

Datenqualität in medizinisch-betriebswirtschaftlichen Informationssystemen

MedConf 2013

Endler Gregor, 16.10.2013



FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

TECHNISCHE FAKULTÄT

Warum Datenqualität?



2002, USA:
600.000.000 \$

Y2k weltweit:
1.500.000.000 \$

Kosten



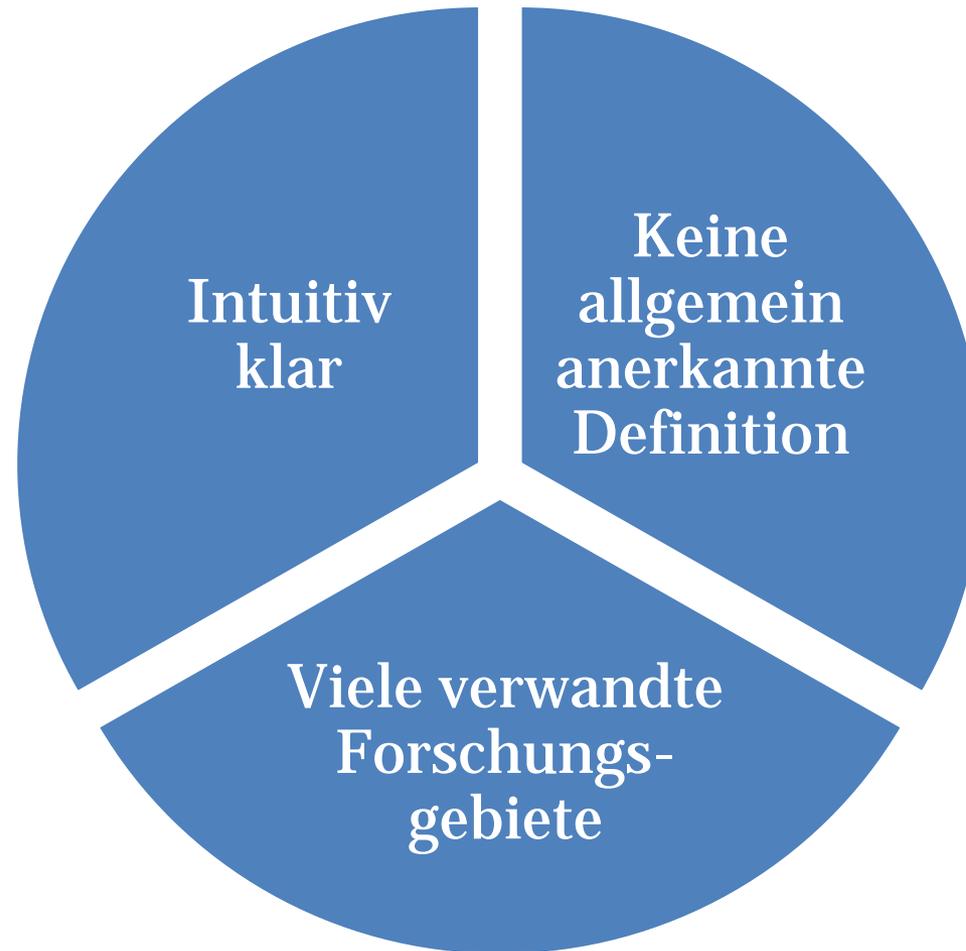
44.000 – 98.000
Todesfälle durch
Fehler



Fehlende Info:
bis 81% der Fälle



Was ist Datenqualität?





Generisch: „Fitness for Use“



1. Datenqualität ist subjektiv

Felder vertauscht

ungenau

fehlende Werte

Duplikate

PID	Name	Vorname	Geburtsjahr	Telefon	MgrZulage
9462	Hans	Müller	1984	1234	
3819	Müller	Hans	1984	1234	
9406	Mustermann	Susanne	1978		
78365	Merkel	Angela	1900		340
2643	Becker	Bris	2015	8374	

Tippfehler

nicht plausibler Wert

unmöglicher Wert

falscher Wert



Viele unterschiedliche Arten von „Qualität“

1. Datenqualität ist subjektiv



2. Datenqualität ist multidimensional





Korrektheit

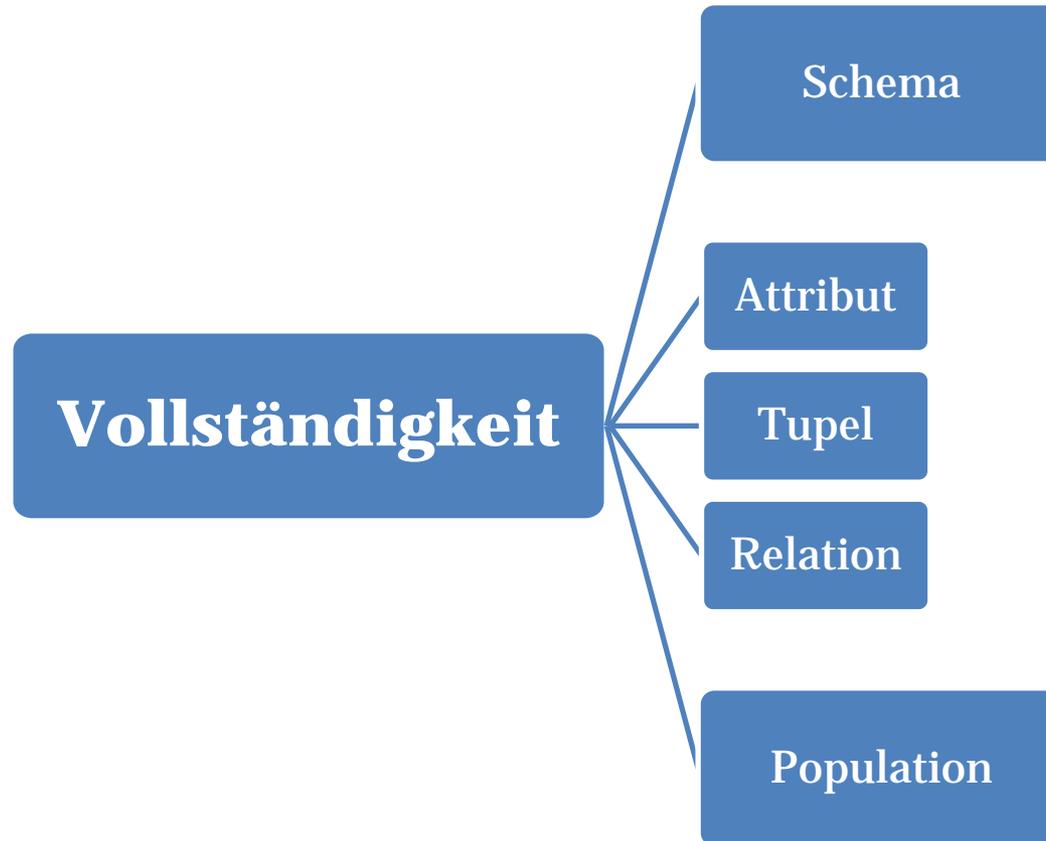
- Übereinstimmung
Datenwert - Realwelt
- Syntaktisch
vs.
Semantisch

