Annual growth in global data volume Seagate forecasts

- ⚲ Unstoppable growth of data volume;

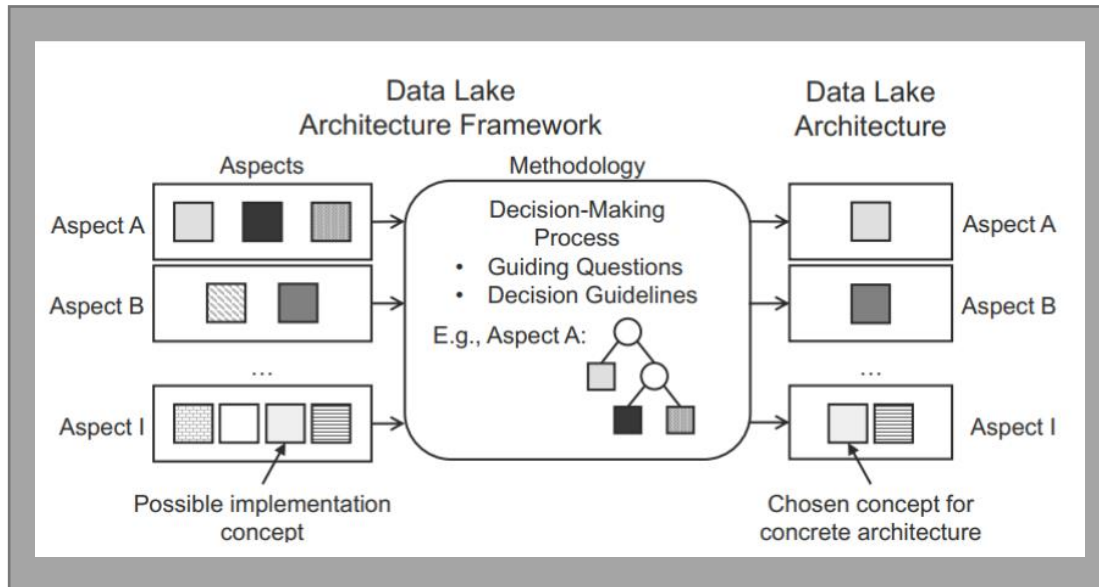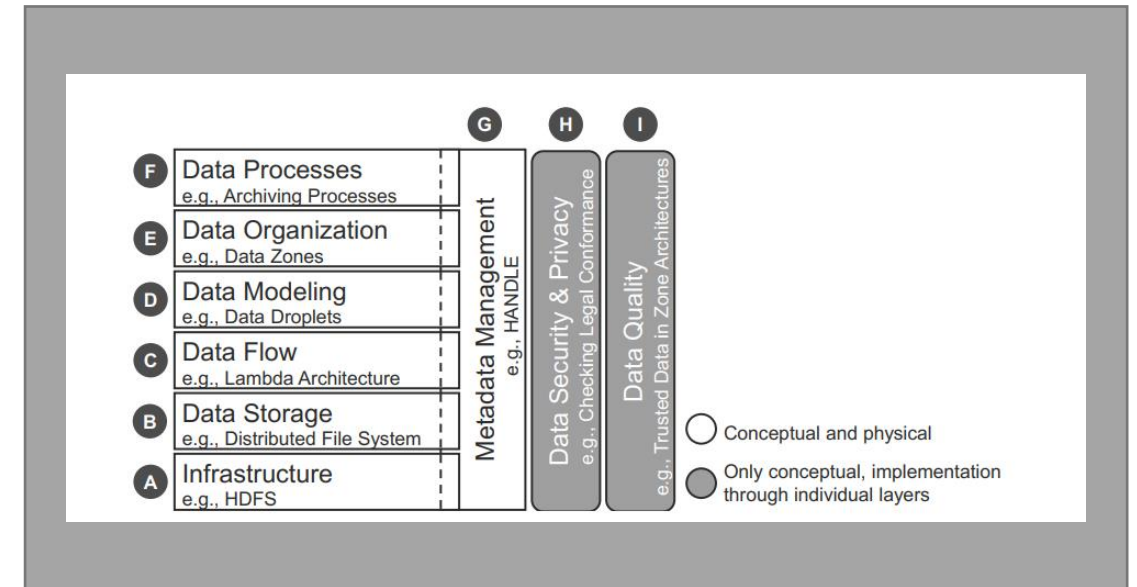- ⚲ High speed of progress of technologies for storing and processing ultra-large data;

