

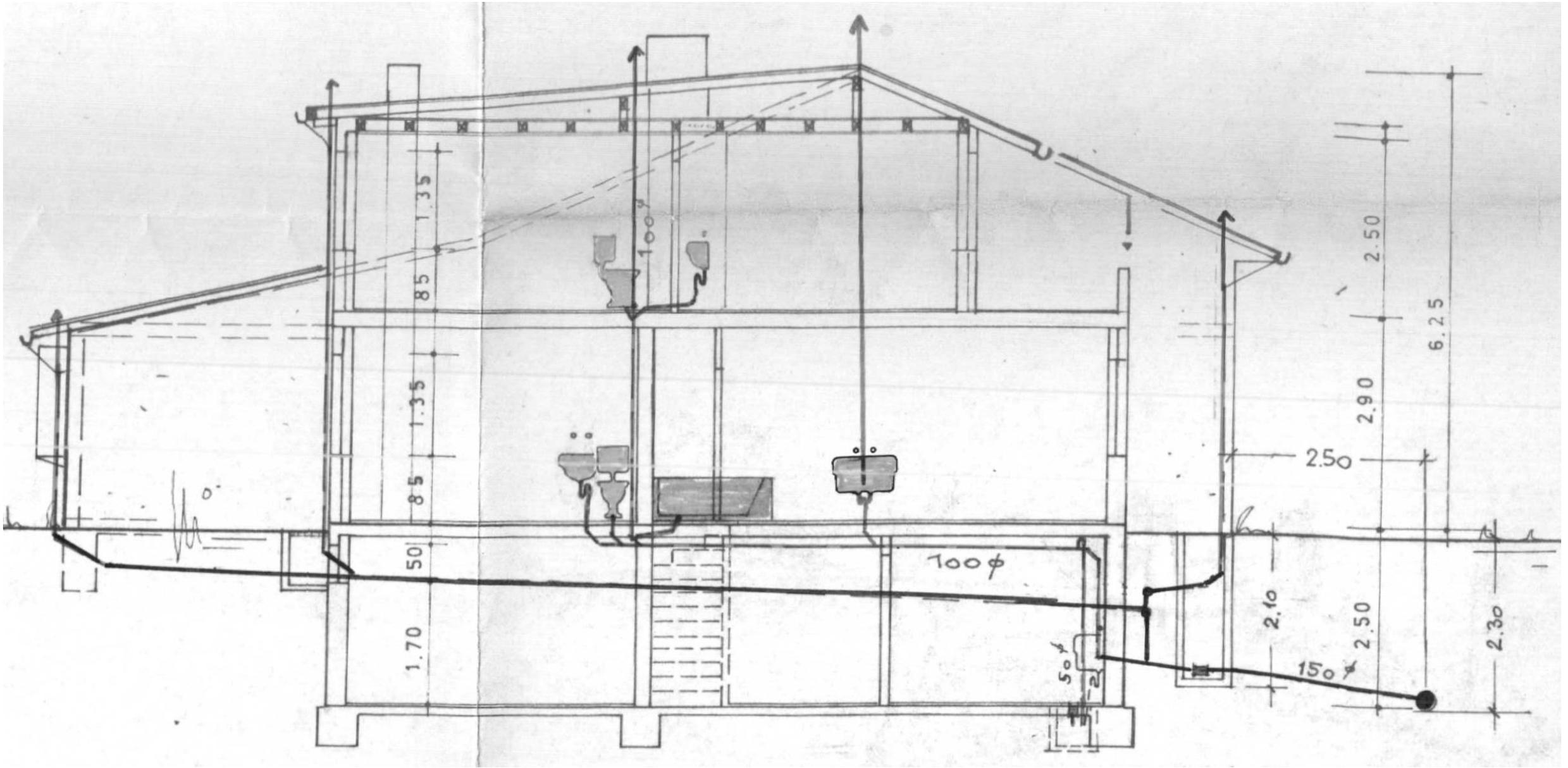
Data Vault Pipeline Description

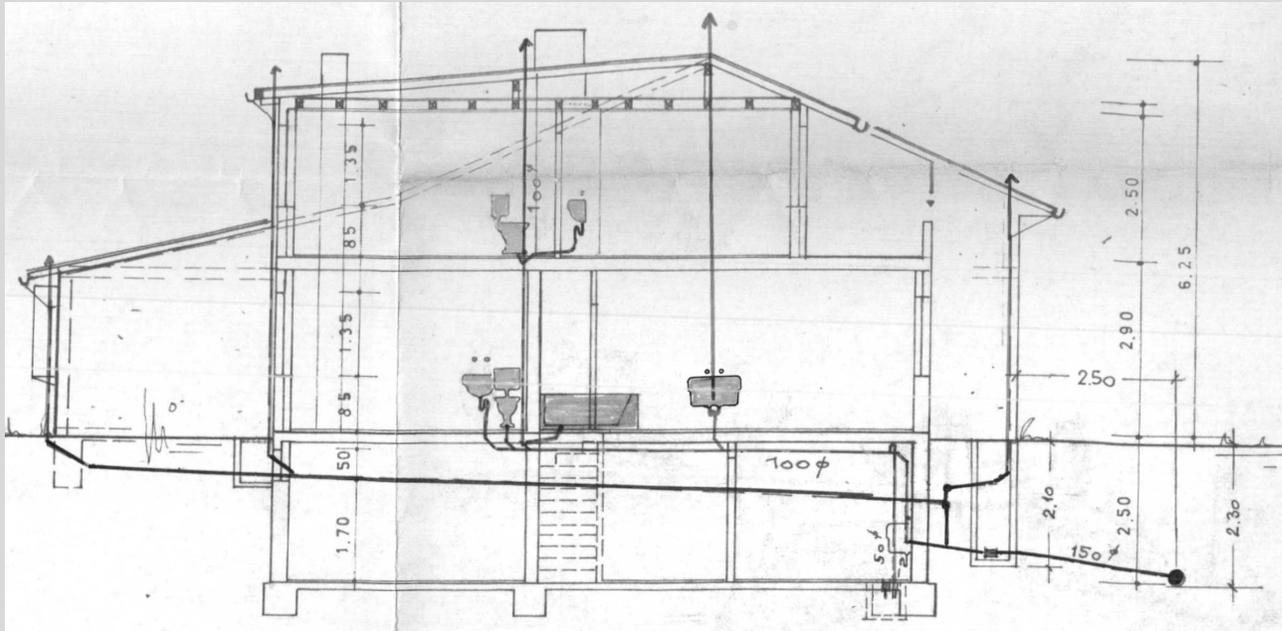
Dokumentformat zur Entkopplung von
DWH Entwicklungswerkzeugen

Matthias Wegner, cimt ag









Motivation

Konkreter Blick

Demo

Designprinzip

Github Projekt

Nutzung



Motivation



Der Data Warehouse Entwicklungsprozess

Anforderung

- Ziel & Informationsbedarf
- Herleitungs/ Transformations- beschreibung
- Datenherkunft

Modellierung

- Profiling der Quelldaten
- RAW Vault Modell
- Mart Modell
- Business Vault Modell

Implementierung

- Technische Anbindung
- Inkrementverwaltung
- Fetch
- Parse
- Stage
- Load
- Export/Send

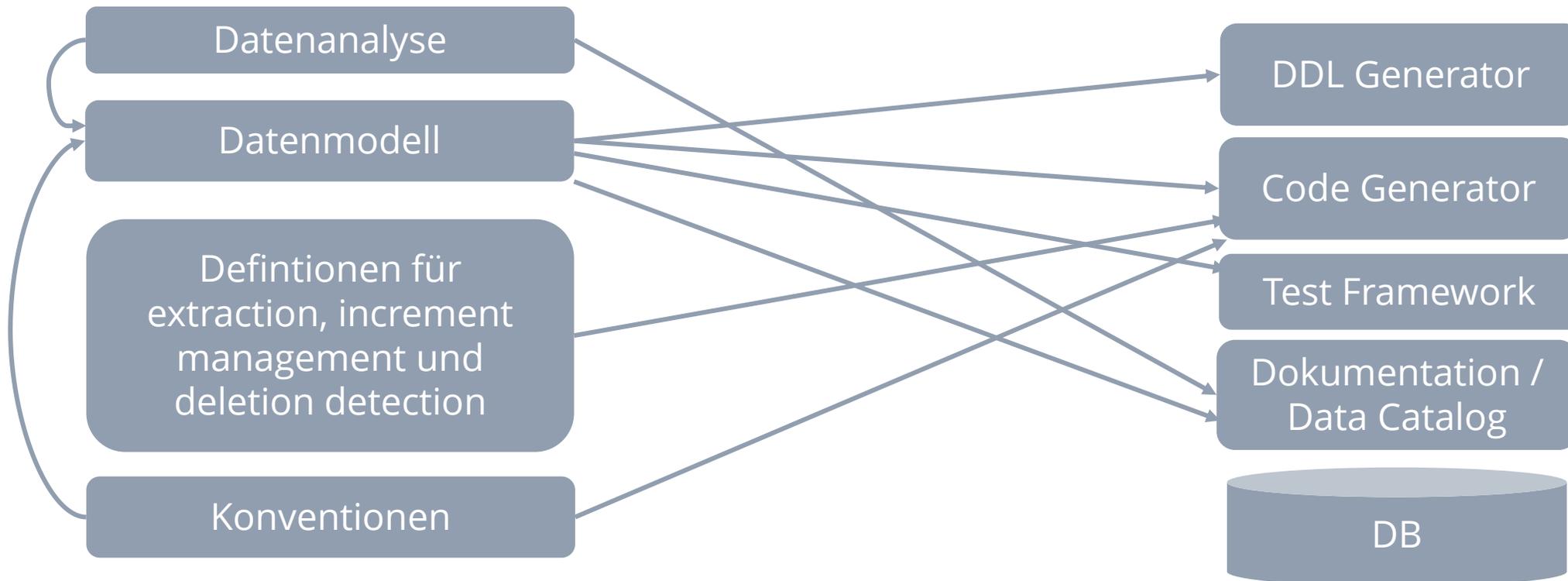
Betrieb

- Deployment
 - DB Objekte
 - Pipeline/Job
- Scheduling
- Monitoring
- Störungs- management

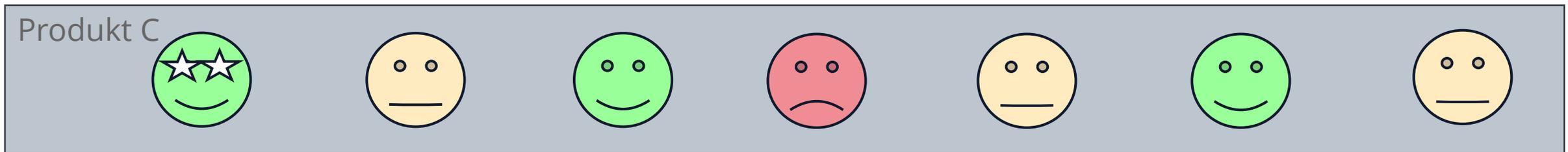
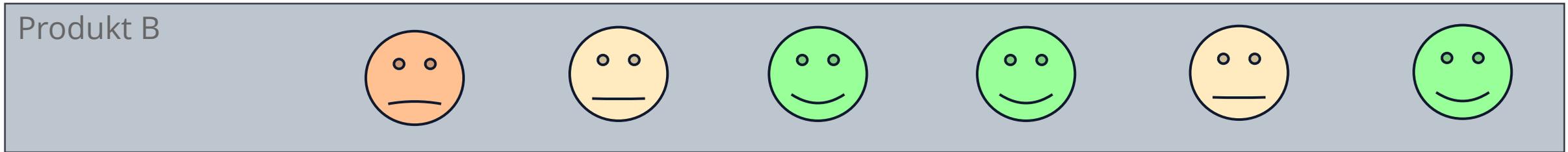
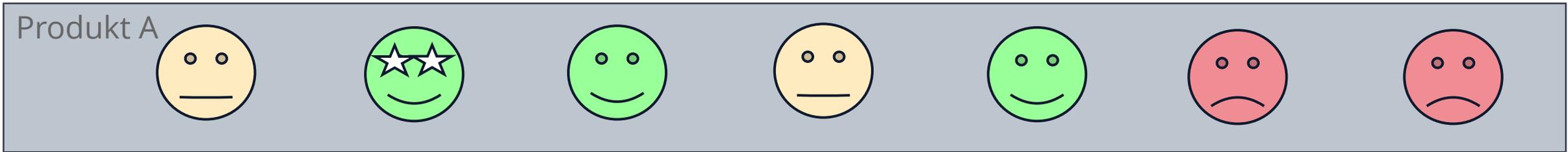
Werkzeugkasten der Entwicklungsteams