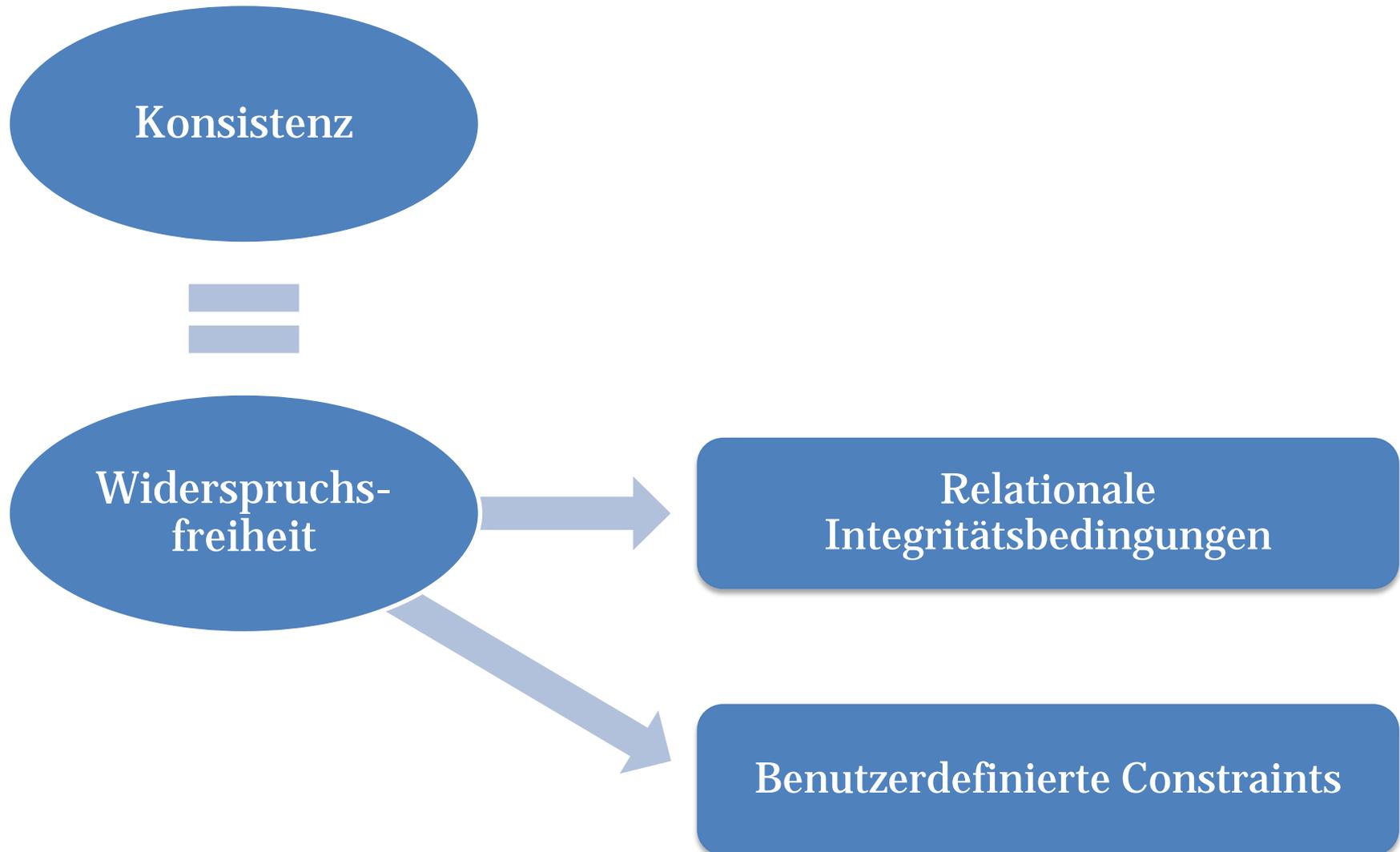
kontext*un*abhängig

Genauigkeit

- Abstand
Datenwert - Realwelt

kontext*ab*hängig







Aktualität

Daten veraltet?