- ⚲ The lag of methods and means of protection of the above-mentioned technologies from the development of the technologies themselves.
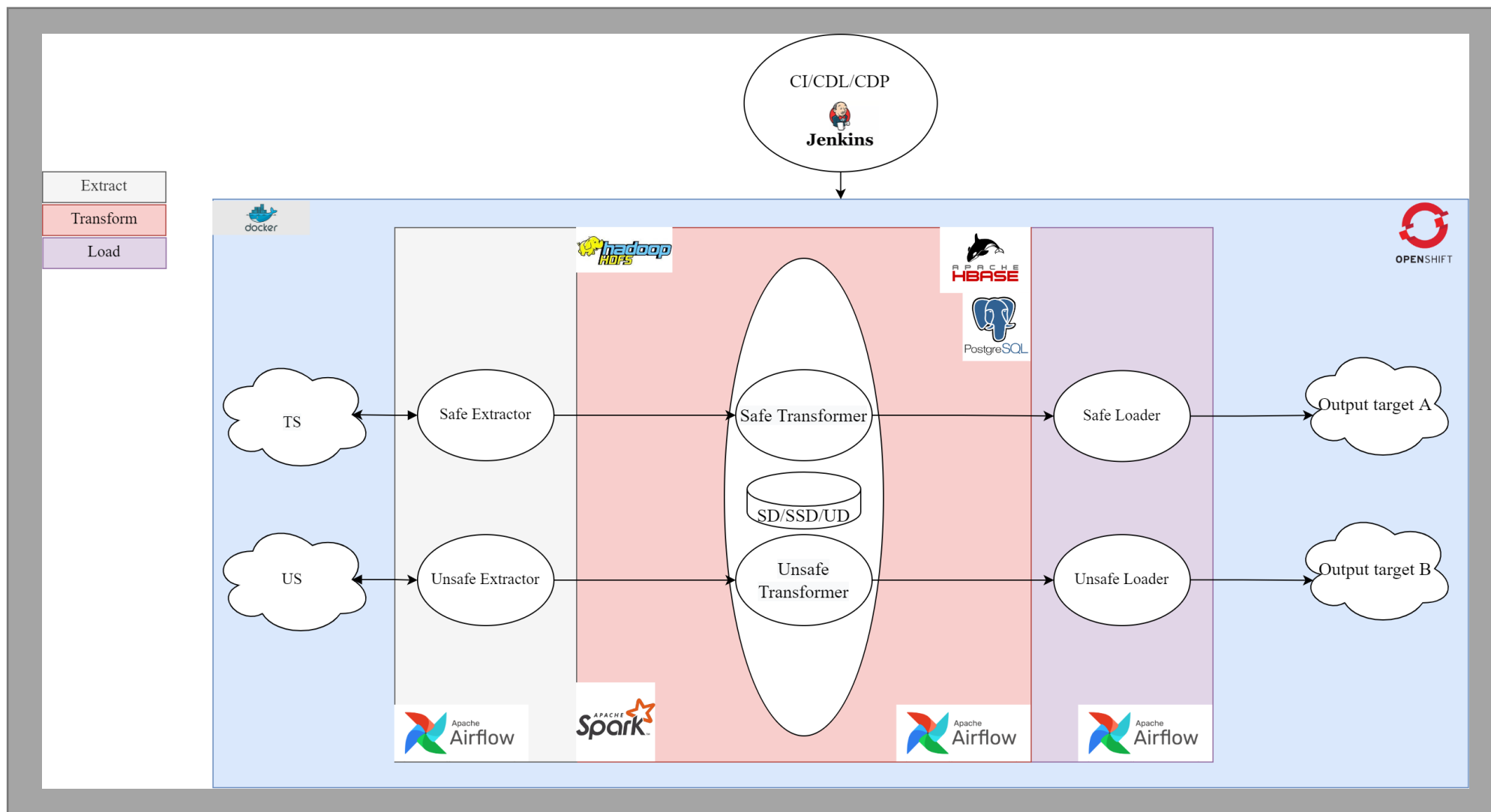


