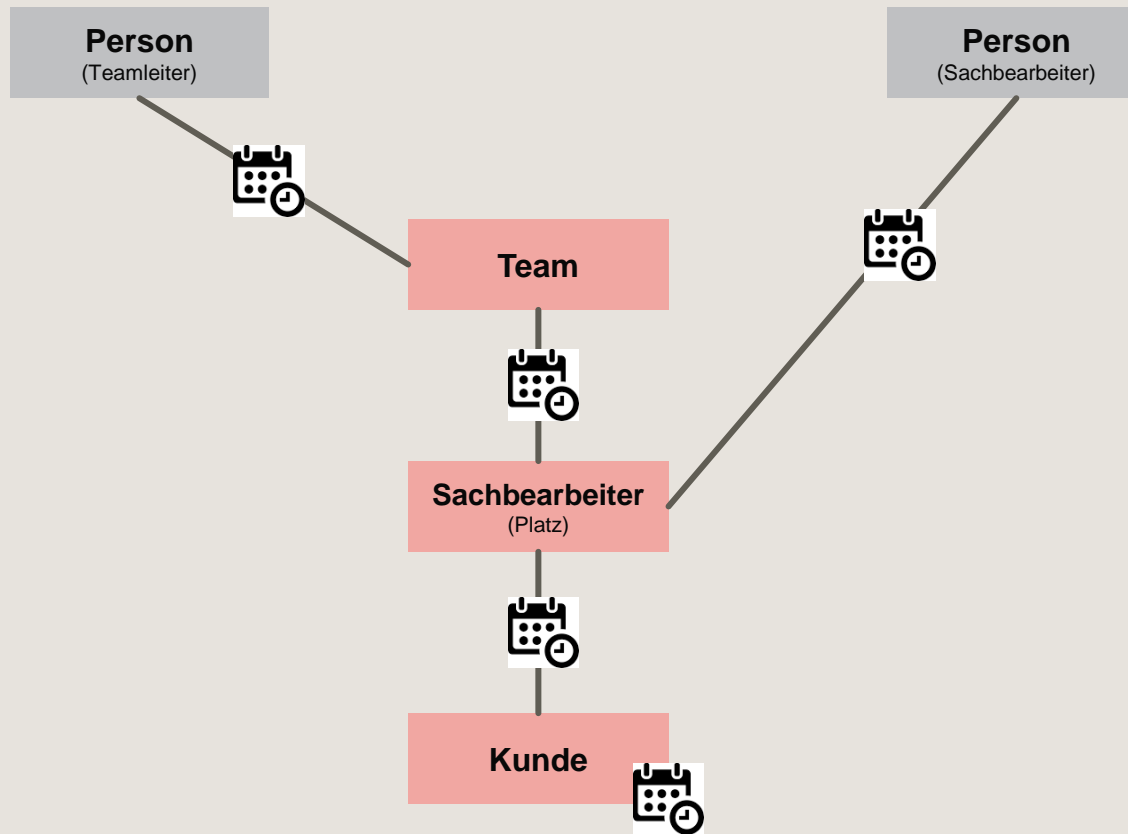




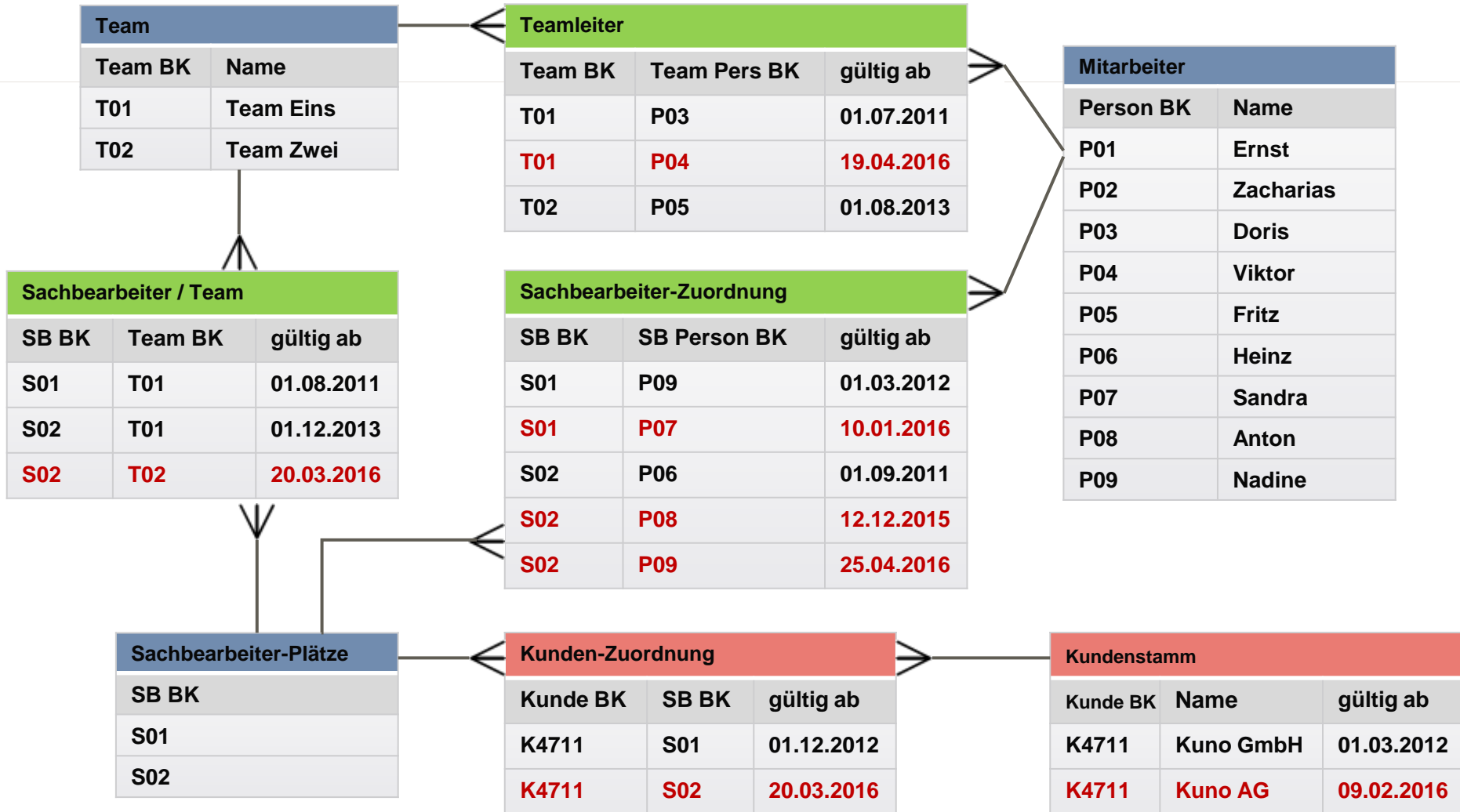
SQL-basierte SCD2-Versionierung hierarchischer Strukturen

Meik Truschkowski
nobilia-Werke J. Stickling GmbH & Co. KG | Verl
Projektleiter Business Intelligence und Data Warehousing

Problem: Historisierung (Versionierung)



Änderungen in den Quelldaten



Zeitleiste

Nov 2015	Dez 2015	Jan 2016	Feb 2016	Mrz 2016	Apr 2016	Mai 2016	Jun 2016	Jul 2016
----------	----------	----------	----------	----------	----------	----------	----------	----------

Team: Teamleiter	T02: P05 (Fritz)							
	T01: P03 (Doris)					T01: P04 (Viktor)		
SB-Platz: Team	S02: T01				S02: T02			
	S01: T01							
SB-Platz: SB-Person	S02: P06 (Heinz)		S02: P08 (Anton)			S02: P09 (Nadine)		
	S01: P09 (Nadine)		S01: P07 (Sandra)					
Kunde: SB-Platz	K4711: S01				K4711: S02			
Kunde (Inhalt)	K4711: Kuno GmbH				K4711: Kuno AG			
Version = Kunde SK	1		2	3	4	5		

Ergebnismenge

kleinstmögliche Zeitscheiben

KUNDE SK	KUNDE BK	KUNDE_NAME	SB BK	SB PERSON BK	TEAM BK	TEAM PERSON BK	GÜLTIG AB	GÜLTIG_BIS
1	K4711	Kuno GmbH	S01	P09	T01	T03	01.01.2000	09.02.2016
2	K4711	Kuno GmbH	S01	P07	T01	T03	10.01.2016	08.02.2016
3	K4711	Kuno AG	S01	P07	T01	T03	09.02.2016	19.03.2016
4	K4711	Kuno AG	S02	P08	T02	P05	20.03.2016	24.04.2016
5	K4711	Kuno AG	S02	P09	T02	P05	25.04.2016	31.12.9999

KIMBALL!

