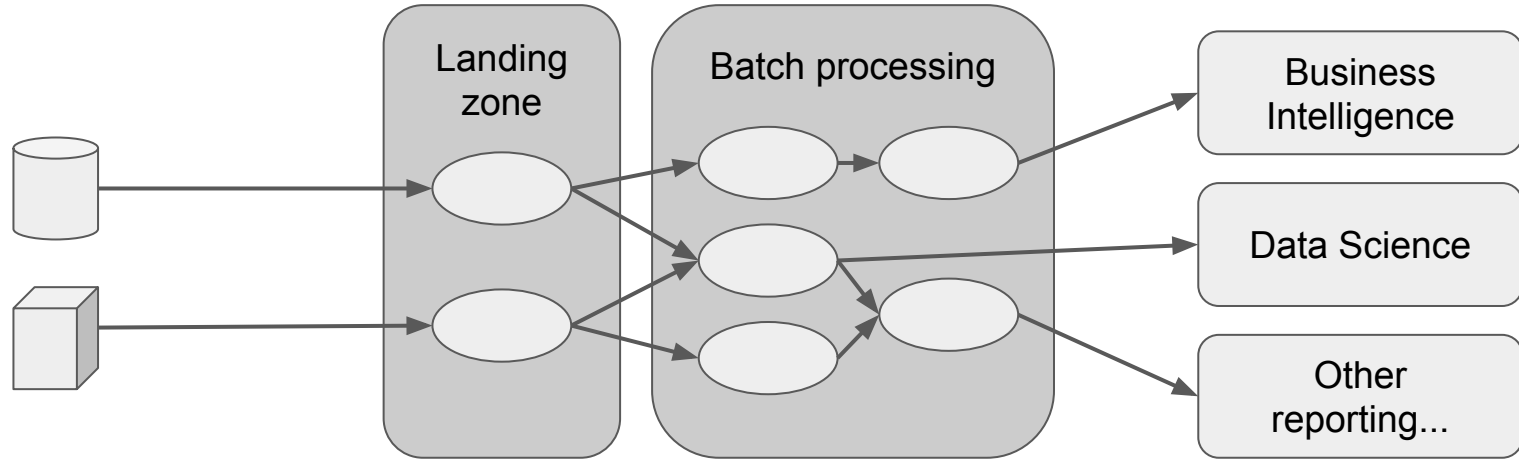# Building a data lake at a bank
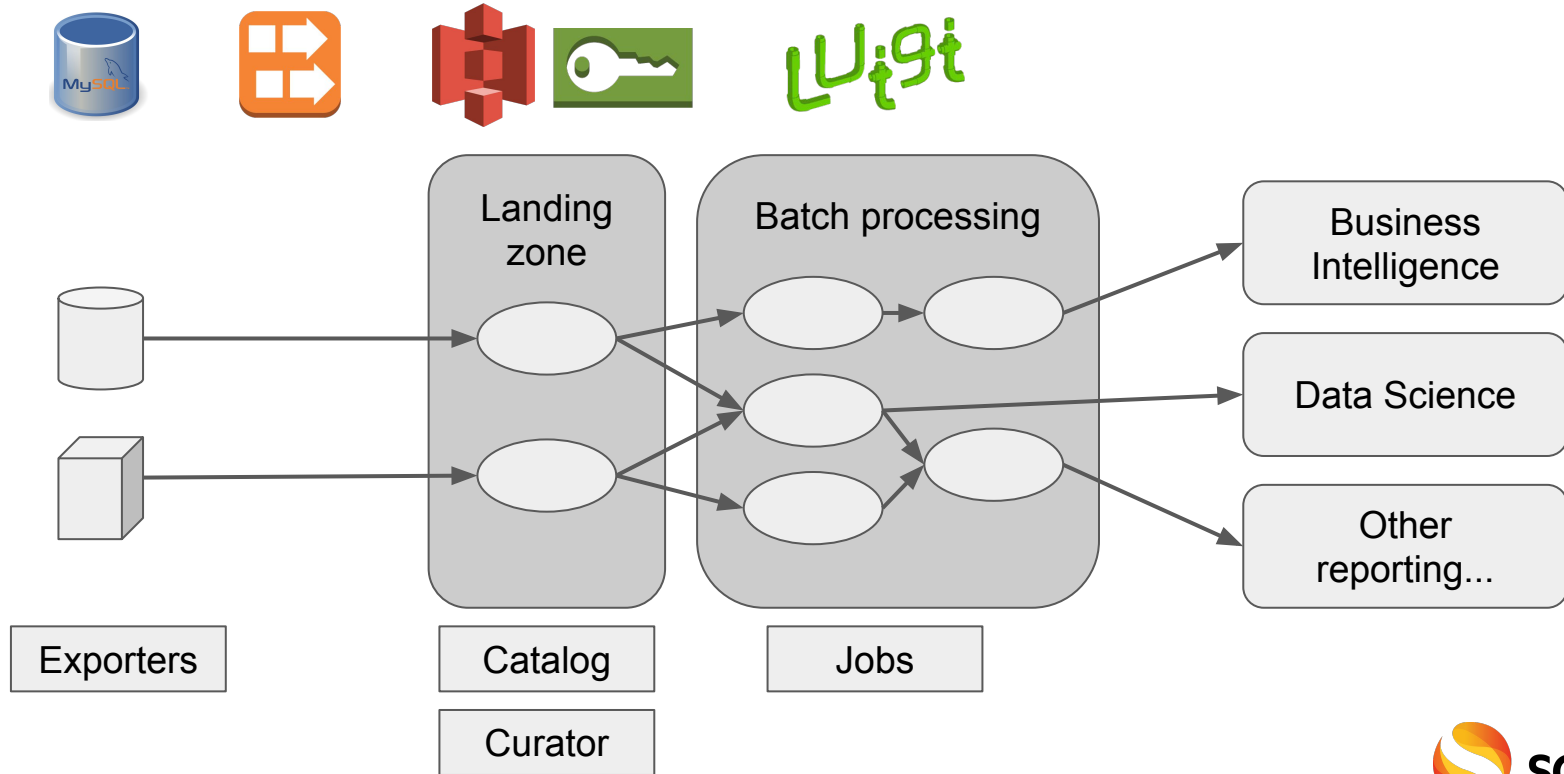
**solarisBank**

