# Big Data Vilnius

## Agile Data Architecture

# About me

Search kaggle

Competitions   Datasets   Kernels   Discussion   Learn   •••   Sign In

## Gerard Toonstra

Rotterdam, Netherlands
Joined 3 years ago · last seen 2 months ago

Followers 1

Competitions Master

Home   Competitions (18)   Kernels (8)   Discussion (149)   Followers (1)   Contact User   Follow User

**KEEP YOUR SOFTWARE SIMPLE**

What does it mean to say "software is complex"?

How do we deal with this complexity?

Gerard Toonstra

## ETL Best Practices with airflow 1.8

1.8

Search docs

- ETL principles
- Gotcha's
- What makes Airflow great?
- ETL example
- Hive example
- Data vault
- Monitoring
- Building your own ETL platform
- Ingesting files

Docs » ETL best practices with Airflow documentation site

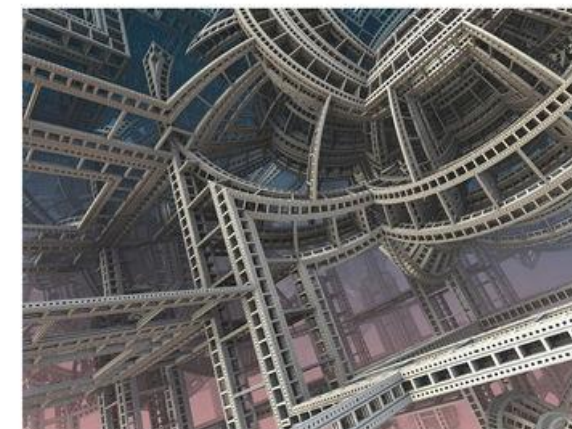# ETL best practices with Airflow documentation site

**❶ Important**

**Disclaimer**: This is not the official documentation site for Apache airflow. This site is not affiliated, monitored or controlled by the official Apache Airflow development effort. If you are looking for the official documentation site, please follow this link:

Official Airflow documentation

What you will find here are interesting examples, usage patterns and ETL principles that I thought are going to help people use airflow to much better effect.
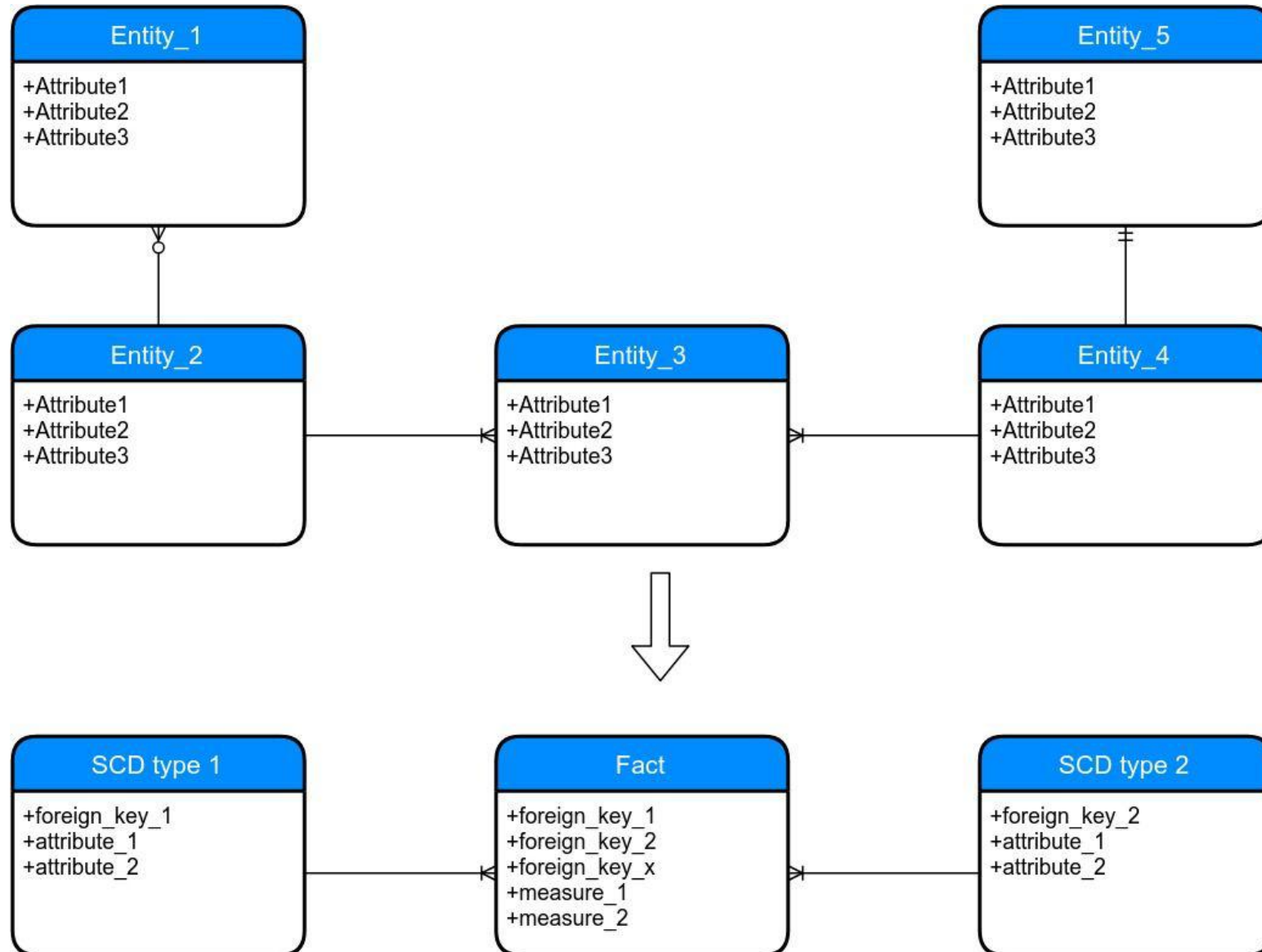
# What is an agile data warehouse?

# Design "lock–in"

# Design "lock-in"

## Slowly Changing Dimension Type 1

| primary_key | customer_name | email |
|---:|---|---|
| 1 | Joe Smith | joe.smith@aol.com |
| 2 | Kerry Jones | kerry.jones@gmail.com |
| 3 | Mary Woods | mary.woods@hotmail.com |

## Slowly Changing Dimension Type 2

| primary_key | customer_name | email | start_date | end_date | is_active |
|---:|---|---|---:|---:|---|
| 1 | Joe Smith | joe.smith@aol.com | 2017-01-05 | 2017-12-24 | N |
| 4 | Joe Smith | joe.smith@gmail.com | 2017-12-24 | NULL | Y |
| 2 | Kerry Jones | kerry.jones@gmail.com | 2017-06-05 | 2017-09-02 | N |
| 5 | Kerry Smith | kerry.smith@gmail.com | 2017-09-02 | NULL | Y |