Zeitnähe

**Zeitgerechte
Bereitstellung?**





Wechselwirkungen zwischen Dimensionen

1. Datenqualität ist subjektiv

2. Datenqualität ist multidimensional

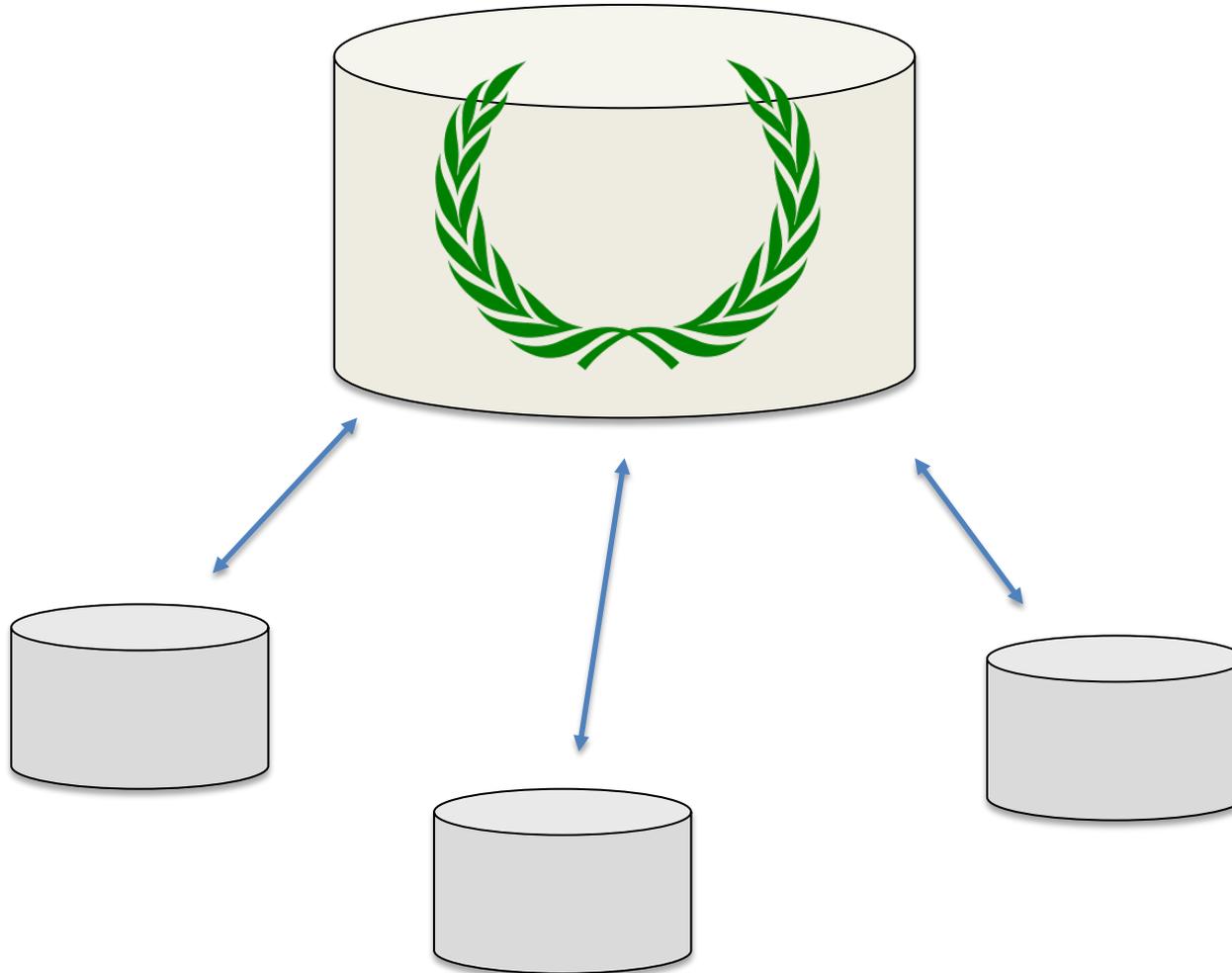
 3. DQ-Dimensionen sind nicht unabhängig



Messen von Datenqualität



FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG
TECHNISCHE FAKULTÄT



Korrektheit

- Syntaktisch: Wertebereich, Rechtschreibung
- Semantisch: Realweltvergleich



Genauigkeit

- Realweltvergleich
- Spezialfall Distanzmessung an Realobjekt



Konsistenz

$$\frac{\# \text{ Tupel, die alle Integritätsbedingungen erfüllen}}{\# \text{ Tupel}}$$

Zeitnähe

- Fehlende Daten für Arbeitsschritt?
- Prozessmonitoring



Vollständigkeit

- Attribut, Tupel, Relation: Anteil *NULL*
- Schema: Schema- & Bedarfsanalyse
- Populationsvollständigkeit: Expertenwissen, Realweltvergleich

