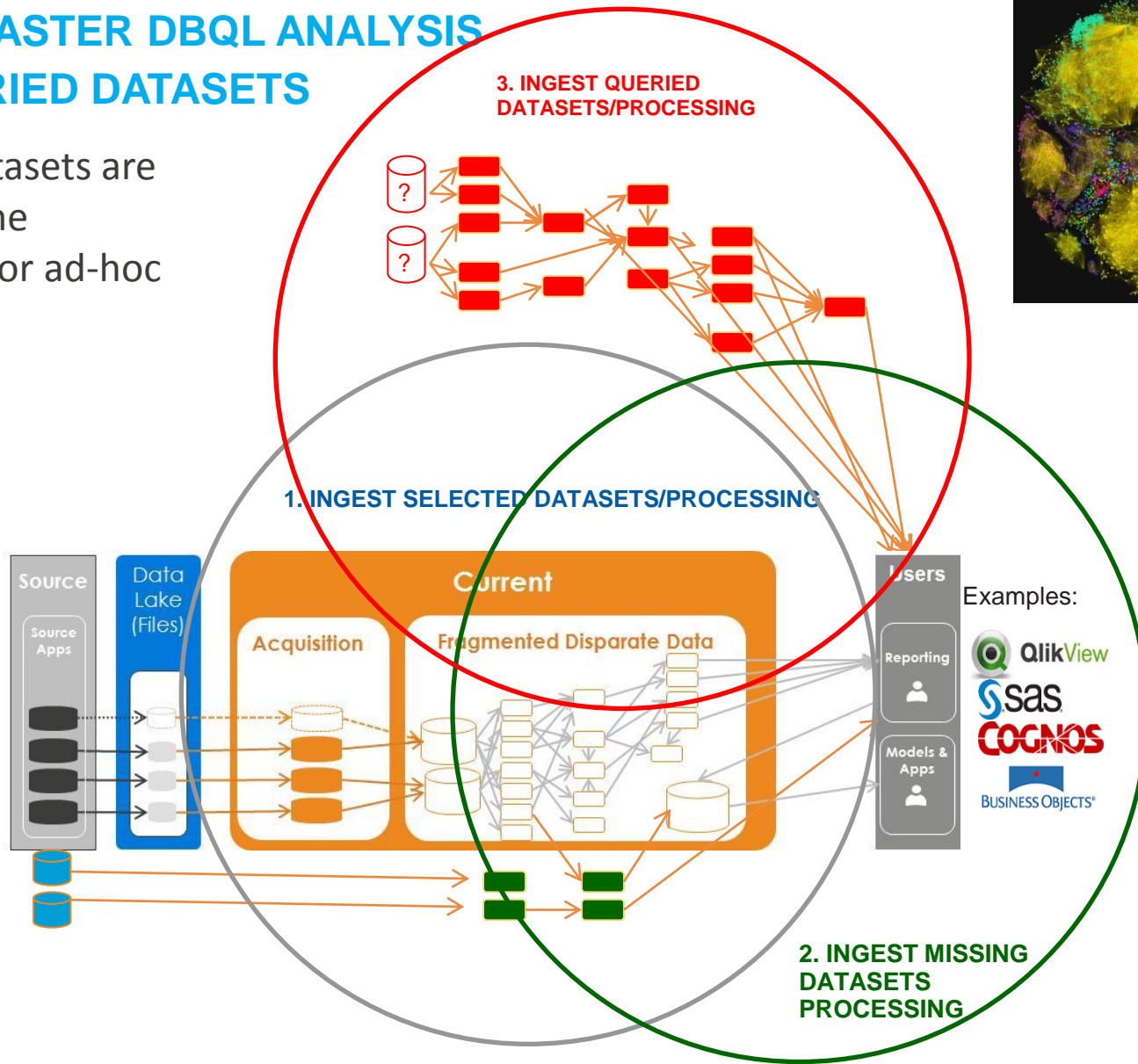
Step B: Right to left Analysis

- Formal Reports to Sources – needle in a haystack
- One business entity populated from multiple places.
- Multiple business entities populated from the same place



STEP C: ASTER DBQL ANALYSIS OF QUERIED DATASETS

Which datasets are used by the business for ad-hoc queries?



