

# Modern Data Engineering in the Cloud

**Dr. Christian Dollfus**  
Dozent

T direkt +41 41 228 22 54  
christian.dollfus@hslu.ch

Luzern

**Dr. Pavlin Mavrodiev**  
Dozent

T direkt +41 76 733 61 66  
pavlin.mavrodiev@alumni.ethz.ch

PART 1 : FOUNDATIONS IN DATA ENGINEERING

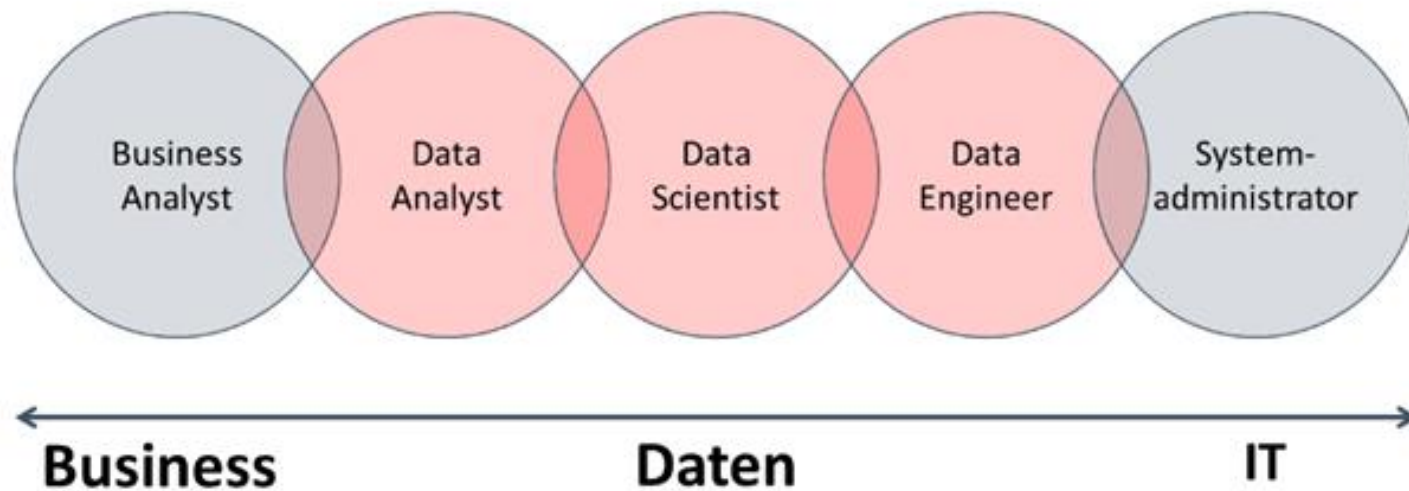
# Introduction and Motivation

- Introductory words
- What is Data Engineering?
- Motivation and Value Proposition - a historic overview
- How does Data Engineering look like in many companies?
- What are features and advantages of workflow-based ETL?
- Motivation for the use of ETL tools.
- Data Engineering and Business Process Automation - similarities and differences

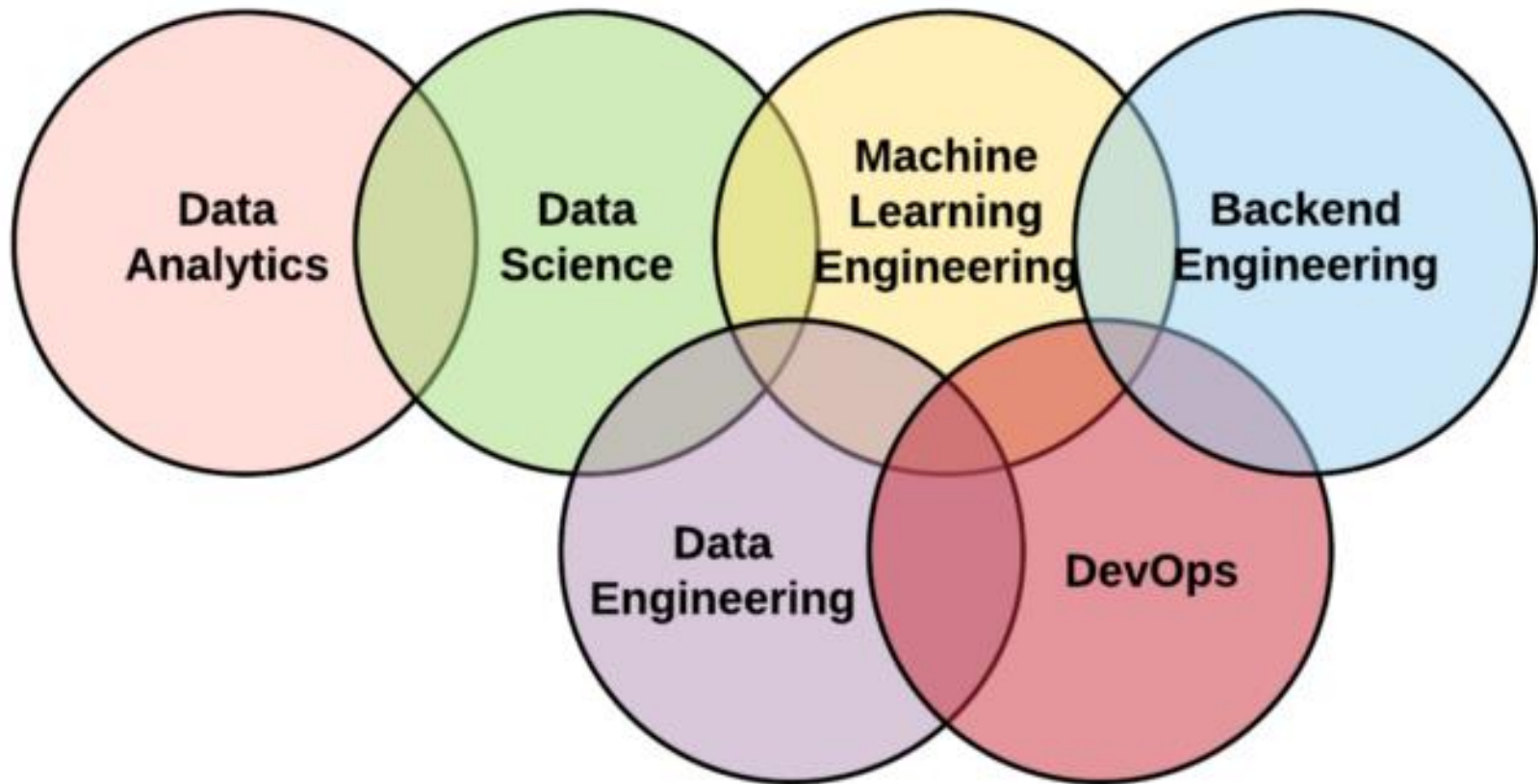
## Learning Goals

- After the 5 days you are able to **engineer automatic batch based and realtime Data-Pipelines** with PDI and KAFKA
- You have hands-on experience with 2 most used techniques in DE
- You have an **overview and foundation** of the huge field of **Data Engineering** namely on the processual part (storage is also a topic of Data Lab 1 and 2)
- You have the **possibility to dive into the topic** if you use it. You got information there.