Spezifikation & Modellierung

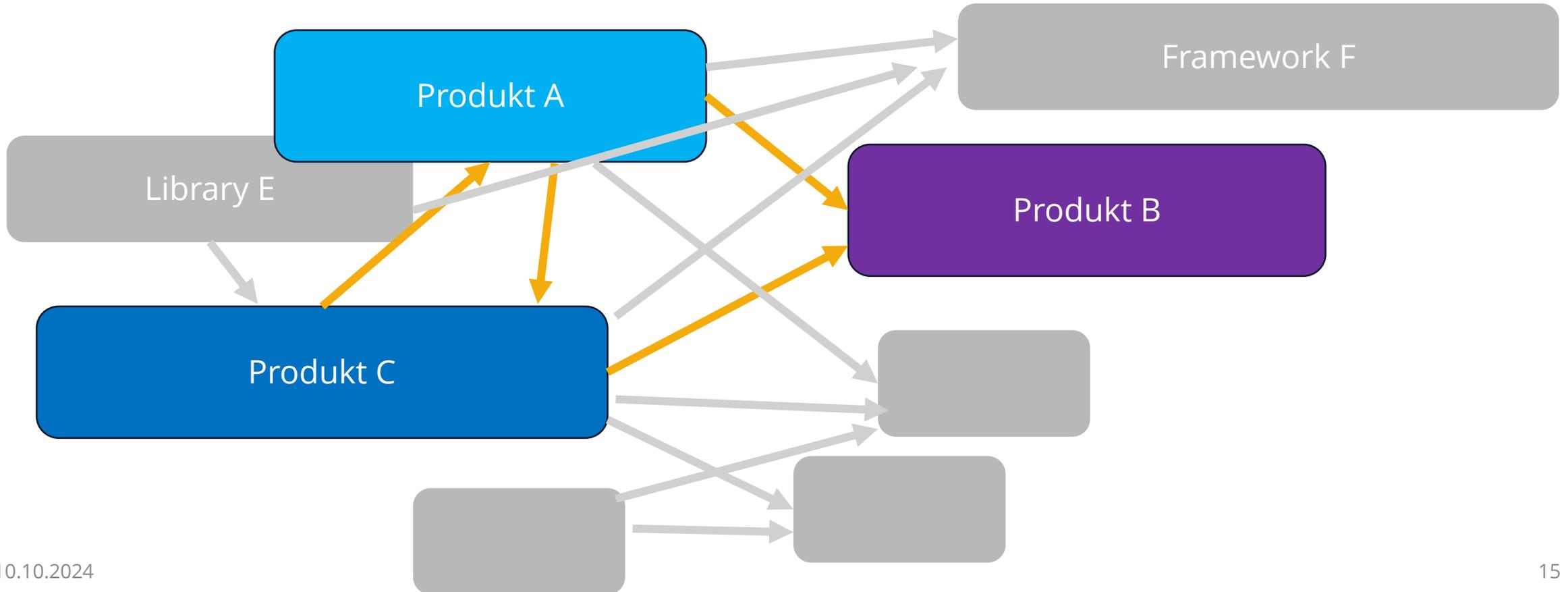
Implementierung



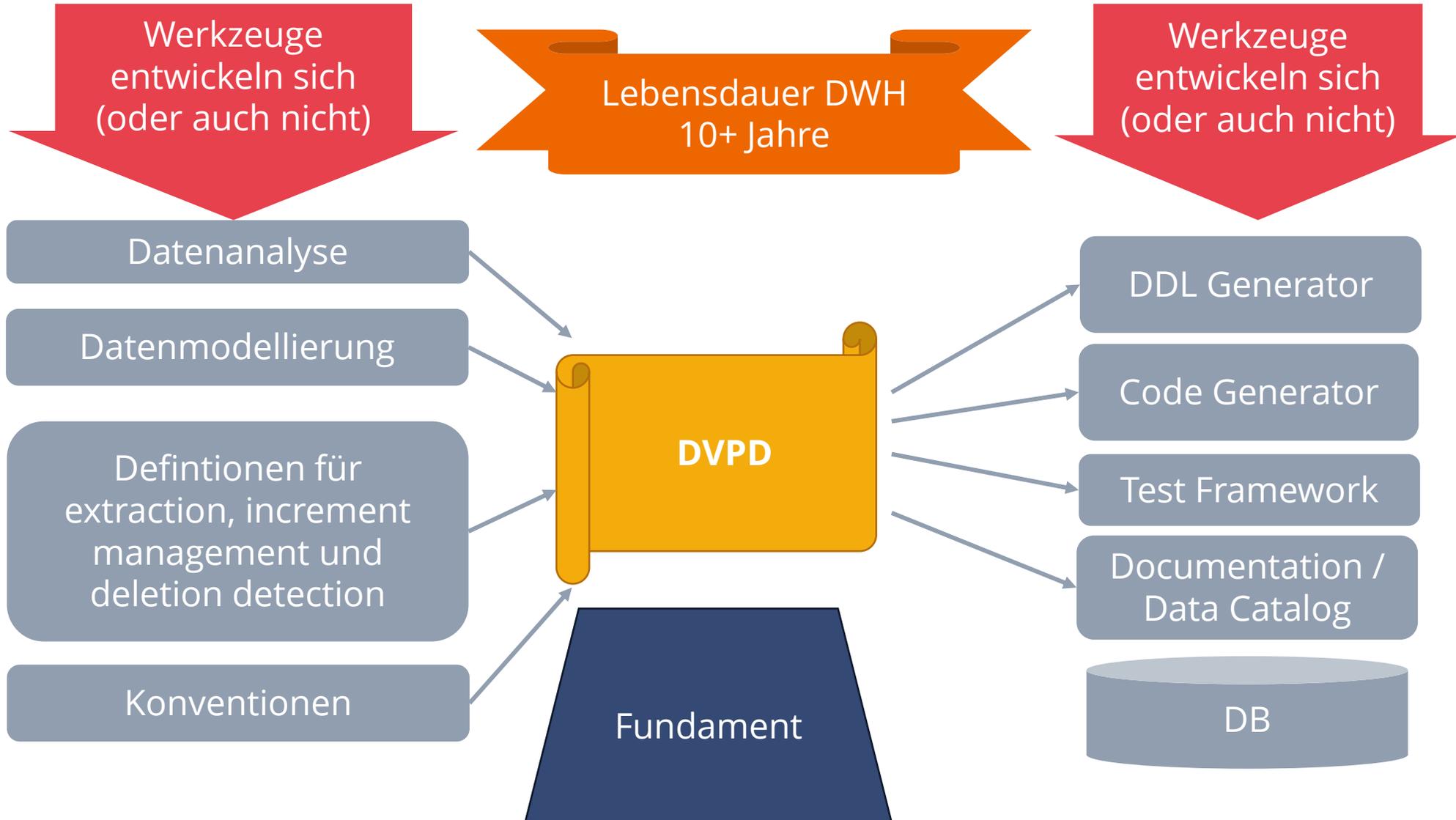
„Data Vault Produkte“ - Funktionsabdeckung



Produktzusammenstellung



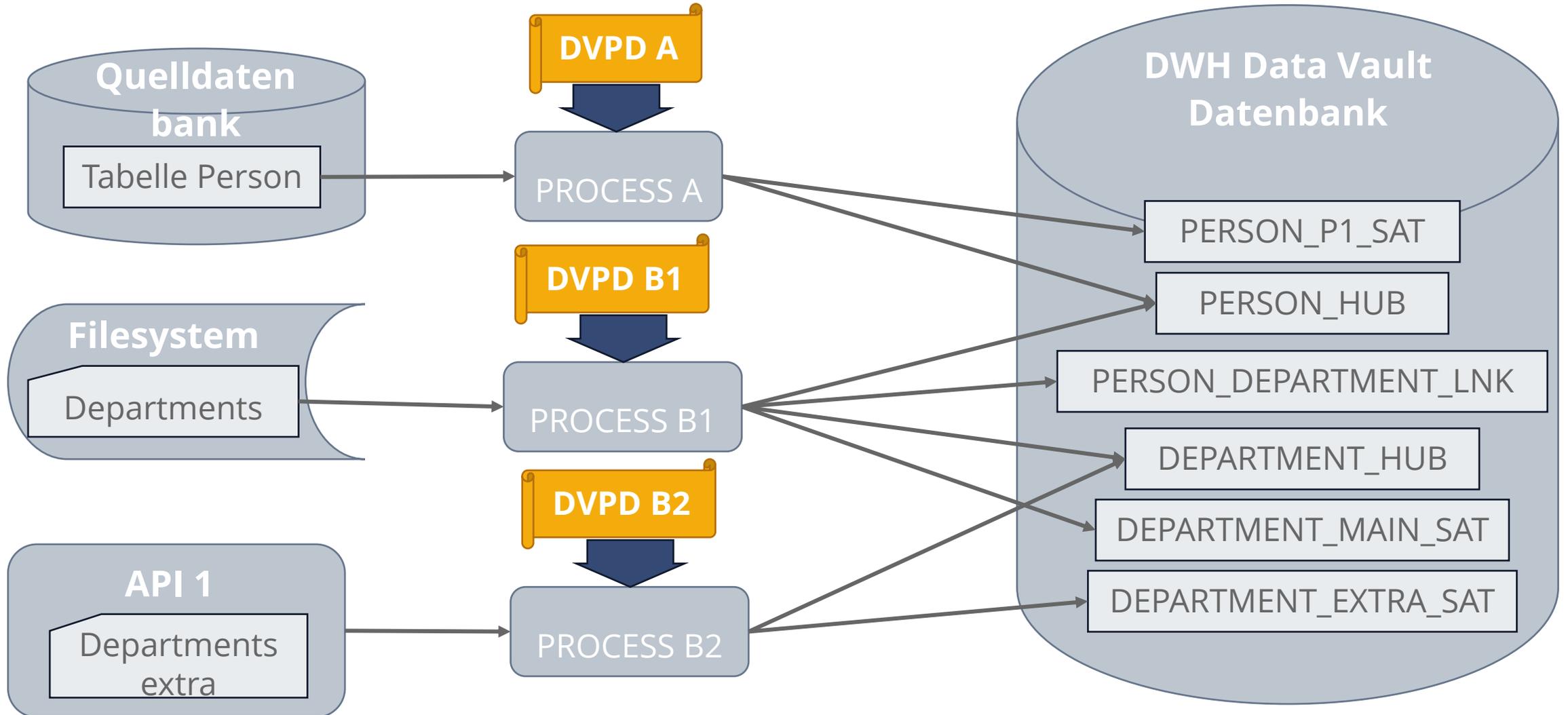
Data Warehouses leben lange





Wie sieht eine DVVPD konkret aus ?

Eine DVPD beschreibt den konkreten physischen Ladeprozess



Einfache Modellbeschreibung ...