...nun zur Lösung in SQL...

Moment mal – in **SQL**???

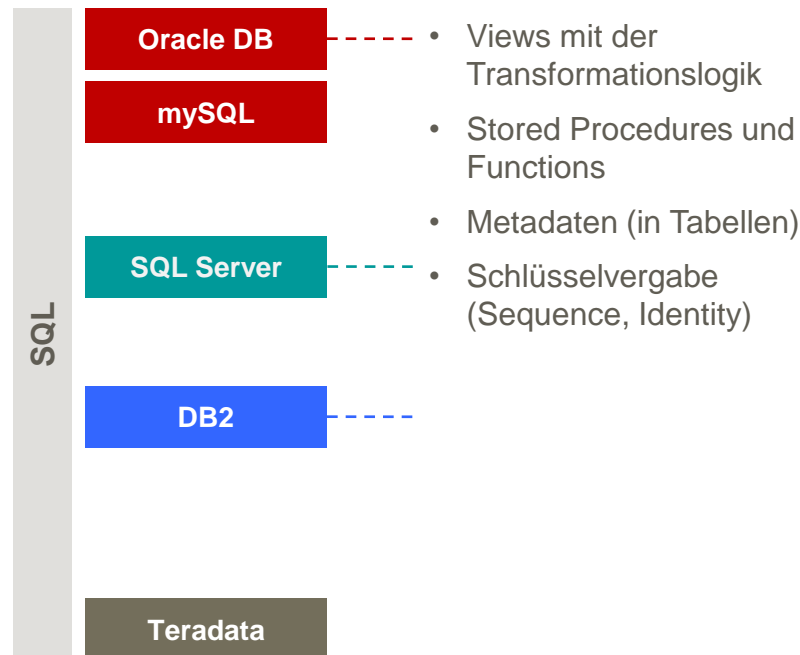
Als Skript?

Es gibt doch moderne **ETL**-Tools!?

Hybrides ~~ETL~~ ELT

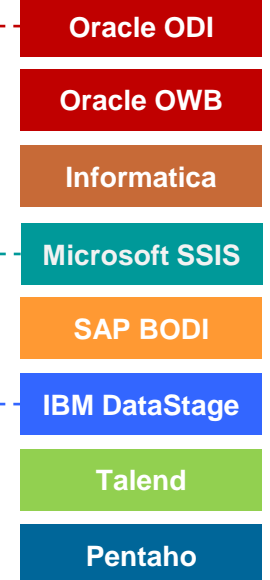
Was heißt hier eigentlich Plattform-unabhängig?

zentrales Data Warehouse (RDBMS)



ETL- bzw. ELT-Werkzeug

- Ablaufsteuerung
- Extraktion und Laden von externen Quellen in das DWH (Bulk Insert)
- DML-Aktionen
 - SCD2-Handling
 - Merge/Upsert-Handling
- Monitoring / Logging



Argumente für und gegen hybrides ELT

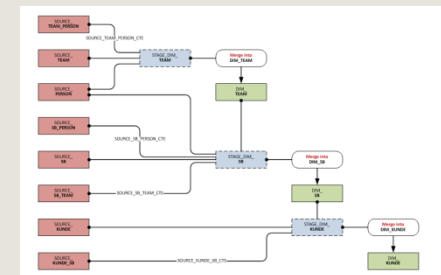
Warum das so toll ist... oder auch nicht...



- SQL-Code allgemein verständlich
- Rapid Prototyping
- ANSI-Standard
- Einflussnahme auf die Performance
- Stored Procedures



- SQL-Kenntnisse notwendig
- kein zeilenbasiertes ETL
- „Magic Numbers“
- Dokumentation (Data Lineage) manuell



Zurück zur Lösung

Voraussetzungen

Gültigkeitsdatum:

- beide Datumsspalten für Zuordnung zu Fakten notwendig
- einheitliches Start-Datum der ersten Versionen (01.01.2000)
- einheitliches End-Datum der gültigen Version (31.12.9999)
- keine Lücken
- keine Überlappungen

In der Demo werden die Gültigkeitsregeln im SQL-Skript sichergestellt!