# Complexity: doing too many things at once

```
INSERT INTO customer_dim
SELECT source_cust_id, first_name, last_name, eff_date, end_date, current_flag
FROM
( MERGE customer_dim cm
  USING customer_source cs
  ON cm.source_cust_id = cs.source_cust_id
  WHEN NOT MATCHED THEN
    INSERT VALUES (cs.source_cust_id, cs.first_name, cs.last_name, convert(char(10), getdate()-1, 101),
'12/31/2199', 'y')
  WHEN MATCHED AND cm.current_flag = 'y' and cm.last_name <> cs.last_name THEN
    UPDATE SET cm.current_flag = 'n', cm.end_date = convert(char(10), getdate()- 2, 101)
  OUTPUT $Action action_out, cs.source_cust_id, cs.first_name, cs.last_name, convert(char(10),
getdate()-1, 101) eff_date, '12/31/2199' end_date, 'y' current_flag
) AS merge_out
WHERE merge_out.action_out = 'UPDATE';
```
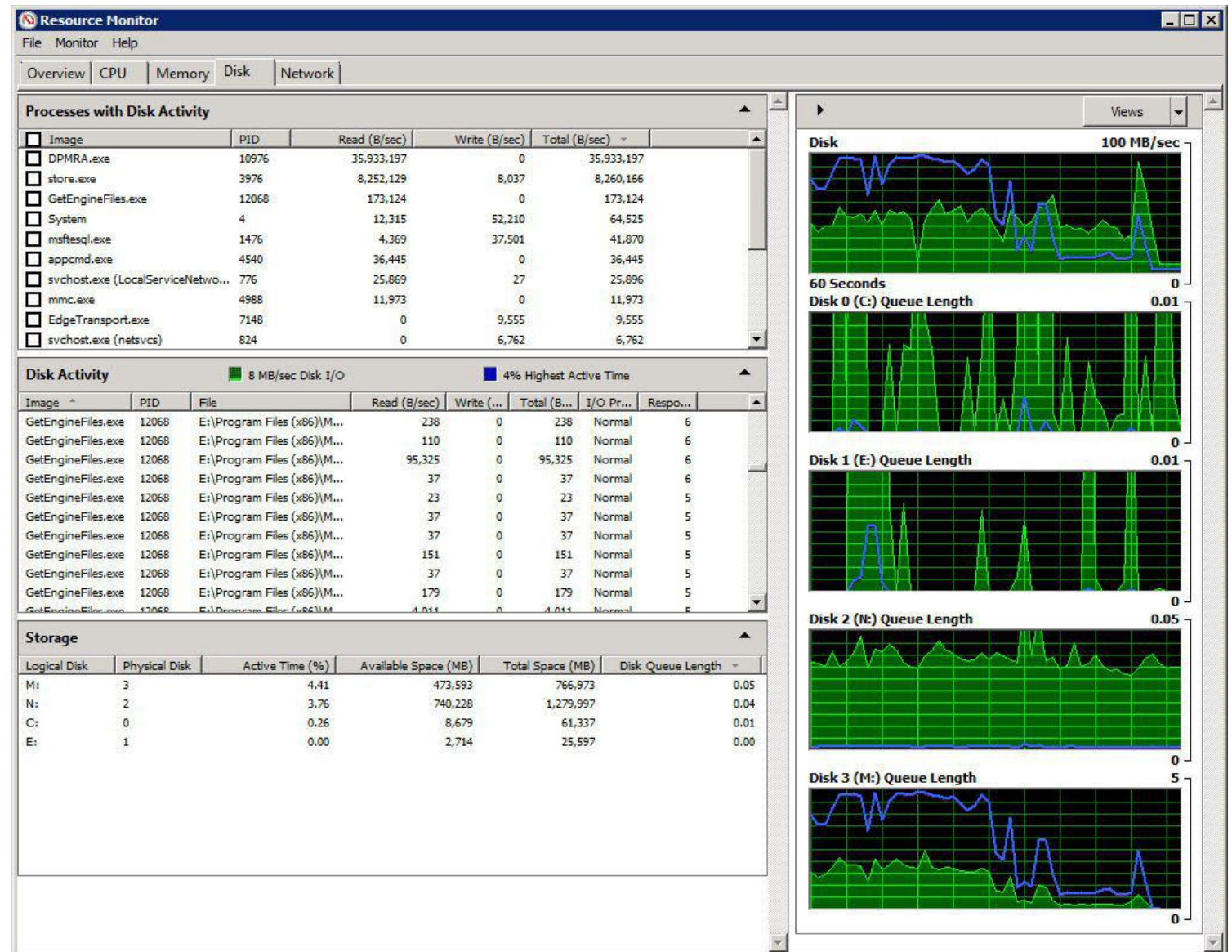
bigdata REPUBLIC
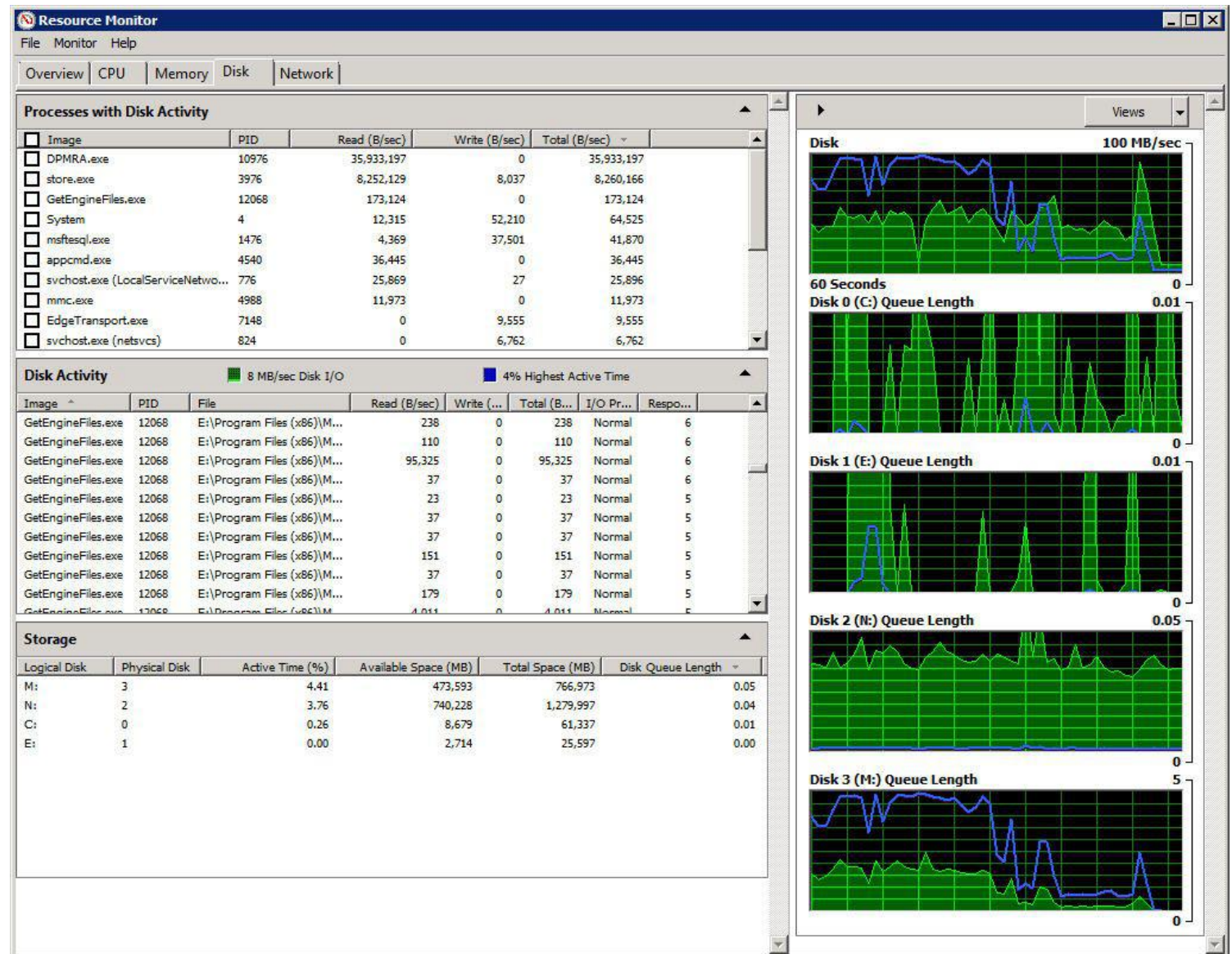
# Data volumes

- Failure to meet SLA

# Data volumes

- Failure to meet SLA

- Long user wait times

# Data volumes

- Failure to meet SLA

- Long user wait times

- **Reports generate high strain → slow**

# Reproducibility

- Rerun parts of your data pipeline without thinking?