# Overview

# Overview

# Regulation / GDPR

Typical technical requirements

solarisBank

# Regulation / GDPR

Typical technical requirements

- Encryption
    - At rest
    - In transit

solarisBank

# Regulation / GDPR

Typical technical requirements

- Encryption
  - At rest
  - In transit
- Strong authentication (MFA)

solarisBank

# Regulation / GDPR

Typical technical requirements

- Encryption
    - At rest
    - In transit
- Strong authentication (MFA)
- Auditability, automation and version control

# Regulation / GDPR

Typical technical requirements

- Encryption
    - At rest
    - In transit
- Strong authentication (MFA)
- Auditability, automation and version control
- Roles, separation of concerns

solarisBank

# Regulation / GDPR

Typical technical requirements

- Encryption
    - At rest
    - In transit
- Strong authentication (MFA)
- Auditability, automation and version control
- Roles, separation of concerns
- Pseudonymization/masking of data

solarisBank

# Regulation / GDPR

Typical technical requirements

- Encryption
    - At rest
    - In transit
- Strong authentication (MFA)
- Auditability, automation and version control
- Roles, separation of concerns
- Pseudonymization/masking of data
- Other IT security measures

solarisBank

# Regulation / GDPR

Typical technical requirements

- Encryption
    - At rest
    - In transit
- Strong authentication (MFA)
- Auditability, automation and version control
- Roles, separation of concerns
- Pseudonymization/masking of data
- Other IT security measures

## "Easy" part

**solarisBank**

# Data protection by design and by default (Art. 25)

solarisBank

# Data protection by design and by default (Art. 25)

Bla bla bla …

- implement appropriate technical and organisational measures …
- which are designed to implement data-protection principles …
- and to integrate the necessary safeguards into the processing in order to meet the requirements of this Regulation