PERSONS

company
person id
name
job

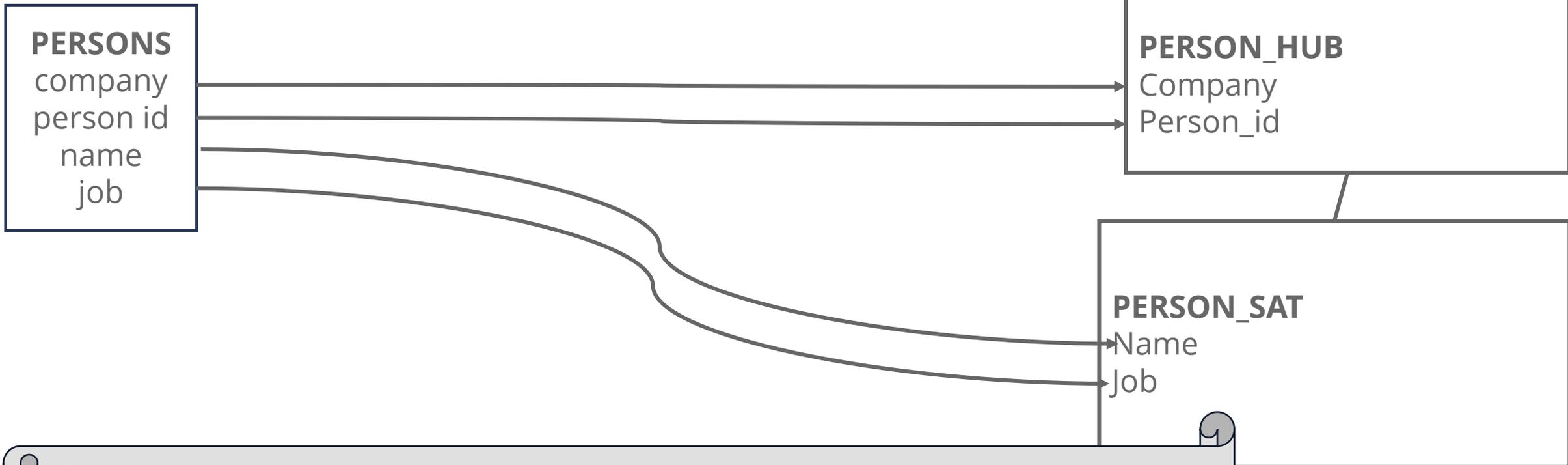
```
"fields":[  
  {"field_name":"company", "field_type":"varchar(5)"},  
  {"field_name":"person_id", "field_type":"varchar(10)"},  
  {"field_name":"name", "field_type":"varchar(250)"},  
  {"field_name":"job", "field_type":"varchar(250)"}  
]
```

PERSON_HUB

PERSON_SAT

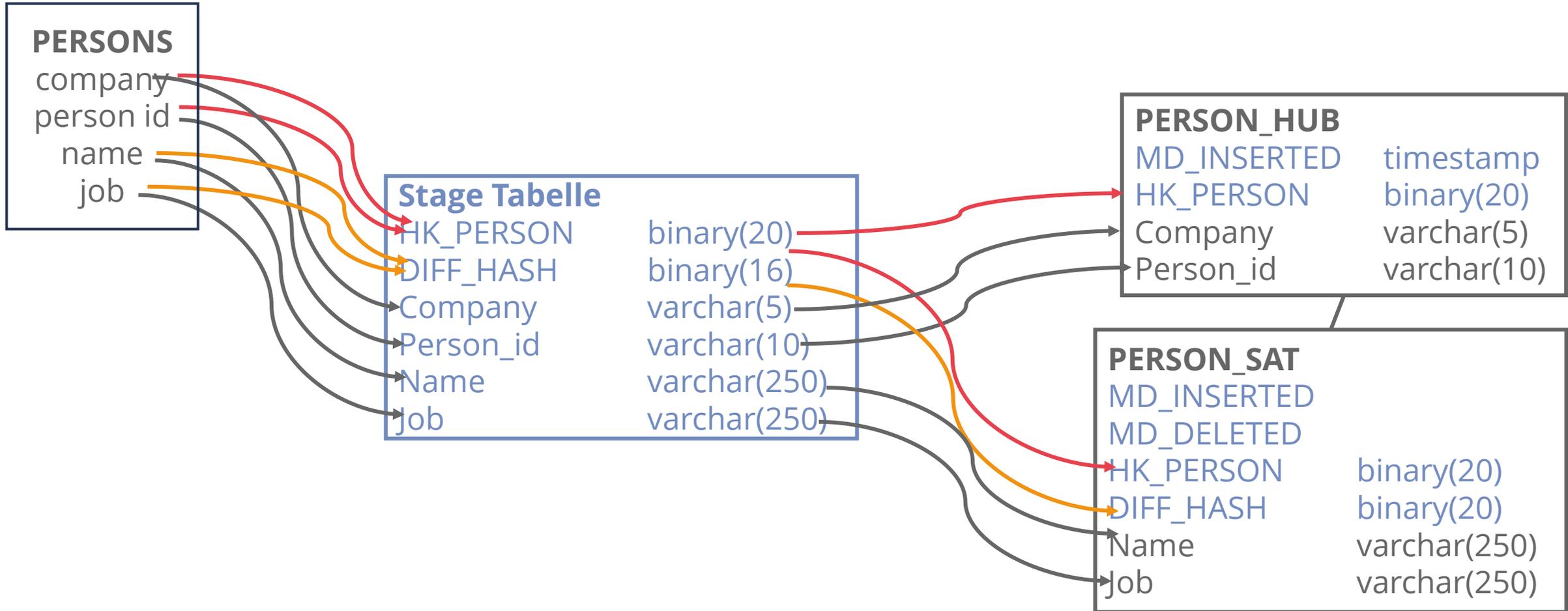
```
"data_vault_model":  
  {"table_name":"person_hub ", "table_stereotype":"hub"},  
  {"table_name":"person_sat ", "table_stereotype":"sat",  
   "satellite_parent_table":"person_hub"},  
  ]
```

Einfache Mappingbeschreibung



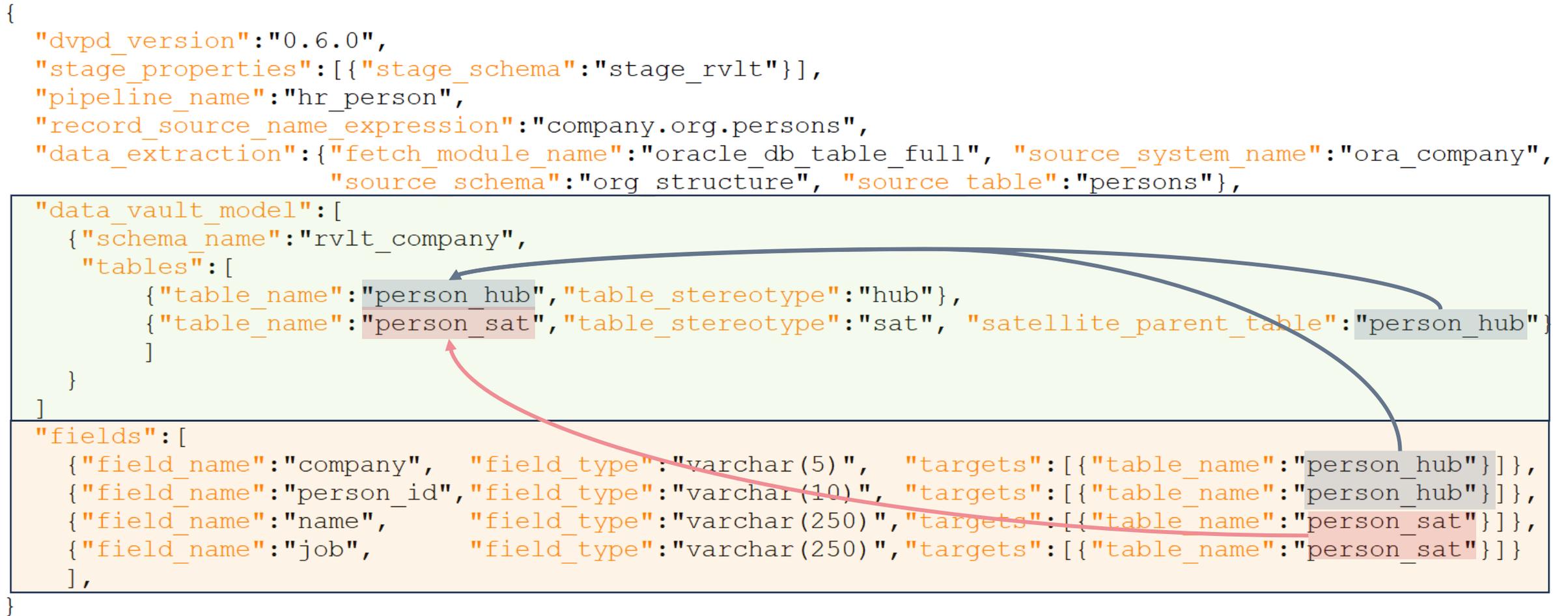
```
"fields":[
  {"field_name":"company", "field_type":"varchar(5)","targets": [{"table_name":"person_hub"}]},
  {"field_name":"person_id", "field_type":"varchar(10)","targets": [{"table_name":"person_hub"}]},
  {"field_name":"name", "field_type":"varchar(250)","targets": [{"table_name":"person_sat"}]},
  {"field_name":"job", "field_type":"varchar(250)","targets": [{"table_name":"person_sat"}]}
]
```

Aus der einfachen Beschreibung prinzipbedingte Schlussfolgerungen ziehen



JSON Syntax

```
{
  "dvpd_version": "0.6.0",
  "stage_properties": [{"stage_schema": "stage_rvlt"}],
  "pipeline_name": "hr_person",
  "record_source_name_expression": "company.org.persons",
  "data_extraction": {"fetch_module_name": "oracle_db_table_full", "source_system_name": "ora_company",
    "source_schema": "org structure", "source_table": "persons"},
  "data_vault_model": [
    {"schema_name": "rvlt_company",
      "tables": [
        {"table_name": "person_hub", "table_stereotype": "hub"},
        {"table_name": "person_sat", "table_stereotype": "sat", "satellite_parent_table": "person_hub"}
      ]
    }
  ],
  "fields": [
    {"field_name": "company", "field_type": "varchar(5)", "targets": [{"table_name": "person_hub"}]},
    {"field_name": "person_id", "field_type": "varchar(10)", "targets": [{"table_name": "person_hub"}]},
    {"field_name": "name", "field_type": "varchar(250)", "targets": [{"table_name": "person_sat"}]},
    {"field_name": "job", "field_type": "varchar(250)", "targets": [{"table_name": "person_sat"}]}
  ],
}
```



Alles (!) was eine Beladung beschreibt

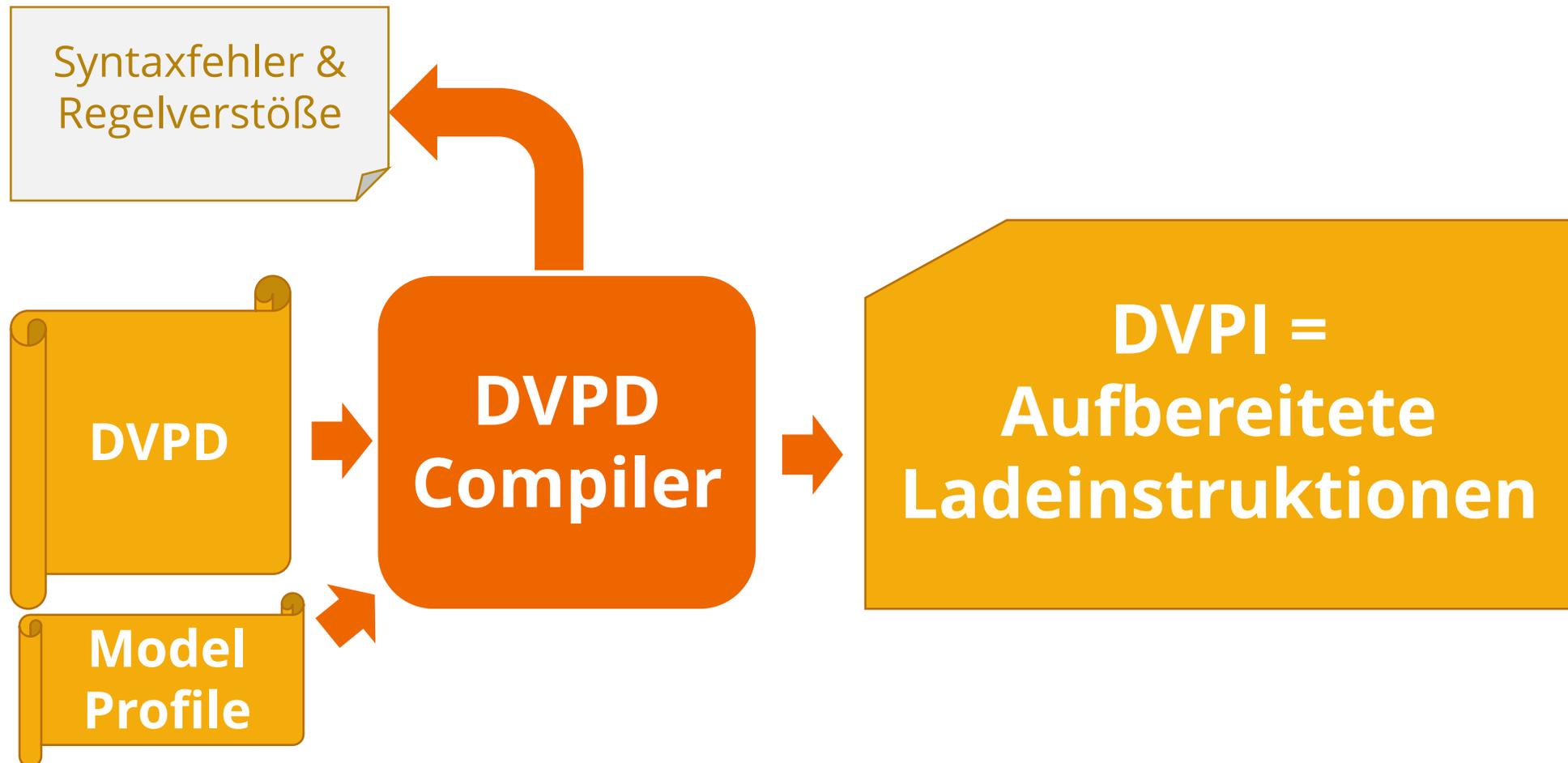
Von Quelle bestimmt

- Abrufverfahren und Transport aus der Quelle
- Inkrementmethode
- Erkennung von gelöschten Daten
- Datenformat und Zerlegung in Felder (Parsing)