Example ASTER DBQL analysis



STEP D: PROFILE/DQ ANALYSIS

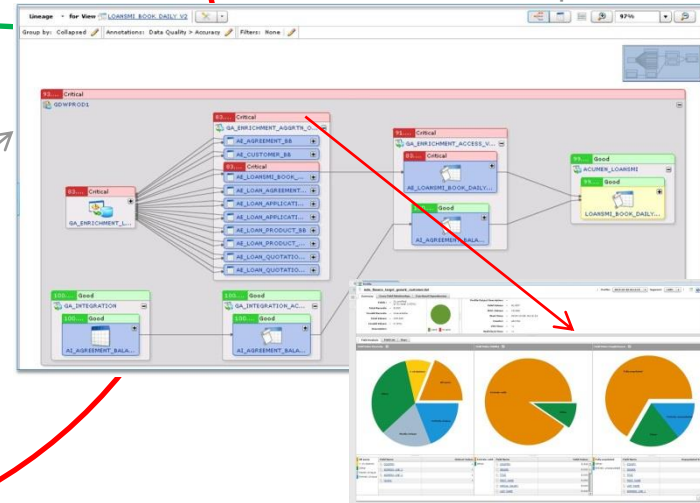
Only executed on datasets that are used and have uncertain DQ

Step C: ASTER DBQL analysis of queried datasets

3. INGEST QUERIED DATASETS/PROCESSING

Steps A, B, C: incremental ingest to populate integrated end to end Metadata

Example DQ



STEP D: PROFILE/DQ ANALYSIS (ON SELECTED DATASETS)

Step A: Left to right ingest

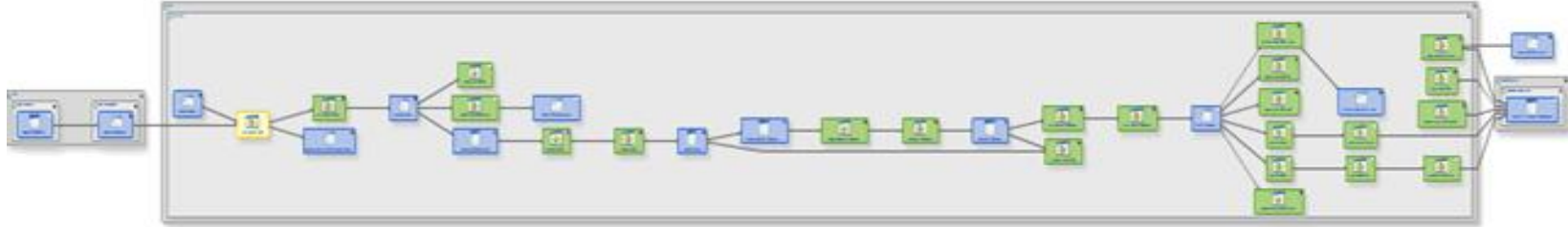
1. INGEST SELECTED DATASETS/PROCESSING

Step B: Right to left analysis

2. INGEST MISSING DATASETS/PROCESSING

Results: Source to Target (Data Lineage Samples)