… bla bla bla

solarisBank

# GDPR Principles

High level GDPR (Art. 5)

- Lawful processing
    - Purpose: you must have a good reason for processing
    - Legal basis: you also need a reason why your purpose is legal
- Data minimization
    - Collection and processing: only use what you really need
    - Deletion: once you are done, you must delete
- Accountability
    - onus of proof is on the company

solarisBank

# Lawful processing
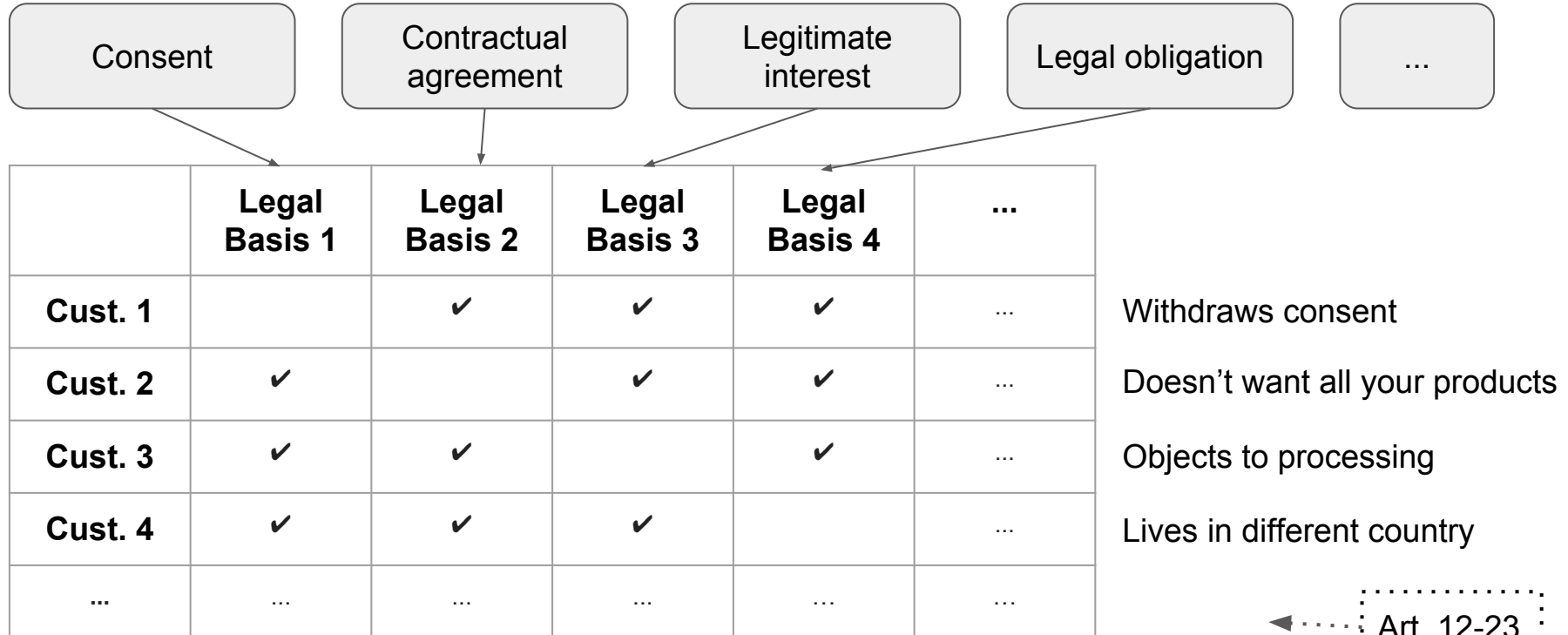
Art. 6

# Lawful processing
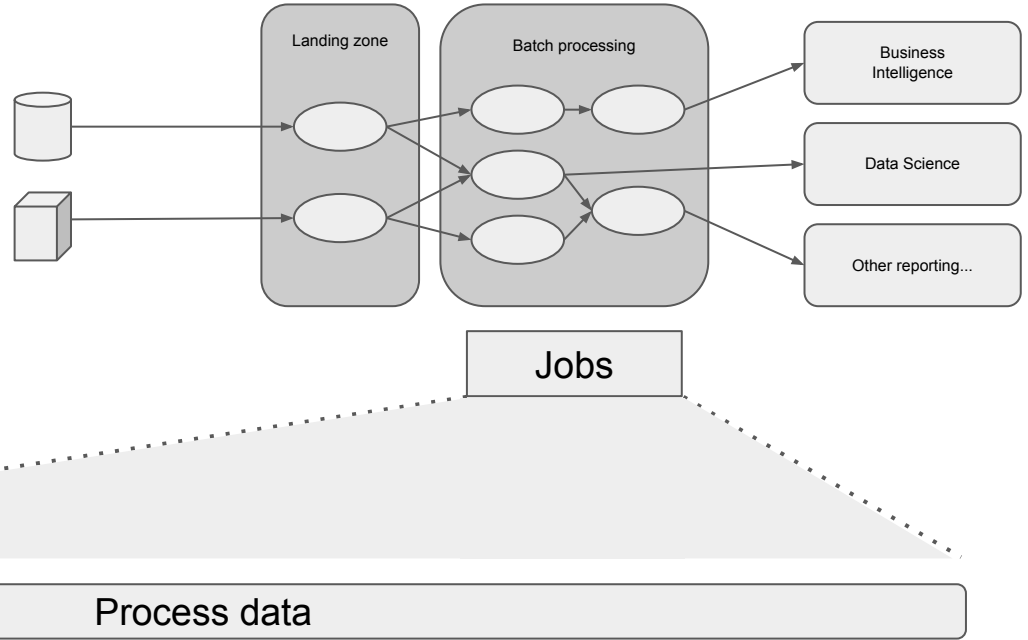
