

Guidelines for Designing New Data Architectures



Rick F. van der Lans Industry analyst Email rick@r20.nl Twitter @rick_vanderlans www.r20.nl

Copyright © 2025 R20/Consultancy B.V., The Netherlands. All rights reserved. No part of this material may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photographic, or otherwise, without the explicit written permission of the copyright owners.

///AdeptEvents

Rick F. van der Lans



Rick F. van der Lans is a highly-respected independent analyst, consultant, author, and internationally acclaimed lecturer specializing in data warehousing, business intelligence, big data, and database technology. He is managing director of R20/Consultancy BV.

He has presented countless seminars, webinars, and keynotes at industry-leading conferences. Rick helps clients worldwide to design their data warehouse, big data, and business intelligence architectures and solutions and assists them with selecting the right products. He has been influential in introducing the new logical data warehouse architecture worldwide which helps organizations to develop more agile business intelligence systems.

He is the author of several books on computing, including his new *Data Virtualization: Selected Writings* and *Data Virtualization for Business Intelligence Systems*. Some of these books are available in different languages. Books such as the popular *Introduction to SQL* is available in English, Dutch, Italian, Chinese, and German and is sold world wide. He also authored numerous whitepapers for vendors.

In 2018 he was selected the sixth most influential BI analyst worldwide by onalytica.com.

You can get in touch with Rick van der Lans via:Email:rick@r20.nlWebsite:www.r20.nlLinkedIn:http://www.linkedin.com/pub/rick-van-der-lans/9/207/223

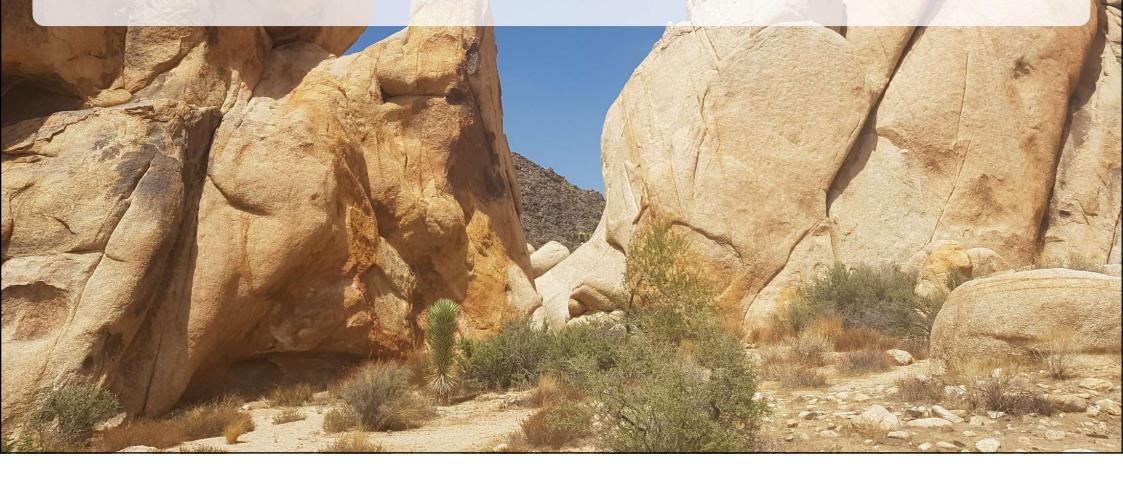


Agenda and Subjects



- **1. Introduction: Why a New Data Architecture?**
- 2. Introduction to Data Architectures
- 3. Steps 1-3: Setting the Stage
- 4. Step 4: Analyze New Technologies for Data Storage, Processing, and Analytics
- 5. Step 5: Architectural Design Principles
- 6. Step 6: Reference Data Architectures
- 7. Step 7-8: Designing New Data Architectures
- 8. Steps 9-10: Final Steps
- 9. Closing Remarks

Part 1: Introduction: Why a New Data Architecture?





600000

Digital transformations are even more difficult than traditional change efforts to pull off. But the results from the most effective transformations point to five factors for success.



Home About the ELI News & Event

Projects > Data Economy

Principles for a Data Economy (with the ALI)

With the rise of an economy in which data is a tradeable asset globally, more certainty is needed with regard to the legal rules that are applicable to the transactions in which data is an asset. Critical questions arise such as who has which right with regard to the data generated by connected devices? They need to be answered urgently, as lack of clarity in this field potentially hinders innovation and growth and, more importantly, troubles consumers, data-driven industries, and start-ups.

For an overview of past and upcoming meetings of this project, please click here.

ELI Members, who are interested in actively contributing to the development of this project are invited to

Data is a business asset beyond imagination – here is why (and where)

It has almost become embarrasing to say that data is a business asset and should be treated as one *(the same goes for information).*

The 'data is an asset' or a 'data is a business asset' message is not new. It goes back over two decades. However, despite the fact that so many people have said it so often before, we still see that there is a difference between preach and practice.

It's not that organizations fail to understand the data, information and actionable intelligence (

4 Steps to Help You Become a Data Driven Company in 2019

CLEMENT RENÉ

Media Monitoring - Social Media - Agencies

3 Ways Your Company Should Be Like Google, and 1 Way it Definitely Shouldn't

TECHNOLOGY

Inc.

Trends

Digital Marketing

INC. 5000 CONFEREN

You don't have to be the world's most popular search engine in order to innovate. Here are four things you can learn from contact us BLOG Q =

Popular Posts

Future Horizons in Review - 5 Tech

Conversations That Ruled in 2018

What to Expect from Immersive

Data Disruption is Here

in Future horizons posted on January 11, 2018 by Mike Parsons @3 minutes Ø0 Comments

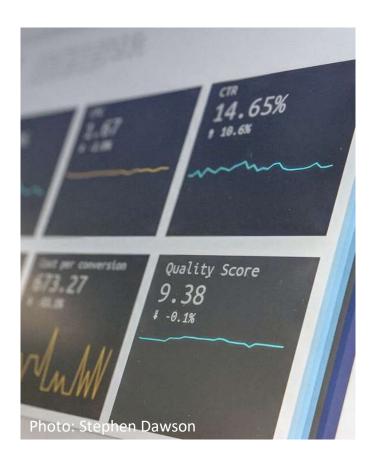
If software is eating the world, then data is swallowing it. The four highest valued companies in the US are tech heavyweights Apple, Google, Microsoft, and Amazon. Facebook is close behind in 8th position.

QUALITANCE

Most Popular

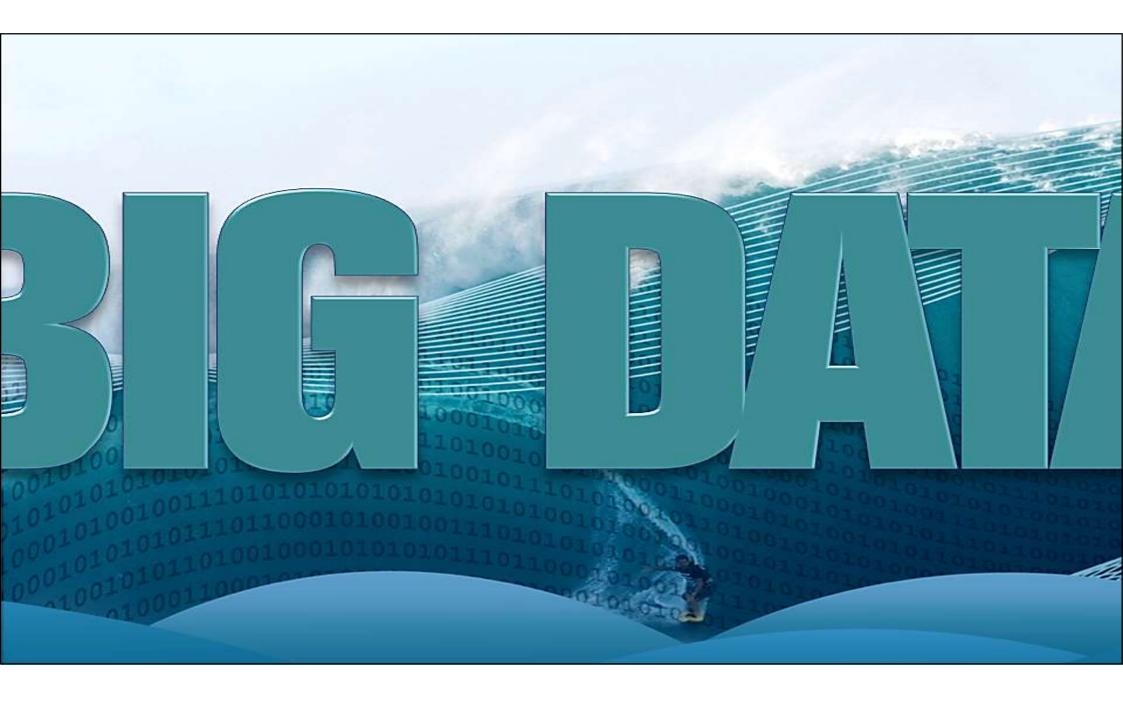
Technologies in 2019 Top 5 Publicly Traded Companies (by Market Cap) Tech Other 2006 EXON 54468 53338 Torxt 53278 5278 Collection 52788 53278 53278 11 Important Truths About How 11 Important Truths About How

Examples of 'Doing More' with Data



- Enabling self-service reporting, dashboards, and for analytics to work with (near) real-time data
- Combining internal with external data to enrich analytical capabilities
- Accelerating AI/machine learning initiatives to discover new patterns or trends in the data
- Simplifying deployment of IoT technology to generate more data on the machines and business processes
- Offering edge analytics in which real-time data is analyzed continuously and near the place where the data is produced by the sensor or business process

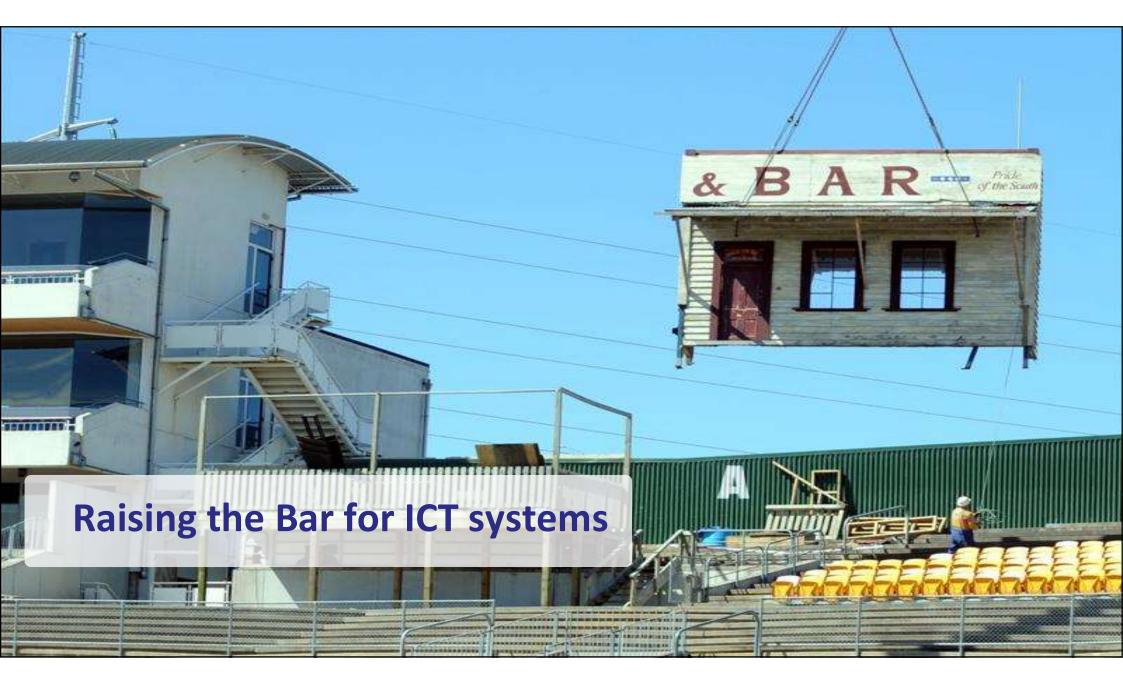




Data hasn't changed, it's just more of the same

Data consumption has changed

Self-service BI Embedded BI Supplier- and Customer-driven BI Applied AI in Text, Image, Video Analysis Edge Analytics Data Marketplace Data Science Automated decisions



How Good is our Track Record? Photo: Samuelk Blanck

Software Development Hall of Shame

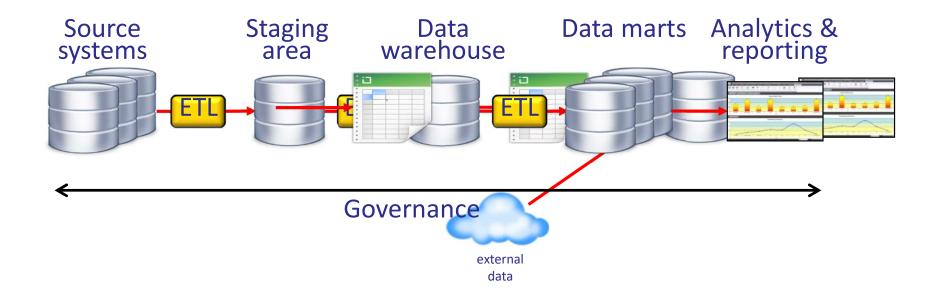
YEAR	COMPANY	OUTCOME (COSTS IN US \$)				
2005	Hudson Bay Co. [Canada]	Problems with inventory system contribute to \$33.3 million* loss.				
2004-05	UK Inland Revenue	Software errors contribute to \$3.45 billion* tax-credit overpayment.				
2004	Avis Europe PLC [UK]	Enterprise resource planning (ERP) system canceled after \$54.5 million [†] is spent.				
2004	Ford Motor Co.	Purchasing system abandoned after deployment costing approximately \$400 million.				
2004	J Sainsbury PLC [UK]	Supply-chain management system abandoned after deployment costing \$527 million.*				
2004	Hewlett-Packard Co.	Problems with ERP system contribute to \$160 million loss.				
2003-04	AT&T Wireless	Customer relations management (CRM) upgrade problems lead to revenue loss of \$100 million.				
2002	McDonald's Corp.	The Innovate information-purchasing system canceled after \$170 million is spent.				
2002	Sydney Water Corp. [Australia]	Billing system canceled after \$33.2 million [†] is spent.				
2002	CIGNA Corp.	Problems with CRM system contribute to \$445 million loss.				
2001	Nike Inc.	Problems with supply-chain management system contribute to \$100 million loss.				
2001	Kmart Corp.	Supply-chain management system canceled after \$130 million is spent.				
2000	Washington, D.C.	City payroll system abandoned after deployment costing \$25 million.				
1999	United Way	Administrative processing system canceled after \$12 million is spent.				
1999	State of Mississippi	Tax system canceled after \$11.2 million is spent; state receives \$185 million damages.				
1999	Hershey Foods Corp.	Problems with ERP system contribute to \$151 million loss.				
1998	Snap-on Inc.	Problems with order-entry system contribute to revenue loss of \$50 million.				
1997	U.S. Internal Revenue Service	Tax modernization effort canceled after \$4 billion is spent.				
1997	State of Washington	Department of Motor Vehicle (DMV) system canceled after \$40 million is spent.				
1997	Oxford Health Plans Inc.	Billing and claims system problems contribute to quarterly loss; stock plummets, leading to \$3.4 billion loss in corporate value				

Enkele Mislukte ICT-Projecten Bij NL Overheid

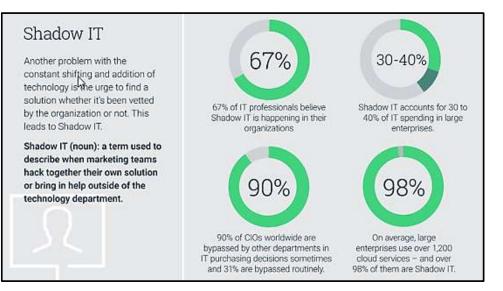
proj wor	ecte	en: e	haa stru		e	
ICT Ove			orojecter wee bek			

- Minister Carola Schouten (Landbouw, CU) [besluit April 2019] te stoppen met de vernieuwing van de ICT-systemen bij de Nederlandse Voedsel- en Warenautoriteit. Na 65 miljoen euro te hebben uitgegeven, bleek er maar weinig te werken.
- De Belastingdienst, dat sinds 2005 probeerde een systeem te bouwen dat alle transacties van de fiscus zou verwerken. Na negen jaar en 203 miljoen euro gaven ze het op.
- Defensie, waar ze sinds 2002 bouwden aan 'Speer'. Na volgens eigen zeggen 433 miljoen euro te hebben uitgegeven – de Algemene Rekenkamer kwam op 900 miljoen euro uit – gaf het ministerie het in 2015 op. Speer was nog lang niet af, en werkte niet zoals bedoeld.
- Het nieuwe bevolkingsregister BRP. Daarvan werd de ontwikkeling in 2017 stopgezet, na tien jaar bouwen en 100 miljoen aan uitgaven.
- Digitalisering van de rechtspraak, die in april 2018 na zes jaar en ruim 200 miljoen euro (oorspronkelijk werden de kosten op 7 miljoen euro geschat) werd stopgezet.

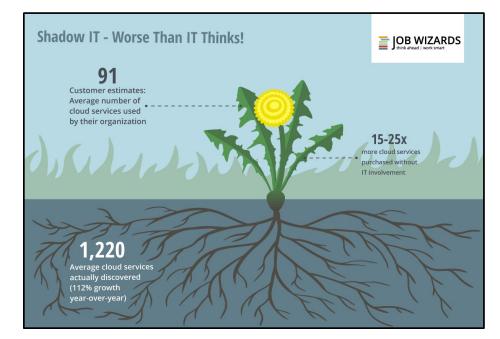
Is This Really the Entire Data Architecture?



Shadow IT



Source: https://www.emailvendorselection.com/ building-a-consolidated-tech-stack/

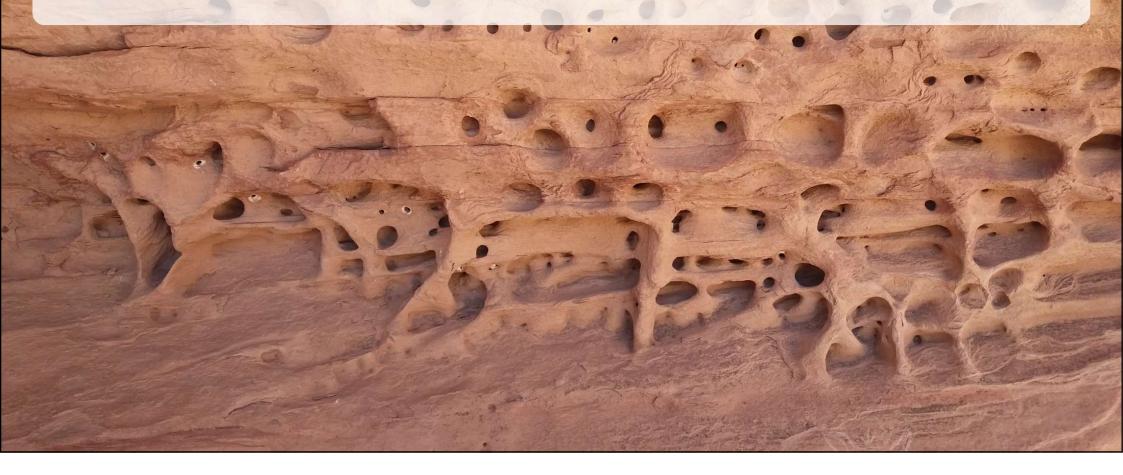


Source: https://job-wizards.com/en/ shadow-it-the-hidden-menace-for-every-company/



Copyright $\ensuremath{\textcircled{C}}$ 2025 R20/Consultancy B.V., The Netherlands

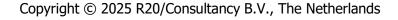
Part 2: Introduction to Data Architectures



What is a Data Architecture?



- Wikipedia: A data architecture is composed of models, policies, rules or standards that govern which data is collected, and how it is stored, arranged, integrated, and put to use in data systems and in organizations.
- Examples of data architectures:
 - Data warehouse architecture
 - Data streaming architecture
 - Transactional system

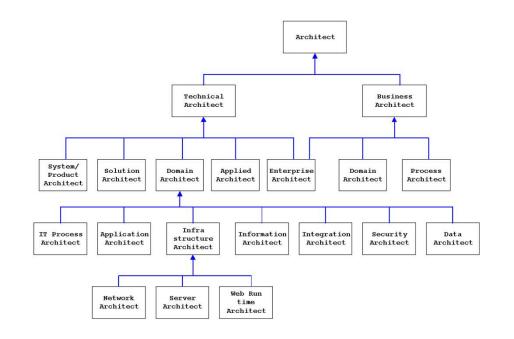


Data Architects versus Solutions Architects

Data Architects	Solutions Architects
focus on how information moves across the system from one application to another	look at the overall technological environment of the company
collaborate with clients to determine the specifications of the project, as well as the business goals that will align with the collected data	meet with their clients and establish their specific technology needs based on their business objectives
design the data model for the organization; where to store the customer data, how to retrieve the data; who can read the data	has a more technical point of view. Do we select a cloud solution, or on premise? What will the network look like? How will everything be connected without failures?



So Many Different Types of Architects

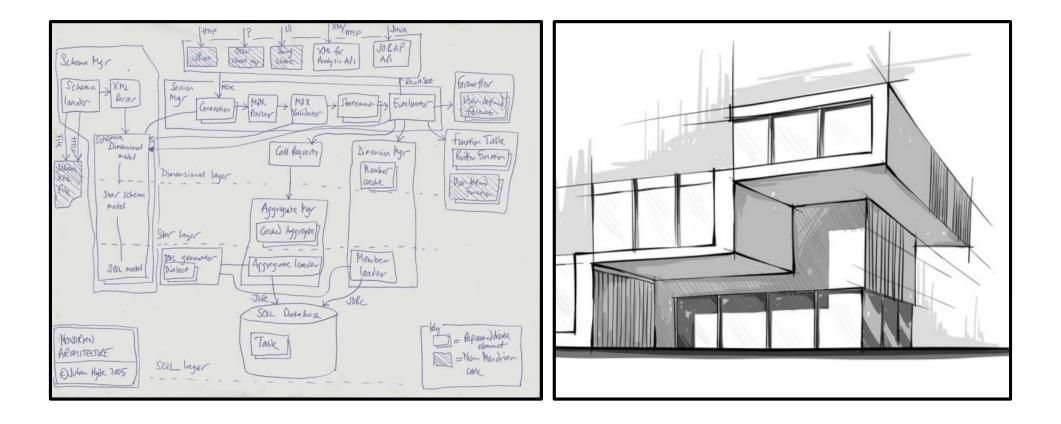


Paul Catalin Tomoiu : In [the] IT world, an architect is a person with enough knowledge to find a high level solution to a challenge in an IT environment.
The challenge type, defines the nature of the architecture role.
We can speak about an Enterprise Architecture, a Business Architecture, Data Architecture, Solution Architecture, IT

Technical Architecture, Application Architecture, Software Architecture, Security Architecture, etc.

Source: https://blog.prabasiva.com/2008/08/21/different-types-of-architects/

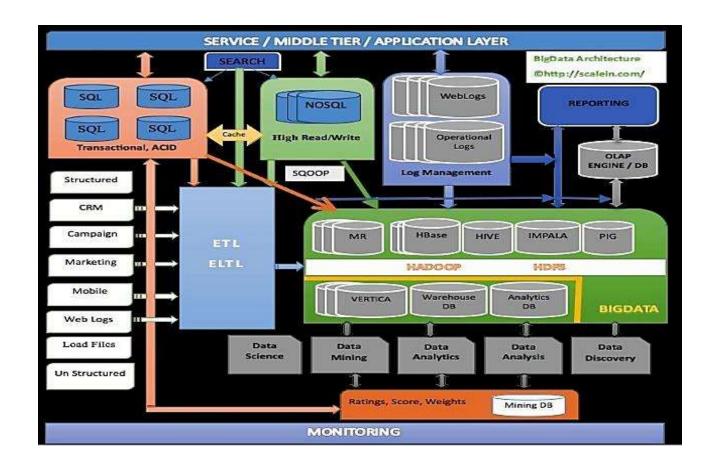
The Birth of a Data Architecture





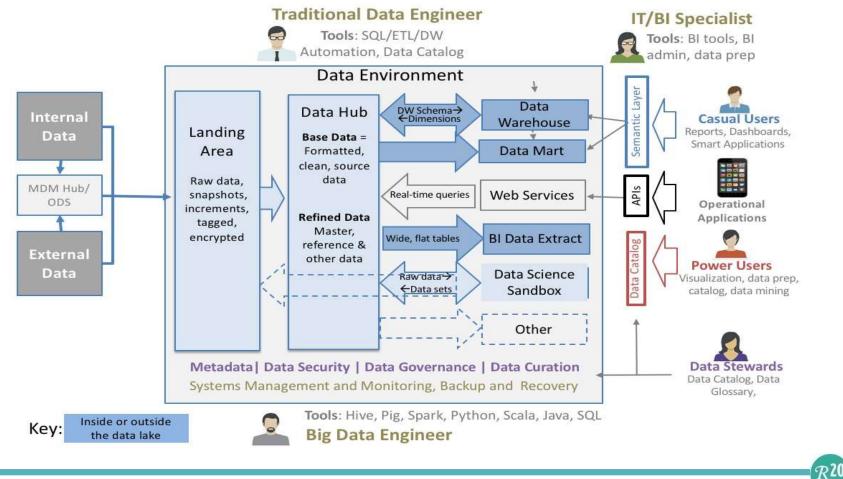
Copyright $\ensuremath{\textcircled{O}}$ 2025 R20/Consultancy B.V., The Netherlands

Quiz: What Type of Architecture is This?

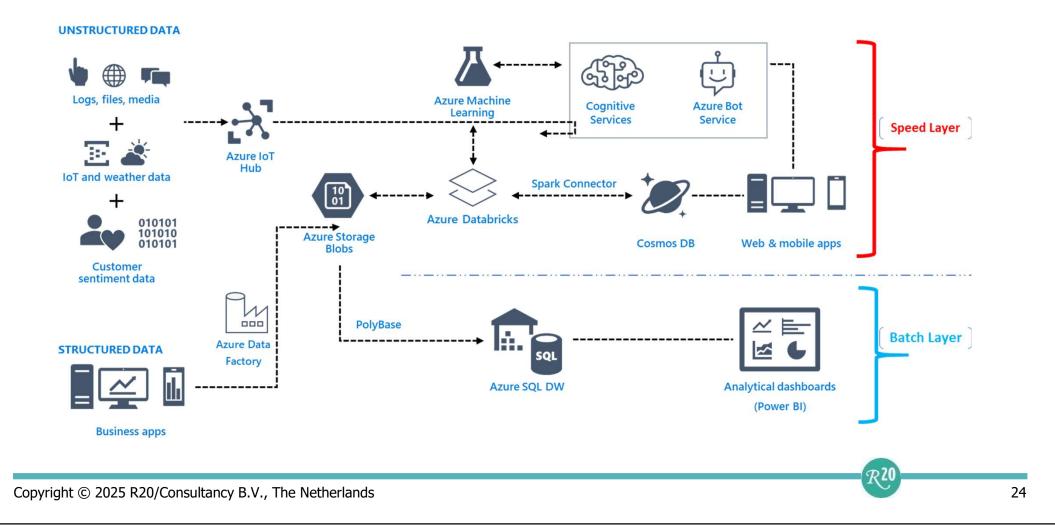




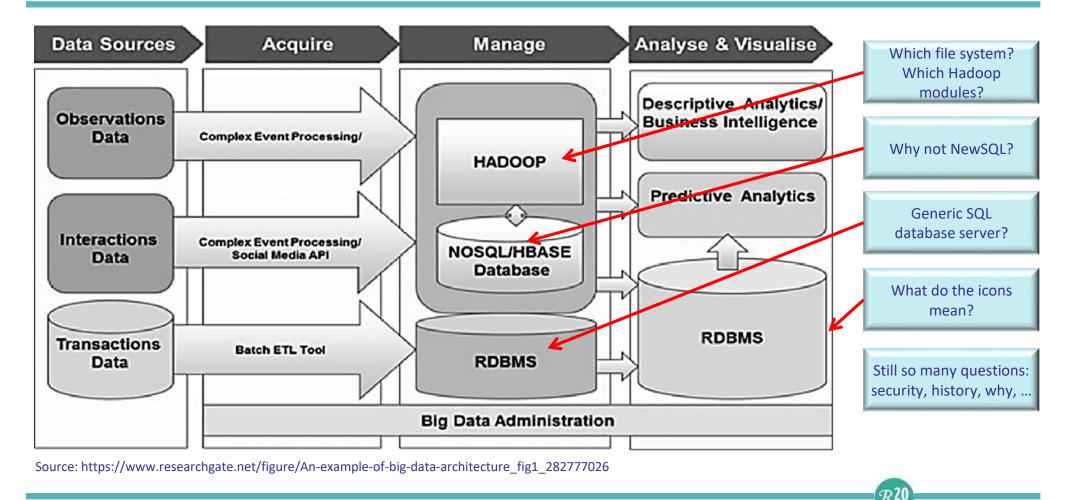
Quiz: What Type of Architecture is This?



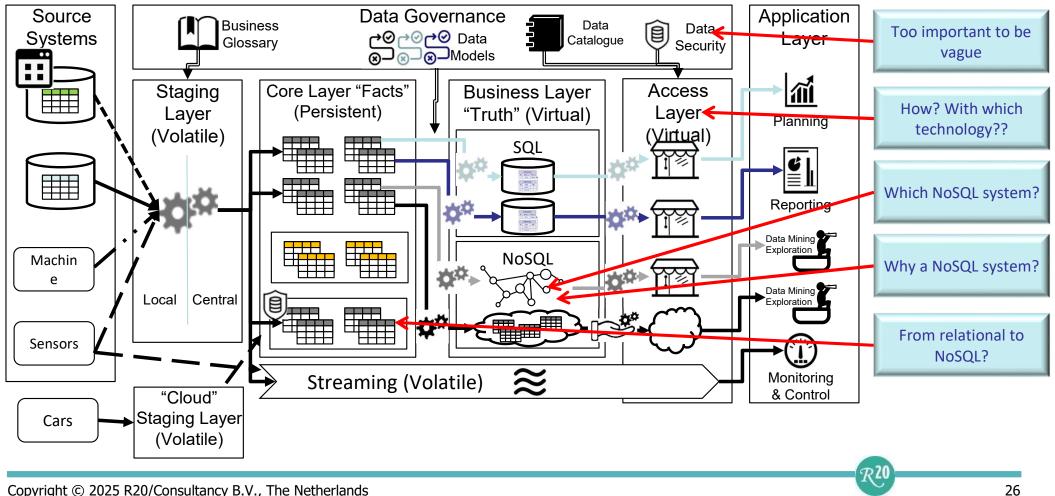
Quiz: What Type of Architecture is This?



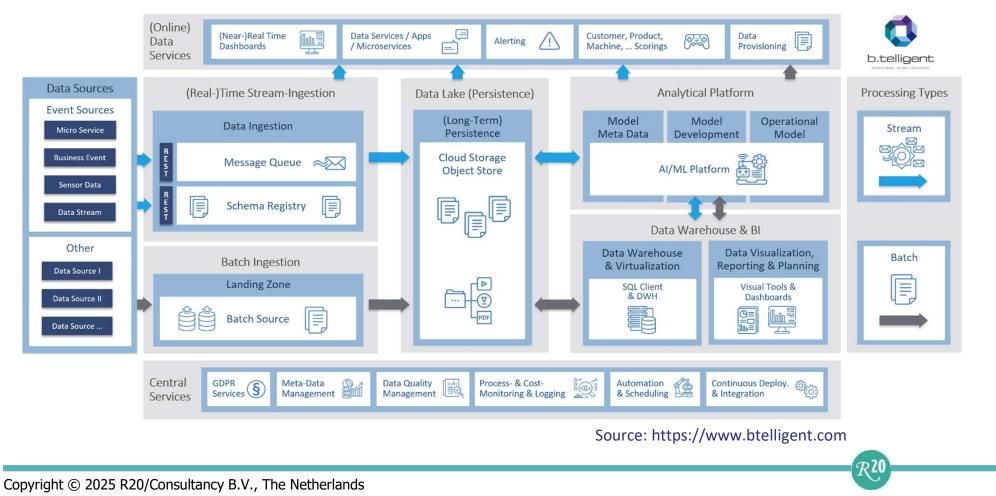
Data Architecture Example 1



Data Architecture Example 2



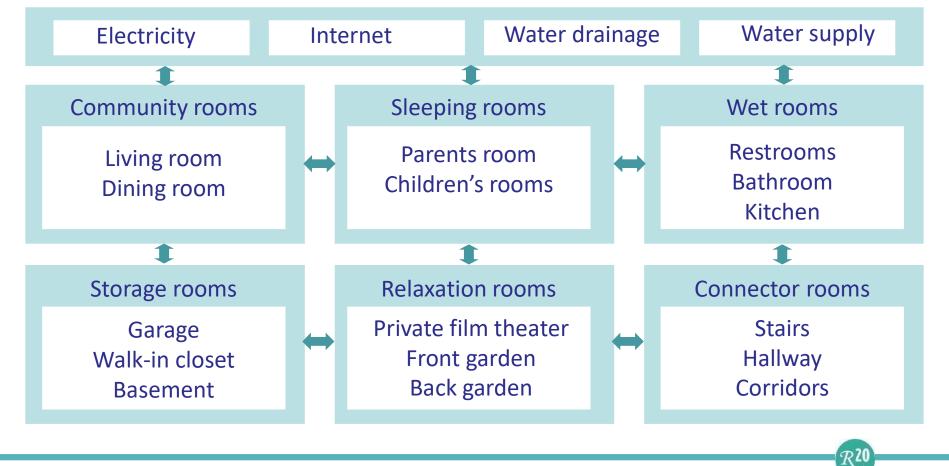
Data Architecture Example 3



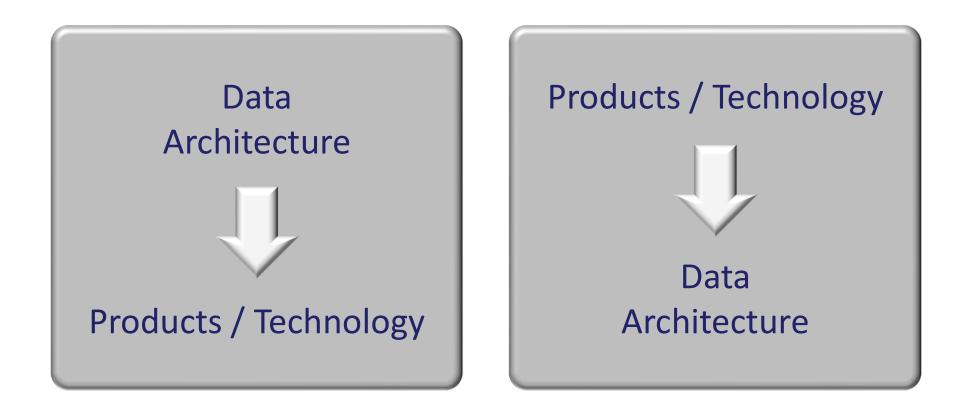
Architecture of a House



Architecture of a House (IT Style)



What Comes First?





Roadmap for Designing Data Architectures

1. Determine business motivations 2. Determine new requirements 3. Analyze the existing environment 4. Study new products and technologies 5. Define architectural design principles 6. Select a reference data architecture 7. Design the new data architecture 8. Determine the Implementation approach 9. Select new products and technologies 10. Introduce the data architecture within the organization

Part 3: Steps 1-3: Setting the Stage



Part 3.1: Step 1: Determine Business Motivations



Poor Examples of Business Motivations



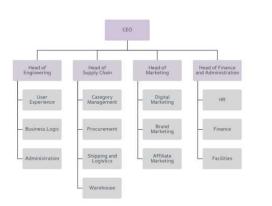
- Change insights and requirements
- Deployment of self-service BI
- Optimization of existing data architecture
- The platform on which the current BI system is hosted externally is old and needs to be replaced
- Move to the cloud
- Data science is not very well supported by current data warehouse environment
- We want to do more with the data we have, but it's hard to get to it

Proper Business Motivations



- Competitive improvement
 - Improving reaction speed to customer requests
- Support for customer journey and valuestream
- New business model
 - Allow customers real-time access to data
- Lower costs of specific business processes to improve margin
- Organization under threat
 - New competitor
- Comply with new laws and regulations
 - E.g. GDPR, CCPA, PSD2

Business Strategy and Data Strategy

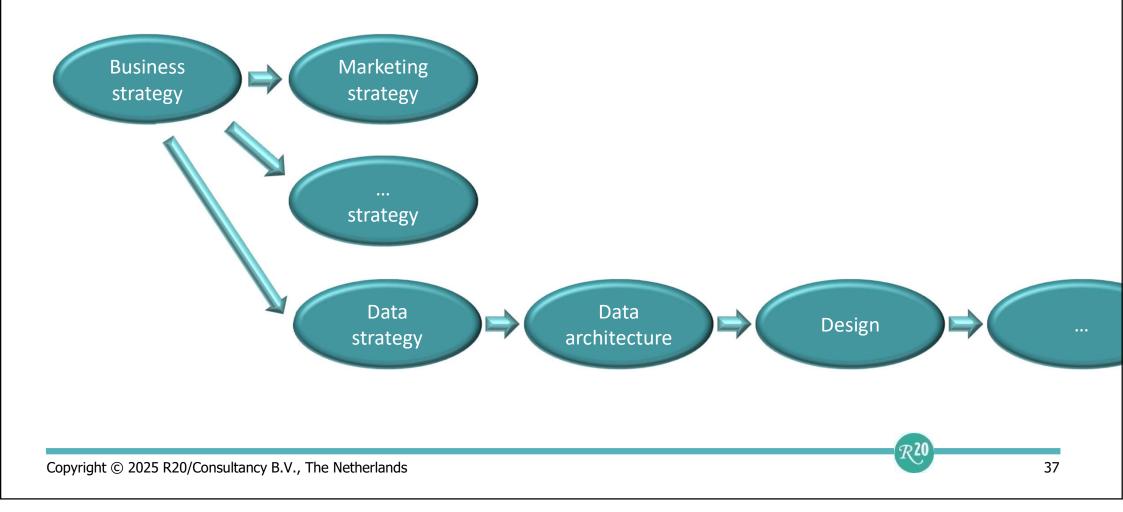


Business Strategy

- The challenges of top executives
 - New regulations, competitors, ...
- The main concerns for current business processes
- Future business developments
 - New business domains
- Data Strategy
 - New data architecture has to "fit" the data strategy
 - New demands for data delivery

36

From Strategy to Data Architecture and Onwards



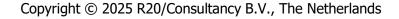
Part 3.2: Step 2: Determine New Requirements



Determine New Requirements



- New analytical functionality
- Lower latency for reports
- More users
- More access to metadata
- Migration to cloud platform
- More data
- More transparency of architecture
- Better security
- Deployment of data science



Determine Constraints



- Laws and regulations
 - GDPR, CCPA, PSD2, ...
- Budget restrictions
- Software limitations
 - One-stop shopping, open source preferred, company preferences, ...
- Hardware limitations
 - No easy processing, memory or storage scalability, ...
- Current legacy systems
 - Mainframe-based, proprietary applications, plain old, out-ofdate/obsolete development environments
- Internal ICT skills



Part 3.3: Step 3: Analyze the Existing Environment

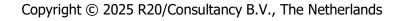


Determine Current ICT Bottlenecks



- Performance
- Report latency
- Productivity backlog
- Functionality
- Costs too high
- Business ICT cooperation
- Non-professional IT organization
- Not IT savvy

. . .



Analyze Existing Data Stores and Data Producers

Types of Data Stores

Transactional databases

Data warehouses

Data marts

Cubes

Data Lakes

Log files

Master databases

Private data files

External data sources

(Cloud) applications

Characteristics of Data Stores		
Size		
Workload		
Technology		
Limitations with respect to consumption		
Known performance problems		
How is data loaded?		
Security aspects		



Analyze Existing Applications

Data producers

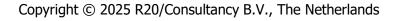
- Can we access the database directly or through an API?
- Current workload?
- Homemade or application?
- Data transformers and transporters
 - Home-made or professional (e.g. ETL, bus, data virtualization)?
 - Implementation style?
- Data Consumers
 - Homemade?
 - Internal or external?



Quantitative Information



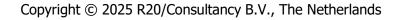
- Current and future of each
- Data consumption: number of users, transaction rate, complexity of queries, types of data consumption (batch, real-time), peaks, SLA's, ...
- Data size: database size (per table), data overlap, data growth, ...
- Data types: words, codes, numbers, dates/times, text, audio, video, geo/gis, multi-structured, ...
- Data update frequency: average, max, ...
- Data retention: minimum/maximum/preferred time of storage



Analyze Technology and Products in Use



- Selected products and versions
- Selected (cloud) platforms
- License costs
- Infrastructure
- Potential migration challenges





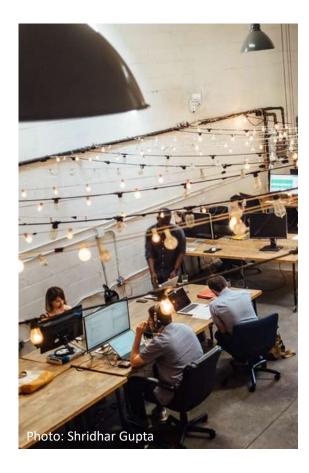
Determine the Culture of the IT Organization

Head of Engineering	Head of Supply Chain	Head of Marketing	Head of Finance and Administration
User Experience	Category Management	Digital Marketing	- HR
Business Logic	- Procurement	Brand Marketing	— Finance
- Administration	Shipping and Logistics	Affiliate Marketing	Facilities

- Traditional?
- Risk evasive?
- No experience with modern technologies?
- Cynical towards new developments?

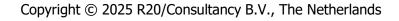


Determine IT Maturity Level of Organization (1)

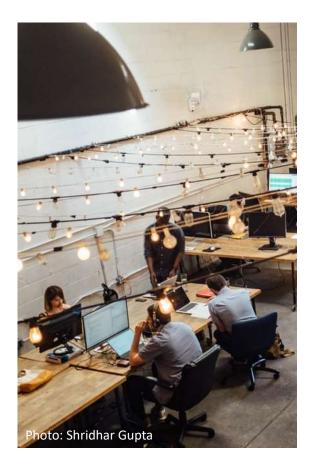


Data processing checks

- Is data primarily stored to support business processes and to conform to reporting regulations?
- Can DBAs see the data?
- Are ETL processes started manually?
- Is ETL crash automatically fixed?
- Are data processing specifications scattered across all modules?
- Is metadata available and kept up to date?
- Are "old" reports reproducible?

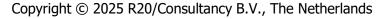


Determine IT Maturity Level of Organization (2)

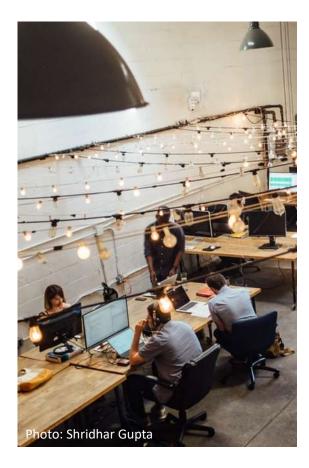


Data consumption

- Do reports primarily show what has happened within business processes?
- High data latency?
- Do they use predictive analytics to optimize business processes and decision-making processes?
- Data management
 - Ownership of data assigned?
 - Is there focus on data quality?
 - Are there procedures in place to fix incorrect data?

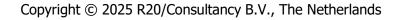


Determine IT Maturity Level of Organization (3)

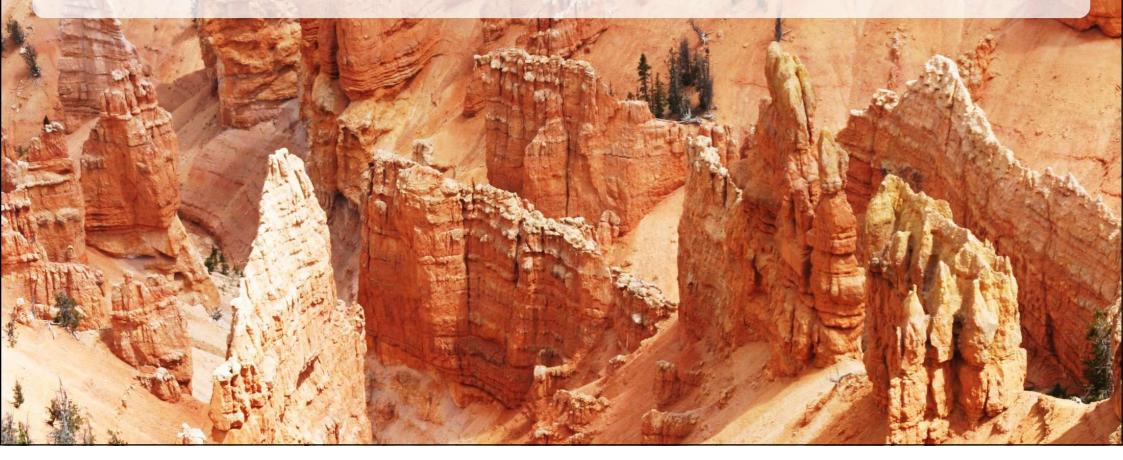


ICT skills

- All development outsourced?
- Many tool-jockeys?
- Performance anxiety?
- Minimal knowledge of new technologies?