# Reproducibility

- Rerun parts of your data pipeline without thinking?
- Can you regenerate your entire warehouse (in principle)?

# Reproducibility

- Rerun parts of your data pipeline without thinking?
- Can you regenerate your entire warehouse (in principle)?
- → Easy to solve bugs

# Limitations of ETL tooling

- Focused on a specific database

# Limitations of ETL tooling

- Focused on a specific database
- **Not scalable**

# Limitations of ETL tooling

- Focused on a specific database
- Not scalable
- Difficult to synchronize with other scheduled pipelines

# Limitations of ETL tooling

- Focused on a specific database
- Not scalable
- Difficult to synchronize with other scheduled pipelines
- **Not built from a functional philosophy**

# Limitations of ETL tooling

- Focused on a specific database
- Not scalable
- Difficult to synchronize with other scheduled pipelines
- Not built from a functional philosophy
- Not extendable as a platform

Engineering is a methodological process of stages requiring:

- engineering skills

- knowing what to do when

- and what NOT to do when

# The concept

# Contextualization
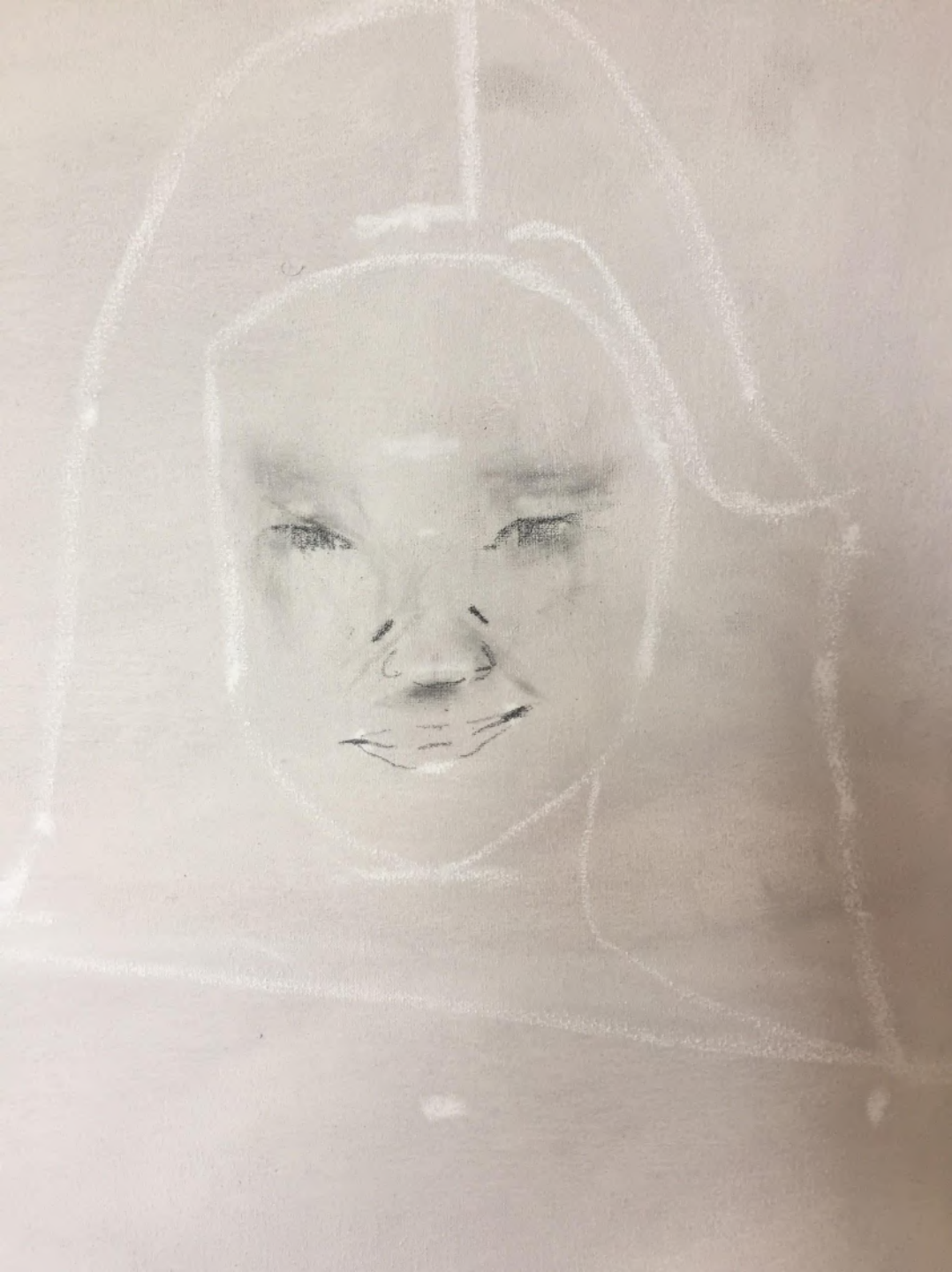
# The underlayer

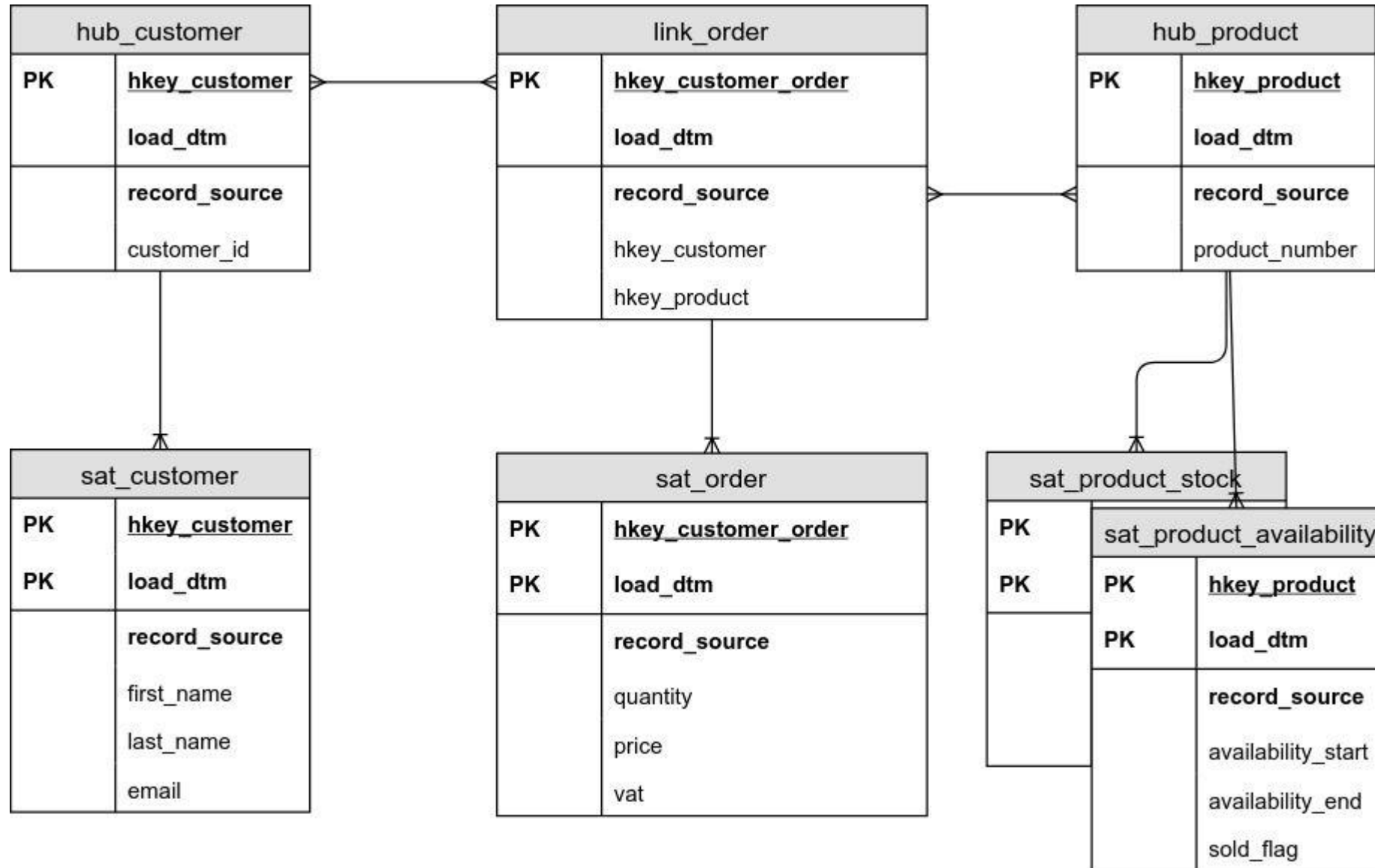# Finished underlayer

# Painting the final layer

# Reductive as well as additive

# Add color

# Agile Data Architecture

- Choose good ETL tools with generalized, reusable components
- Apply design decisions at the right stage, not too early
- Build 'functional' data pipelines
- Make pipelines reproducible → immutable partitions
- Decouple the source schema from your warehouse schema
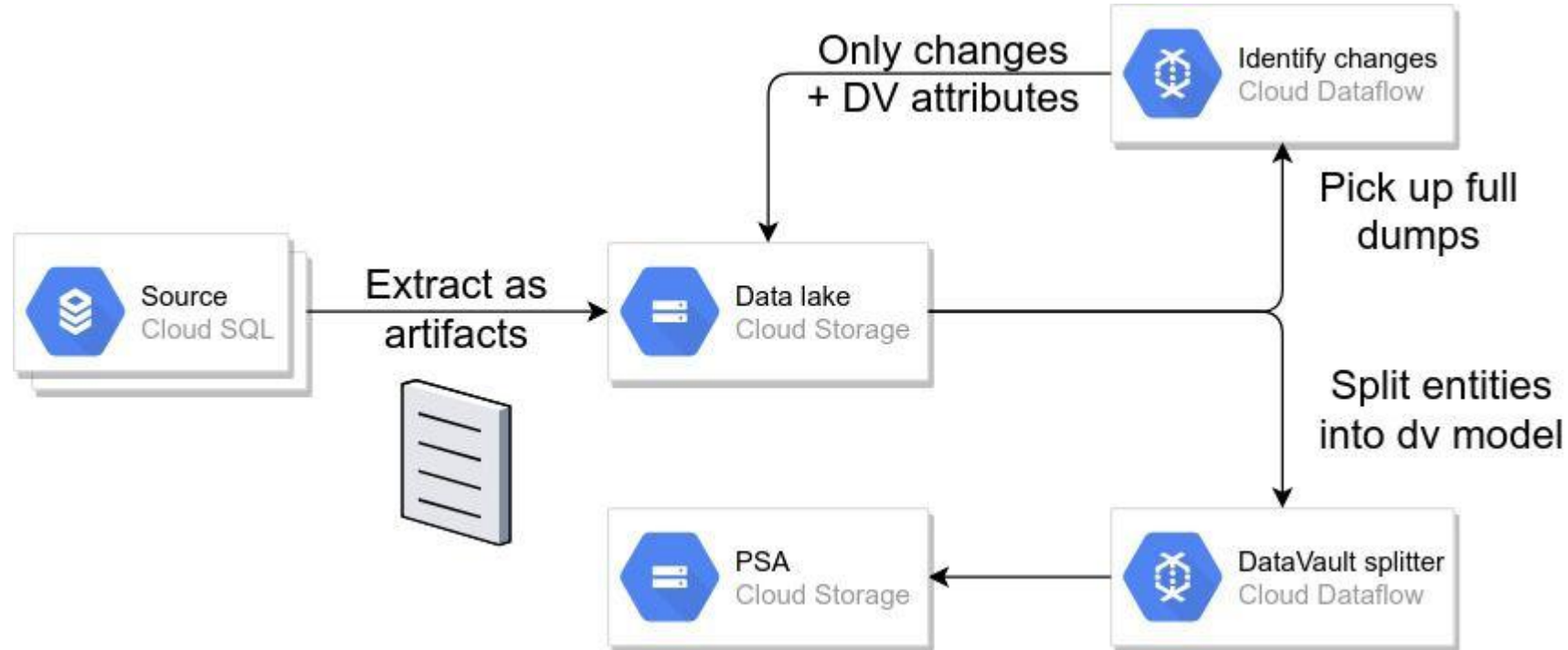- Don't (immediately) resist more persistent copies of the data in pipeline