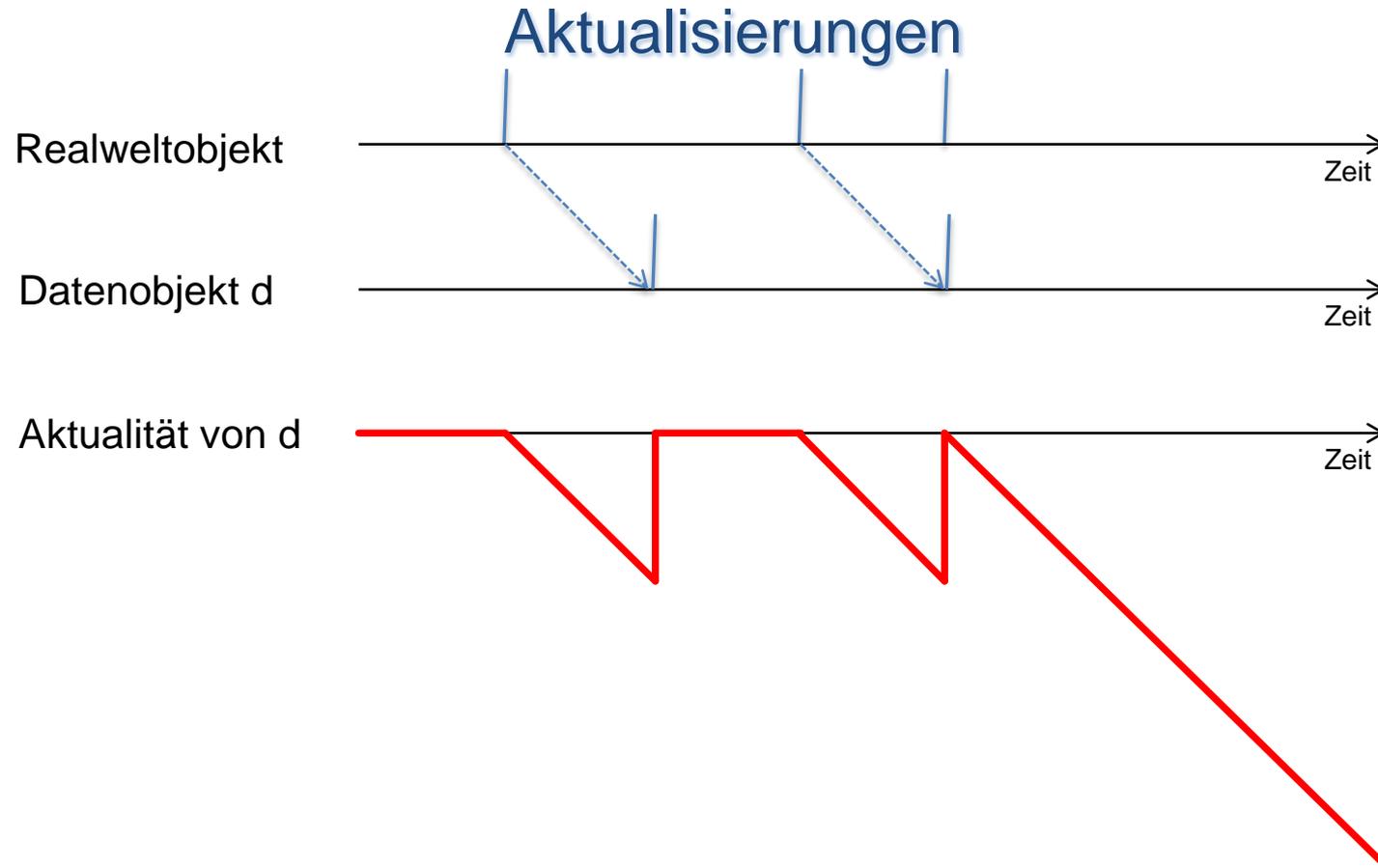


NULL

Wert existiert nicht

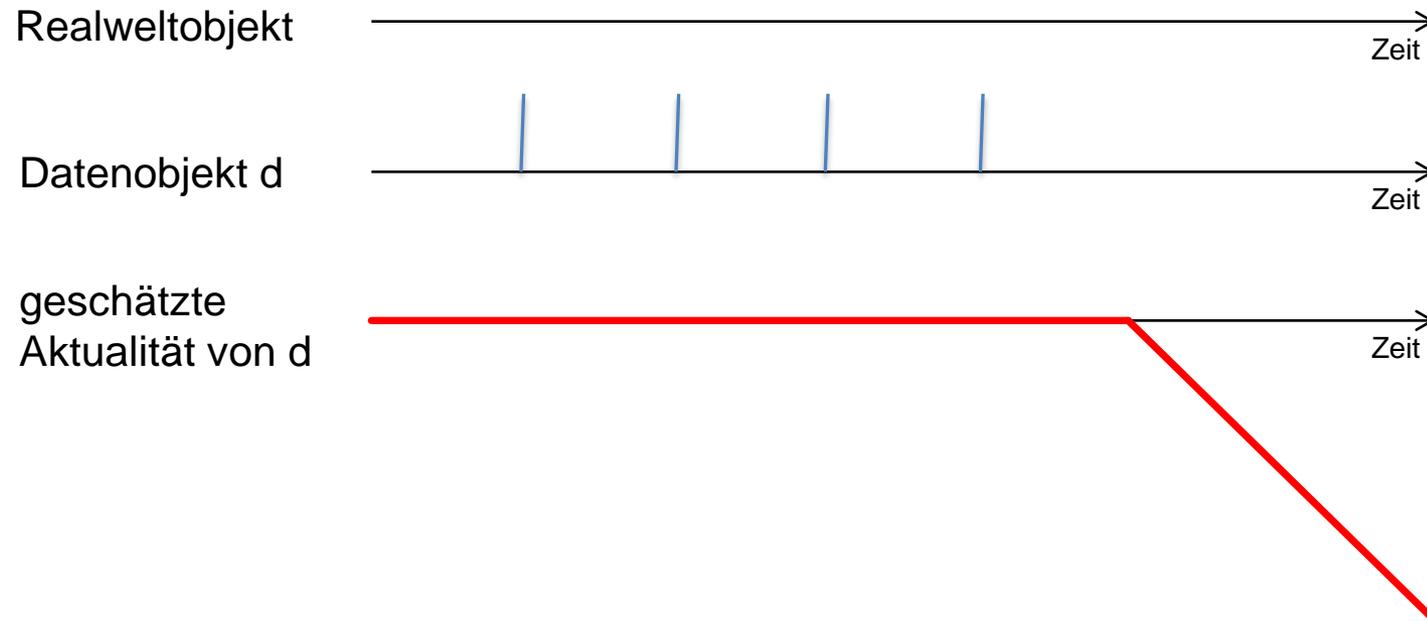
Wert existiert, ist aber nicht bekannt

Nicht bekannt, ob Wert existiert





Volatilität ?



