

# Risikobasiertes Testen von Daten

Ing. Harald Foidl, PhD

Das ATB Fluchtachterl

Raus aus dem Alltag,  
rein ins Gespräch

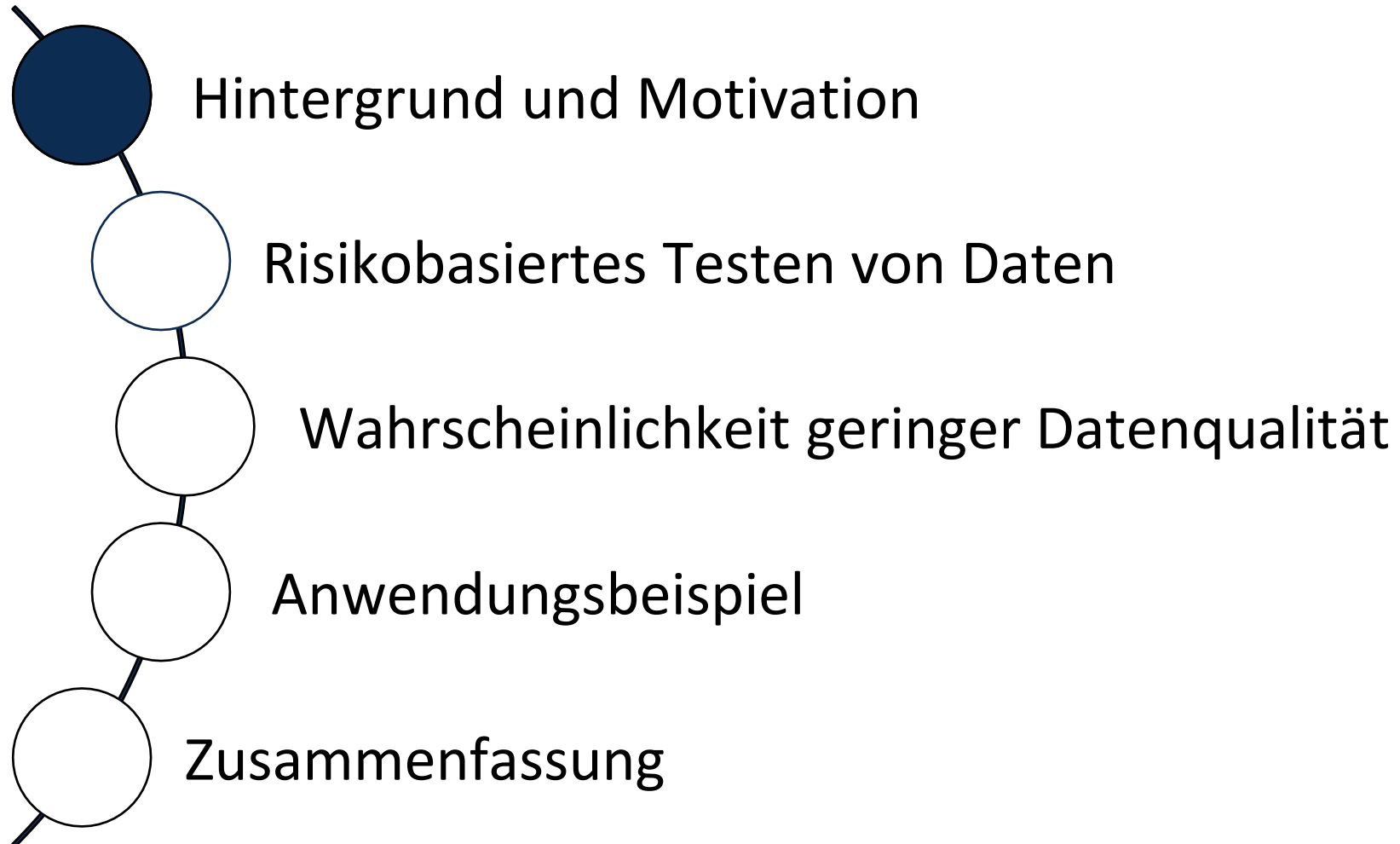
ATB Fluchtachterl

18.06.2025



# Agenda

---







# Zuverlässigkeit von Datenintensiven Systemen

Fehlentscheidungen solcher Systeme können schwerwiegende Folgen haben!

## How a Self-Driving Uber Killed a Pedestrian in Arizona

By TROY GRIGGS and DAISUKE WAKABAYASHI UPDATED MARCH 21, 2018



## Report: IBM Watson delivered 'unsafe and inaccurate' cancer recommendations

JULY 25, 2018 BY FINK DENSFORD — LEAVE A COMMENT



➔ **Datenqualität entscheidend**

# Lösung: Testen von Daten (Datenvalidierung)

- Kontinuierliche Prüfung/Überwachung von Daten zur Qualitätsbewertung
- Definierte Schemata, Wertebereiche und weitere Kriterien (z. B. Eindeutigkeit, Datentypen)
- Techniken, z.B.:



Hypothesentests



Korrelationsanalysen



Aggregationen



Lageparameter (z.B. Mittelwert)

# Jedoch ...

---



Datenvalidierung beeinflusst die Systemperformance



Manuell, zeitaufwendig, erfordert Expertenwissen



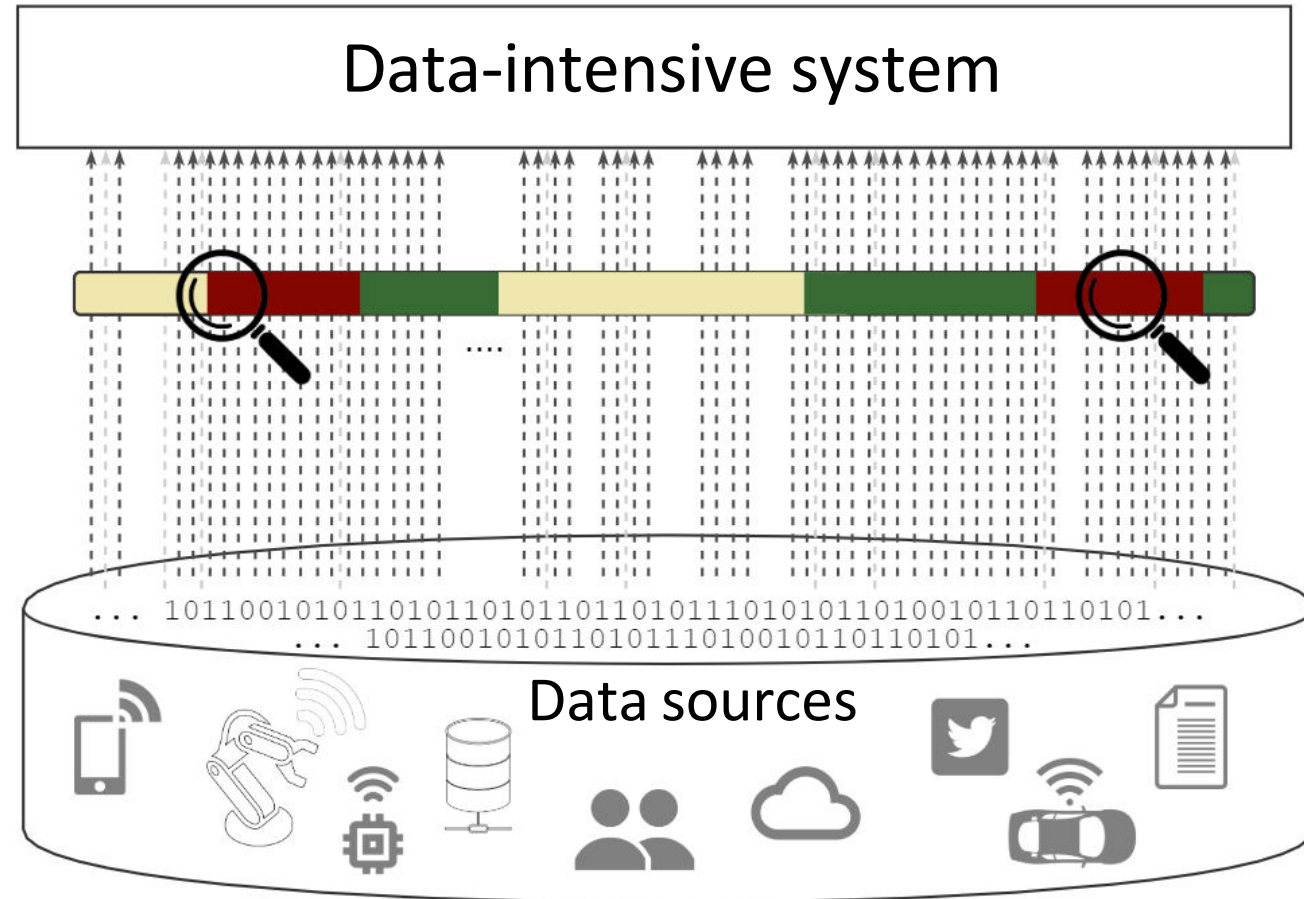
Problem der Validierungstiefe/Strenge der Validierung



Vollständige Datenvalidierung ist nicht möglich

# Problemstellung

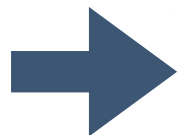
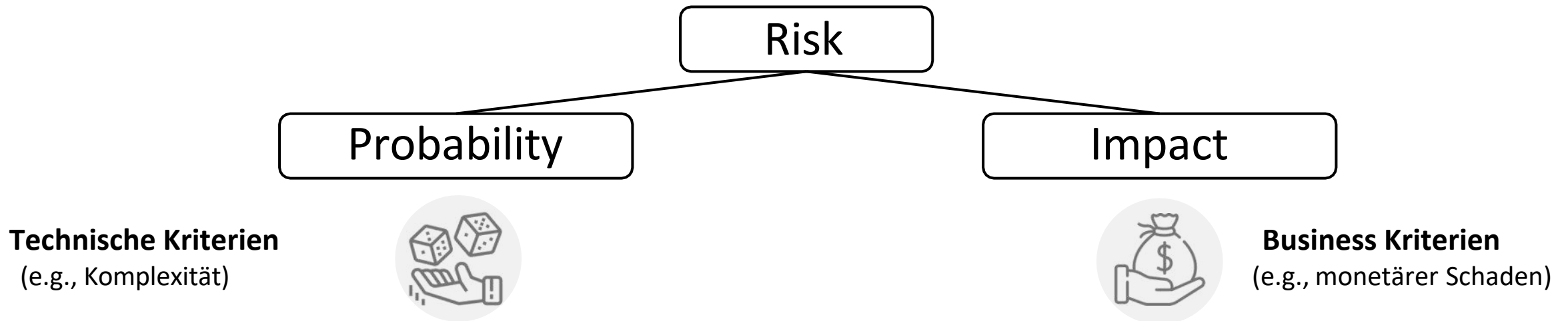
SELEKTIVES  
Testen von Daten



*Welche* Daten sollen validiert werden – und wie *gründlich*?