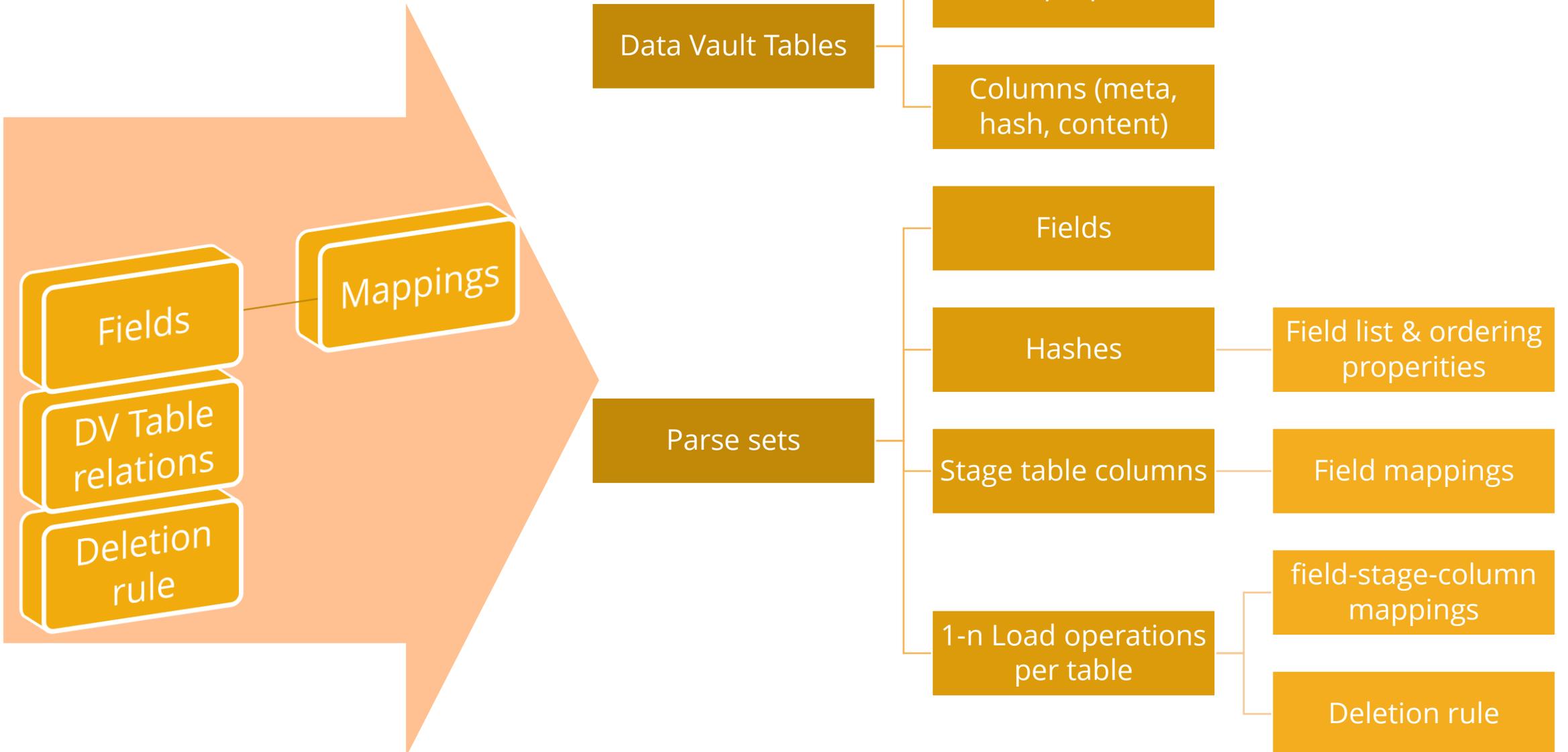
Von Data Vault Methode
und Zielplattform
bestimmt

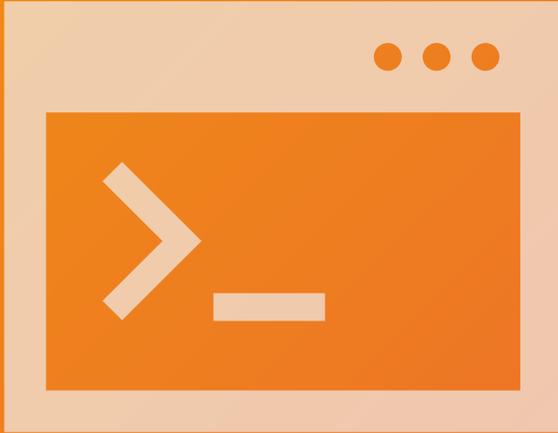
- Zieldatenmodell
- Mapping

Compiler und Data Vault Pipeline Instruction (DVPI)



Struktur der DVPI

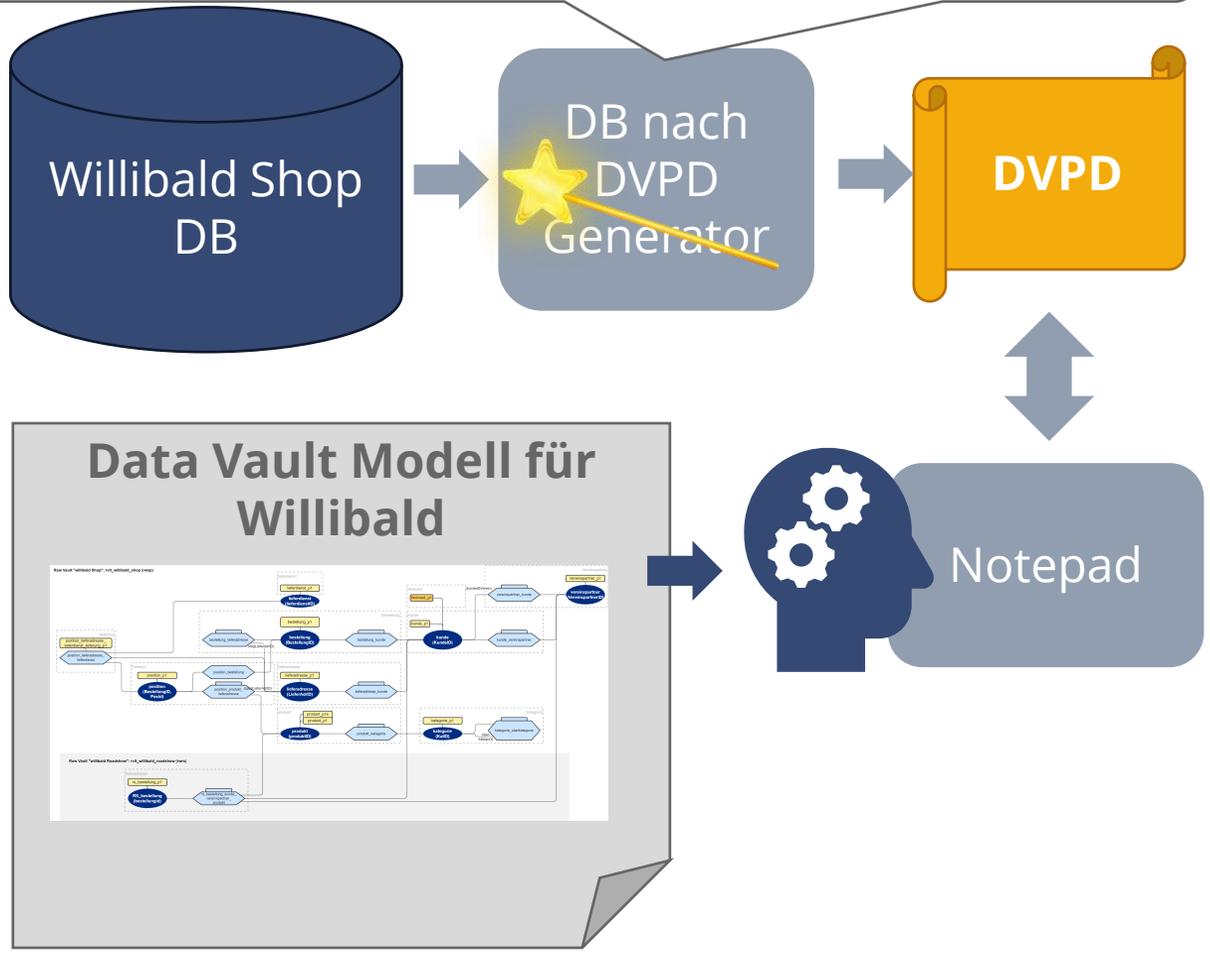




Demonstration

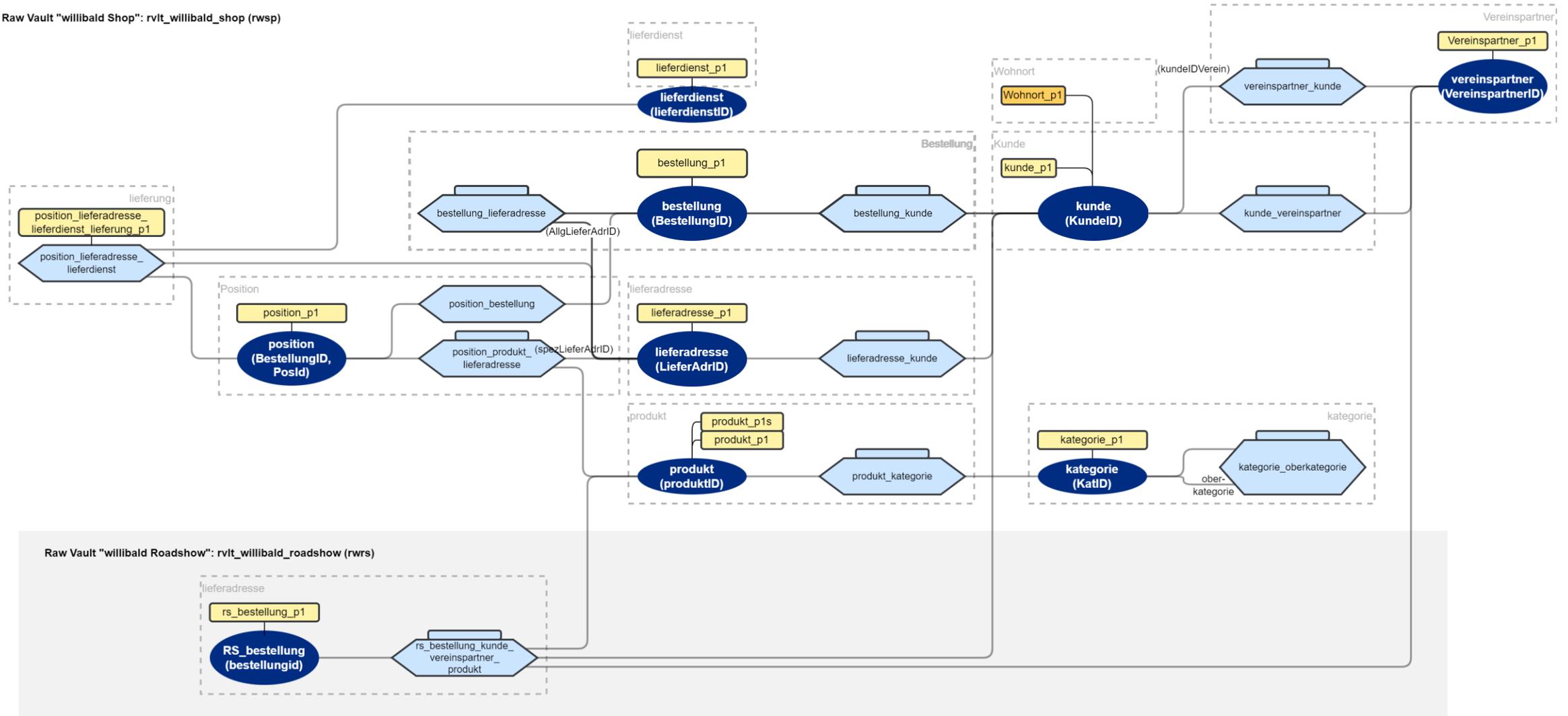
„Generieren einer DVPD aus der Quellstruktur“

```
dvpd_generate_from_db willibald_shop_p1 bestellung
```