# The overall data pipeline
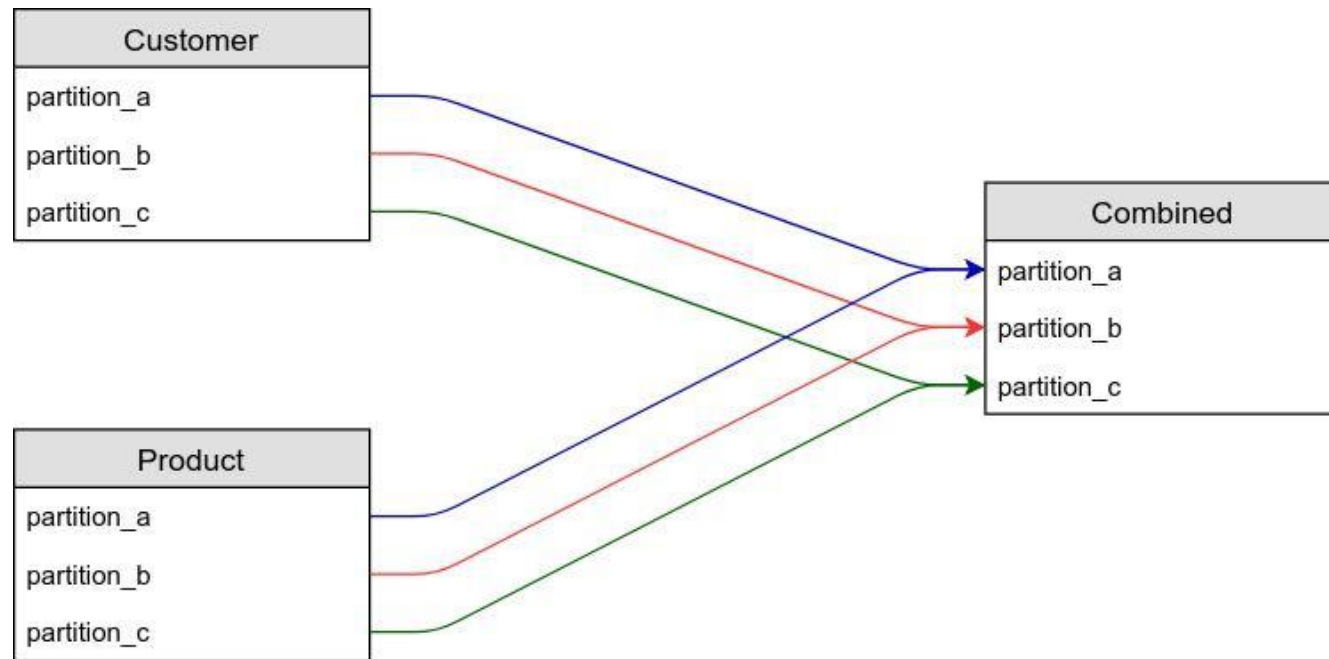
# What is a (raw) data vault 2.0?

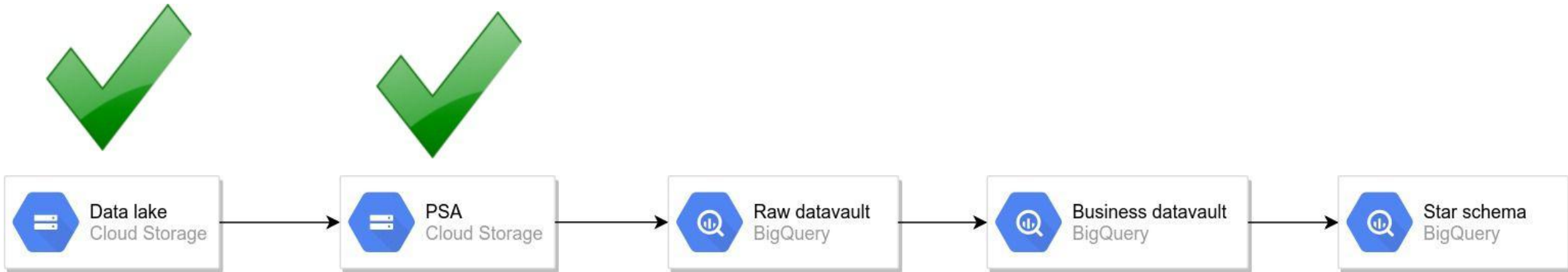# Step 1: Break data into reloadable partitions
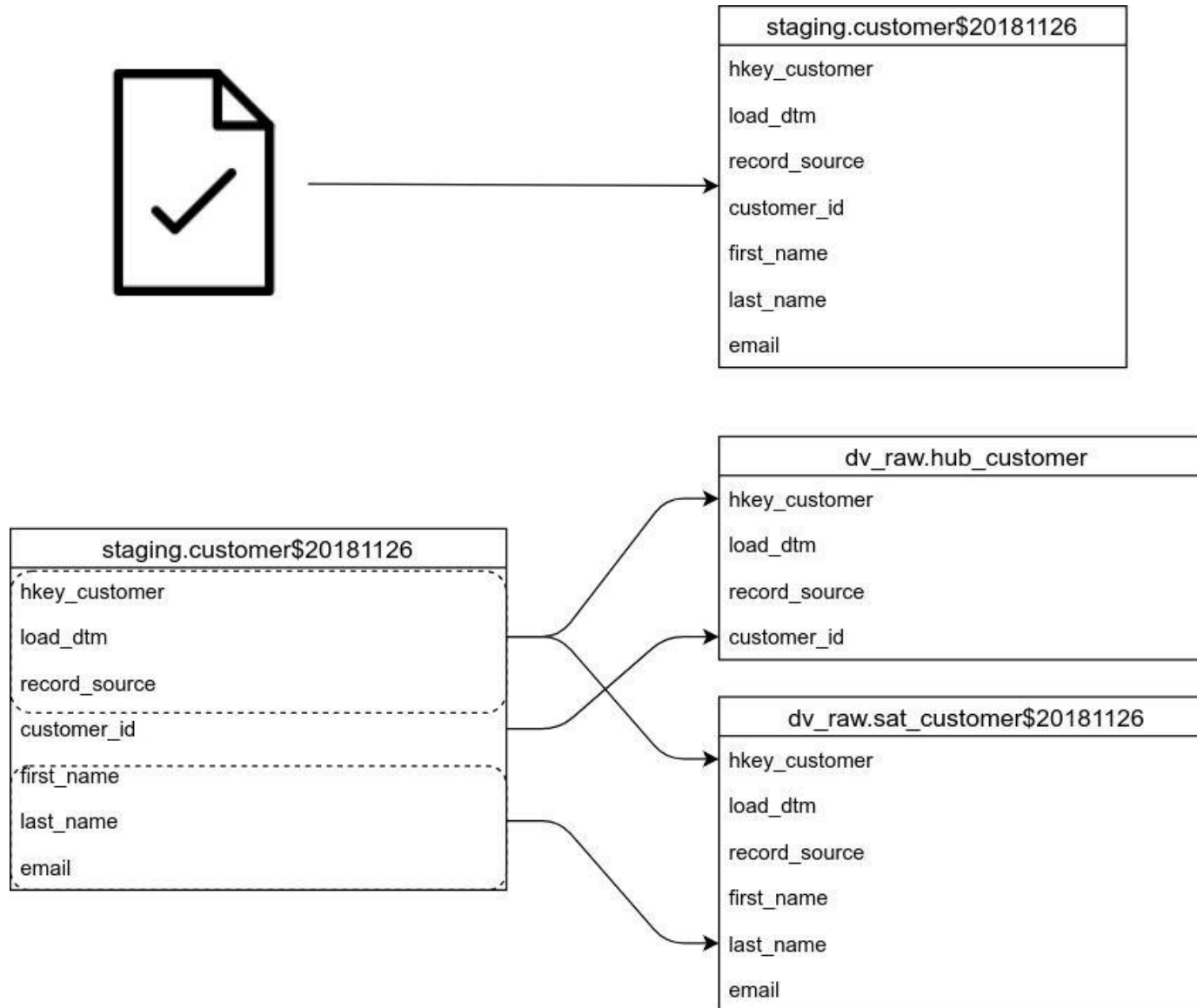
# Step 1: Break data into reloadable partitions