# Ansatz

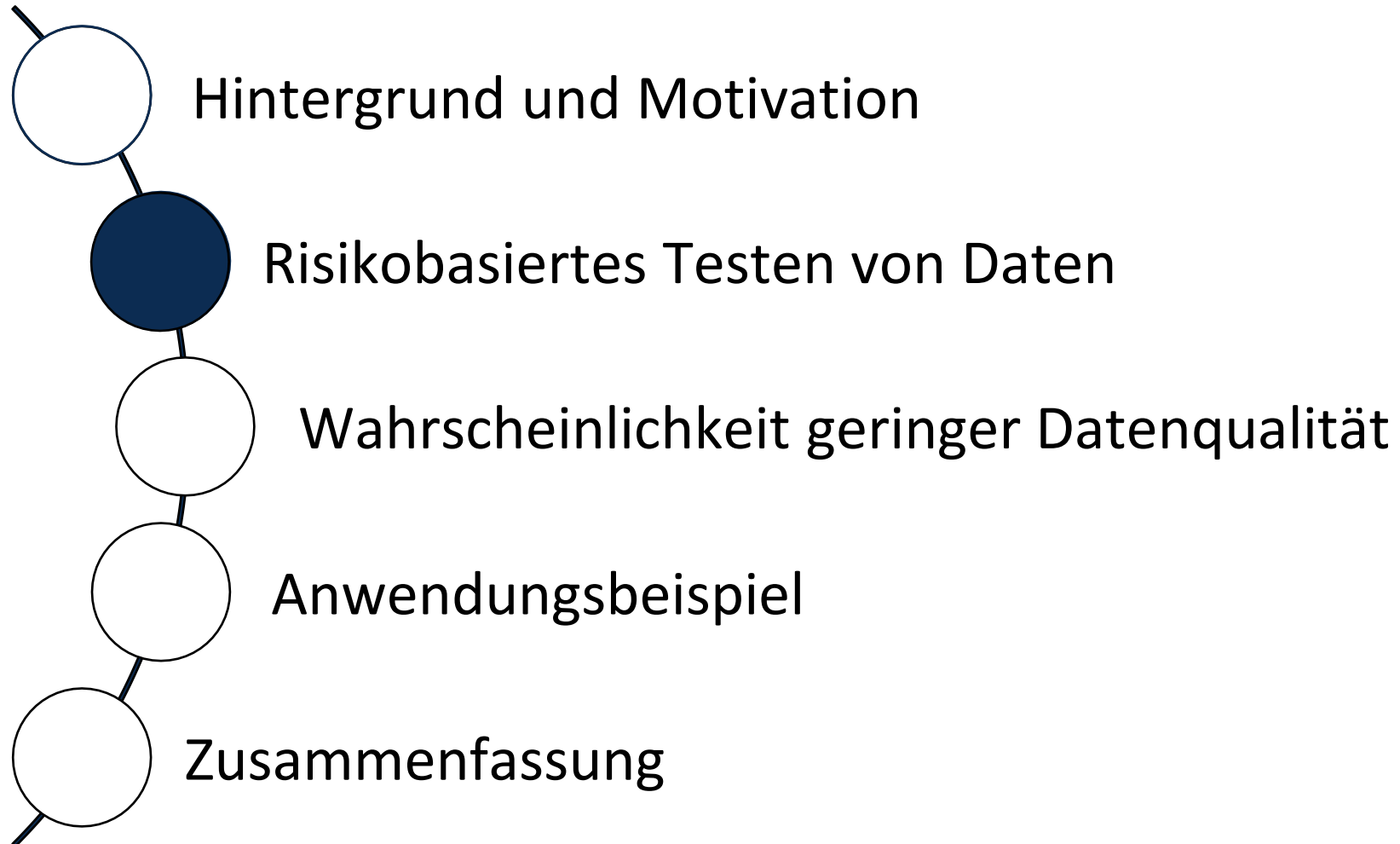
- Softwaretesten: risikobasierte Strategie für vergleichbare Entscheidungsprobleme
- Risikobasiertes Testen: Risiken eines Softwareprodukts steuern den Testprozess



**Nutzung des Konzepts des risikobasierten Softwaretestens**

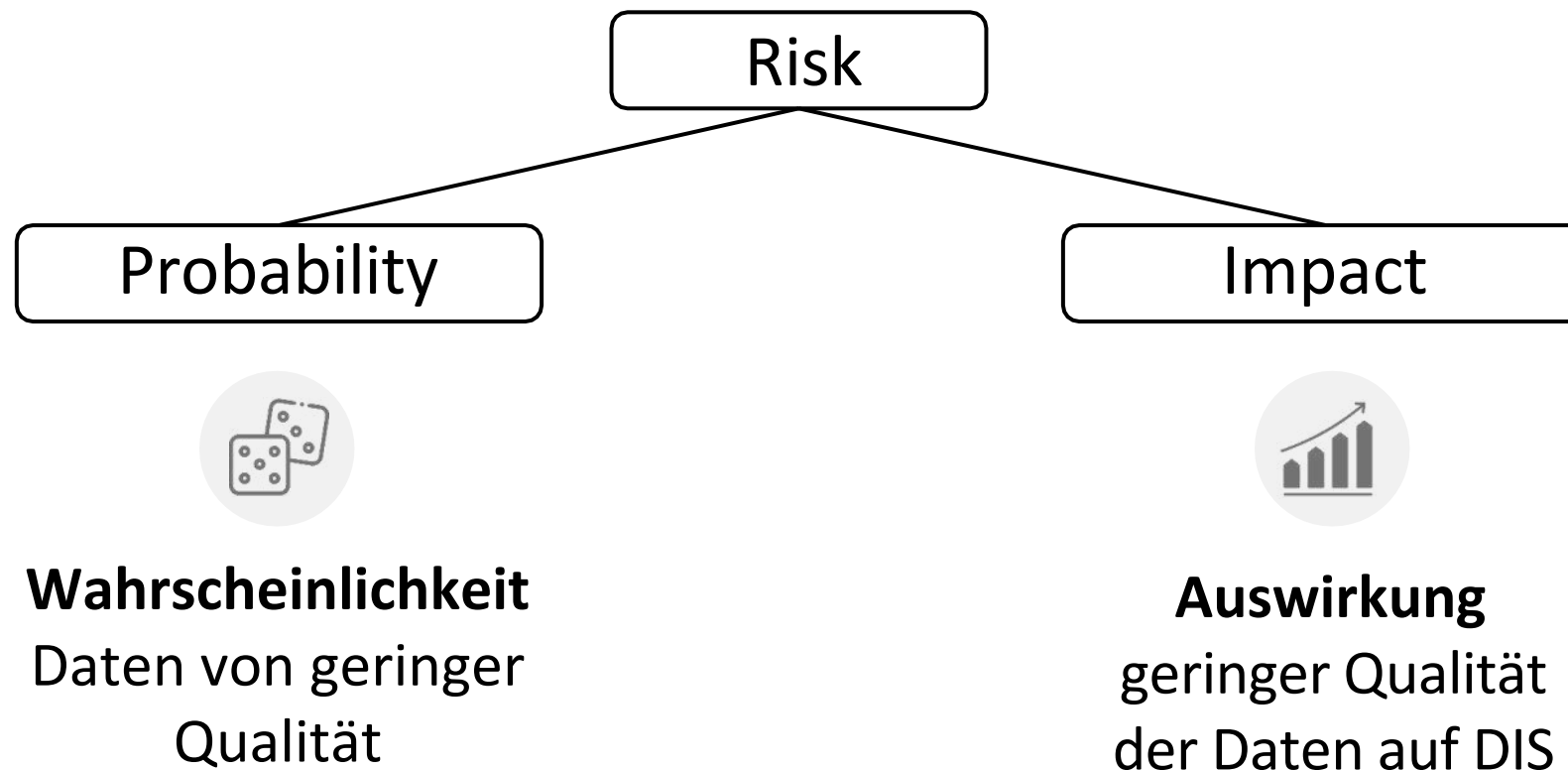
# Agenda

---

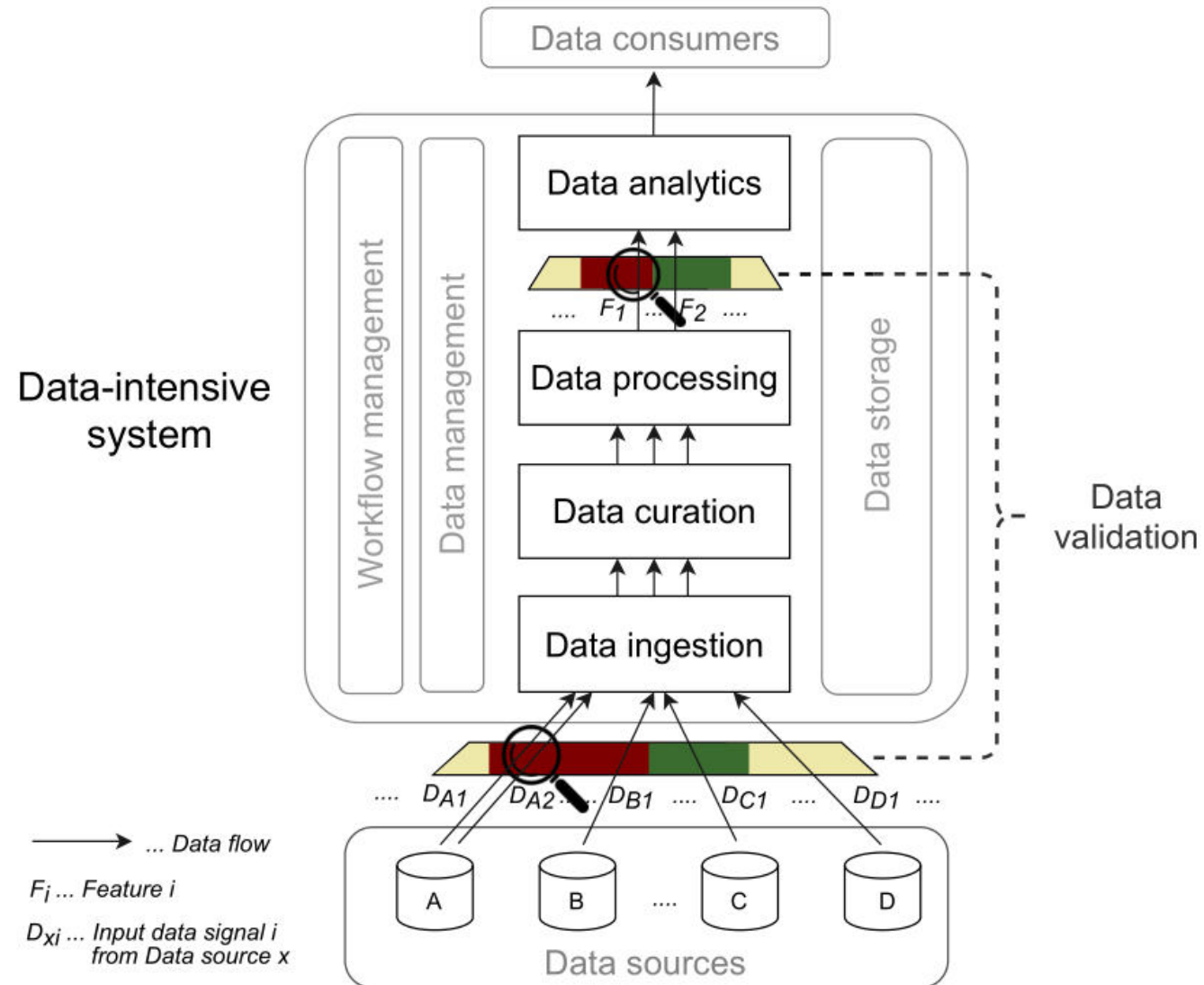


# Idee des risikobasierten Testens von Daten

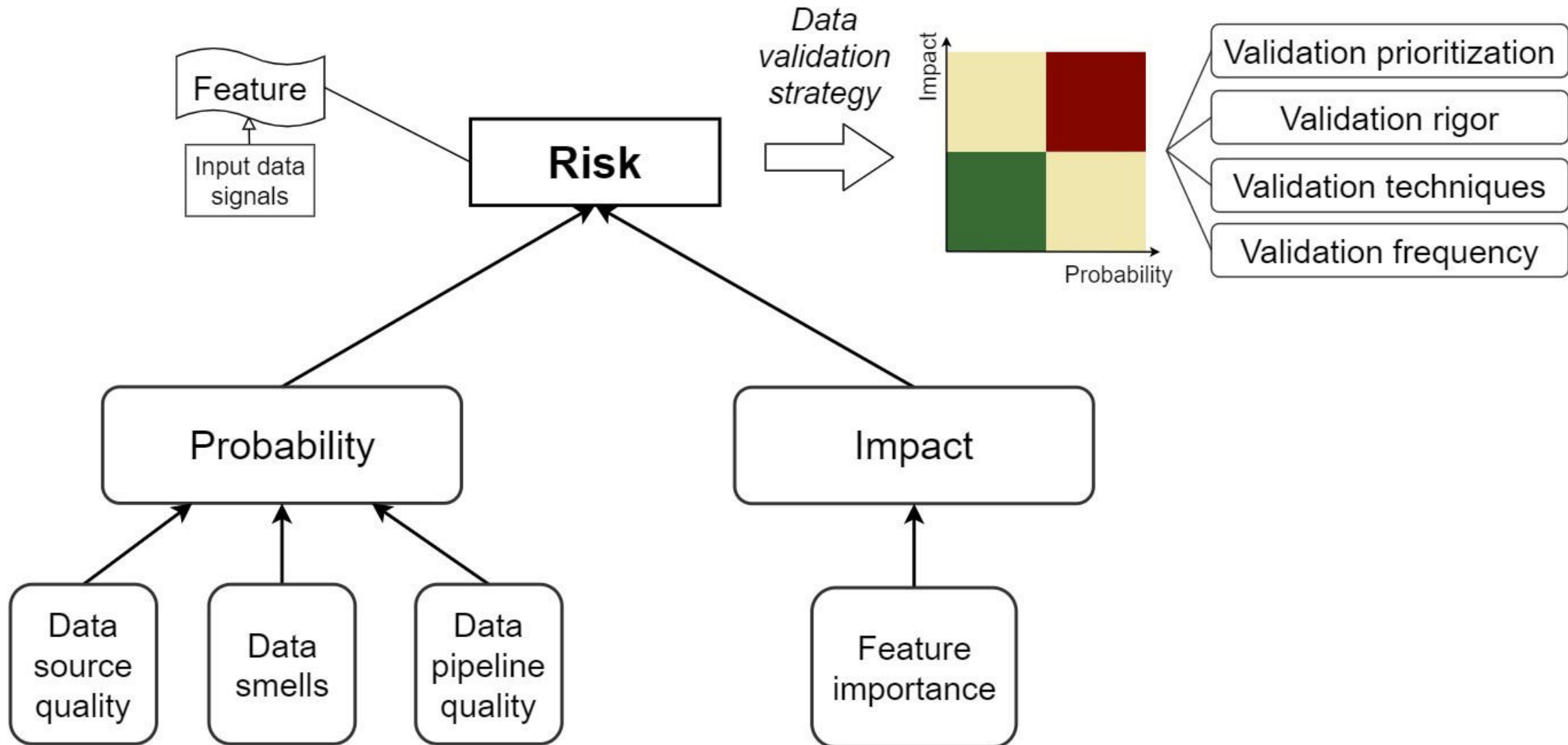
Bestimme **Risiko von geringer Datenqualität**,  
um Testen zu steuern