Part 4: Step 4: Analyze New Technologies for Data Storage, Processing, and Analytics

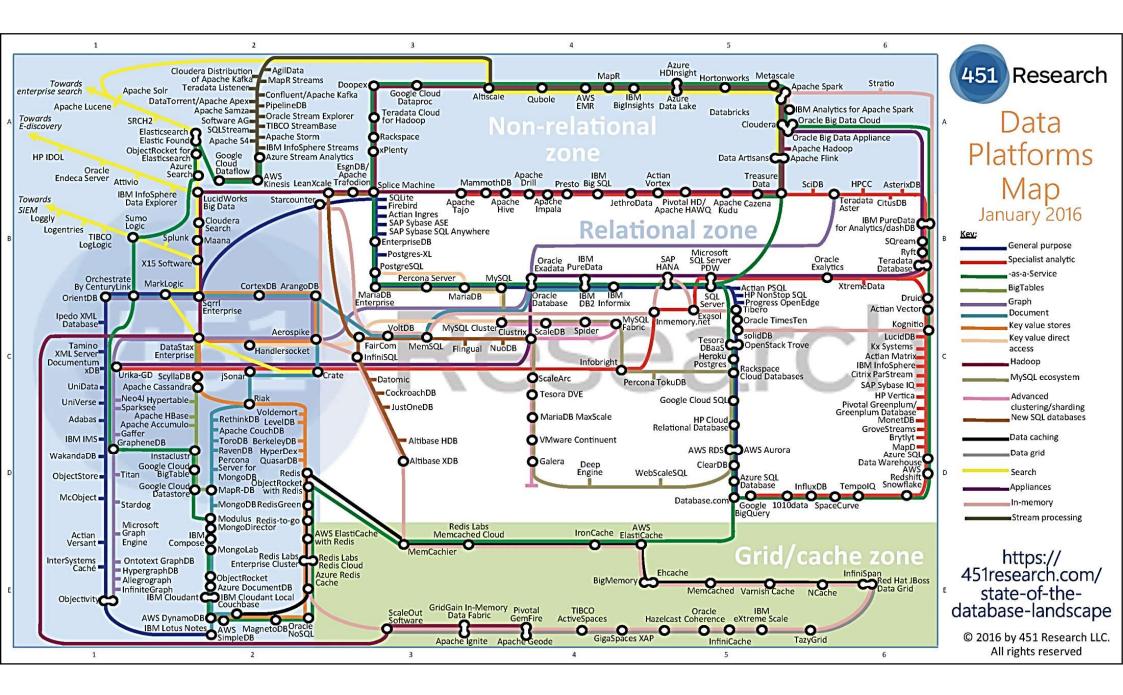


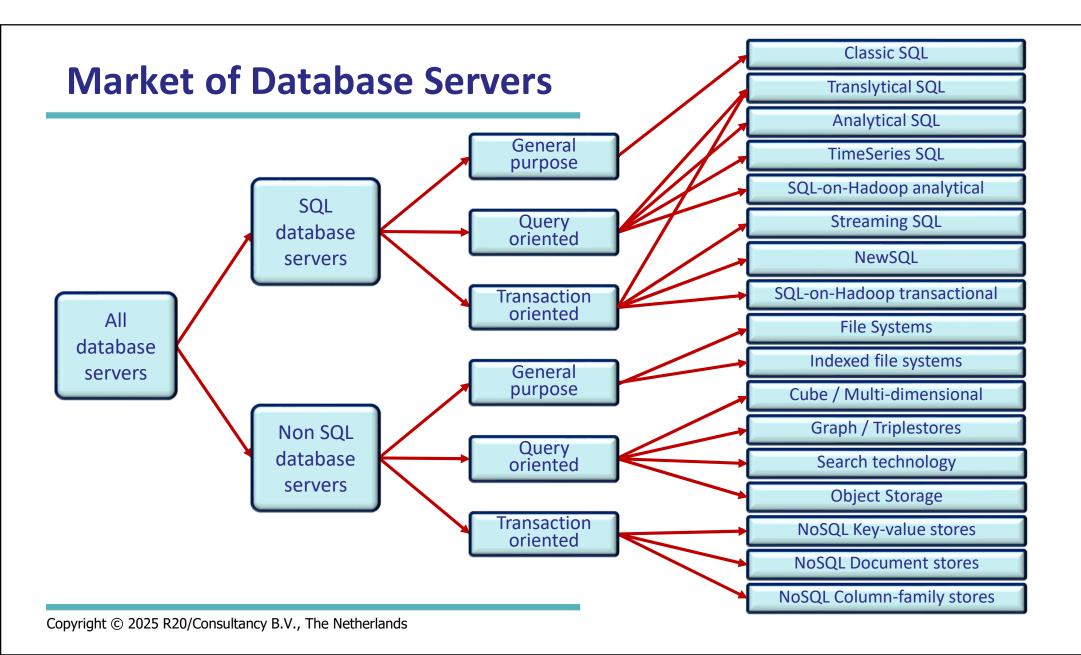
Part 4.1: **Data Storage**

Trends in Database Market

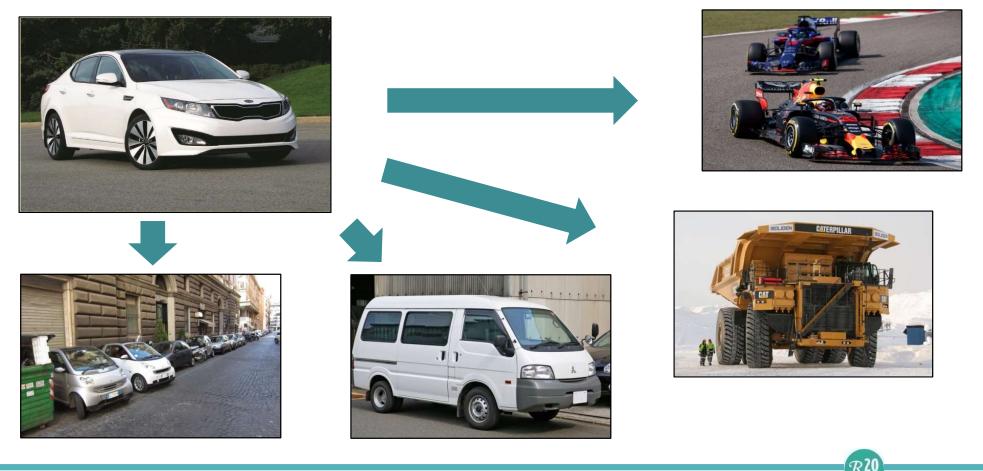


- More data, more queries, more transactions, more concurrent users, more complex queries, ...
- Coming and going of products
- Less standards
 - Less portability
- No interchangeable products
- Specialization of products
 - Limited use cases



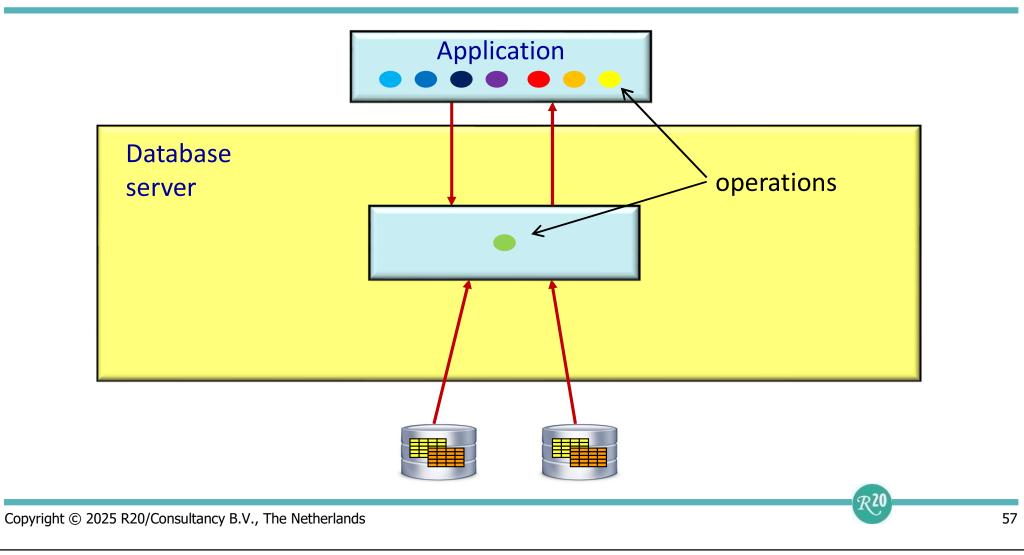


Specialization of Cars

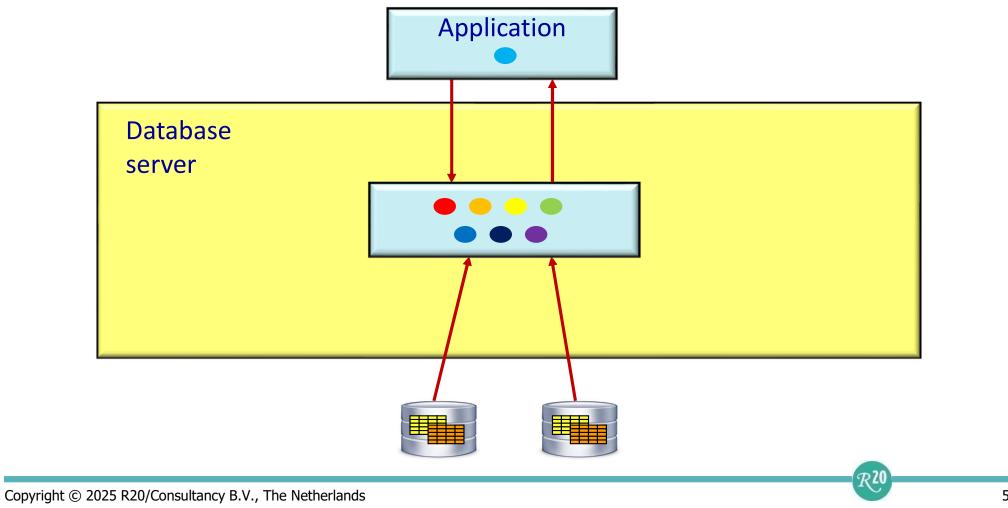




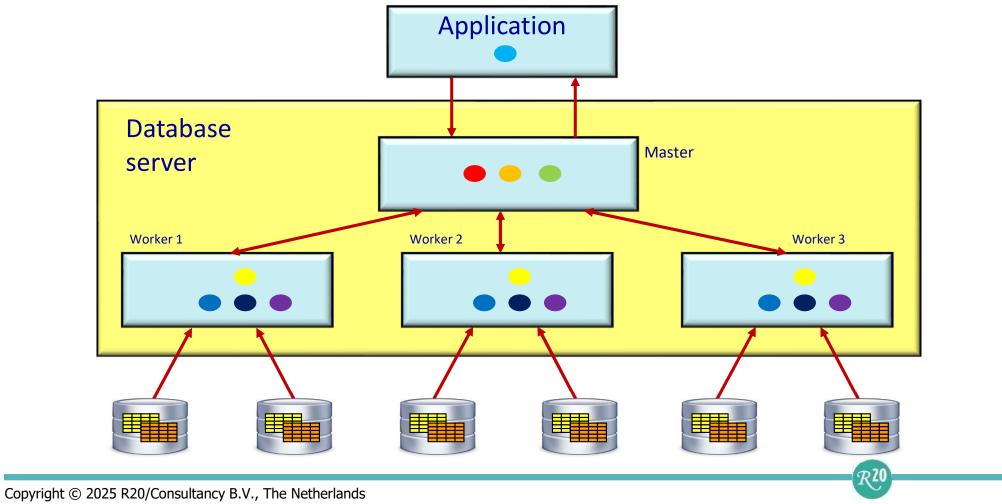
Application-based Analytics



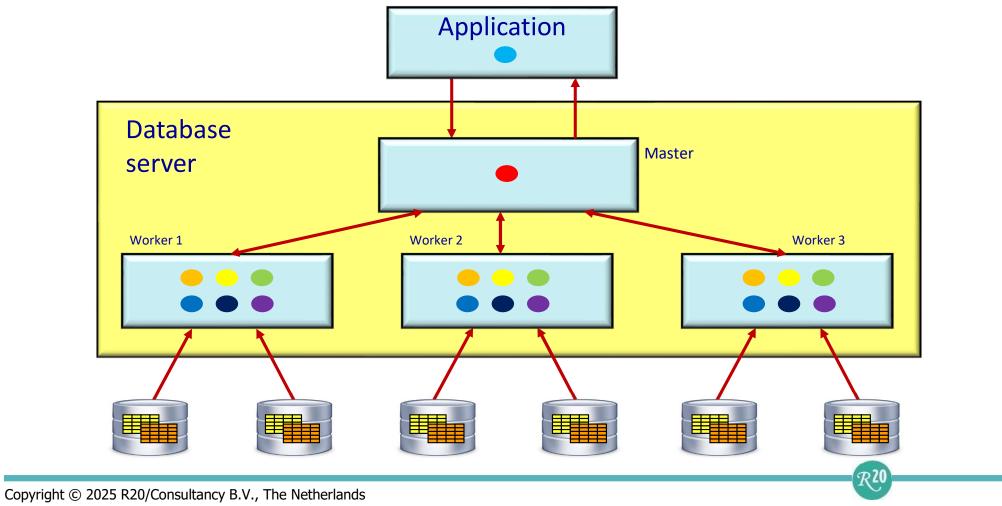
In-Database Analytics



Partial Parallel Analytics

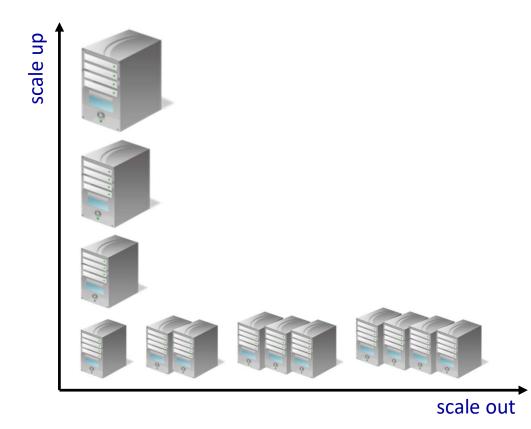


Full Parallel Analytics



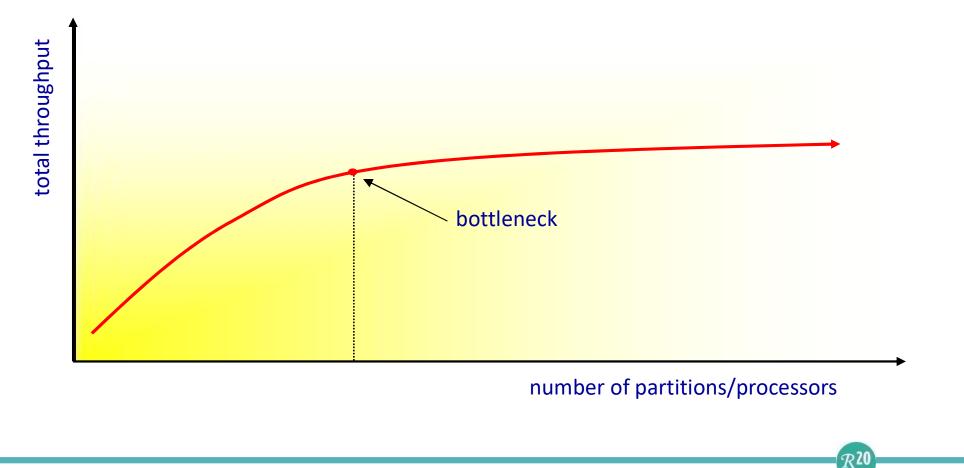
60

Scale Up versus Scale Out

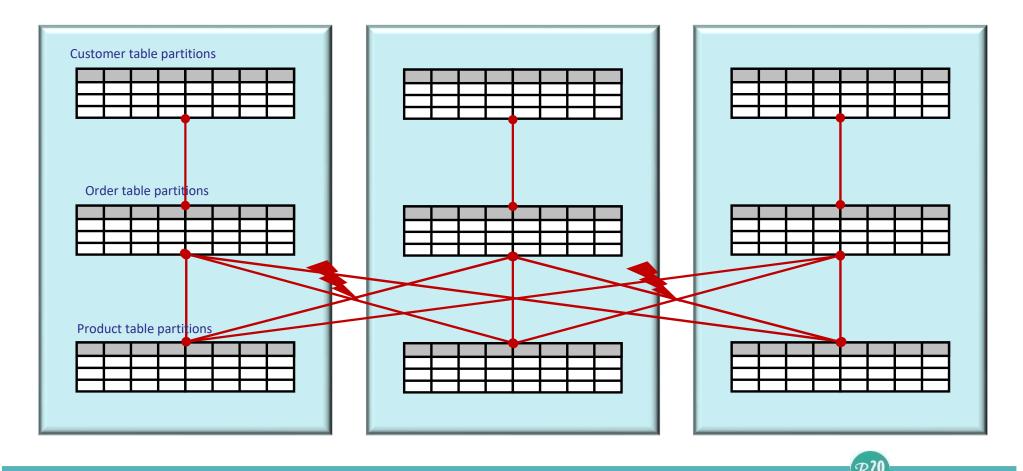


- Scale up (vertical scaling) means adding more resources to one node in a system
- Scale out (horizontal scaling) means adding more nodes to a system
 - Continuous availability/redundancy
 - Cost/performance flexibility
 - Contiguous upgrades
 - Geographical distribution

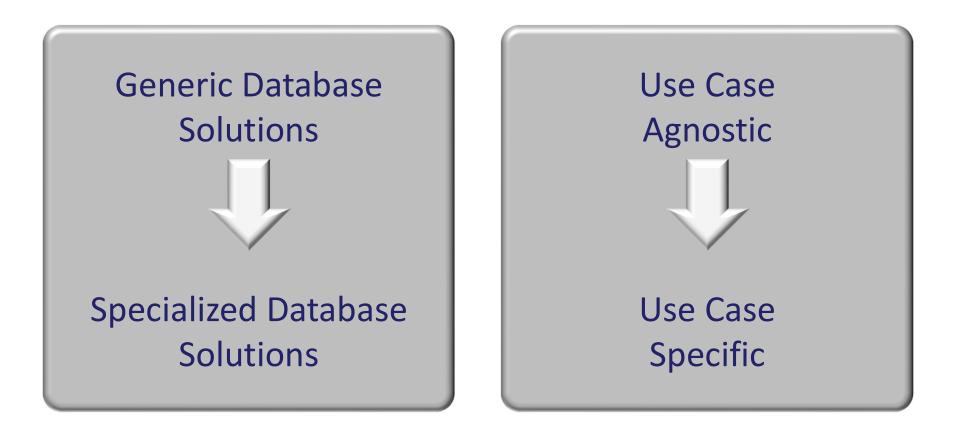
Effect of Partitions on Query Response



Optimizing Distributed Joins



Database Technology has Changed





NoSQL Database Servers

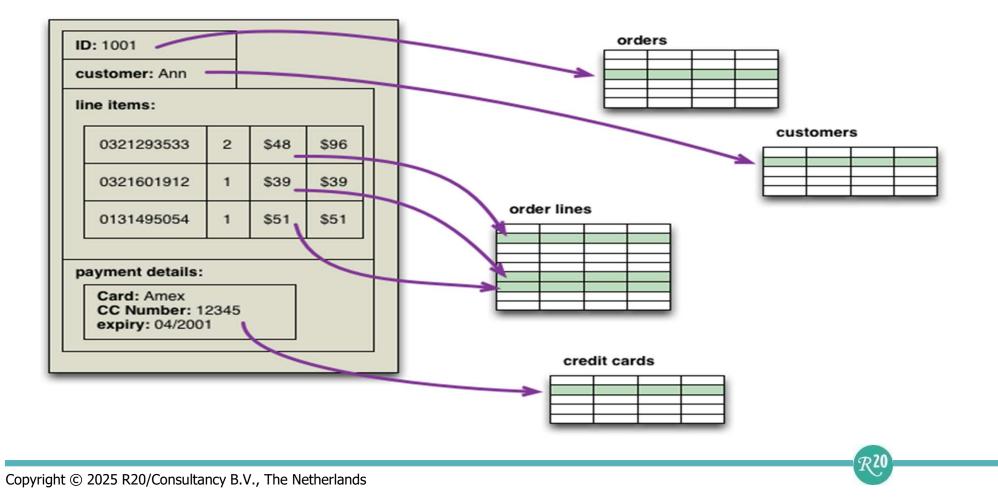


Tricks to Improve Performance



- Aggregate data model
 - To remove the impedance mismatch
- Design architecture to scale-out
 - Sharding
- Reduce functionality (security, query power, data integrity, ...)
- Lower consistency
- Give developers full control over internal processing
- Push down" complex operations

NoSQL: Aggregate Data Model



Typical NoSQL Use Cases

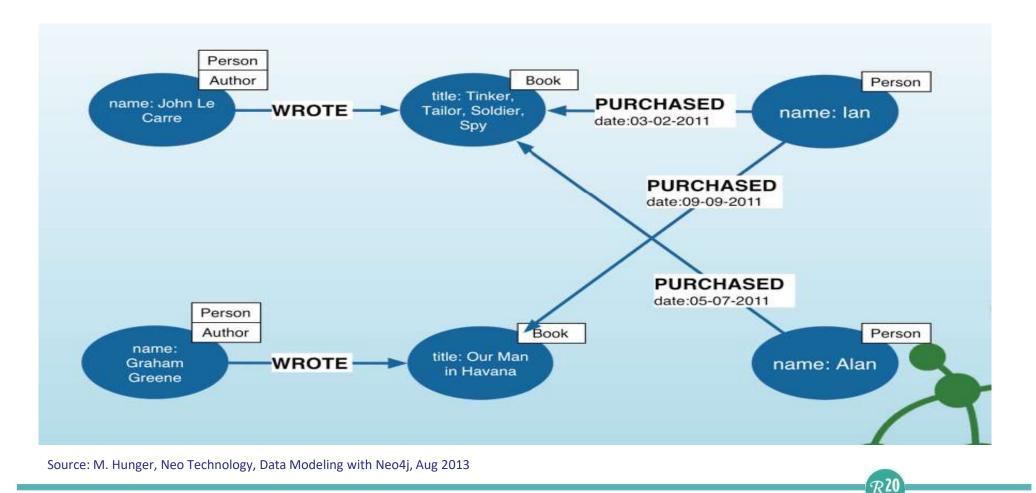


Transactional

- Big transactional workload
- Single record/document transactions
- Massive data ingestion
- Simple reporting point queries
- Dynamic data structures
- Complex data structures
- Narrow" data model



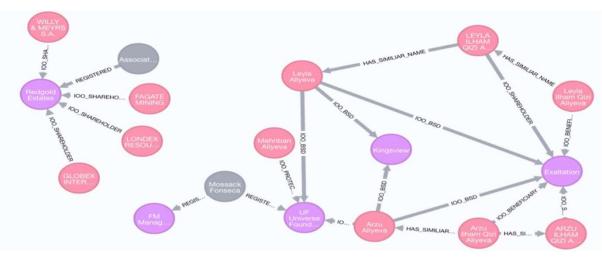
Graph Database Servers



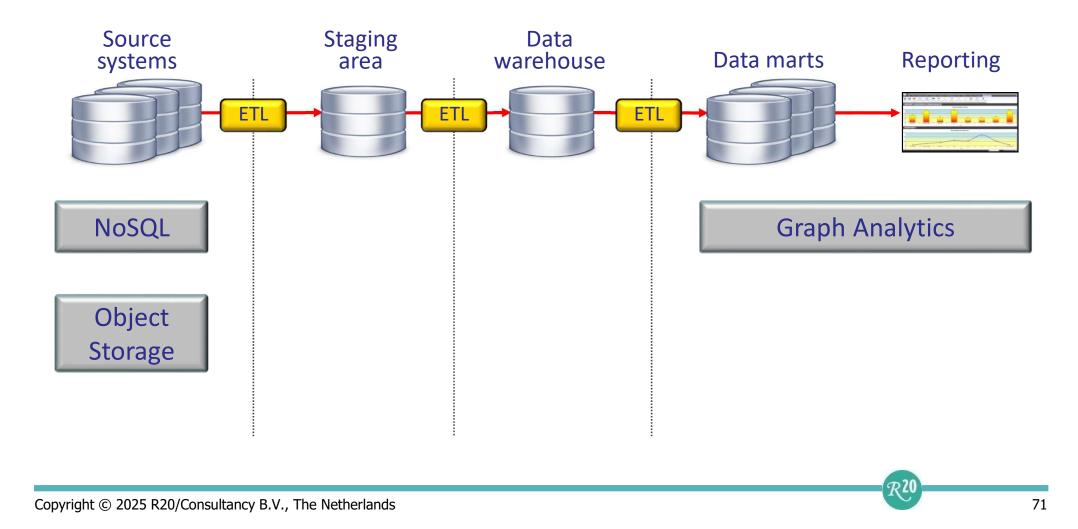
69

Application Areas Graph Databases

- Social network analysis
- Network impact analysis
- Optimal route determination
- Internet retail recommendations
- Logistics
- Fraud analysis
- Securities and debts
- "Panama papers"
- And many more



Deploying NoSQL, Graph and Object Storage



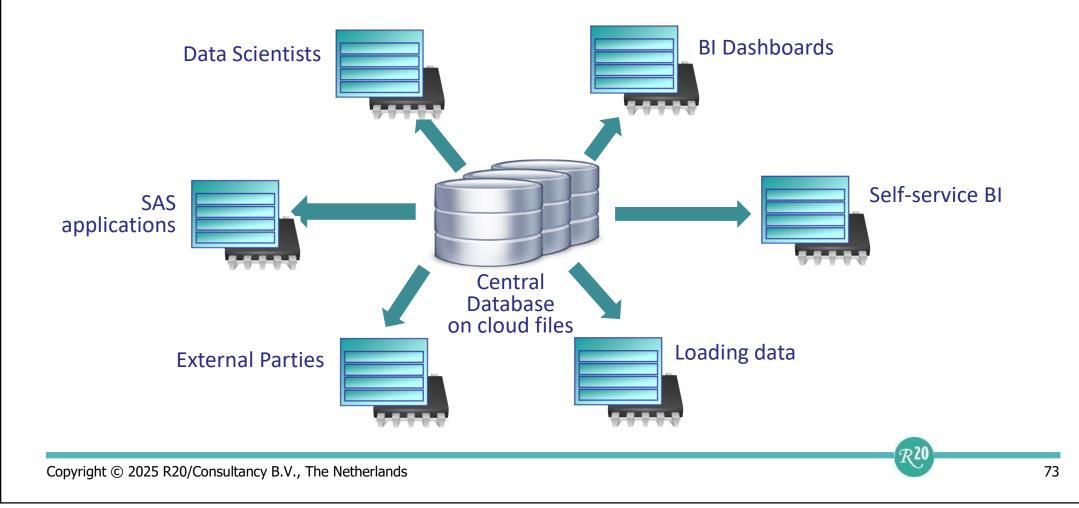
Analytical SQL Database Servers

Examples		
1010Data	Microsoft Azure Synapse	
Amazon Redshift and Athena	OmniSci (MapD)	
Apache HAWQ	Oracle Database In-Memory	
BlazingSQL	SAP HANA	
CitusDB (PostgreSQL)	SAP Sybase IQ	
ClickHouse SQL	SnowflakeDB	
Databricks Delta Lake	Splice Machine	
Edge Intelligence	SQream	
Exasol	Starburst (Trino formerly Presto)	
Google BigQuery	Teradata Vantage	
Greenplum	XTremeData dbX	
IBM DB2 Warehouse on Cloud	Several SQL-on-Hadoop engines	
Ignite InfoBright DB	Vertica	
Kinetica	And many others	
Kognitio WX2		

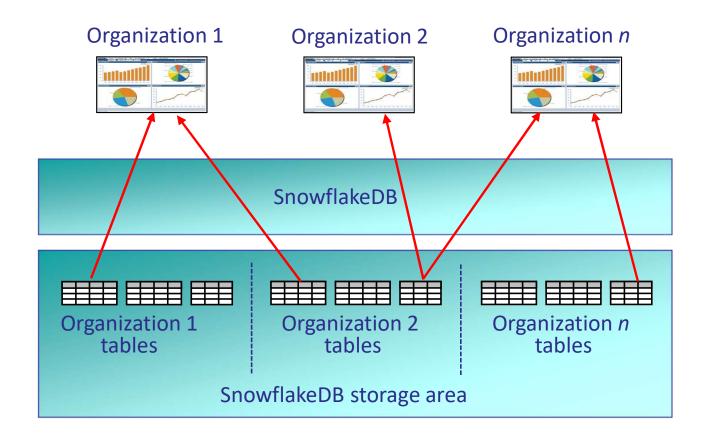
Copyright © 2025 R20/Consultancy B.V., The Netherlands

 \mathcal{R}^{20}

Example 1: Snowflake

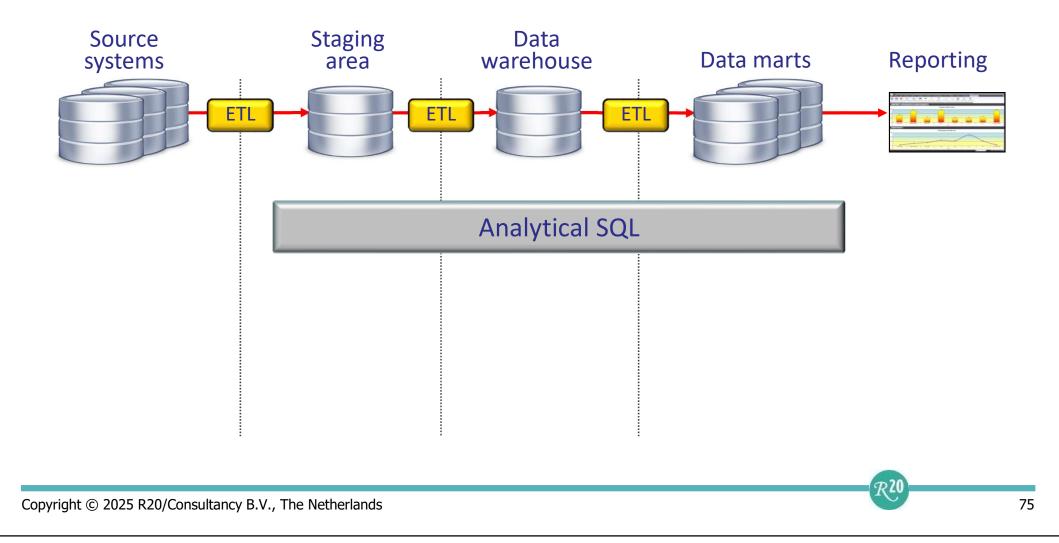


Example 1: Snowflake

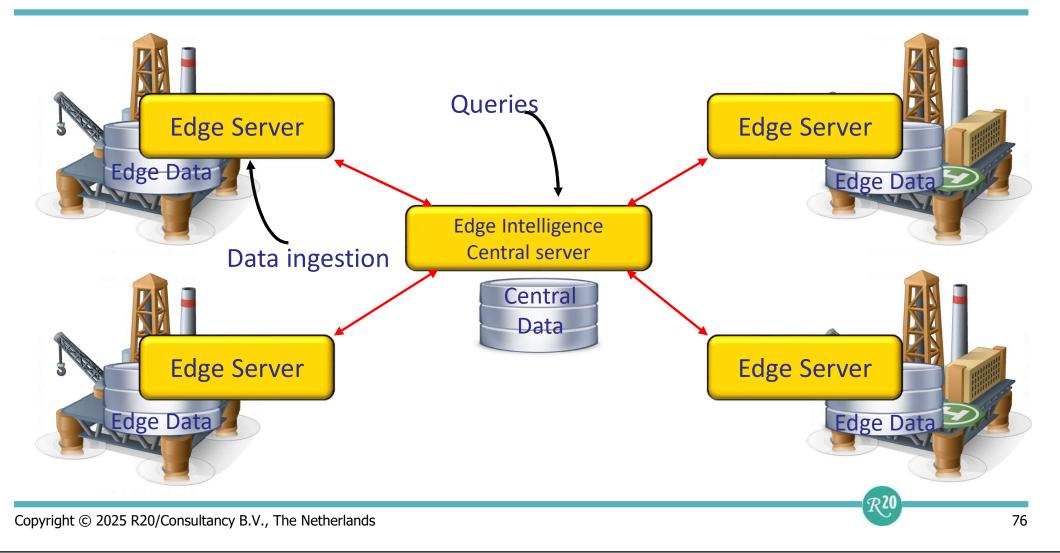




Deploying Analytical SQL Databases

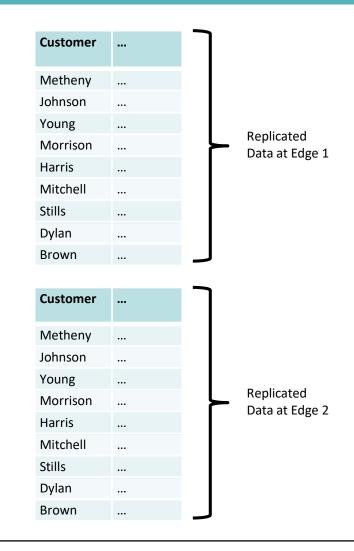


Example 2: Edge Intelligence



Example 2: Edge Intelligence

ſ	Store_id	Customer	Datetime	
	1	Metheny	2017-11-01 12:00:08	
Non-replicated	1	Johnson	2017-11-01 12:10:18	
Data at Edge 1	1	Young	2017-11-01 12:12:33	
	1	Morrison	2017-11-01 12:50:09	
	1	Harris	2017-11-01 12:55:45	

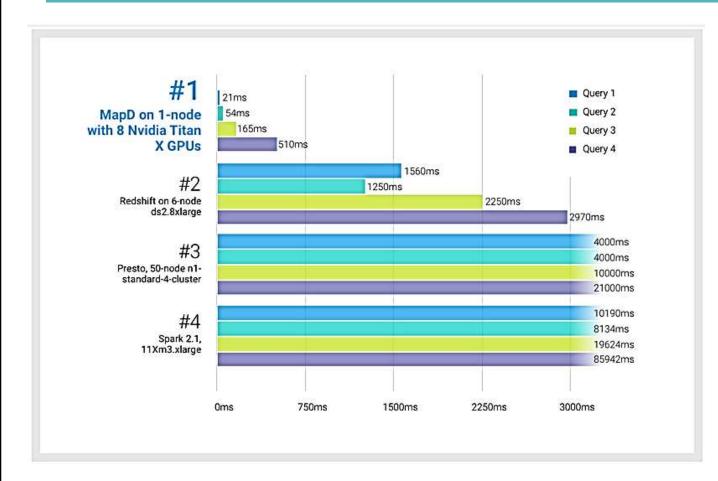


Non-replicated Data at Edge 2	Store_id	Customer	Datetime	
	2	Metheny	2017-11-01 12:01:32	
	2	Mitchell	2017-11-01 12:05:42	
	2	Stills	2017-11-01 12:11:39	
	2	Dylan	2017-11-01 12:12:30	
	2	Brown	2017-11-01 12:40:19	

NVIDIA TITAN V: GPU With More Than 5,000 Cores



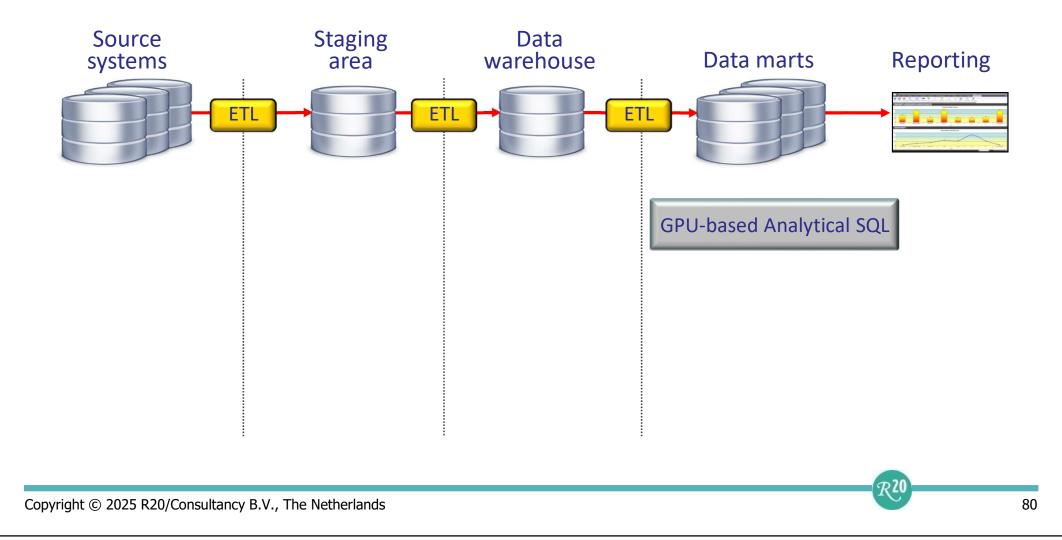
Comparison of Analytical SQL Database Servers



- Products:
 BlazingSQL, Kinetica,
 OmniSci (MapD),
 SQream
- They make use of the parallel power of GPU's
- Long-term data persistency is not their core business

R20

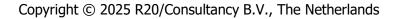
Deploying GPU-based SQL Databases



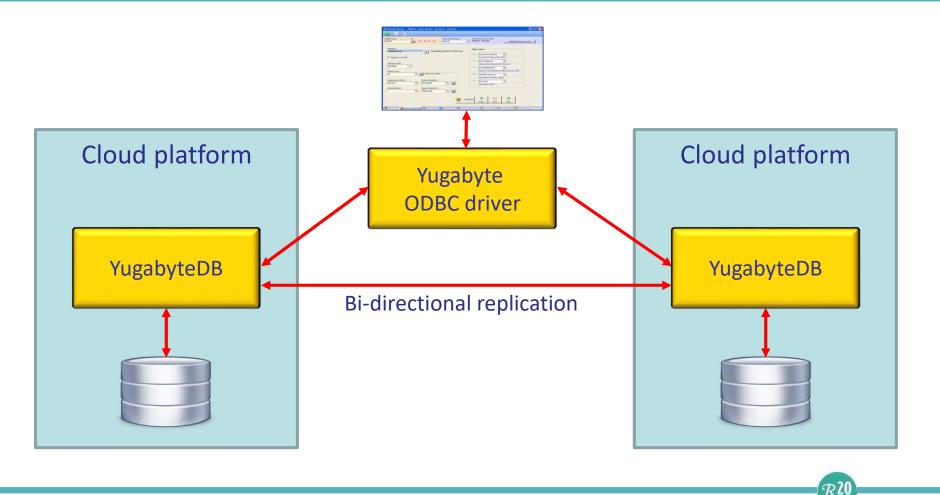
Transactional SQL Database Servers



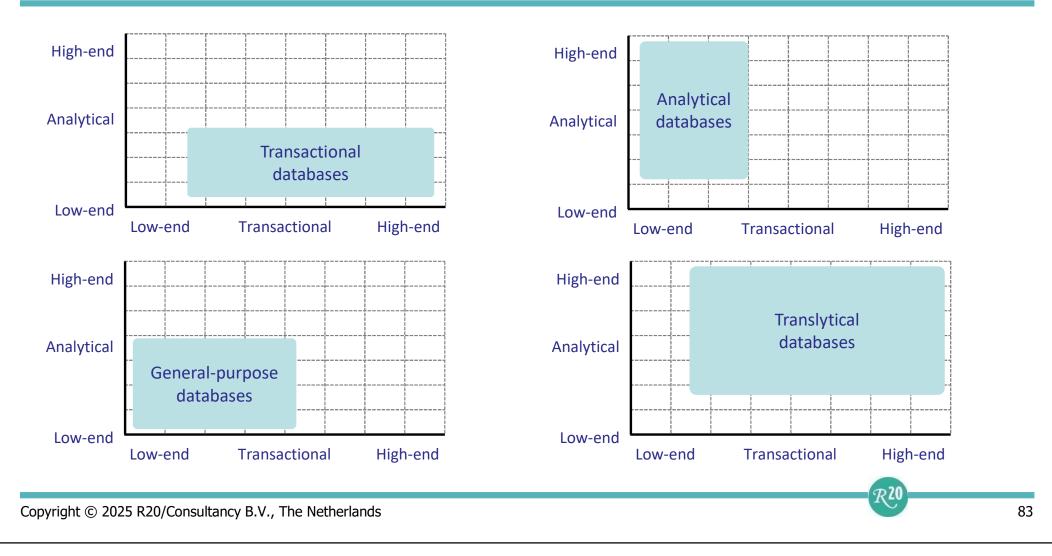
- Examples: Clustrix, DataBricks Delta lake, SingleStore (MemSQL), Splice Machine, Pivotal GemFire XD (SQLFire), VoltDB, and YugabyteDB
- NewSQL is not a new SQL dialect
 - The internal architectures are different from classic SQL database servers
- High scalability with respect to transactions
- Full-blown SQL high level of data independence
- ACID-compliant = 100% consistency
- Exploitation of low-cost clusters



Example: YugabyteDB

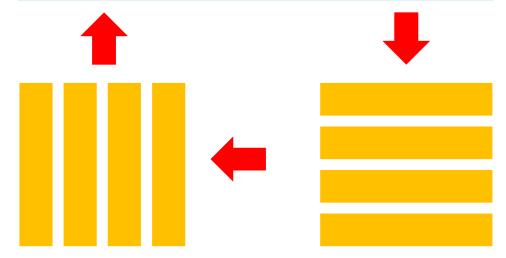


Four Categories of SQL Databases



Example: SingleStore (translytical)

ID	Name	Initials	Date Entered	City	State
12345	Young	Ν	Aug 4, 2008	San Francisco	CA
23324	Stills	S	Sep 10, 2009	New Orleans	LA
57657	Furay	R	Oct 16, 2010	Yellow Springs	ОН
65461	Palmer	В	Nov 22, 2011	Boston	MA

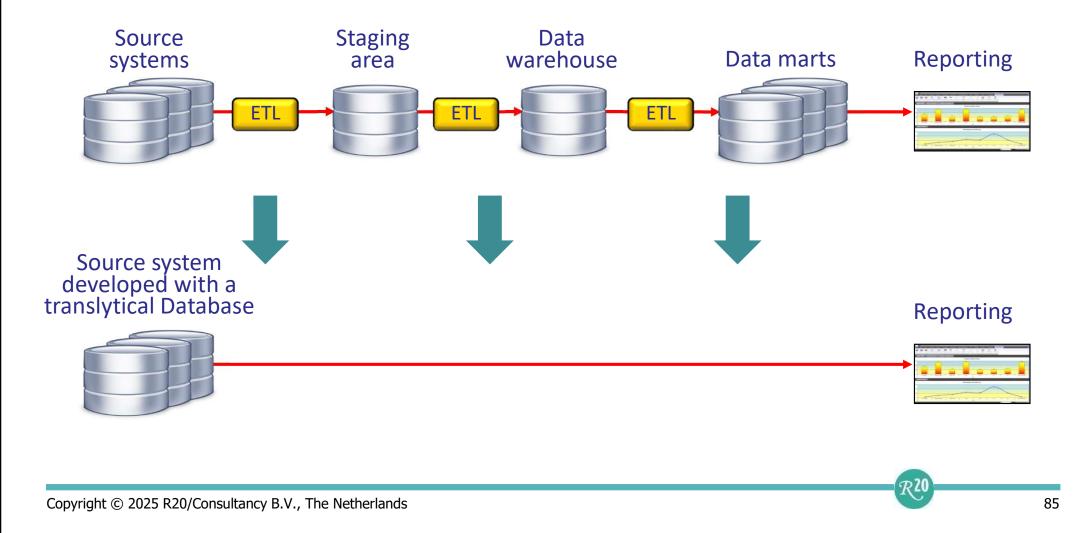


Copyright © 2025 R20/Consultancy B.V., The Netherlands

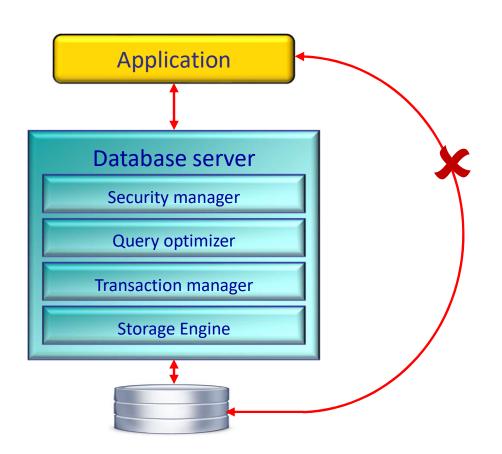


R20

Zero Data-Latency Architectures



Most Database Servers Use Proprietary Files

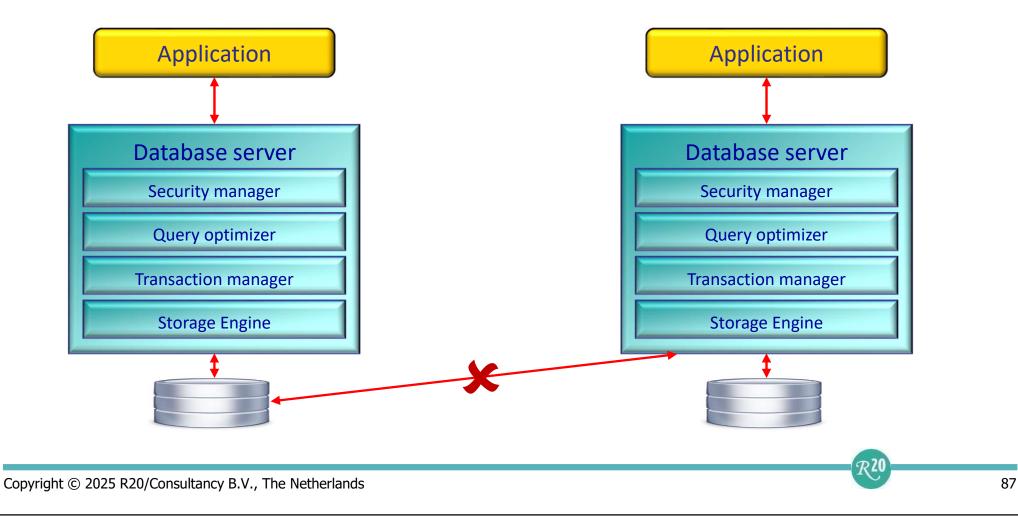


- In many database servers a proprietary storage format is used
- Data can only be accessed via database server
- Data needs to be copied for other database servers