Projected data volume by 2025

2

DLAF schema containing possible implementation concepts and configurations
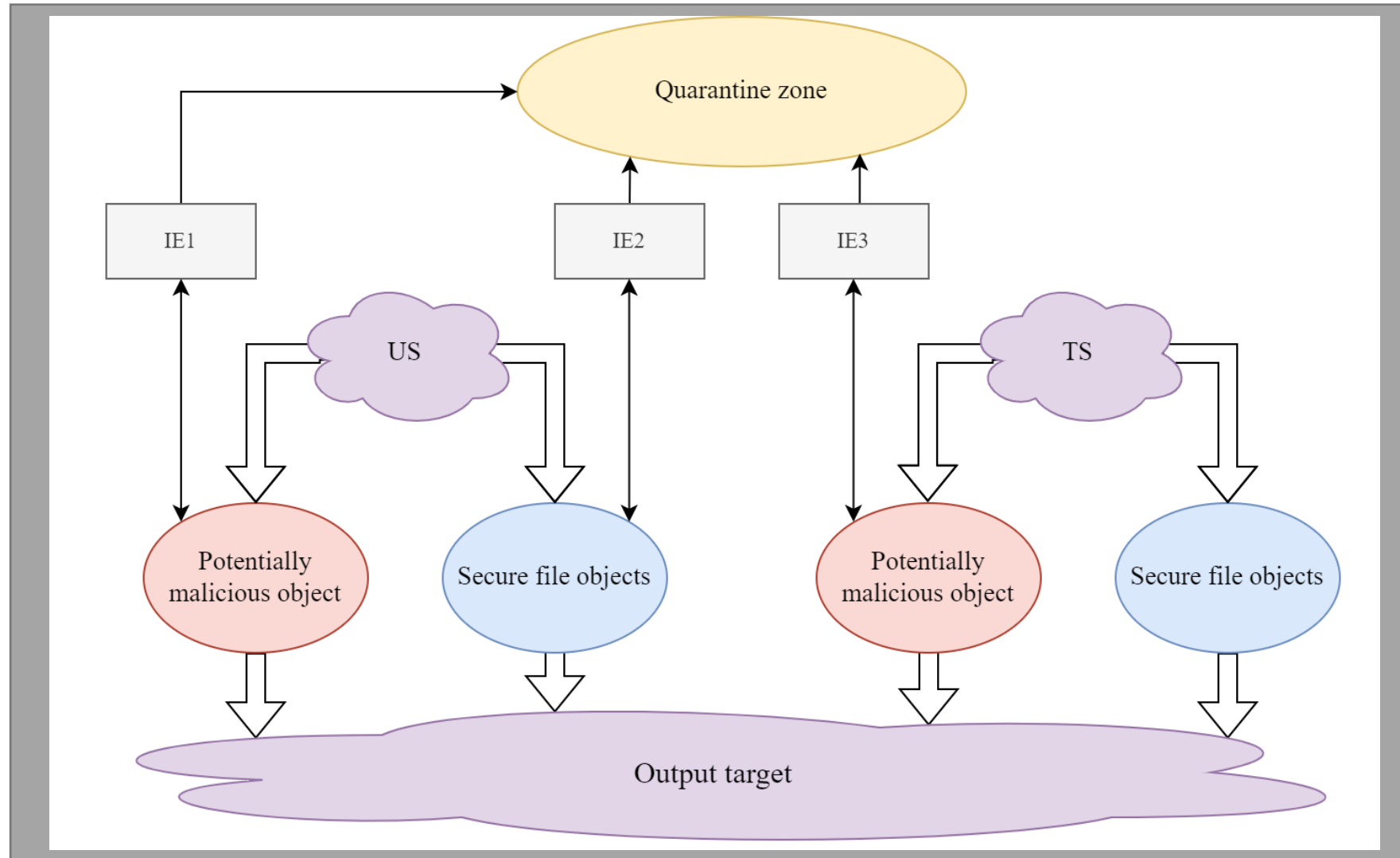