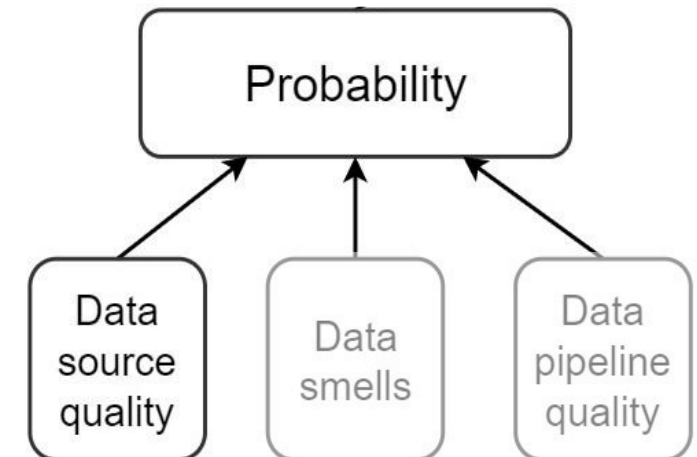
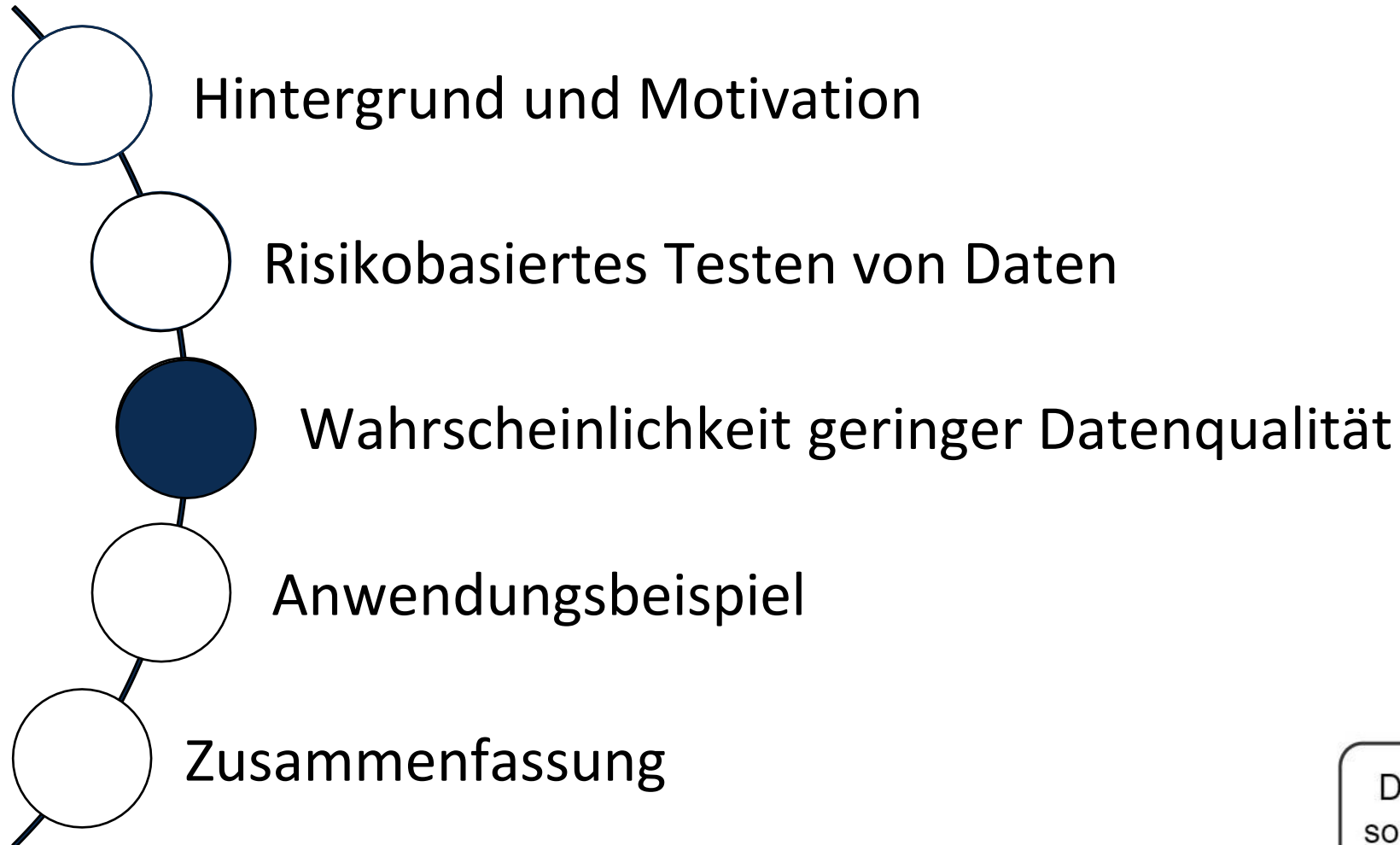
# Risikobasiertes Testen von Daten



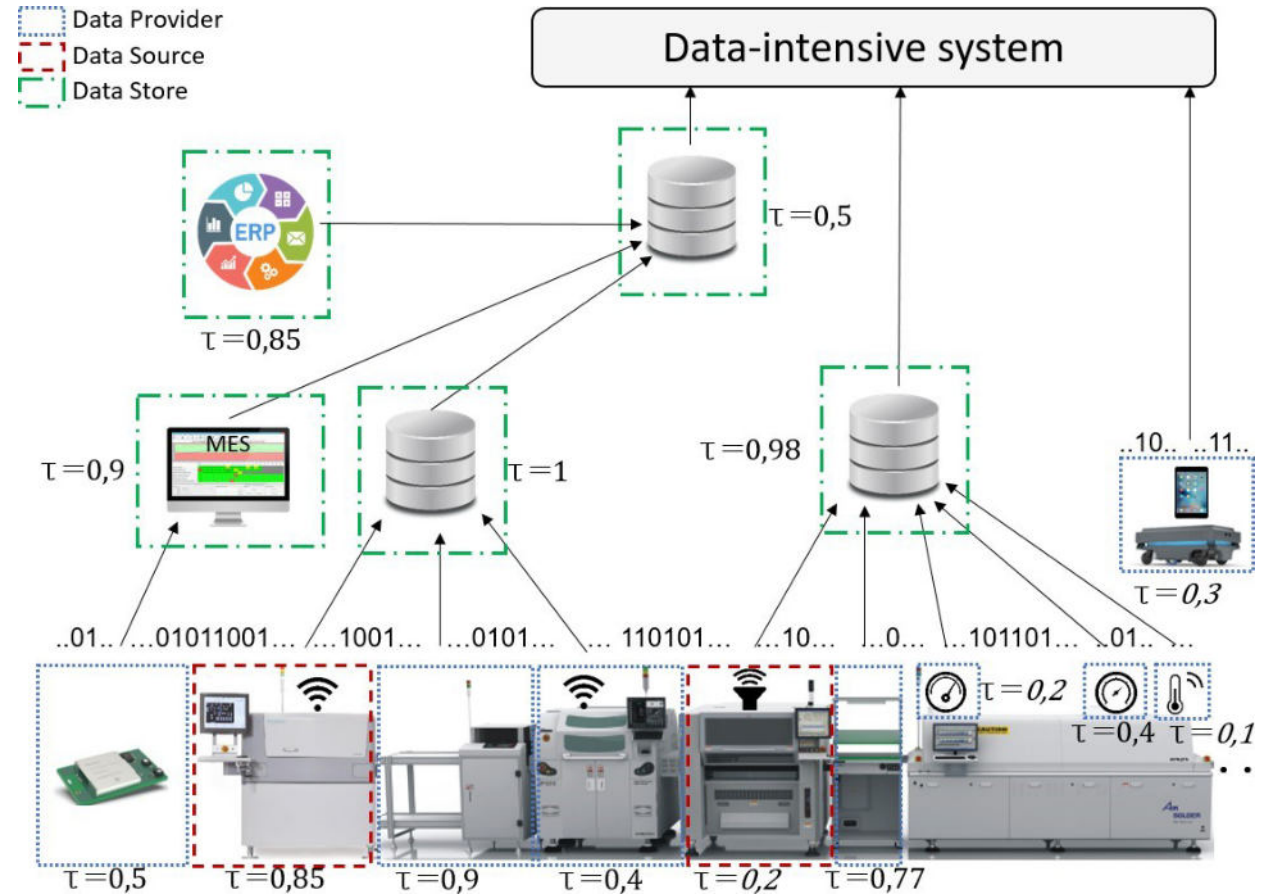
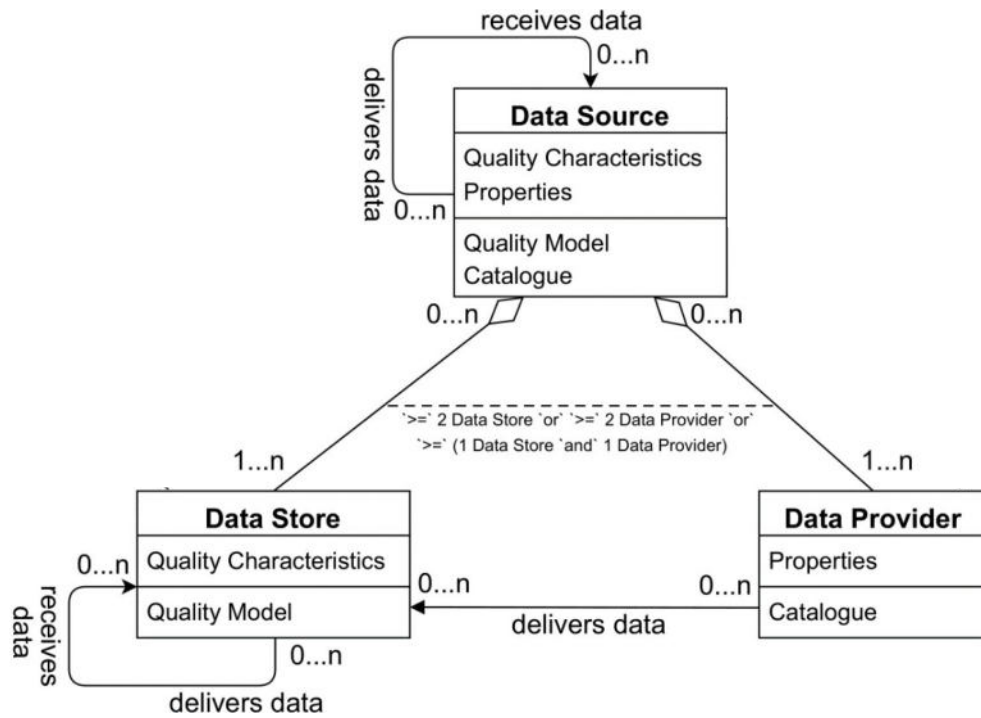
# Risikobasiertes Testen von Daten - Framework