Teradata → SQL Server → Teradata



Oracle → Teradata



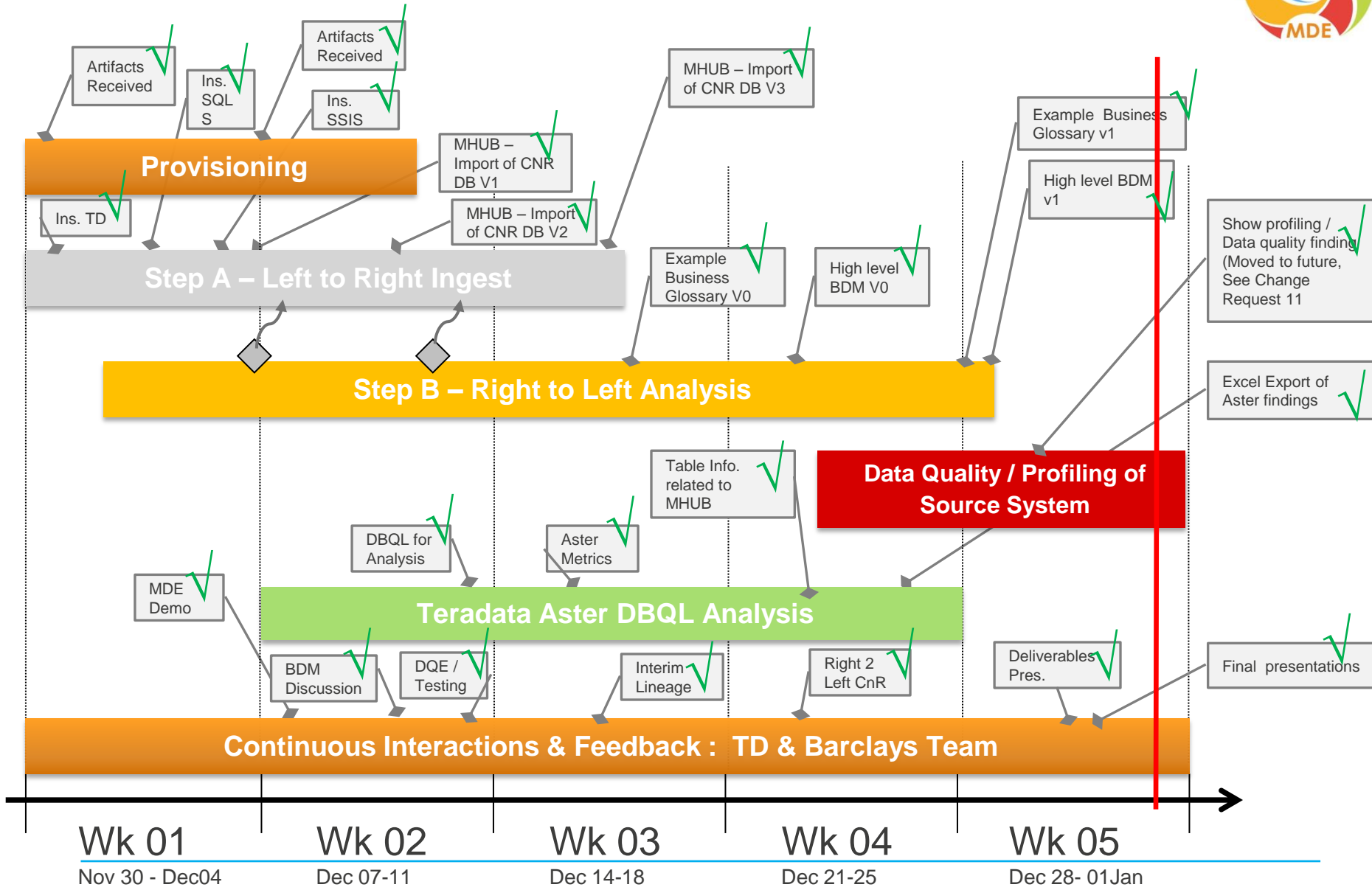
Source is a file System feed into SQL Server

IBM Mainframe → Oracle



Data Lineage	Estate Clean Up	New Estate Build
Helps to identify end to end source and targets	✓	✓
Acts as a trusted source of truth for data		✓
Capable of showing transformations & business logic details.		✓
Helps perform impact analysis of upstream/ downstream Sys.	✓	✓
Data quality overlay possible		✓

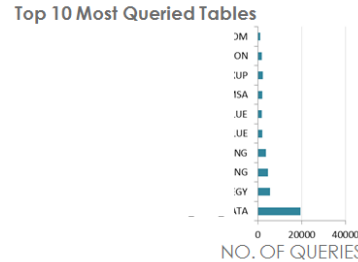
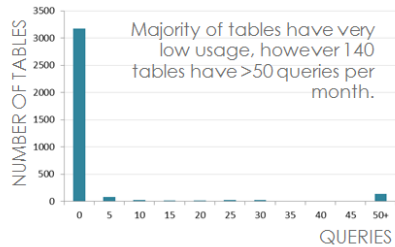
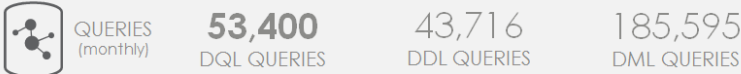
Phase I - All done in a 5 Week Timeline