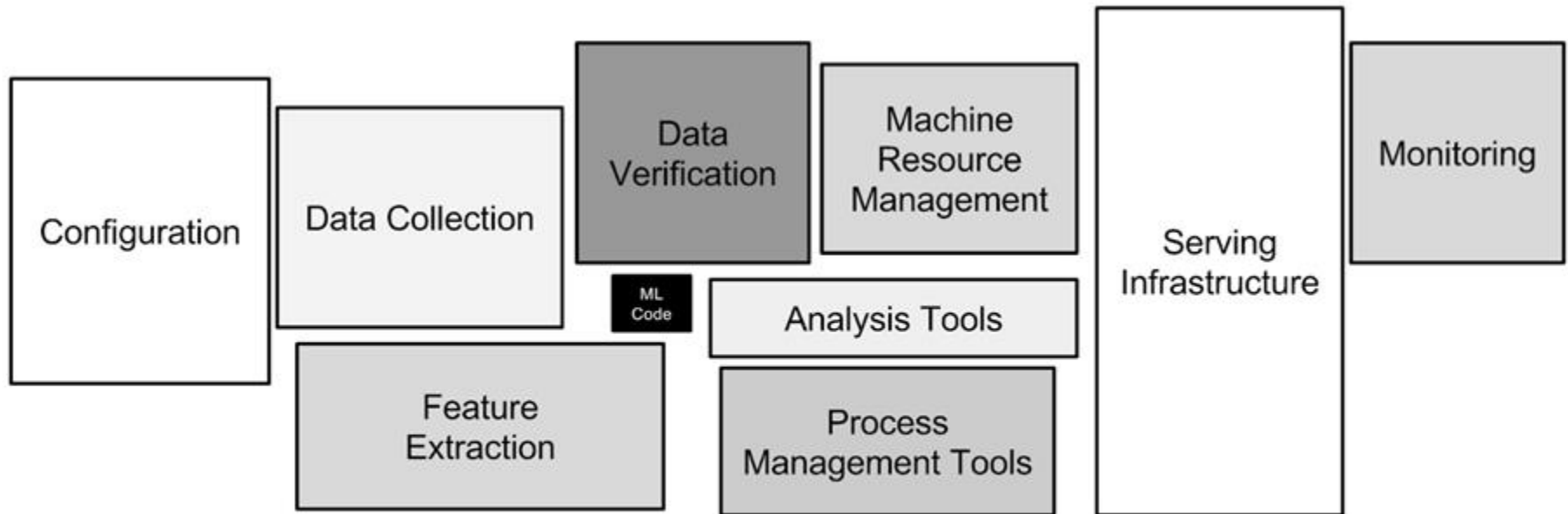
## Introduction and Motivation: Roles



## Introduction and Motivation: Areas

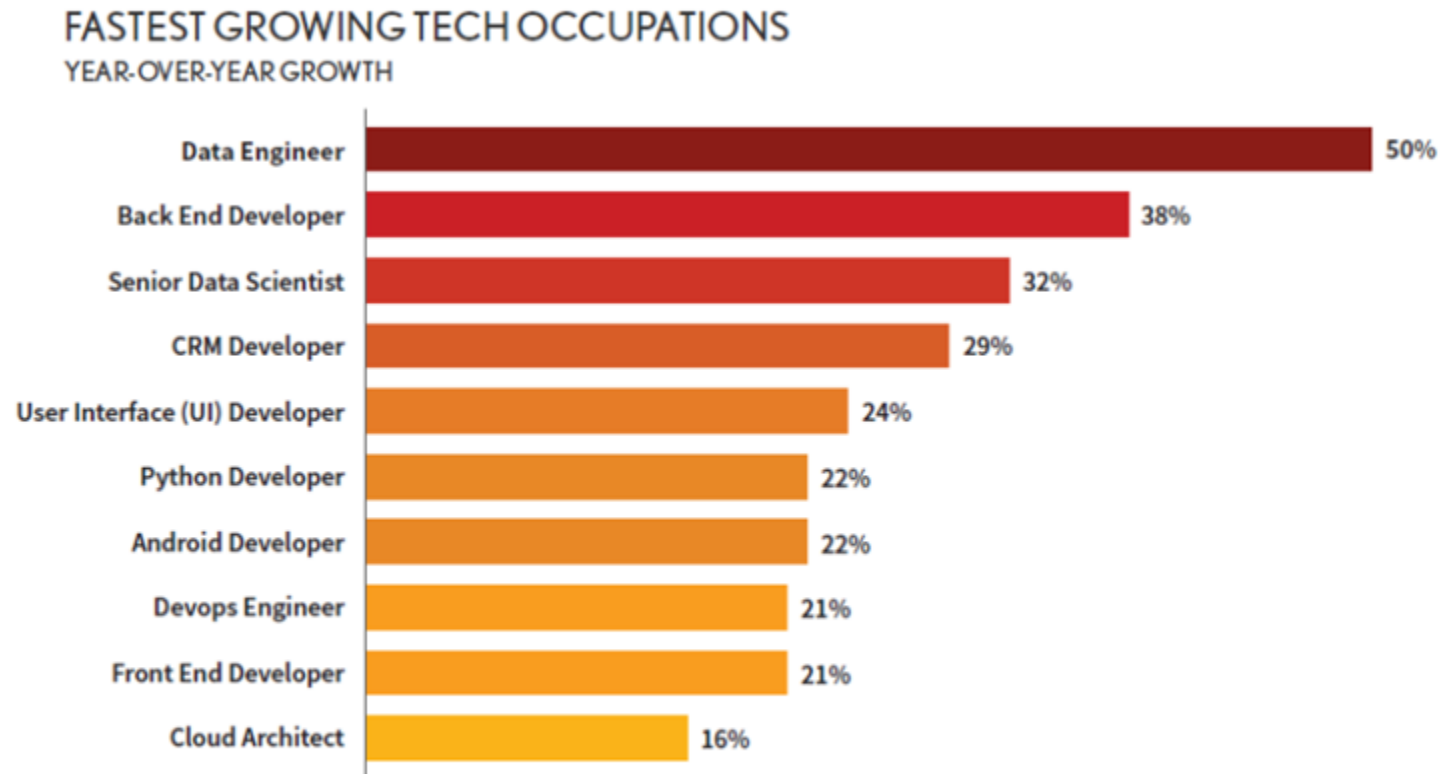


Source: <https://towardsdatascience.com/data-science-is-boring-1d43473e353e>



Source: Sculley, D., Holt, G., Golovin, D. et al. Hidden Technical Debt in Machine Learning Systems