# Agenda



# Qualität von Datenquellen



# Qualitätsmodell für Datenspeicher

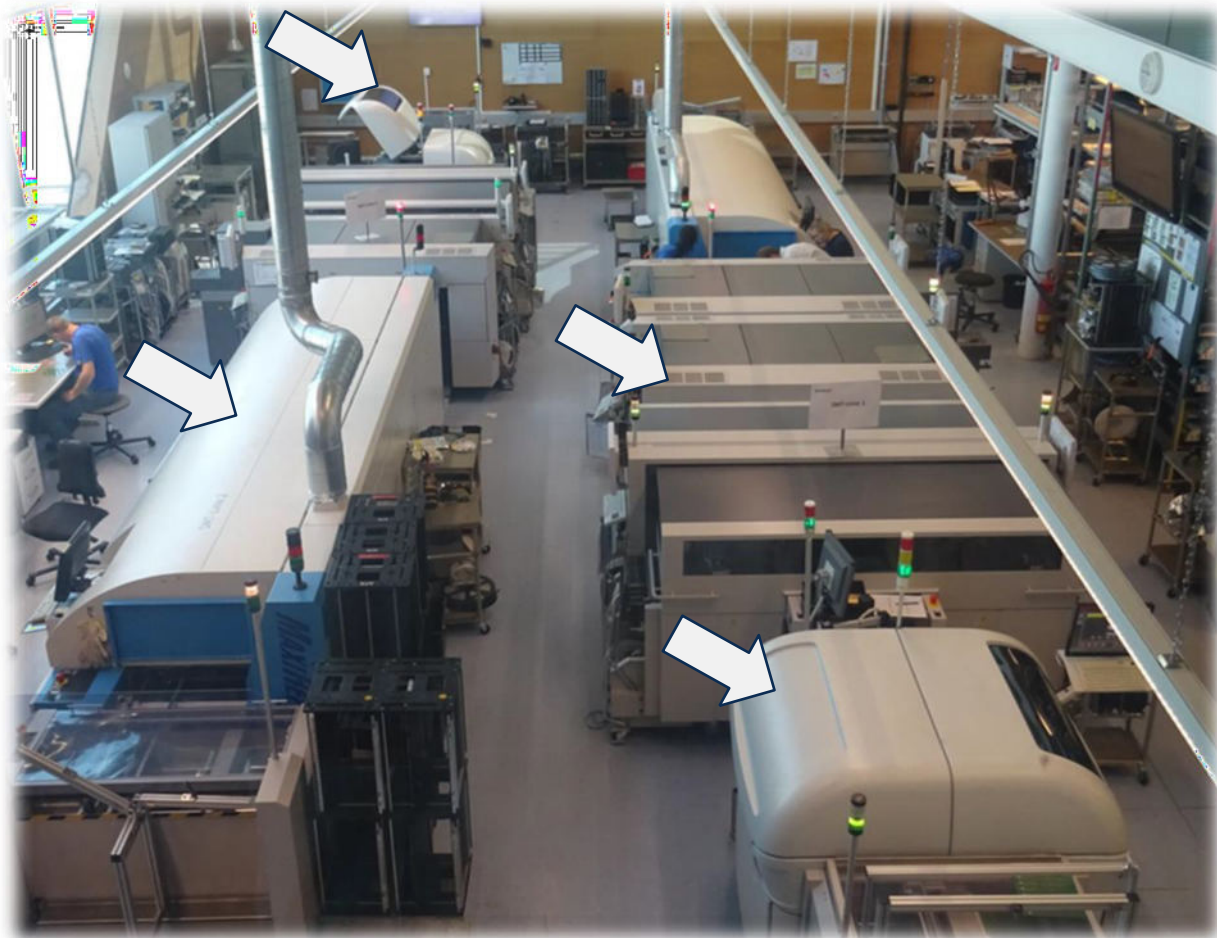
Data Quality Category	Data Quality Characteristic	Description	Properties and Sub-characteristics
Representational Data Store Quality	<b>Representational Adequacy</b>	Degree to which a data store presents data in a <i>concise and organised way</i> .	Schema Minimality Schema Normalisation Schema Pertinence
	<b>Representational Consistency</b>	Degree to which a data store presents data <i>always in the same format and compatible with previous data</i> .	Data Format Variety Data Type Variety Schema Change Proneness
	<b>Understandability</b>	Degree to which <i>users can understand the data</i> provided by a data store.	Data Format Complexity Documentation Degree Metadata Quality Schema Readability
Dynamical Data Store Quality	<b>Accessibility</b>	Degree to which data are <i>easily and quickly retrievable</i> .	Access Maturity Operability Retrievability
	<b>Availability</b>	Degree to which <i>data are available</i> from a data store.	Durability Fault Tolerance Recoverability Scalability Uptime
	<b>Security</b>	Degree to which <i>access to data for unauthorised persons is restricted</i> by a data store.	Authentication/Encryption Authorisation Policy
	<b>Timeliness</b>	Degree to which a data store <i>provides up-to-date data in a timely manner</i> .	Refresh Rate Response Time
Statical Data Store Quality	<b>Completeness</b>	Degree to which a data store is able to represent <i>every meaningful state of the real world</i> .	Schema Completeness Schema Correctness
	<b>Contactability</b>	Degree to which a data store provides <i>contact information for further inquiries</i> .	Support Degree Complexity Data Governance
	<b>Trustworthiness</b>	Degree to which a <i>data store can be trusted</i> .	Maturity Verifiability

# Eigenschaftskatalog von Datenanbietern

Category	Property	Option
General	Age	old
		middle
		new
	Provider type	persons involved only sensor/device/machine
Location	Environment	harsh, dynamically changing
		moderate, seldom changing
		mild, never changing
	Mobility	always
sometimes never		
Energy	Power source	battery
		hardwired
	Usage	low
		medium high
Connectivity	Mechanism	wireless
		physical
		wired

Category	Property	Option
Hardware	Quality	low
		medium
		high
	Maintenance/calibration	never
		sporadic
		regularly
	Ease of replacement/repair	hard
		medium
		easy
	Constraints	severe limitations
moderate limitations		
no limitations		
Data pre-processing		high
		medium
		low/no

# Fallstudie - Elektronikfertigung



SMT assembly machine



wearable barcode scanner

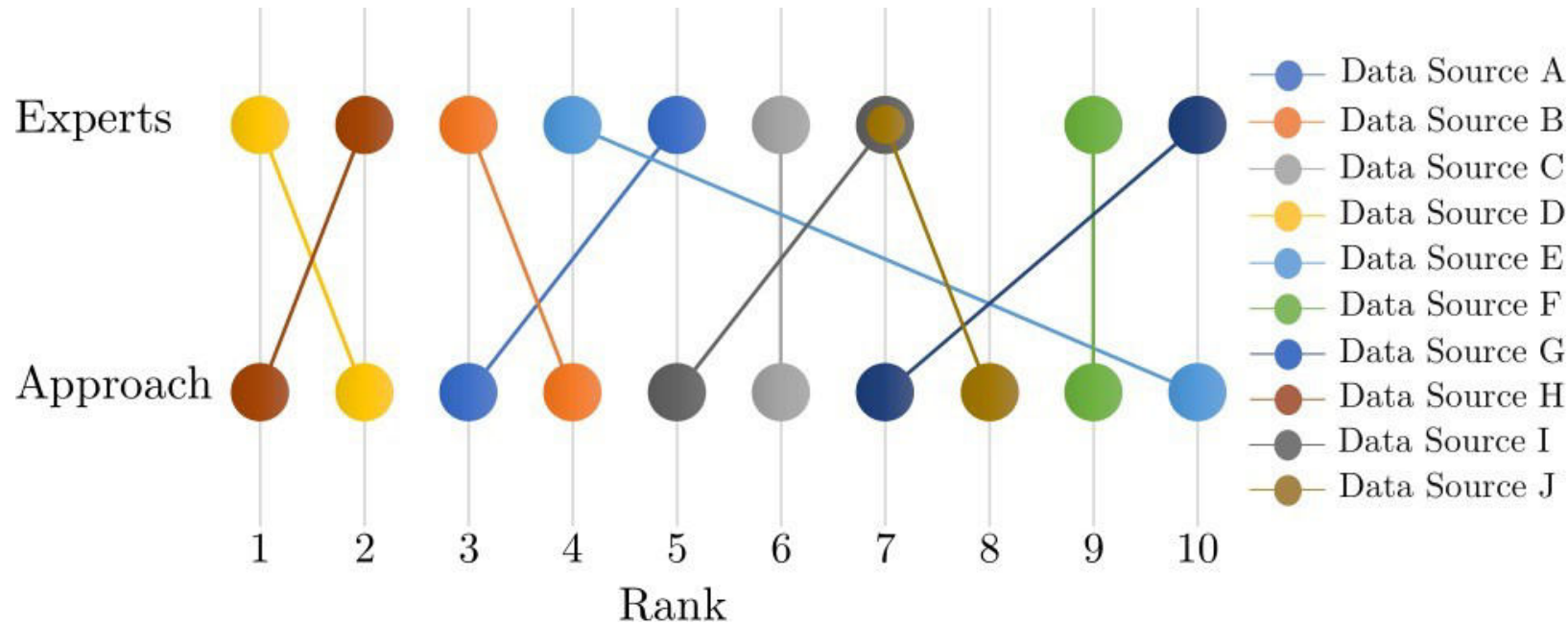


nitrogen sensor

**kontron**

The Power of IoT

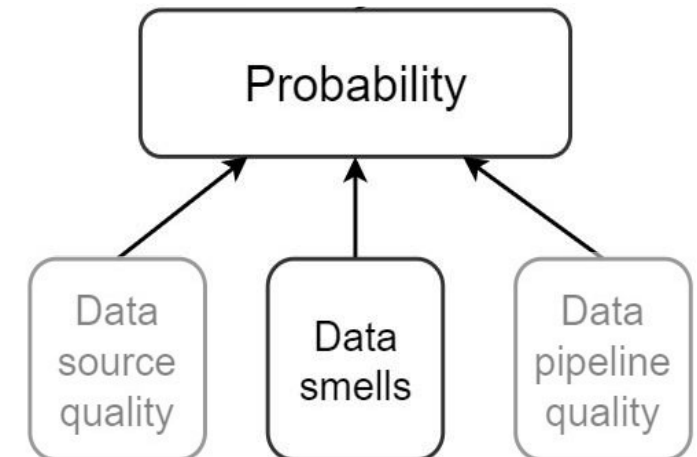
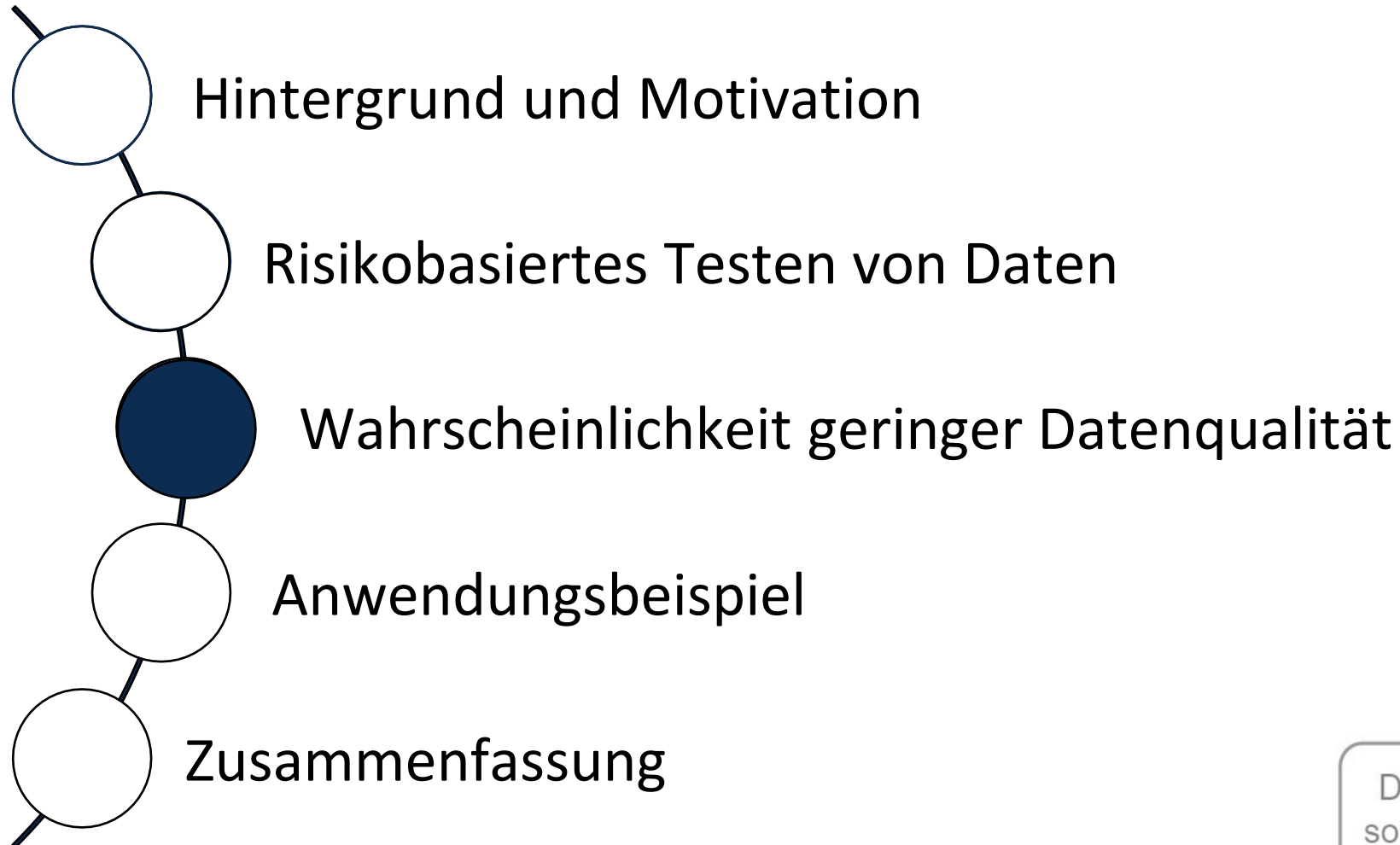
# Fallstudie - Ergebnis



