



MONITORING EVALUATION FOR LARGE-SCALE ENVIRONMENTS AND OPTIMIZING DATA SYSTEM HEALTH

Supervisor: Apostolos Karvelas

Anastasiia Petrovych

AGENDA OVERVIEW

01

PROJECT BACKGROUND

02

PROBLEM STATEMENT

03

FRAMEWORK

04

METHODOLOGY

05

PIPELINE

06

PIPELINE COMPONENTS

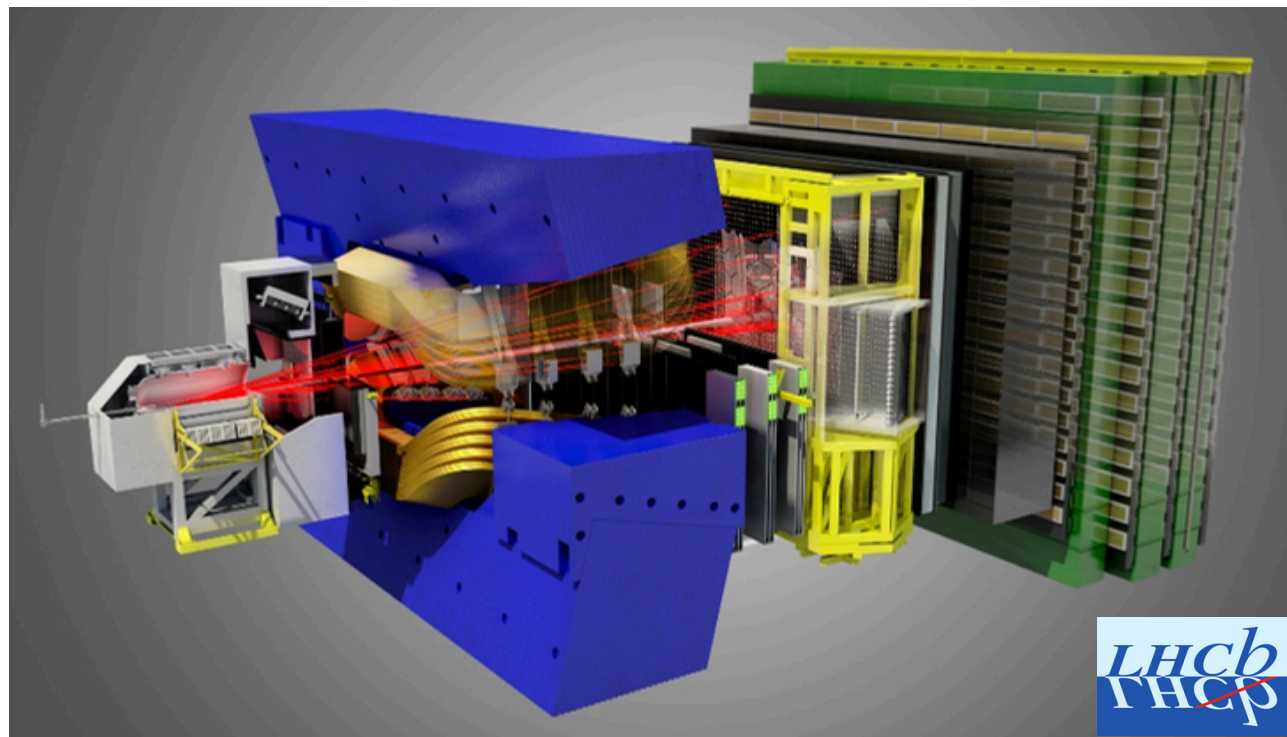
07

VISUAL REPRESENTATION

08

CONCLUSION

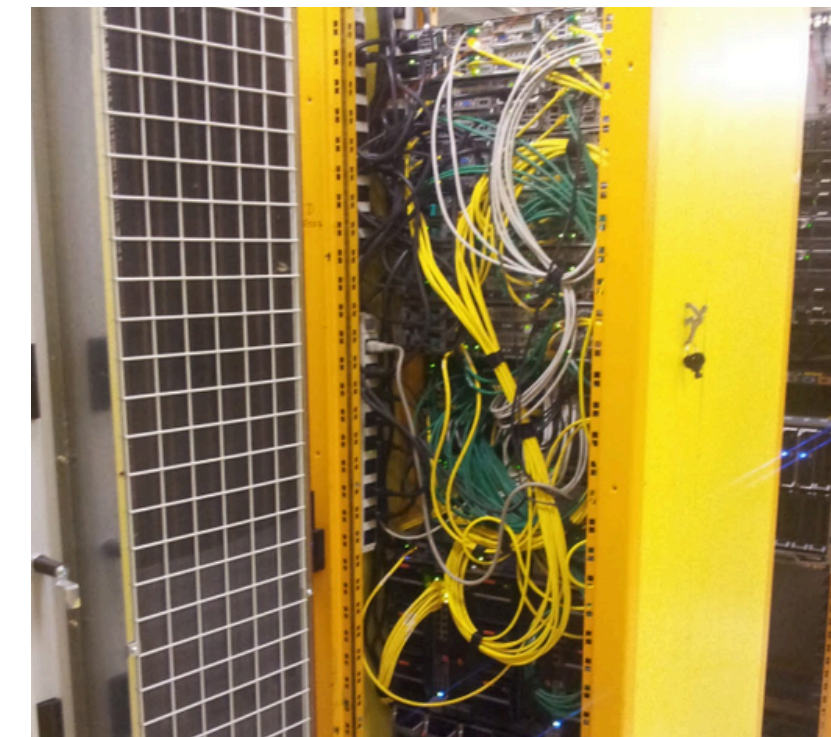
PROJECT BACKGROUND



LHCb experiment