Art. 6

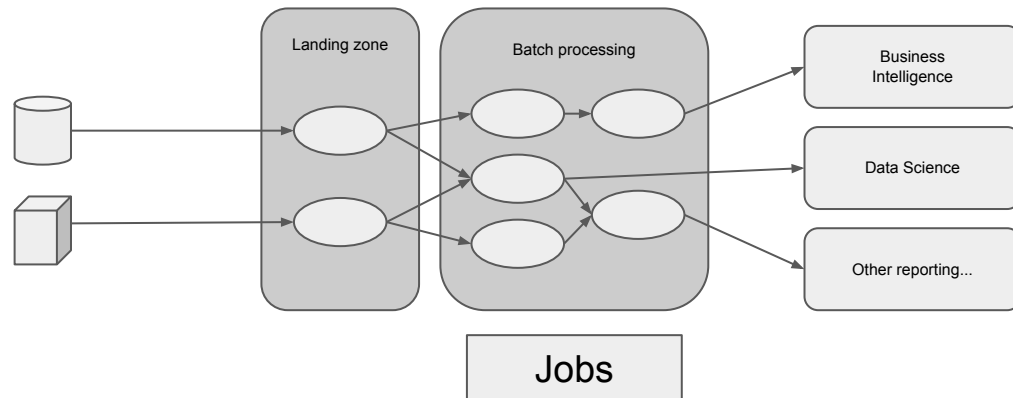| Consent | Contractual agreement | Legitimate interest | Legal obligation | ... |

solarisBank

# Lawful processing

| | **Legal Basis 1** | **Legal Basis 2** | **Legal Basis 3** | **Legal Basis 4** | **...** |
|---|---|---|---|---|---|
| **Cust. 1** | | ✔ | ✔ | ✔ | ... |
| **Cust. 2** | ✔ | | ✔ | ✔ | ... |
| **Cust. 3** | ✔ | ✔ | | ✔ | ... |
| **Cust. 4** | ✔ | ✔ | ✔ | | ... |
| ... | ... | ... | ... | ... | ... |

Consent

Contractual agreement

Legitimate interest

Legal obligation

...