*Moderate positive Korrelation zwischen Experten Bewertung  
und entwickelten Ansatz.  
(Ranking von Datenquellen)*

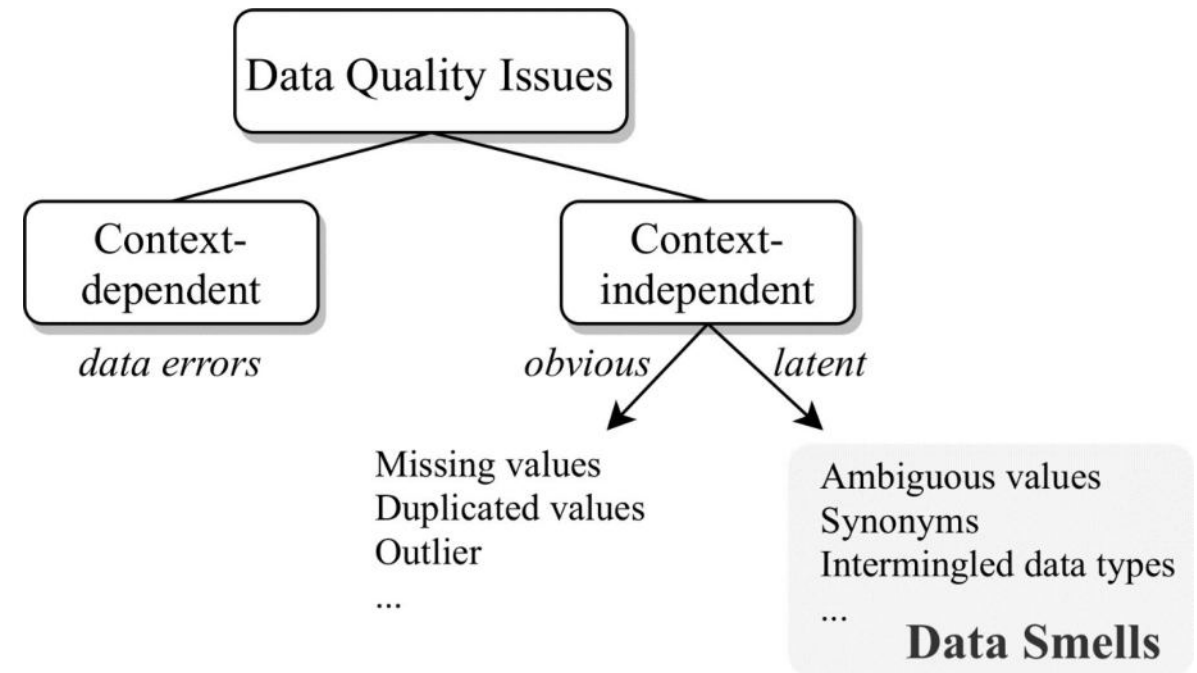
Spearman's Rho: 0.64 (p-value = 0.04)  
Inter-rater Agreement of Experts: 0.79 (Krippendorff's Alpha)

# Agenda

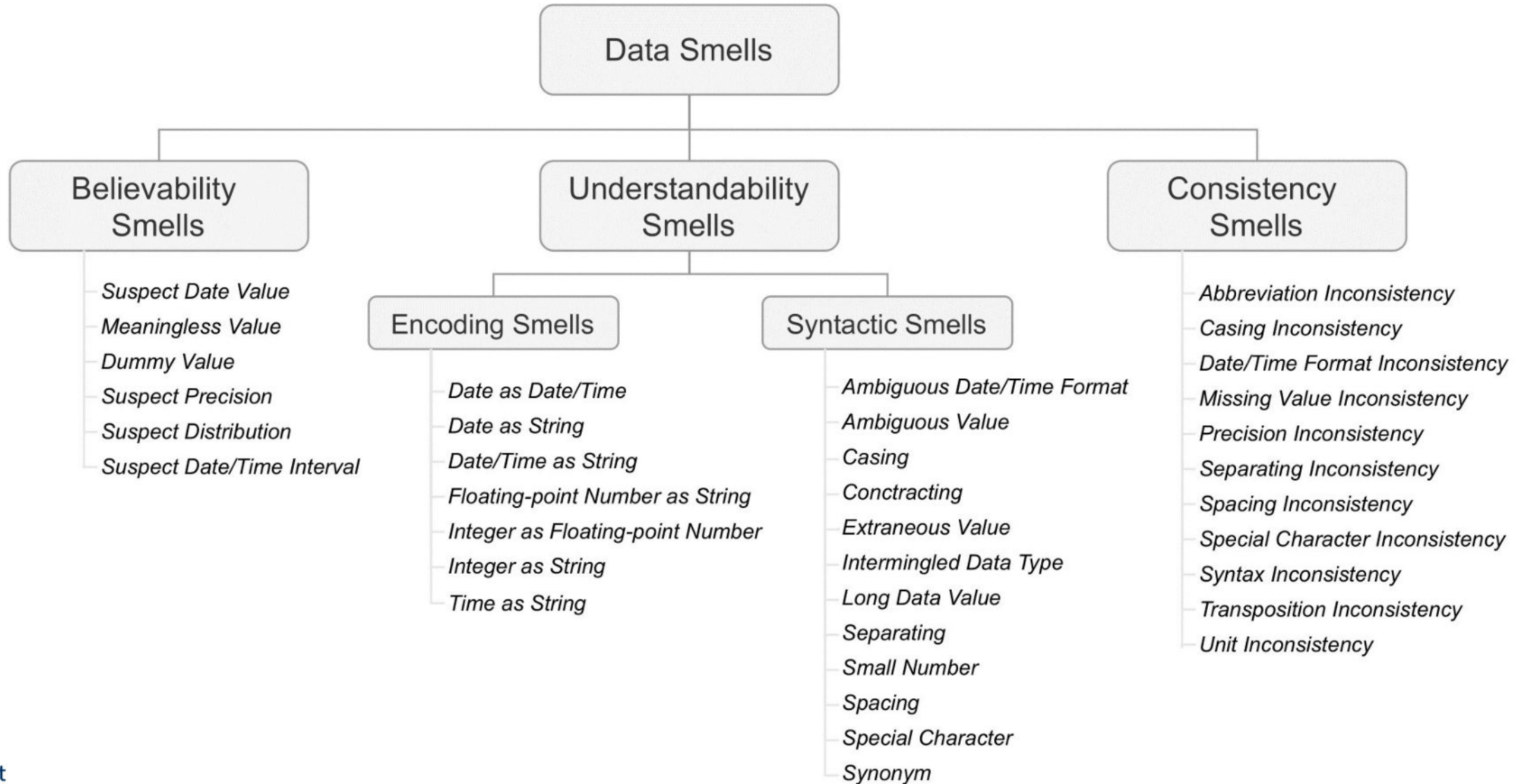


# Data Smells

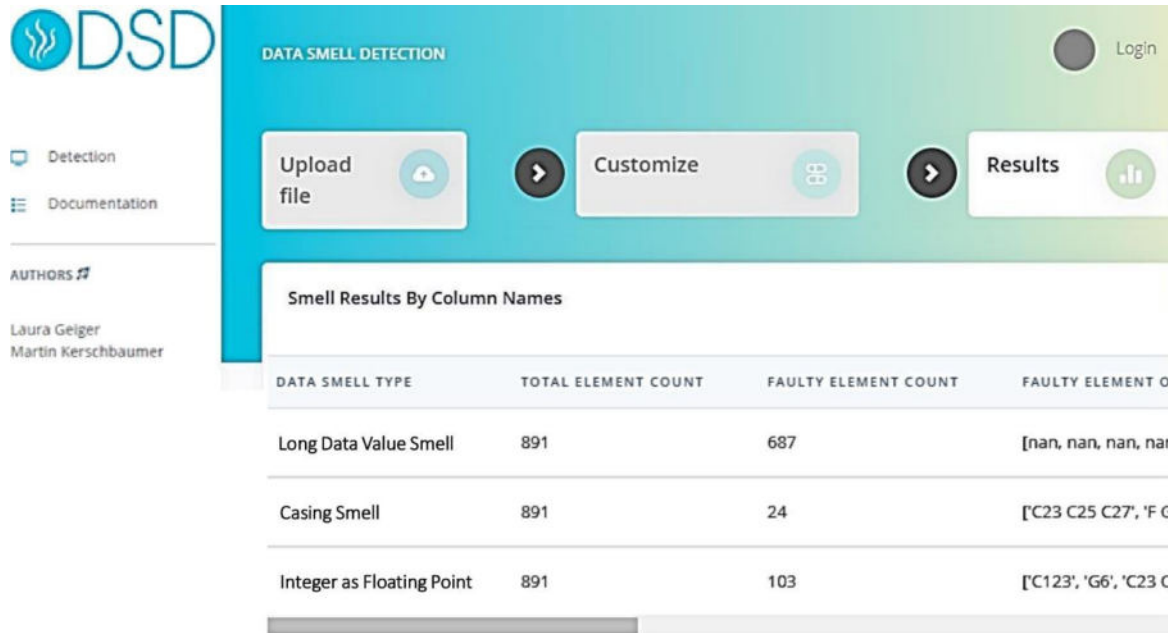
kontextunabhängige, datenwertbasierte Indikatoren für latente Datenqualitätsprobleme, die durch schlechte Praktiken verursacht werden und zu zukünftigen Problemen führen können.



# Data Smells Katalog



# Data Smells Erkennung



DATA SMELL TYPE	TOTAL ELEMENT COUNT	FAULTY ELEMENT COUNT	FAULTY ELEMENT O
Long Data Value Smell	891	687	[nan, nan, nan, nar
Casing Smell	891	24	['C23 C25 C27', 'F C
Integer as Floating Point	891	103	['C123', 'G6', 'C23 C

Regelbasierte Erkennung

## Data Smell Detection with Machine Learning

### LSTM Detection Results

**Agent:** Sample: Date Smell Classification using LSTMs

**Dataset:** Sample Dataset: LSTM Date Classification

[Download Results](#)

### LSTM Classification

Shown below is the class distribution of your data as well as examples for each class. By default, the classes are labeled by the corresponding data smells according to the research paper. This can be disabled on the analyze page. If the class distribution does not match up your expectations, please download the corresponding dataset to further inspect the classification.

#### Class DateTime as String Smell

This class contains **48** total samples, or **3.45%** of the total data.

See some examples below for what data has been classified in this class.