Data centers



Data farms

PROBLEM STATEMENT

How to detect anomalies in servers?

Problems

Hardware Failures

Software Issues

Resource Constraints

Environmental Factors



Consequences

Data Loss

Downtime

Data Inconsistency

Inefficient use of
resources

FRAMEWORK

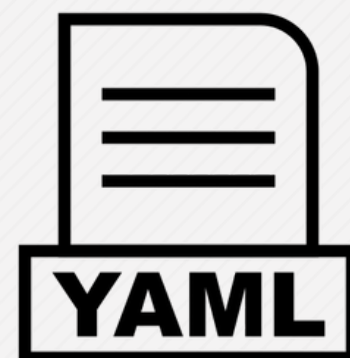


Kubeflow

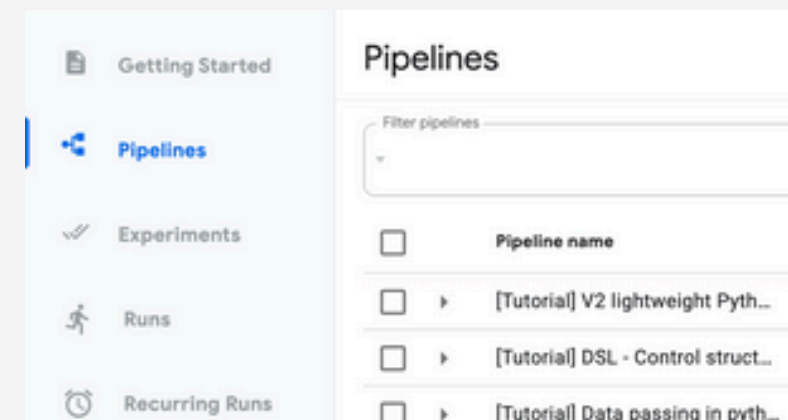
an open-source platform for machine learning and MLOps on Kubernetes