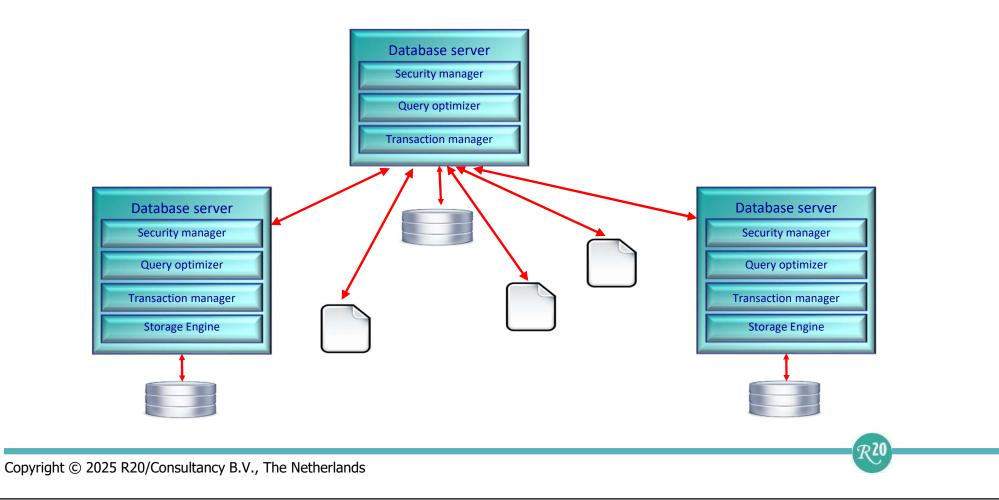


Copyright $\ensuremath{\textcircled{C}}$ 2025 R20/Consultancy B.V., The Netherlands

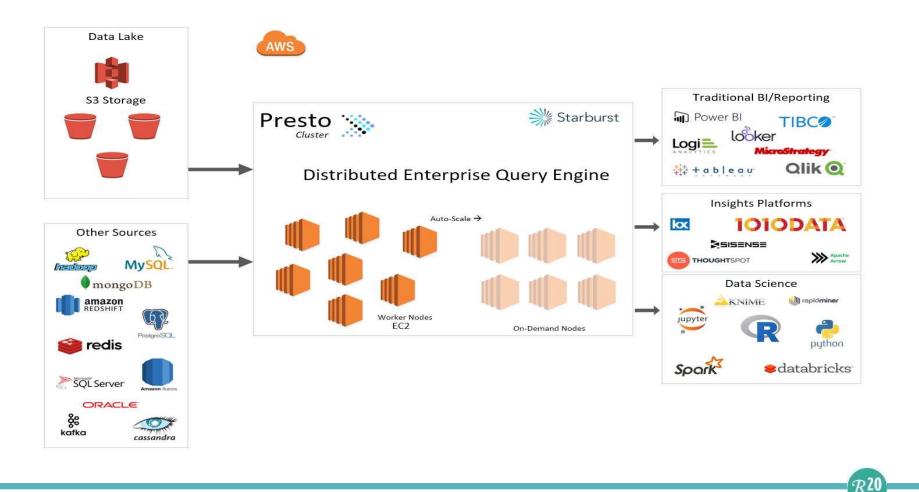
Database Servers Can't Share Data



Accessing "External" Data



Example: Starburst (based on Trino)

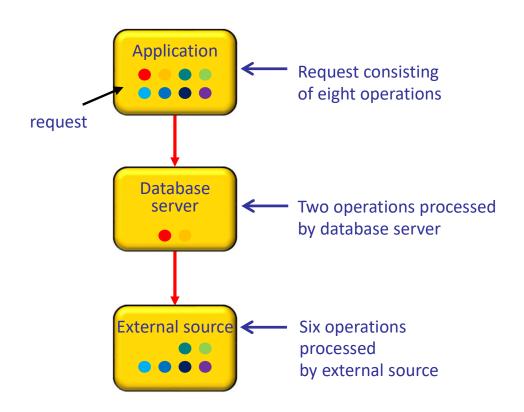


Example: Amazon Athena

```
CREATE EXTERNAL TABLE employee
  ID string,
 NAME string,
 AGE string,
  GEN string,
 CREATE DATE bigint,
 PROCESS NAME string,
 UPDATE DATE bigint
STORED AS AVRO
LOCATION 's3://my-bucket/staging/employees'
TBLPROPERTIES (
'avro.schema.literal'='
4
    "type" : "record",
    "name" : "AutoGeneratedSchema",
    "doc" : "Sqoop import of QueryResult",
    "fields" : [ {
      "name" : "ID",
      "type" : [ "null", "string" ],
      "default" : null,
      "columnName" : "ID",
```

```
}, {
  "name" : "NAME",
  "type" : [ "null", "string" ],
  "default" : null,
  "columnName" : "NAME",
  "salType" : "12"
}, {
  "name" : "AGE",
  "type" : [ "null", "string" ],
  "default" : null,
  "columnName" : "AGE",
  "salType" : "2"
}, {
  "name" : "GEN",
  "type" : [ "null", "string" ],
  "default" : null,
  "columnName" : "GEN",
  "sqlType" : "12"
}, {
  "name" : "CREATE DATE",
  "type" : [ "null", "long" ],
  "default" : null,
  "columnName" : "CREATE DATE",
```

Accessing External Data and Query Pushdown

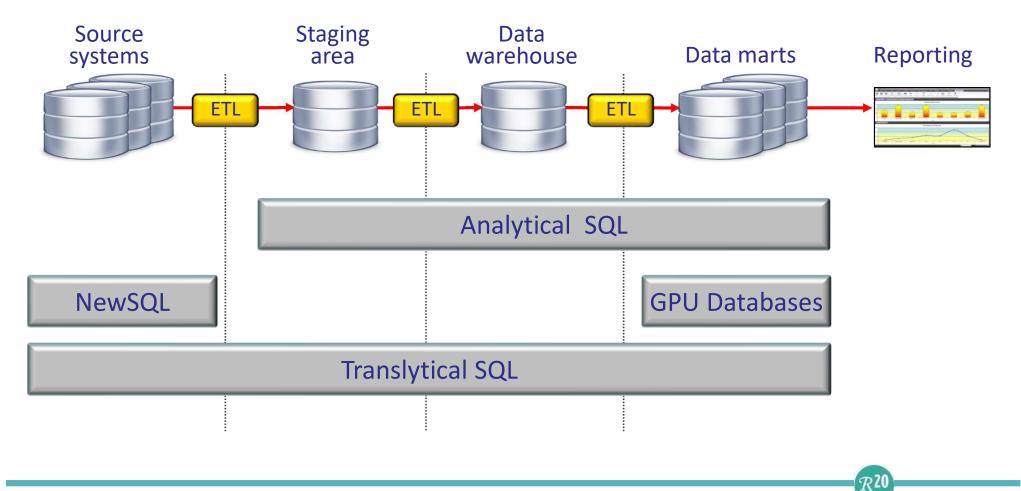


- Push processing to the external source
- Minimize network traffic
- Exploit the full power of the external source
- Optimize distributed joins
- Deal with datatype differences
- From structure-less data to structure-rich data

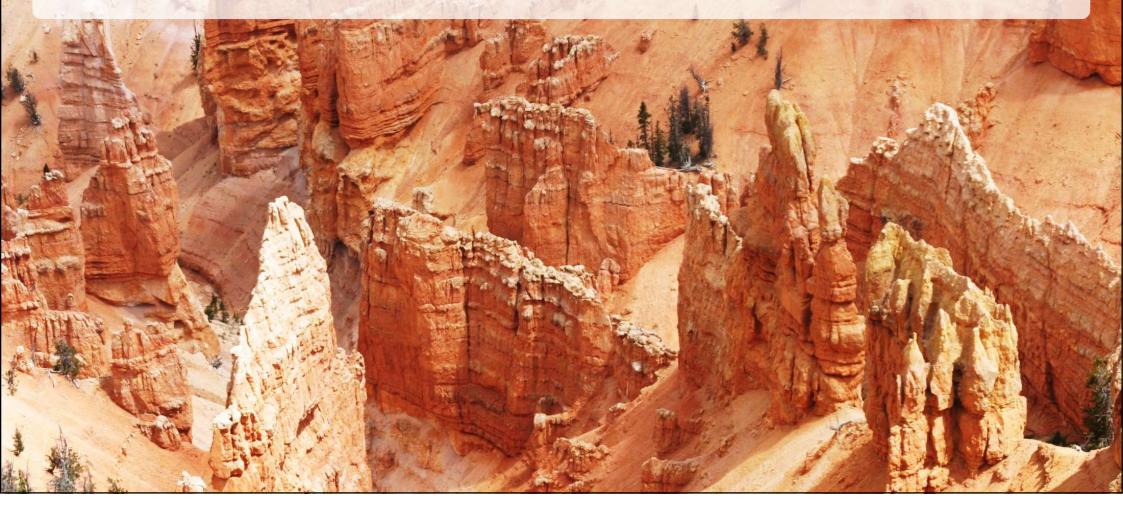


Copyright $\ensuremath{\mathbb{C}}$ 2025 R20/Consultancy B.V., The Netherlands

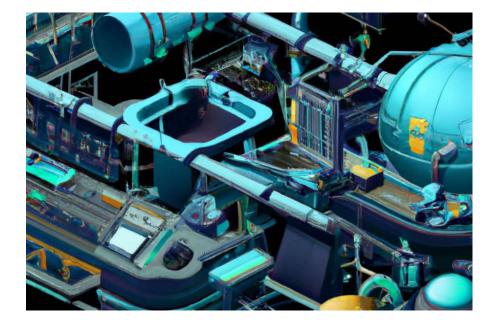
Deploying Specific SQL Databases



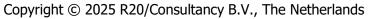
Part 4.2: Data Processing



Categories for Data Processing



- ETL (Extract Transform Load)
- Data Replication (Change Data Capture)
- ESB (Enterprise Service Bus)
- Data Virtualization



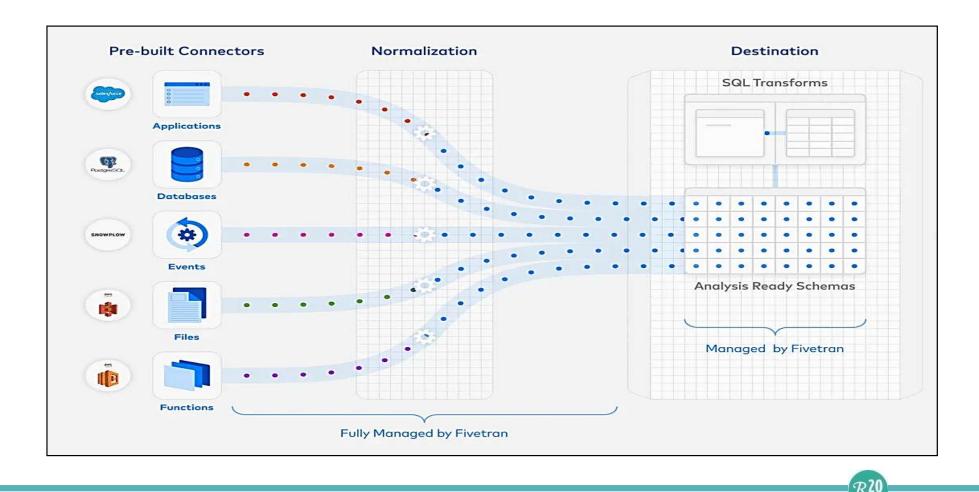


ETL = Extract Transform Load

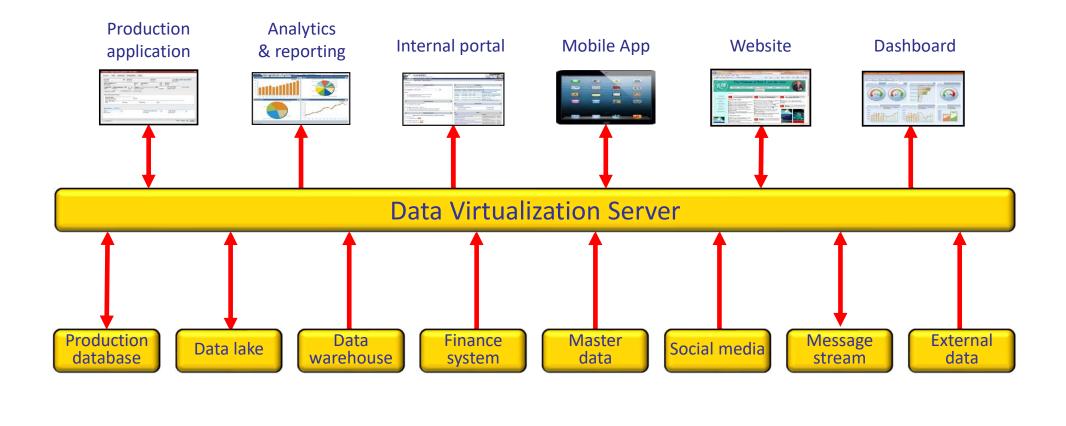
- Transforming of data structures
 - To a data structure suitable for reporting and analysis
- Cleansing of data
- Integration of data from production systems
- Transforming data
 - Filtering, aggregating, projecting, joining, splitting, ...
- Scheduling the ETL process
 - Batch-oriented
- Managing the ETL process



Example: Fivetran

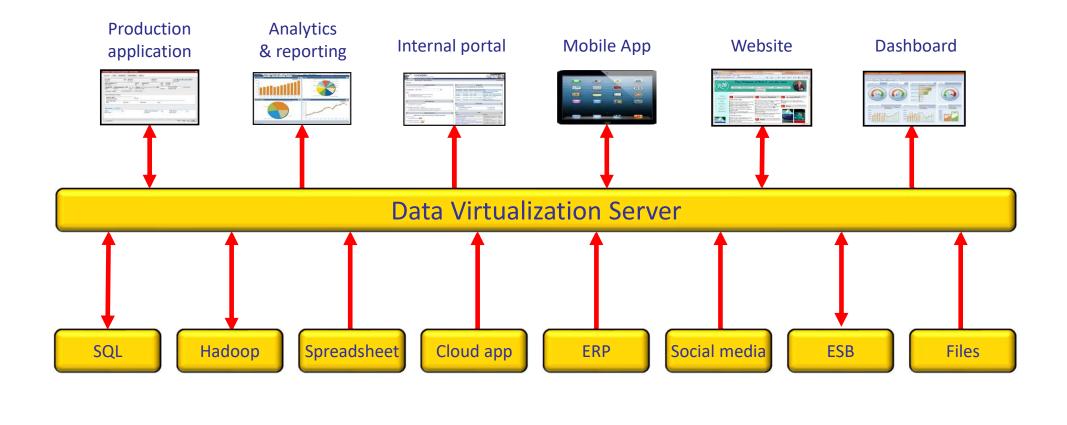


Data Virtualization Overview (1)



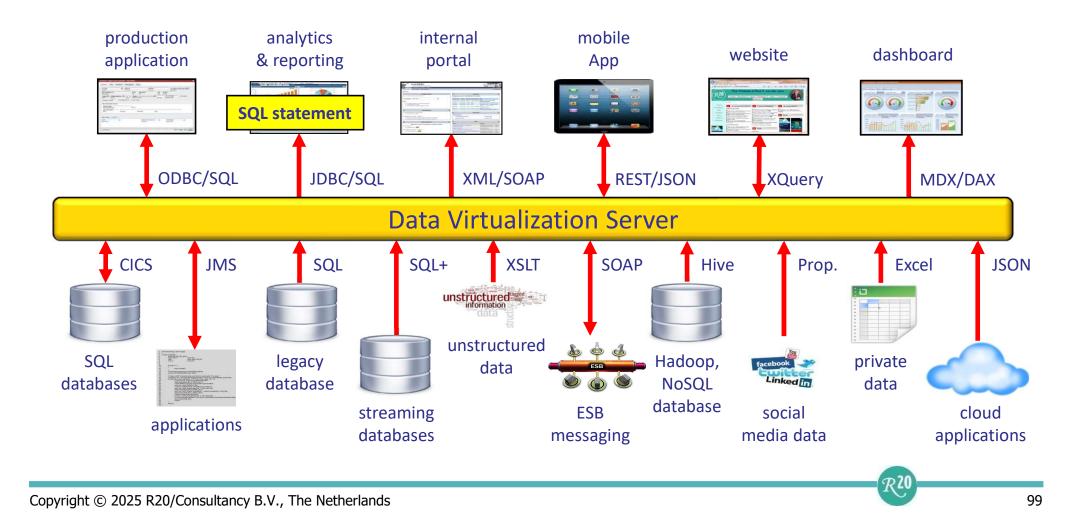
R20

Data Virtualization Overview (2)

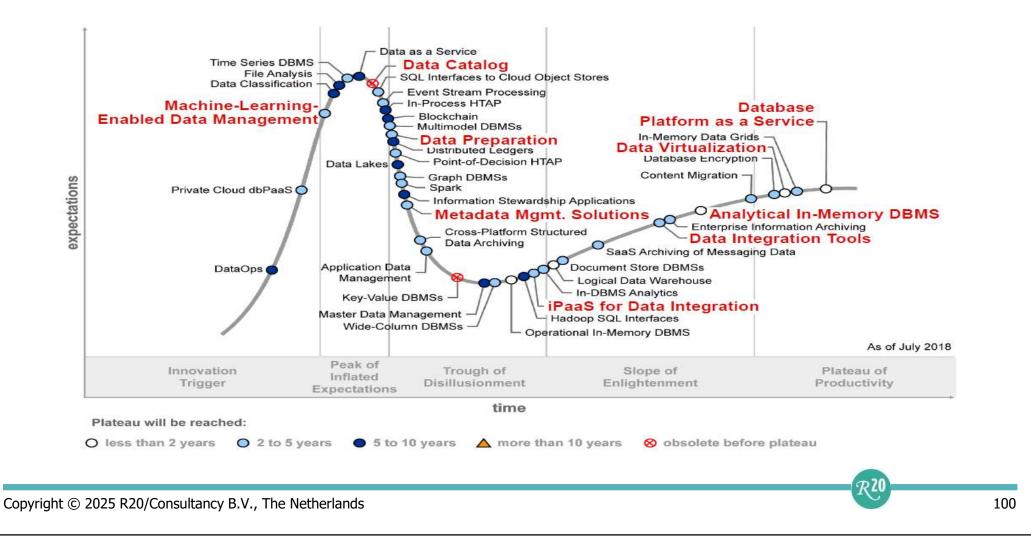


 \mathcal{R}^{20}

Data Virtualization Overview (3)



Gartner Gives Data Virtualization its Highest Maturity Rating



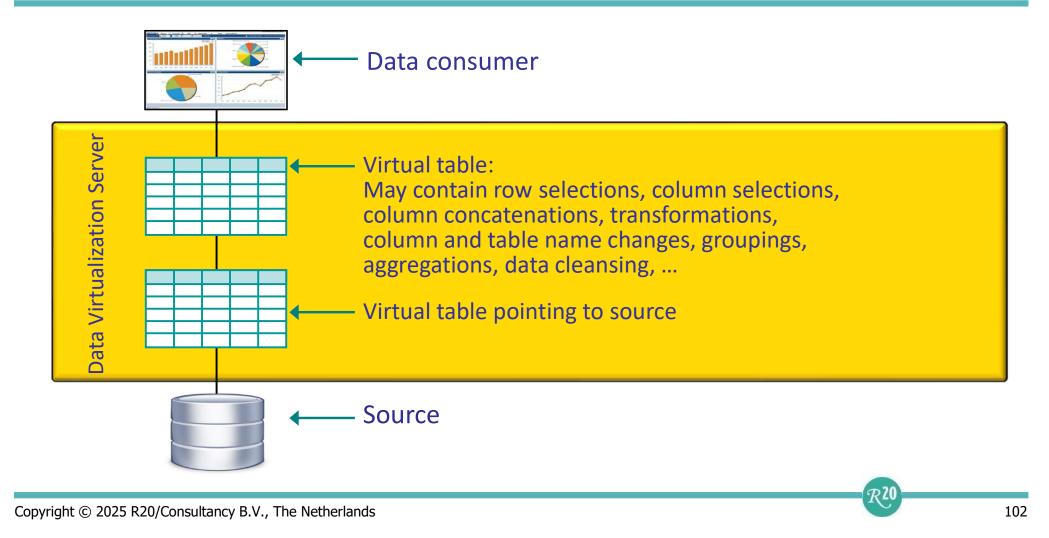
The Market of Data Virtualization Servers



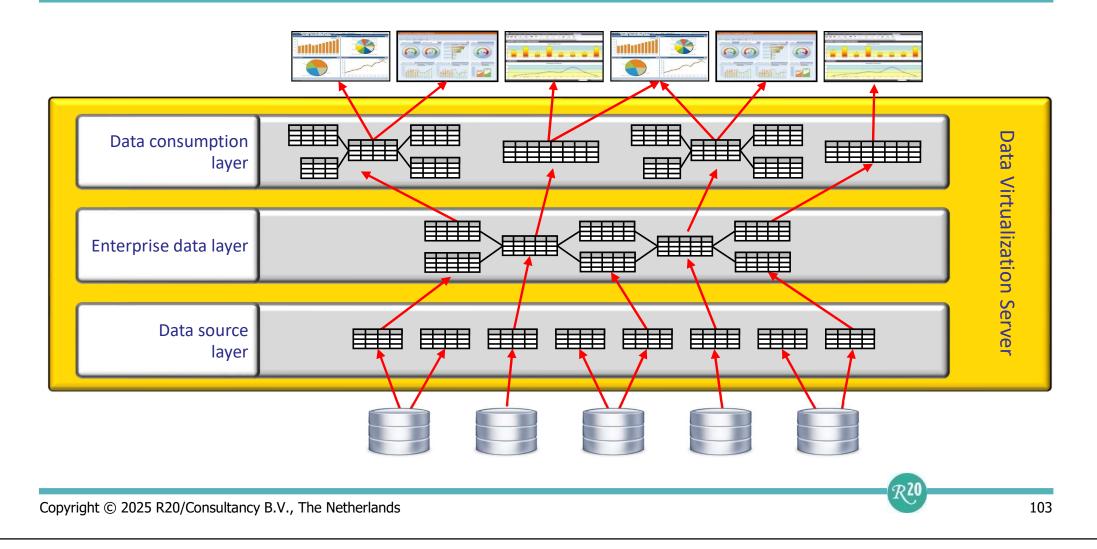
- AtScale
- DataVirtuality (Pipes, Pipes Prof, LDW)
- Denodo Platform
- Dremio
- Fraxses
- IBM InfoSphere Federation Server & IBM Data Virtualization Manager for z/OS (formerly Rocket Data Virtualization)
- Red Hat JBoss Data Virtualization (Teiid) ??
 - Stone Bond Enterprise Enabler Virtuoso
- TIBCOData Virtualization (formerly Cisco & Composite)
- Zetaris
- And many more ...



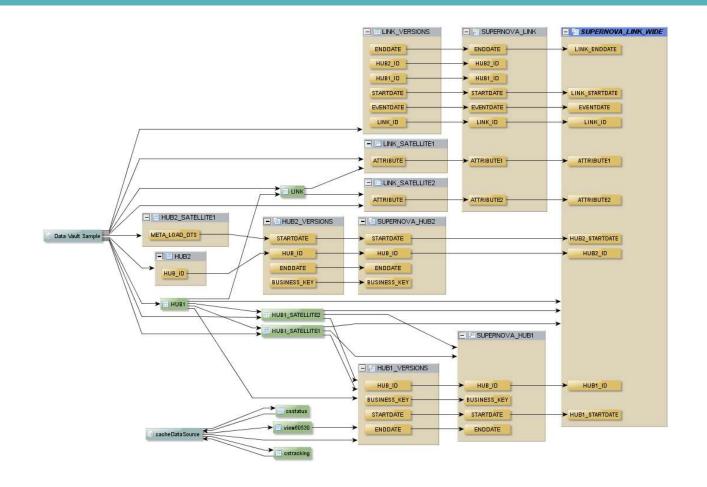
Developing Virtual Tables



Layers of Virtual Tables



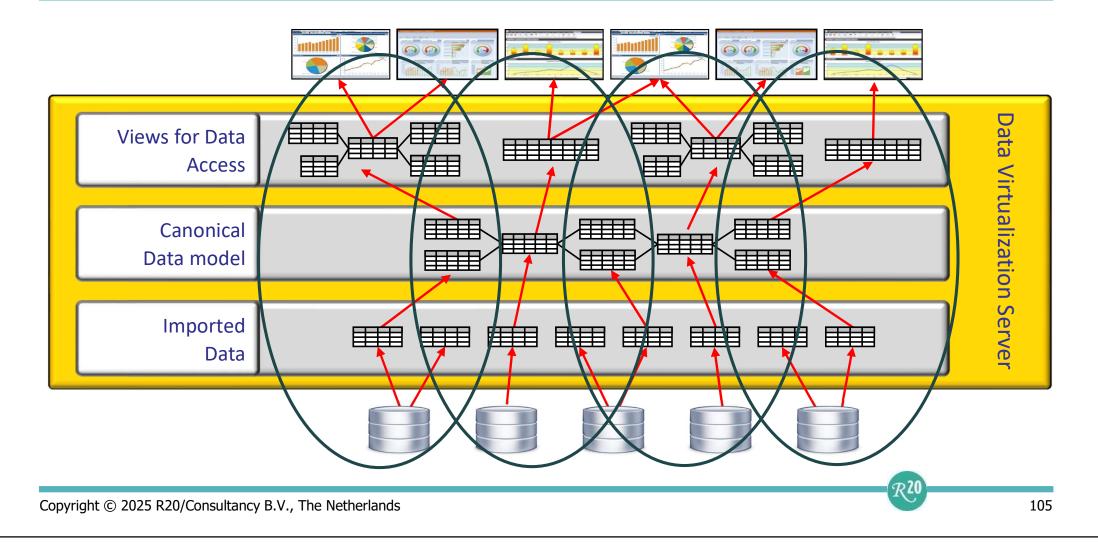
Lineage and Impact Analysis



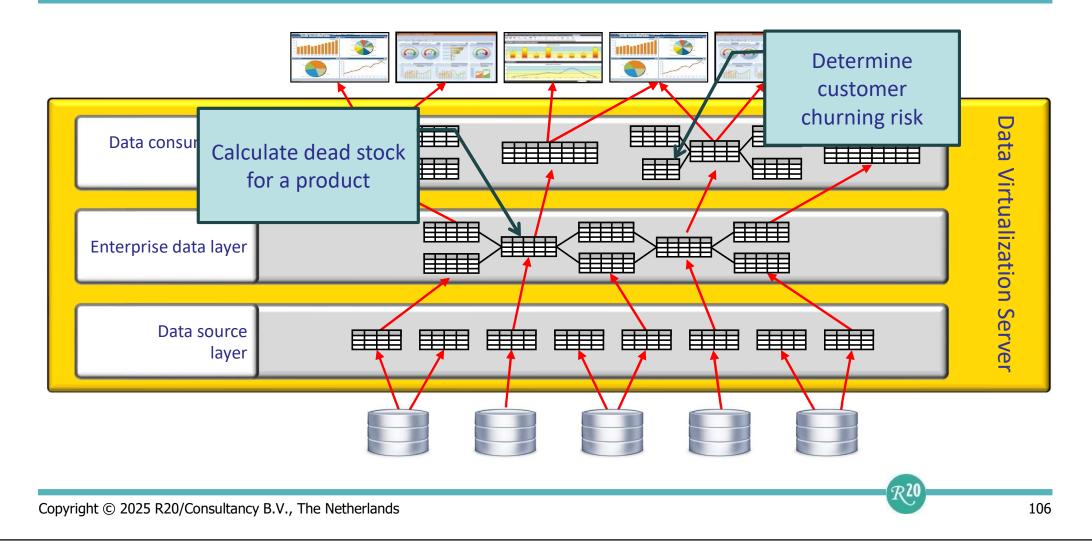
Copyright © 2025 R20/Consultancy B.V., The Netherlands

 \mathcal{R}^{20}

Evolutionary Development Approach



Improved Productity Through Sharing



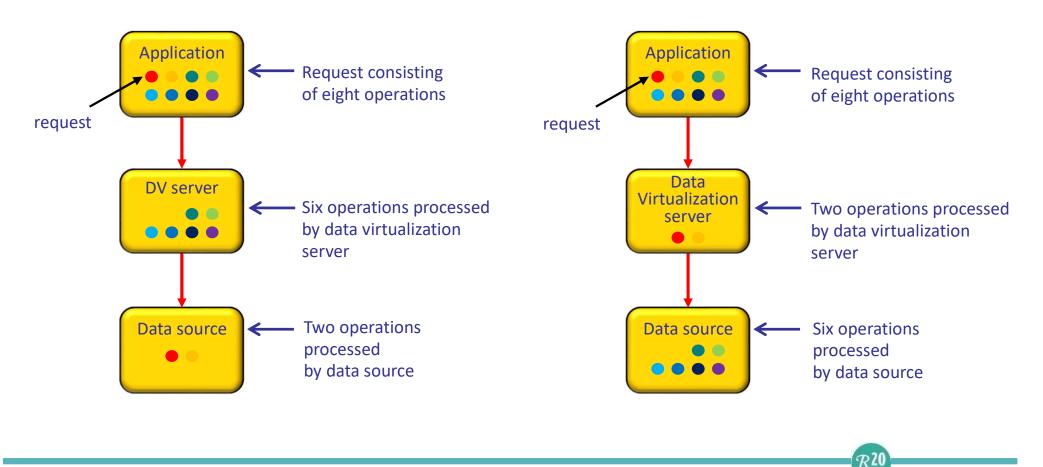
Performance Improving Features



- Easy-to-optimize queries
- Environment setup
- Query optimization
- Parallel processing and parallel pushdown
- Caching virtual tables
- The network
- Efficient drivers and connectors

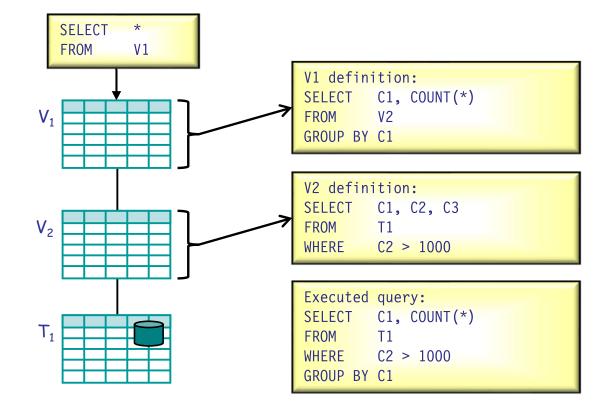


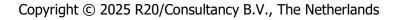
Improved Performance Through Query Pushdown



108

Many Levels and One Query







Push Down Query Processing

ation	(1) Incoming Query:	(3) Executed Query:
Data Virtualization Server	SELECT C2, CONCAT(C5, C6) FROM Virtual table WHERE C1 = > 1000 AND C4 BETWEEN 10 AND 20	SELECT C2, CONCAT(C5, C6) FROM Result

e	(2) Executed Query:
Data Source	SELECT C2, C5, C6 FROM Table WHERE C1 = $>$ 1000 AND C4 BETWEEN 10 AND 20

Copyright © 2025 R20/Consultancy B.V., The Netherlands

Accessing Files

Virtualization Server	(1) Incoming Query:	(3) Executed Query:
Data Virtual Serve	SELECT C2, CONCAT(C5, C6) FROM Virtual table WHERE C1 = > 1000 AND C4 BETWEEN 10 AND 20	SELECT C2, CONCAT(C5, C6) FROM Result WHERE C1 = $>$ 1000 AND C4 BETWEEN 10 AND 20

Data Source	(2) Executed Query: SELECT C1, C2, C4, C5, C6 FROM File

Copyright © 2025 R20/Consultancy B.V., The Netherlands

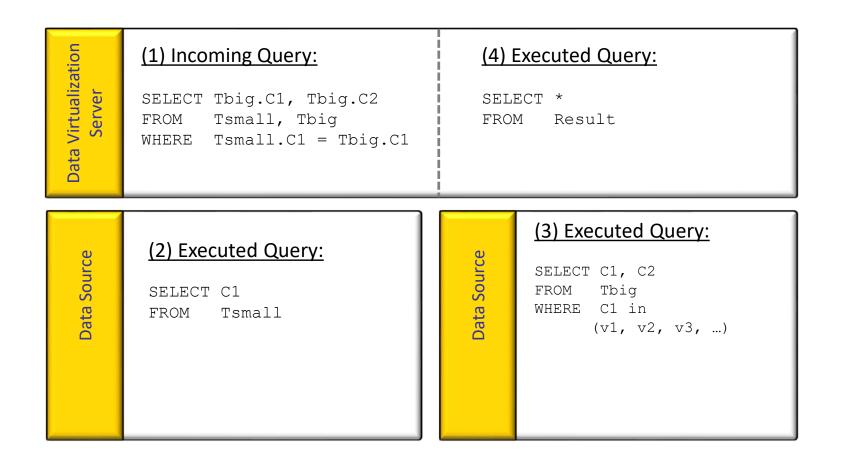
Inferred Filters

ation	(1) Incoming Query:	(3) Executed Query:
Data Virtualization Server	SELECT T1.C1, T1.C2, T2.C2 FROM T1, T2 WHERE T1.C1 = T2.C1 AND T1.C1 => 1000 AND T2.C2 BETWEEN 10 AND 20	SELECT T1.C1, T1.C2, T2.C2 FROM T1, T2 WHERE T1.C1 = T2.C1

U	(2a) Executed Query:	ce	(2b) Executed Query:
Data Source	SELECT C1, C2 FROM T1 WHERE C1 => 1000	Data Sourc	SELECT C1, C2 FROM T2 WHERE T2.C1 => 1000 AND T2.C2 BETWEEN 10 AND 20

Copyright © 2025 R20/Consultancy B.V., The Netherlands

Join of Data Sources with Query Injection



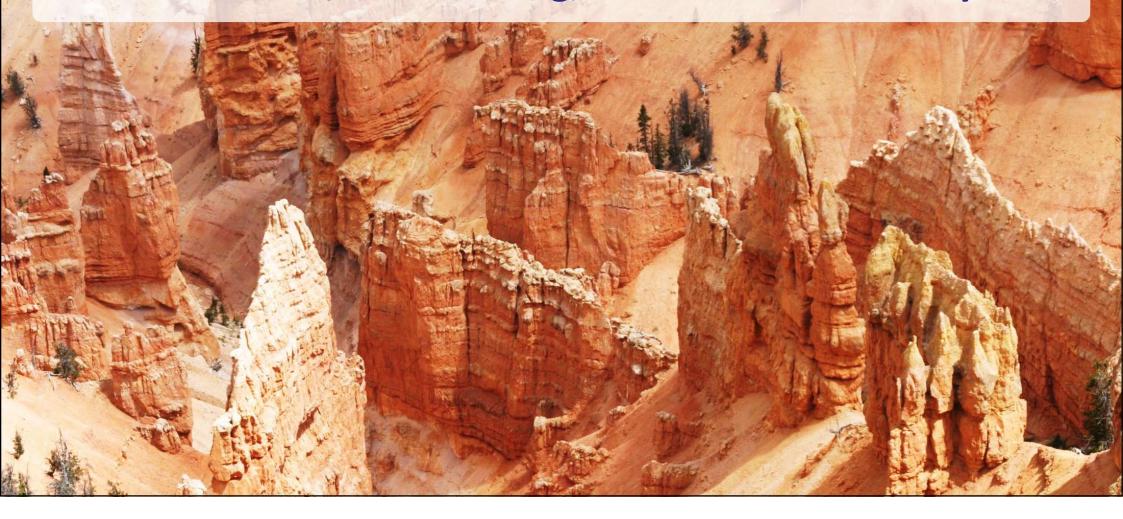


Join of Data Sources with Ship Join

Data Virtualization Server	(1) Incoming Query: SELECT T1.C1, T1.C2, T2. C2 FROM Tbig1, Tbig2 WHERE Tbig1.C1 = Tbig2.C1	<u>(4) Executed Query:</u> SELECT * FROM Result				
Data Source	(2) Executed Query: SELECT C1, C2 FROM Tbig1	(3) Executed Query: CREATE TEMP TABLE TEMP; INSERT INTO TEMP SELECT * FROM Tbig1; SELECT TEMP.C1, TEMP.C2, Tbig2.C2 FROM TEMP, Tbig2 WHERE TEMP.C1 = Tbig2.C1; DROP TEMPORARY TABLE TEMP;				



Part 4.3: Metadata, Data Catalog, and Business Glossary



Types of Metadata

Metadata on Data

- Textual definition
- Description
- Annotations by users and IT specialists
- Data lineage including transformations
- Retention information
- Qualifications: trustworthiness, completeness, data quality, ...
- Value descriptions
- Original or masked/anonymized
- Ontology
- Owner and support

...

Metadata on Metadata

- Retention information on metadata
- History of metadata
- Qualifications: trustworthiness, completeness, data quality, ...
- Metadata value descriptions
- Original or masked/anonymized metadata
- Owner and support
-

Copyright © 2025 R20/Consultancy B.V., The Netherlands



116

Metadata in 1976

DE DATA DICTIONARY/DIRECTORY (DD/D)

door L. Delport

1.1 Wat is een DD/D? (Data Dictionary/Directory) (gegevenskataloog)

Zeer algemeen kunnen we een DD/D als volgt bepalen:

Een DD/D is een katalogus die de omschrijving bevat van alle informatie-elementen die in een bedrijf bestaan.' Deze bepaling laat natuurlijk de weg open

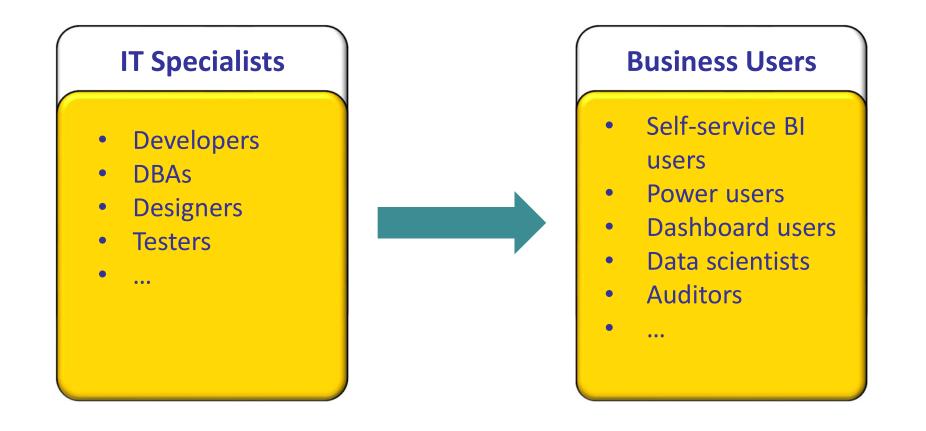
1.2 Waarom hebben we een DD/D nodig?

Centralisatie, het vermijden van dubbele gegevens en het opzoeken en vinden van informatie langs allerlei wegen en door middel van allerlei sleutels zijn technieken eigen aan M.I.S. (7) en systemen van gegevensbanken.

Informatie jaargang 18 nr. 7/8 pag. 430 t/m 491 Amsterdam juli/augustus 1976



A Target Audience Shift of Metadata



Numerous Solutions that Manage Metadata

- Home made metadata systems
- Professional data catalogs and business glossaries: Alation, Apache Atlas, Collibra, Informatica, TIBCO EBX, ...
- Scraping and linking: ASG, Manta, SQLdep (Collibra), ...
- Data warehouse automation: Astera, Attunity Compose, BiGenius, TimeXtender, WhereScape, ...
- BI tools: semantic layers
- ETL tools
- Data profiling tools
- Data quality tools
- Data virtualization servers: Data Virtuality, Denodo, Fraxses, Tibco DV, ...
- And many more ...

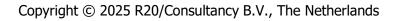


Study Metadata Needs of Data Consumers (1)

- Metadata on which objects are required by which user
 - Metadata on data elements, reports/dashboards, data science models, business rules, ...
- How do they want to access metadata?
 - Instant metadata (integrated within BI dashboard)
 - Via a service interface
 - Search interface on metadata
 - Do they need an ontology?
- Need for adding personal annotations to metadata
 - Annotation on tables and columns
 - Annotations on individual data values
 - Annotations on aggregated/derived values

Instant Metadata

Stock overview Cohelion															
Note: This is an aggregated re	eport. Most da	ta is first calo	culated at the	e native resol	ution (per gr	oup <mark>and poo</mark> l	I), and the res	sults are <mark>s</mark> un	med or aver	aged for an a	aggregated re	eport.			
	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec	Jan	Feb	Mar	Apr	May	Trend vs last year	Overall
Total demand	128,213	127,098	117,923	116,033	115,431	108,253	114,681	117,610	115,882	115,554	115,122	116,106	116,849		117,289
Total inventory	145,401	144,654	144,753	143,668	142,977	142,224	141,258	139,011	138,414	137,690	137,108	136,660	136,191		140,770
Total available inventory	111,562	111,719	113,504	114,175	114,342	112,682	113,106	112,081	110,908	109,726	108,631	107,575	106,545		111,274
Available Bycicles	100,653	101,344	102,974	103,659	103,766	101,010	101,180	100,155	98,982	97,800	96,705	95,649	94,619		99,884
Rented Bycicles	10,909	10,375	10,530	10,516	10,576	11,672	11,926	11,926	Available	Bycicles	C. Haller	,926	11,926		11,389
Unservicable Bycicles	6,490	6,528	7,793	6,931	7,094	7,716	7,990	7,040	Jan			,059	7,158		7,142
In repair Bycicles	2,445	2,508	2,635	2,809	2,607	2,749	2,499	2,744	Forecast 98	8,982 tered manua	lly Ecroport	,688	2,667		2,652
Storage Bycicles	1,223	1,652	1,526	3,836	4,205	4,568	4,934	2,515	is the last m	ionth repeate	ed, but adjust		3,072		2,981
Lost / Stolen Bycicles	23,681	22,247	19,295	15,917	14,729	14,509	12,729	14,631	of that mont	th. Formula: /	and decreas Available Byc s - Sales Car	icles 388	16,749		16,721
Rebrand In	0	0	0	0	0	0	0	24		values are ta	s - Sales Cap aken from the		ξ ₁		0
Rebrand Out	0	0	0	0	0	0	0		previous m				×,		0
Lease In		0	0	20	0	25	116								20



R20

Study Metadata Needs of Data Consumers (1)

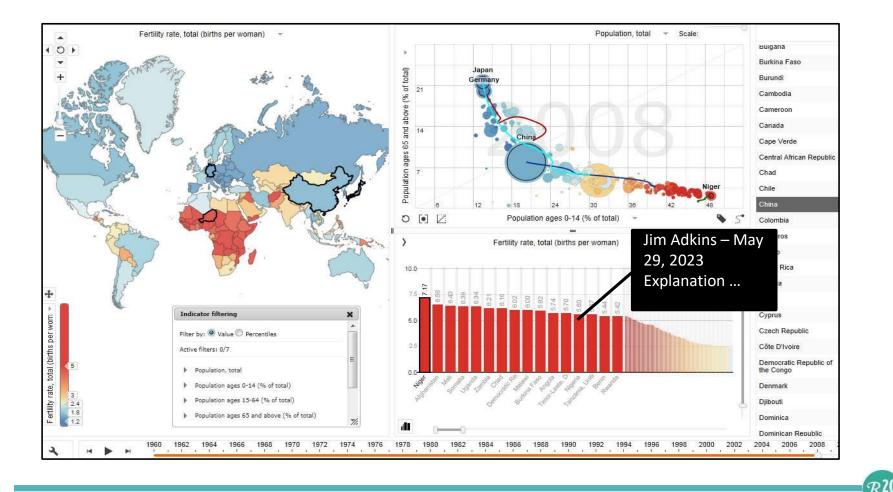
- Metadata on which objects are required by which user
 - Metadata on data elements, reports/dashboards, data science models, business rules, calculations, retention specifications, ...