SQL-Dialekt:

- Common Table Expressions (CTE)
- Greatest() und Least()-Funktionen
- Case When Then Else bzw. IIF oder ähnliches Konstrukt

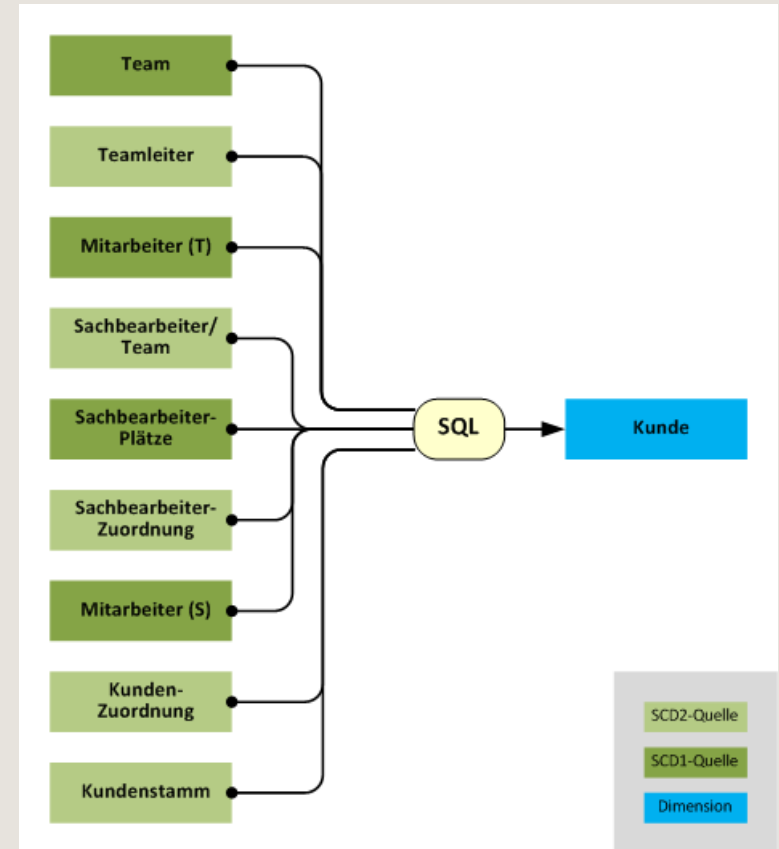
weiter mit der Lösung

Variante 1

Eine einzige Ergebnis-Tabelle

- Ein Transformationsschritt (ein SQL-Statement)
- Referenzierung nur auf Version des Kunden möglich
- keine Referenzierung auf Version von SB-Platz oder Team möglich

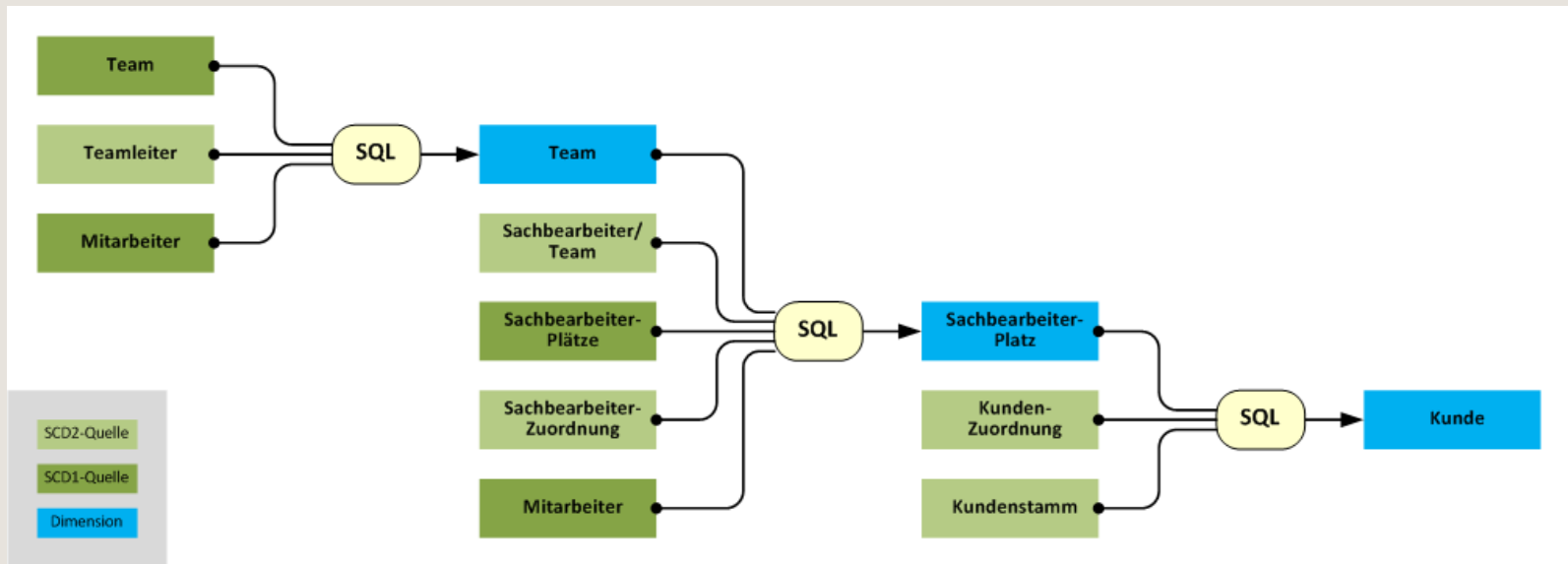
Gute Lösung, wenn keine Hierarchien zwischen den Dimensionen bestehen!