# Introduction and Motivation

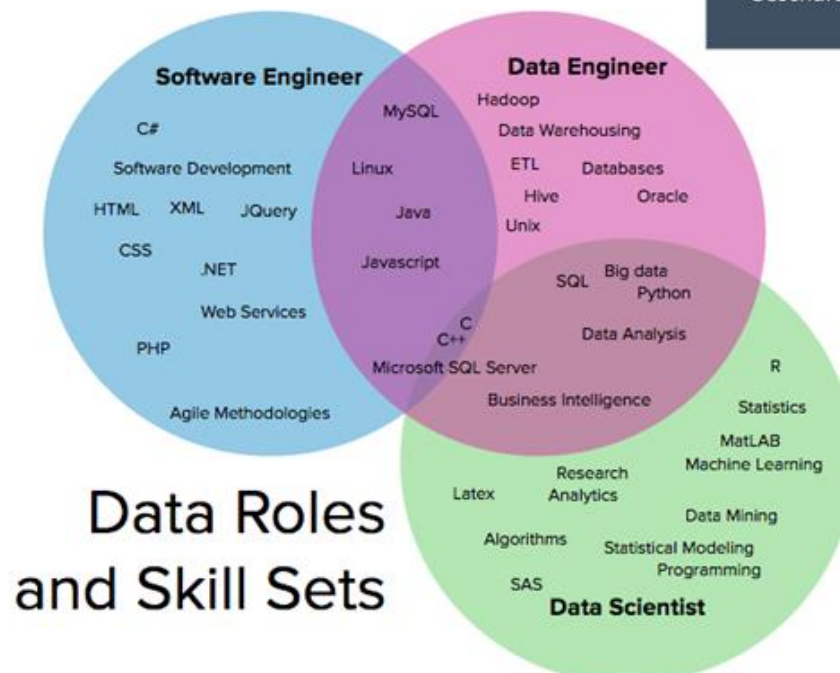


Source: [The Dice 2020 Tech Job Report](#)

# Introduction and Motivation

Data Analyst	Data Scientist	Data Engineer
<b>Schwerpunkt:</b> BWL	<b>Schwerpunkt:</b> Mathematik	<b>Schwerpunkt:</b> Informatik
<b>Wichtigste Skills:</b> <ul style="list-style-type: none"> <li>• Excel, ggf. SQL</li> <li>• BI-Anwendungen</li> <li>• Kommunikation</li> </ul>	<b>Wichtigste Skills:</b> <ul style="list-style-type: none"> <li>• Statistik</li> <li>• R oder Python</li> <li>• Datenbanken</li> <li>• Datenvisualisierung</li> </ul>	<b>Wichtigste Skills:</b> <ul style="list-style-type: none"> <li>• Data Warehouse</li> <li>• ETL (Extraction, Transformation &amp; Load)</li> <li>• Datenbank-Design</li> </ul>
<b>Aufgaben:</b> <ul style="list-style-type: none"> <li>• Reporting</li> <li>• Geschäftsentwicklung</li> </ul>	<b>Aufgaben:</b> <ul style="list-style-type: none"> <li>• Insights</li> <li>• Predictive Analytics</li> </ul>	<b>Aufgaben:</b> <ul style="list-style-type: none"> <li>• Data Pipeline</li> <li>• Datenarchitektur</li> </ul>

[jobs-karriere/data-scientist-oder-data-engineer-was-ist-der-unterschied](https://jobs-karriere/data-scientist-oder-data-engineer-was-ist-der-unterschied)



Quelle : <https://i.pinimg.com/originals/aa/b6/98/aab6987a5979683edada36d01e13862d.png>



If data science is the discipline of making data *useful*, then you can think of data engineering as the discipline of making data *usable*.

Data engineers are the heroes who provide behind-the-scenes infrastructure support for Data Scientists

Data scientists are the data-wrangers, while data engineers are the data-pipeline-wrangers



Quelle : [https://www.linkedin.com/posts/hamed-zitoun-machine-learning-freelance\\_datascience-dataengineering-machinelearning-activity-6730899989536505857-MCDo/](https://www.linkedin.com/posts/hamed-zitoun-machine-learning-freelance_datascience-dataengineering-machinelearning-activity-6730899989536505857-MCDo/) and <https://towardsdatascience.com/data-science-without-any-data-6c1ae9509d92>