How do they want to access metadata?

- Instant metadata (integrated within BI dashboard)
- Via a service interface
- Search interface on metadata
- Do they need an ontology?
- Need for adding personal annotations to metadata
 - Annotation on tables and columns
 - Annotations on individual data values
 - Annotations on aggregated/derived values



Adding Annotations on Data Points



Study Metadata Needs of Data Consumers (2)

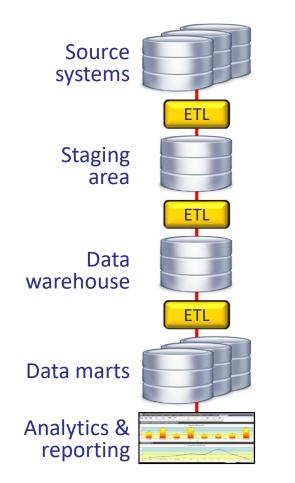
Metadata tagging

- Retention period, real or anonymized, responsible data steward, trustworthiness, completeness, data quality, ...
- Taggings added by business users
- Versioning of metadata
- Automatic notifications
 - a retention period of frequently used data is about to expire
 - a definition has changed
 - a piece of legislation is about to change
- Fuzzy boundaries between metadata and master data
 - Are all the state codes metadata or master data?



Copyright $\ensuremath{\mathbb{C}}$ 2025 R20/Consultancy B.V., The Netherlands

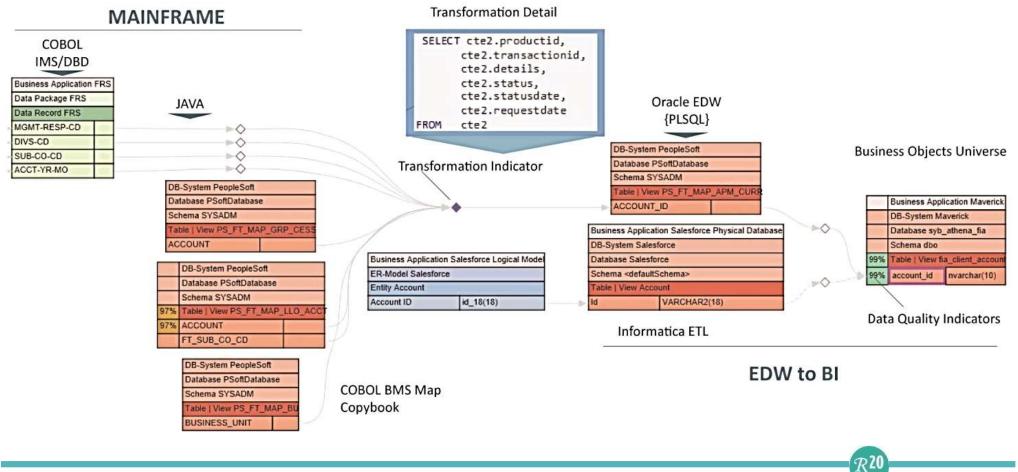
Metadata is Dispersed



- Metadata dispersed across many systems
 - In database servers (system tables)
 - In integration tools
 - In documentation
 - In reporting tools (semantic layer)
 - In spreadsheets
 - In application code
 - And many more ...
- Most is technical and not business metadata
- Not integrated no clear relationships between metadata elements



Example: Lineage Through Scraping by ASG



Example: Lineage Through Scraping by SQLdep (Collibra)



Part 4.4: **Master Data**

Master Data Management: Two Customer Tables

Customer table in **Sales** System

ID	Name	Initials	Date Entered	City	State
12345	Young	Ν	Aug 4, 2008	San Francisco	CA
23324	Stills	S	Sep 10, 2009	New Orleans	LA
57657	Furay	R	Oct 16, 2010	Yellow Springs	ОН
65461	Palmer	В	Nov 22, 2011	Boston	MA

Customer table in Finance System

ID	Name	Initials	Date Entered	City	State
C5729	Young	Ν	Sep 16, 2007	San Francisco	СА
LA781	Stils	S	Dec 8, 2010	New Orleans	LA
J7301	Furay	R	Jan 10, 2008	Yellow Springs	ОН
К8839	Palmer	В	Feb 11, 2009	New York	NY

Copyright $\ensuremath{\textcircled{O}}$ 2025 R20/Consultancy B.V., The Netherlands



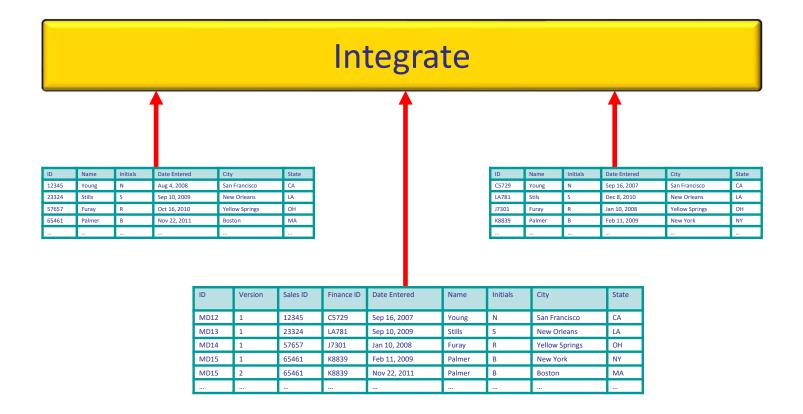
The Master Customer Table

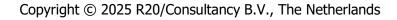
Master Customer table

ID	Version	Sales ID	Finance ID	Date Entered	Name	Initials	City	State
MD12	1	12345	C5729	Sep 16, 2007	Young	Ν	San Francisco	CA
MD13	1	23324	LA781	Sep 10, 2009	Stills	S	New Orleans	LA
MD14	1	57657	J7301	Jan 10, 2008	Furay	R	Yellow Springs	ОН
MD15	1	65461	K8839	Feb 11, 2009	Palmer	В	New York	NY
MD15	2	65461	K8839	Nov 22, 2011	Palmer	В	Boston	MA

Copyright $\ensuremath{\textcircled{O}}$ 2025 R20/Consultancy B.V., The Netherlands

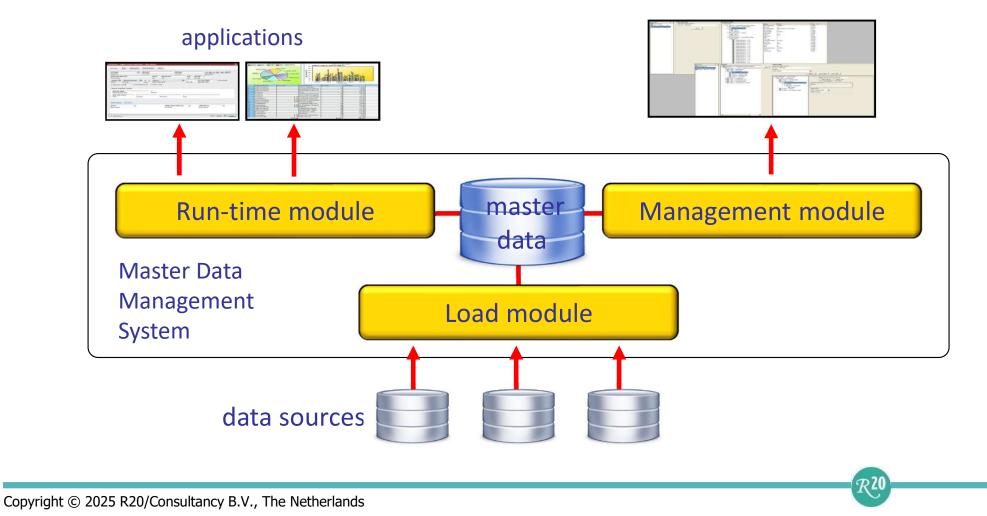
Joining Needs Master Data



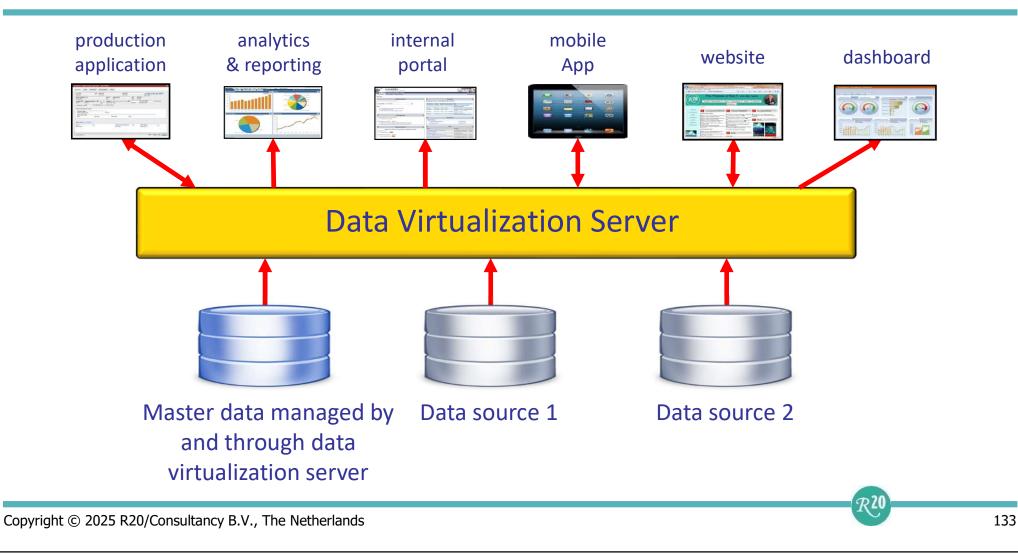




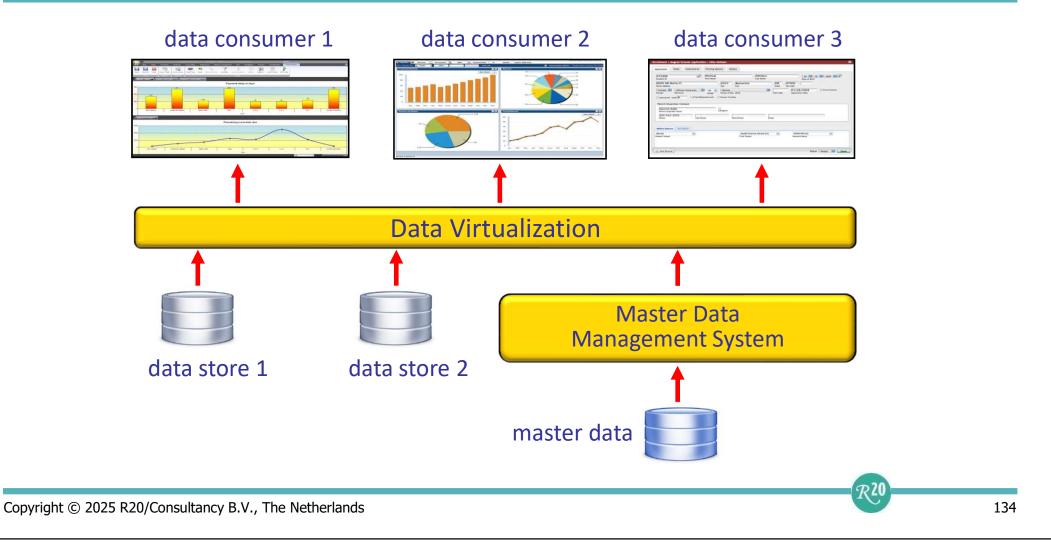
Overall Architecture of an MDM System



Lightweight Solution for Master Data



MDM as Source for Data Virtualization

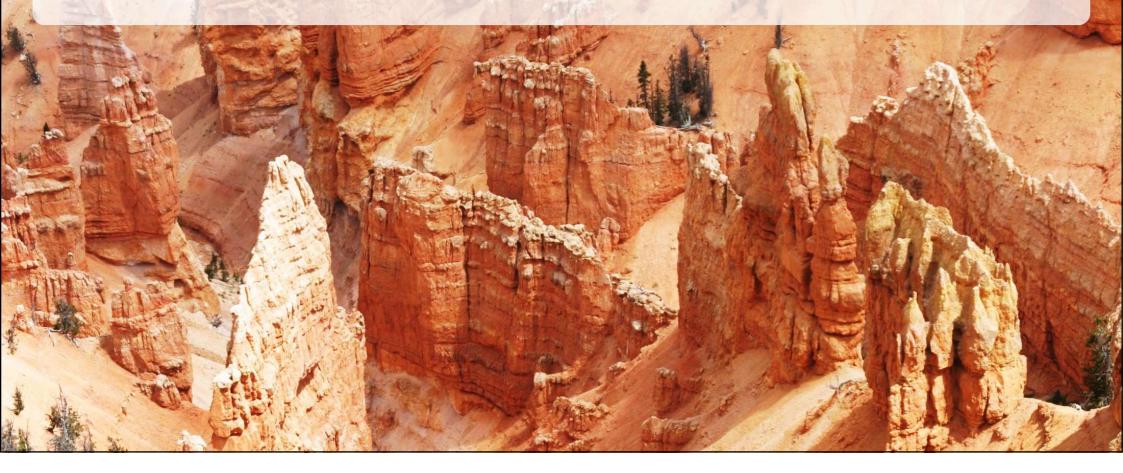


Example: Cohelion (Master Data-driven BI)

Accounts Payable	30,809	35,972	37,370	37,668	37,346	33,360	34,141	42,505
Profit - Loss Accounts								
Gross profit	79,150	56,450	61,467	77,134	Actual Manual SalesForce - Feed SAP - Feed		185,254.00 20,248.00	60,400
EBIT	11,327	-732	-1,247	11,101			30,855.00 134,151.00	-374
Sales Income	201,736	167,789	170,810	222,655	199,000	185,254	7,600	131,000
Rental Income	0	0	0	0	0	0	0	0
Other income	0	o	0	0	0	0	0	0
Expense Accounts								
Office Expense	812	681	790	532	960	372	1,010	455

Copyright © 2025 R20/Consultancy B.V., The Netherlands

Part 4.5: Incorporating Cloud in the Data Architecture



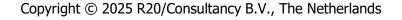
Cloud Platforms are Becoming the New Mainframes



Mainframe = Lock In



- Proprietary operating systems
- Proprietary system management software
- Proprietary database servers
- Proprietary security systems
- Proprietary development environments
- Proprietary JCLs
- Proprietary ...



Cloud Platform = Lock In?



- Proprietary operating systems
- Proprietary management software
- Proprietary database servers
 - E.g. Amazon: RDS, RedShift (SQL), S3, ...
- Proprietary security systems
- Proprietary development environments
 - E.g. Microsoft Azure: Reporting Services, Analytics services, Data Management Services, ...

Proprietary ...



Data Storage Technologies Available on Cloud Platforms

Cloud Platform	Data Storage Technology
Amazon AWS	Aurora DocumentDB DynamoDB Elasticache for Redis RDS Redshift S3 Timestream
Google	BigQuery Cloud Bigtable Cloud Firestore Cloud Spanner Cloud SQL
Microsoft Azure	Cache for Redis Cosmos DB Data Lake SQL Database Synapse Analytics

Copyright © 2025 R20/Consultancy B.V., The Netherlands



Stay Cloud Platform Independent

Design to Migrate

Watch Out For Egress Costs!

Public Cloud	Typical Data Egress Charge (per GB)		Discounted Data Egress Charge (per GB) for 100 TB	Cost to move 100 TB, per month
Azure	\$0.08	\$800	\$0.07	\$7,000
AWS	\$0.02	\$200	\$0.02	\$2,000
Google Cloud Platform	\$0.11	\$1100	\$0.08	\$8,000
Oracle	Free up to 10TB	free	\$0.0085-\$0.050 depending on geography	\$850 to \$5,000

Source: https://www.factioninc.com/blog/it-challenges/egress-charges-how-to-prevent-costs/



 \mathcal{R}^{20}

Copyright $\ensuremath{\textcircled{O}}$ 2025 R20/Consultancy B.V., The Netherlands

Cloud Platform Fees



- Fees can have influence on data architectureExample:
 - SnowflakeDB: pay for data usage (queries)
 - Store more derived data
 - Exasol: pay for environment size (queries for free)
 - Work with views in stead of physical data marts
- How well can the technology exploit the cloud platform?
 - E.g. cloud is endless MPP, what about the database server?
- Pushing processing into the cloud, close to where data is produced

Part 5: Step 5: Define Architectural Design Principles



Forget Old Architectural Design Principles



- No reporting on the production database
 - Reporting and transaction workloads clash
- Physical data marts are needed to improve reporting performance
- Data marts need a star schema design to speed up analytical queries
- ETL is used to transform data
 - Batch oriented
- When SQL databases are used
 - Indexes are required to improve query performance
 - Use locking for concurrency management
 - Not ideal for MPP
 - Need constant tuning by DBA

Examples of Architectural Design Principles



- Centralized and active data processing specifications
 - Searchable definitions and descriptions for technical and business users
 - Lineage and impact analysis
- One universal architecture for all forms of data consumption
 - Standard reporting, self-service BI, apps, data science, ...
- Data storage and access technology agnostic
 - Hadoop, SQL, cubes, ...
 - Abstraction
- Push the processing to the data, not the data to the processing
 - Decentralized data production
 - Edge analytics
 - Hyper-decentralized data production and storage
- Generator-driven

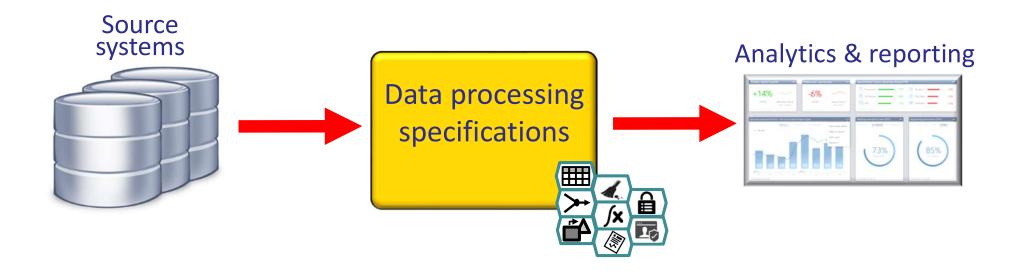


146

Part 5.1: Data Processing Specifications



The Data Processing Specifications





Data structure specifications Integration specifications Transformation specifications Data security specifications



Data cleansing specifications Analytical specifications Visualization specifications Data privacy specifications



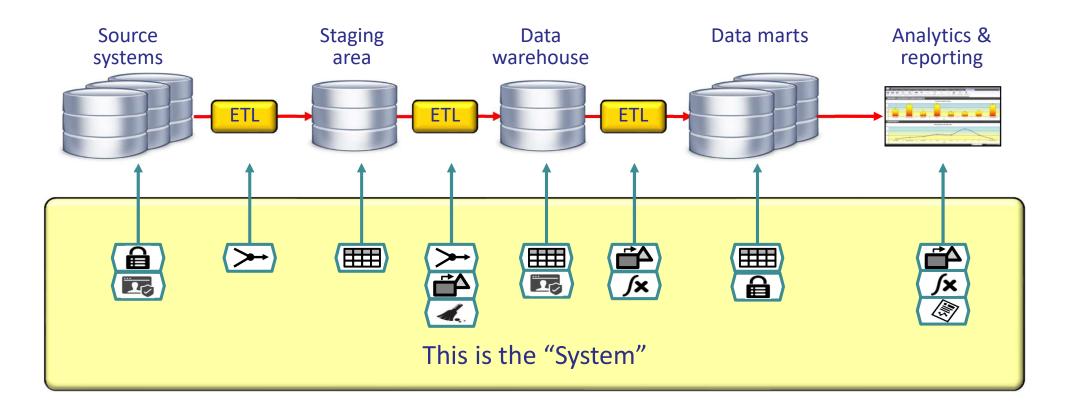
Examples of Data Processing Specifications



- Data value transformations
- Data structure transformations
- Aggregations
- Filters
- Calculations
- Integrations
- Technical corrections
- Functional corrections
- Anonymizations
- Authorizations and authentications
- Historizations
- Metadata-related specifications
- ...

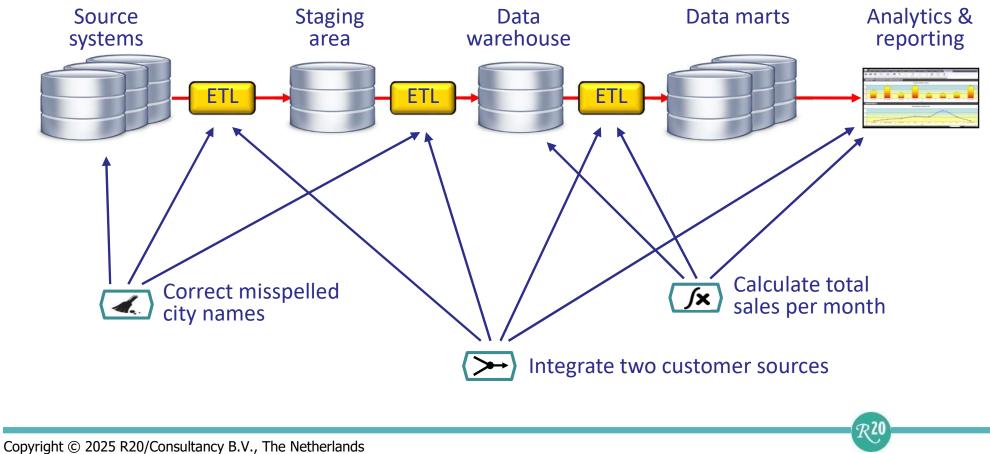


Data Processing Specifications

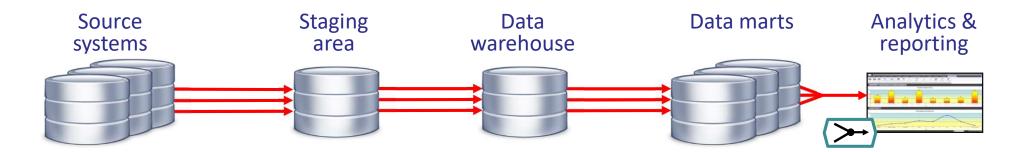




Where to Implement Data Processing Specifications?

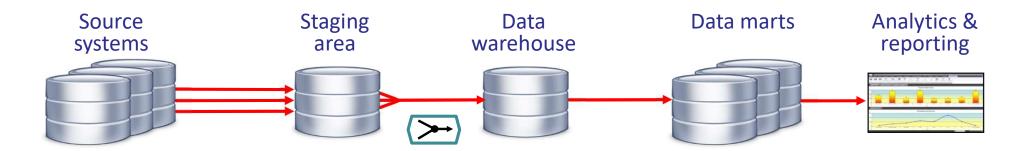


Example: Integration and Aggregation (1)



- Slow processing of integration logic in reports
- Complex queries in reports
- No sharing of integration logic across reports and tools
- Potential errors and inconsistencies in reports
- Fast copying and lower data latency
- Data structure of source database determines all data structures
- Use of original raw data possible
- What about integration errors?

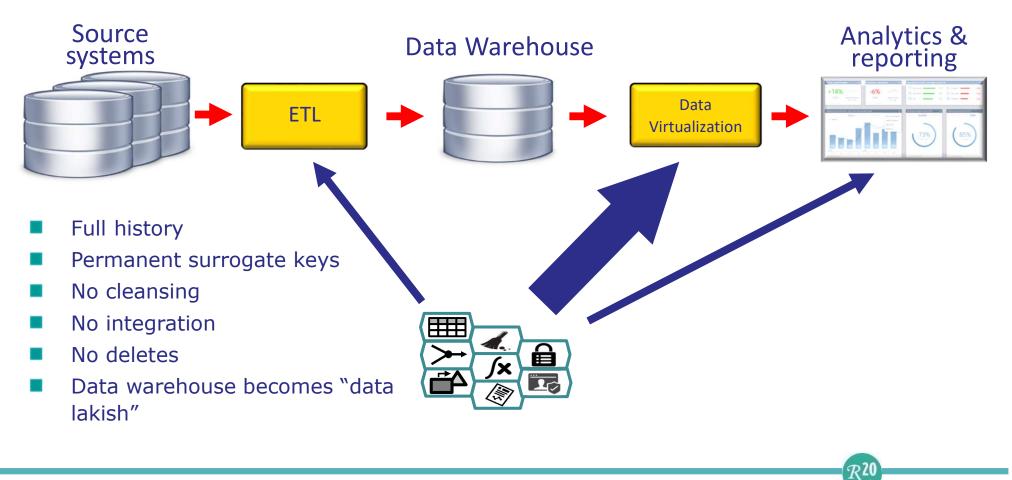
Example: Integration and Aggregation (2)



- Fast processing of integration logic in reports
- Simpler queries in reports
- Sharing of integration logic across reports and tools
- Potential errors and inconsistencies in ETL
- Slower copying and higher data latency
- Data structure of source database does not determine all data structures
- Use of original raw data possible
- Integration errors easier to fix



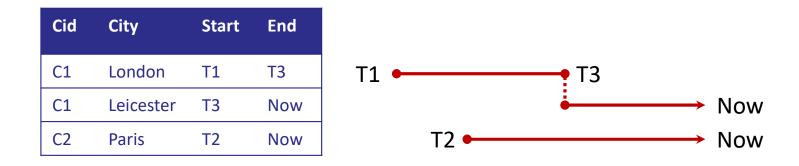
Implementing the Data Processing Specifications



Part 5.2: Dealing with History

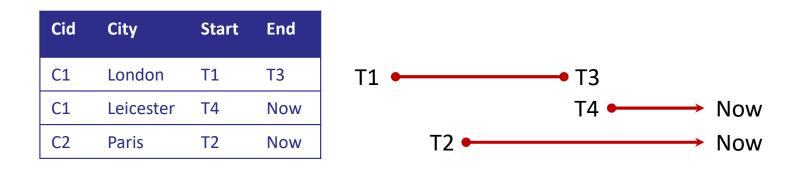


Modeling History: Simple History for Updates



- No gaps in history
- Only one value for an object on a specific datetime
- Supports following queries:
 - What is the current value Where End = Now
 - What was the value on a specific datetime Where date between Start and End

Modeling History: Simple History for Updates with Gaps

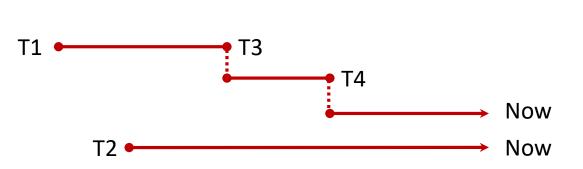


- No gaps in history
- Only one value for an object on a specific datetime
- Supports following queries:
 - What is the current value Where End = Now
 - What was the value on a specific datetime Where date between Start and End may return no values

157

Modeling History: Simple History for Updates Without Gaps

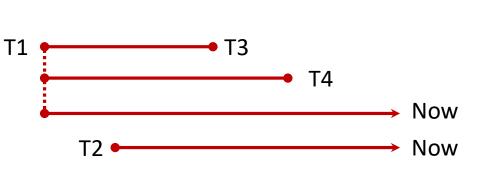
Cid	City	Start	End
C1	London	T1	Т3
C1	Gap	Т3	T4
C1	Leicester	T4	Now
C2	Paris	T2	Now



- No gaps in history
- Only one value for an object on a specific datetime
- Supports following queries:
 - What is the current value Where End = Now
 - What was the value on a specific datetime Where date between Start and End always returns a value; sometimes nothing

Modeling History: Corrections

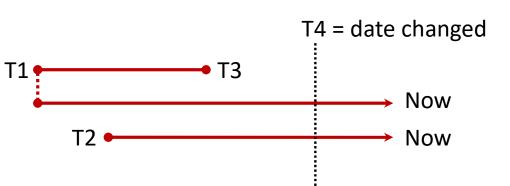
Cid	City	Start	End
C1	Londdon	T1	Т3
C1	Londn	T1	T4
C1	London	T1	Now
C2	Paris	T2	Now



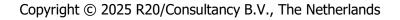
- No gaps in history
- Multiple values for an object on a specific datetime
- Supports following queries:
 - What is the current value Where End = Now
 - What was the value on a specific datetime Get the oldest where date between Start and End

Modeling History: Delayed Corrections

Cid	City	Start	End	Changed
C1	Londdon	T1	Т3	T4
C1	London	T1	Now	
C2	Paris	T2	Now	



- Extra column required
- No gaps in history
- Multiple values for an object on a specific datetime
- Supports following queries:
 - What is the current value Where End = Now
 - What was the value on a specific datetime Get the oldest where date between Start and End



Modeling History: Logging Updates and Corrections

Cid	City	Start	End	Changed	Insert id	Change id
C1	London	T1	T2		Insert1	Update1
C1	Leicester	T2	Now		Insert2	
C2	Paris	Т3	T4		Insert3	Update2
C2	Lyon	Т4	T5		Insert4	Delete1

Change id	Who	When	Where	•••
Insert1	User1	T1		
Insert2	User2	T2	•••	
Insert3	User1	Т3		
Insert4	User3	Τ4		
Update1	User2	T2	•••	
Update2	User4	T4		
Delete1	User5	T5		

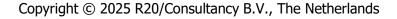
- Log table for auditing purposes
- Batch inserts, updates, and deletes

Modeling Streaming Data: Single Values

Incoming Stream	Кеу	S
T1, S1, Temp=50	1	
T2, S1, Temp=52	2	Т
T3, S2, Temp=51	5	Т
T4, S3, Temp=49	3	
T5, S1, Temp=52;	6	Т
T5, S2, Temp=53	4	

Кеу	Start	End	Sensor	Тетр	Avg Temp
1		T1	S1	50	50
2	T1	T2	S1	52	51
5	T2	T5	S1	52	51,3
3		Т3	S2	51	51
6	Т3	T5	S2	53	52
4		T4	S3	49	49

- Key is unique artificial value
- Measurement is considered as temperature since previous measurement
- Sensor data is arriving in the right order



Modeling Streaming Data: Multiple Values

Incoming Stream
T1, S1, Temp=50
T2, S1, Temp=52
T3, S2, Temp=51
T4, S3, Temp=49
T5, S1, Temp=52; S2, Temp=53

Кеу	Start	End	Sensor	Тетр	Avg Temp
1		T1	S1	50	50
2	T1	T2	S1	52	51
5	T2	T5	S1	52	51,3
3		Т3	S2	51	51
6	Т3	T5	S2	53	52
4		T4	S3	49	49

Key is unique artificial value

- Stream records are flattened
- Measurement is considered as temperature since previous measurement
- Sensor data is arriving in the right order



Modeling Streaming Data: Delta Values

Incoming Stream				
T1, S1, Temp=50				
T2, S1, Temp=+2				
T3, S2, Temp=51				
T4, S3, Temp=49				
T5, S1, Temp=+0				
T5, S2, Temp=+2				

Кеу	Start	End	Sensor	Тетр	Avg Temp
1		T1	S1	50	50
2	T1	T2	S1	52	51
5	T2	T5	S1	52	51,3
3		Т3	S2	51	51
6	Т3	T5	S2	53	52
4		T4	S3	49	49

- Key is unique artificial value
- Measurement is considered as change in temperature
- Sensor data is arriving in the right order

Modeling Streaming Data: Log Data

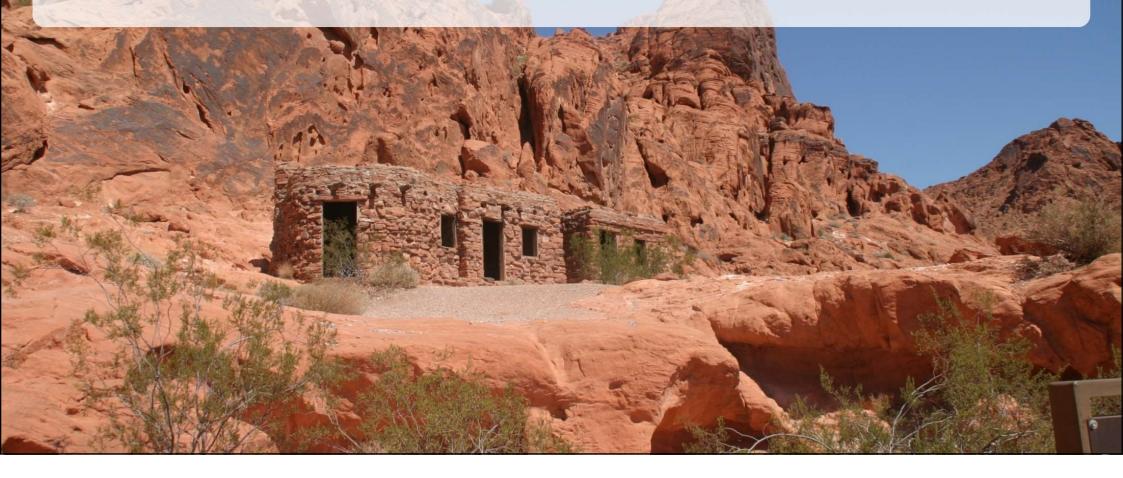
Incoming Stream
T1, Insert, C1, London
T2, Update, C1, Leicester
T3, Insert, C2, Paris
T4, Insert, C3, Berlin
T5, Update, C1, Manchester
T5, Update, C3, Munich

Ci	d C	City	Start	End
C1	. L	ondon	T1	Т2
C1	. L	eicester	Т2	T5
C1	. N	Machester	T5	Now
C2	. P	Paris	Т3	Now
C3	e E	Berlin	T4	T5
C3	S N	Junich	T5	Now

- Business key used
- Stream is seen as data entry
- Careful with parallel inserts; order not unimportant
 - Loading with hashed keys?



Part 5.3: Supporting GDPR

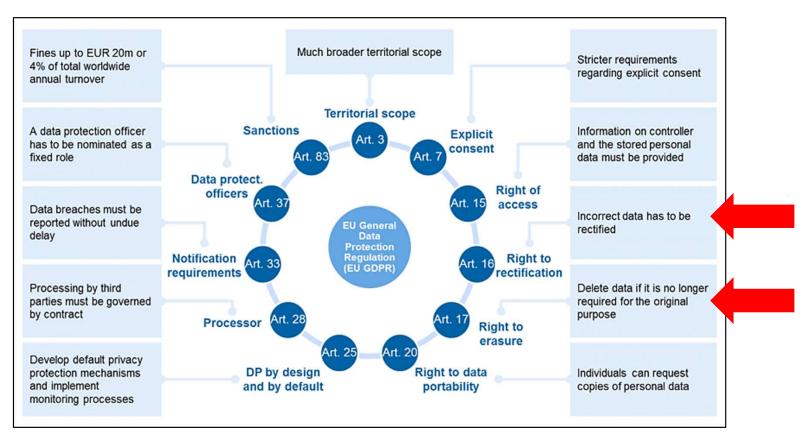


GDPR – The Right to be Forgotten



- Art. 17 GDPR Right to erasure (right to be forgotten)
 - The rule primarily regulates erasure obligations
- According to this, personal data must be erased immediately where the data are no longer needed for their original processing purpose, or the data subject has withdrawn his consent and there is no other legal ground for processing, the data subject has objected and there are no overriding legitimate grounds for the processing, or erasure is required to fulfill a statutory obligation under the EU law or the right of the Member States.
- In addition, data must naturally be erased if the processing itself was against the law in the first place
- A data subject should have the right to have personal data concerning him or her rectified

Requirements of GDPR



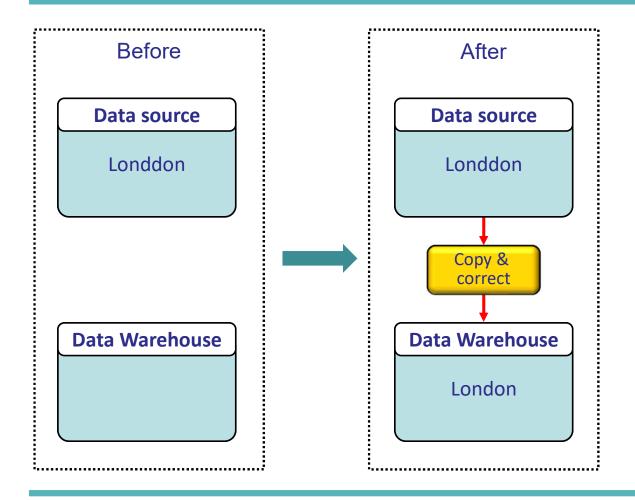
Source: Banking Hub, November 2017; see https://www.bankinghub.eu/banking/finance-risk/gdpr-deep-dive-implement-right-forgotten



Part 5.4: Dealing with Incorrect Data



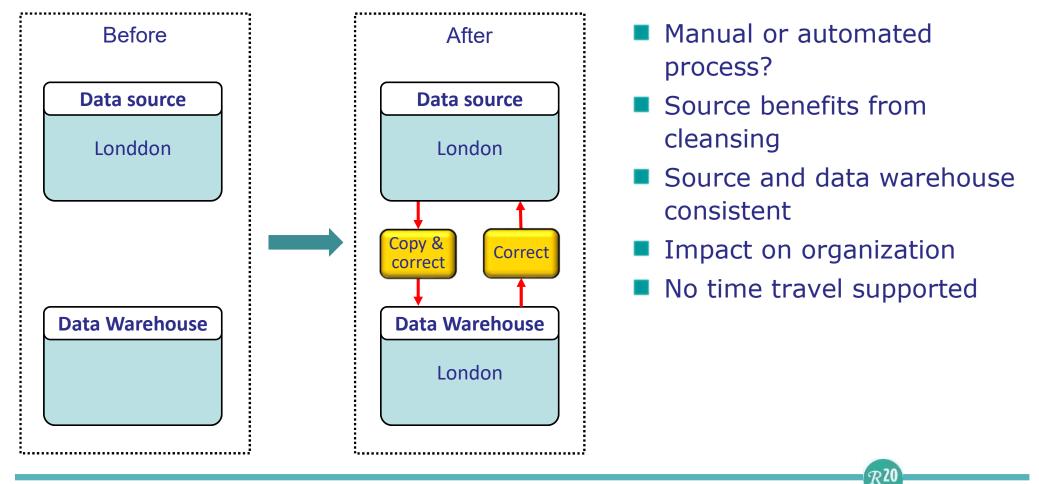
Data Correction Strategies: Simple



- Restricted to programmable cleansing operations
- Easy to implement
- Source doesn't benefit from cleansing
- Source and data warehouse inconsistent
- No impact on organization
- No time travel supported

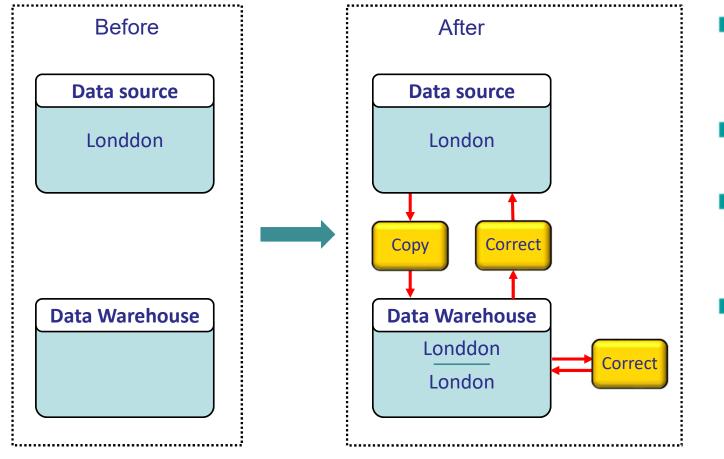
Copyright $\ensuremath{\textcircled{C}}$ 2025 R20/Consultancy B.V., The Netherlands

Data Correction Strategies: Synchronize



171

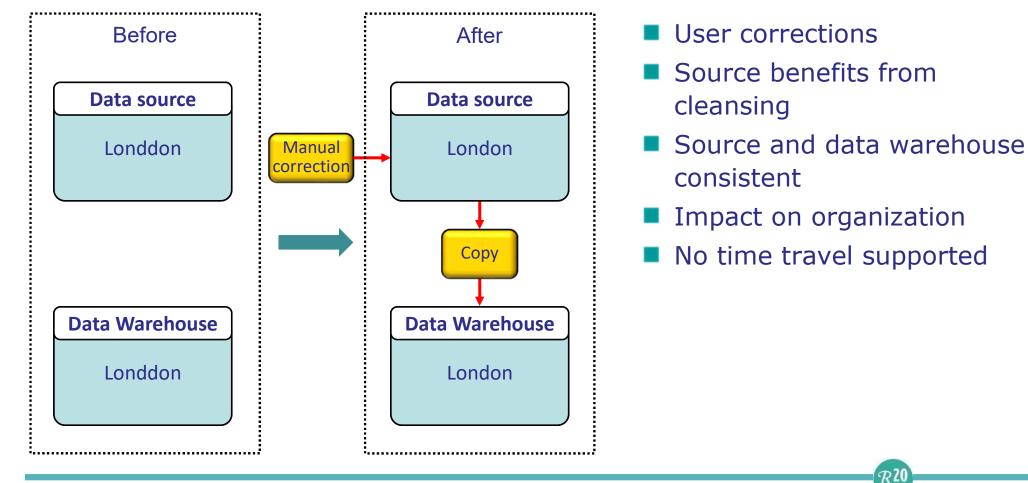
Data Correction Strategies: Time Travel



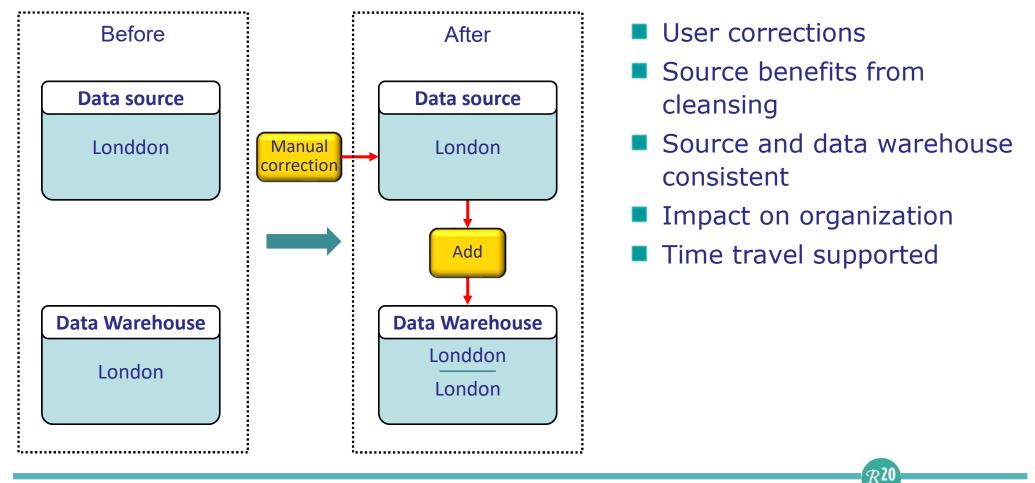
- Restricted to programmable cleansing operations
- Source benefits from cleansing
- Source and data warehouse consistent
- Time travel supported

Copyright $\ensuremath{\textcircled{C}}$ 2025 R20/Consultancy B.V., The Netherlands

Data Correction Strategies: Manual Corrections



Data Correction Strategies: Manual Corrections + Time Travel



Part 6: Step 6: Select a Reference Data Architecture



Roadmap for Designing Data Architectures

1. Determine business motivations 2. Determine new requirements 3. Analyze the existing environment 4. Study new products and technologies 5. Define architectural design principles 6. Select a reference data architecture 7. Design the new data architecture 8. Determine the Implementation approach 9. Select new products and technologies 10. Introduce the data architecture within the organization