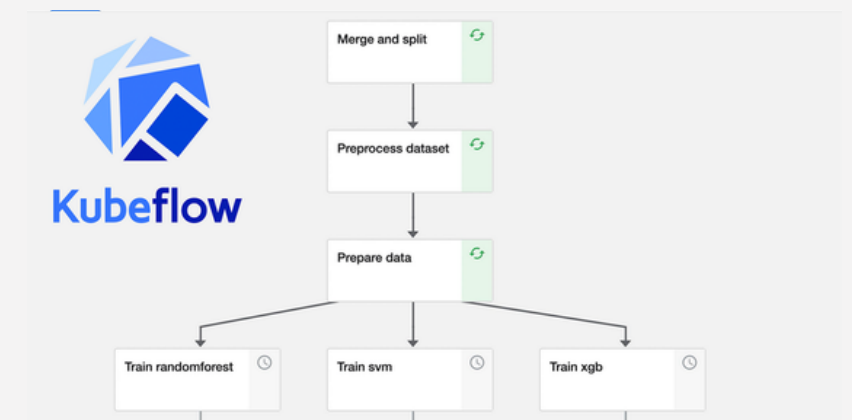
1. Define pipeline components



2. Compile to YAML file

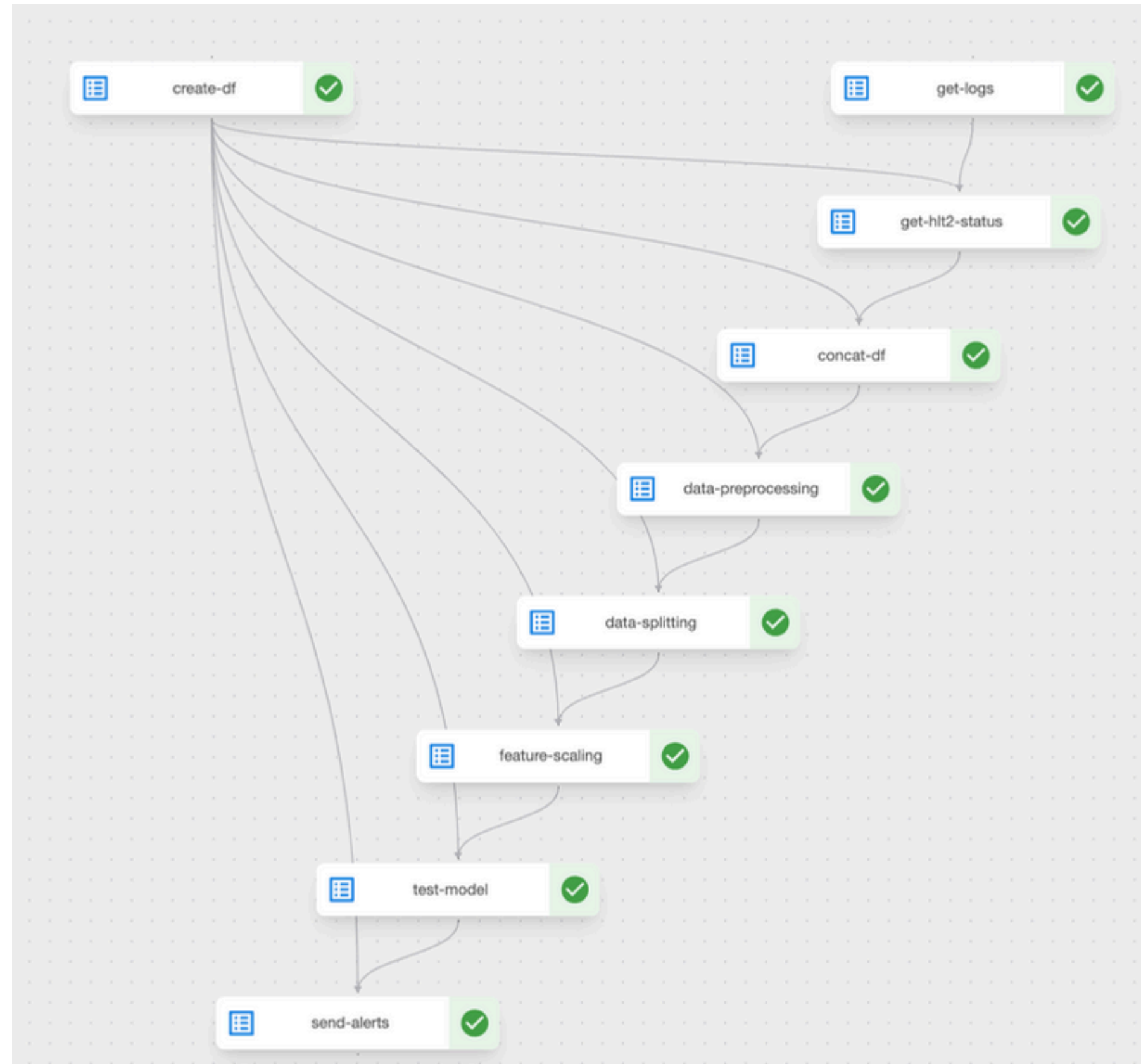


3. Create pipeline in Kubeflow Central Dashboard

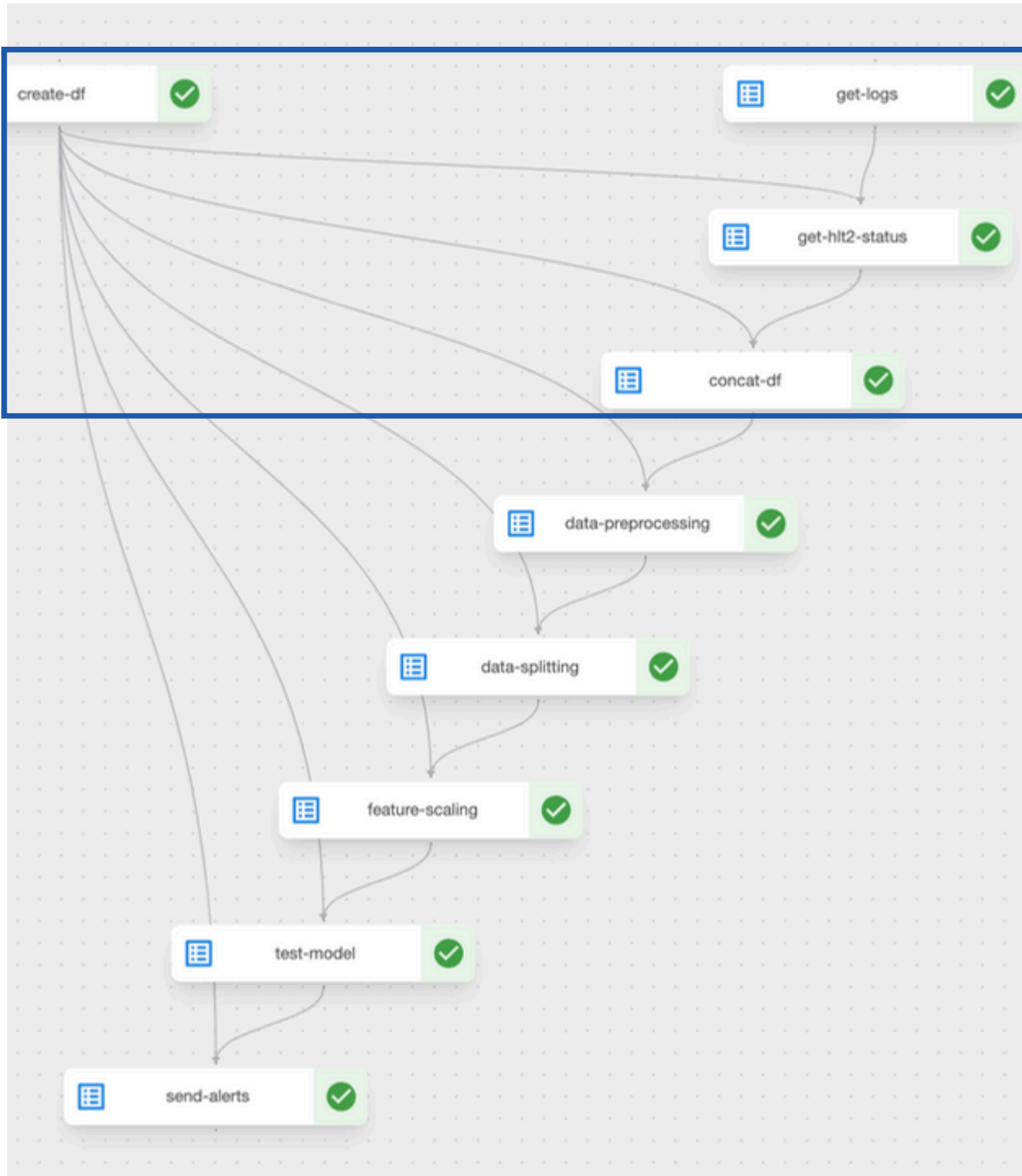


4. Run recurrent calls


PIPELINE



DATA COLLECTION



Prometheus



```
HLT20319_2> in state <PAUSED>
HLT20303_2> in state <PAUSED>
HLT20329_2> in state <PAUSED>
HLT20333_2> in state <PAUSED>
```