# THE DATA SCIENCE HIERARCHY OF NEEDS

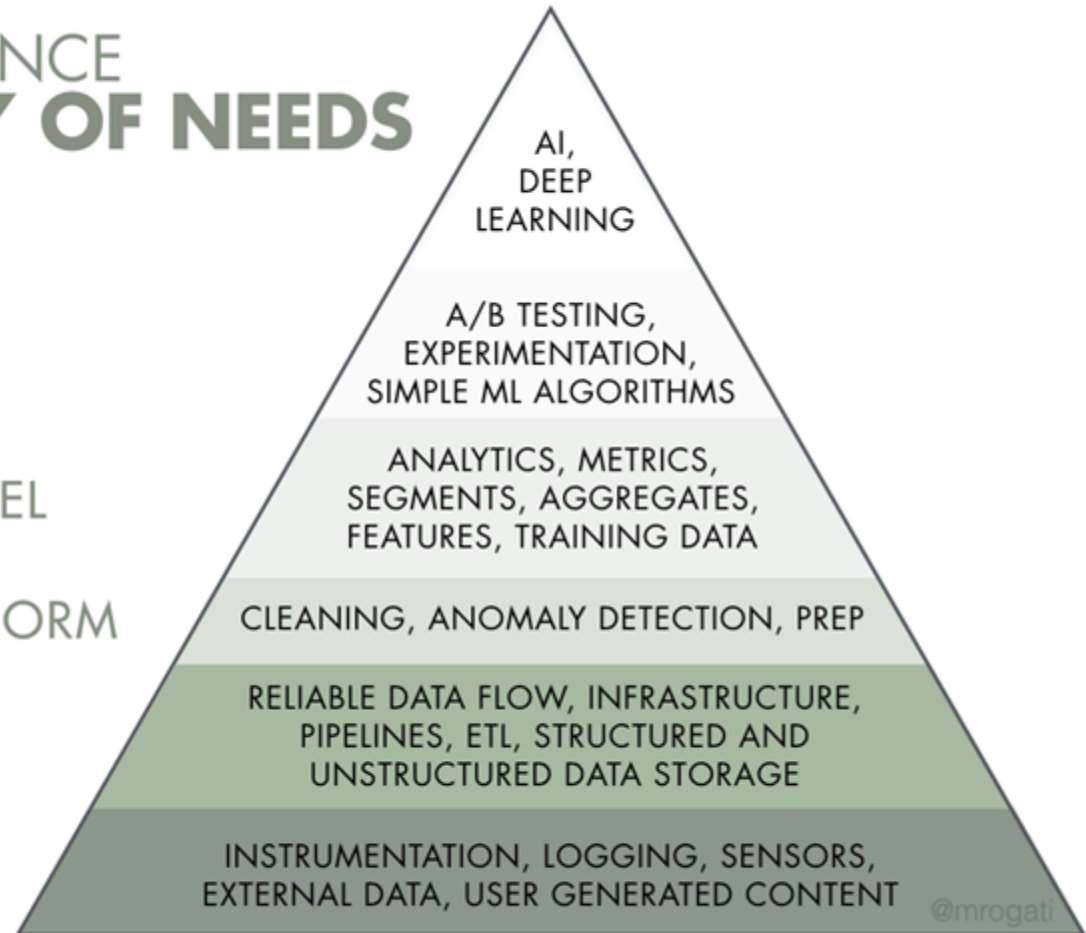
LEARN/OPTIMIZE

AGGREGATE/LABEL

EXPLORE/TRANSFORM

MOVE/STORE

COLLECT



Quelle : <https://hackernoon.com/the-ai-hierarchy-of-needs-18f111fcc007>

## Introduction and Motivation

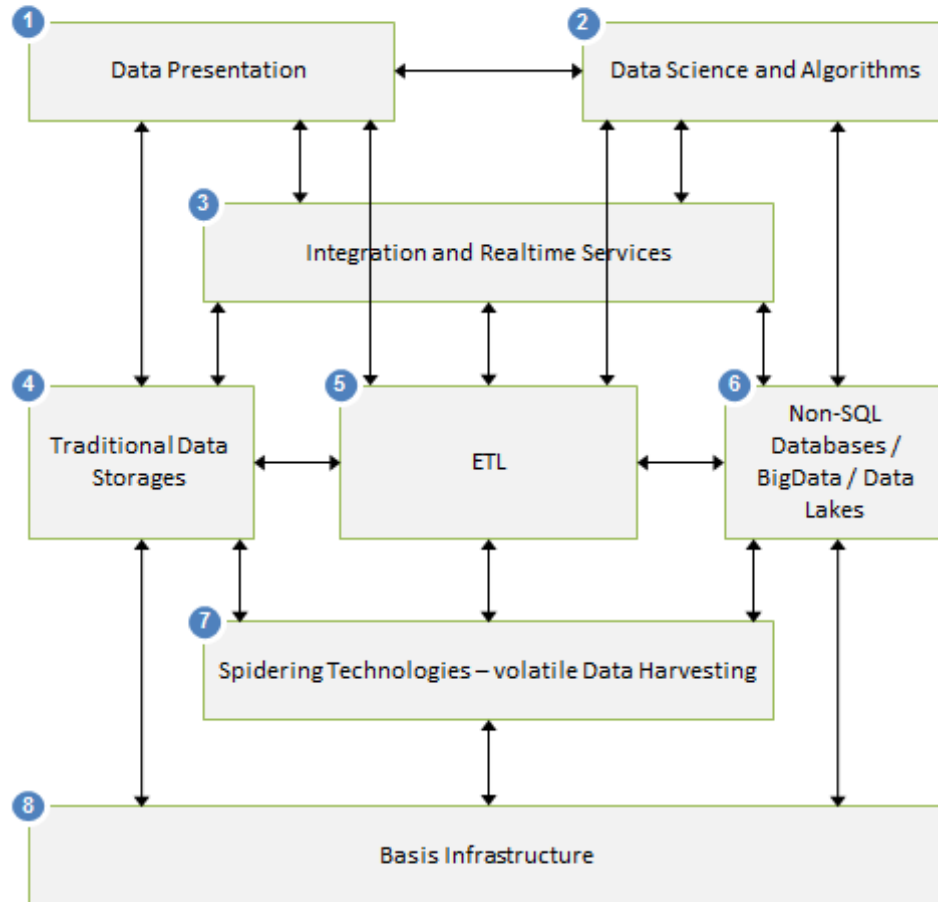
*Volume + Velocity + Variety + Veracity = Value*



Quelle : <https://medium.com/dev-genius/what-you-need-to-know-before-you-become-a-data-engineer-career-advice-503b95e7a3cf>