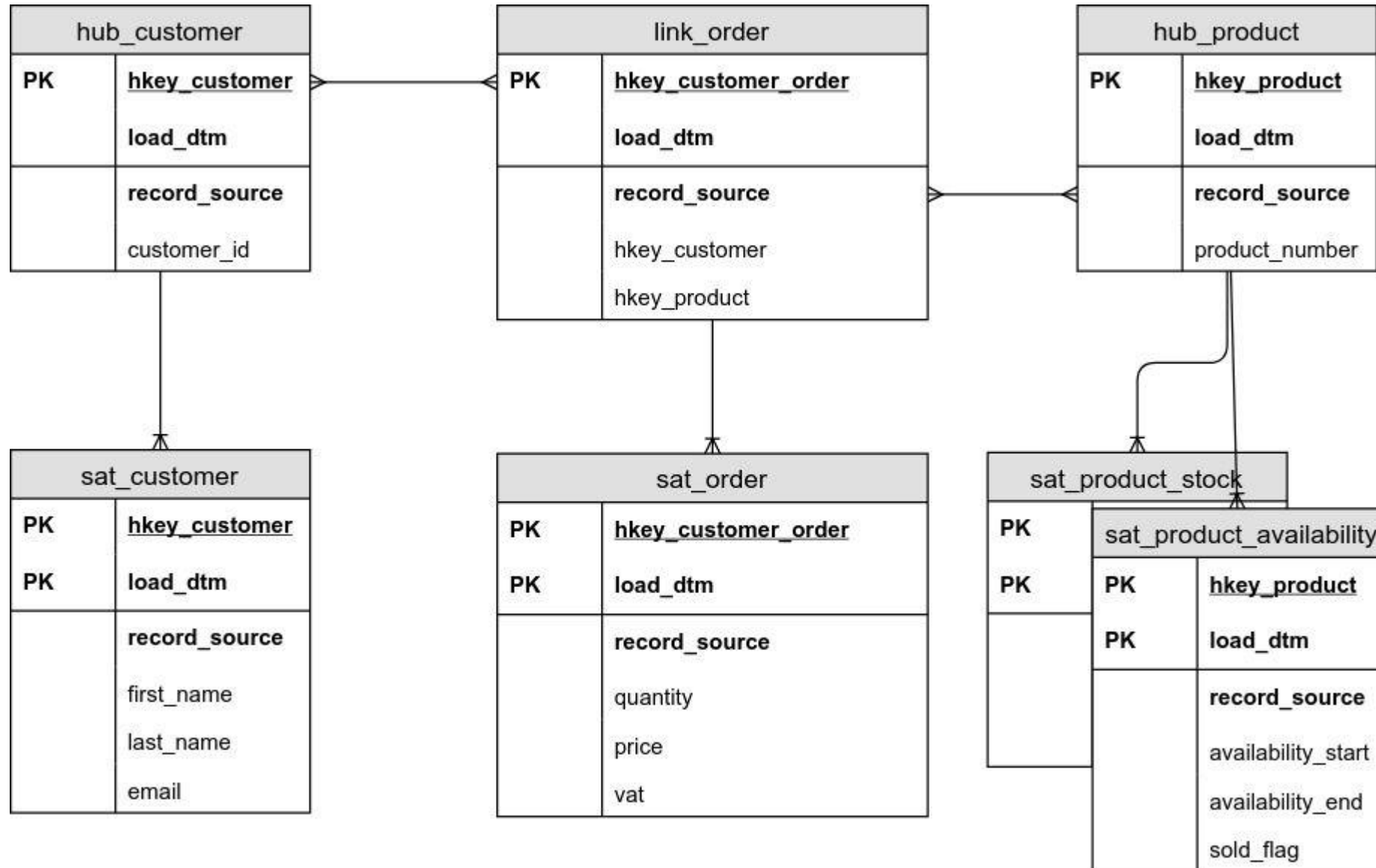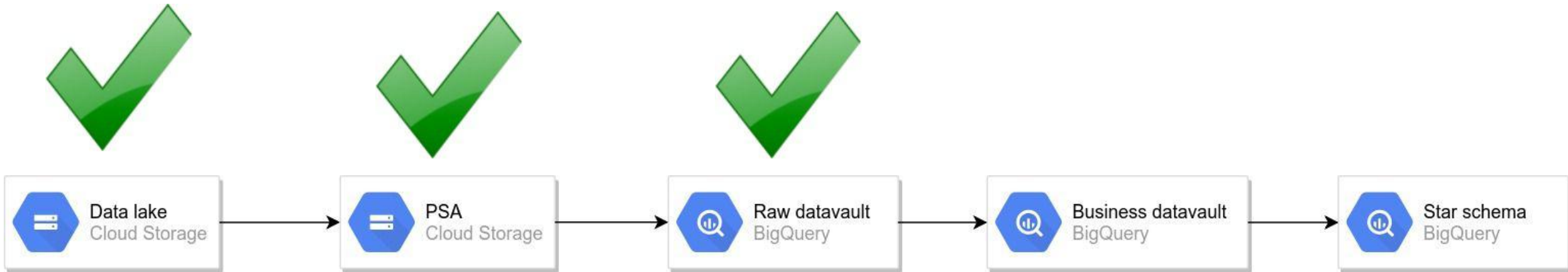| gs://datalake/datavault/psa/crm/customer/2018/11/26/data.avro | partition_a (customer) |
| --- | --- |
| gs://datalake/datavault/psa/crm/customer/2018/11/27/data.avro | partition_b (customer) |
| gs://datalake/datavault/psa/crm/customer/2018/11/28/data.avro | partition_c (customer) |
| gs://datalake/datavault/psa/crm/product/2018/11/26/data.avro | partition_a (product) |
| gs://datalake/datavault/psa/crm/product/2018/11/27/data.avro | partition_b (product) |
| gs://datalake/datavault/psa/crm/product/2018/11/28/data.avro | partition_c (product) |

# The overall architecture
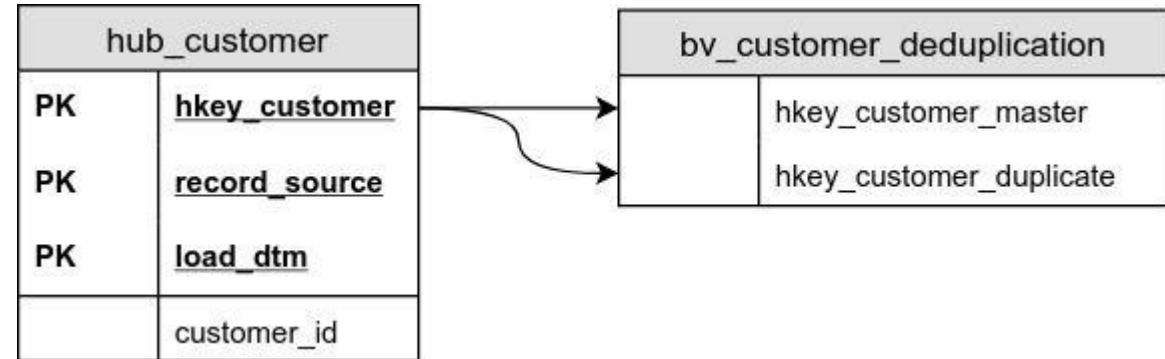
# Step 2: Load raw datavault

# Do this for all entities

# The overall architecture

# The business vault (data interpretation)

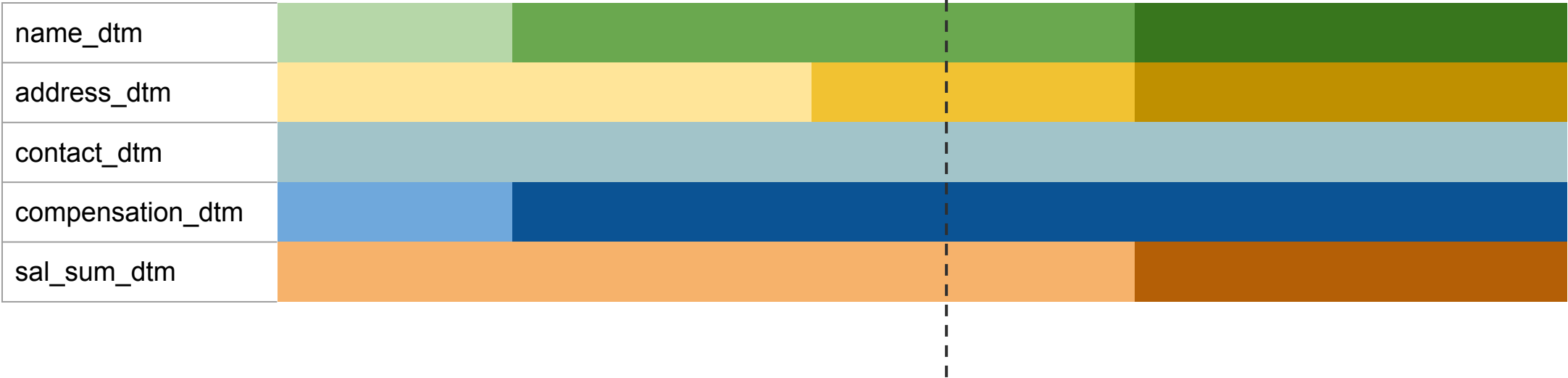Customer deduplication:



Bridge & PIT tables:



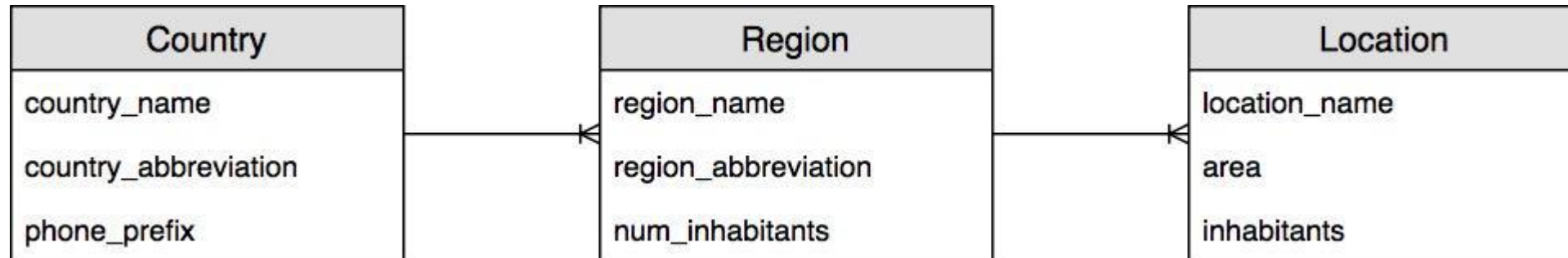Calculations, data cleaning, data quality, soft business rules

# The PIT table (data interpretation)

EMPLOYEE_PIT
- hkey_employee
- pit_load_dtm
- name_load_dtm
- address_load_dtm
- contact_load_dtm
- compensation_load_dtm
- sal_sum_load_dtm

?

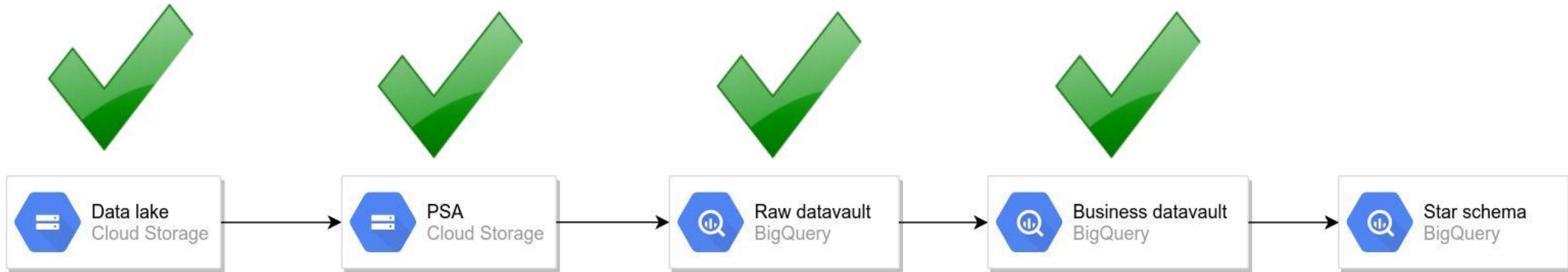| name_dtm | | | |
| --- | --- | --- | --- |
| address_dtm | | | |
| contact_dtm | | | |
| compensation_dtm | | | |
| sal_sum_dtm | | | |

# The Bridge table (data interpretation)

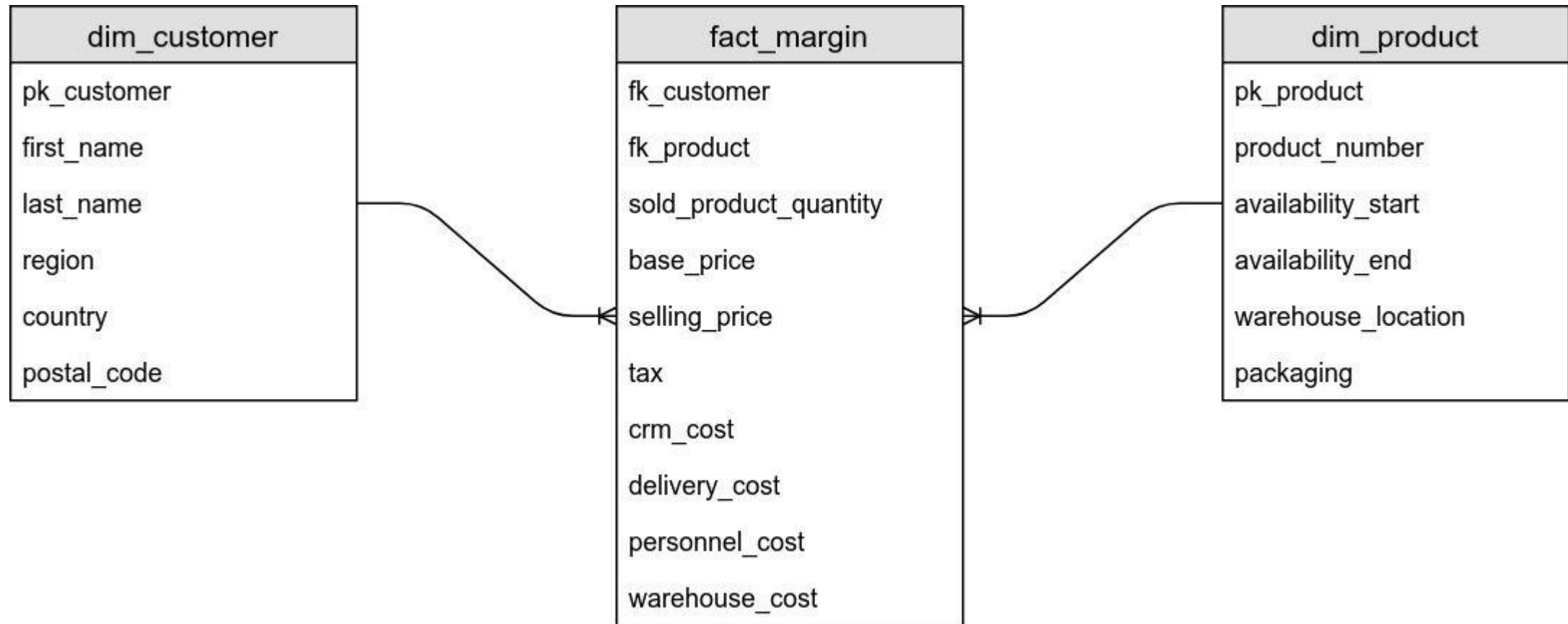**BRIDGE_GEOGRAPHY**

hkey_bridge_geography

bridge_load_dtm

hkey_link_regional_country

hkey_location
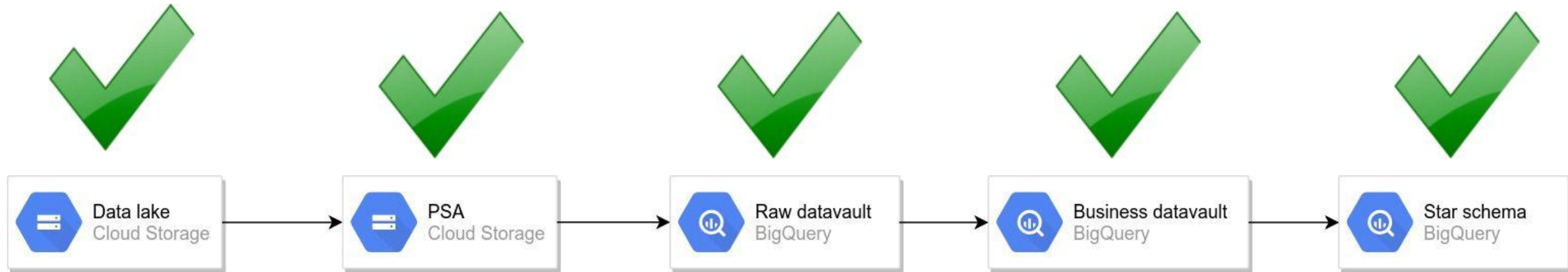
hkey_country

hkey_region

country_name

**Country**

country_name

country_abbreviation

phone_prefix

**Region**

region_name

region_abbreviation

num_inhabitants

**Location**

location_name

area

inhabitants

# The overall architecture



| Data lake Cloud Storage | → | PSA Cloud Storage | → | Raw datavault BigQuery | → | Business datavault BigQuery | → | Star schema BigQuery |