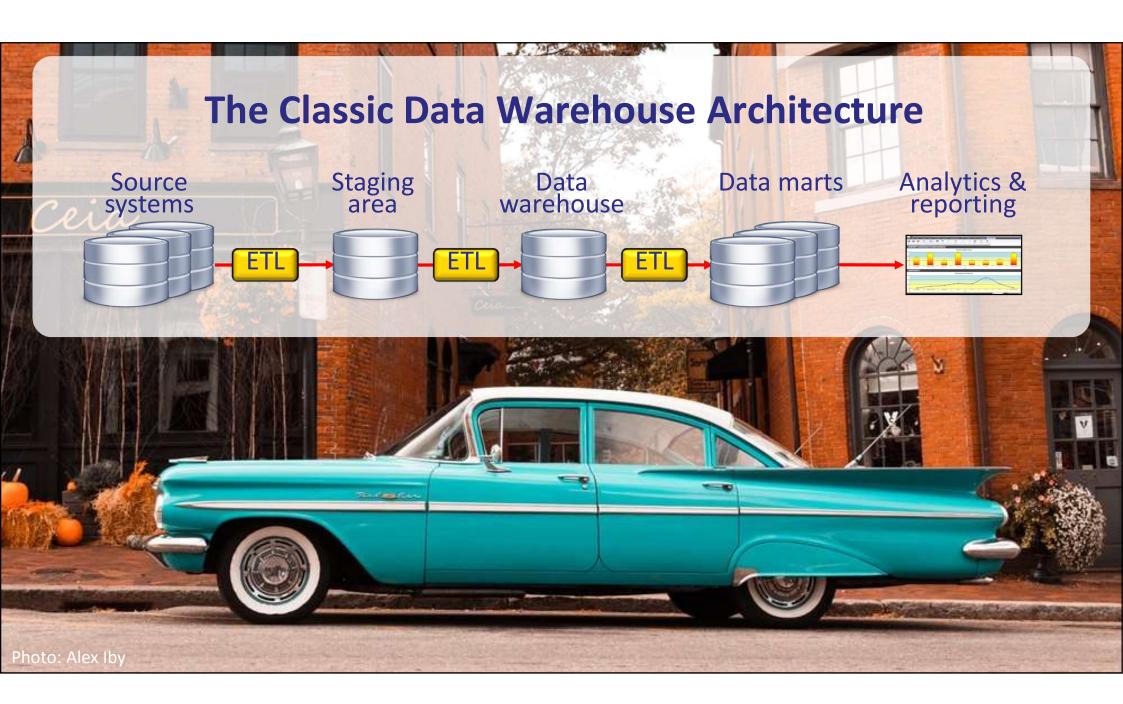
Common Challenges

- Source data must be queryable
- Developers and data consumers can't find data easily
- Every insert, update, delete and query should be logged for reconstruction purposes and transparency
- Horizontal and operational lineage
- CRUD interface for real-time synchronization
- Centralized, reusable, versioned business logic
- Proper authorization when integrating data

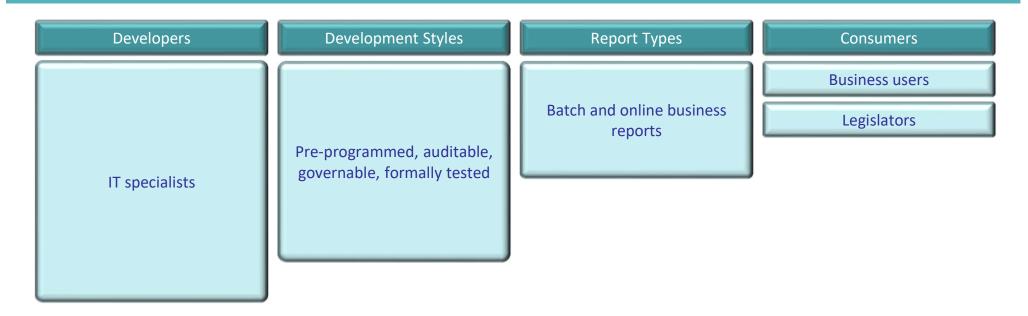


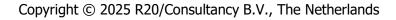
Part 6.1: The Classic Data Warehouse Architecture





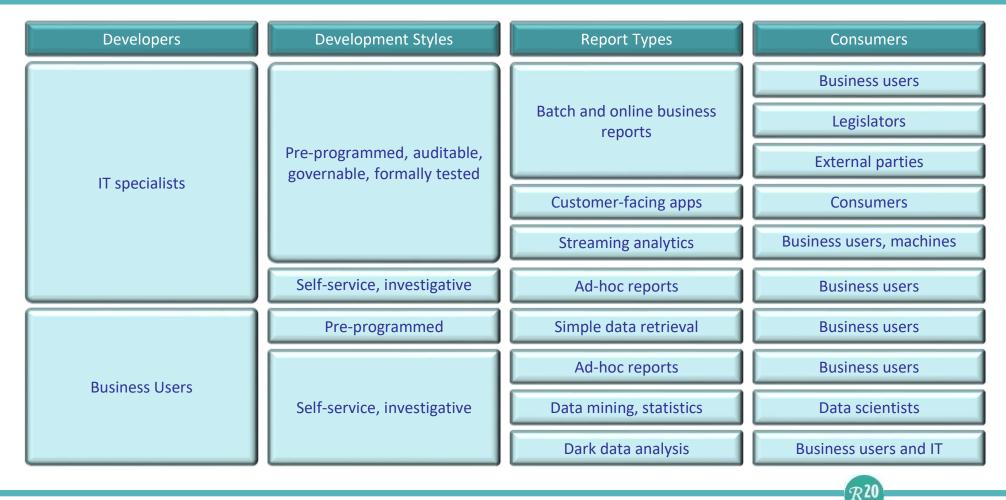
Yesterday: Data Warehouse and Data Consumption







Today & Tomorrow: Data Warehouse and Data Comsumption

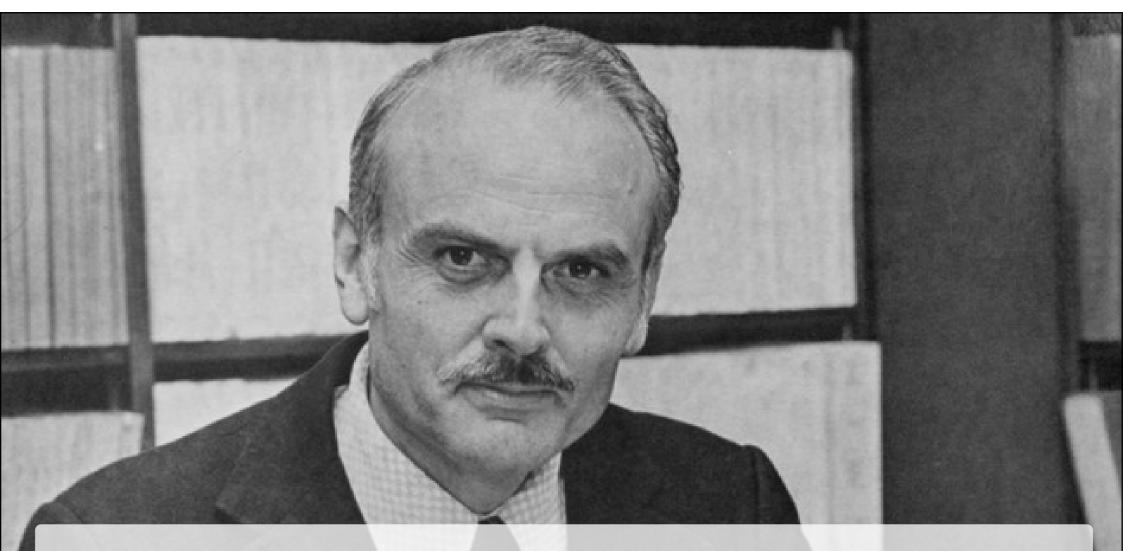


The Classic Data Warehouse Architecture is Like a Rigid Assembly Line



Part 6.2: The Logical Data Warehouse Architecture





Edgar Frank "Tedd" Codd

Ted Codd – June 1970

Future users of large data banks must be protected from having to know how the data is organized (...) application programs should remain unaffected when the internal representation of data is changed

Source: https://cs.uwaterloo.ca/~david/cs848s14/codd-relational.pdf

Ted Codd on Data Independence



The 1981 ACM Turing Award was presented to Edgar F. Codd, an IBM Fellow of the San Jose Research Laboratory, by President Peter Denning on November 9, 1981 at the ACM Annual Conference in Los Angeles, California. It is the Association's foremost award for technical contributions to the computing community.

Codd was selected by the ACM General Technical Achievement Award Committee for his "fundamental and continuing contributions to the theory and practice of database management systems." The originator of the relational model for databases, Codd has made further important contributions in the development of relational algebra, relational calculus, and normalization of relations.

Edgar F. Codd joined IBM in 1949 to prepare programs for the Selective Sequence Electronic Calculator. Since then, his work in computing has encompassed logical design of computers (IBM 701 and Stretch), managing a computer center in Canada, heading the development of one of the first operating systems with a general multiprogramming capability, contributing to the logic of selfreproducing automata, developing high level techniques for software specifica-

tion, creating and extending the relational approach to database management, and developing an English analyzing and synthesizing subsystem for casual users of relational databases. He is also the author of *Cellular Automata*, an early volume in the ACM Monograph Series.

The 1981 ACM Turing Award Lecture Delivered at ACM '81, Los Angeles, California, November 9, 1981

Codd received his B.A. and M.A. in Mathematics from Oxford University in England, and his M.Sc. and Ph.D. in Computer and Communication Sciences from the University of Michigan. He is a Member of the National Academy of Engineering (USA) and a Fellow of the British Computer Society.

The ACM Turing Award is presented each year in commemoration of A. M. Turing, the English mathematician

2. Motivation

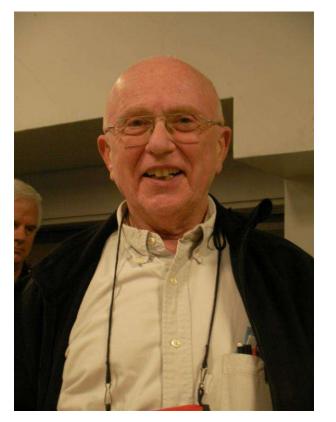
The most important motivation for the research work that resulted in the relational model was the objective of providing a sharp and clear boundary between the logical and physical aspects of database management (including database design, data retrieval, and data manipulation). We call this the *data independence objective*.

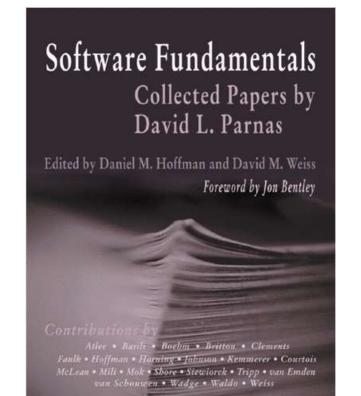
A second objective was to make the model structurally simple, so that all kinds of users and programmers could have a common understanding of the data, and could therefore communicate with one another about the database. We call this the *communicability objective*.

A third objective was to introduce high level language concepts (but not specific syntax) to enable users to express operations upon large chunks of information at a time. This entailed providing a foundation for setoriented processing (i.e., the ability to express in a single statement the processing of multiple sets of records at a time). We call this the *set-processing objective*.



David Parnas - Information Hiding - 1972

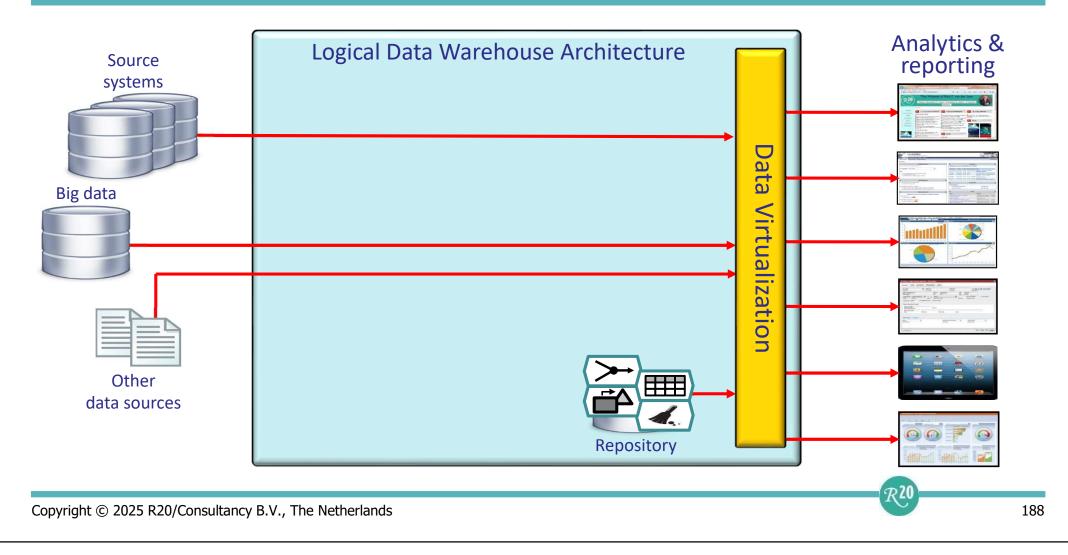




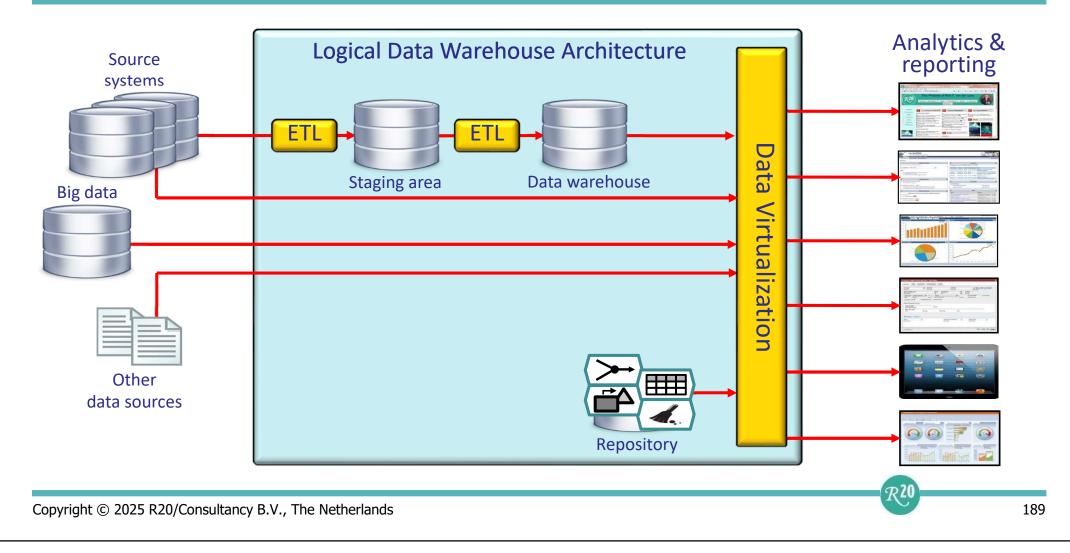
http://www.cs.umd.edu/class/spring2003/cmsc838p/Design/criteria.pdf



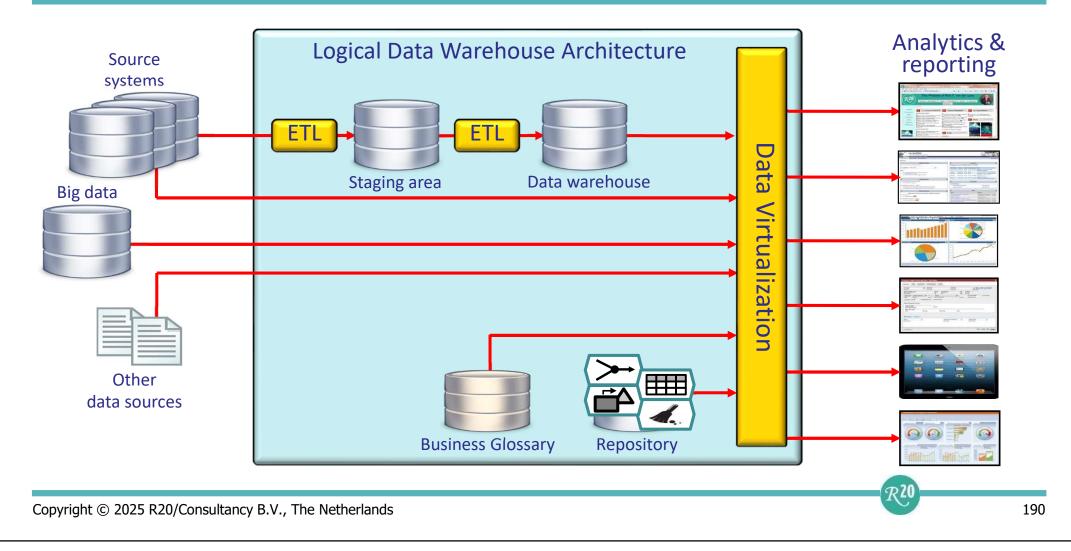
The Logical Data Warehouse Architecture (1)



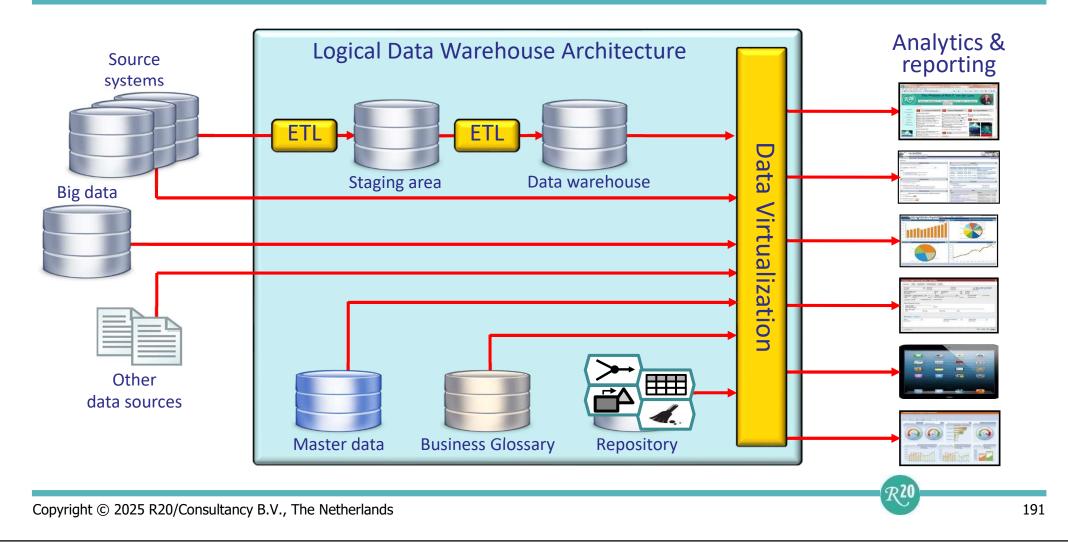
The Logical Data Warehouse Architecture (2)



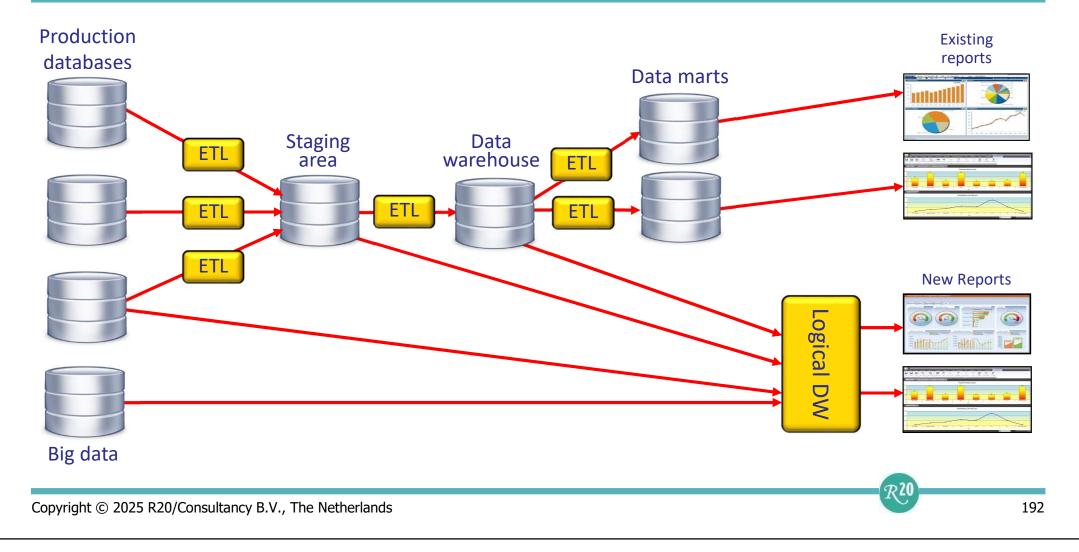
The Logical Data Warehouse Architecture (3)



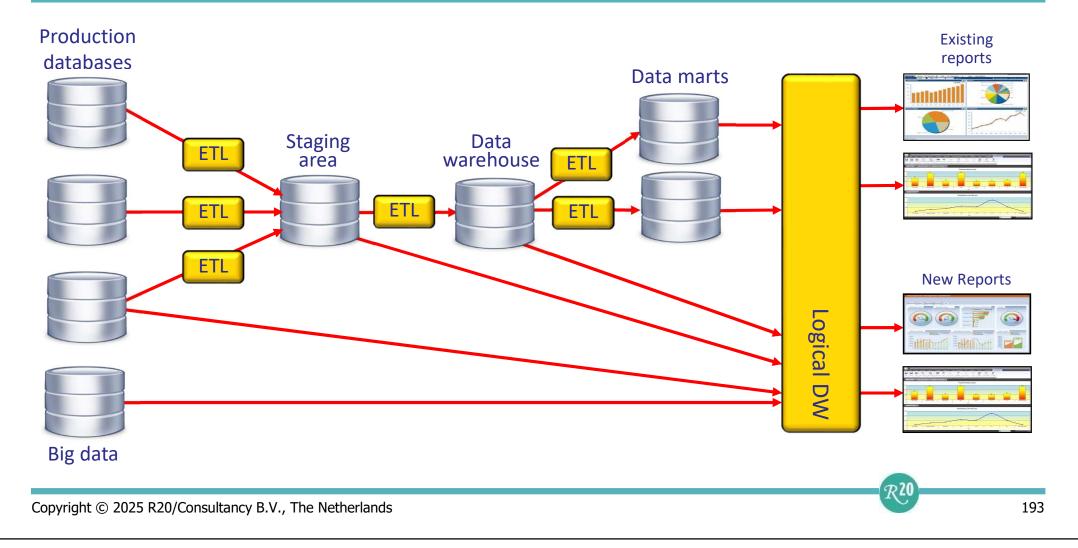
The Logical Data Warehouse Architecture (4)



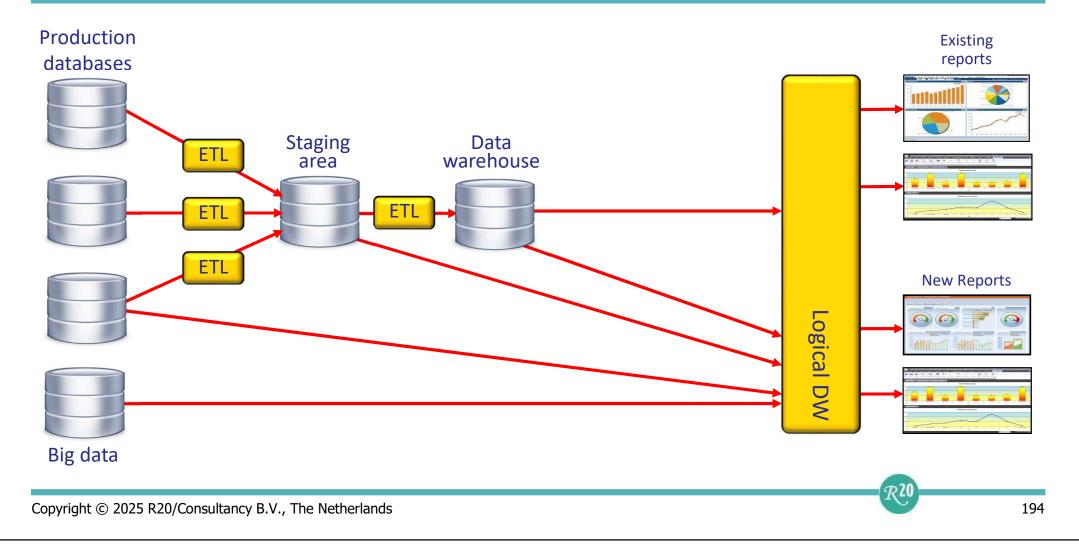
Wrap the Old Data Warehouse (1)

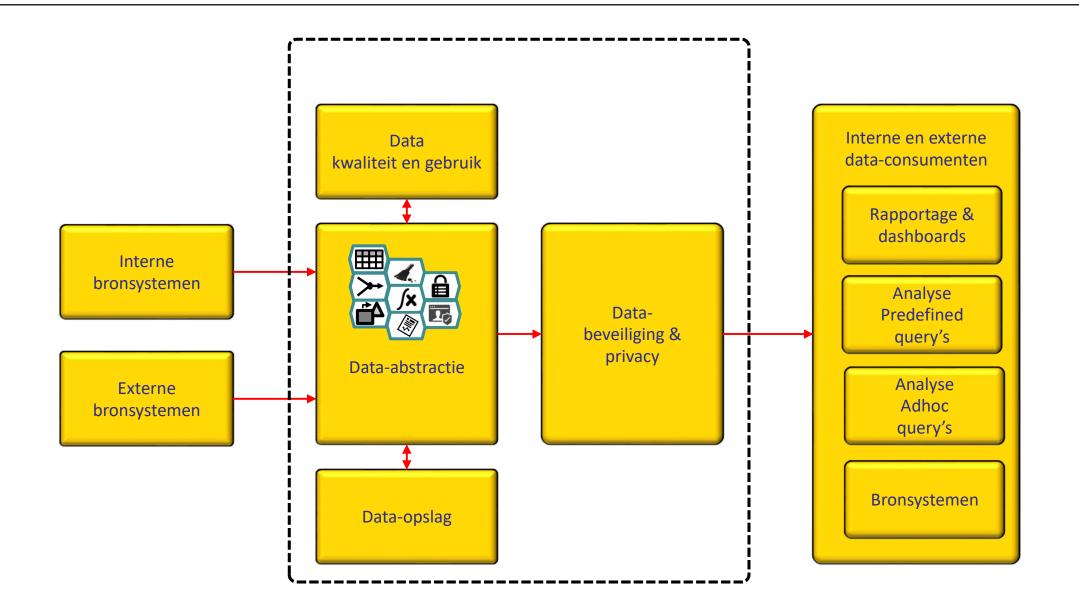


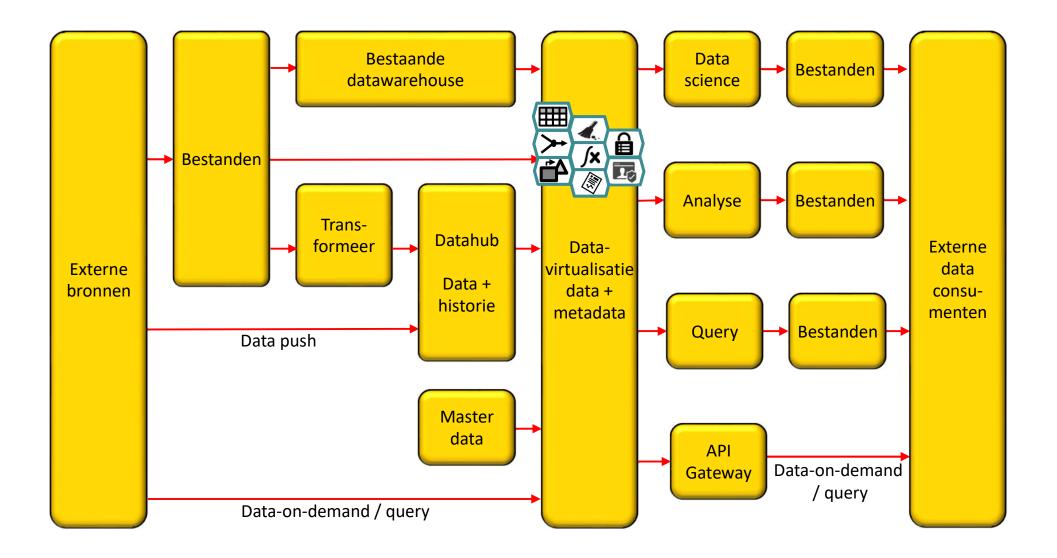
Wrap the Old Data Warehouse (2)

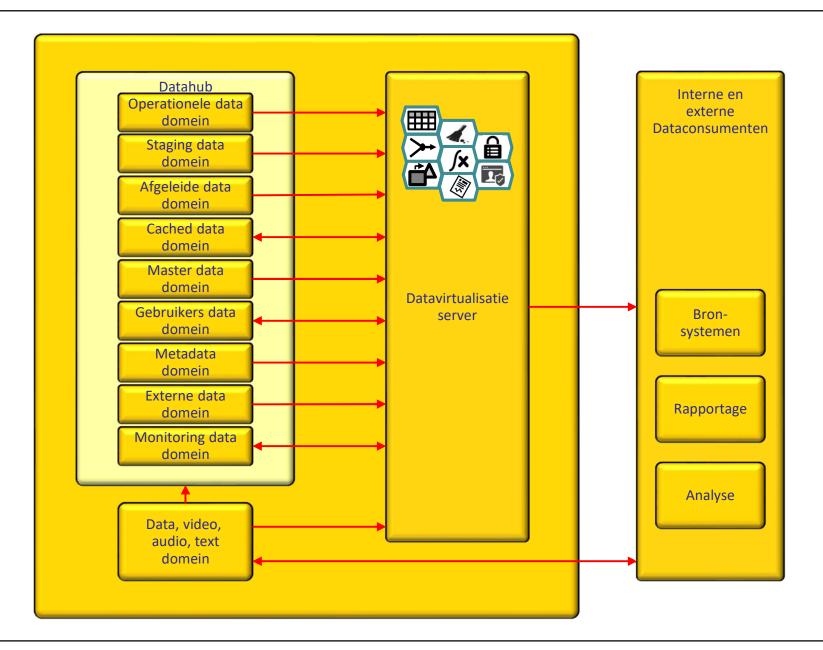


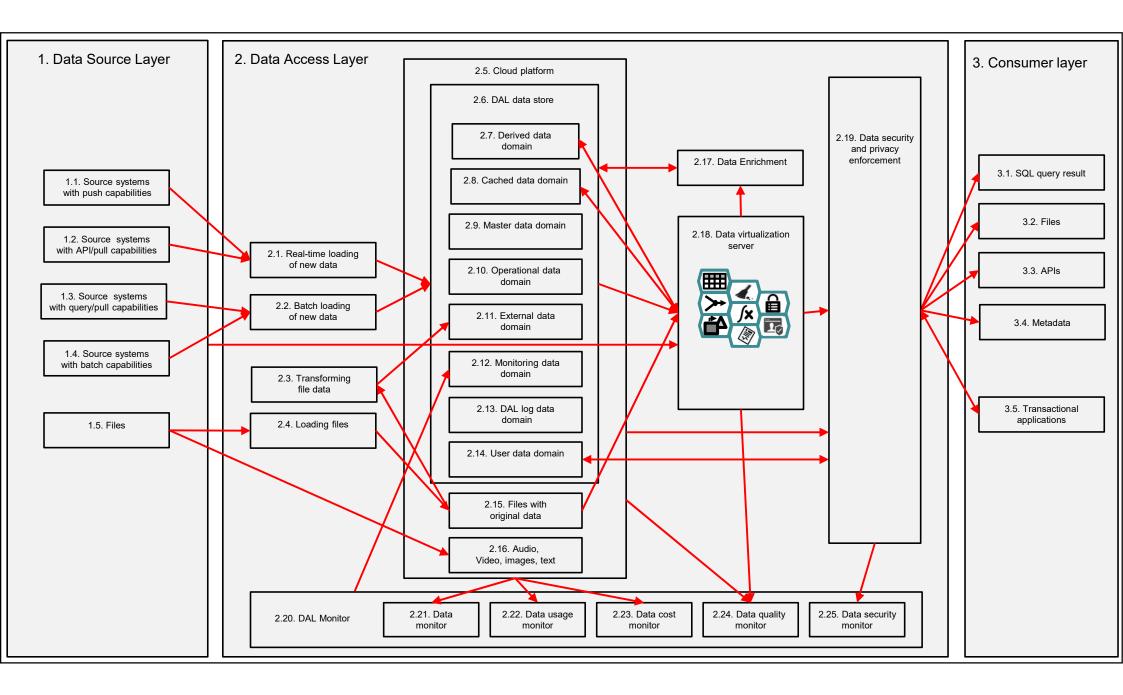
Wrap the Old Data Warehouse (3)







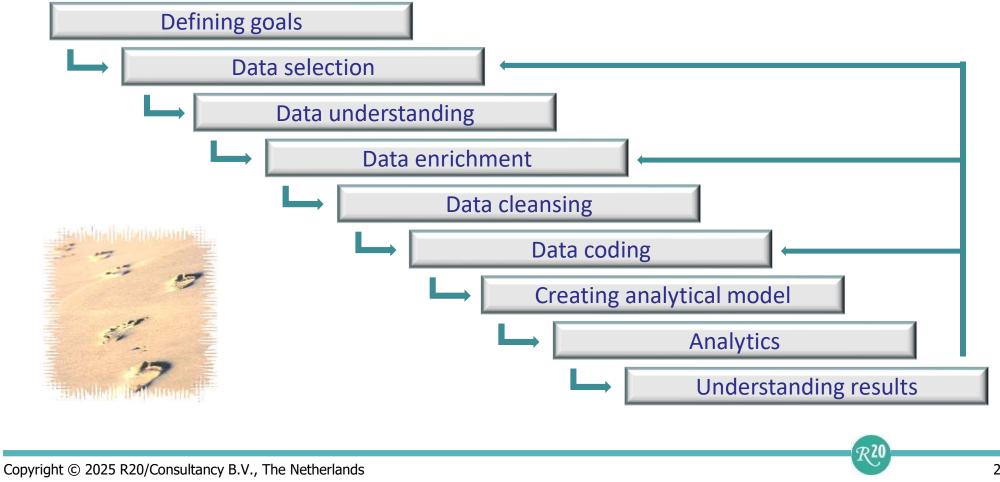




Part 6.3: The Data Lake



Data Science Steps



Data Coding



- Computation
 - examples: divide all monthly salaries by 1000; round all prices
- Grouping continuous values
 - example: all transaction between 08:00 and 10:30 will belong to group 1, all transactions between 10:31 and 12:00 will belong to group 2
 - do groups need equal sizes (with respect to ranges)?
 - do groups need equal numbers of values?
- Scaling
 - most neural networks accept numeric data only in the range 0.0 to 1.0 or -1.0 to 1.0; used for continuous values, such as salary and weight
- Normalizing
 - sum all elements, and divide each element by the sum
 - value represents the percentage of contribution
- Symbolic to numeric transformations
 - example: the string "yes" becomes 1, and "no" becomes 0
- Coding discrete values
 - transform a column with fixed set of values (F) into F columns with yes/no values

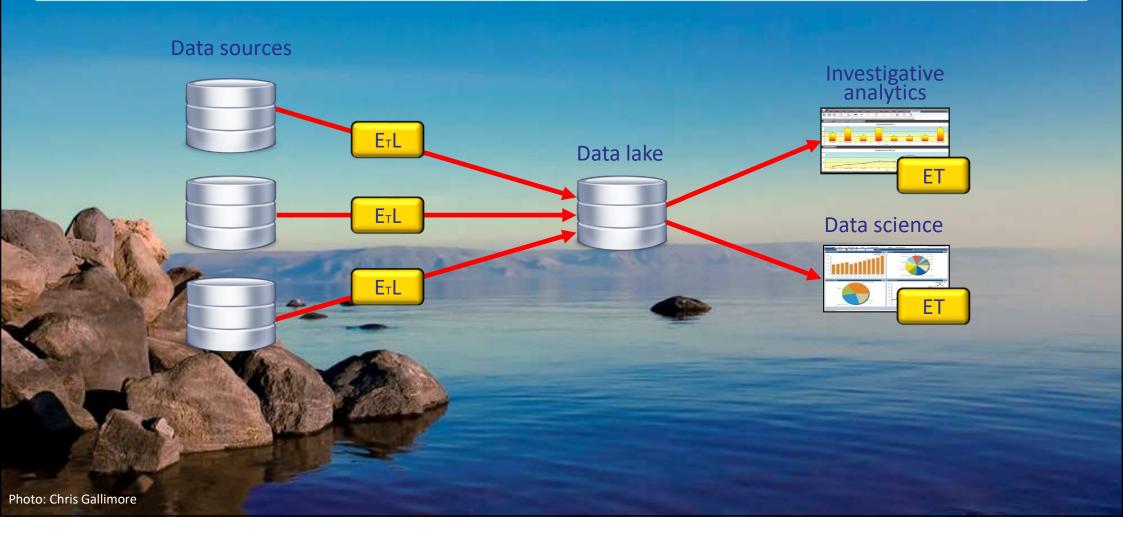
Common Definition of Data Lake

James Serra:

A "data lake" is a storage repository, usually in Hadoop, that holds a vast amount of raw data in its native format until it is needed. It's a great place for investigating, exploring, experimenting, and refining data, in addition to archiving data.

"

The Data Lake

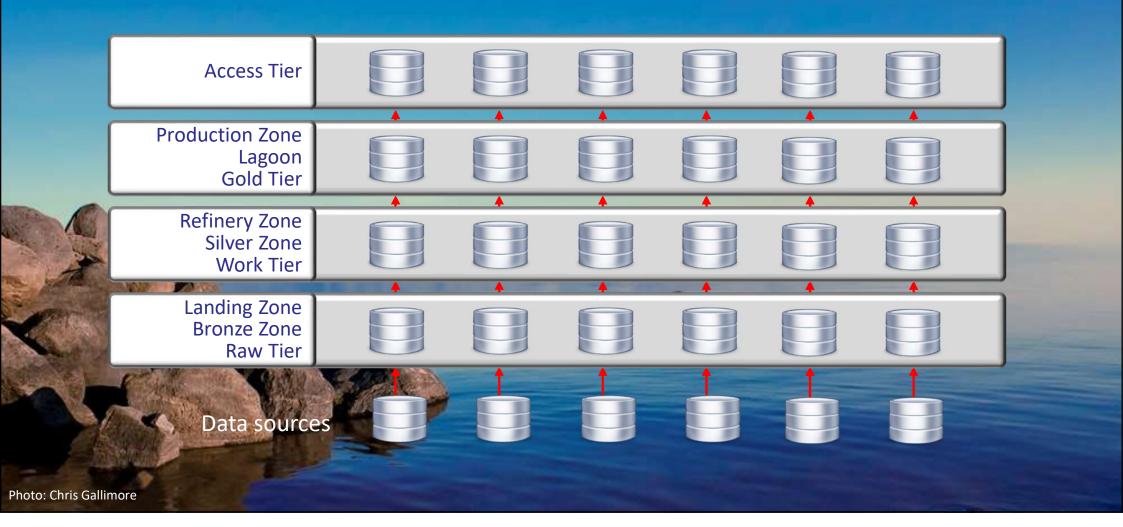


Challenges of a Physical Data Lake

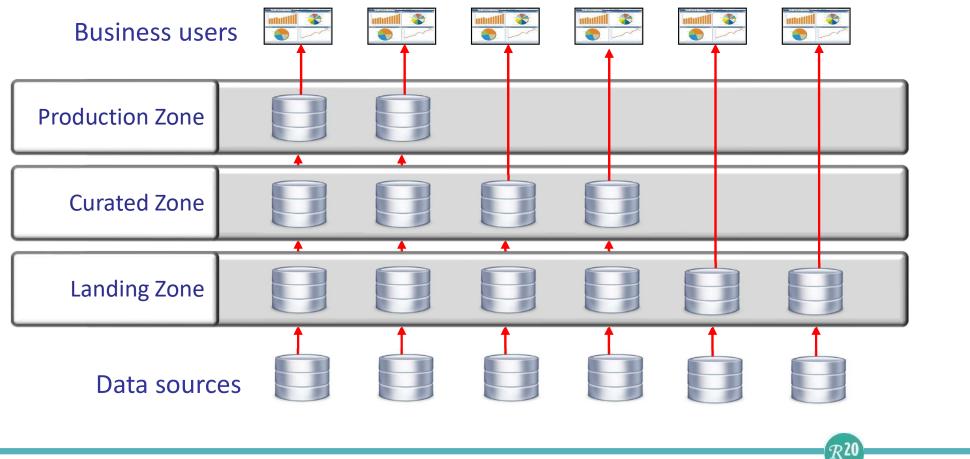


- Big data too big to move
 - Too slow to copy and bandwidth issues
- Complex "T" moved to data consumption
- Company politics
- Data privacy and protection regulations
- Data in data lake is stored outside original security realm
- Metadata to describe data
- Some sources are hard to copy
 - For example, mainframe data
- Refreshing of data lake
- Management of data lake required

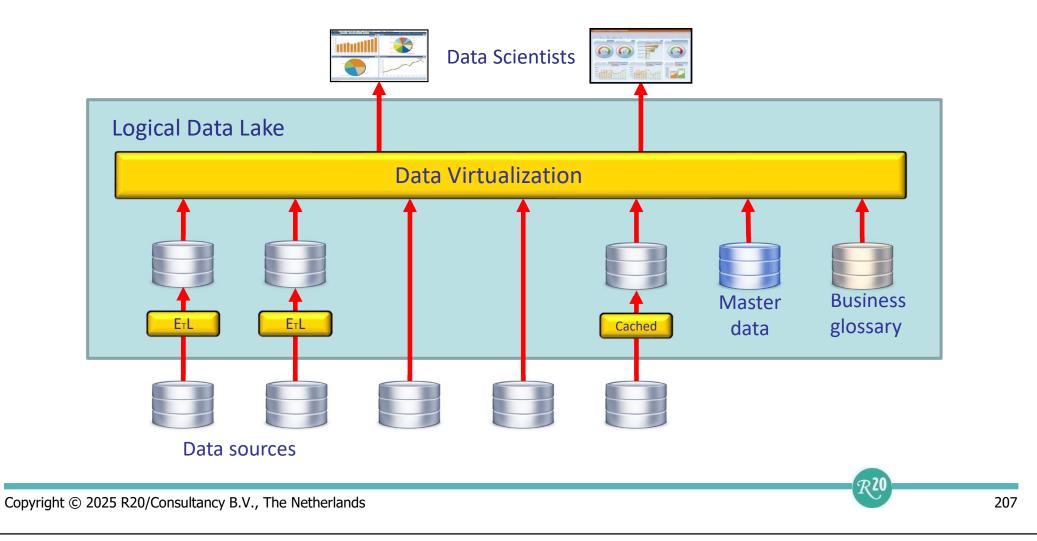
The Business Data Lake



A Data Lake With Multiple Zones



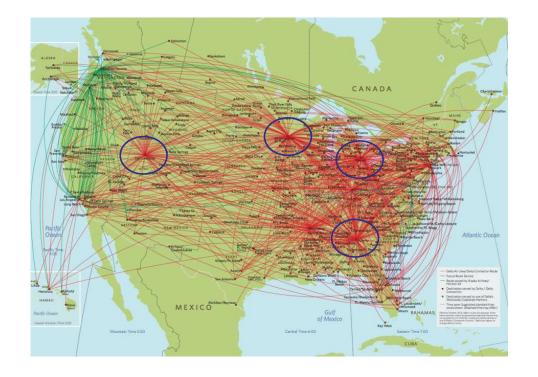
The Logical (Virtual) Data Lake



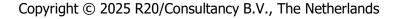
Part 6.4: The Data Hub



What is a Data Hub?