- "03-MAR-1994 09:57PM+06:00" (100.0%)
- "25-APR-2001 00:14+10:00" (100.0%)
- "11-JUN-1977 02:03:44.0051AM +06:00" (100.0%)
- "18-JUN-2015 09:09:25.0453AM +04:00" (100.0%)
- "20-APR-1999 02:51:26AM+03:00" (100.0%)

#### Class Date as DateTime Smell

This class contains **48** total samples, or **3.45%** of the total data.

See some examples below for what data has been classified in this class, and w

- 10-APR-1992 00:00+00:00 (99.93%)
- 15-DEC-1993 00:00+00:00 (99.92%)
- 01-OCT-1999 12:00:00.0000AM +00:00 (100.0%)
- 30-MAR-1995 12:00:00.0000AM +00:00 (100.0%)
- 06-JAN-1976 00:00:00+00:00 (99.93%)

Intelligente Erkennung

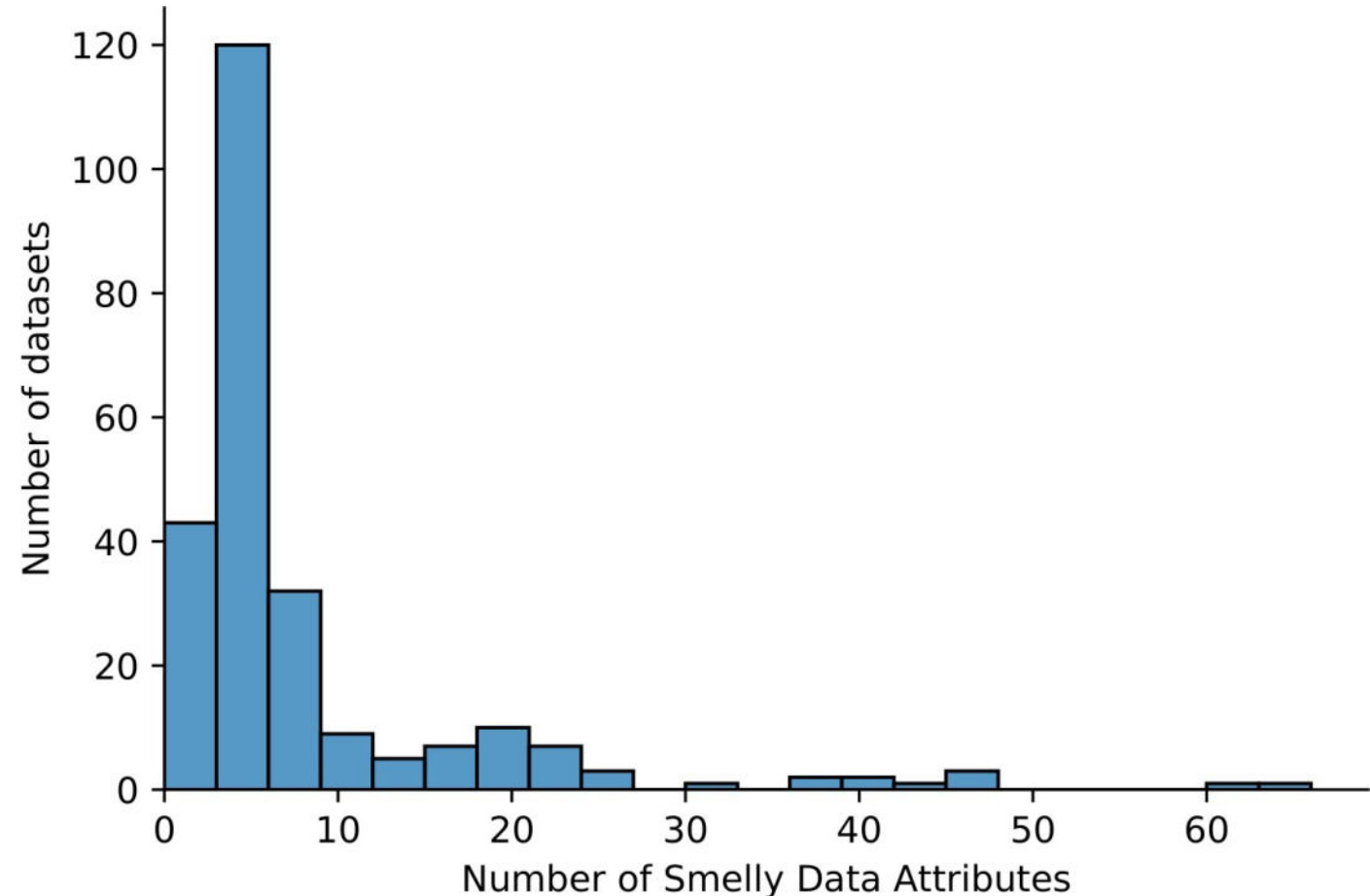
# Data Smells in der Praxis

246 Kaggle Datensätze

> 2000 Spalten

> 42 Mrd. Zeilen

- › Meist drei bis fünf Smells
- › Drei Datensätze keine Smells
- › 45 Datensätze > 10 Smells
- › Ergebnisse handhabbar



# Data Smells Beispiele

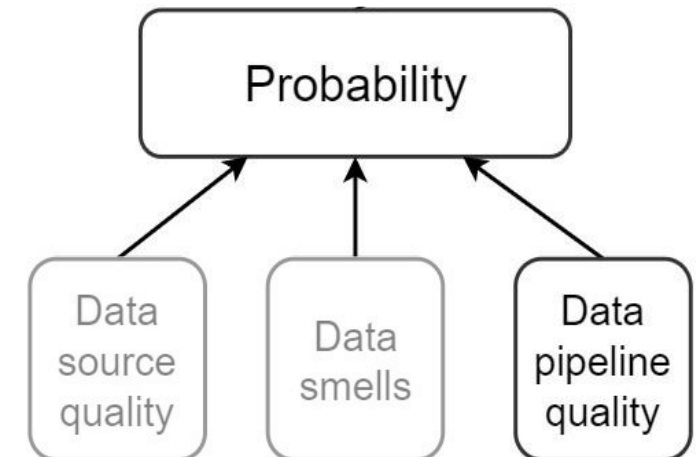
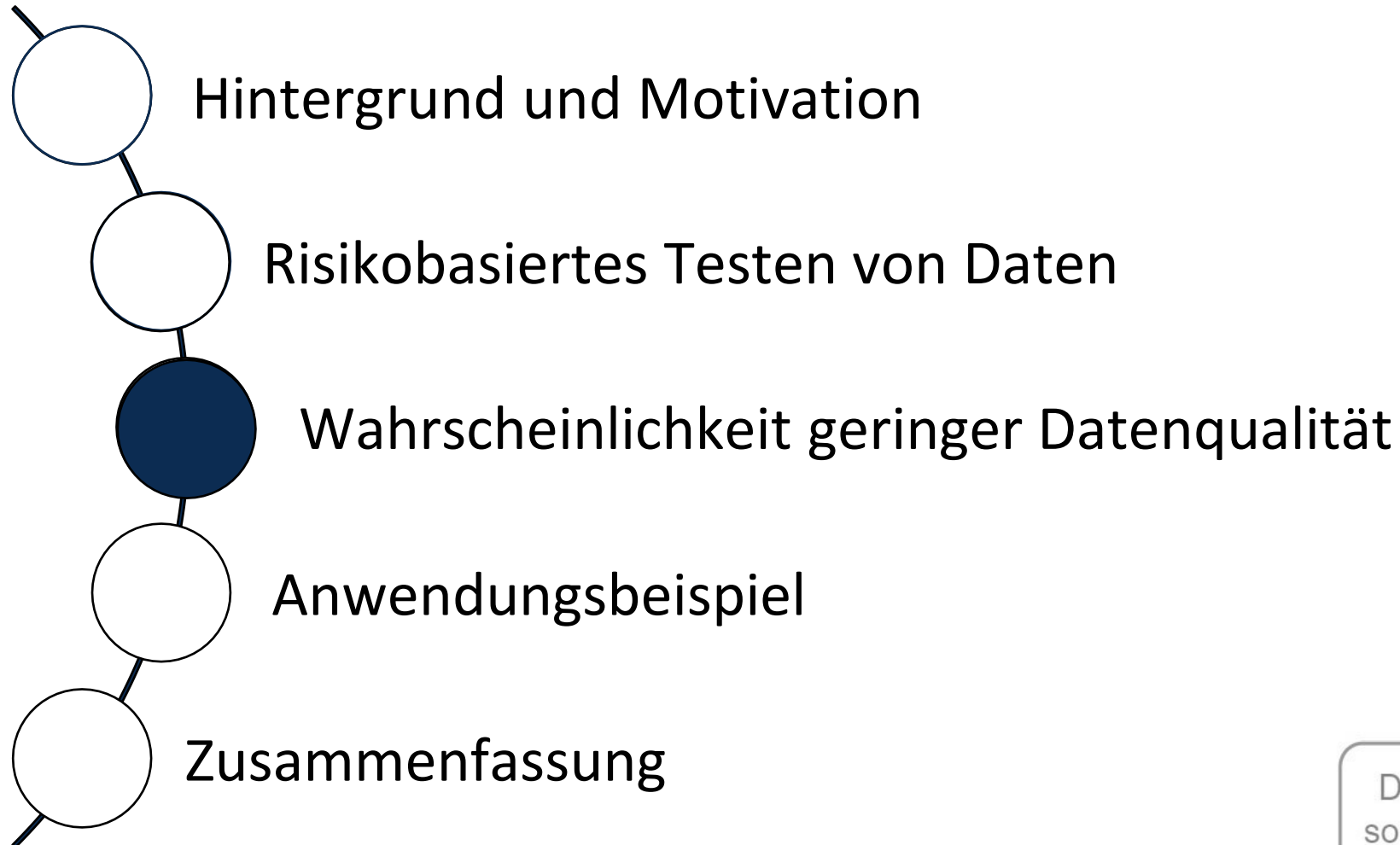
Date	# Daily Confi...	# Total Confi...
11-Feb	0	3
12-Feb	0	3
13-Feb	0	3
14-Feb	0	3
15-Feb	0	3
16-Feb	0	3
17-Feb	0	3
18-Feb	0	3

Ambiguous Date/Time Format

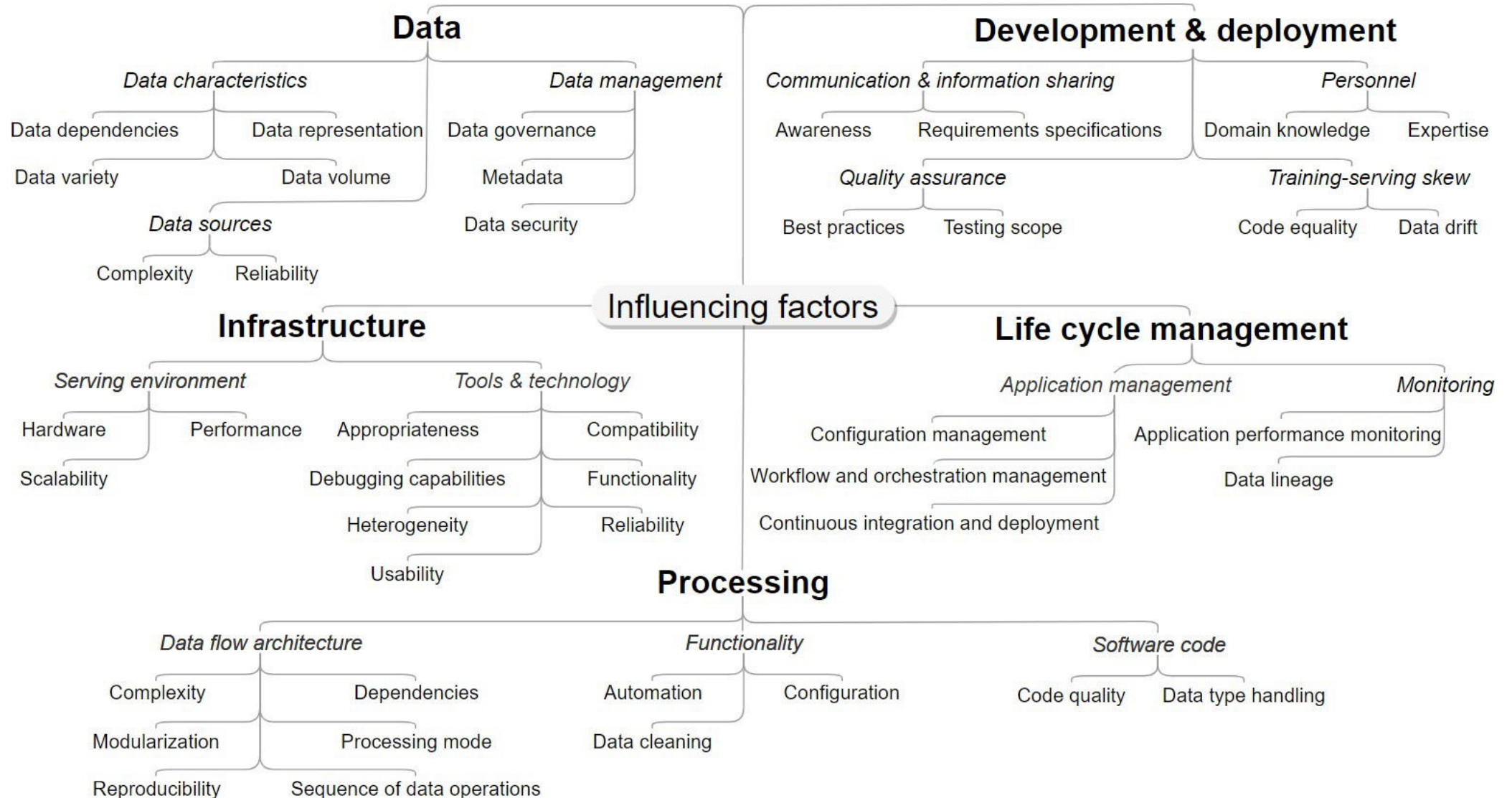
Appointme...	Appointme...	# Age
5638447	2016-04-29T00:00:00Z	21
5629123	2016-04-29T00:00:00Z	19
5630213	2016-04-29T00:00:00Z	30
5620163	2016-04-29T00:00:00Z	29
5634718	2016-04-29T00:00:00Z	22
5636249	2016-04-29T00:00:00Z	28