bigdata REPUBLIC

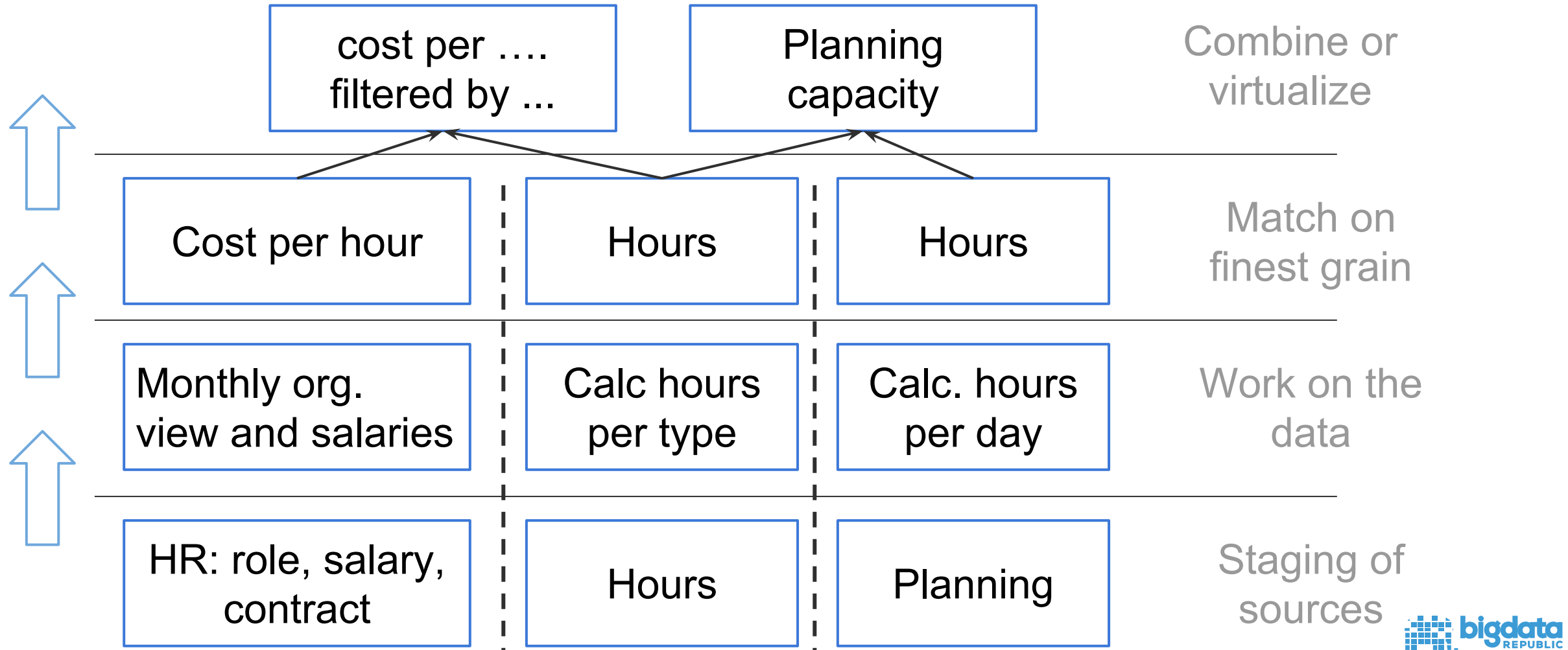# The star schema (data interpretation)
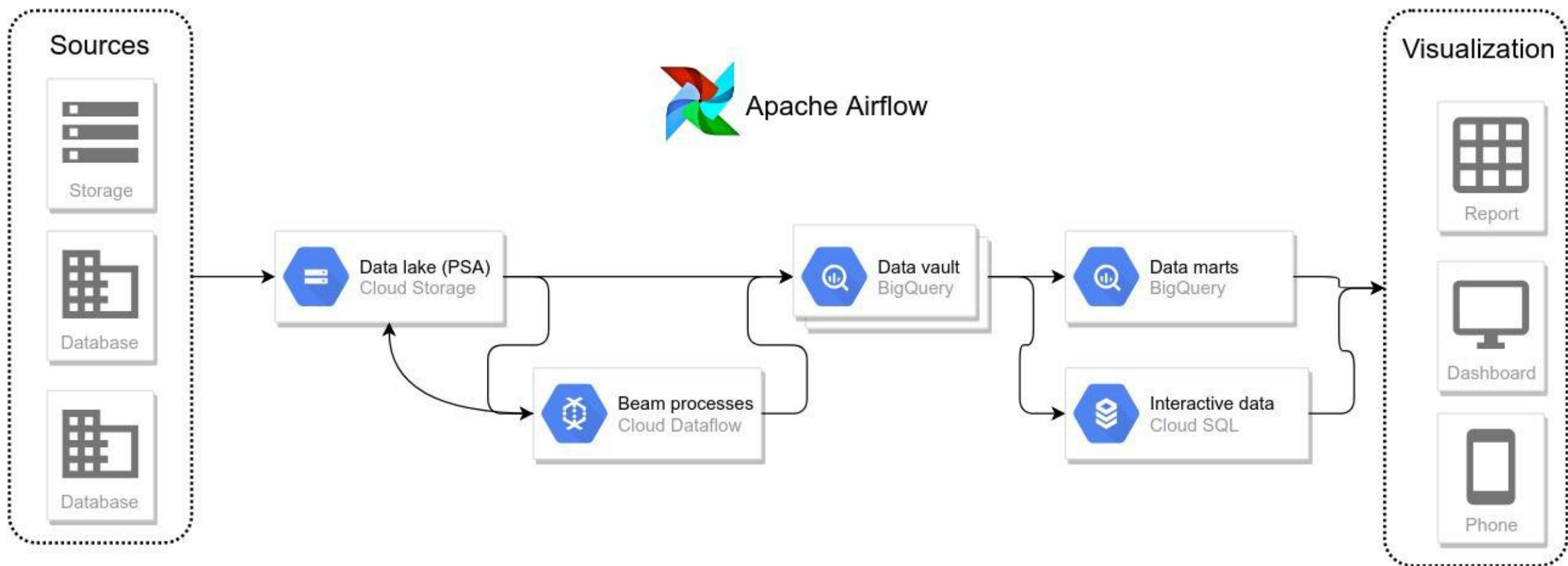
# The overall architecture

# Avoid complication with "data lanes"

Merge data from different sources as late as possible

# A possible architecture

Join at
slido.com: #bigdata2018

# bigdata REPUBLIC

DATA SCIENCE | BIG DATA ANALYTICS | BIG DATA ARCHITECTURES

📞 +31 (0)1 – 68479294

🏠 Coltbaan 4C, Nieuwegein

🐦 @bigdatarep

🌐 www.bigdatarepublic.nl

in /company/bigdata–republic

✉ info@bigdatarepublic.nl

bigdata REPUBLIC