# Introduction and Motivation

## Main building blocks in the DI technology landscape

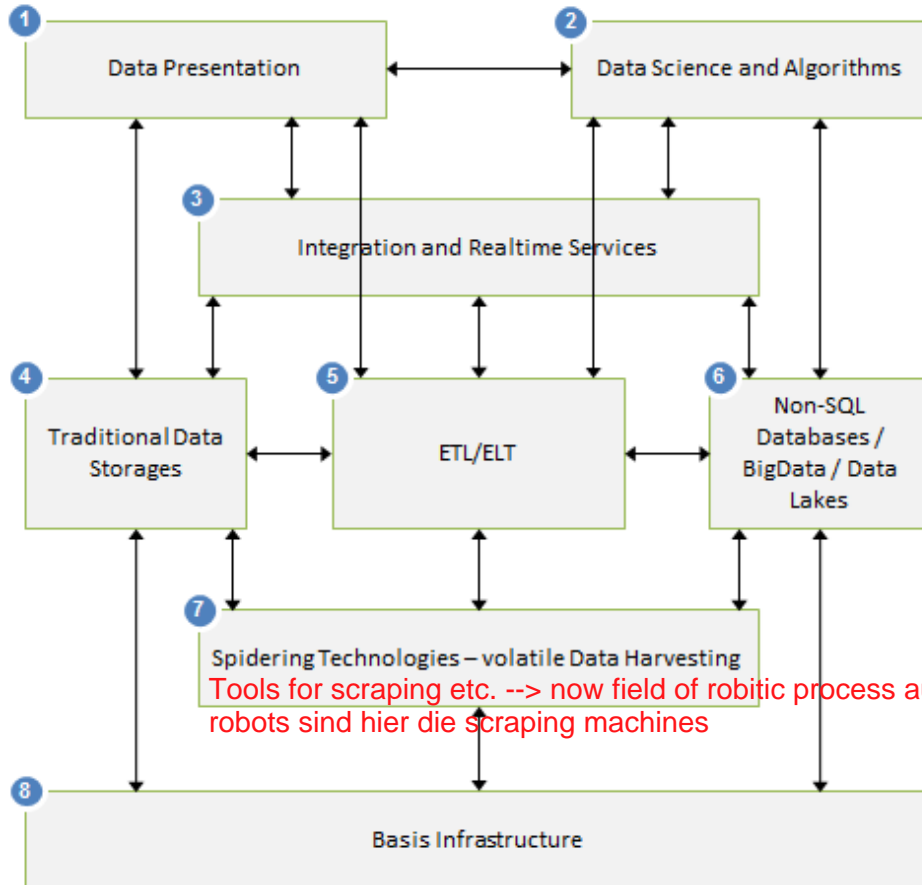


## Building blocks

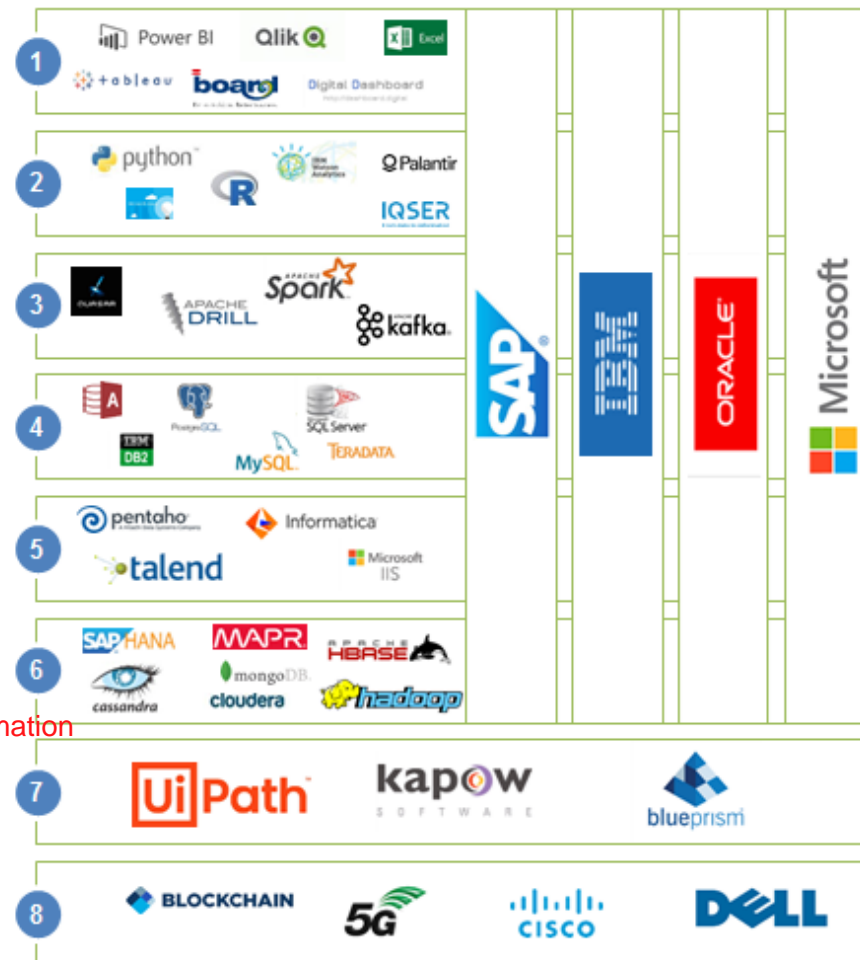
- 1 Structured and unstructured presentation of data such as Reporting, Cockpits, Dashboards and other data visualization.
- 2 Deep learning, artificial intelligence, neuronal networks, machine learning, algorithms (e.g. regression, clustering, decision trees) and bot's.
- 3 Entire underlying technology in order to allow for real-time enterprise data integration.
- 4 Traditional relational and multidimensional data storages such as SQL, DWH and BI technologies including in-memory approaches.
- 5 Data extraction, transformation and loading including robots for data provisioning automation approaches and bot's
- 6 Non-SQL data storages such as big data technologies, MapReduce, Hadoop including in-memory approaches.
- 7 Web Crawling techniques to get unstructured data from different sources such as ontology techniques, interfaces to web-services as well as related bot technologies.
- 8 Entire underlying infrastructure such as hardware, network, OS-layer, orchestration as well as new technologies such as blockchain.

# Introduction and Motivation

## Main building blocks in the DE technology landscape



## Products (excerpt)





## BIG DATA &amp; AI LANDSCAPE 2018





## Introduction and Motivation: How does Data Engineering look like in many companies ?

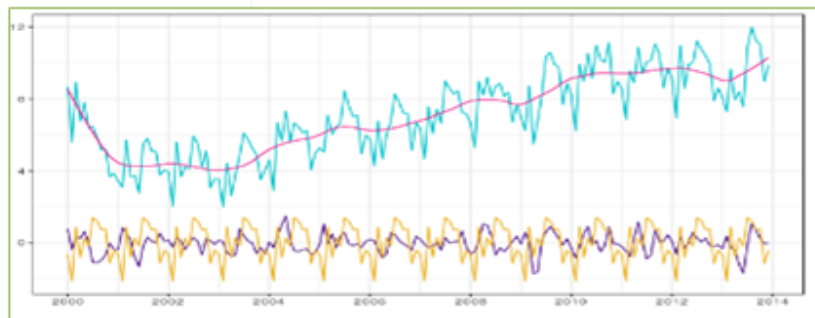


# Introduction and Motivation: How does Data Engineering look like in many companies ?

## Use Excel as Data Scientist only as a presentation-tool

		MONATE												letzte 4 Wochen			
Produkt	Jahr	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Okt	Nov	Dec	9	10	11	12
Besuche aller Webseiten	2014	702041	558186	576874	592	0	0	0	0	0	0	0	0	119092	128916	150799	12
	2013	705818	537505	574537	527663	705955	835130	833372	799296	681926	529972	371207	393298	130256	130387	133314	12
	Verhältnis	99.46%	103.85%	100.41%	0.11%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	91.43%	98.87%	113.12%
Umsatz Pauschalreisen	2014	3706499	3549532	3891273	0	0	0	0	0	0	0	0	0	797622	814920	882998	92
	2013	2838527	2523908	3239348	3142750	3178890	4662563	5571988	4839891	4259929	2595905	1280367	1592120	800863	752441	711269	68
	Verhältnis	130.58%	140.64%	120.13%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	99.60%	108.30%	124.14%
Buchungen Pauschalreisen	2014	1317	1286	1469	0	0	0	0	0	0	0	0	0	312	303	329	
	2013	1014	980	1269	1313	1364	1803	2164	2068	1722	1023	483	507	299	300	261	
	Verhältnis	129.88%	131.22%	115.76%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	104.35%	101.00%	126.05%
Umsatz Linienflüge	2014	231819	151634	200121	0	0	0	0	0	0	0	0	0	43551	40898	50914	5
	2013	254603	131606	162818	210350	243237	183768	214069	199159	215188	188514	137561	113704	47333	39961	31733	2
	Verhältnis	91.05%	115.22%	122.91%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	92.01%	102.34%	160.44%
Buchungen Linienflüge	2014	160	124	166	0	0	0	0	0	0	0	0	0	38	41	33	
	2013	175	104	131	163	201	155	200	151	172	156	96	85	37	28	30	
	Verhältnis	91.43%	119.23%	126.72%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	102.70%	146.43%	110.00%

### Result Report Switzerland



Rolling Forecast			Test	Customers						Pot. Prosp. PB
Curr. SL	Next SL	Quotes	Re	Bu	Active Re	%	Bu	New Re	%	
69'002	249'538	207'500	8	17	30	176.47	9	1	11.11	
100'558	179'475	535'400	38	29	27	93.1	12	7	58.33	
153'775	170'492	315'000	41	49	45	91.84	12	11	91.67	
104'890	187'684	213'750	33	29	28	96.55	9	12	133.33	
94'287	36'775	393'500	20	35	18	51.43	8	6	75	
236'220	311'782	1'071'745	79	68	69	101.47	34	28	82.35	
758'731.64	1'135'746.75	2'736'894.97	219.00	227.00	217.00	95.59	84.00	65.00	77.38	

Total reported consultants

Total consultants' factors for av.