Date as Date/Time

# Agenda

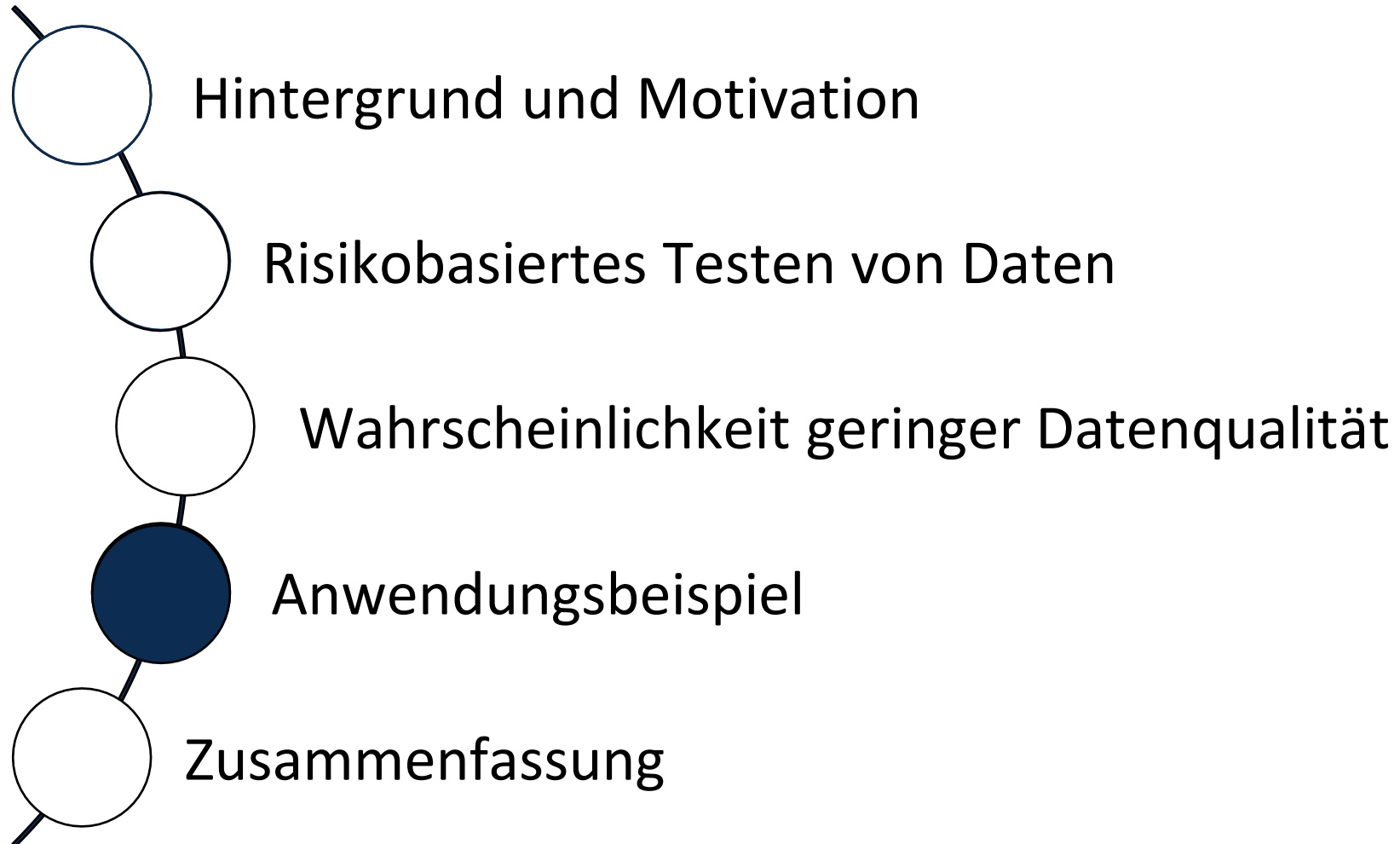


# Datenpipeline Qualität

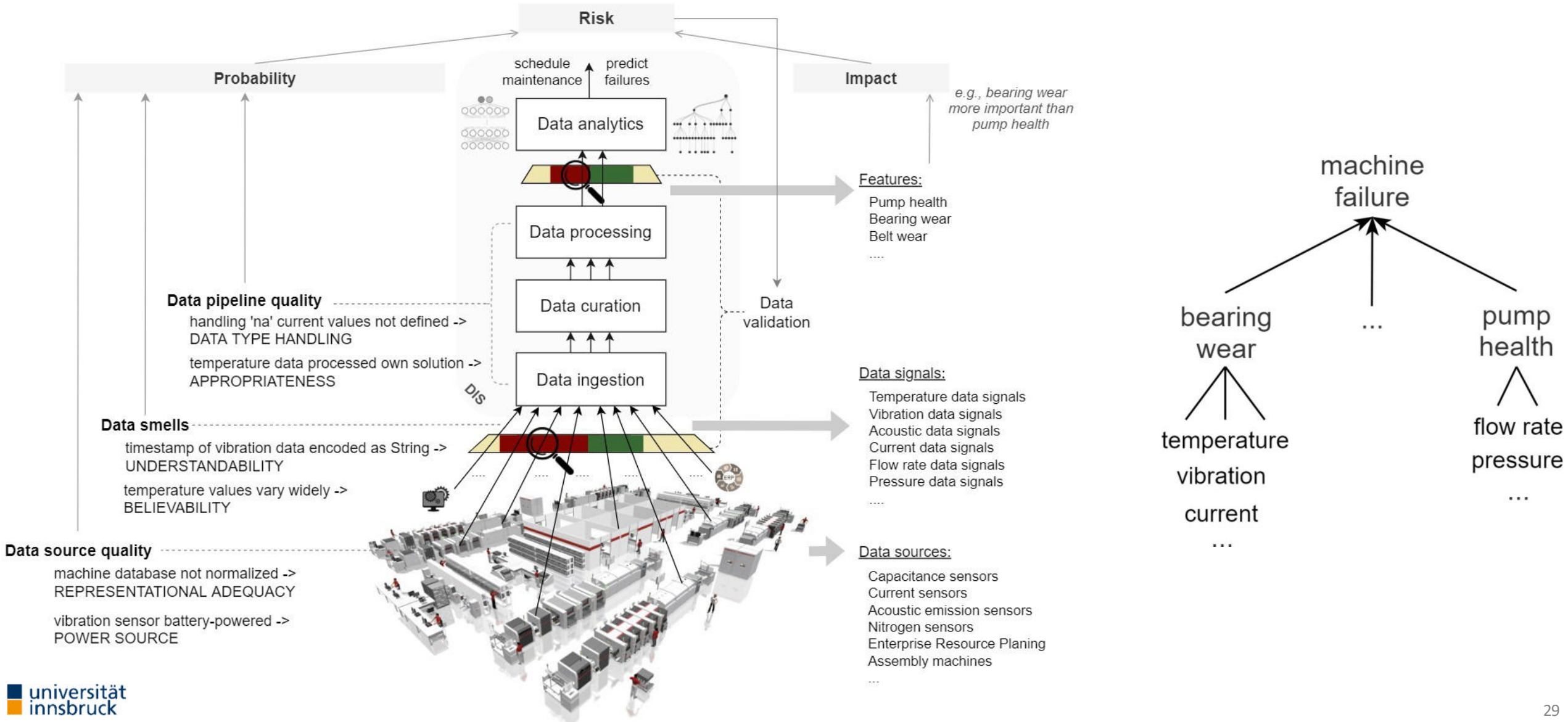


# Agenda

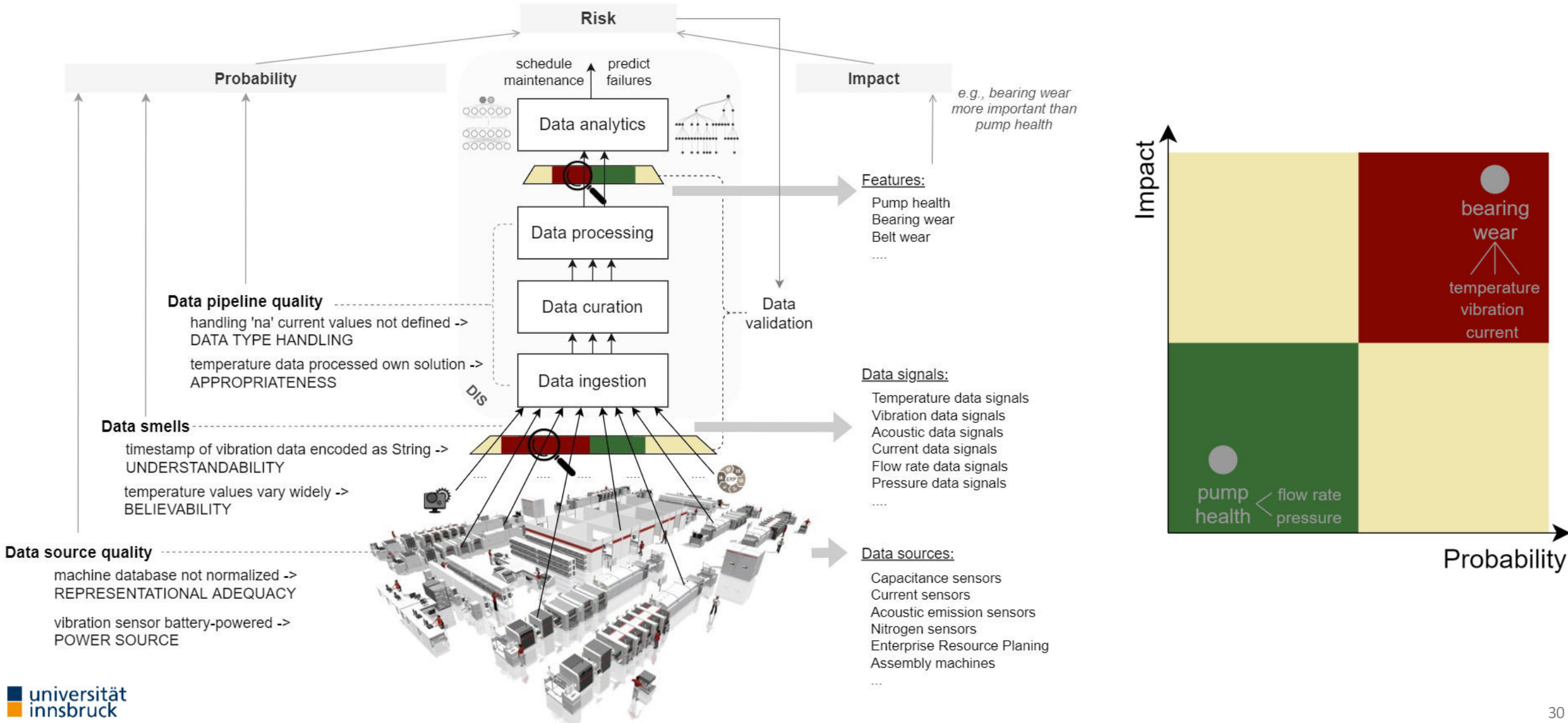
---



# Anwendungsbeispiel

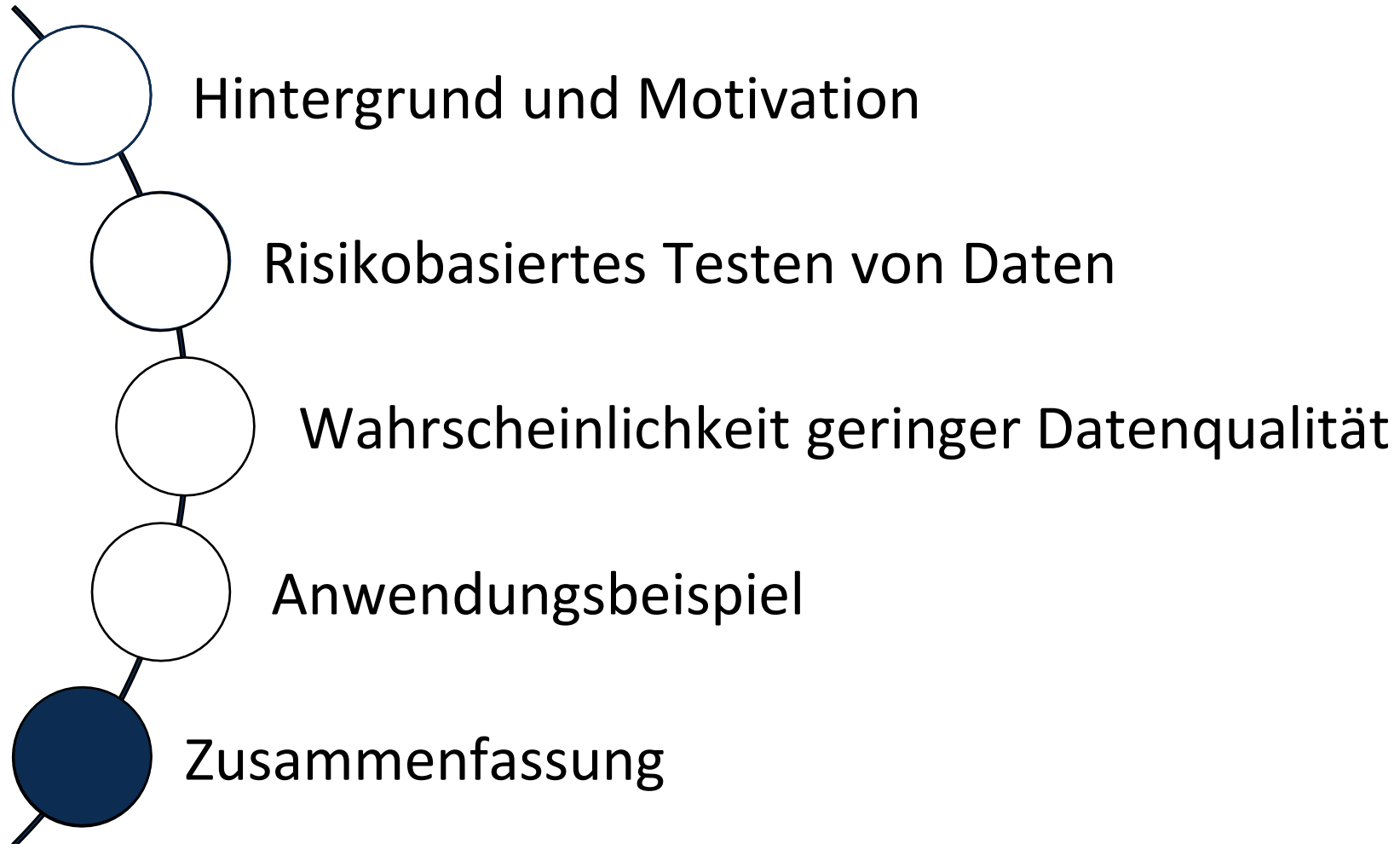


# Anwendungsbeispiel

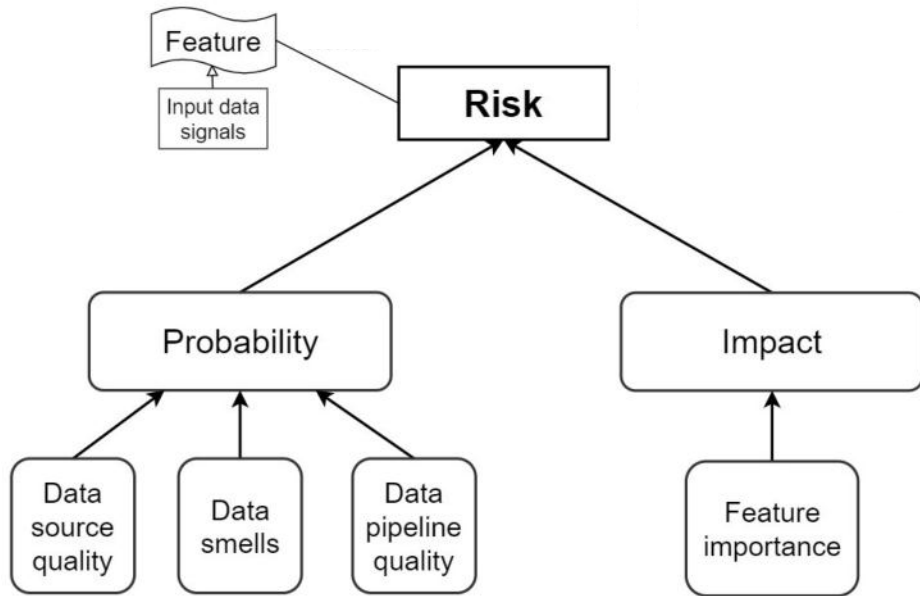


# Agenda

---



# Risikobasiertes Testen von Daten



Ermöglicht selektive Datenvalidierung



Anwendbar in heterogenen Umgebungen



Unterstützt fundierte Entscheidungsfindung



Gezielte Ressourcenzuweisung



Steigerung von Effizienz und Produktivität

# Veröffentlichte Artikel

Foidl, H., & Felderer, M. (2015). Research challenges of industry 4.0 for quality management. In *International Conference on Enterprise Resource Planning Systems*, pp. 121-137. Springer, Cham.

Foidl, H., & Felderer, M. (2016). Data science challenges to improve quality assurance of Internet of Things applications. In *International Symposium on Leveraging Applications of Formal Methods*, pp. 707-726. Springer, Cham.

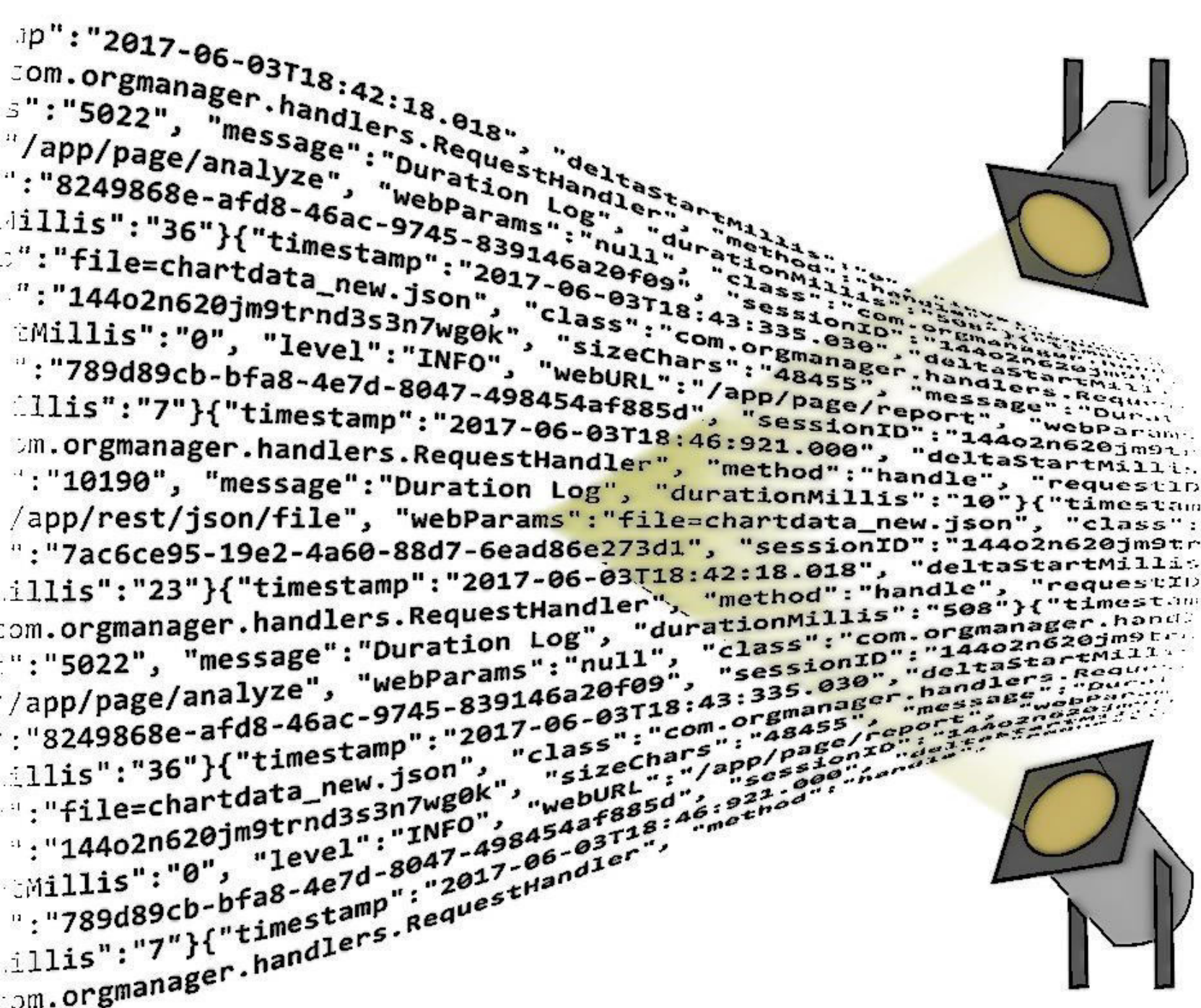
Foidl, H., Felderer, M., & Biffi, S. (2019). Technical debt in data-intensive software systems. In *45th Euromicro Conference on Software Engineering and Advanced Applications (SEAA)*, pp. 338-341. IEEE.