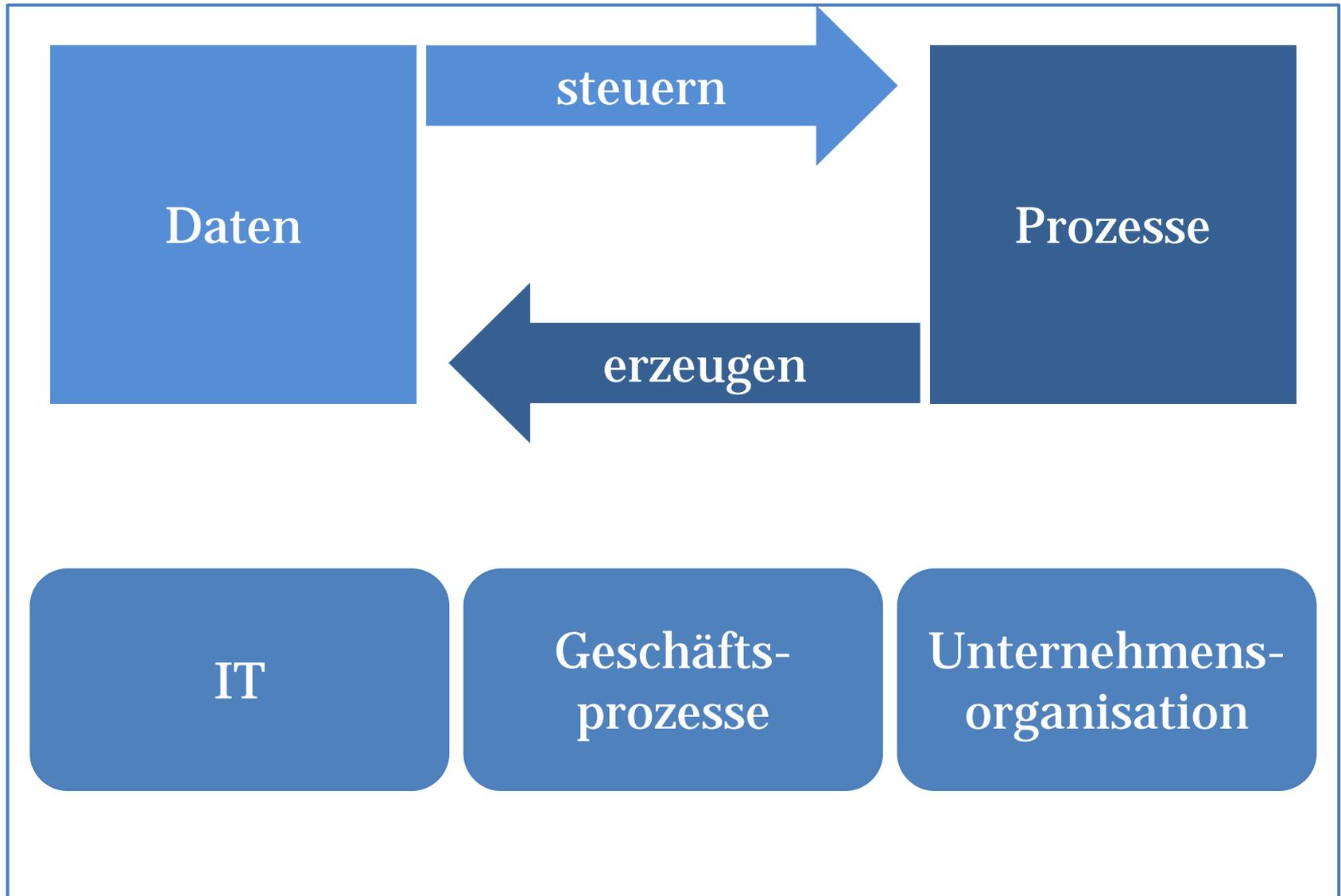


Verbesserung von Datenqualität

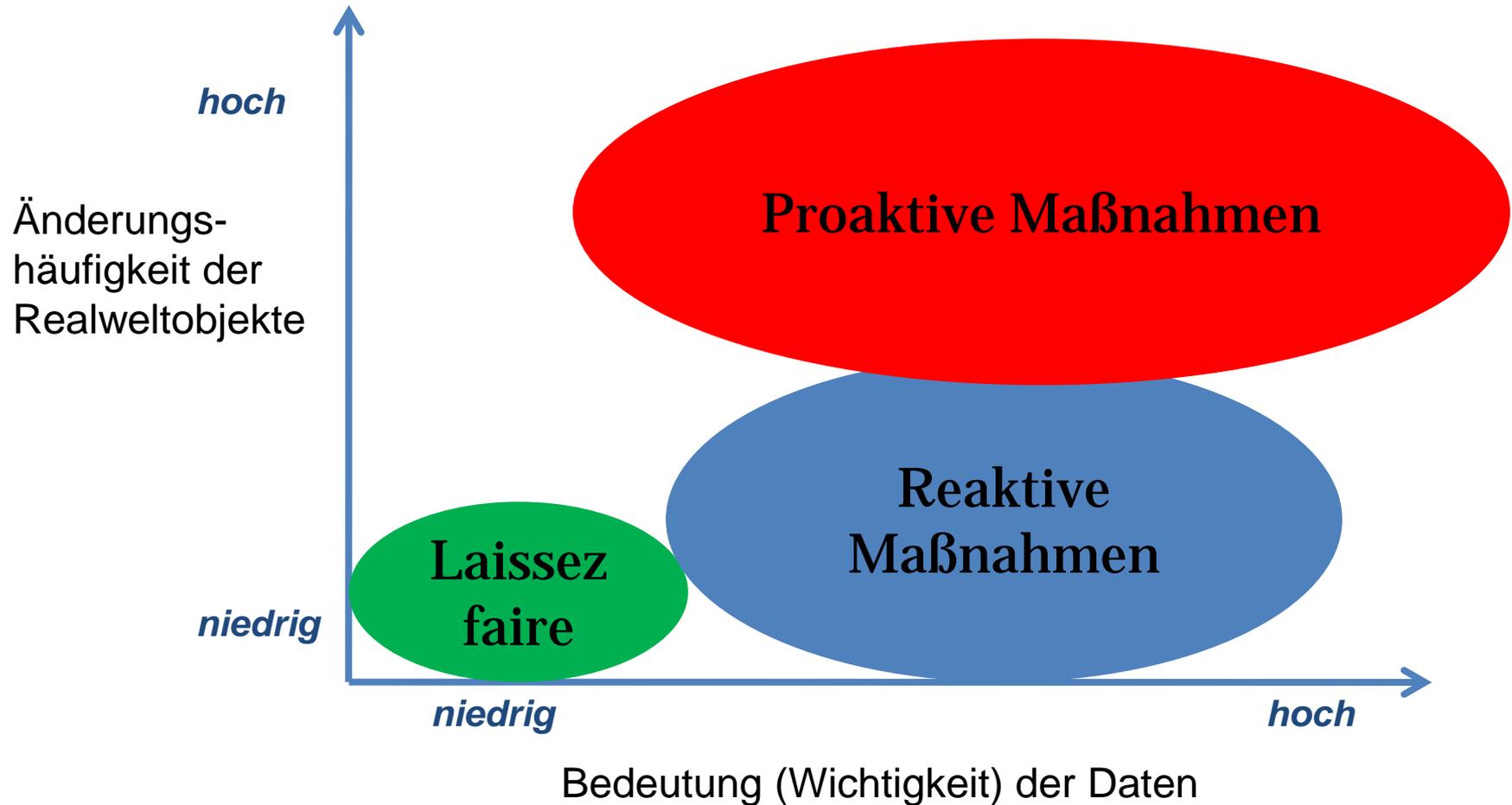