Logs

event monitoring
and alerting

Metrics

CPU Utilization
Disk I/O
Memory
Network
System Overview

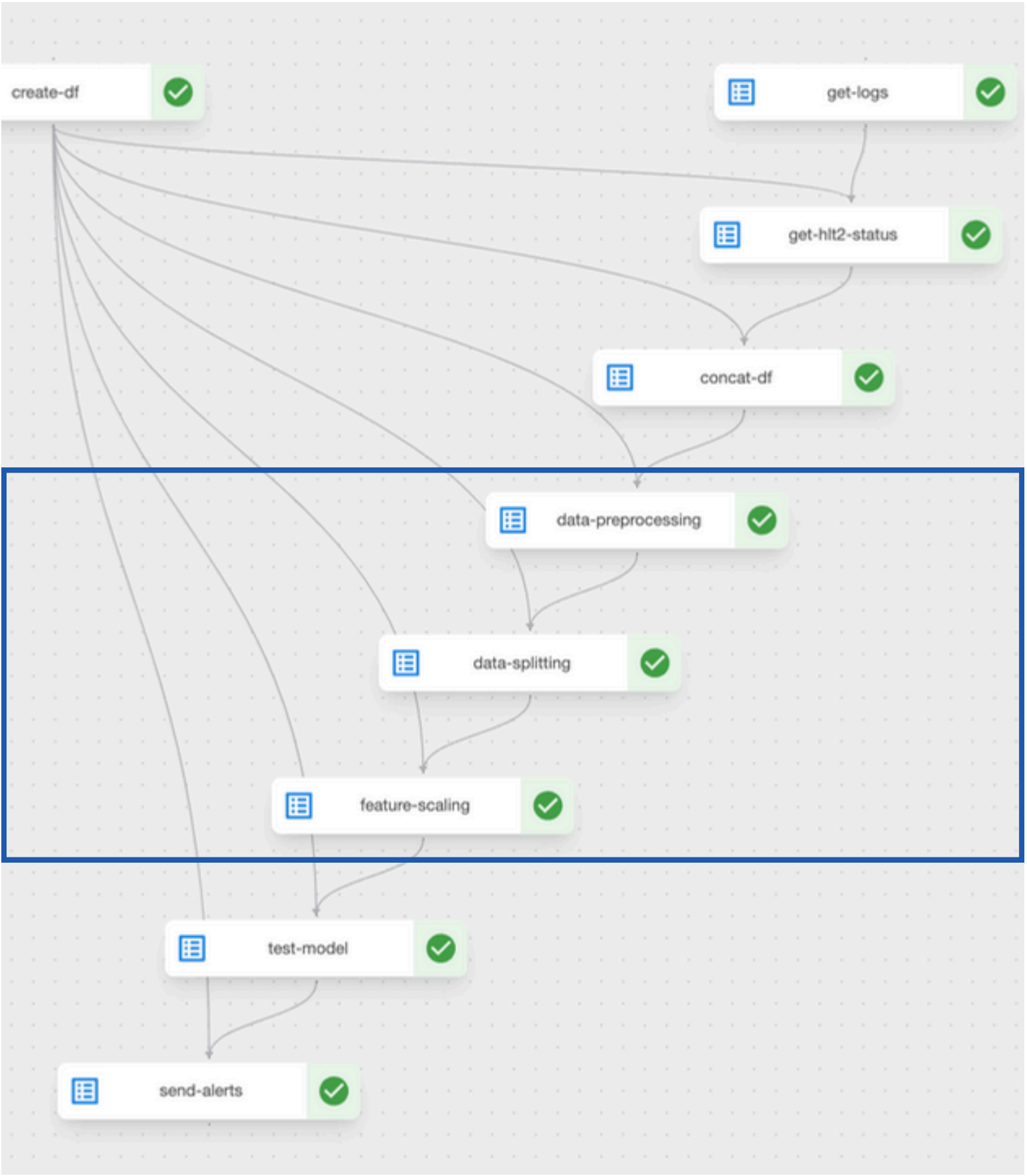
DATA PREPROCESSING

Fill missing values and