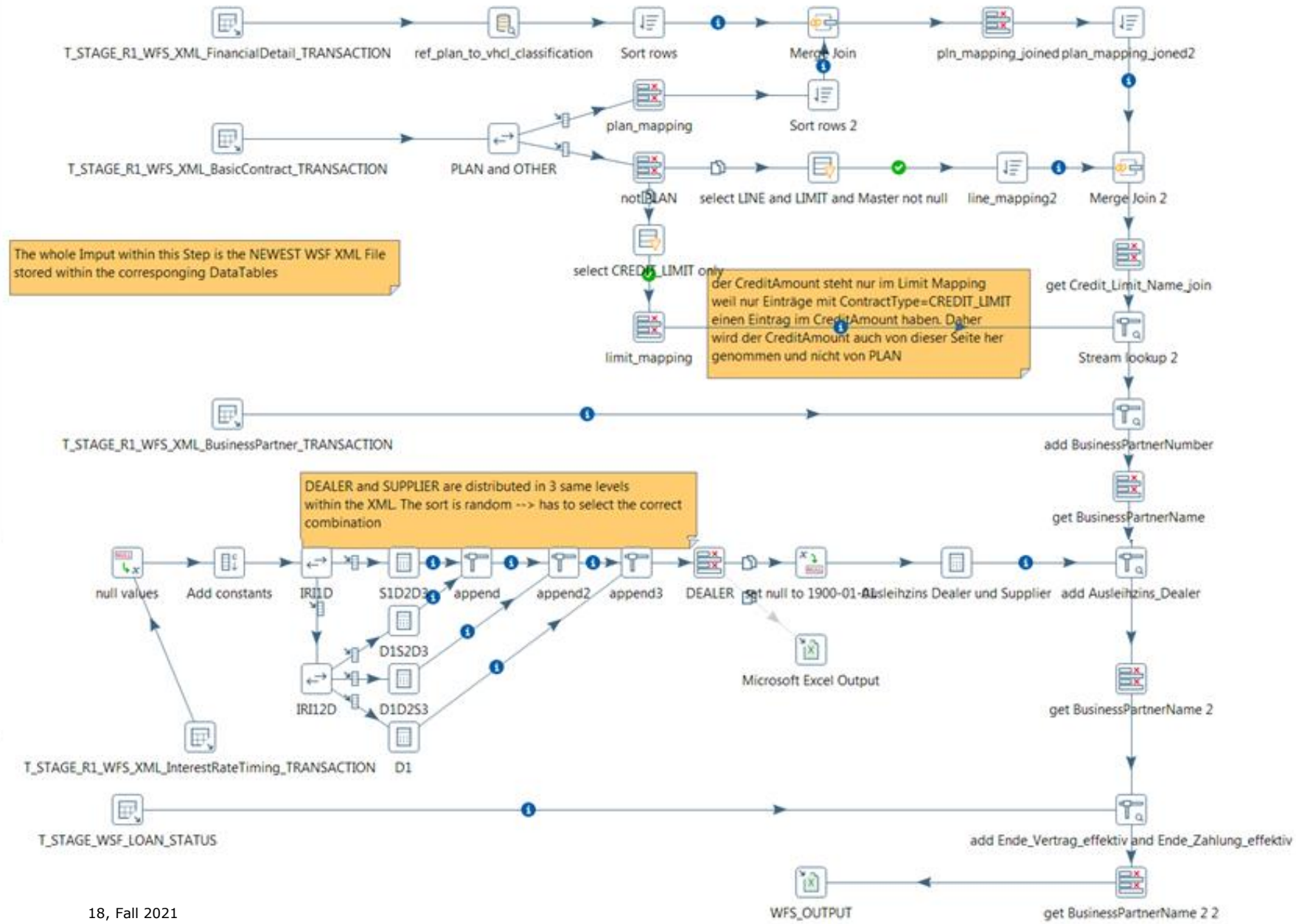
16, Fall 2021 budget on period 92.3076923



## **Introduction and Motivation:** Motivation for the use of ETL tools. Data Engineering and Business Process Automation - similarities/differences

I choose a lazy person to do a hard job. Because a lazy person will find an easy way to do it.  
-Bill Gates

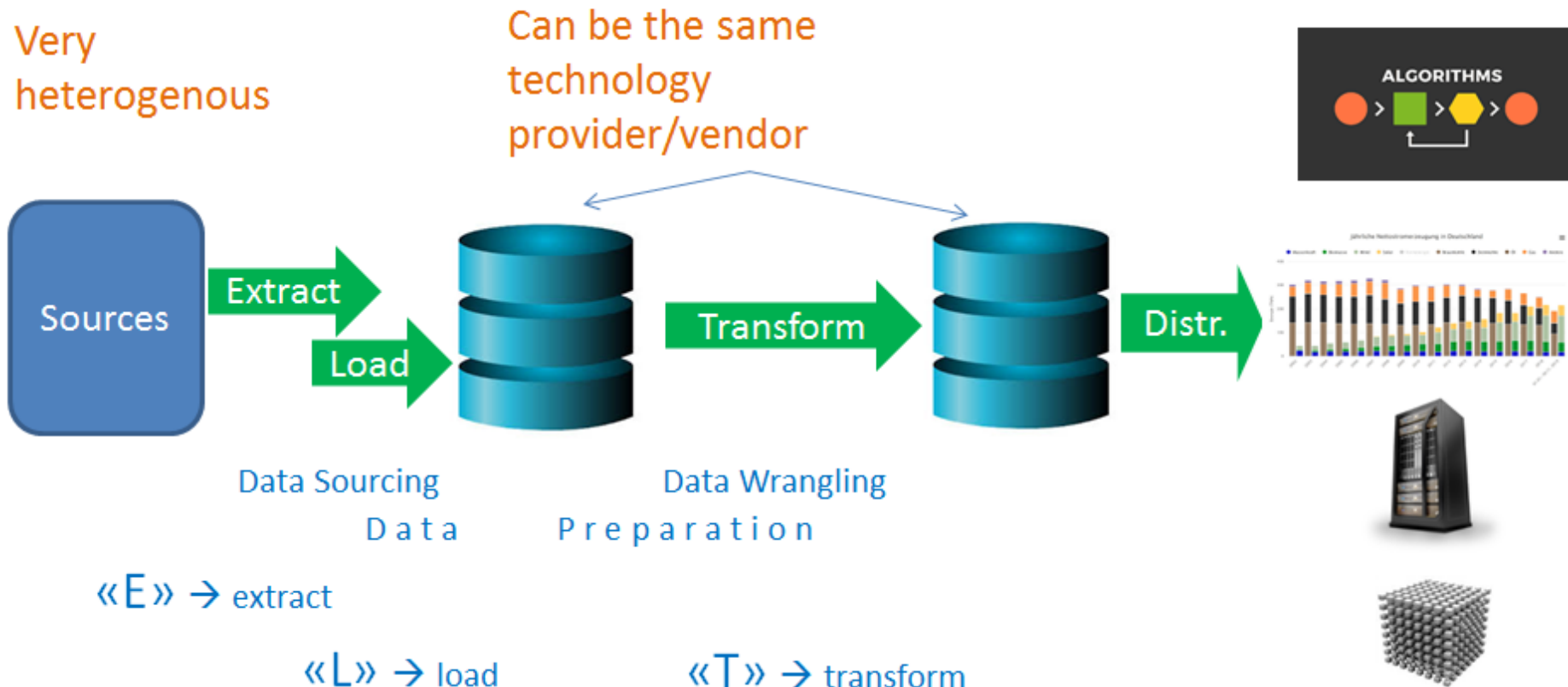
- Velocity and agility
- Cost of development → predefined functions
- «ease of use» without programming knowledge
- Independen from programming language → more people can participate in the process
- Overview of process → Directed acyclic graph
- Maintenance cost
- Documentation
- 1000 rows JAVA Code in 1 workflow (Ex. WFS\_IMPORT\_ROUTINE\_FOR\_DWH\_STORE\_PREPARATION