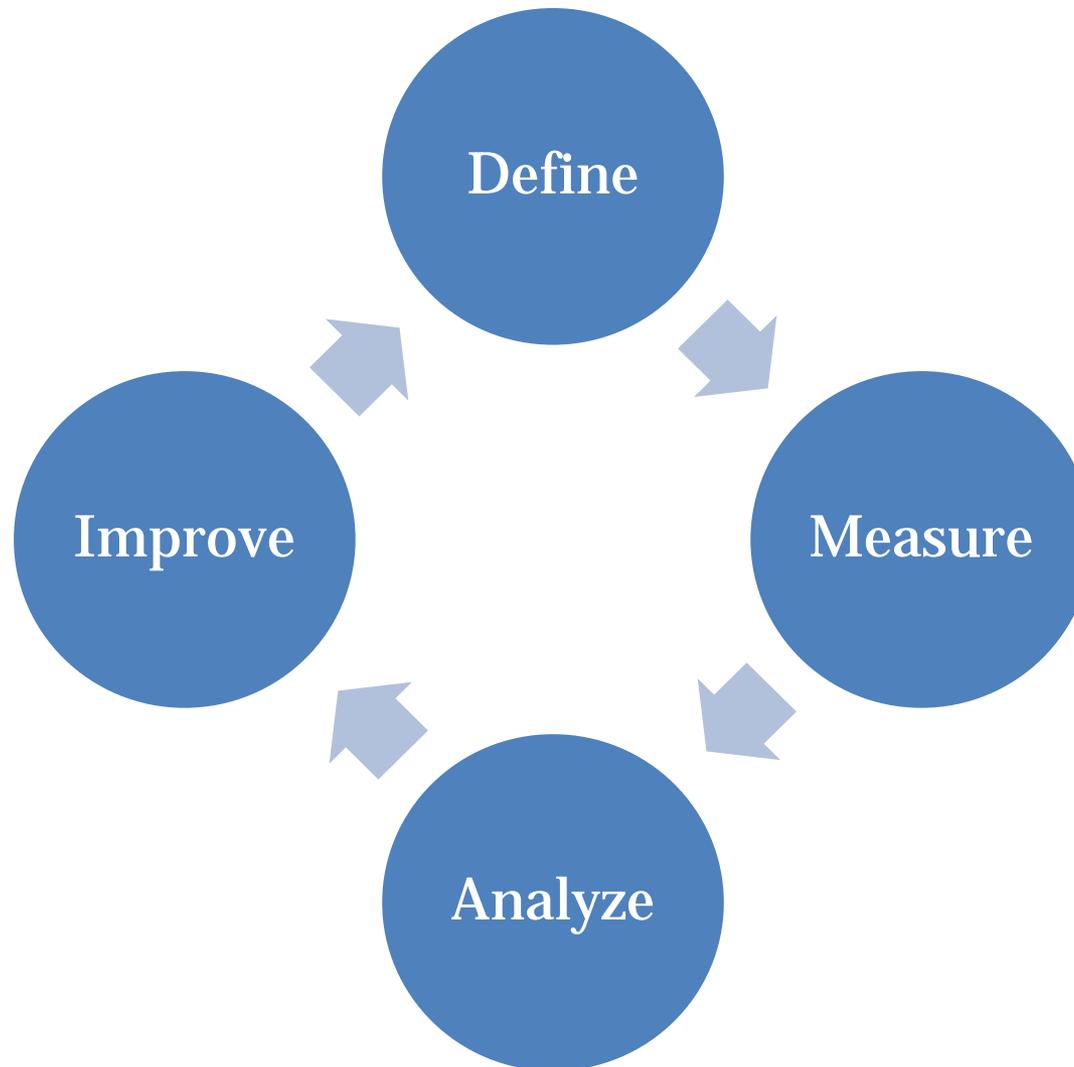
FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