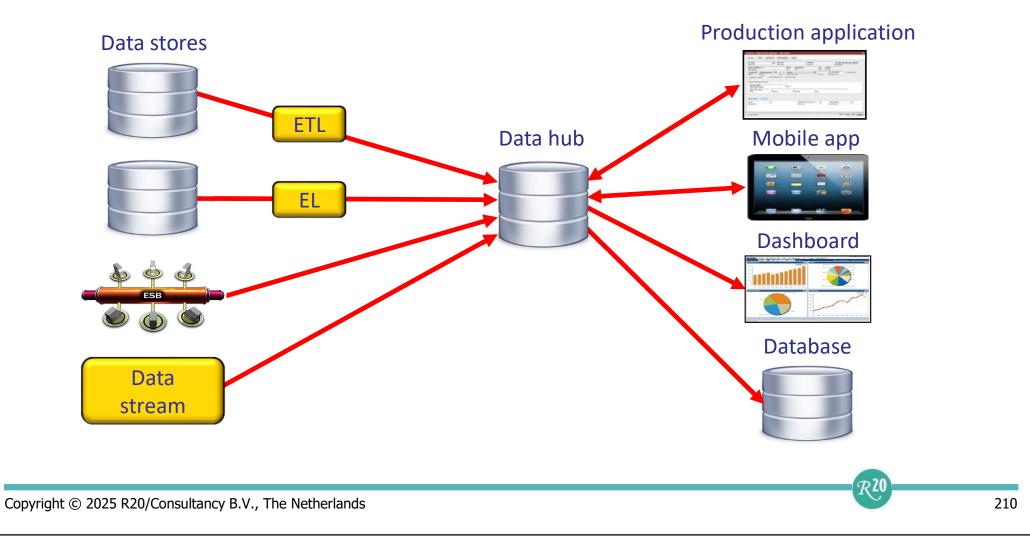
Wikepedia: A data hub is a collection of data from multiple sources organized for distribution, sharing, and often subsetting and sharing. Generally this data distribution is in the form of a hub and spoke architecture



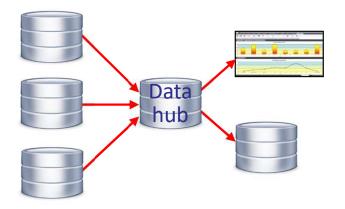


209

The Data Hub (Sharing of Data)



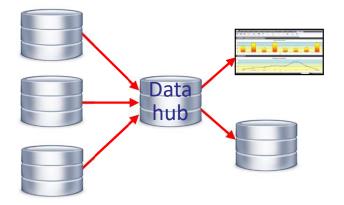
Characteristics of the Data Hub



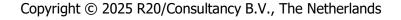
- The main goal of a data hub is to organize data efficiently, storing it in a cost-efficient manner and expose it towards key business functions
- It excels in easy integration, and enables deduplication, security, quality and data standardization
- The data hub can be leveraged to enable data processing activities with the end use-case in mind, and typically has governance capabilities
- Although operationally focused, it can be trusted as an analytical data source



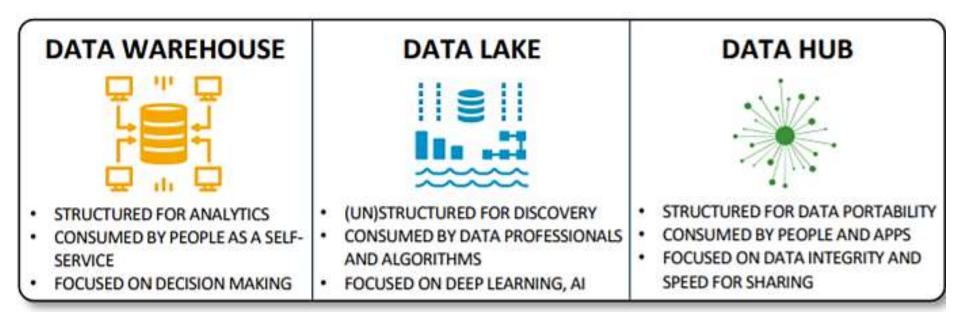
Data Hub Versus the Rest Of the World



- Data hub versus *data warehouse*: a data hub is generally non-integrated and often at different grains
- Data hub versus operational data store: a data hub does not need to be limited to operational data
- Data hub versus *data lake*: a data lake tends to store data in one place for availability, and allow/require the consumer to process or add value to the data
- Data warehouses and data lakes may be endpoints, data hubs are not endpoints, they serve as points of intermediation and data exchange



Comparison of Three Data Storage Environments



Source: Talend; see https://www.talend.com/resources/customer-360-data-hub/



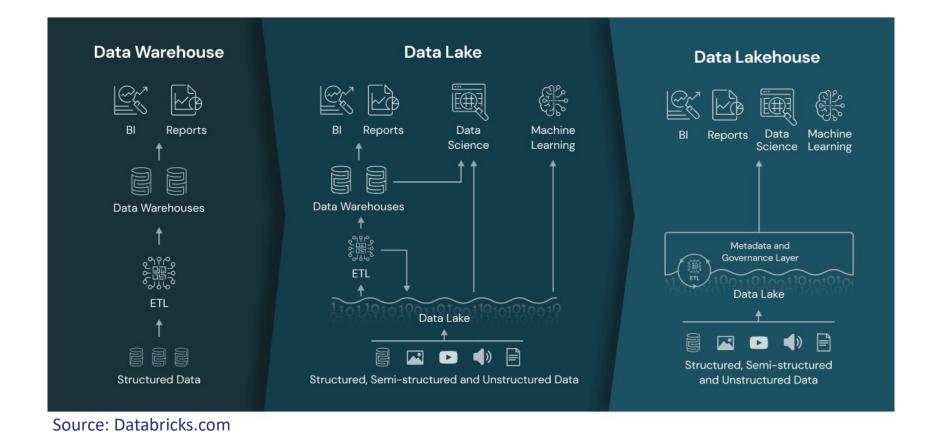
Part 6.5: The Data Lakehouse



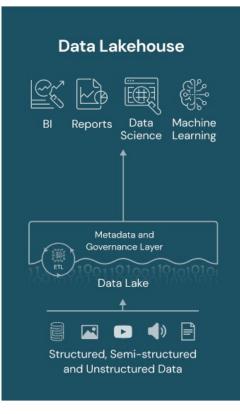
Definitions of Data Lakehouse

- DataBricks: "A data lakehouse is a [...] open data management architecture that combines the flexibility, cost-efficiency, and scale of data lakes with the data management and ACID transactions of data warehouses, enabling business intelligence (BI) and machine learning (ML) on all data."
- Striim, John Kutay: "A data lakehouse is a new, big-data storage architecture that combines the best features of both data warehouses and data lakes. A data lakehouse enables a single repository for all your data (structured, semi-structured, and unstructured) while enabling best-in-class machine learning, business intelligence, and streaming capabilities."
- Dremio, Deepa Sankar: "A [data] lakehouse has the performance and optimization of a data warehouse combined with the flexibility of a data lake."
- Wikipedia: "Databricks develops and sells a cloud data platform using the marketing term lakehouse, a portmanteau based on the terms data warehouse and data lake."

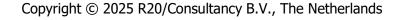
A Comparison of Three Data Architectures



Key Characteristics of a Data Lakehouse

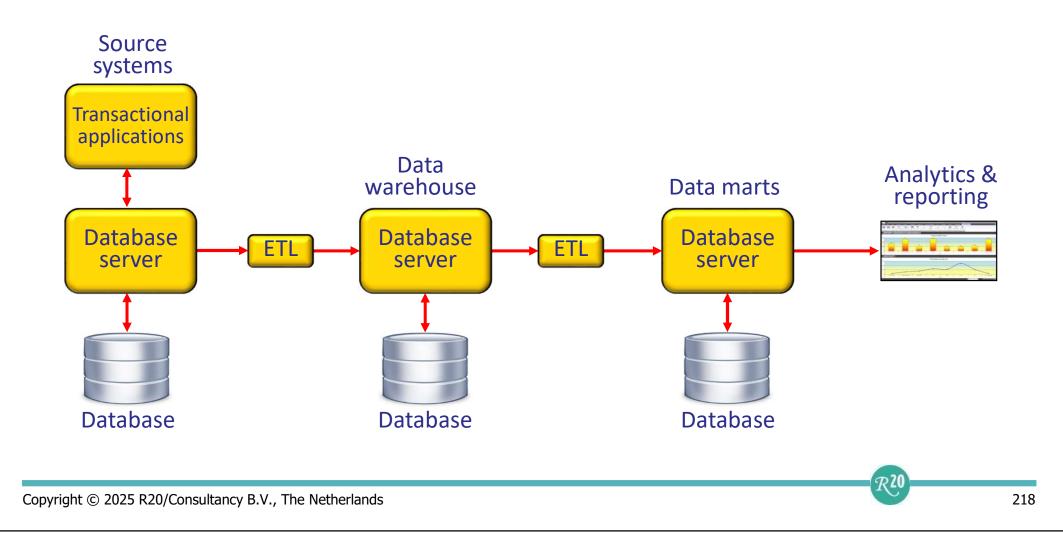


- Two use cases: BI and data science
- Data is stored once
- Supports structured and unstructured data
- Schema enforcement
- Open file formats
- Low-cost data storage
- ACID compliant

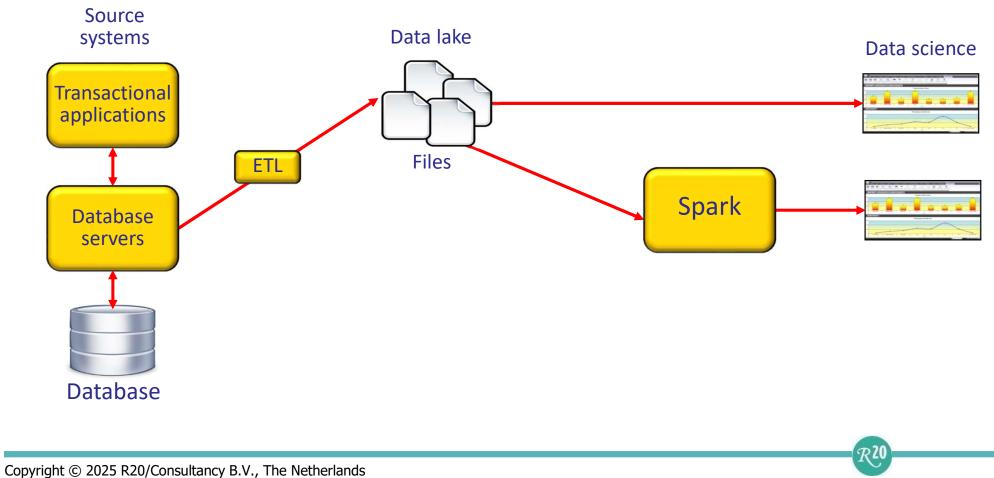




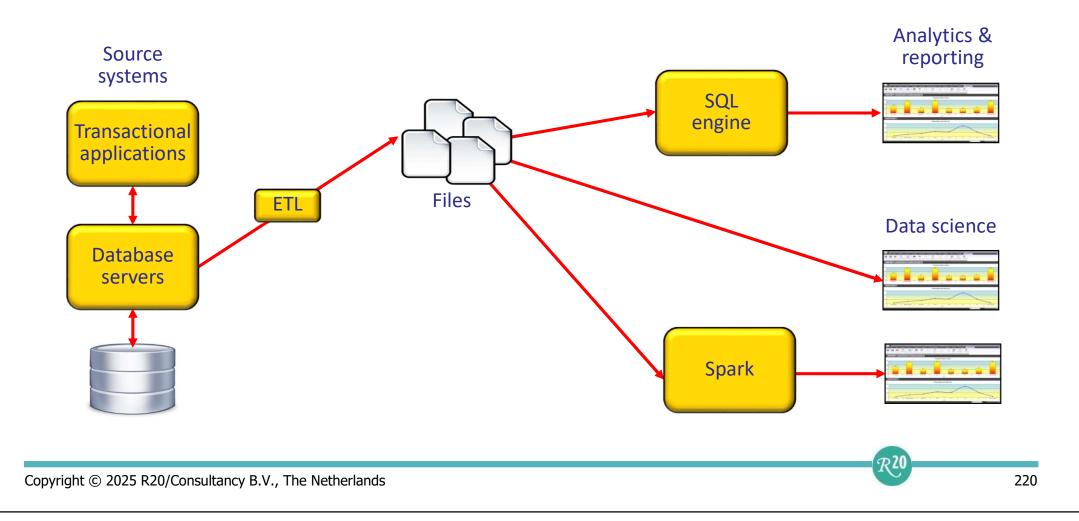
The Data Warehouse Architecture in More Detail



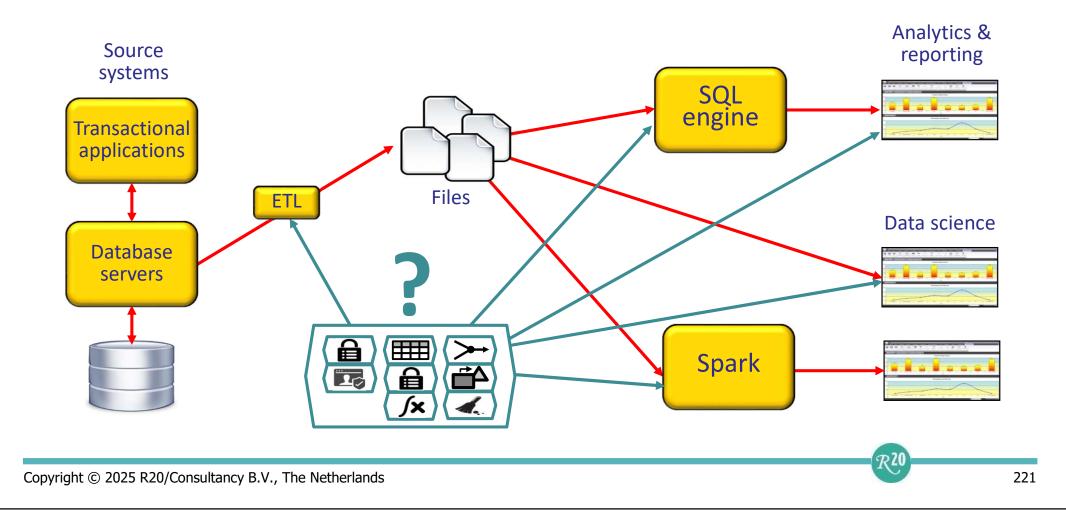
The Data Lake Architecture in More Detail



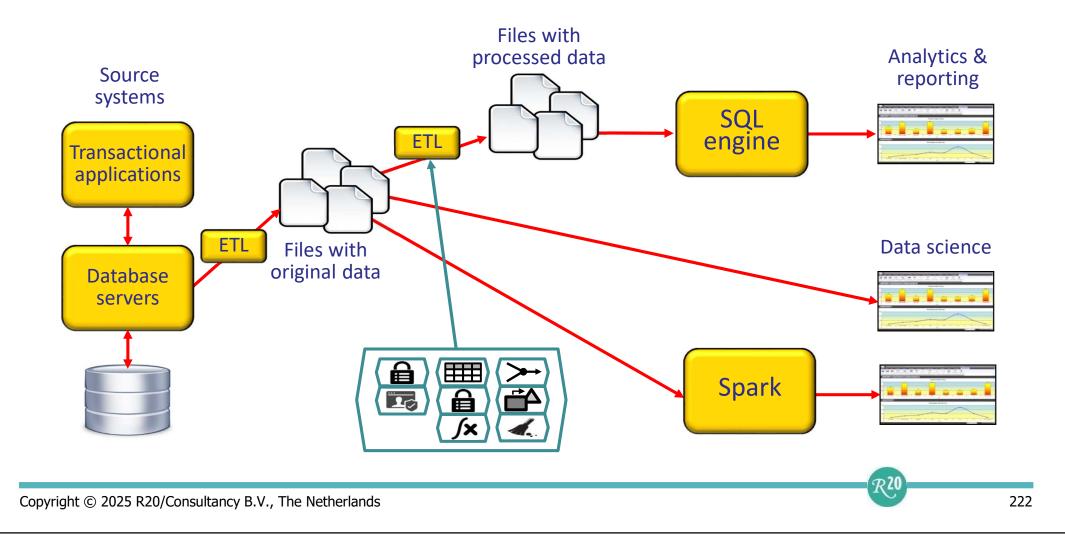
The Data Lakehouse Architecture in More Detail



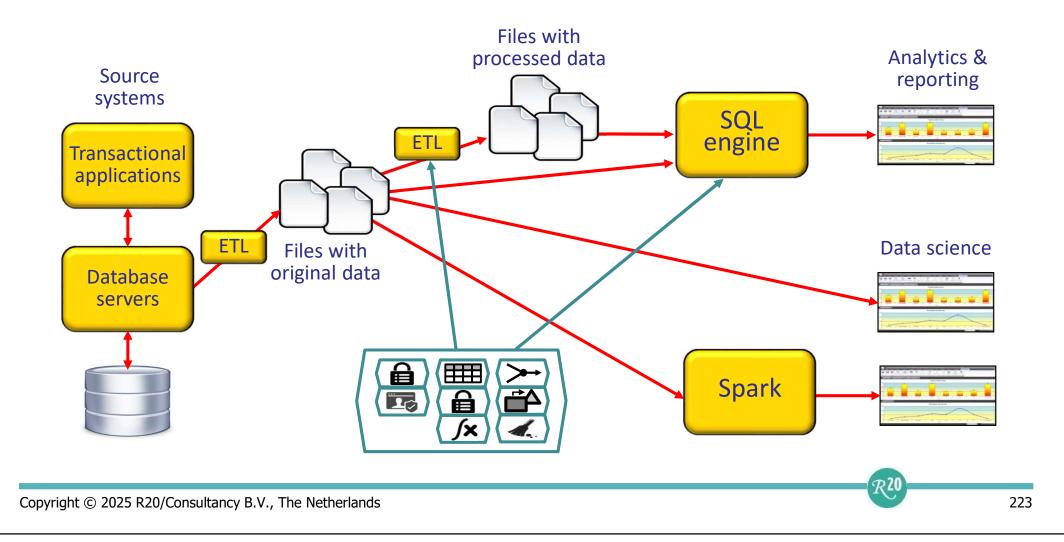
Where to Implement Data Processing Specifications?



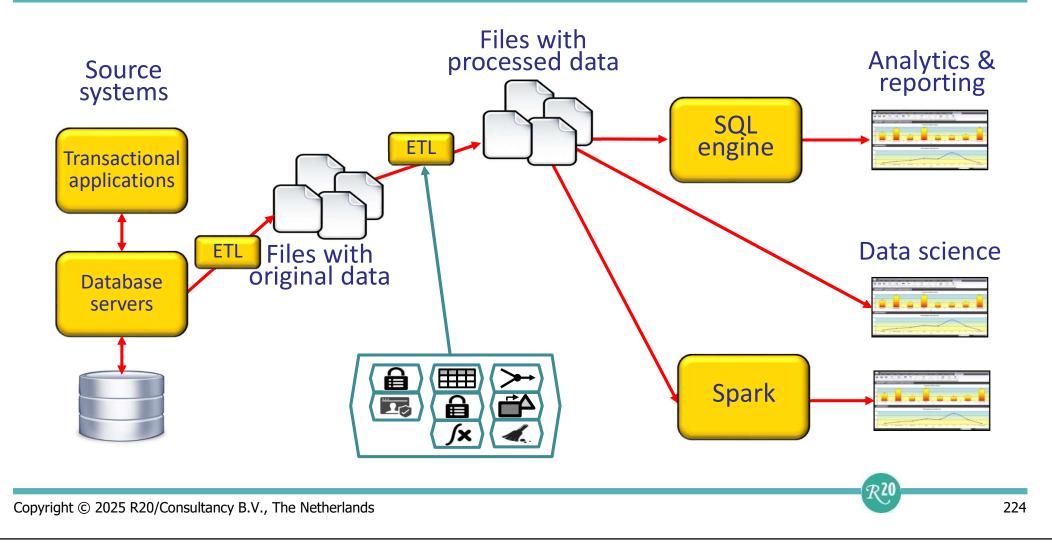
Solution 1: *All* **BI via Processed Data**



Solution 2: Some BI via Processed Data



Solution 3: Both Use Cases via Processed Data



High-Level Comparison of Three Architectures

Characteristic	Data Warehouse	Data lake	Data Lakehouse		
Type of data	Structured	Structured and unstructured	Structured and unstructured		
Use cases	BI, reporting, dashboarding	Experimental, investigative	Both		
Schema enforcement	Yes	Optional	Optional		
Open file format	No	Yes	Yes		
Low-cost data storage	No	Yes	Yes		
ACID-compliant	Yes	No	Yes		
Near real-time data	No	Yes	Yes Based on		
Non-siloed	No	No	Yes assumptions		
Data copies minimal	No	No	Yes		
Anonymization	Yes	Depends	Depends		
Auditable	Yes	Depends	Depends		
Performance optimized for BI	Yes	No	?		
Performance optimized for data science	No	Yes	Yes		
	Source: Various articles in the Internet				

Source: Various articles in the Internet

Copyright $\ensuremath{\textcircled{C}}$ 2025 R20/Consultancy B.V., The Netherlands

 \mathcal{R}^{20}

Part 6.6: The Data Fabric



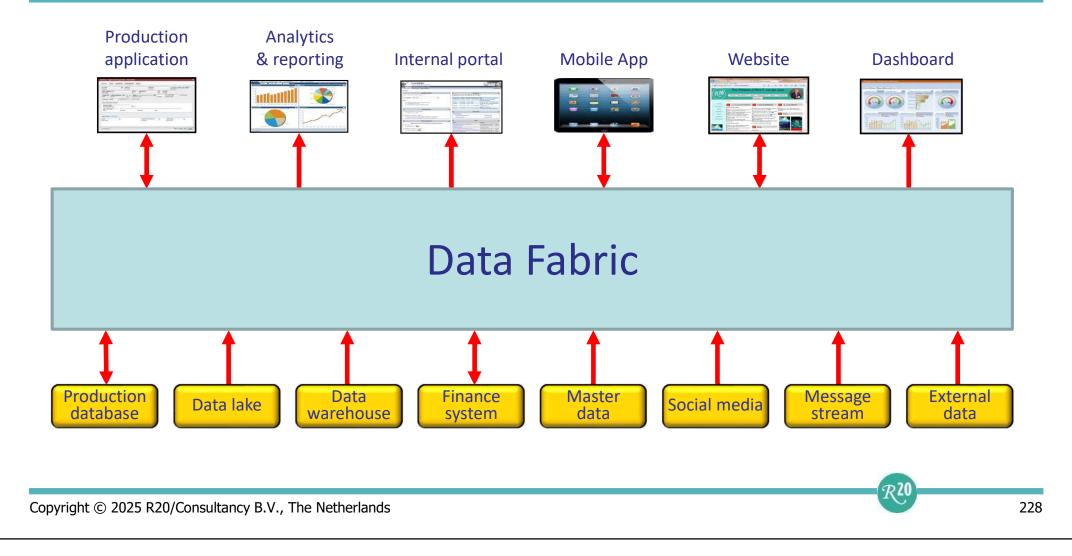
What is the Data Fabric?



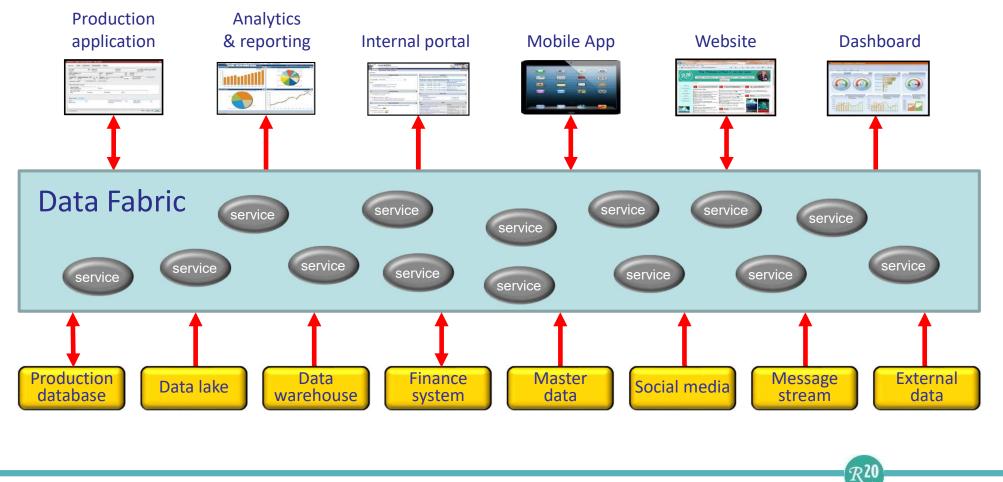
- Gartner: A data fabric is generally a custommade design that provides reusable data services, pipelines, semantic tiers or APIs via combination of data integration approaches in an orchestrated fashion.
- Gartner: Data fabric enables *frictionless access* and sharing of data in a distributed data environment. It enables a *single* and *consistent* data management framework, which allows seamless data access and processing by design across otherwise siloed storage.



The Data Fabric

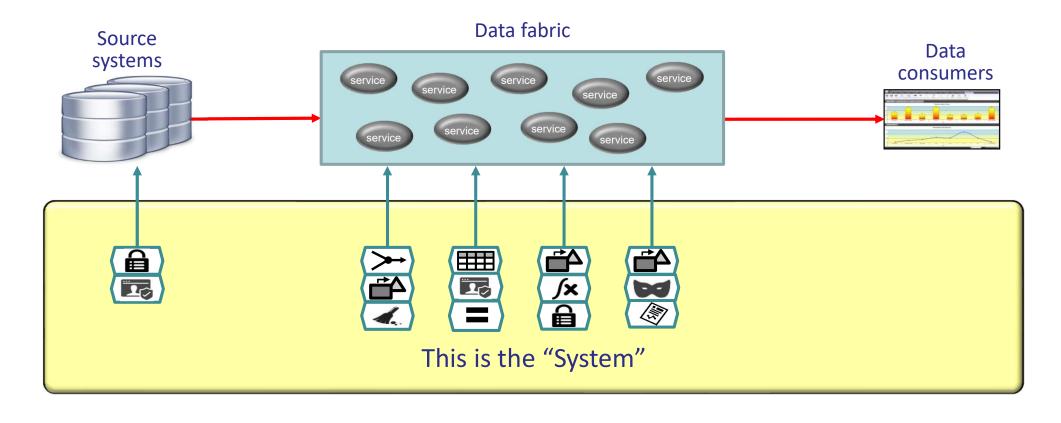


The Services of a Data Fabric



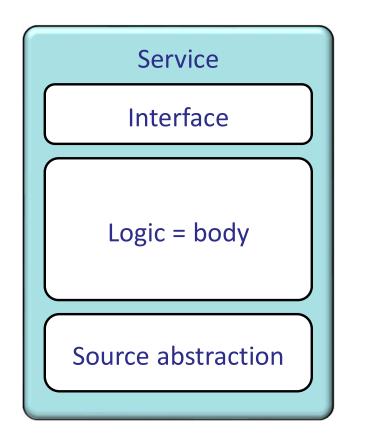
Copyright © 2025 R20/Consultancy B.V., The Netherlands

Data Fabrics and Data Processing Specifications

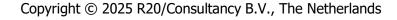




The Components of Services

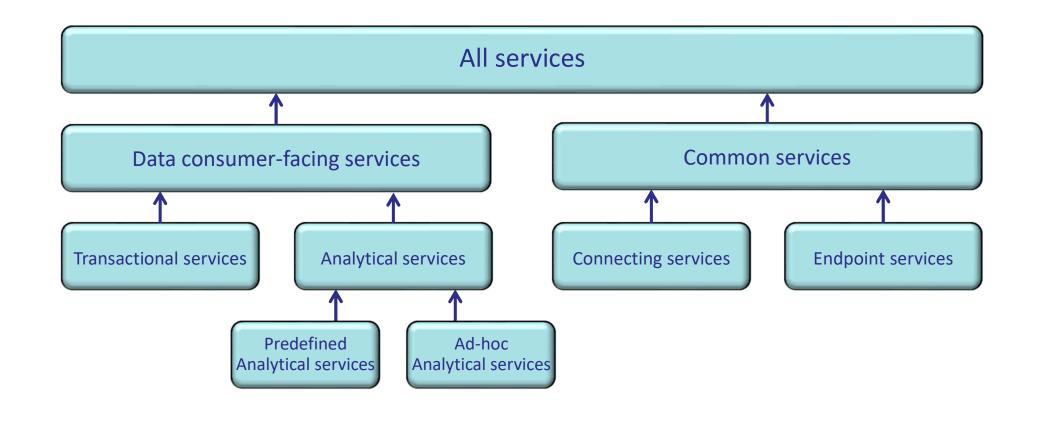


- The interface component is responsible for handling incoming parameters and outgoing results
- The logic of the service form the body
- The body deals with data processing specifications
- Abstraction layer to make it independent of changes to the IT systems





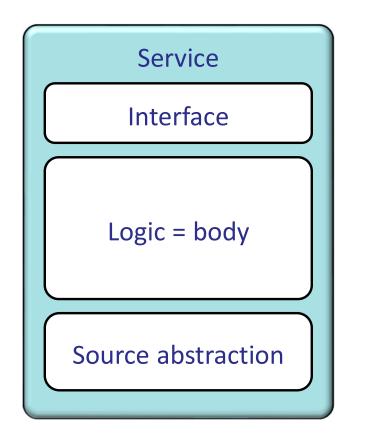
A Data Fabric Consists of Different Types of Services



Copyright © 2025 R20/Consultancy B.V., The Netherlands

 \mathcal{R}^{20}

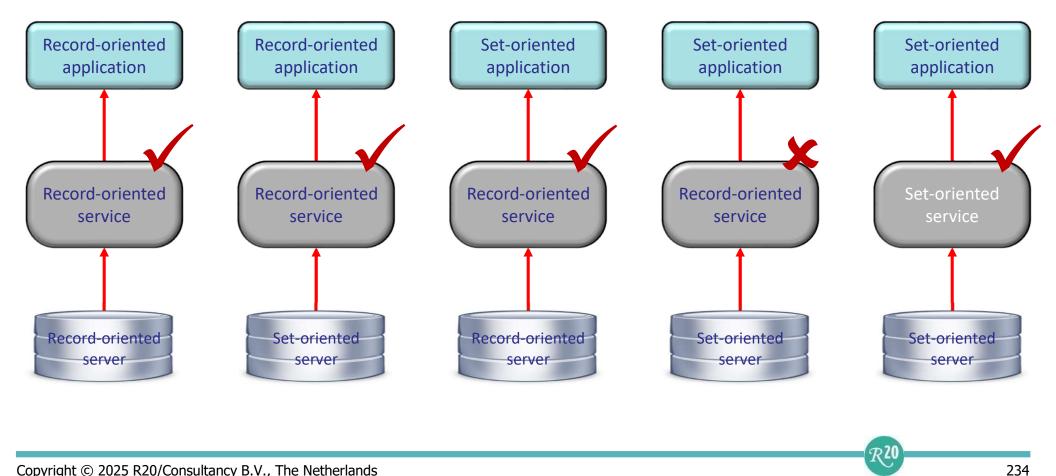
12 Capabilities for Frictionless Data Access



- Data preparation, such as transformations, calculations, aggregations, filters, joins, ...
- Adaptable logic
- High performance
- Data access by many data consumption forms
- Access to all the data sources
- Processing of all types of data
- Data security and privacy
- Real-time data access
- Read and write data access
- Data minimization
- Event processing
- Technical and business metadata management
- Master and reference data management

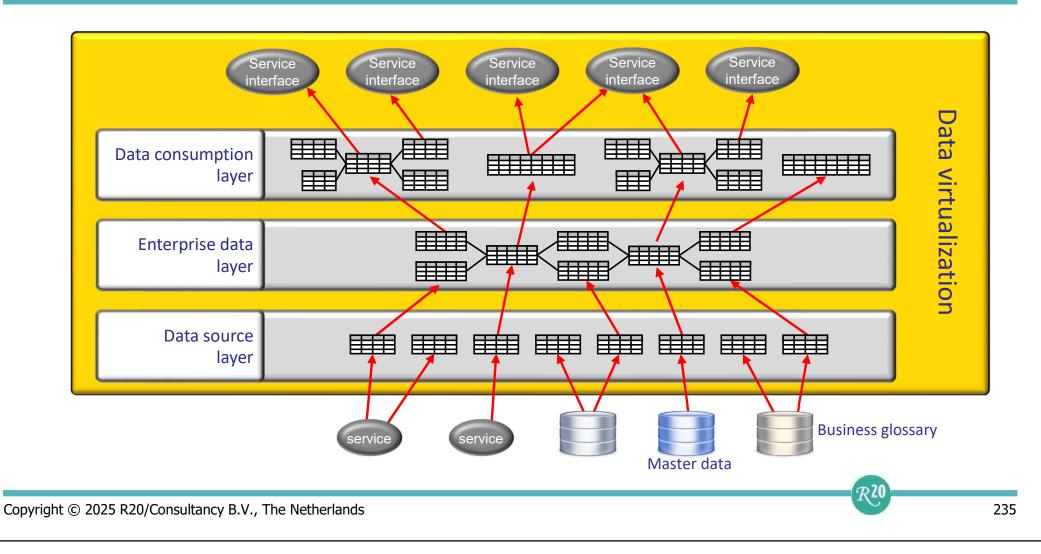


Record-Oriented or Set-Oriented Interface?

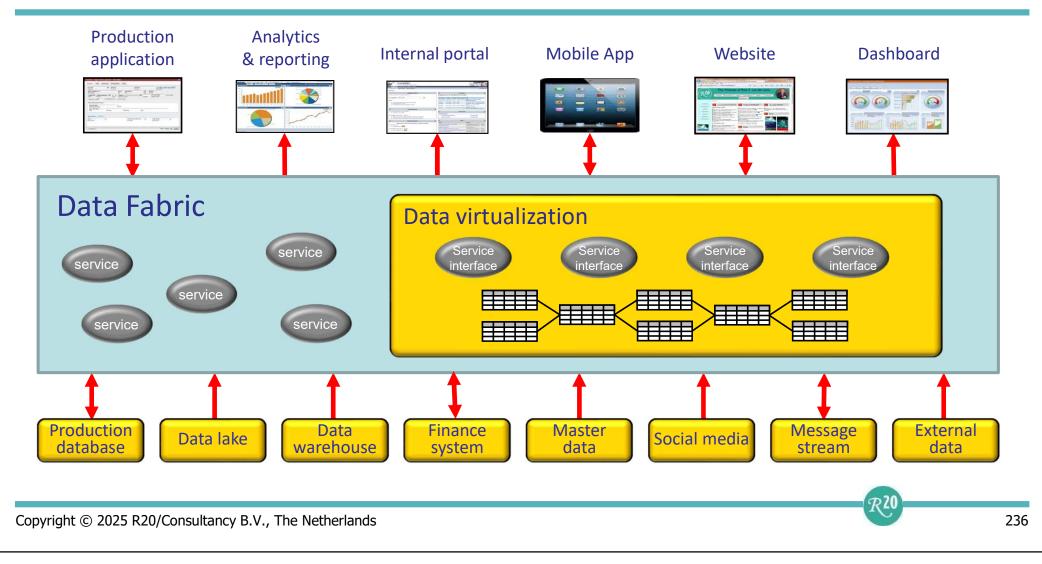


Copyright © 2025 R20/Consultancy B.V., The Netherlands

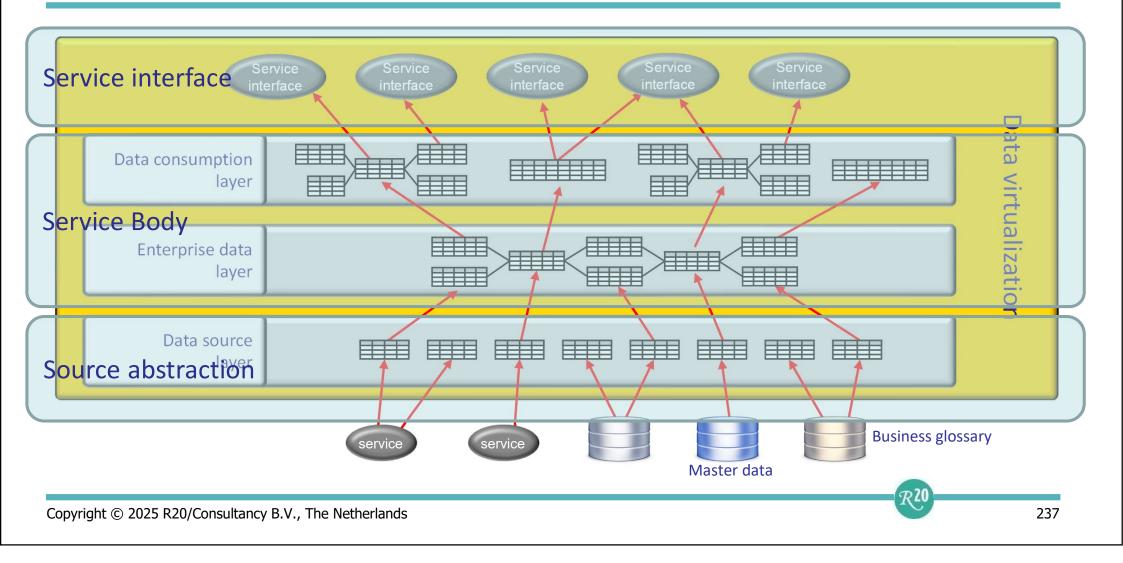
Developing Data Fabric Services with Data Virtualization



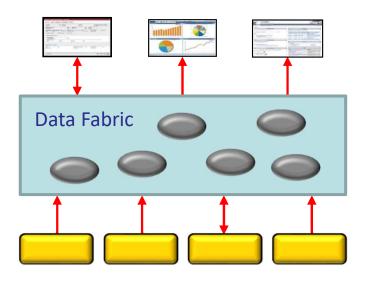
The Data Fabric



Developing Data Fabric Services with Data Virtualization

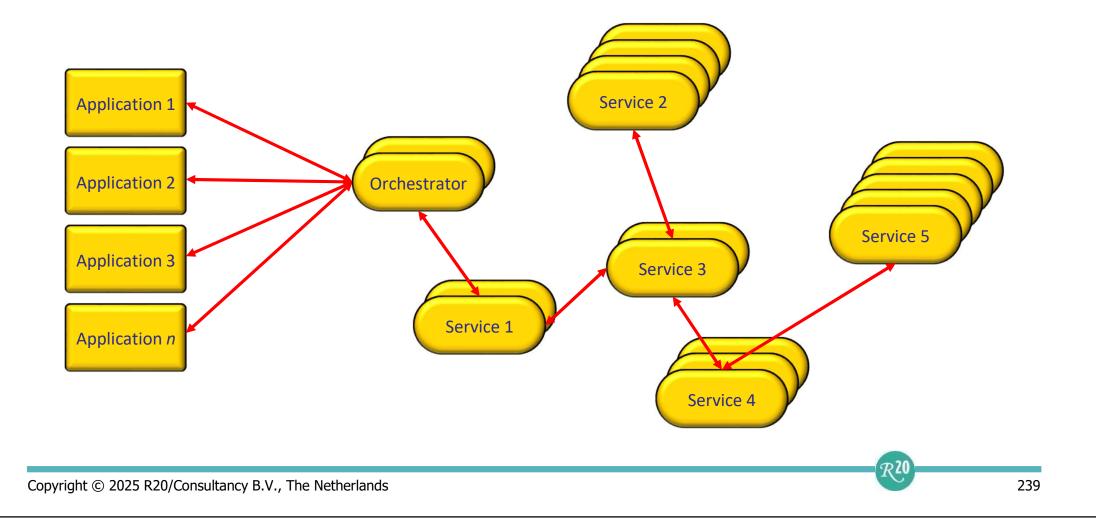


Remarks on Data Fabric

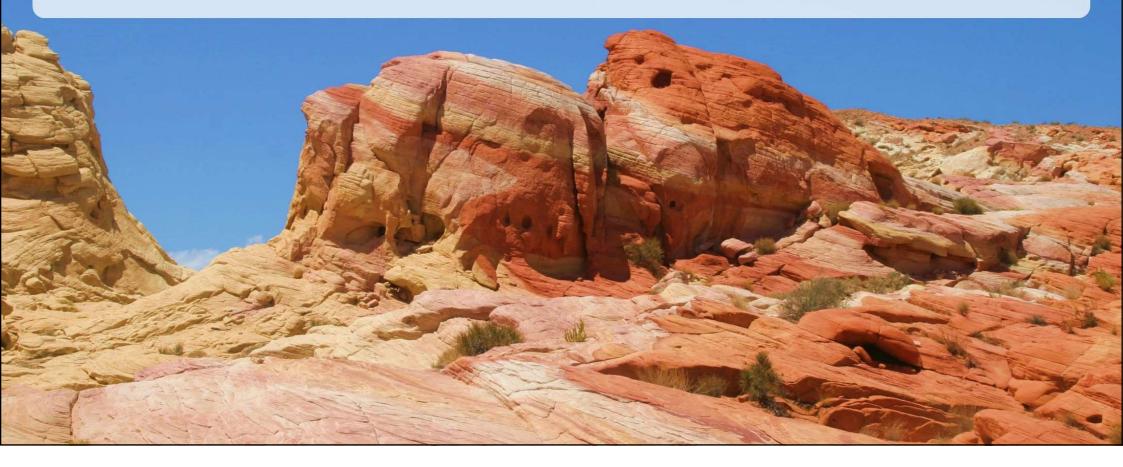


- Poorly defined concept
- Metadata is required
- Master data needed to integrate data correctly
- A data fabric may contain a data warehouse, data lake, or data hub
- Dedicated tool market is small, e.g. Cinchy
- Data fabric must support transactional and analytical workloads

Data Fabric \neq **Micro-Services Architecture**



Part 6.7: The Data Mesh



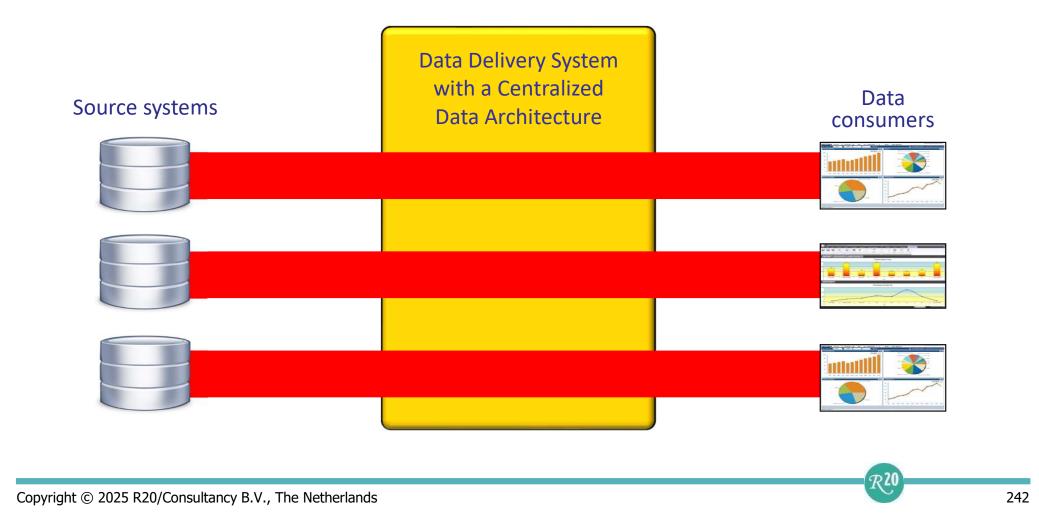
What is a Data Mesh?



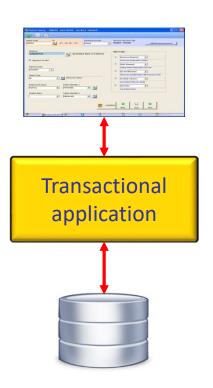
- Introduced by Zhamak Dehghani:
- Data platforms based on the data lake architecture have common failure modes that lead to *unfulfilled promises* at scale.
- To address these failure modes we need to shift from the centralized paradigm of a lake, or its predecessor data warehouse.
- We need to shift to a paradigm that draws from *modern* distributed architecture: considering domains as the first class concern, applying platform thinking to create selfserve data infrastructure, and treating data as a product."



Single-Domain Data Consumers

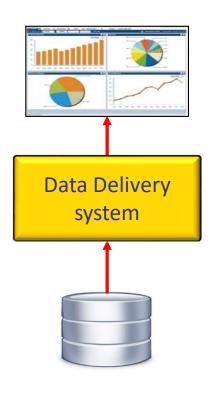


Data Engineers for Transactional Applications



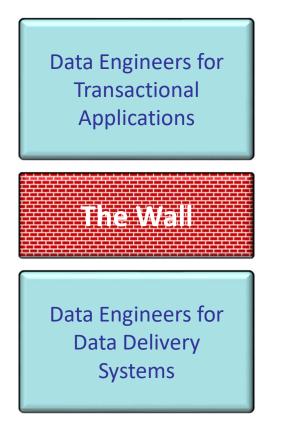
- They are single-domain experts
- They focus only on data requirements of transactional applications
 - Not on other forms of data consumption, such as BI
- They do not make the data easily consumable
- They implement all the business rules
- They know about all the exceptions
- They know when changes are implemented

Data Engineers for Data Delivery Systems



- They need to understand the data of all the domains
 - Hyper-domain experts
- They need to transform the data into consumable data
- They need to work with data not designed to be integrated
- They need to understand all the business rules that need to be applied
 - Complex ETL processes
- They need to understand the data requirements of the data consumers
- They need to deal with SLAs of the data consumers

The Wall between Two Groups of Data Engineers

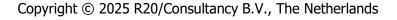


- Different groups of data engineers
- Different tool sets
- Different responsibilities
- Changes of data or data structures not always communicated
- Who owns the data and who is responsible for data quality?
- How to implement "the right to be forgotten/corrected"?