Phase 1 – some example deliverables



5 DB's: DB AxY example



Table's with No Selects

Description	Category	PPT slide	Excel name
Tables with no usage in 3 mths	Data Duplication – Optimise Investments	Ref: Q 2	Tables with No Selects – no usage

Ex: Tables with No Selects

2 unique users
No selects

1.4 TB

115,223 NO QUERIES
24,218 SPACE UTILIZATION
41% (115k) Tables No Queries

objectdatabase	objectname	size_mb	total_queries	num_select	num_dcl	num_dml	num_dml_c	num_dml_n	num_sel	unique_users	free_ddl	freq_week	tbl_size	usage_type	
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated
AAA	AAA	1	0	0	0	0	0	0	0	1	0	1	1.00B33333	0.25	4 isolated

Data Lineage – different areas

Teradata → SQL Server → Teradata (Feed – Dialler)



SQL Server → Teradata (Feed - Total View)



Source is a file System feed into SQL Server

SQL Server → Teradata (Feed - BBC)



Data Lineage	Estate Clean Up	New Estate Build
Helps to identify end to end source and targets	✔	✔
Acts as a trusted source of truth for data	✔	✔
Capable of showing transformations & business logic details.	✔	✔
Helps perform impact analysis of upstream/ downstream Sys.	✔	✔
Data quality overlay possible	✔	✔

User Metrics

Description	Category	PPT slide	Excel name
All User metrics by categories	User categories	Ref: Q 3	All User Metrics

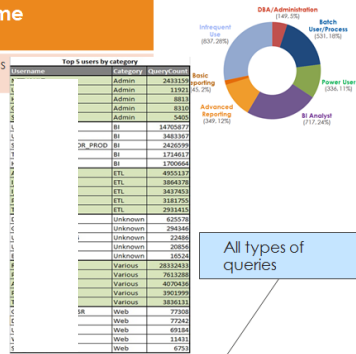
Ex: Advanced Reporting

E01909598
158 queries a day?

Table with columns: user_id, username, segment, freq_day, freq_week, avg_num_obj, avg_num_tbl, avg_num_sbo, num_sel, num_dcl, num_dml, num_dml_c, num_dml_n, num_sel, unique_users, free_ddl, freq_week, tbl_size, usage_type

Ex: BI Analysts

Table with columns: user_id, username, segment, freq_day, freq_week, avg_num_obj, avg_num_tbl, avg_num_sbo, num_sel, num_dcl, num_dml, num_dml_c, num_dml_n, num_sel, unique_users, free_ddl, freq_week, tbl_size, usage_type