TECHNISCHE FAKULTÄT



Maßnahmenportfolio zur Verbesserung der DQ
[Redman 1996]





TDQM: Wang et al: "Data Quality", Kluwer, 2000



... der DQ

... der Datenproduktion

A1, QS1	A2, QS2	A3, QS3
w1, qw1	w2, qw2	...
...

QS2	q1	q2	...
qw2	qind1	qind2	...
...



... der DQ

... der Datenproduktion

Provenance

Why

Where

How



Maßnahmen



Profiling

Attributanalyse

Abhängigkeiten

**Fremdschlüssel-
beziehungen**

Redundanzen

Manuelle Korrektur

...

**(semi-)automatische
Korrektur**

Ausreißer

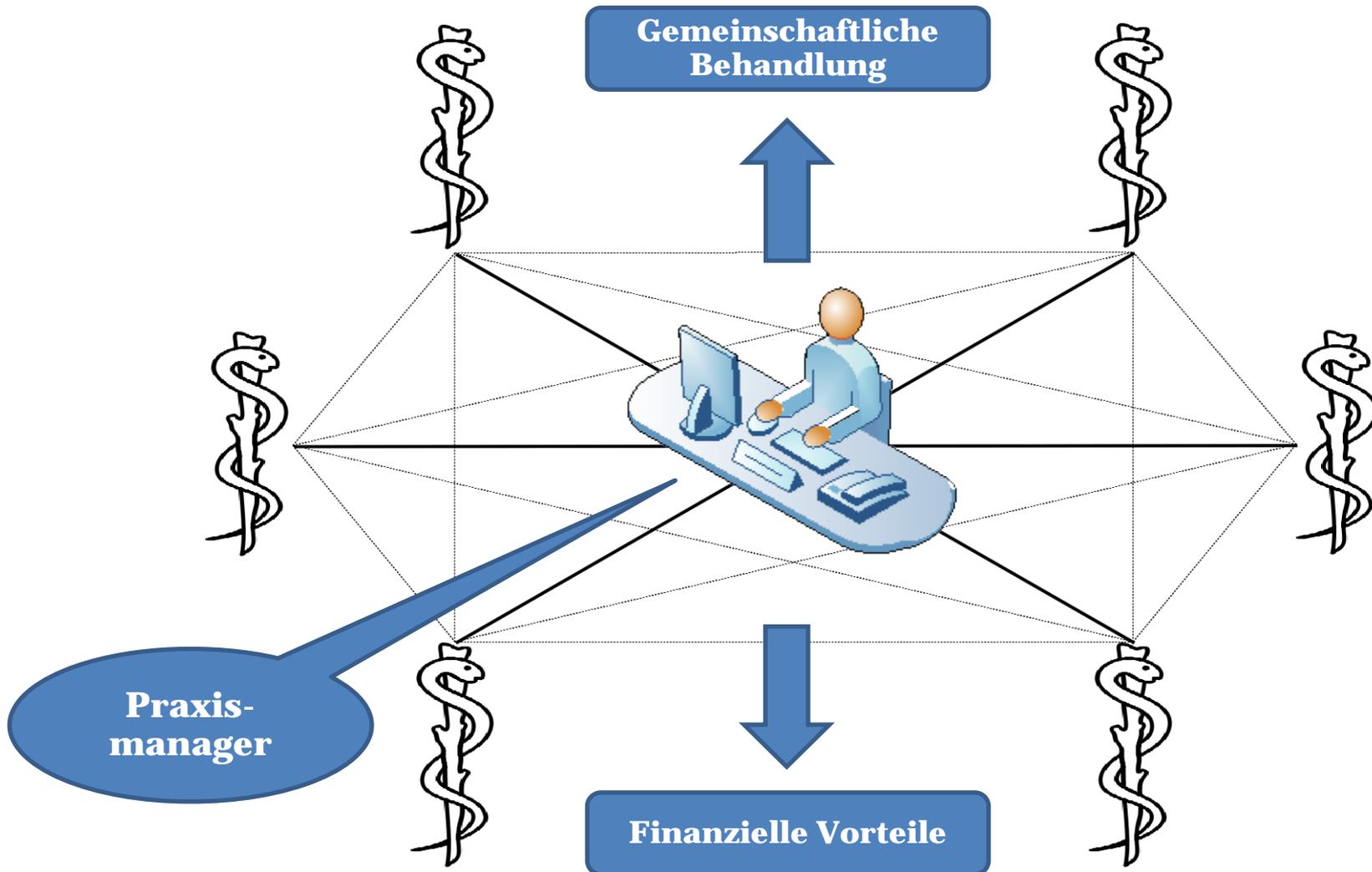
Identity Matching

Record Linkage

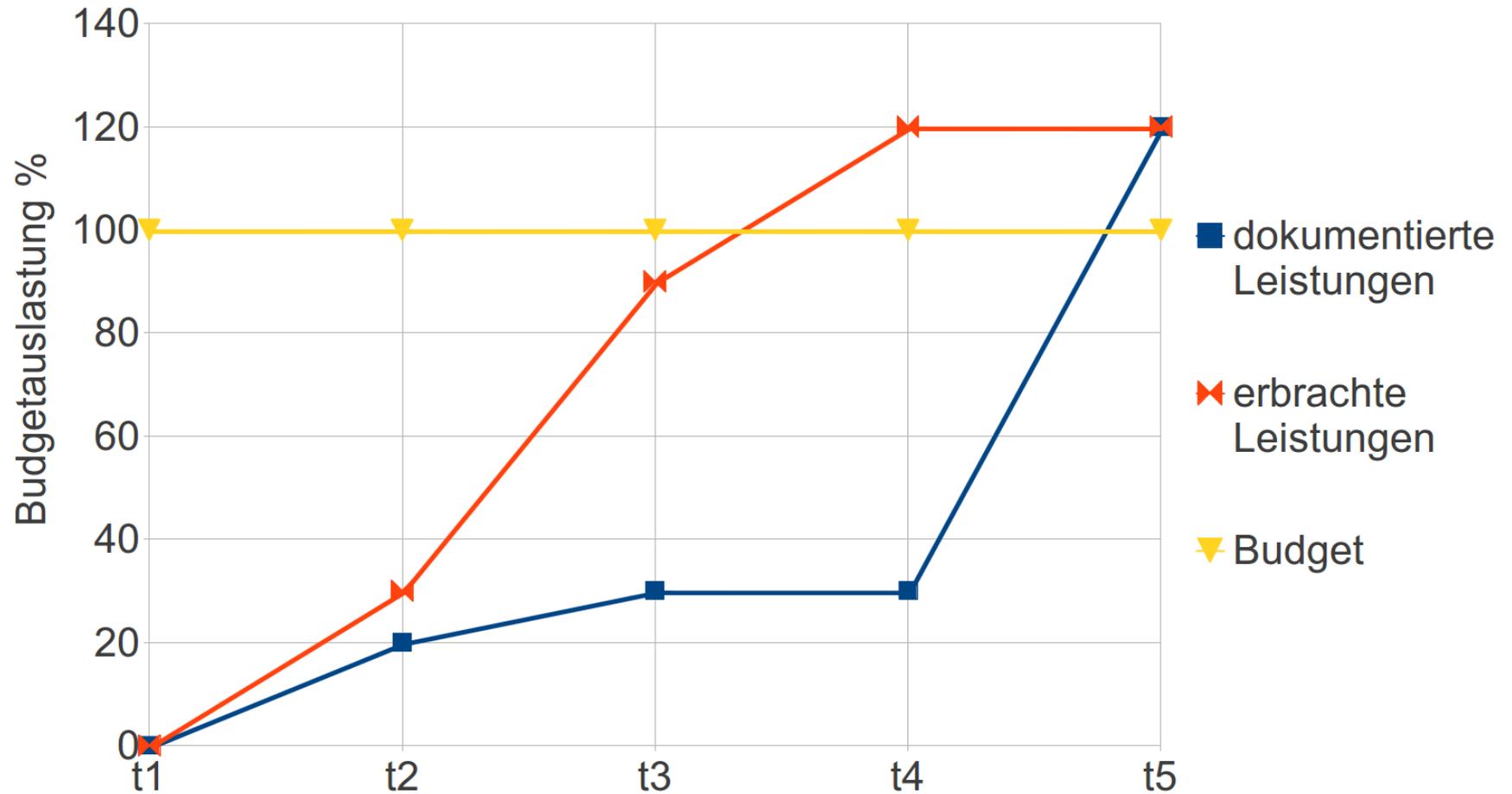


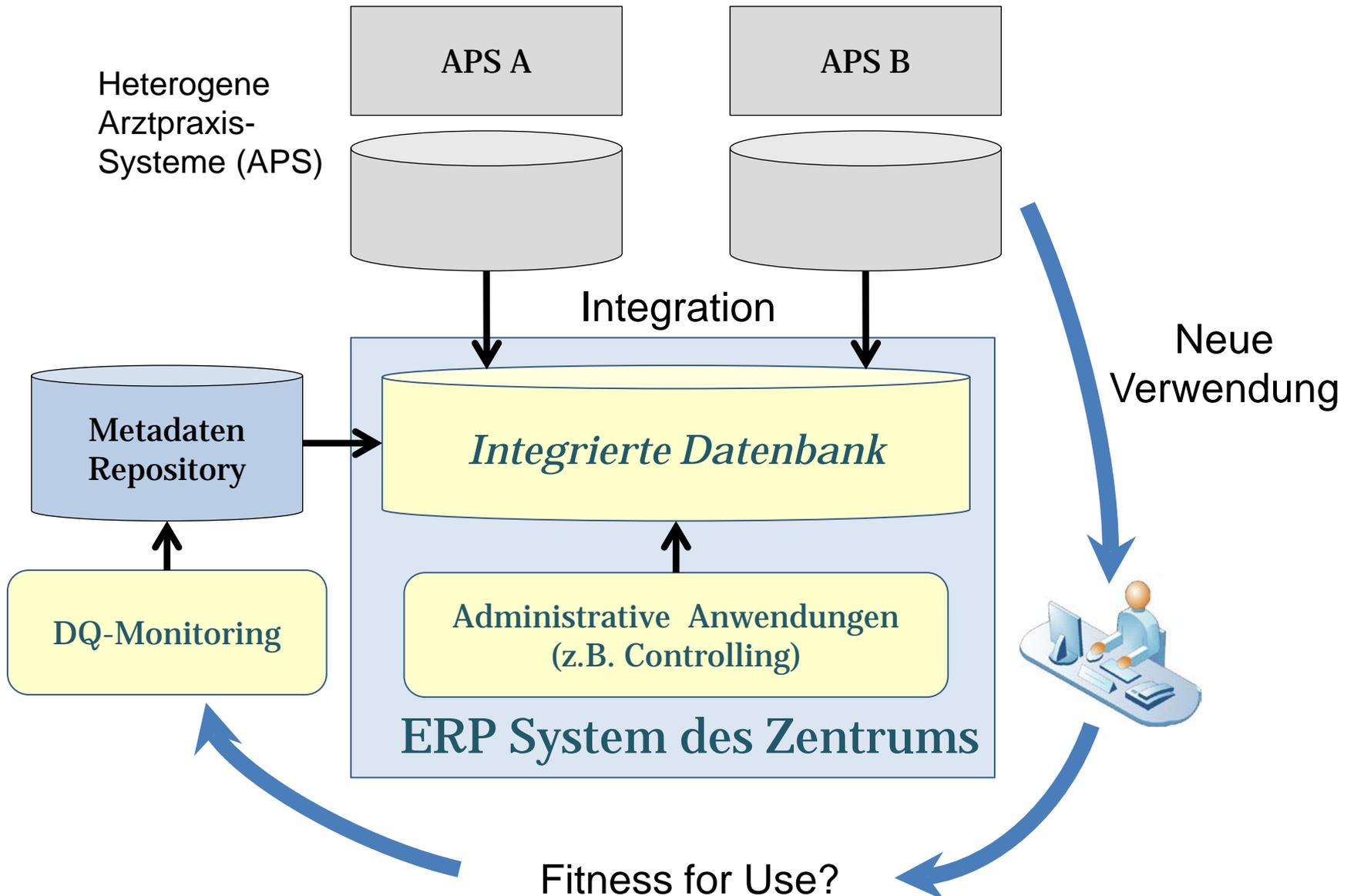
Datenqualität im Projekt MEDITALK

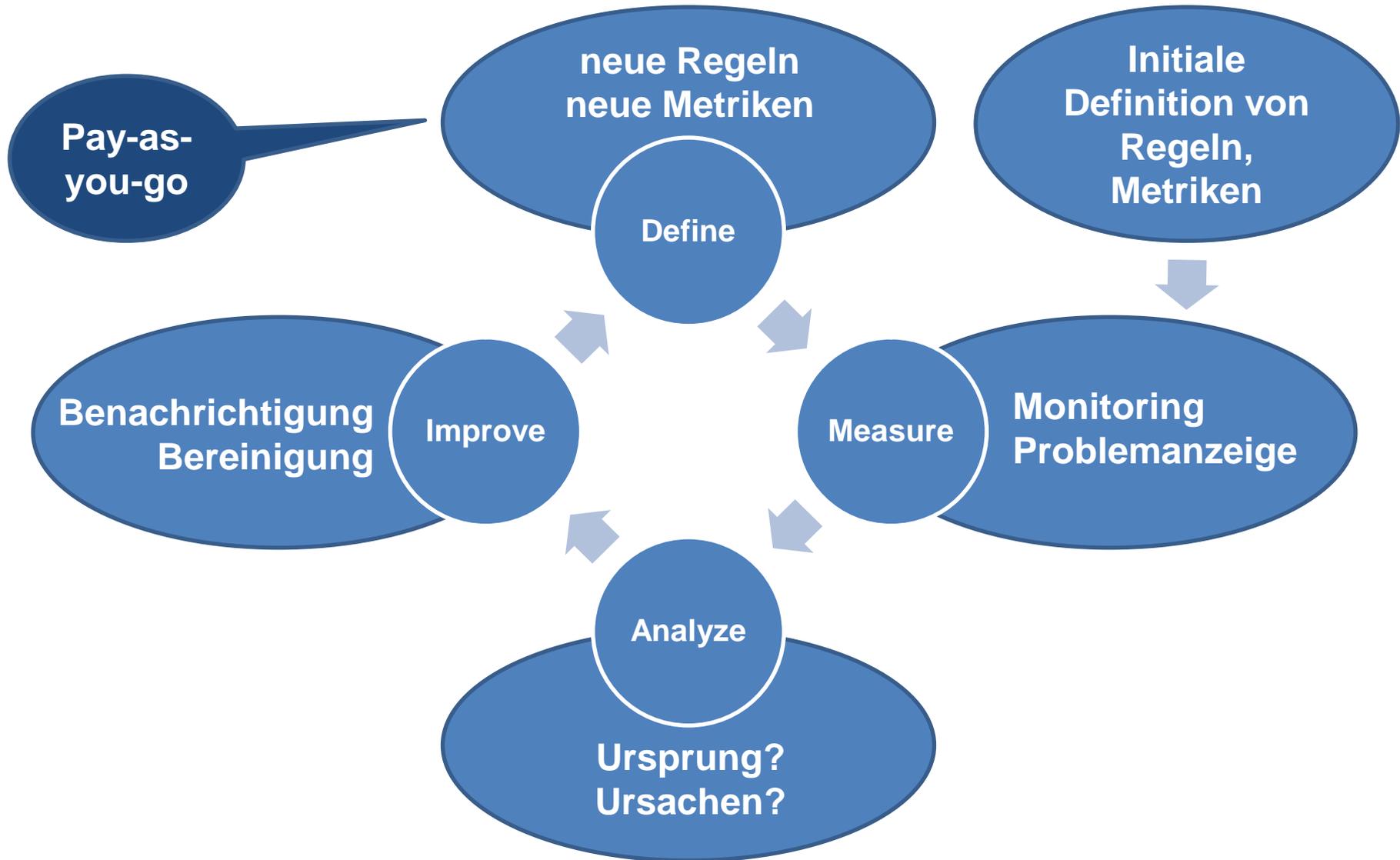




Wo drückt der Schuh?







TDQM: Wang et al: "Data Quality", Kluwer, 2000



Datenqualität ist subjektiv

Datenqualität ist multidimensional

DQ-Dimensionen sind nicht unabhängig

Messen entlang der Dimensionen

Oft kontinuierliches DQ-Management nötig

Kontakt

Gregor Endler
Lehrstuhl für Informatik 6 (Datenmanagement)
FAU Erlangen-Nürnberg

gregor.endler@fau.de

www6.cs.fau.de/people/greg/



FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

TECHNISCHE FAKULTÄT



Batini, C. and Scannapieco, M.:

Data Quality. Concepts, Methodologies and Techniques, Springer, 2006

Eckerson, W.:

Data Quality and the Bottom Line: Achieving Business Success through a Commitment to High Quality Data, The Data Warehouse Institute, Repost Series, 2002

Fisher, C.W. and Kingma, B.R.:

Criticality of Data Quality as Exemplified in Two Disasters, Information Management, 2002

English, L.P.:

Improving Data Warehouse and Business Information Quality, Wiley & Sons, 1999

Institute of Medicine:

IOM Report 1999

IOM Report 2001

Lenz, R.Y.:

Vorlesungsmaterial *Evolutionäre Informationssysteme*, 2012

Miller, D.W., et al.:

Missing prenatal records at a birth center: A communication problem quantified, AMIA Annu. Symp. Proc., 2005

Redman, T.C.:

Data Quality for the Information Age, Artech House, 1996

Wang, R. et.al.:

Data Quality, Kluwer, 2000