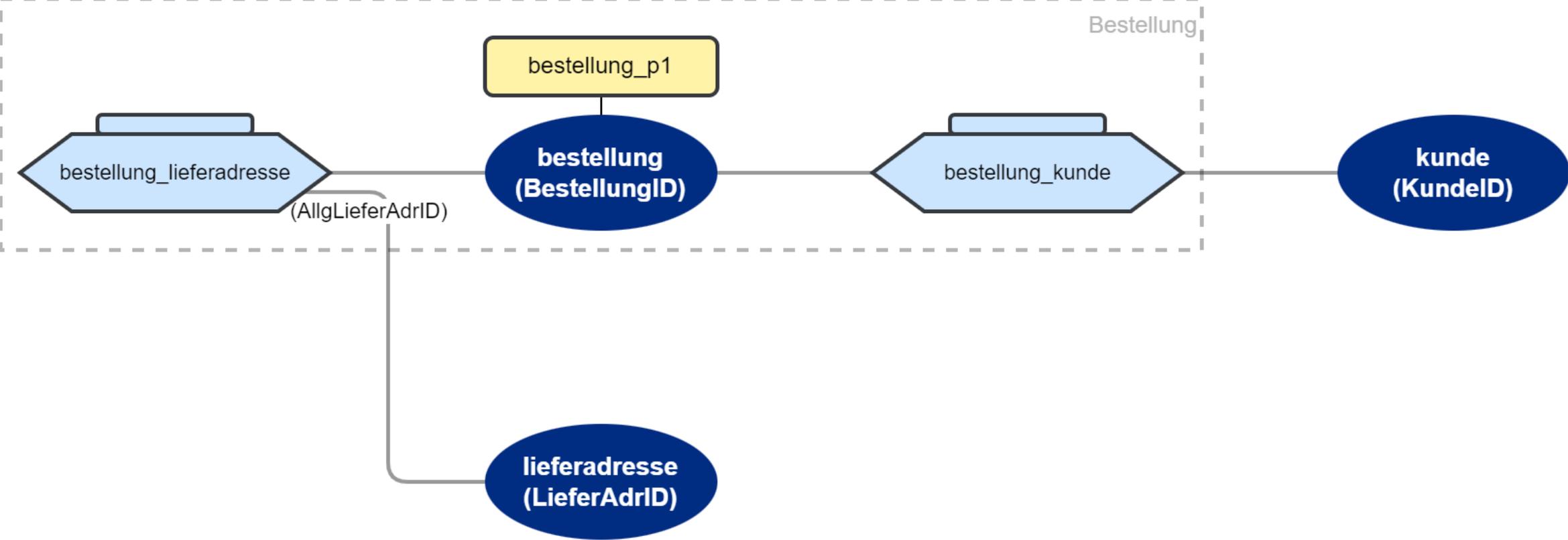


Modell der Pipeline

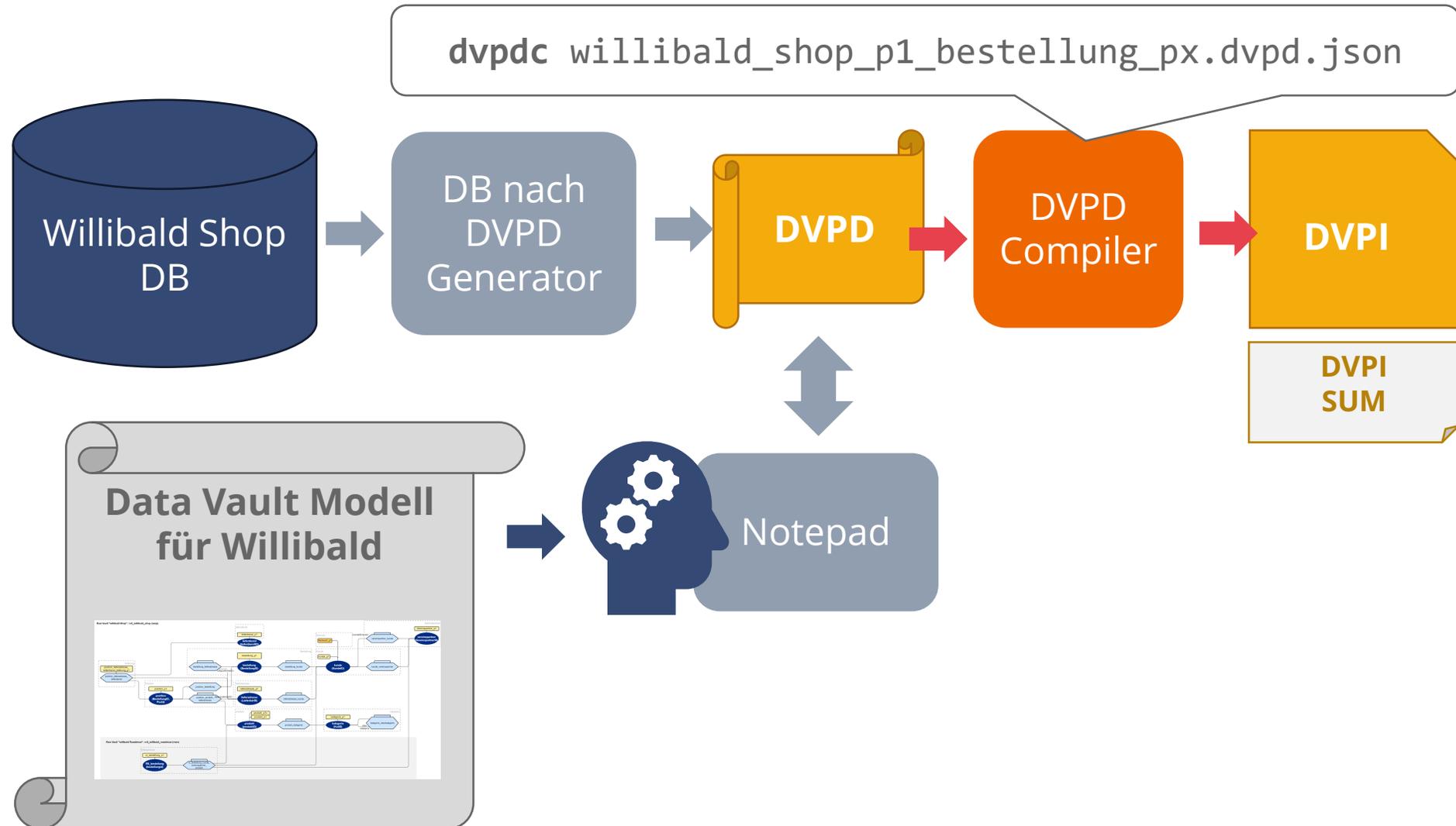
Raw Vault "willibald Shop": rvit_willibald_shop (rwsp)