A DLAF framework consisting of nine aspects of data lakes to consider when creating a comprehensive data lake architecture
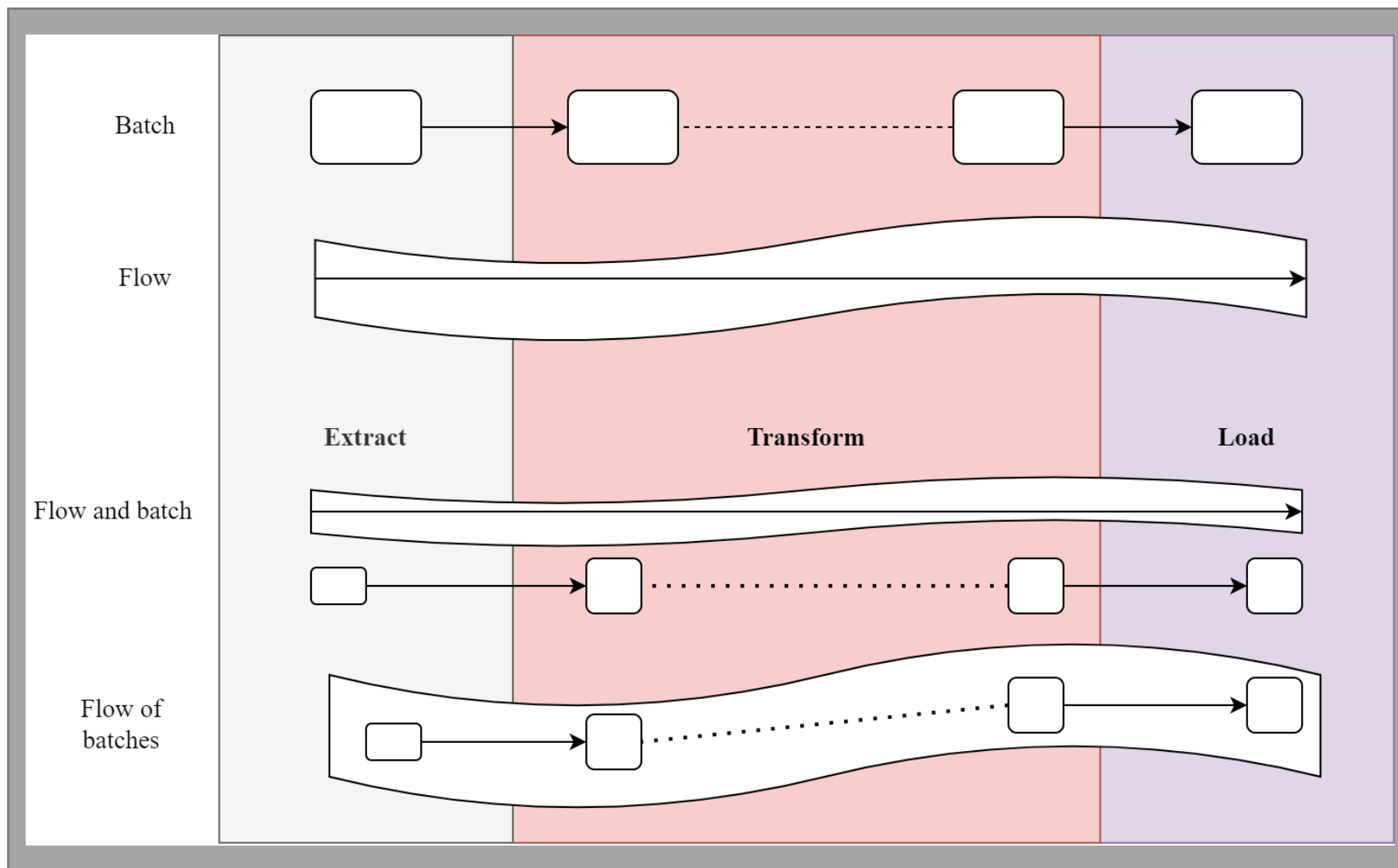
3