Withdraws consent

Doesn't want all your products

Objects to processing

Lives in different country

Art. 12-23

# Lawful processing, by design

Landing zone

Batch processing

Business Intelligence

Data Science

Other reporting...

Jobs

**Before:** Process data

**After:** Identify and filter → Process data

Legal basis table

solarisBank

# Lawful processing, by design



Landing zone

Batch processing

Business Intelligence

Data Science

Other reporting...

Preprocessing necessary !?!

Jobs

**Before:** Process data

**After:** Identify and filter → Process data

Legal basis table

solarisBank

# Data minimization

Only store/use personal data you need, and store it only as long as you need it.

→ Use only the data you really need: very difficult topic… (e.g, machine learning)

→ Implement data retention policies everywhere: deletion!

solarisBank

# Data minimization

Only store/use personal data you need, and store it only as long as you need it.

→ Implement data retention policies everywhere

| Consent | Contractual agreement | Legitimate interest | Legal obligation | ... |

Triggers: every legal basis may an expiry date

- Consent can be withdrawn
- Contracts may end
- Legitimate interests may be subject to objection (Art. 21)
- Legal obligations usually have a time-frame
- ...
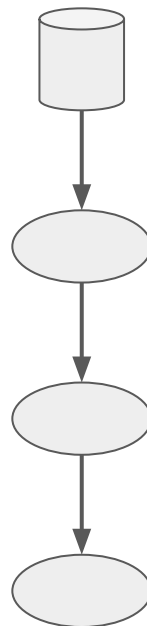
→ Definitely stop processing

→ Maybe also delete

solarisBank

# Data minimization - deletion

Transactional systems:
- Hopefully easy, if backed by something like RDBMS
- Look at DB schema, get and/or delete all data about a person

Data Pipelines
- Could be problematic
    - Ideal: "Functional style"
    - Immutable intermediate results, often materialized
    - Deletion by "marking as deleted" does not count
        ...

solarisBank
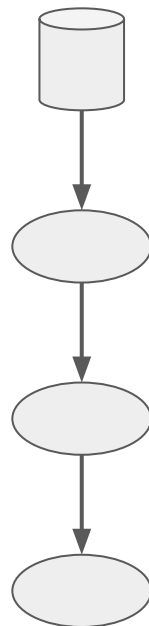
# Data minimization - deletion, by design

Transactional systems:
- Hopefully easy, if backed by something like RDBMS
- Look at DB schema, get and/or delete all data about a person

Data Pipelines
- Could be problematic
    - Ideal: "Functional style"
    - Immutable intermediate results, often materialized
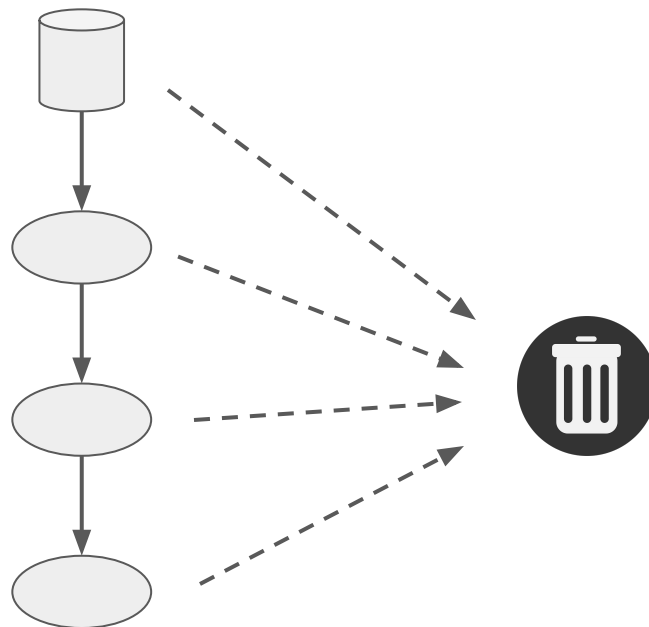    - Deletion by "marking as deleted" does not count ...
- Only 3 ways out:

# Data minimization - deletion, by design

Transactional systems:
- Hopefully easy, if backed by something like RDBMS
- Look at DB schema, get and/or delete all data about a person

Data Pipelines
- Could be problematic
  - Ideal: "Functional style"
  - Immutable intermediate results, often materialized
  - Deletion by "marking as deleted" does not count ...
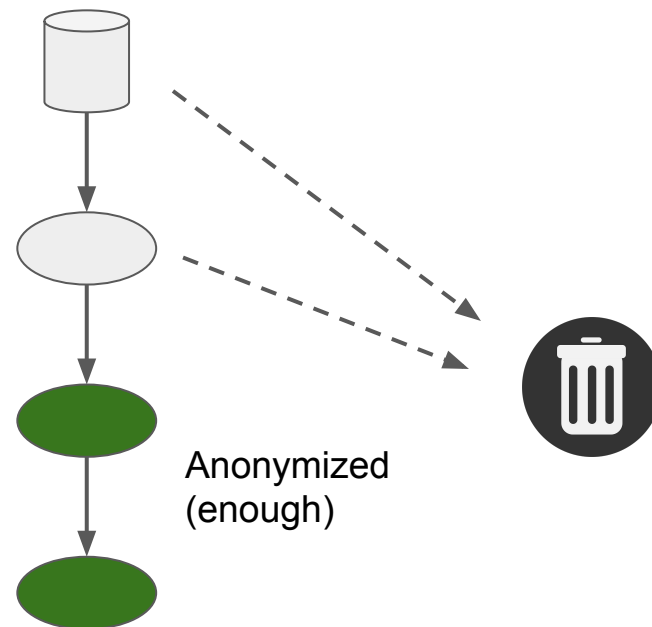- Only 3 ways out:
  - Low enough retention by default

solarisBank

# Data minimization - deletion, by design

Transactional systems:
- Hopefully easy, if backed by something like RDBMS
- Look at DB schema, get and/or delete all data about a person

Data Pipelines
- Could be problematic
  - Ideal: "Functional style"
  - Immutable intermediate results, often materialized
  - Deletion by "marking as deleted" does not count ...
- Only 3 ways out:
  - Low enough retention by default
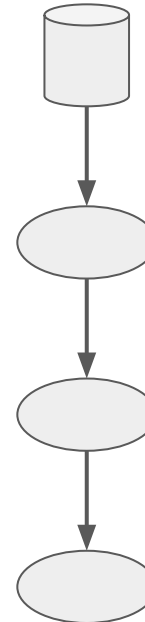  - Make data anonymized

Anonymized (enough)

solarisBank

# Data minimization - deletion, by design

Transactional systems:
- Hopefully easy, if backed by something like RDBMS
- Look at DB schema, get and/or delete all data about a person

Data Pipelines
- Could be problematic
  - Ideal: "Functional style"
  - Immutable intermediate results, often materialized
  - Deletion by "marking as deleted" does not count ...
- Only 3 ways out:
  - Low enough retention by default
  - Make data anonymized
  - Clean up after you

solarisBank

# Data minimization - deletion, by design

**Clean up after you**

Transactional systems:
- Hopefully easy, if backed by something like RDBMS
- Look at DB schema, get and/or delete all data about a person

Data Pipelines
- Could be problematic
  - Ideal: "Functional style"
  - Immutable intermediate results, often materialized
  - Deletion by "marking as deleted" does not count ...
- Only 3 ways out:
  - Low enough retention by default
  - Make data anonymized
  - Clean up after you