Modell der Pipeline



Prüfung der DVPD mit dem Compiler



Demo „Ergebnisse aus der DVPD“

dvpd_ddl_render
willibald_shop_p1_bestellung_px.dvpi.json

DDL
Generator

DDL Script für
Kunde_hub
Bestellung_sat für
Bestellung_hub

dvpd_devsheet_render
willibald_shop_p1_bestellung_px.dvpi.json



„Dev Sheet
Generator“

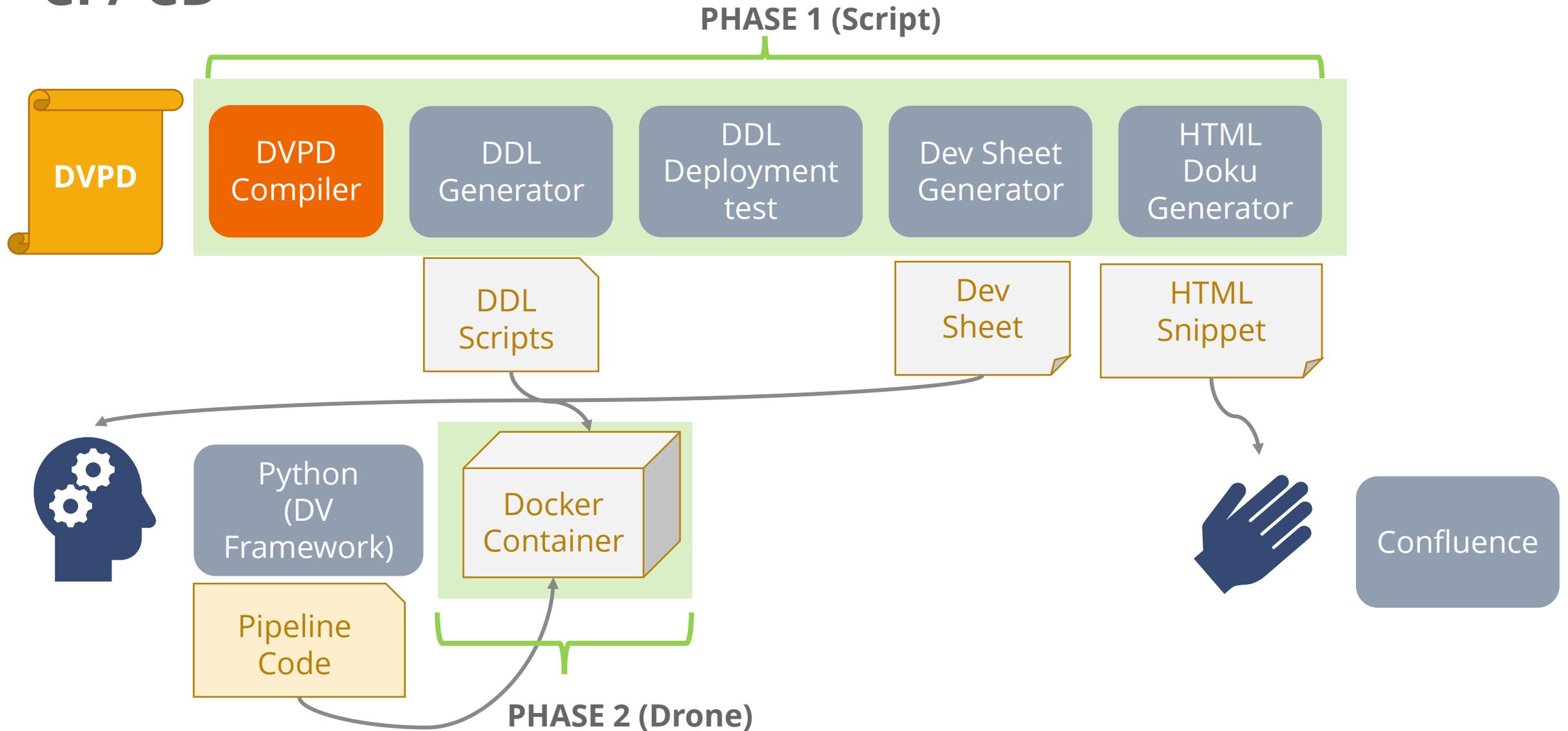
Doku für Entwickler

dvpd_doc_render
willibald_shop_p1_bestellung_px.dvdpd.json

HTML
Doku
Generator

HTML Doku für
Confluence

CI / CD





Designprinzipien der DVVPD Syntax

„Den ganzen Weg gehen“

DVPD =



Komplette Beschreibung eines Datenladeprozesses

Vollständige Unterstützung
der Data Vault Methode

Standpunkt der Syntax ist
die Modelltransformation

DVPD =



Komplette Beschreibung
eines Datenladeprozesses

**Vollständige Unterstützung
der Data Vault Methode**

Standpunkt der Syntax ist
die Modelltransformation

DVPD =



Komplette Beschreibung
eines Datenladeprozesses

Vollständige Unterstützung
der Data Vault Methode

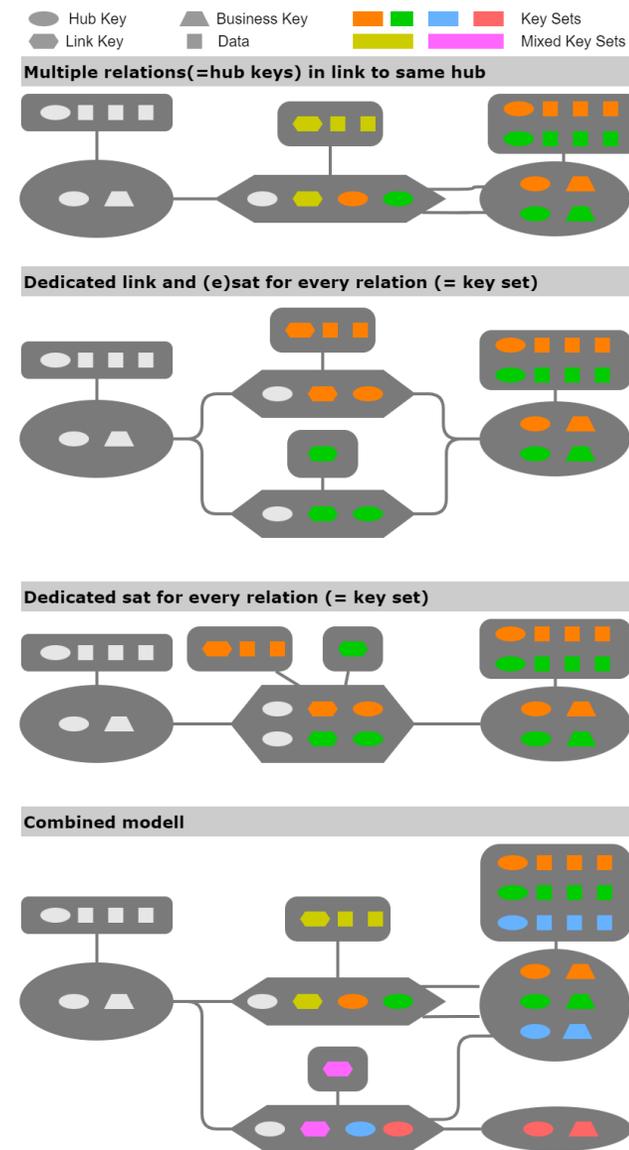
**Standpunkt der Syntax ist
die Modelltransformation**

Standpunkt der Syntax ist die Modelltransformation

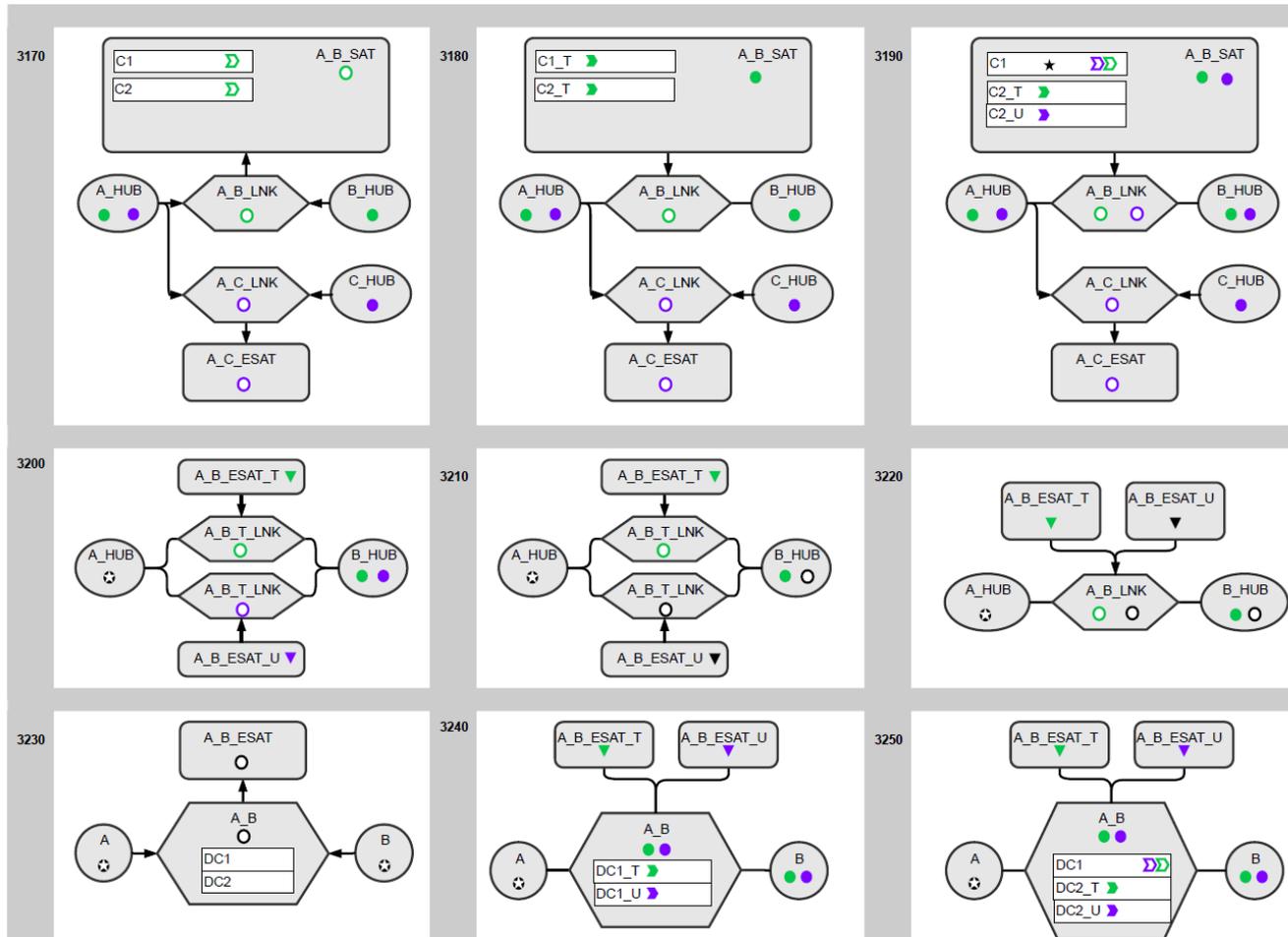


Vollständige Unterstützung der Data Vault Methode

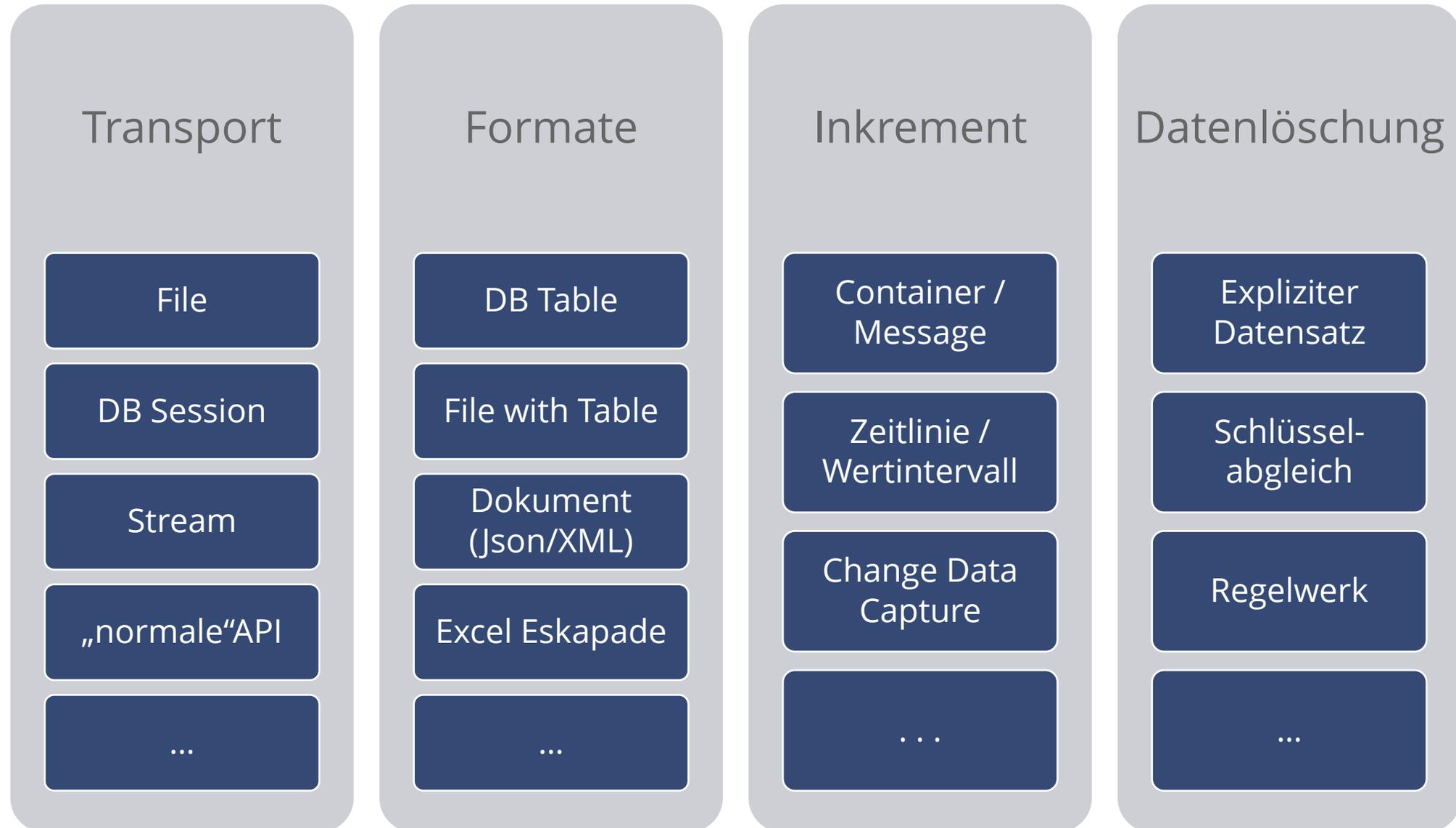
- „Building a Scalable Data Warehouse with Data Vault 2.0“ von Linstedt/Olschimke ist die Grundlage (Data_Vault_method_coverage_and_syntax_examples)
- Kombinatorische Analysen und Ableitungen
 - Varianten von Data Vault Modellen
 - Varianten von Datenmappings
 - Begrenzt auf 1 Datenquellobjekt
- Ergänzt um Erfahrungen aus Projekten und der DV Community
 - Varianten bzgl. Modelldetails