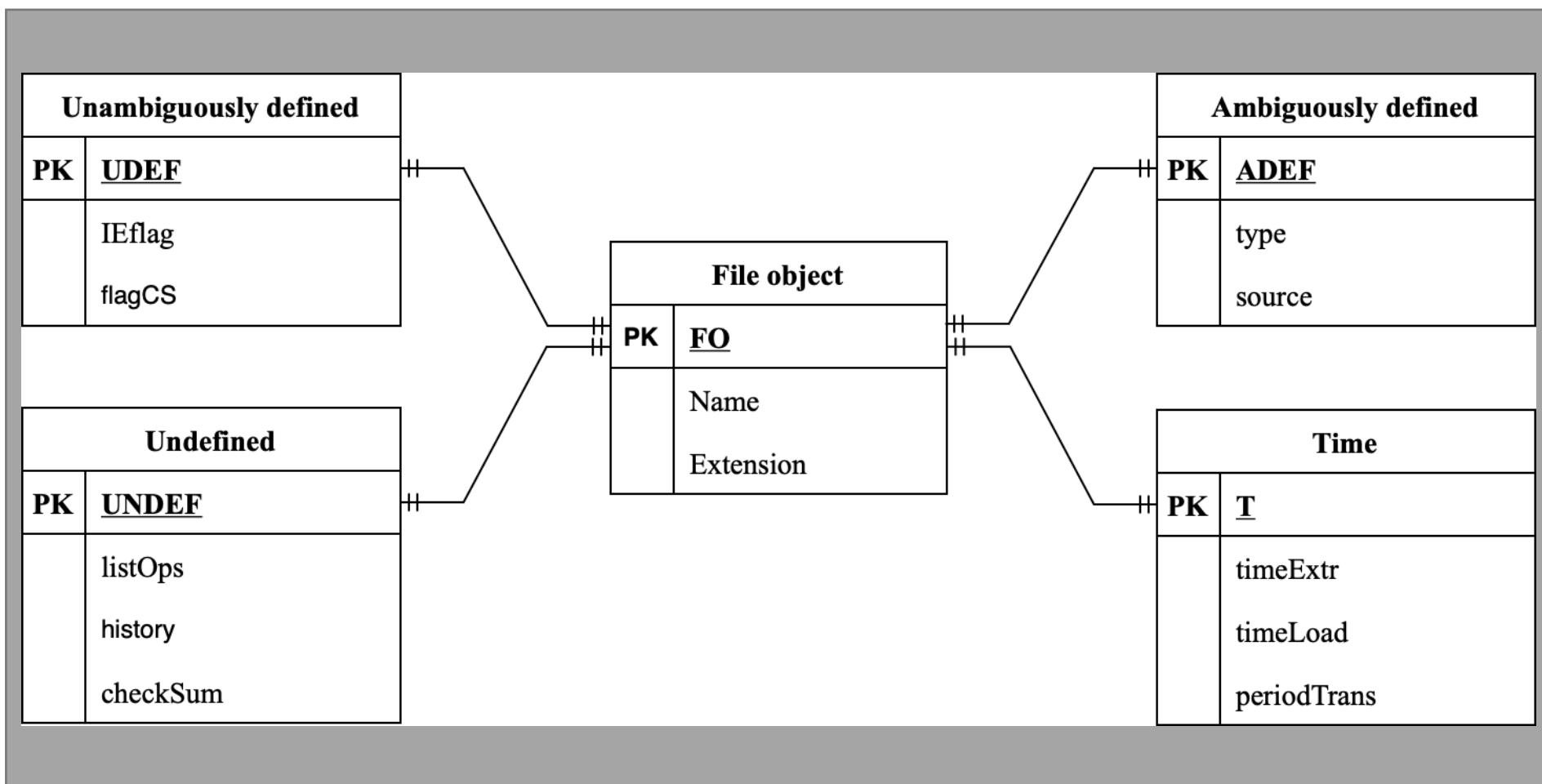**Secure data lake infrastructure**

Secure data lake repository

Secure data lake flow

Secure data lake model

Secure data lake organization

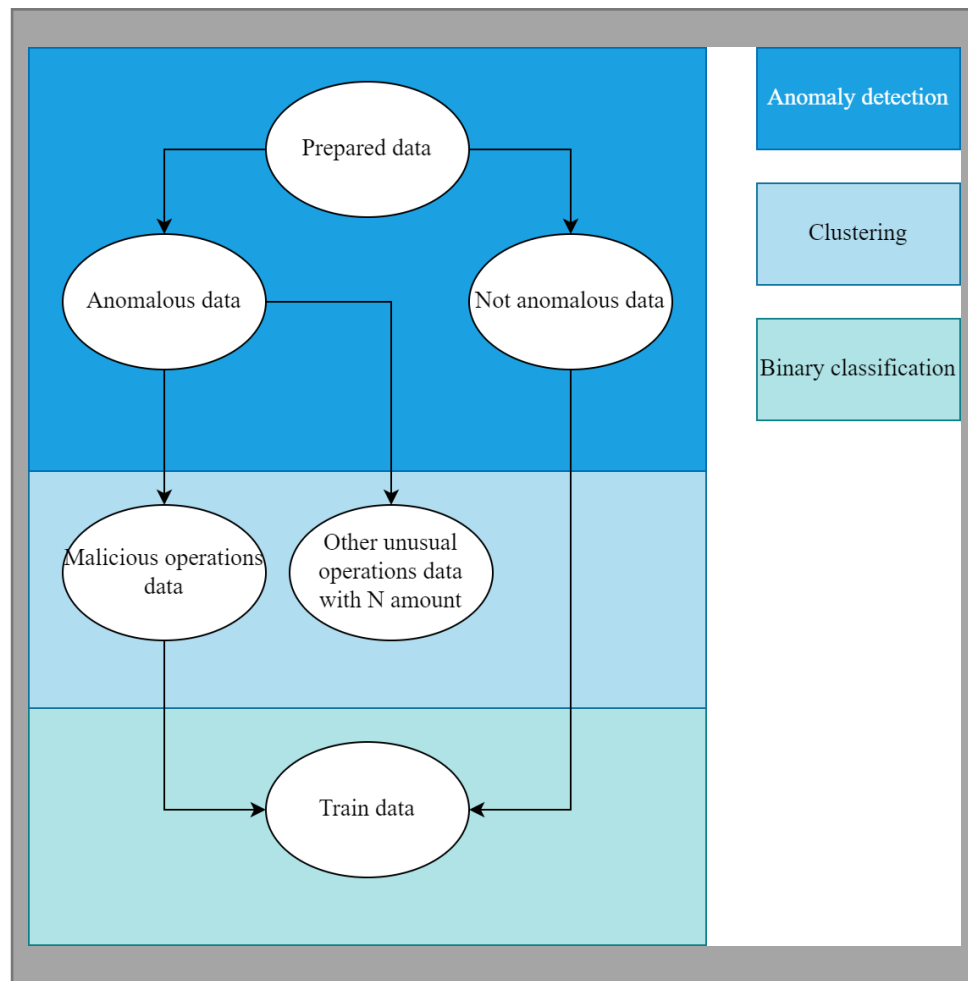**Secure data lake processing**

**Secure data lake metadata management**

There are **four** security models to choose from:
- ⬨ **DRBAC** - an implemented modified security model RBAC with an embedded Clark-Wilson discretionary integrity control model;

- ⬨ **MRBAC** - an implemented modified RBAC security model with Ken Beeb's inbuilt mandate model of integrity control;

- ⬨ **MABAC** - modified security model with ABAC;

- ⬨ **XACML** - generally accepted standard for ABAC implemented in OPC without any modifications.

All the above MBs interact with the **two** main components of a secure OPC architecture:
- ⬨ **GMT** - Global Monitoring Tool;

- ⬨ **CSC** - Check Sum Controller - check sum handler.

Prepared data

Anomalous data

Not anomalous data

Malicious operations data