# Introduction and Motivation

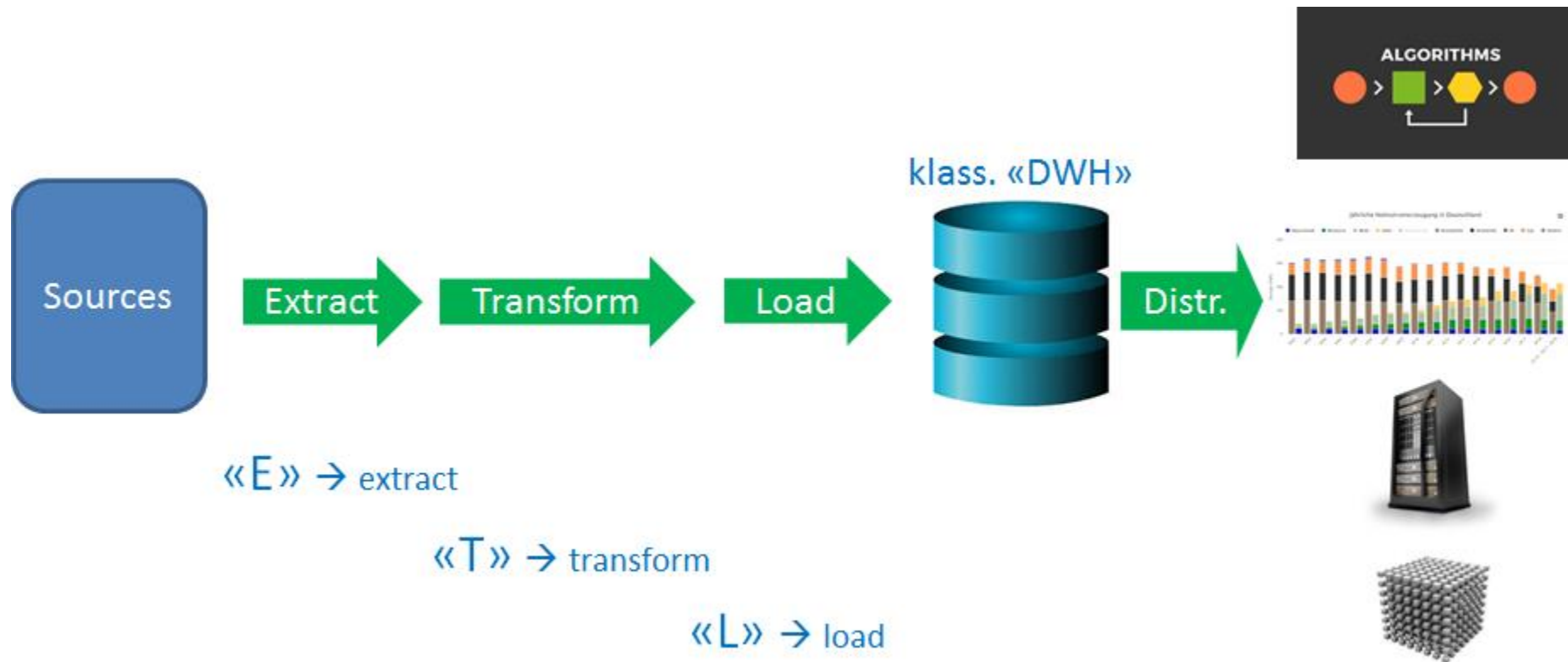
## Elements of a Data Pipeline:

An ETL/ELT-Tool deals with all Tasks

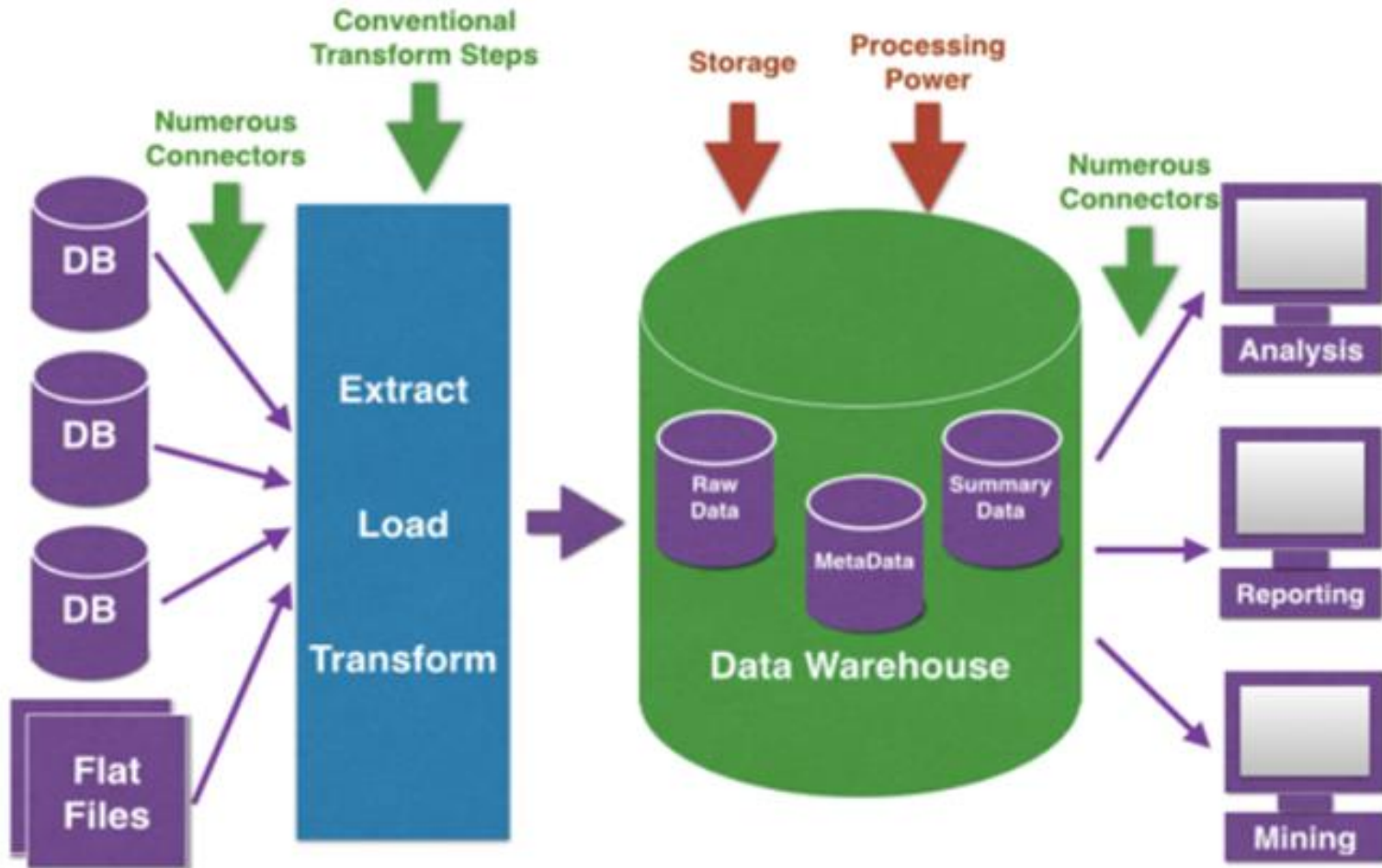


ELT is better than ETL if you have unreliable data --> everything is extracted and loaded and then the process is under control of you and you can check if there for example is a data type problem or so