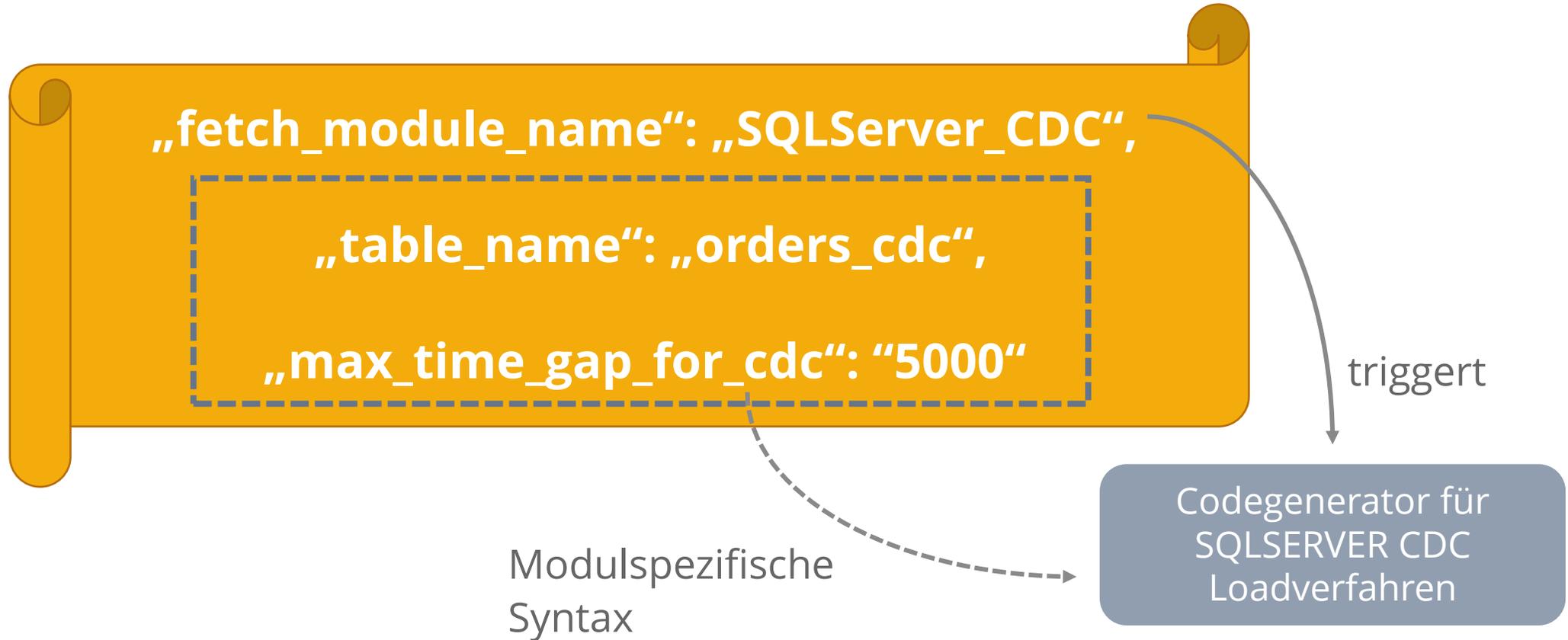
Kombinatorische Analysen (Ausschnitt)



Vielfalt der Quellformate und Verhalten



Modulprinzip + offene Syntax = Flexibilität



Das optimale Verhältnis zwischen hartem Programmcode und generischen Komponenten hängt vom Projekt ab

Quellbezogene Parameter



Abrufmodul,
hart codiert für diese
Technologie und
Inkrementmethode

```
{
  "dvpd_version":"0.6.0",
  "stage_properties":[{"stage_schema":"stage_rvlt"}],
  "pipeline_name":"hr_person",
  "record_source_name_expression":"company.org.persons",
  "data_extraction":{"fetch_module_name":"oracle_db_table_full", "source_system_name":"ora_company",
    "source_schema":"org_structure", "source_table":"persons"},
  "data_vault_model":[
    {"schema_name":"rvlt_company",
      "tables":[
        {"table_name":"person_hub","table_stereotype":"hub"},
        {"table_name":"person_sat","table_stereotype":"sat", "satellite_parent_table":"person_hub"}
      ]
    }
  ]
}
"fields":[
  {"field_name":"company", "field_type":"varchar(5)", "targets":[{"table_name":"person_hub"}]},
  {"field_name":"person_id", "field_type":"varchar(10)", "targets":[{"table_name":"person_hub"}]},
  {"field_name":"name", "field_type":"varchar(250)", "targets":[{"table_name":"person_sat"}]},
  {"field_name":"job", "field_type":"varchar(250)", "targets":[{"table_name":"person_sat"}]}
],
}
```

Parameter für das Abrufmodul

Dies sind auch Parameter für
den Datenabruf (Parsing ist bei
DB Anbindung nicht notwendig)

Warum JSON ?

weit verbreitetes Verfahren

Interpretierbar durch
Menschen und Maschinen



Das DVPD Entwicklungsprojekt

[https://github.com/cimt-ag/
data_vault_pipelinedescription](https://github.com/cimt-ag/data_vault_pipelinedescription)

Inhalte des Konzepts

Dokumentation

- Konzept
 - Anforderungsrahmen
 - Design
 - **Regelwerk**
 - Workflow Beispiele
- **Syntaxreferenz**
- Begleitende Erkenntnisse
 - Data Mapping Taxonomie
 - Deletion Detection Szenarien



DVPD
DVPI

Referenzimplementierung

- Compiler
 - Syntax- und Konsistenzprüfung
 - Ableitungsverfahren
 - Ergebnisbereitstellung
- Automatisierter Test
 - Testfälle
 - Referenzergebnisse
 - Vergleich
- Generatorbeispiele



Bisher
232



DDL
DBT(!)

Lizenz und Weiterverwendung



Dokumentation der Spezifikation

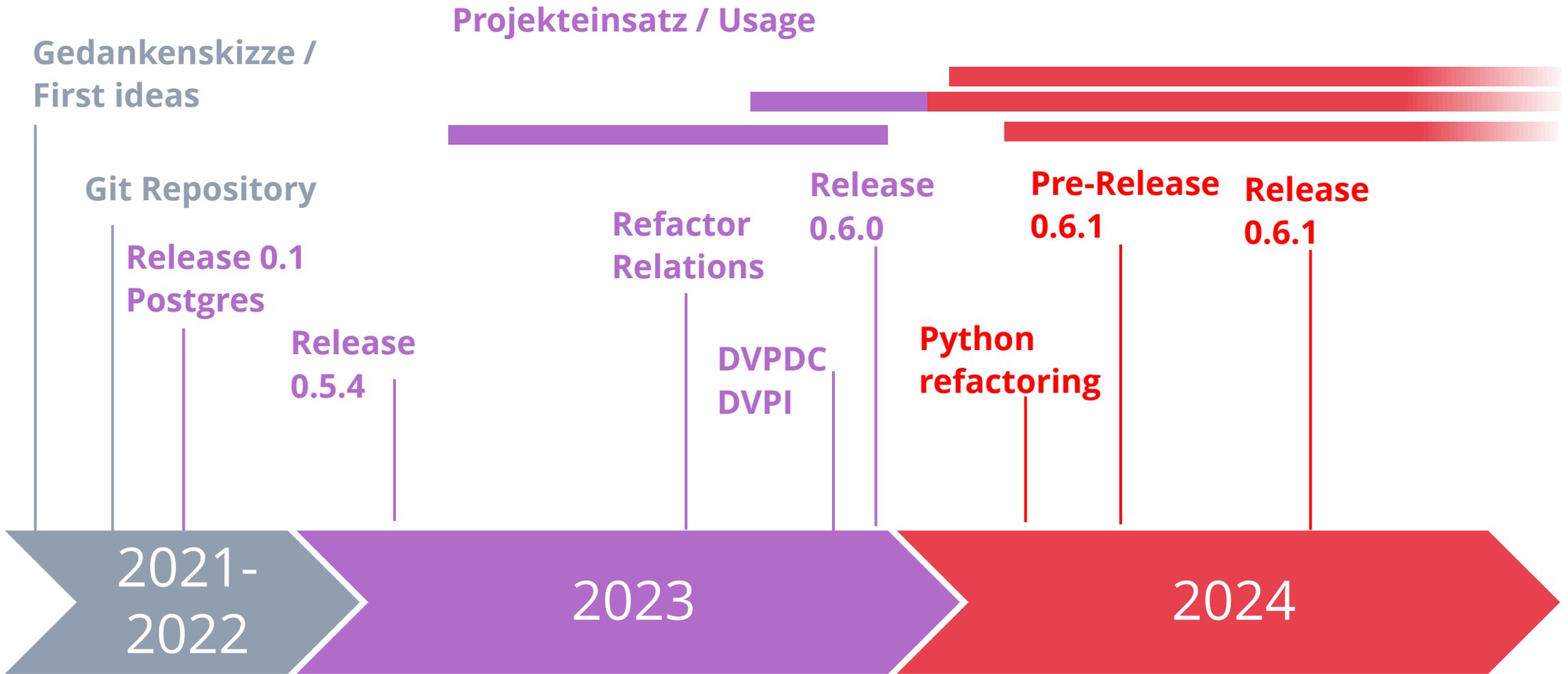
- Creative Commons BY-ND 4.0
 - Attribution
 - NoDerivatives

Referenzimplementierung

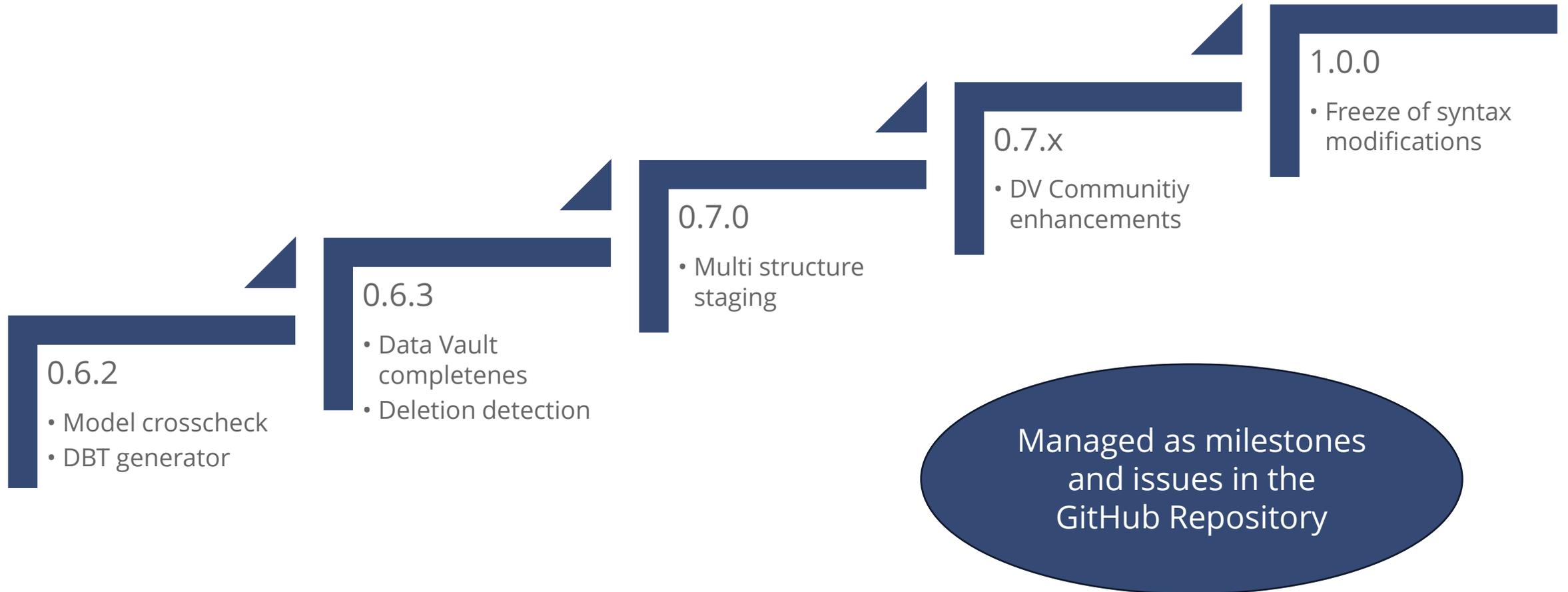
- Apache License, Version 2.0

Weiterentwicklung des Konzeptes erfolgt unter der redaktionellen Führung der cimt ag, Matthias Wegner.
Neue Versionen werden unter den gleichen Bedingungen veröffentlicht.

Historie/History



Roadmap





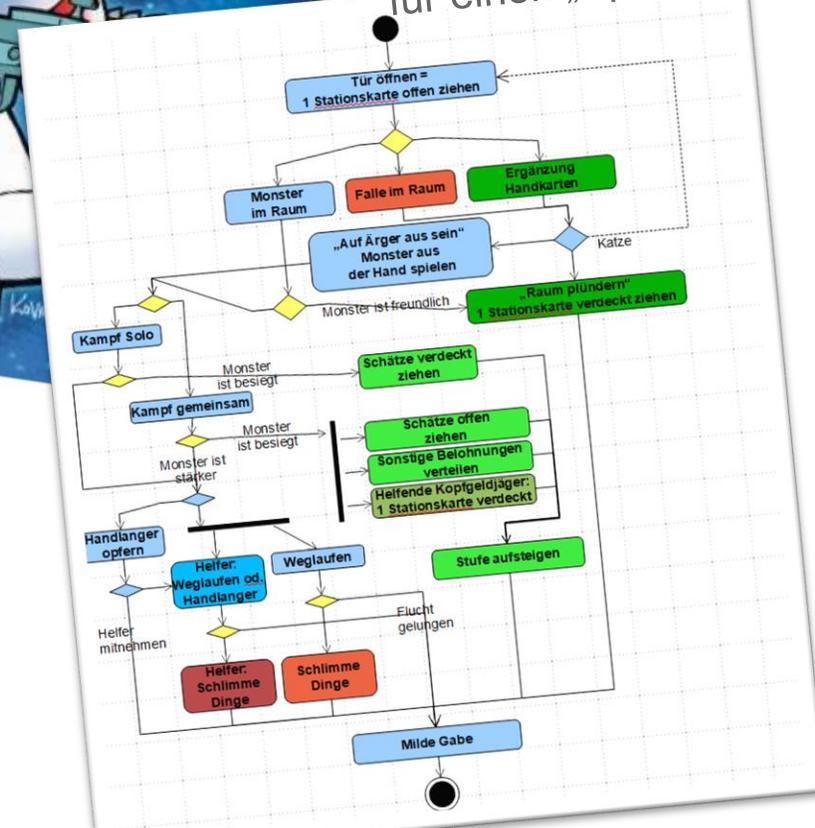
Wie kann man es nutzen ?