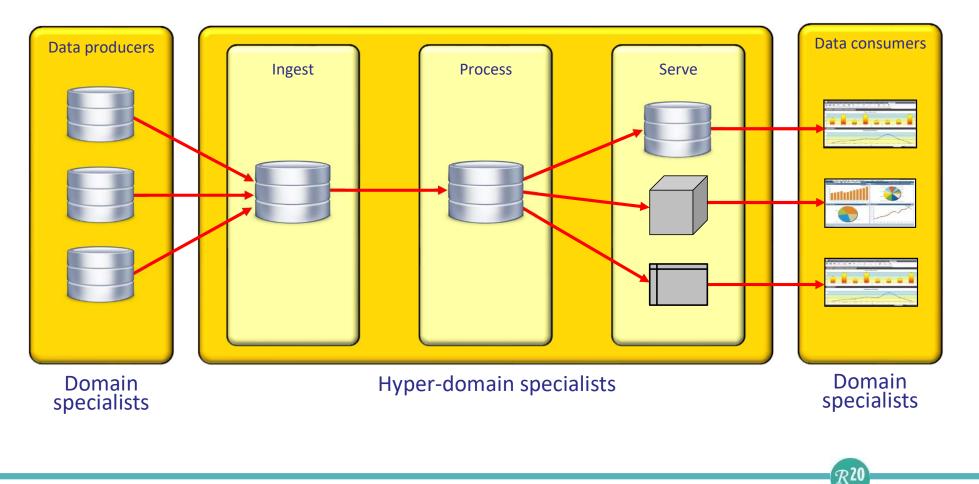
Interfaces of Transactional Applications



- Most are not developed to offer an interface
- Interfaces that do exist, are commonly developed for record-oriented access
- Direct database access complex or not always allowed
- Rarely support for historic data
- Risky to bypass multi-tenant systems

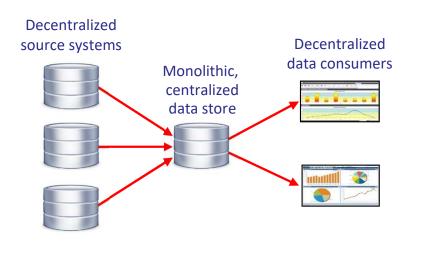


Current Centralized, Monolithic Data Architectures



Copyright \odot 2025 R20/Consultancy B.V., The Netherlands

Traditional Centralized, Monolithic Architectures

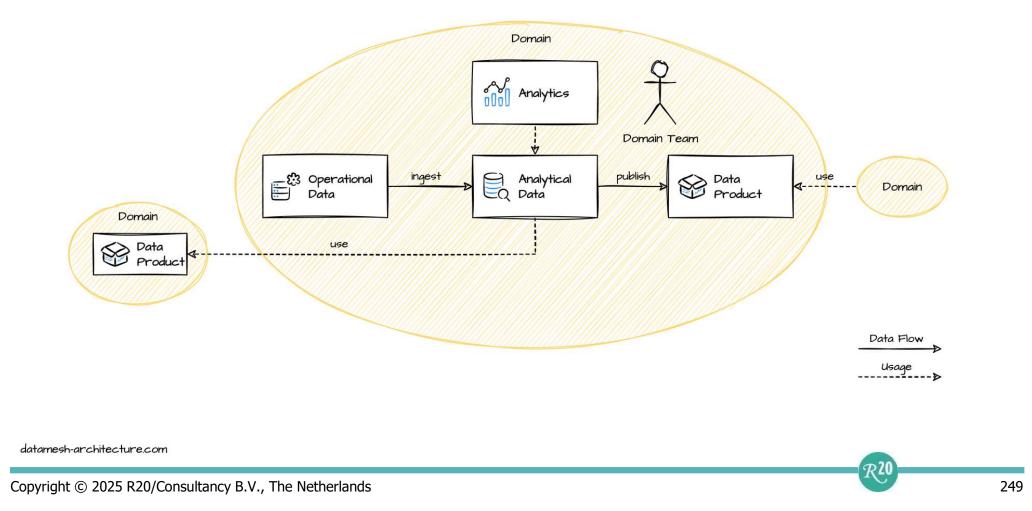


- Large and complex monolithic solutions
- Single-domain versus multi-domain experts
 - Application developers = single-domain experts
 - Data engineers = multi-domain experts
 - Data consumers = single-domain experts
- Data engineers need to understand all the business logic
- Who owns the data in the central database?
- Storage of redundant data



Copyright © 2025 R20/Consultancy B.V., The Netherlands

A Domain-Oriented Architecture

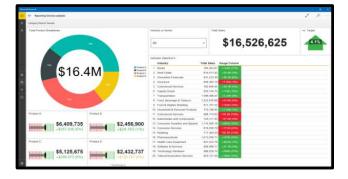


Potential Data Products

Data as file

Set As			Filter		Merge		
SampleID	Common_Name	Description	KeyHand	DigitHand	Hand	Individual	14
	keyboard	Akey	Left	NA	Left	M9	
M98key217,141032	keyboard	Bkey	Ambiguous	NA	Ambiguous	M9	
	keyboard	Ckey	Left	NA	Left	M9	
M9Dkey217.1410		Dikey	Left	NA	Left	M9	
M9Ekey217.141044		Ekey	Left	NA	Left	M9	
M9Enter217.1410		Ente	Right	NA	Right	M9	
M9Fkey217.141065		Fkey	Left	NA	Left	M9	
M9Gkey217.1410		Gkey	Left	NA	Left	M9	
	keyboard	Hkey	Right	NA	Right	M9	
M9Indi217.141066		finger_tip	NA	Left	Left	M9	
M9Indr217.140998		finger_tip	NA	Right	Right	M9	
	keyboard	Kkey	Right	NA	Right	M9	
M9Midt217.141043		finger_tip	NA	Left	Left	M9	
M9Midr217.141060		finger_tip	NA	Right	Right	M9	
	keyboard	Mkey	Right	NA	Right	M9	
	keyboard	Nkey	Right	NA	Right	M9	
	keyboard	Okey	Right	NA	Right	M9	
	human_skin	finger_tip	NA	Left	Left	M9	
M9Pinr217.141002		finger_tip	NA	Right	Right	M9	
M9Pkey217.141096		Pkey	Right	NA	Right	M9	
M9Qkey217.1410		Qkey	Left	NA	Left	M9	
	human_skin	finger_tip	NA	Left	Left	M9	
M9Rinr217.141080	human_skin	finger_tip	NA	Right	Right	M9	
M95key217.141004	keyboard	Skey	Left	NA	Left	M9	
M95pace217.141	keyboard	Space bar	Ambiguous	NA	Ambiguous	M9	
M9Thml217,1410	human skin	finger_tip	NA	Left	Left	M9	
	human_skin	finger_tip	NA	Right	Right	A49	
M9Vkey217.1410	keyboard	Vkey	Left	NA	Left	M9	
	keyboard	Wkey	Left	NA	Left	M9	
	keyboard	Xkey	Left	NA	Left	M9	
M9Ykey217,141029	keyboard	Ykey	Right	NA	Right	M9	

Report



Service



Data via SQL

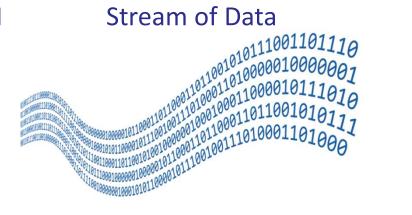




Embeddable KPI



Stream of Data

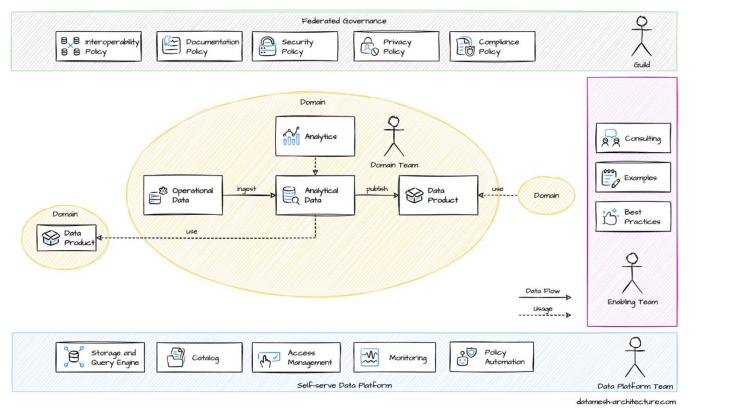


R20

Copyright © 2025 R20/Consultancy B.V., The Netherlands



The Periphery of a Domain

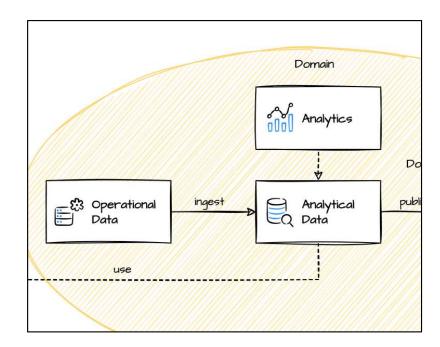


Data Mesh Architecture

Copyright © 2025 R20/Consultancy B.V., The Netherlands

 \mathcal{R}^{20}

What is Ingest?



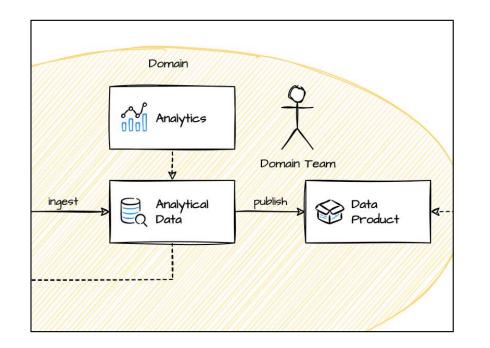
Potential technologies:

- ETL
- Data replication (Change Data Capture)
- ESB (Enterprise Service Bus)
- Messaging
- Database triggers
- Data virtualization
- Involves data processing specifications
 - Transforming data values and structure
 - Masking
 - Cleansing
 - Calculations
 - ...



Copyright © 2025 R20/Consultancy B.V., The Netherlands

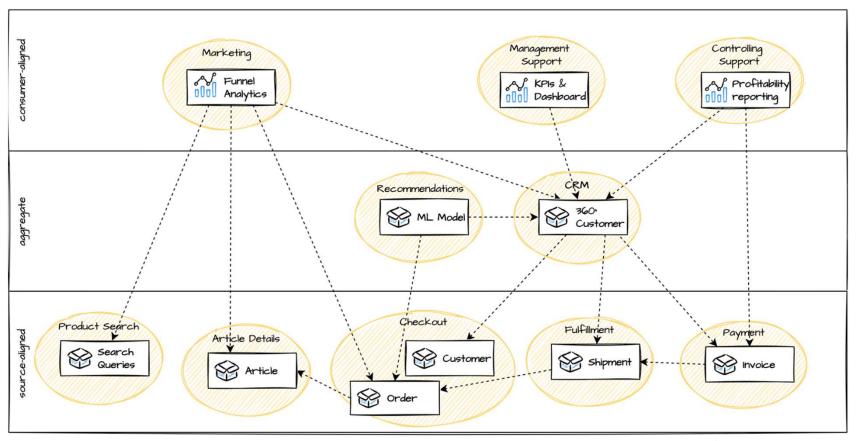
Storing Analytical Data



- One database-solution per domain
 - Multiple databases for different use cases
- Different architectural solutions
 - Data warehouse
 - Data warehouse + data marts
 - Data lake
 - Data lakehouse
 - Data hub
 - Translytical database



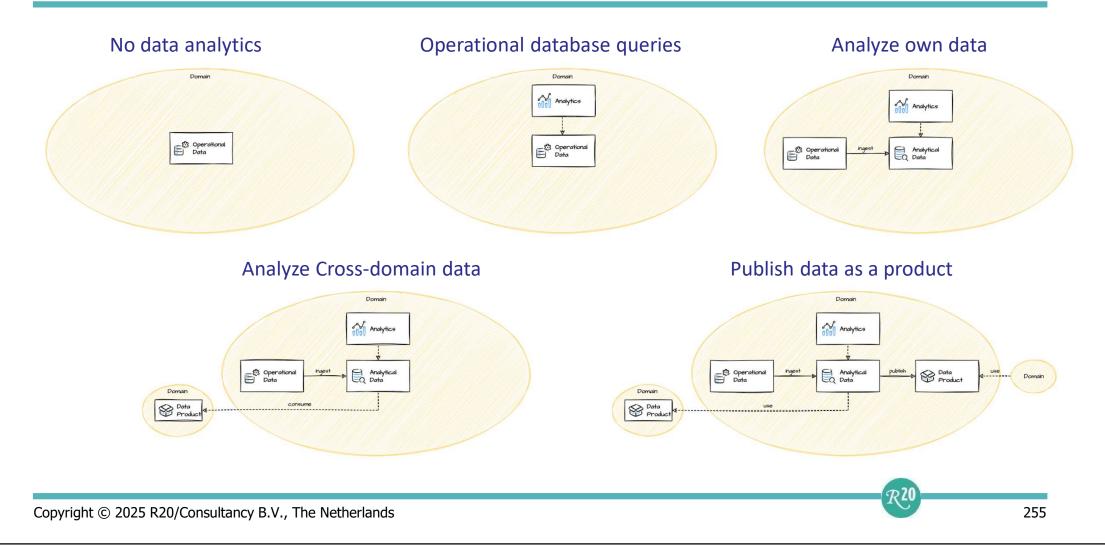
The Mesh Itself



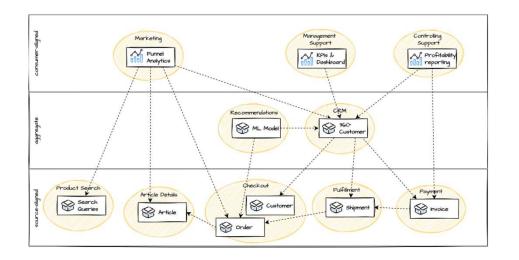
datamesh-architecture.com

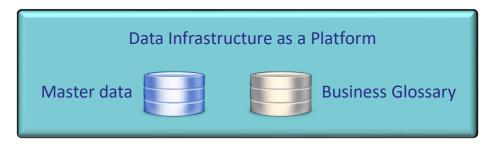
 \mathcal{R}^{20}

The Domain Team's Journey



The Foundation: Data Infrastructure as a Platform





- Scalable polyglot big data storage
- Encryption for data at rest and in motion
- Data product versioning
- Data product schema
- Data product de-identification
- Unified data access control and logging
- Data pipeline implementation and orchestration
- Data product discovery, catalog registration and publishing
- Data governance and standardization
- Data product lineage
- Data product monitoring/alerting/log
- Data product quality metrics (collection and sharing)
- In memory data caching
- Federated identity management
- Compute and data locality
 -

R20

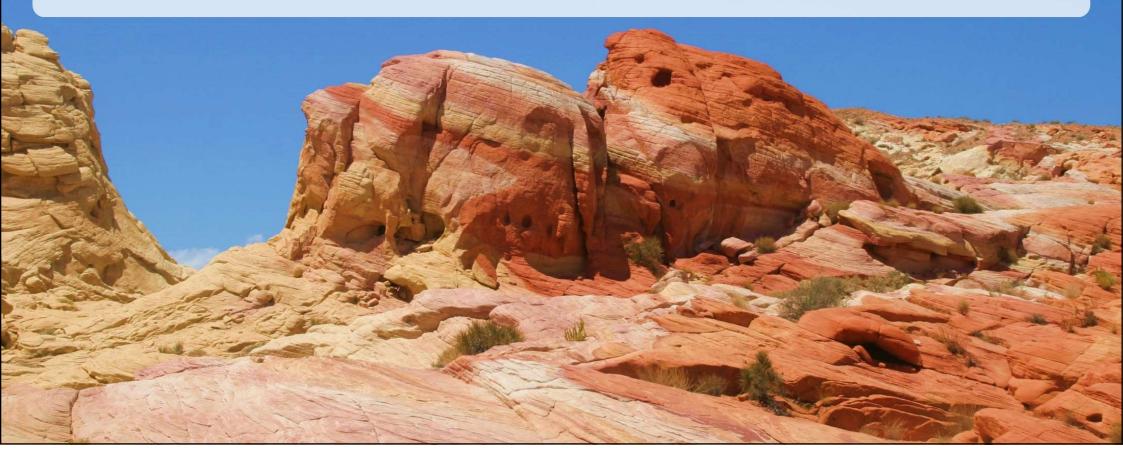
Data Mesh Challenges

produces and produces	Marrael Sold Par	mar 1	Persymmet Support Dombissered	Controlling Support
ารรู้เรียง		Rourres M.	and a state of the	
1 1 1	search	e Dout	Ruinsen	Payrigent Brances

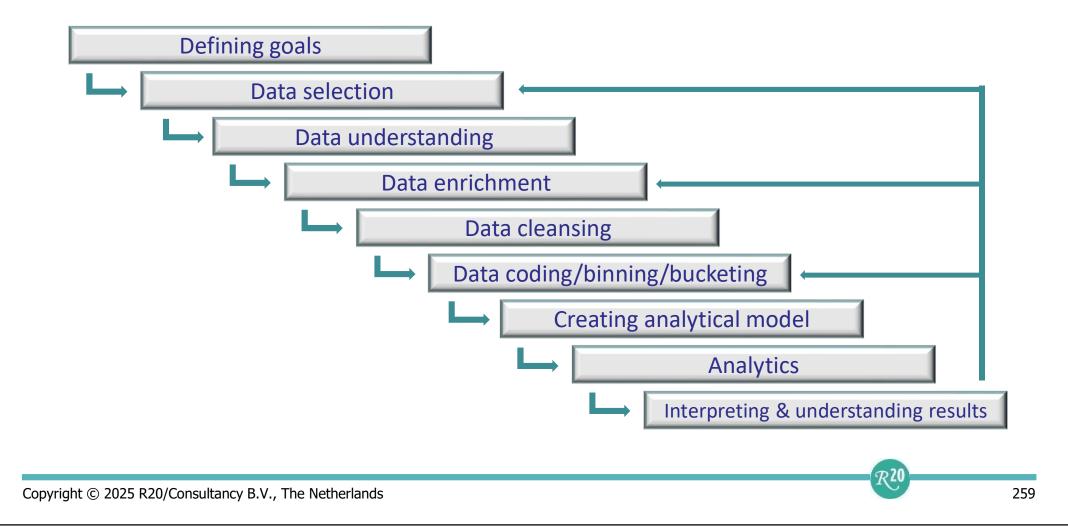
- Massive organizational change for IT
- Standardization of interfaces required
- How clear are BI requirements when new transactional applications are developed?
- Development of transactional applications more complex
- What about transactional applications that are bought?
 - Do they need to be wrapped?
- What about multi-domain transactional applications?
- Distribution of knowledge of data delivery technologies



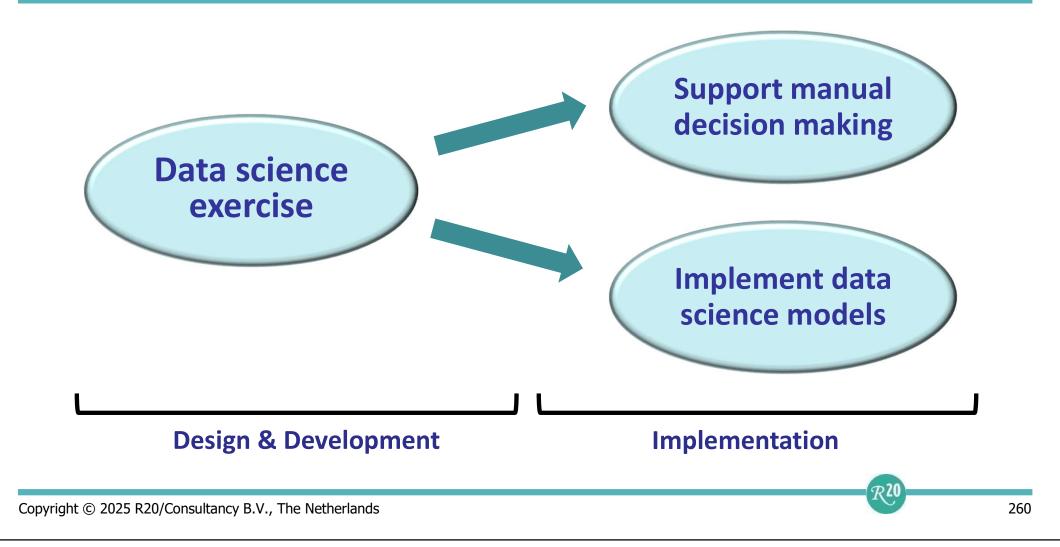
Part 6.8: Embedding Data Science Models in Data Architectures



Development Steps for Data Science



Operationalization of Data Science Models



Operationalization of Data Science Models

Type of Decision	Example	
Singular manual decision	What will be the impact on total sales of acquiring company X?	
Repeatable manual decision	Does a specific location have the right characteristics for opening a new shop?	
Partial automated decision	In a call center: What is the churn risk for a customer? What should we offer?	
Full automated decision without automated reaction	When credit card payment is dubious, send message to operator	
Full automated decision with automated reaction	When sensor indicates the component is heating up too fast, switch off machine	

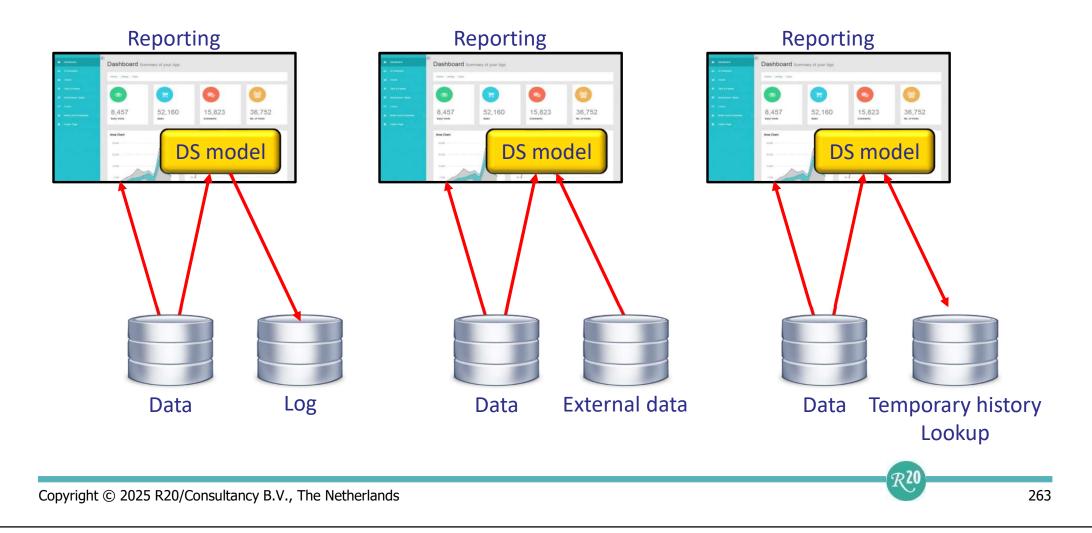
 \mathcal{R}^{20}

Requirements for a Supporting Data Architecture

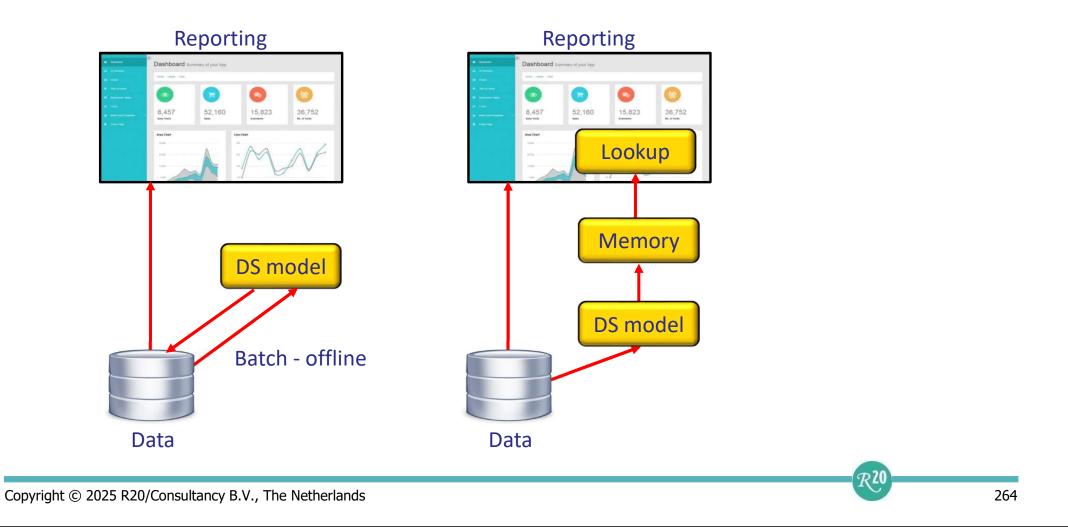
- Versioning of data science models
 - Immutable models
- Auditable data science models
 - Reproducible data for reproducible models
 - Transparency of models
- Different codings must be easy and quick to apply
- Self-learning models or not?
- Delivering metadata
 - Descriptions, definitions, tags, relationships, searchable
- Fast evaluation of models

- Max time to execute model, SLAs
- And many more ...

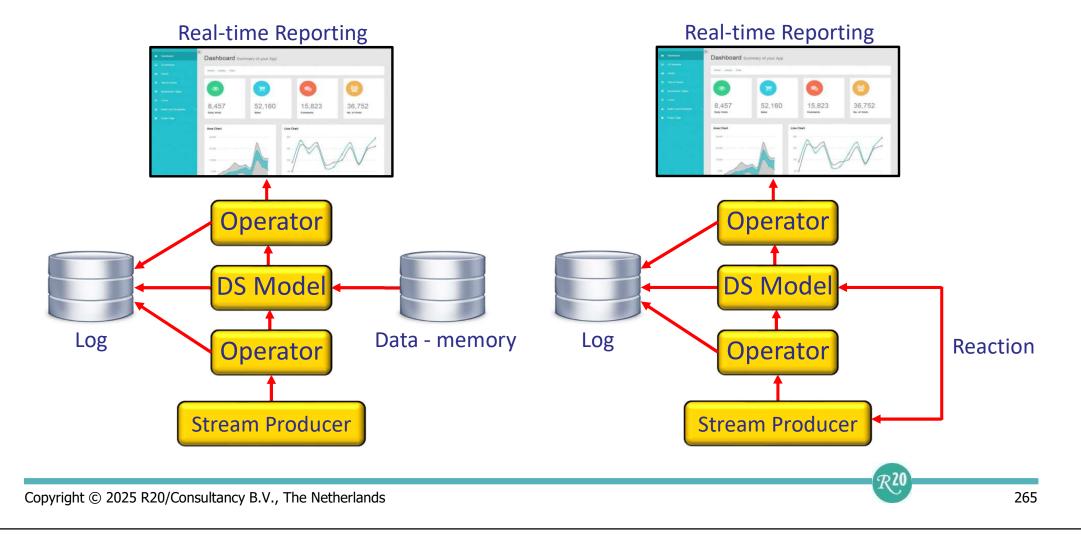
Architectural Aspects (1)



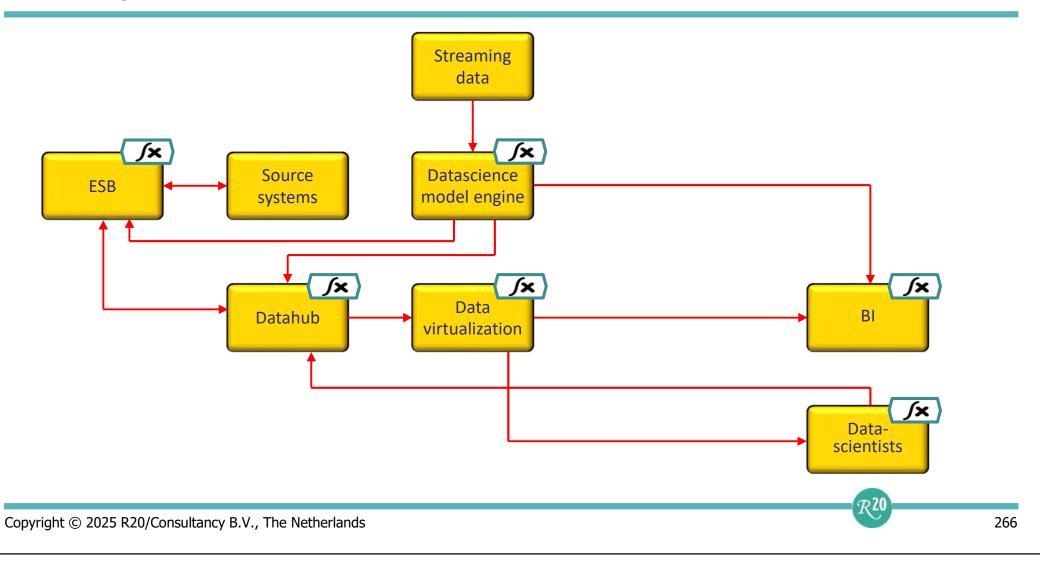
Architectural Aspects (2)



Architectural Aspects (3)



Example of a Data Architecture



Part 6.9: Netflixing Your Data



From Video-by-Sneakers to Video-on-Demand









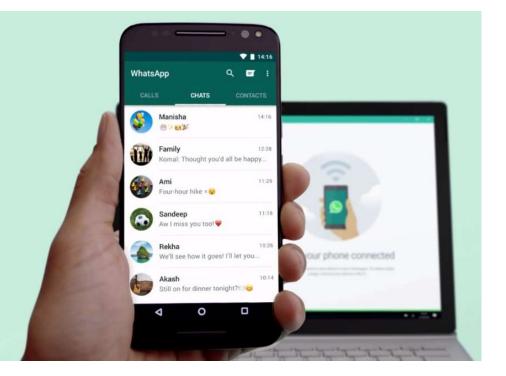
From Music-by-Sneakers to Music-on-Demand





From Message-by-Pigeon to Message-on-Demand









 \mathcal{R}^2

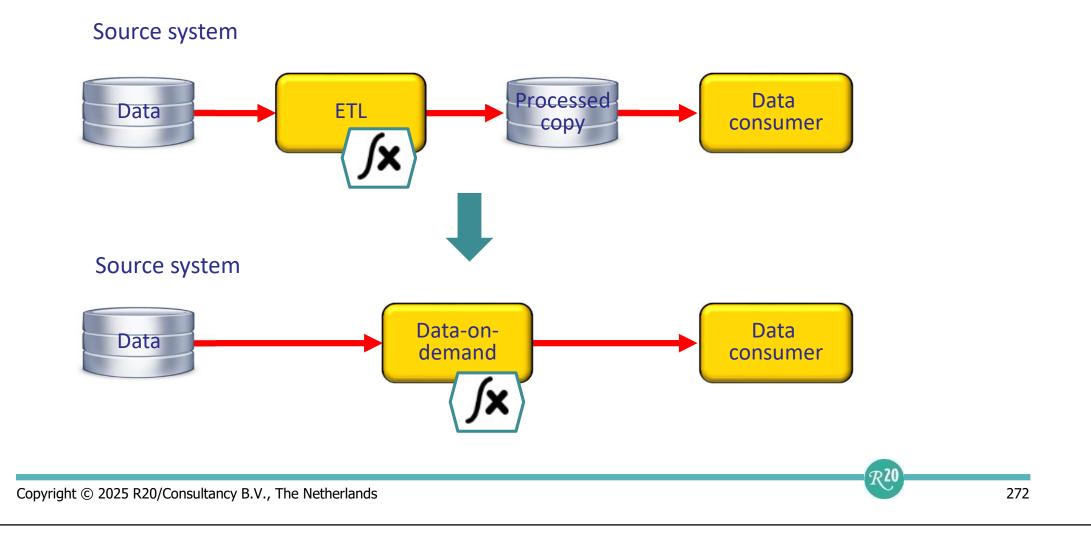
Video-on-Demand

Music-on-Demand

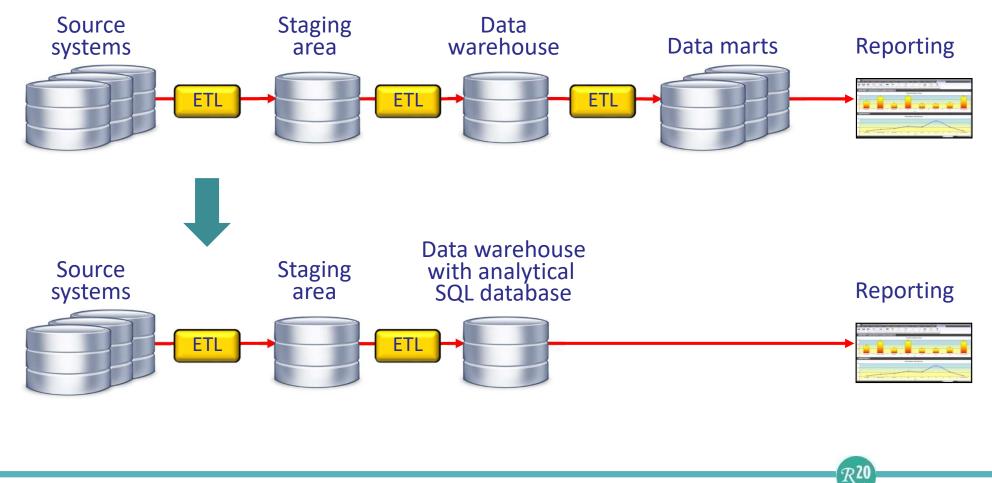
Message-on-Demand

Data-on-Demand?

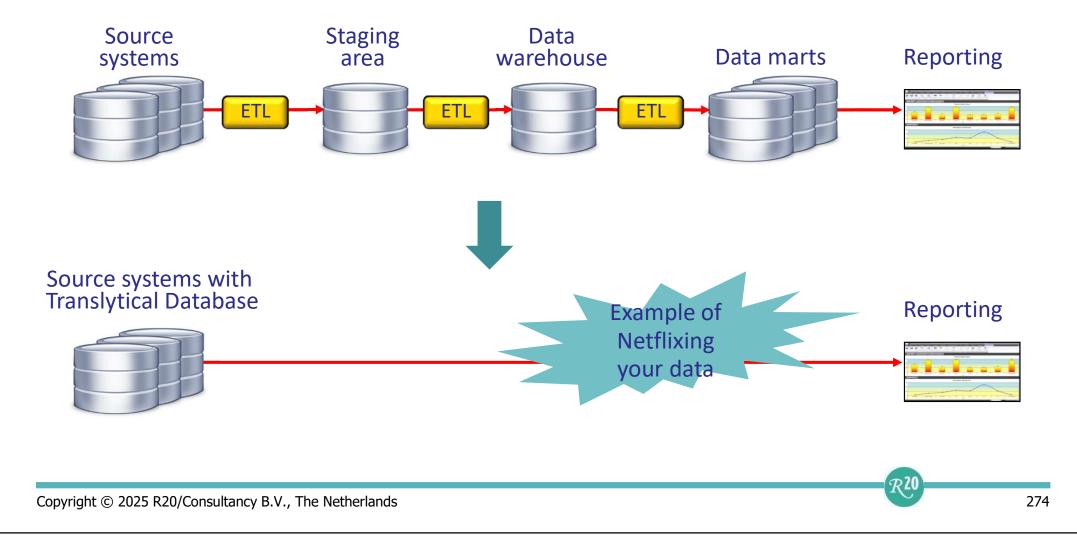
Replace Derived Data by Original Data (1)



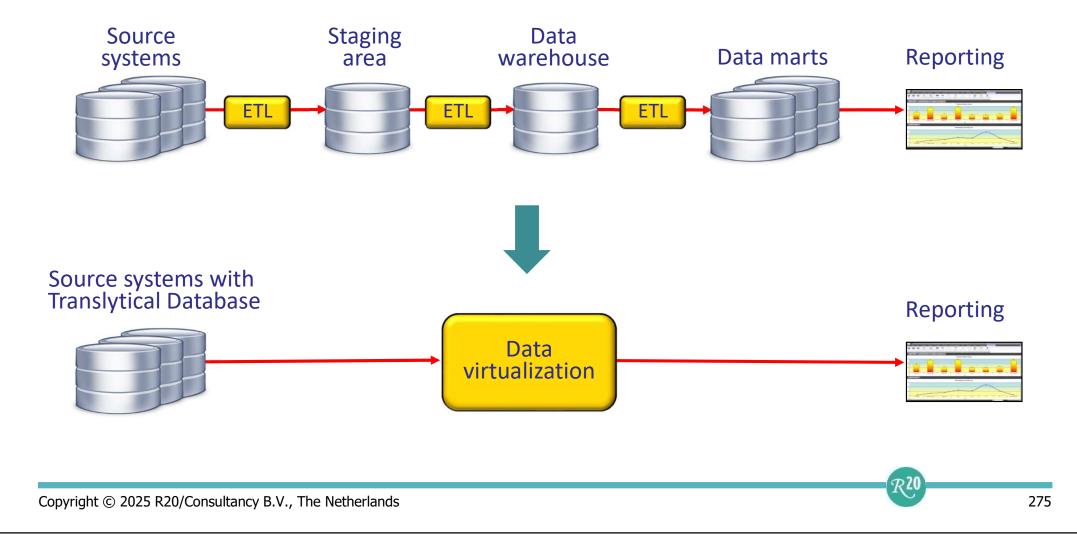
Replace Derived Data by Original Data (2)



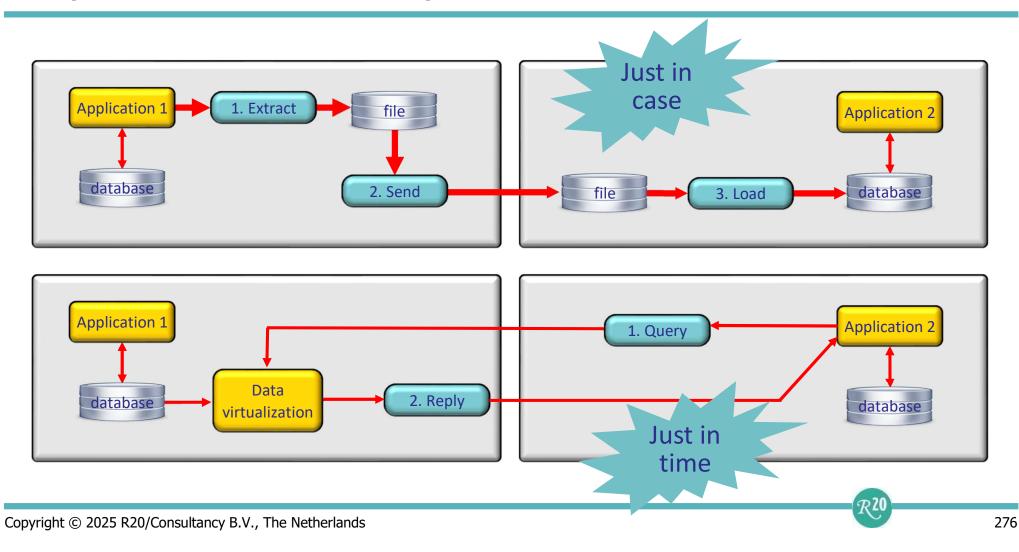
Replace Derived Data by Original Data (3)



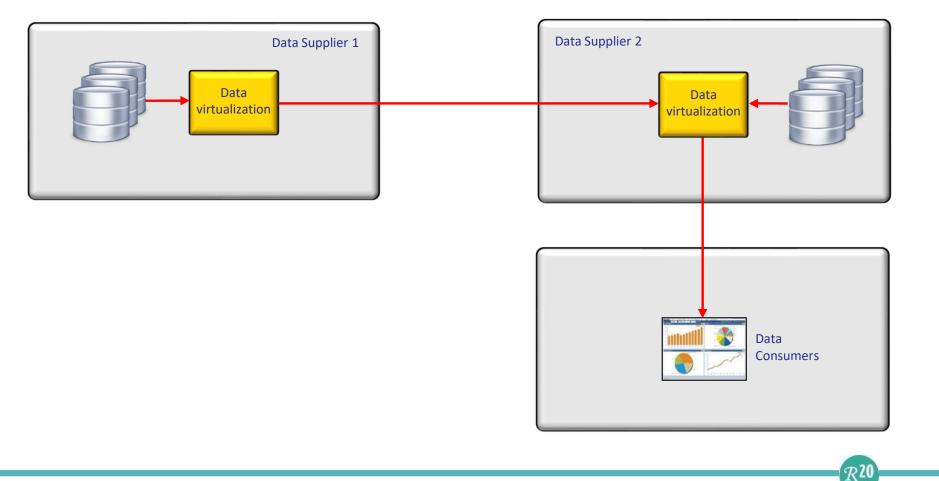
Replace Derived Data by Original Data (4)



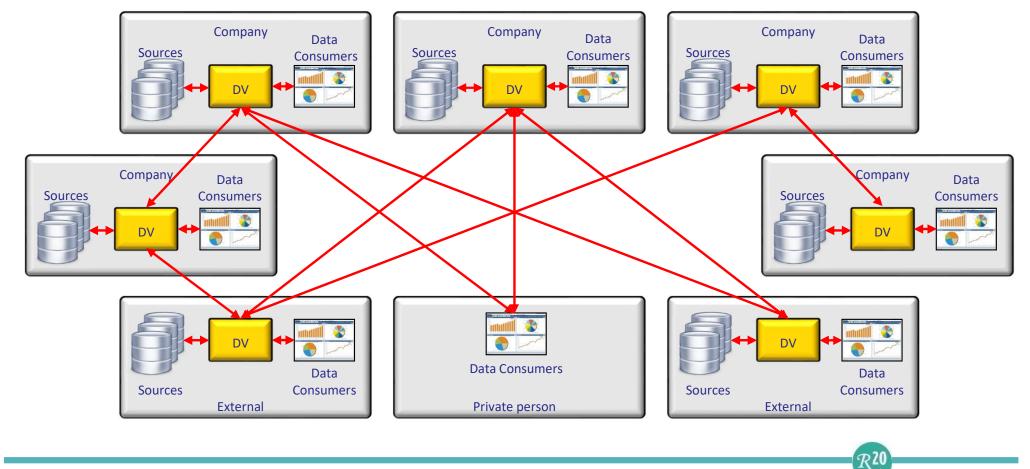
Replace File Transfer by Data-on-Demand



Global Data Architecture Based on Data Minimization



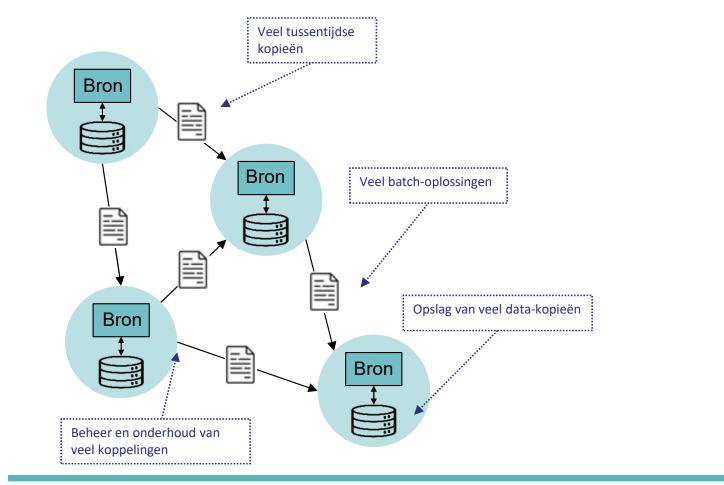
Global Data Architecture Based on Data Minimization



Part 6.10: The Delta Architecture Under Development



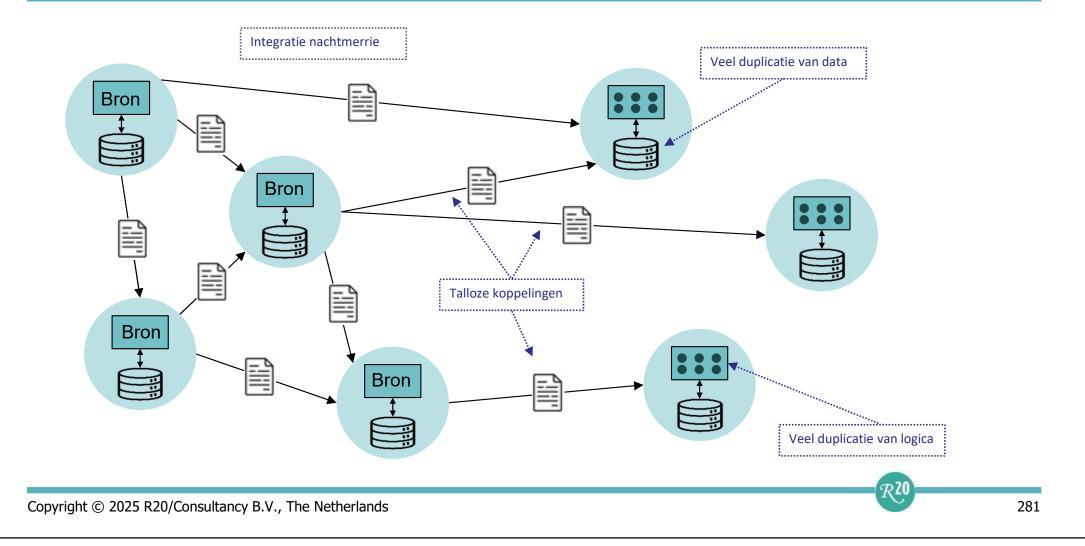
Huidige situatie: Data-uitwisseling tussen bronsystemen