Recomputation
- Clean up source
- Re-run downstream jobs

solarisBank

# Data minimization - deletion, by design

Transactional systems:
- Hopefully easy, if backed by something like RDBMS
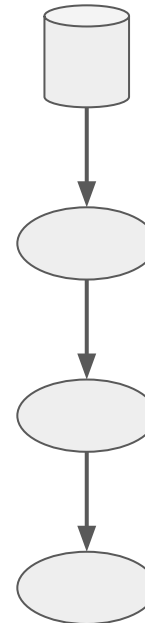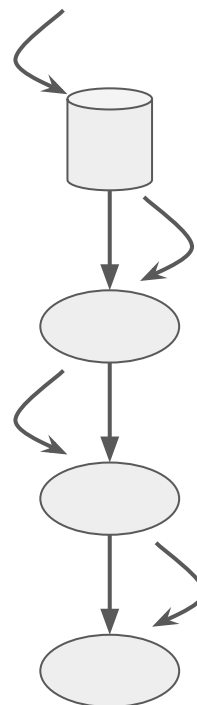- Look at DB schema, get and/or delete all data about a person

Data Pipelines
- Could be problematic
  - Ideal: "Functional style"
  - Immutable intermediate results, often materialized
  - Deletion by "marking as deleted" does not count ...
- Only 3 ways out:
  - Low enough retention by default
  - Make data anonymized
  - Clean up after you

**Clean up after you**

Recomputation
- Clean up source
- Re-run downstream jobs

Cleaning pipelines
- Additional job cleans up source data and all intermediate/end results

solarisBank

# Data minimization - deletion, by design

Transactional systems:
- Hopefully easy, if backed by something like RDBMS
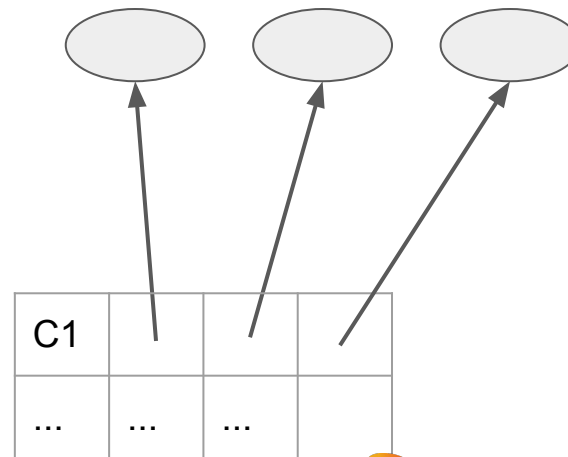- Look at DB schema, get and/or delete all data about a person

Data Pipelines
- Could be problematic
  - Ideal: "Functional style"
  - Immutable intermediate results, often materialized
  - Deletion by "marking as deleted" does not count ...
- Only 3 ways out:
  - Low enough retention by default
  - Make data anonymized
  - Clean up after you

**Factor out**

Join table
- Assign a pseudonym per dataset and person
- Maintain one table to join all pseudonyms



solarisBank

# Data minimization - deletion, by design

Transactional systems:
- Hopefully easy, if backed by something like RDBMS
- Look at DB schema, get and/or delete all data about a person

Data Pipelines
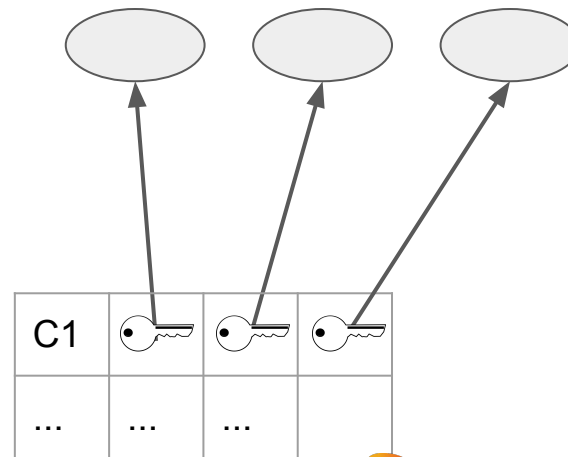- Could be problematic
  - Ideal: "Functional style"
  - Immutable intermediate results, often materialized
  - Deletion by "marking as deleted" does not count ...
- Only 3 ways out:
  - Low enough retention by default
  - Make data anonymized
  - Clean up after you

**Factor out**

Encryption key join table
- "Lost key pattern"
- Assign an encryption key per dataset and person
- Maintain one table to join all keys



solarisBank

# Compliance as code

Smaller Word documents = world is a better place