Other unusual operations data with N amount
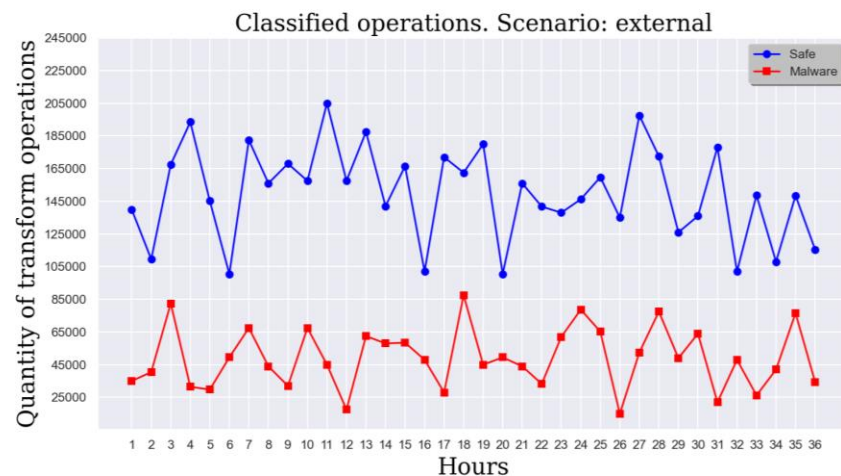
Train data

Anomaly detection

Clustering

Binary classification

**The process of generating a dataset**



Feature Maps

Feature Maps

Feature Maps

Feature Maps

D

Convolution + ReLu

Pooling

Convolution + ReLu

Pooling

Fully Connected Layers

Output Layer

**Architecture of CNN
classification of log code sequences**

Python

Caffe2    C++

Linux daemon

**GMT model development process**

12

Dashboard of the MABAC model in batch mode and emulation of external attack scenarios

The developed Secured Data Lake Architecture Framework (SDLAF) provides the ability to design secure data lakes **without losing key benefits** and while **maintaining flexibility** for a wide range of business requirements.

The next stage of work – **GMT 2.0**:

- Data trait engineering of monitoring and journaling systems;

- Automatic anomaly detection;

- Adaptation to **any log sequence** regardless of system.

Thanks for your attention