Phase 2: Remaining 2 Petabytes in 5 Week Timeline

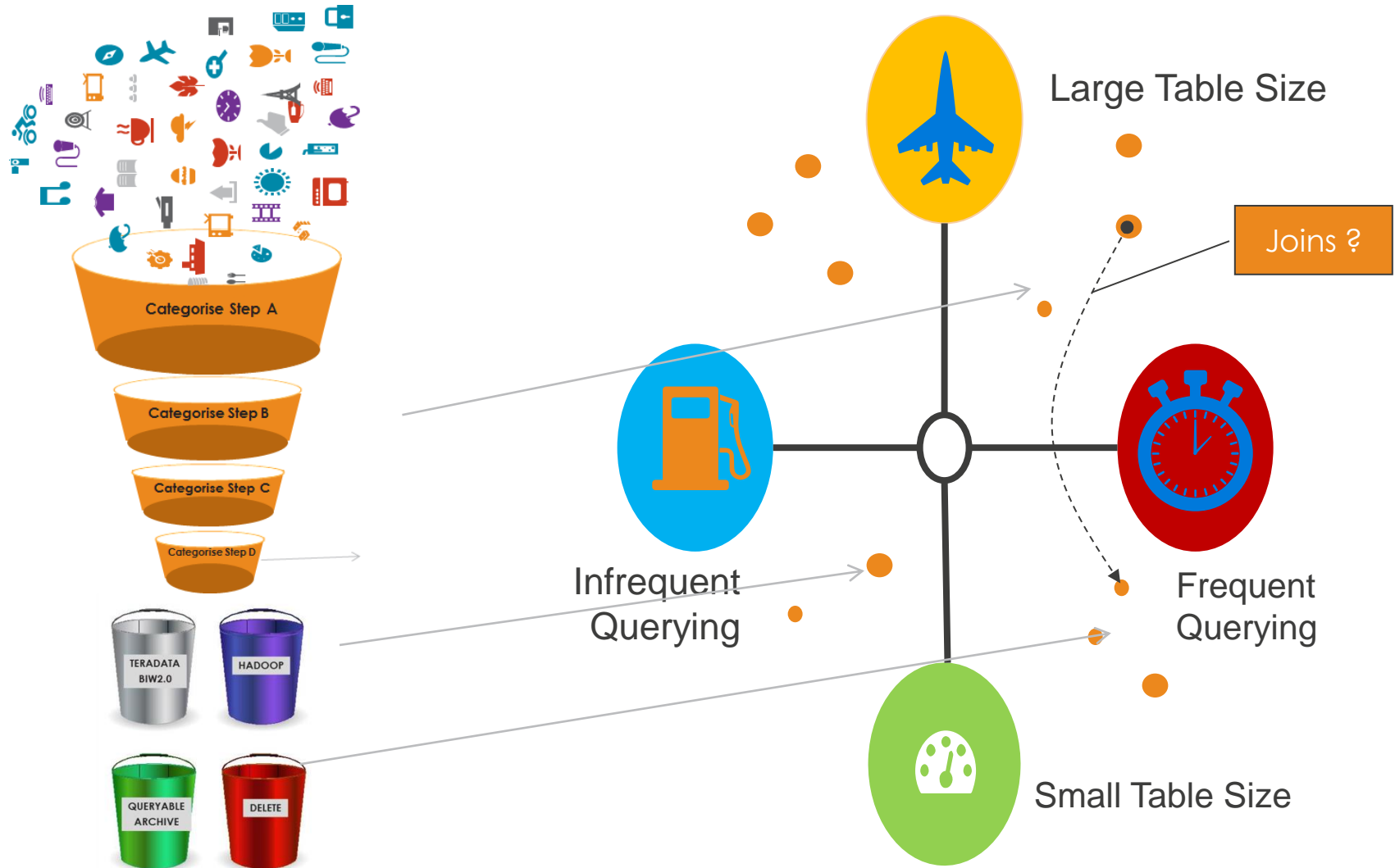


REPEATABLE PROCESS TO APPLY POLICY



CATEGORISATION / SEGMENTATION

by Area to support User Meetings / Adoption

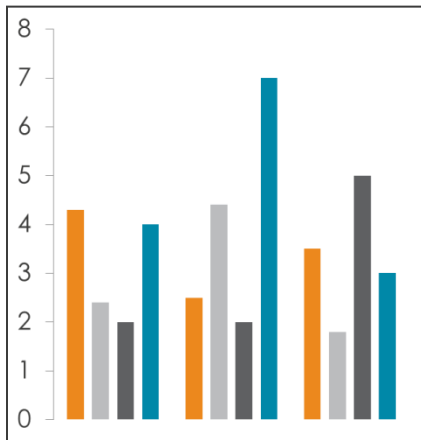


ACTION STEPS: Apply rules, Document, Review, Decide

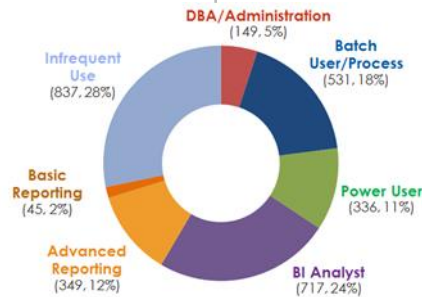
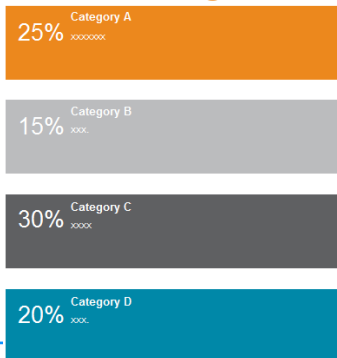
Category