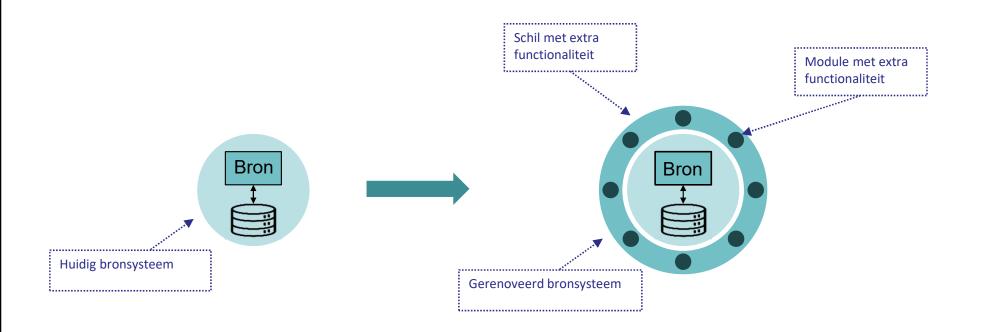
Copyright © 2025 R20/Consultancy B.V., The Netherlands

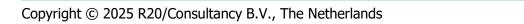
 \mathcal{R}^2

Huidige situatie: Data-uitwisseling tussen bron- en compensatiesystemen



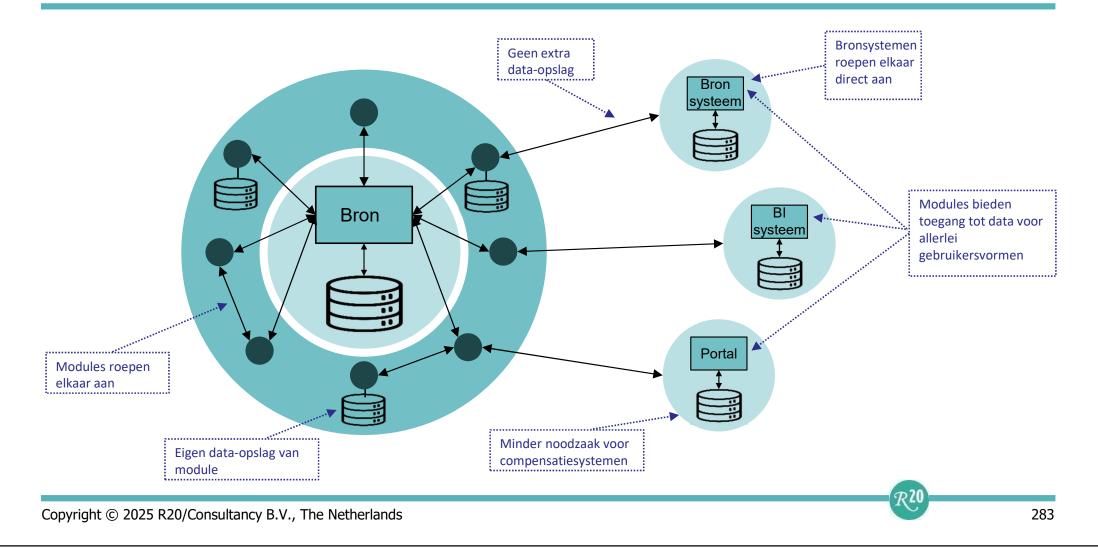
Inkapselen en uitbreiden van bronsystemen





 \mathcal{R}^2

Inzoomen op een ingekapseld bronsysteem



Voorbeelden van nieuwe modules

Data-opvraag

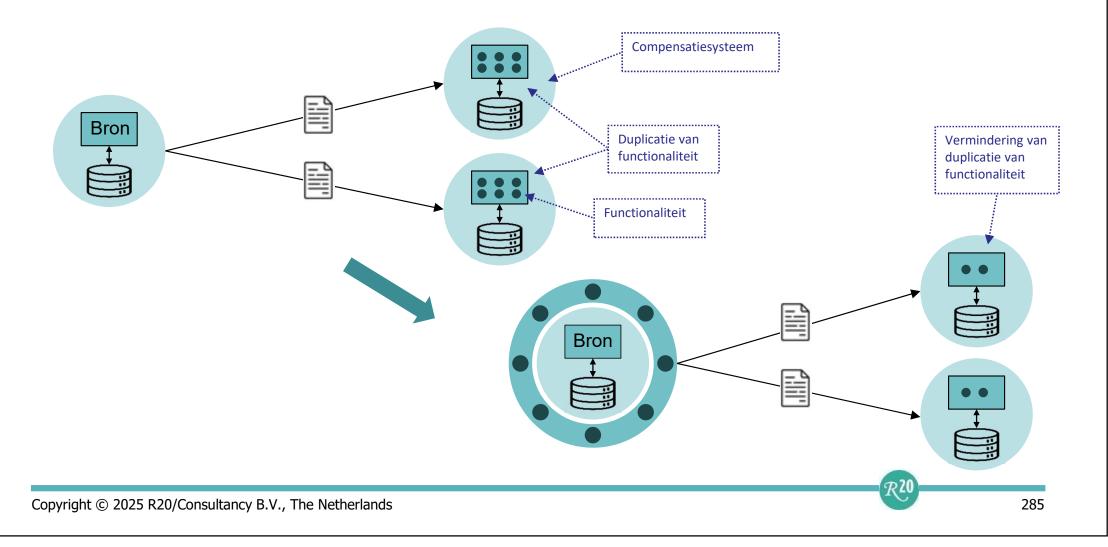
- Opvragen van verzamelingen van business-objecten
- Tijdreizen
- Conform het Enterprise datamodel

CRUD

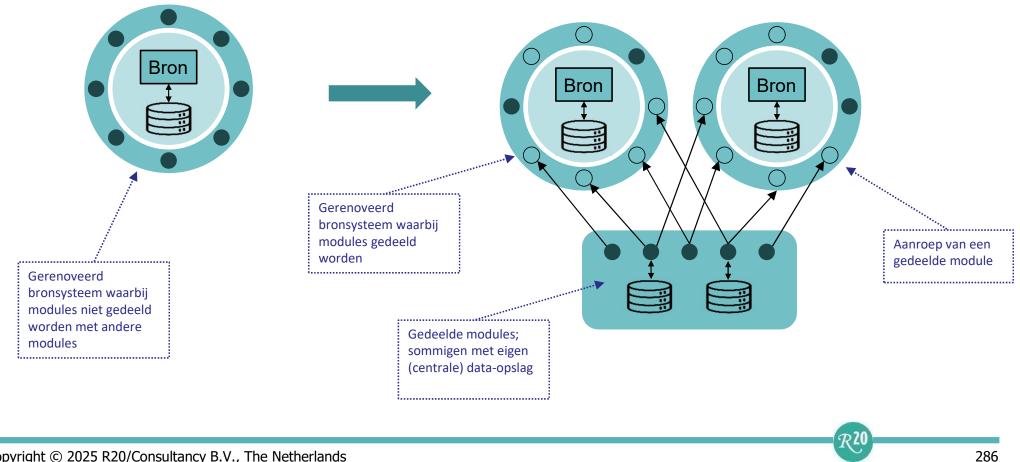
- Bewerken en opvragen van individuele business-objecten
- Conform het Enterprise datamodel
- **Business** logica
 - Bijvoorbeeld: complexe berekeningen, controles en beslissingen
- Databeveiliging
- Datakwaliteit
 - Actief en passief
- Logging
 - Benaderbaar voor analyses
- Metadata



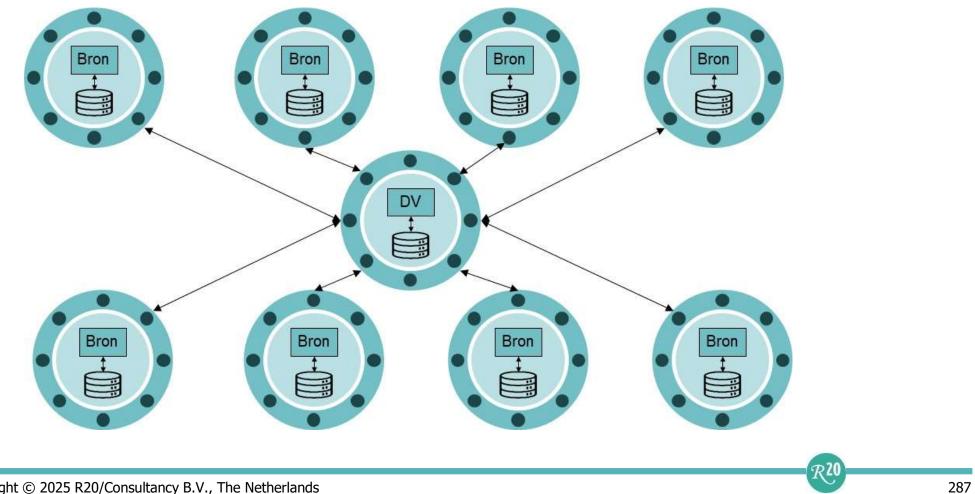
Vermindering duplicatie van functionaliteit



Gebruik van algemene voorzieningen



Overall architectuur voor integratie

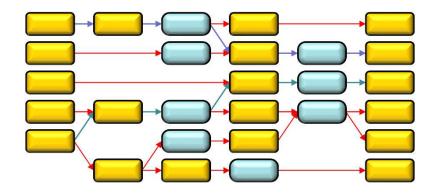


Part 7: Steps 7-8: Design the New Data Architecture, Determine the Implementation Approach



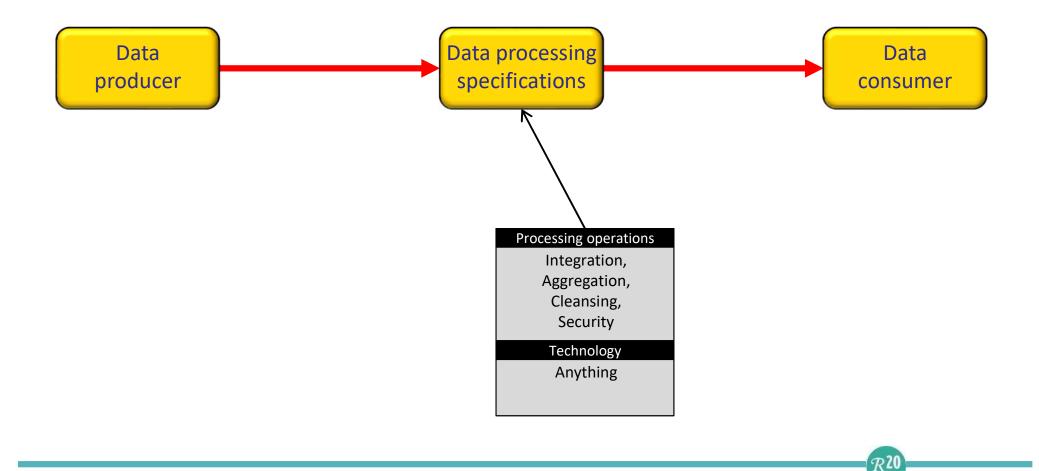
What is a Data Track?

- A data track indicates how data "flows" from data producers to data consumers, and specifies the data processing specifications to be applied and by which module.
- Multiple data consumers can share one data track.
- Data tracks may merge and split.

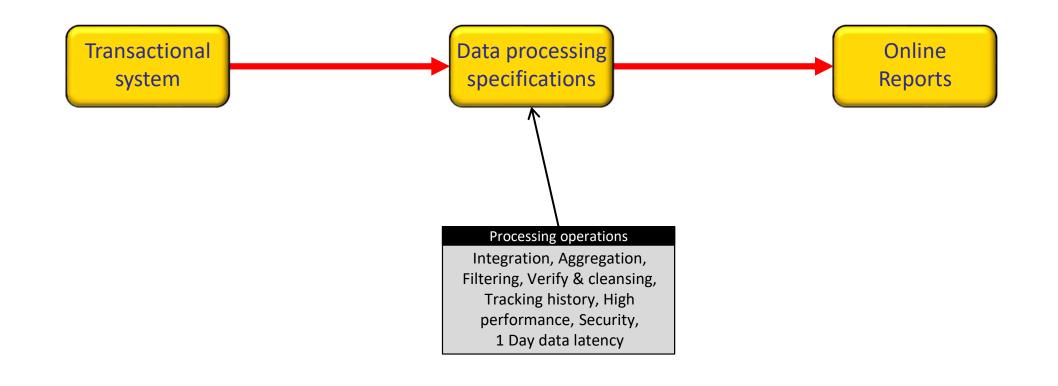


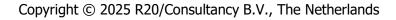


A Data Track Diagram



Data Track Example: Standard Online Reporting (1)







R20

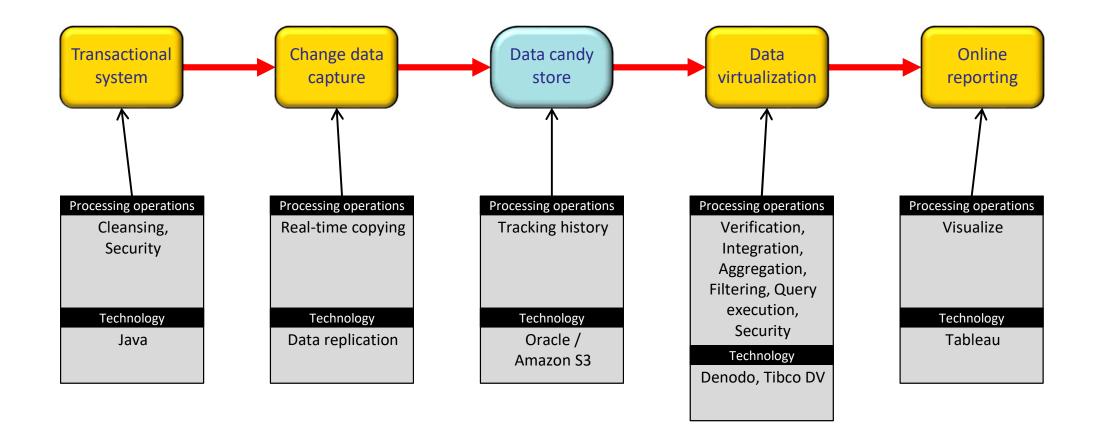
Data Track Example: Standard Online Reporting (2)

Transactional system	Change data capture	Data hub	ETL	Data mart	Online reporting
Processing operations	Processing operations	Processing operations	Processing operations	Processing operations	Processing operations
Cleansing, Security	Real-time copying	Tracking history	Verification, Integration, Aggregation, Filtering, Every	Query execution, Security	Visualize
Technology Java	Technology Data replication	Technology Oracle / Amazon S3	night Technology Talend	Technology GPU database	Technology Tableau



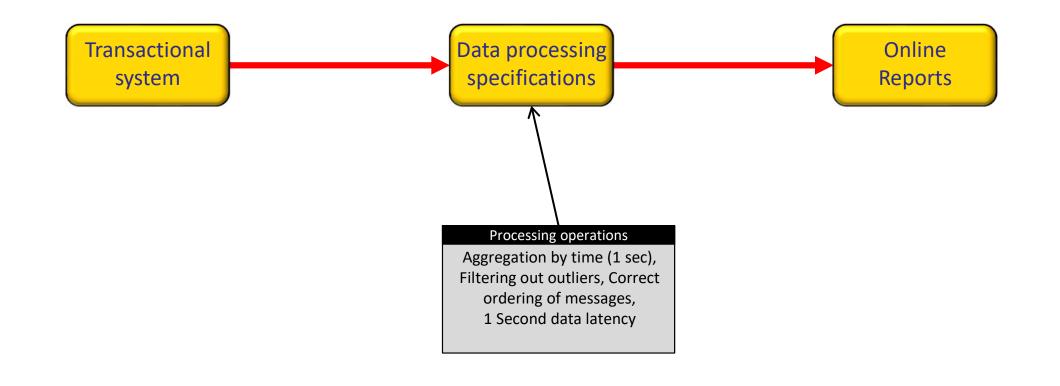
 \mathcal{R}^{20}

Data Track Example: Standard Online Reporting (3)



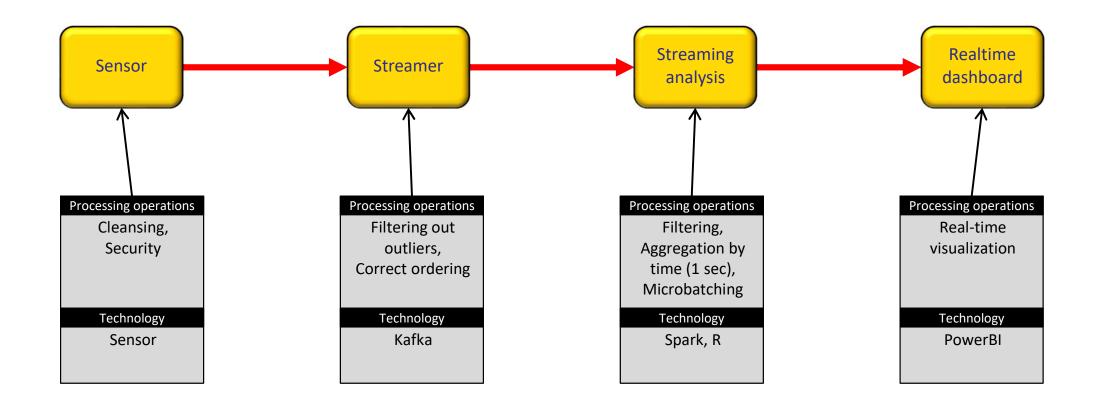
R20

Data Track Example: Streaming Real-time Dashboard (1)



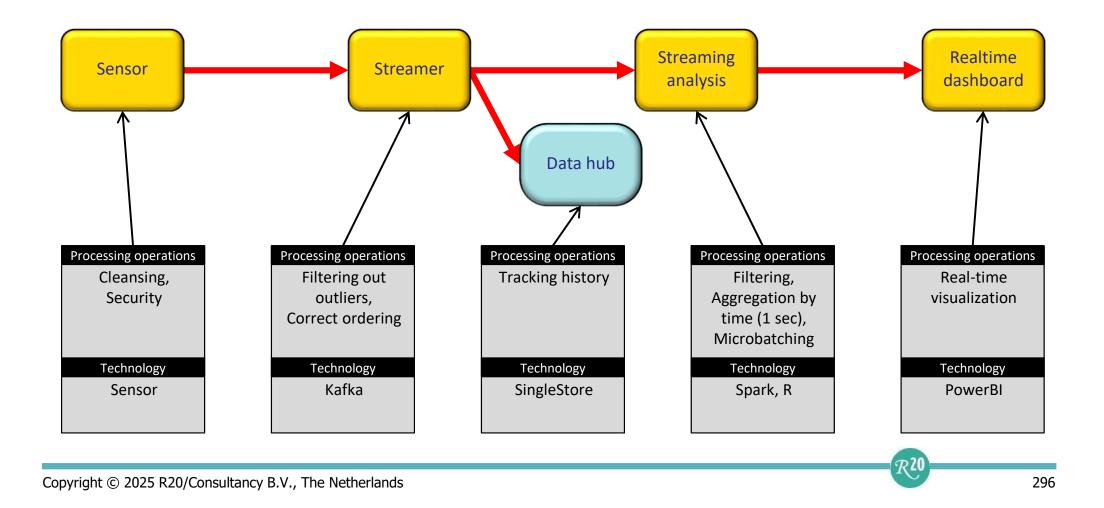
R20

Data Track Example: Streaming Real-time Dashboard (2)

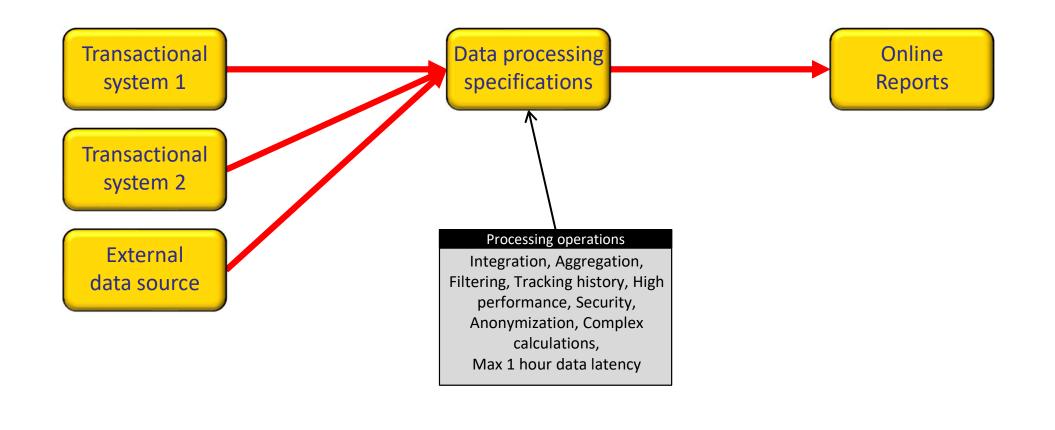




Data Track Example: Streaming Real-time Dashboard (3)



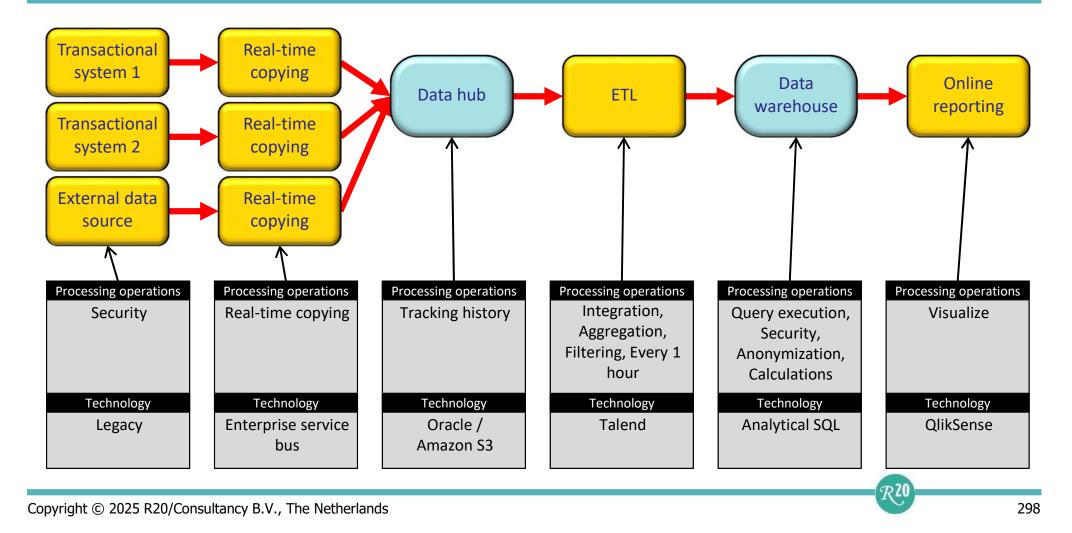
Data Track Example: Integrated Online Reporting (1)



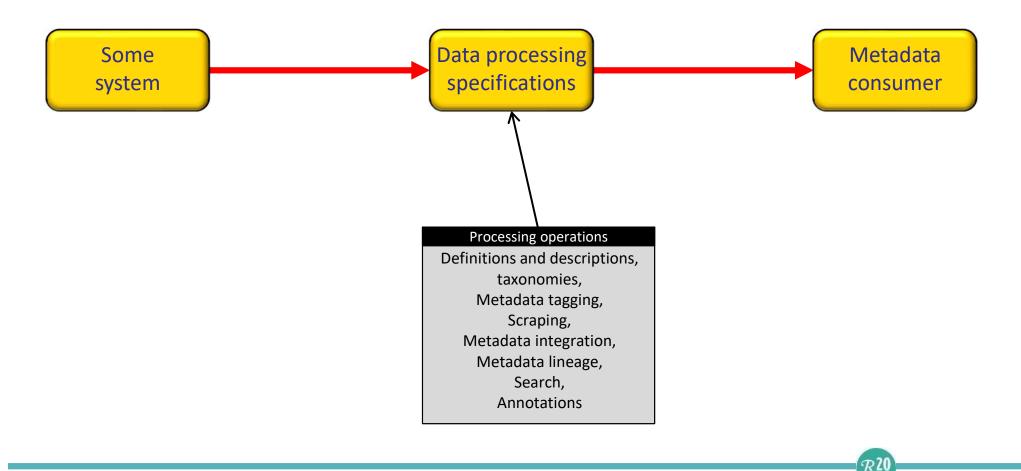
Copyright © 2025 R20/Consultancy B.V., The Netherlands

D?

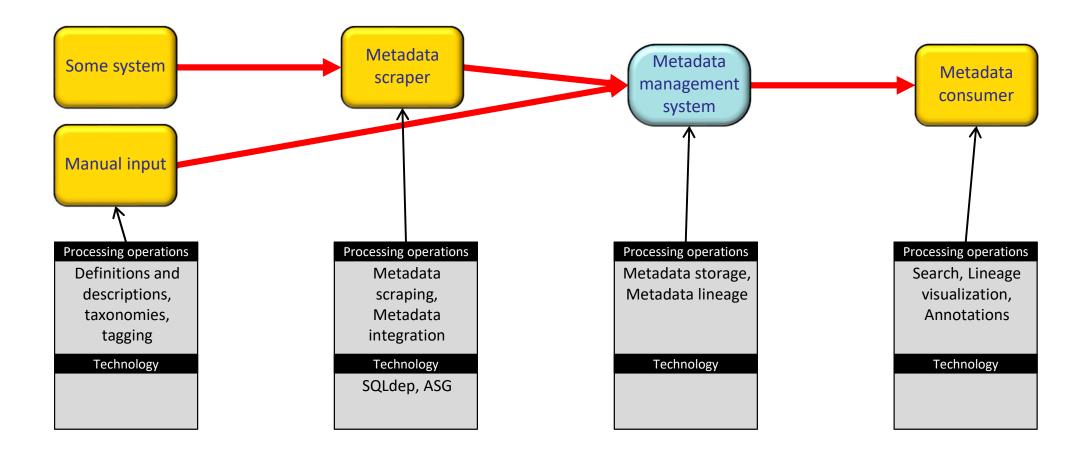
Data Track Example: Integrated Online Reporting (2)

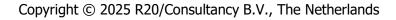


Data Track Example: Metadata Delivery (1)



Data Track Example: Metadata Delivery (2)

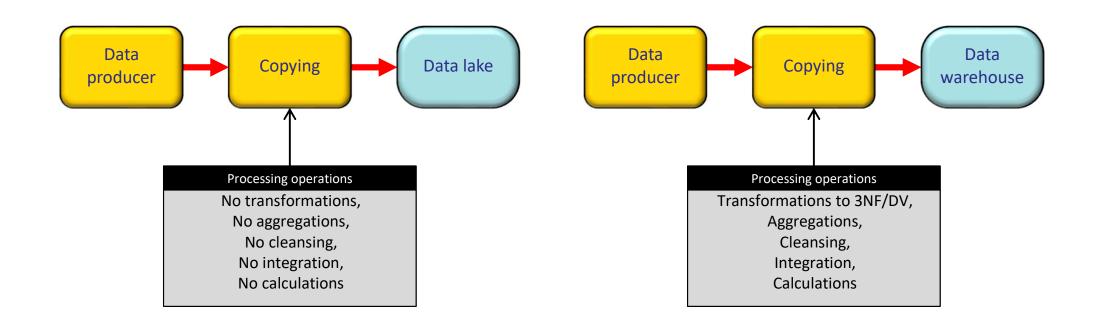






R20

What's in a Name? (1)

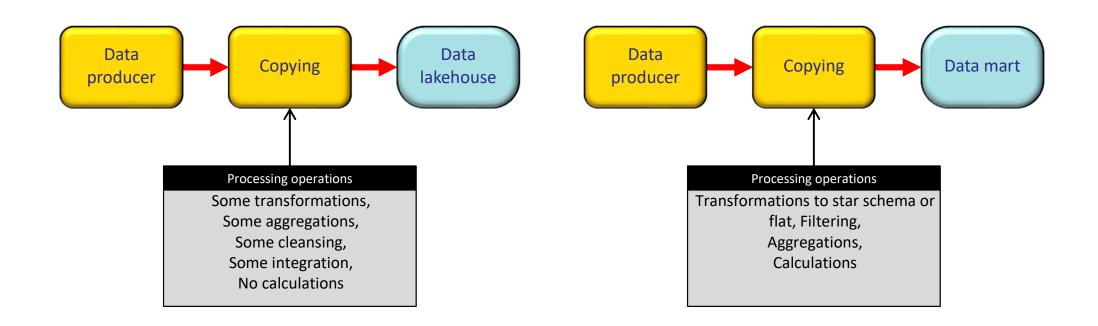


Copyright © 2025 R20/Consultancy B.V., The Netherlands



 \mathcal{R}^{20}

What's in a Name? (2)

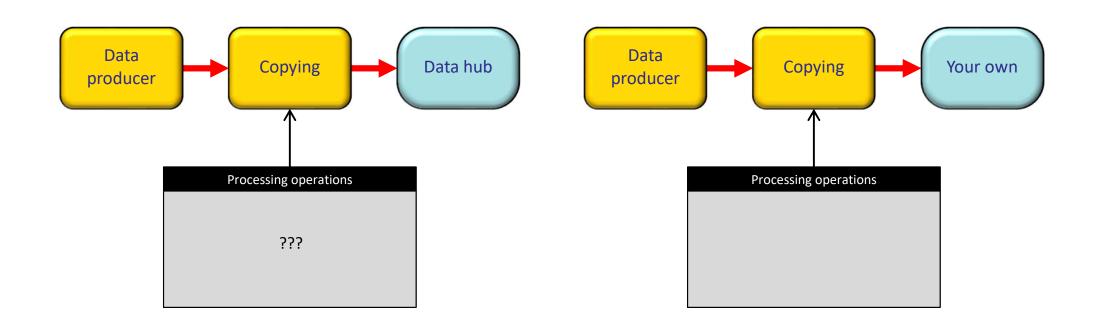


Copyright © 2025 R20/Consultancy B.V., The Netherlands



R20

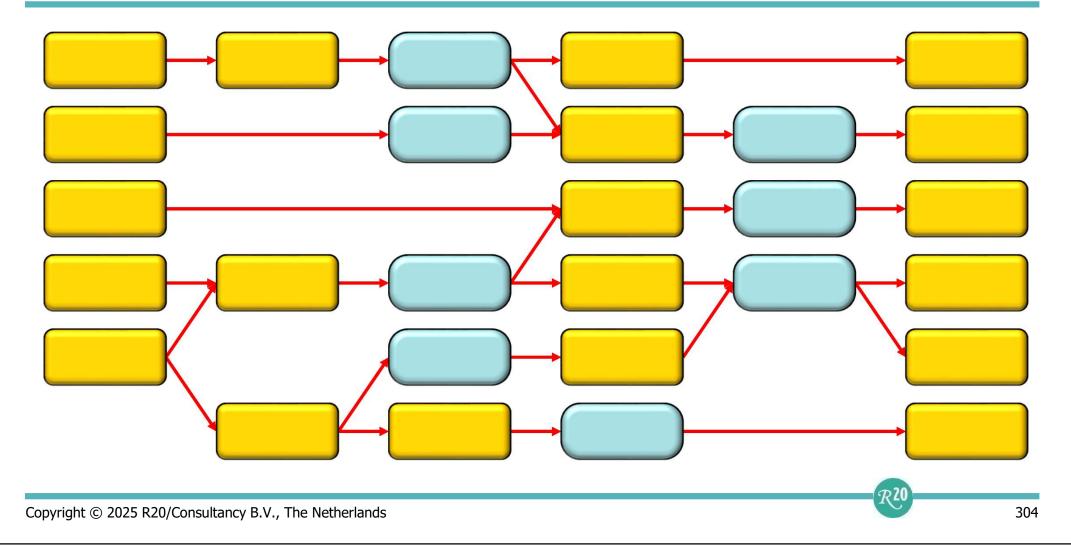
What's in a Name? (3)



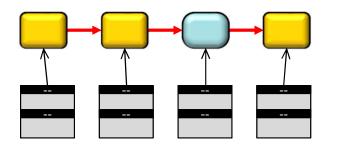


 \mathcal{R}^{20}

High-Level View of the Tracks



Recommendations for Designing Data Tracks



- Design them "backwards" (from consumer to producer)
- Identify data processing specifications first, before assigning of the specs to modules
- One data track can support many comparable data consumers
- Define data track patterns!
 - Architectural design principle?
- Don't solve specific problems

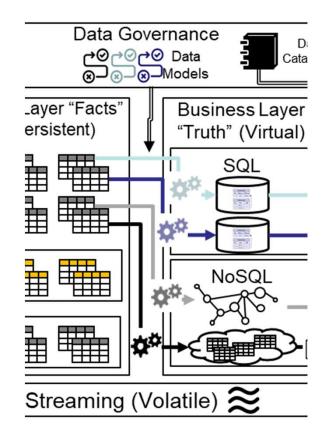


Some More Guidelines for Data Architectures

- 1. Treasure your data processing specifications
- 2. Centralize implementation of data processing specifications
- 3. Centralize technical and business metadata
- 4. Implement abstraction / decoupling
- 5. Make plug and play of technology possible
- 6. Store all data
- 7. Minimize stored data redundancy compute over store
- 8. Choose productivity over performance
- 9. Minimize design exceptions
- 10. Implement cross checks
- 11. Don't send data, let them get it
- 12. Source systems responsible for data quality
- 13. Deploy a holistic design approach



Determine the Intermediate Diagrams



- The current data architecture diagram
 - The new data architecture diagram
 - The dream
 - Will never be reached
- The intermediate data architectures
 - The path from current to new
 - Make the steps as small as possible
 - Preferred: Each step leads to business value
- Think big, act small



Part 8: Steps 9-10: Final Steps



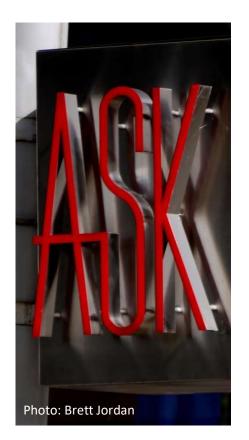
Roadmap for Designing Data Architectures

1. Determine business motivations 2. Determine new requirements 3. Analyze the existing environment 4. Study new products and technologies 5. Define architectural design principles 6. Select a reference data architecture 7. Design the new data architecture 8. Determine the Implementation approach 9. Select new products and technologies 10. Introduce the data architecture within the organization

Part 8.1: Step 9: Select New Products and Technologies



Developing the Request for Information



- Get access to in-depth technological know-how
 - To ask the right questions
- Use distinguishing questions
- Weighs of requirements explain why
- Closed questions easier for comparisons
- Deliver sufficient info to let vendor provide details for pricing and products
- Are extra modules/versions required?
- Remember the new requirements!



Product and Vendor Evaluation (1)



Evaluation of products

- Features, performance, costs, market share
- Local support and partners
 - Experience?
- Extra software required
 - Master data management for complex integration
 - Data cleansing
 - Database server for reference tables and caches
 - Data security
 - Special connectors/drivers



Product and Vendor Evaluation (2)



- Products need to "fit" the architecture
- The intended use cases of the products must match use cases of organization
- Do you need the best tool?
 - Remember Betamax and quadrophonic records
- One-stop shopping or best-of-breed?
 - Minimize number of vendors
 - Never independent of zero vendors

Product and Vendor Evaluation (3)

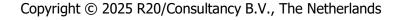


Standardize for back-end tools

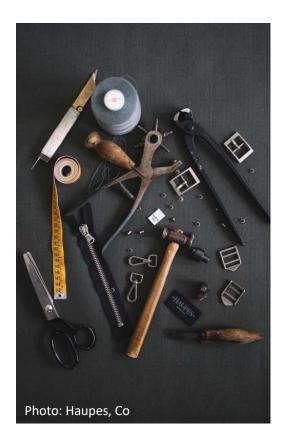
- If use cases allow
- Not one BI tool for all forms of data consumption

Open Source software

- Open source ≠ Non-proprietary software
- Standards = Non-proprietary software
- Study how active the development group is
- What if open source vendor goes commercial?
 - MySQL, Revolutionary Analytics (R), ...



The Proof of Concept/Pilot



- The PoC must be representative of the new system
 - Data: size, characteristics (value distribution, uniqueness), anonymized data?
 - Applications and reports must have a representative complexity
- Performance
 - Multi-user tests
 - Experts required
- Tough SLAs must be tested!
 - Can be expensive
- Invite vendors to install and optimize software themselves
- Select ICT personnel for PoC
 - Developers who enjoy working with new technology and are willing to stay over the weekend



Part 8.2: Step 10: Introduce the Data Architecture Within the Organization



Introduction Within the Organization (1)

Head of Engineering	Head of Supply Chain	Head of Marketing	Head of Finance and Administration
User Experience	Category Management	Digital Marketing	HR
Business Logic	- Procurement	Brand Marketing	— Finance
Administration	Shipping and Logistics	Affiliate Marketing	Facilities

- Three approaches for introduction
 - Via business
 - Via IT bottom up (Trojan horse)
 - Via IT top down
- Identify resistance
 - DBA, source owners, ...
- Educational/missionary program for everyone
 - From programmers to C-level management
 - De-mystify
 - Refute the mythical performance problem
 - Sell the new data architecture
 - Find a champion

Introduction Within the Organization (2)

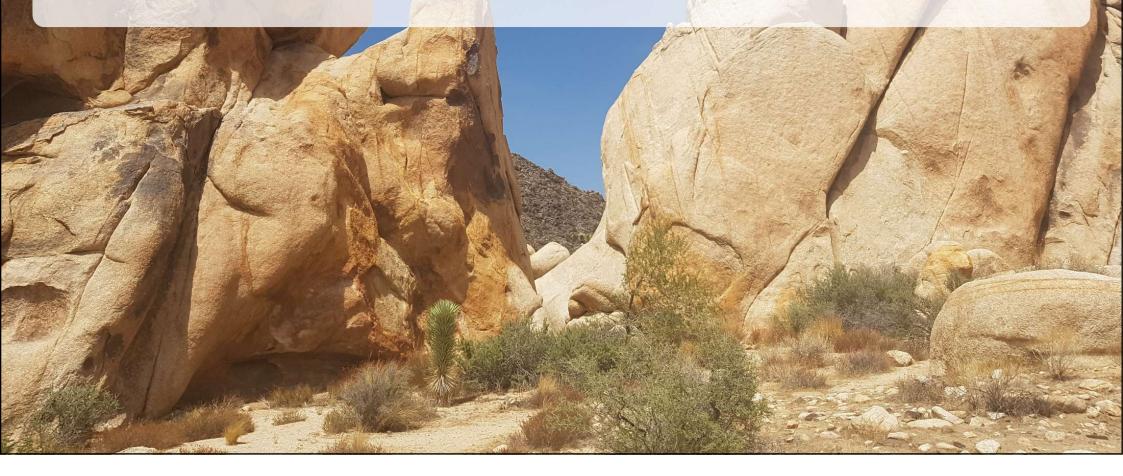
			1	
Head of Engineering	Head of Supply Chain	Head of Marketing	Head of Finance and Administration	
User Experience	Category Management	Digital Marketing	- HR	
Business Logic	- Procurement	 Brand Marketing 	— Finance	
- Administration	Shipping and Logistics	Affiliate Marketing	Facilities	

Impact of new data architecture on organization

- New roles and new responsibilities, examples
 - New roles related to data stewardship
 - Ownership of data
 - Data science models to support business users cooperating
 - New BI tools
- Training

R20

Part 9: Closing Remarks



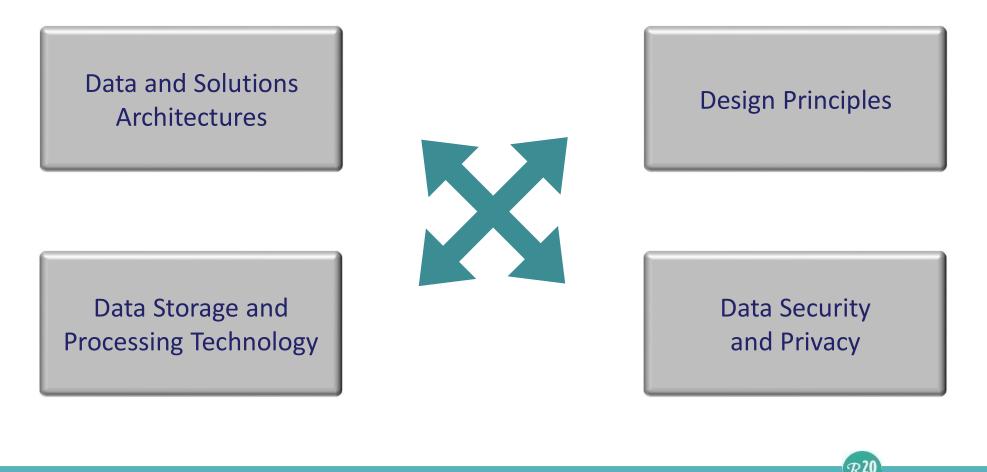
Time to Start!



- New data architectures are required
- Focus on data processing specifications, before drawing the storage "boxes"
- Architects must be familiar with the strengths, weaknesses, and use cases of data storage and processing technologies
 - Without this knowledge:
 - Unnecessarily complex architecture
 - Incorrect use of technology
 - Not able to use the full power of a technology
- Design guidelines impact architecture



From a Linear to a Holistic Approach



Whitepapers by Rick van der Lans – www.r20.nl

VVI	lite	P	ipe	r
		= =		
		ΞΞ		
		ΞΞ		
		ΞΞ		
		ĒĒ		
		= =		

- Streamlining Self-Service Business-Driven Development, November 2022
 - Data Virtualization in the Time of Big Data, April 2022
- Master Data Management for and by Business Users, December 2021
- Logical Data Fabric to the Rescue: Integrating Data Warehouses, Data Lakes, and Data Hub, November 2021
- SingleStore: Simplifying Data Architectures with a Unified Database, July 2021
- Data Fabrics for Frictionless Data Access, April 2021
- Raising the Bar for Data Virtualization, September 2020
- Overcoming Cloud Data Silos with Data Virtualization, June 2020
- Modernizing Data Architectures for a Digital Age Using Data Virtualization, October 2019
- The Business Benefits of Data Virtualization, May 2019
- The Fusion of Distributed Data Lakes Developing Modern Data Lakes, February 2019
- Unifying Data Delivery Systems Through Data Virtualization, October 2018
- Architecting the Multi-Purpose Data Lake With Data Virtualization, April 2018
- The Next Wave of Analytics At the Edge, December 2017
- Data Virtualization in the Time of Big Data, December 2017
- Developing a Bi-Modal Logical Data Warehouse Architecture Using Data Virtualization, October 2016
- Designing a Logical Data Warehouse, February 2016
- Designing a Data Virtualization Environment; A Step-By-Step Approach, January 2016
- How Drill Enriches Self-Service Analytics; The Added Value of a SQL-on-Everything Engine; November 2015; sponsored by MapR Technologies
- Strengthening Self-Service Analytics With Data Preparation and Data Virtualization, September 2015

Articles by Rick van der Lans – www.r20.nl

			N
_			
		/	-
_	and and		_
1		_	/
-			

- The Data Lakehouse: Blending Data Warehouses and Data Lakes, April 2022
- Use the Cloud More Creatively, January 2022
- Sustainable Data Architectures Through Data Architecture Automation, October 2021
- Do Unified Databases Make Polyglot Persistence Irrelevant?, October 2021
- From Data Warehouse Automation to Data Architecture Automation, June 2021
- Data Minimization as Design Guideline for New Data Architectures, May 2021
- Data Fabrics Need to Coexist with data Warehouses and Other Database-Centric Solutions, March 2021
- A Decentralized Master Data Solution using Data Virtualization, February 2021
- Streamlining External Data Access to Enrich Analytics, December 2020
- The Data Mesh, the New Kid on the Data Architecture Block, December 2020
- Developing a Data Fabric, December 2020
- Making Big Data Easy with Data Virtualization, December 2020
- Data Herding Is Not Data Integration!, November 2020
- Benefits of Data Virtualization to Data Scientists, October 2020
- Becoming a Data-driven Organization Requires a Cultural Change, September 2020
- Cohelion, an All-in-one Data Warehouse Factory, July 2020
- How Do You Design New Data Architectures, July 2020
- Eight Data Virtualization Features to Help an Organization Become Data-Driven, June 2020
- New Data Architectures are too Data-Store-Centric, February 2020
- Data Virtualization and SnowflakeDB: A Powerful Combination, January 2020

....