Erster „Vollkontakt“ mit DVPD

- Dokumentation ist ein harter Einstiegspfad
 - Komplexe Konzepte dominieren die Dokumentation, werden aber selten notwendig
- Bei der Anwendung wird's plötzlich einfach und einleuchtend
- Coaching kann den Weg über die Lernkurve erleichtern

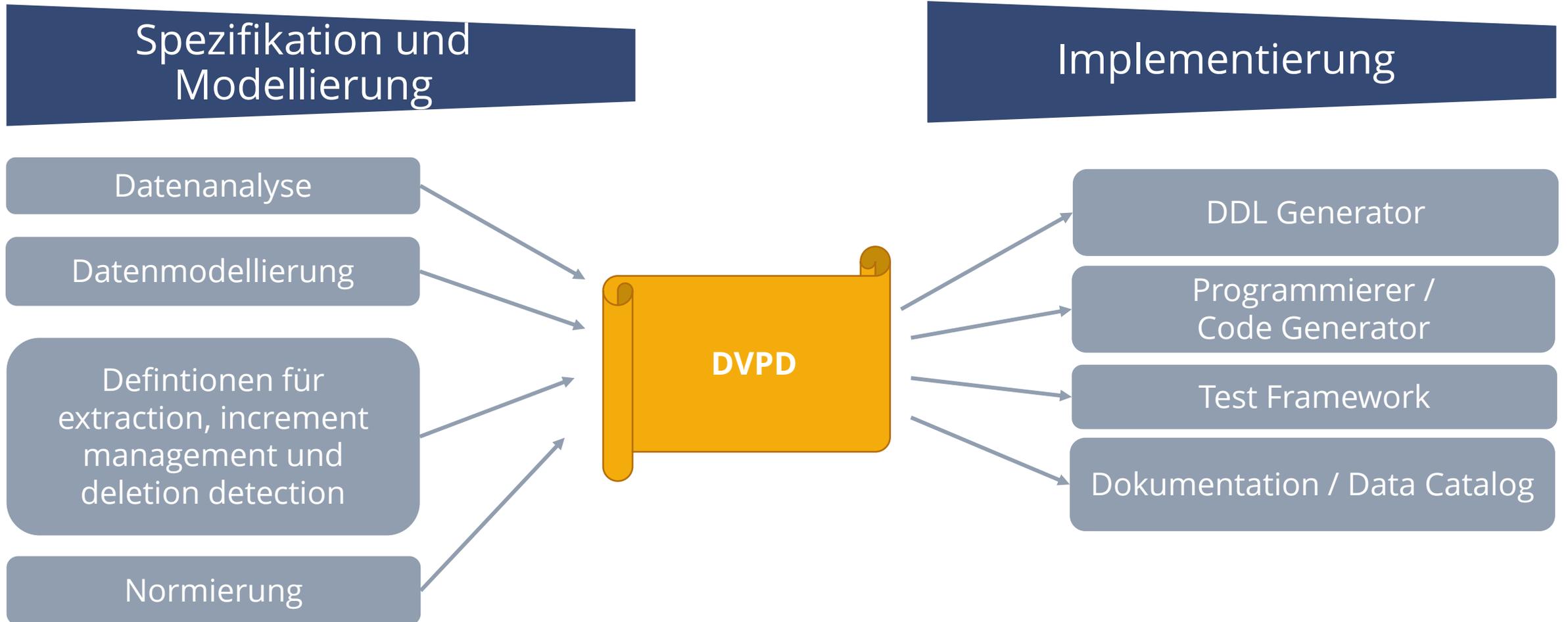


Ablaufplan für einen „Spielzug“

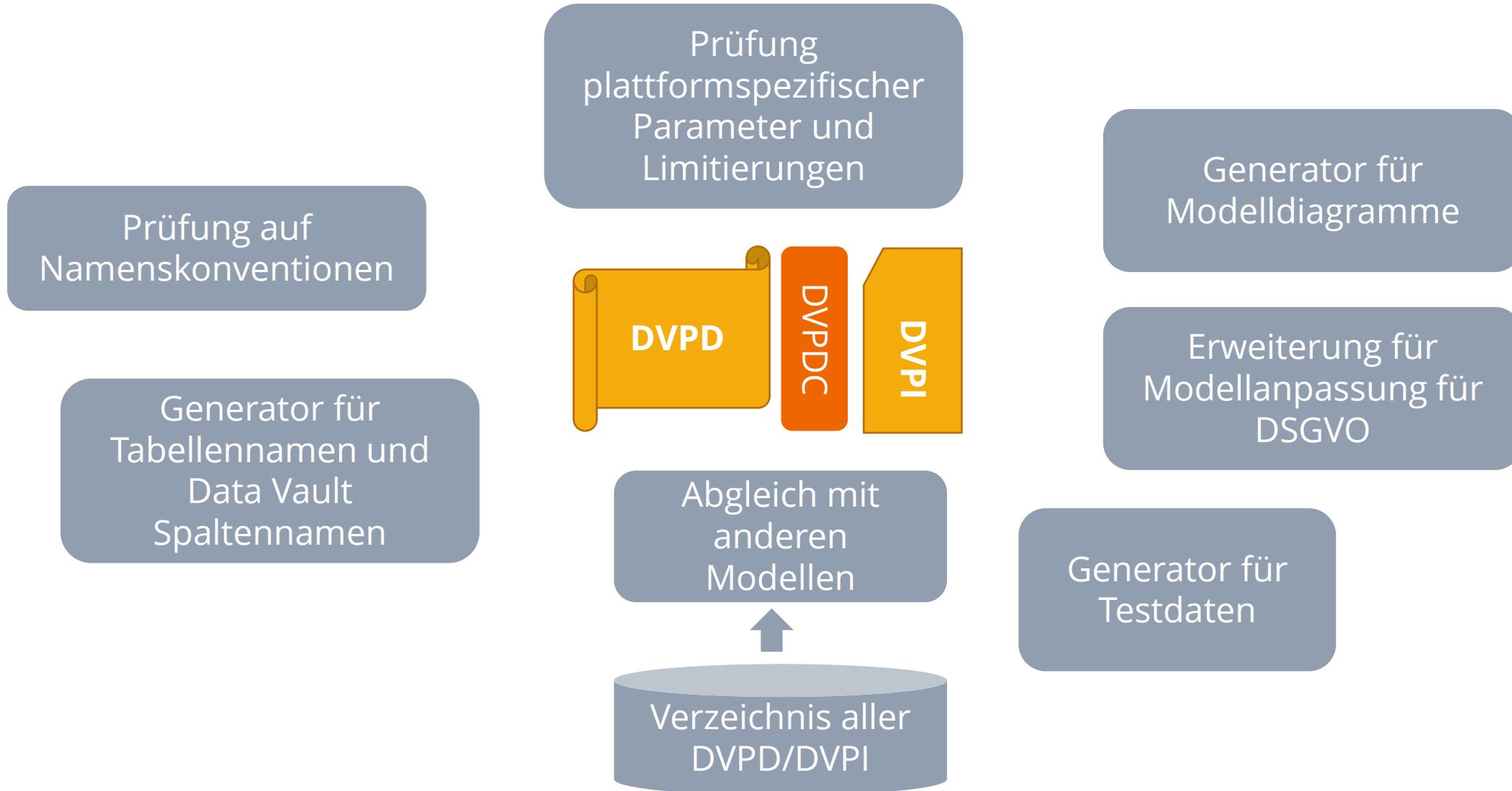




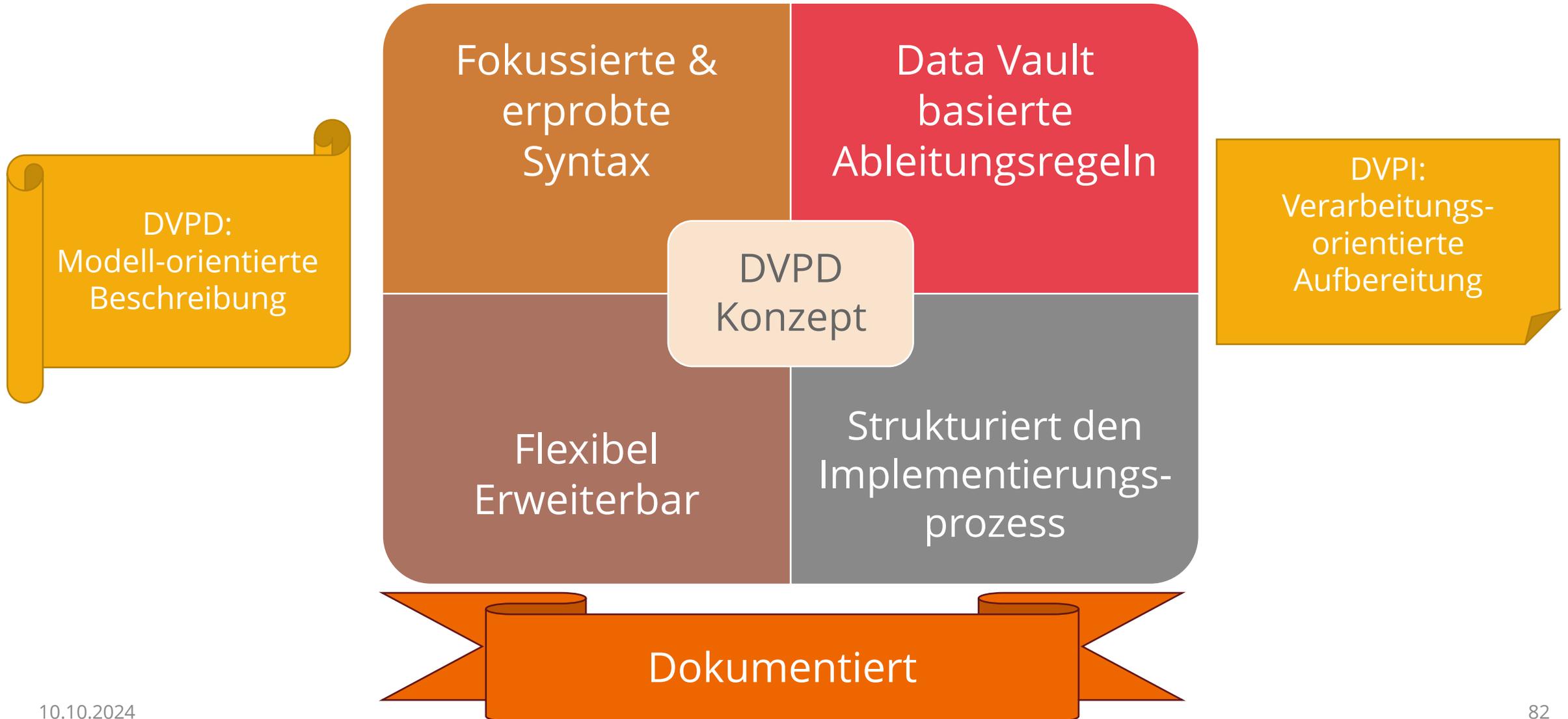
Flexibilität der Toolchain ermöglichen



Freiheit für Erweiterungen



DVPD als Leitlinie & Werkzeug





Q & A



Vielen Dank

Kontakt:

matthias.wegner@cimt-ag.de

LinkedIn: Matthias Wegner

[https://github.com/cimt-ag/
data_vault_pipelinedescription/discussions](https://github.com/cimt-ag/data_vault_pipelinedescription/discussions)