# Introduction and Motivation



## Introduction and Motivation



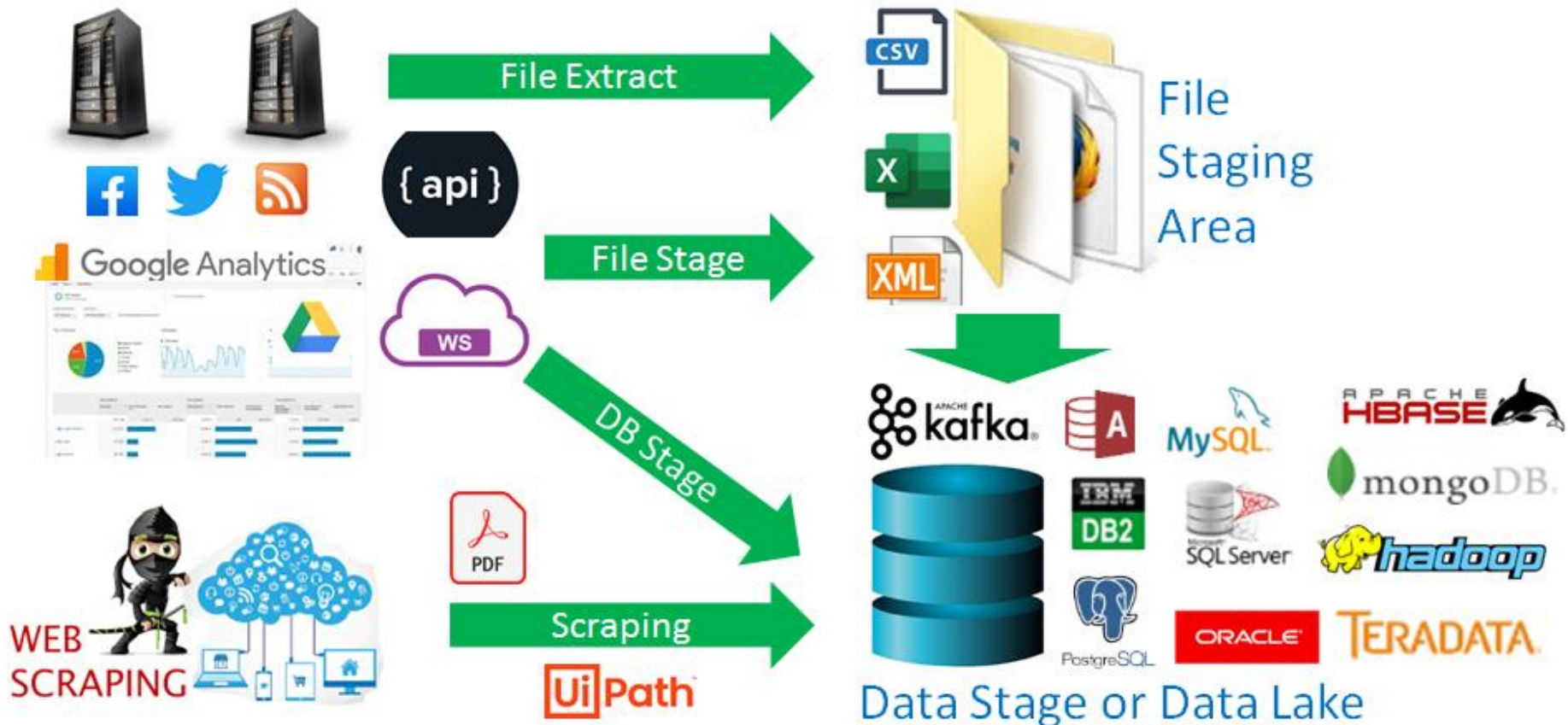
*Figure 1 - Traditional Data Integration*



# Introduction and Motivation

## Elements of a Data Pipeline: Part 1

## Data Sourcing



## Introduction and Motivation

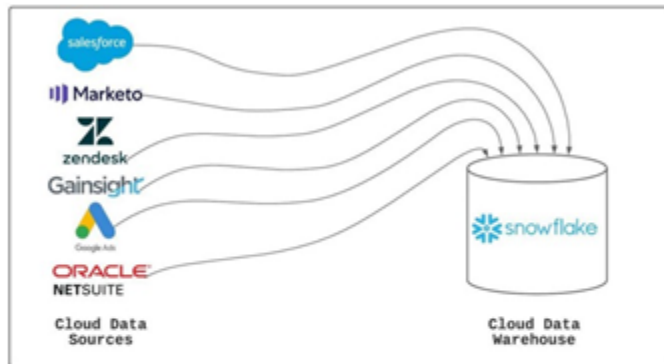
Don't underestimate the managerial and psychological tasks in the Data Sourcing part:

- Information about the elements of KPI's
- Which System holds which data ?
- Who can tell me details in case of missing Data Dictionary ?
- Whow can i access them ?
- Do i really get what i effectively want ? -> testing !

→ don't underestimate the amout of time you use here !

# Introduction and Motivation

## ETL



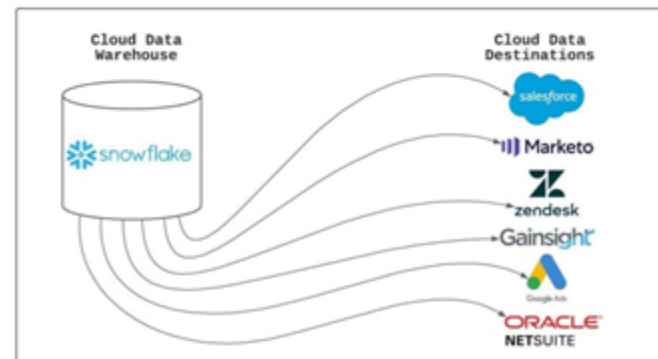
### Definition:

process of copying data from the data warehouse into systems of record across a company.

There are 3 primary use cases for Reverse ETL:

- 1. Operational analytics** — feeding insights from analytics to business teams in their usual workflow so they can make more data-informed decisions
- 2. Data automation** — not all data problems are so glamorous. “Can I get a CSV to issue some invoices?”, your finance team asks. Reverse ETL poses a simple solution.
- 3. Data infrastructure** — with a growing number of data sources, reverse ETL is emerging as a general-purpose pattern in software engineering.

## Reverse ETL



Source: <https://tejasmanohar.medium.com/what-is-reverse-etl-6f228a14f6ec>



## Introduction and Motivation

Elements of a Data Pipeline : Part 2 und 3

ETL/ELT und Ziel



# Introduction and Motivation



Figure 2 - Big Data Integration

# Introduction and Motivation

## Robotic Process Automation:

--> enhance ETL tool functionality



## Introduction and Motivation

DEMO: VIDEO



Example of Web Scraping

[https://www.youtube.com/watch?v=m\\_PkWeHSOrE](https://www.youtube.com/watch?v=m_PkWeHSOrE)