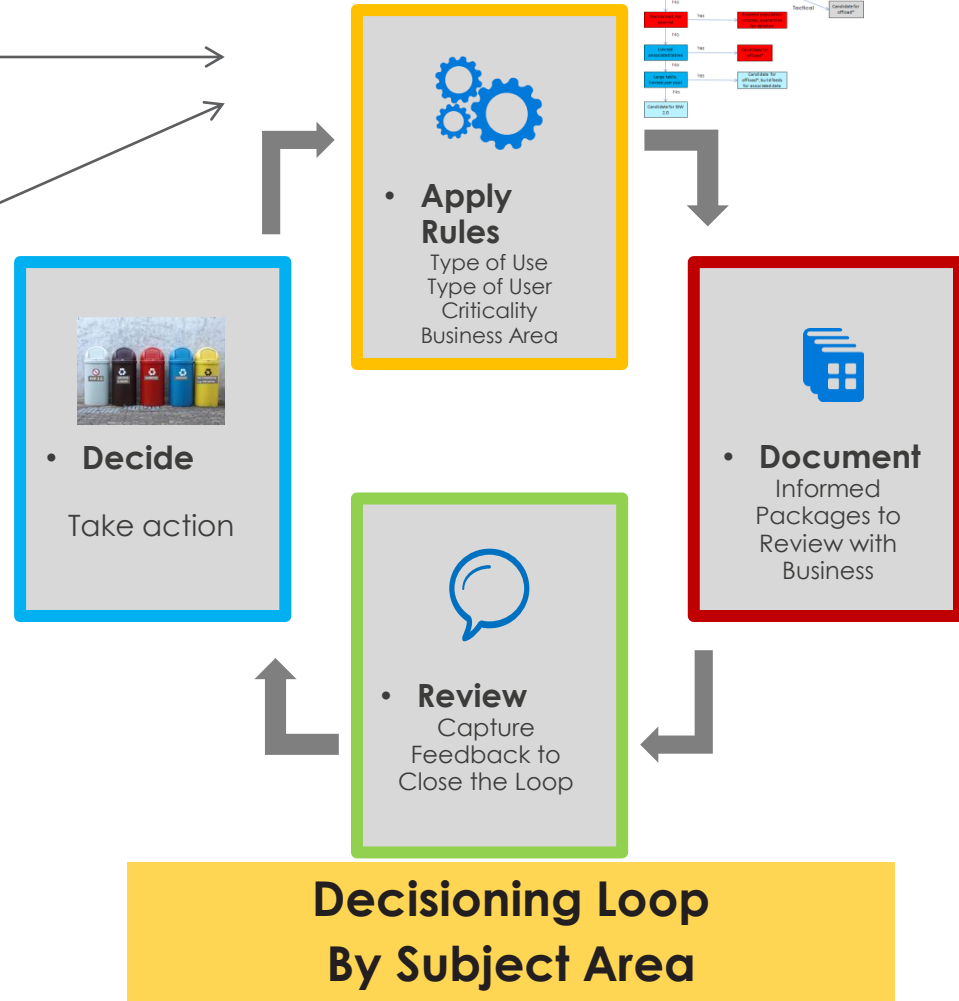
Segment



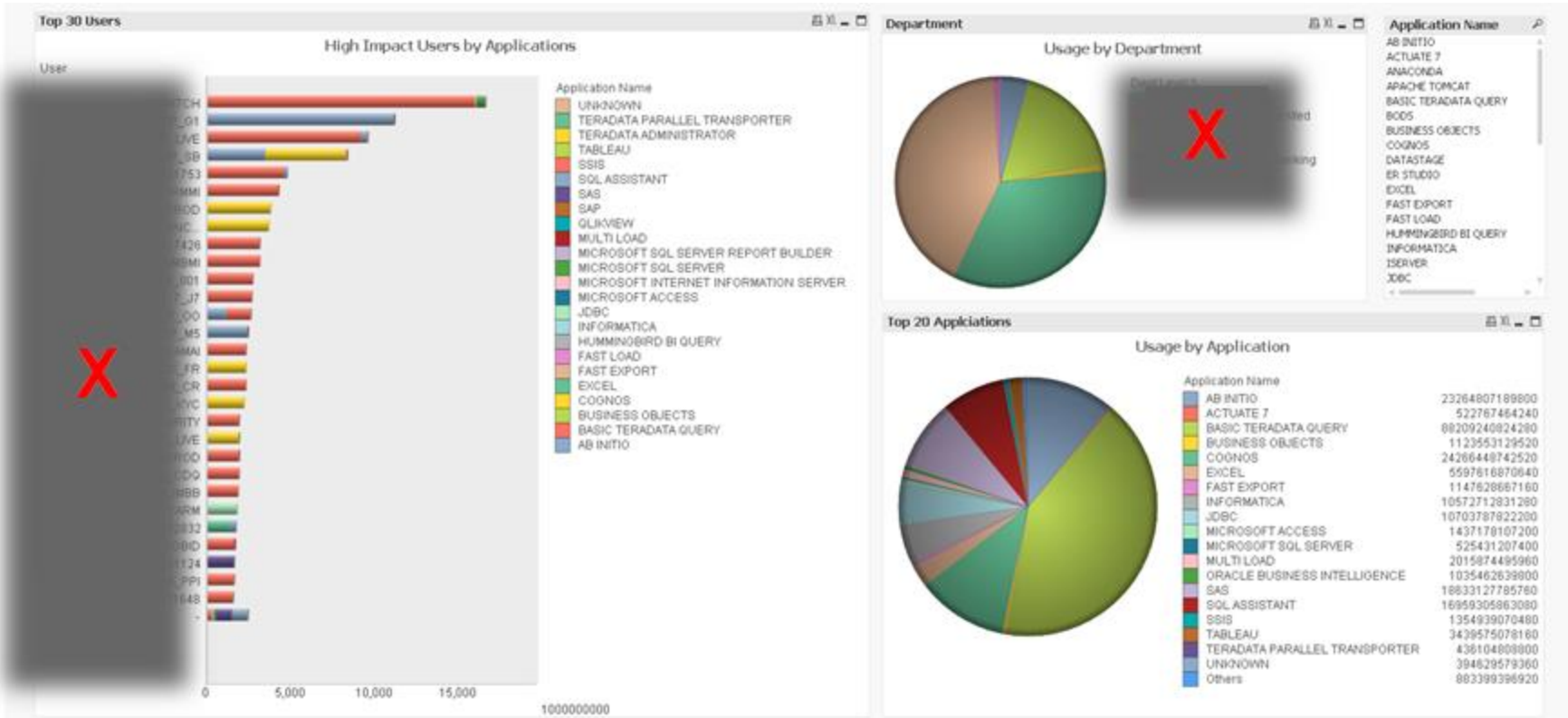
Rules by Segment



Decision Tree



Phase 2 – some example deliverables



So why does the business care?



- Accurate datasets with end to end data lineage baked-in (as required by regulator e.g. BCBS)
- Understand business usage/value of data assets
- Maintaining a “general ledger” of data assets
- “Non-productive” footprint removed to make space for new datasets
- Overnight batch reduced (not populating “non-productive” footprint)



Thank you for your time



For more information on Metadata Driven Estate visit Teradata and Ab Initio at the Ab Initio stand, or attend the Innovation Hub:

Presentation Date & Time:	Tuesday 19 April, 16:30-16:55
Presentation Zone:	Zone D
Presentation Title:	Using Metadata to Accelerate Delivery
Speakers	Elaine Fletcher, Partner, Teradata UK Professional Services Damian Worsdall, Technical Account Manager, Ab Initio Software