weiter mit der Lösung

Variante 2

- Drei Transformationsschritte (drei SQL-Statements)
- jeweils einzelne versionierte Tabellen für Team und SB-Platz
- Referenzierung in der Kunden-Tabelle auf deren SKs
- Star-Schema-Tabelle als View

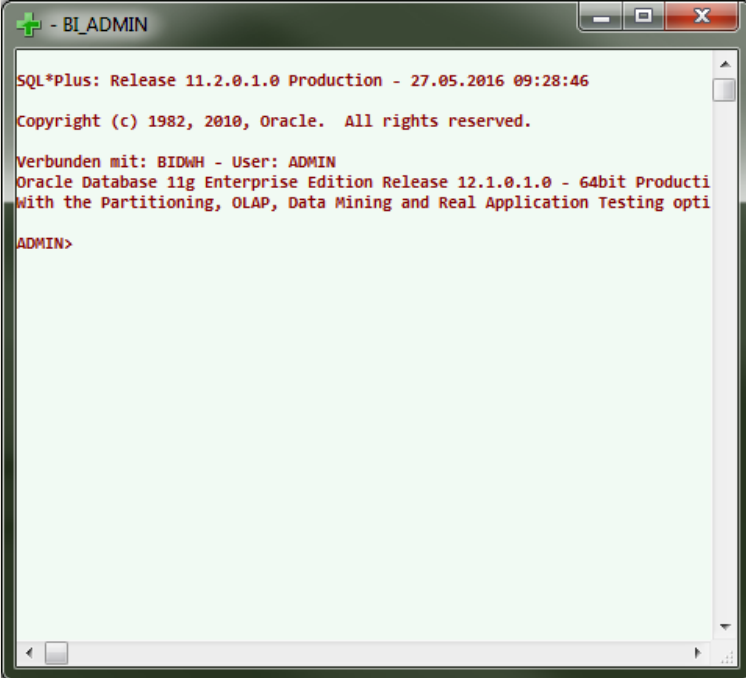


Demonstration der Lösung

in Variante 2

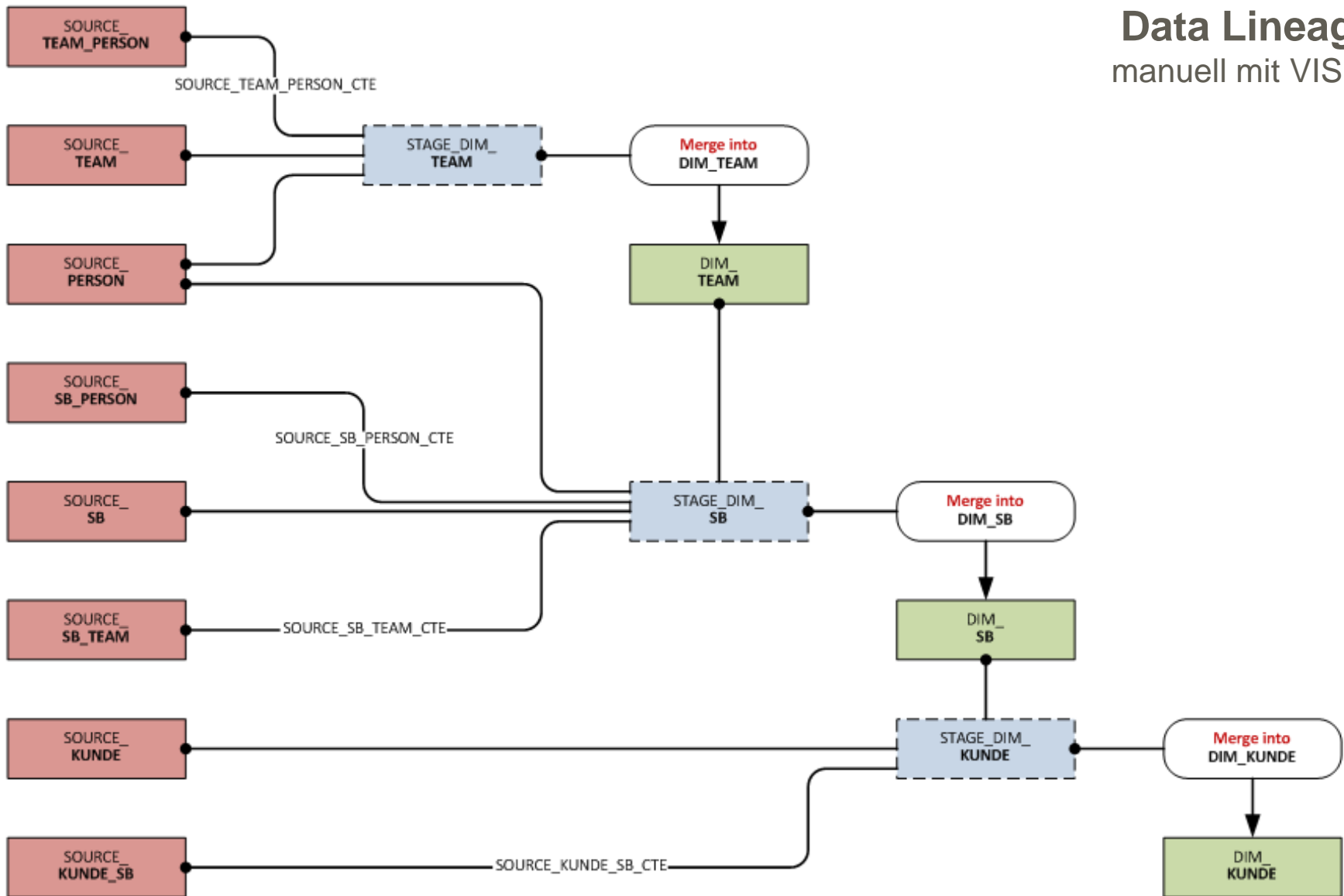
- reines SQL (Oracle)
- kein ETL/ELT-Tool im Beispiel
- alle Tabellen/Views im Schema „TEST“
- Demo im SQL*Plus

Die gezeigte Lösung funktioniert übrigens auch mit Parent-Child-Hierarchien. Dabei muss allerdings rekursives SQL (mit verschachtelten CTEs) angewendet werden!




```
BI_ADMIN
SQL*Plus: Release 11.2.0.1.0 Production - 27.05.2016 09:28:46
Copyright (c) 1982, 2010, Oracle. All rights reserved.
Verbunden mit: BIDWH - User: ADMIN
Oracle Database 11g Enterprise Edition Release 12.1.0.1.0 - 64bit Producti
With the Partitioning, OLAP, Data Mining and Real Application Testing opti
ADMIN>
```

Data Lineage manuell mit VISIO



Probleme und Lösungen?

Wie bekommt man seine Versionen sauber?

Wenn...?	...dann!
Komponente ohne Versionen	eine Version mit dem ersten Gültig-ab- und finalen Gültig-bis-Datum erzeugen
Lücken oder Überlappungen	eine der Datumsspalten als „führend“ bestimmen, danach die fehlende Spalte in SQL errechnen
nur eine der beiden Datumsspalten vorhanden	fehlende Spalte in SQL errechnen (Sub-Select)
erste Parent-Version ist erst später gültig als erste Child-Version	einheitliches erstes Gültig-ab- und finales Gültig-bis-Datum verwenden
„Late arriving Facts“	bei der jeweils letzten Version immer das finale Gültig-bis-Datum verwenden
Gültigkeiten ändern sich nachträglich und rückwirkend (vor allem bei Zuordnungen) 	a) bitemporale Datenhaltung in Dimensionen und Fakten (noch nicht ausprobiert) b) Dimension neu aufbauen, Fakten neu schlüsseln (Informationsverlust)

...noch Antworten?