Foidl, H., & Felderer, M. (2019). Risk-based data validation in machine learning-based software systems. In *Proceedings of the 3rd ACM SIGSOFT international workshop on machine learning techniques for software quality evaluation*, pp. 13-18.

Foidl, H., Felderer, M., & Ramler, R. (2022). Data smells: categories, causes and consequences, and detection of suspicious data in AI-based systems. In *1st International Conference on AI Engineering - Software Engineering for AI (CAIN)*, pp. 229-239. IEEE/ACM.

Foidl, H., & Felderer, M. (2023). An approach for assessing industrial IoT data sources to determine their data trustworthiness. In *Internet of Things*, 22, 100735.

Foidl, H., Golendukhina, V., Ramler, R., & Felderer, M. (2024). Data pipeline quality: influencing factors, root causes of data-related issues, and processing problem areas for developers. In *Journal of Systems & Software*, 207, 111855.

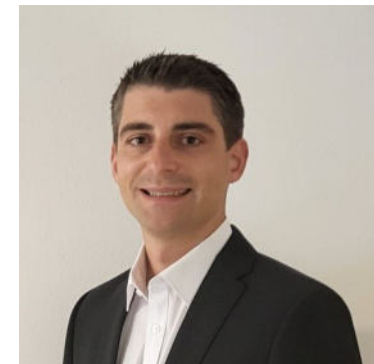


# Danke

## Risikobasiertes Testen von Daten

Ing. Harald Foidl, PhD

Universität Innsbruck



 [harald.foidl@gmail.com](mailto:harald.foidl@gmail.com)

 <http://www.linkedin.com/in/harald-foidl-7616b790>