choose important
feature

Split dataset into train
and validation

Feature scaling

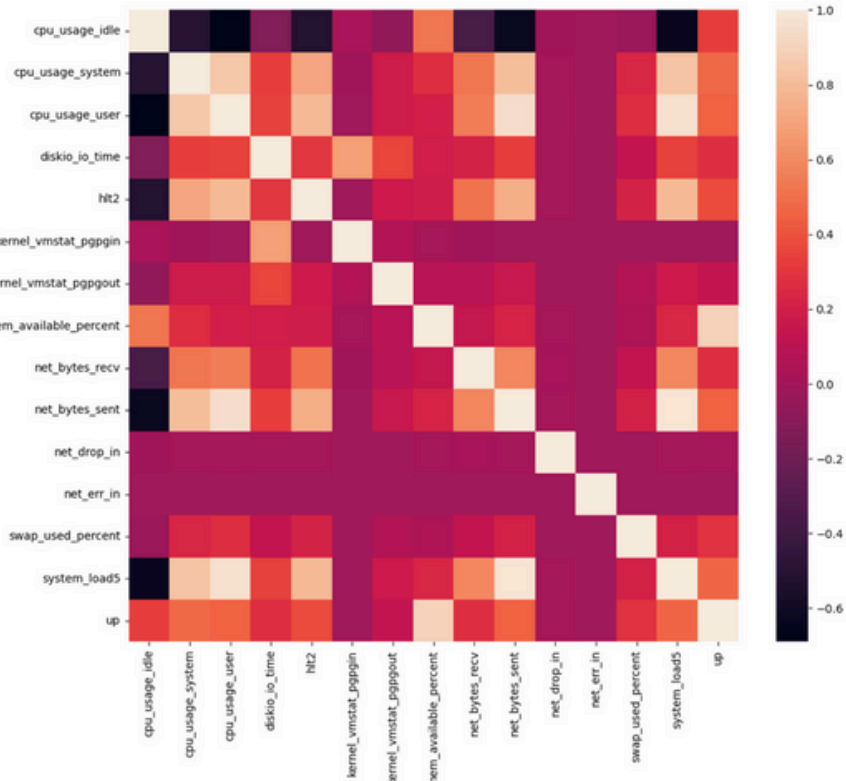
DATA PREPROCESSING



Fill missing values and choose important feature

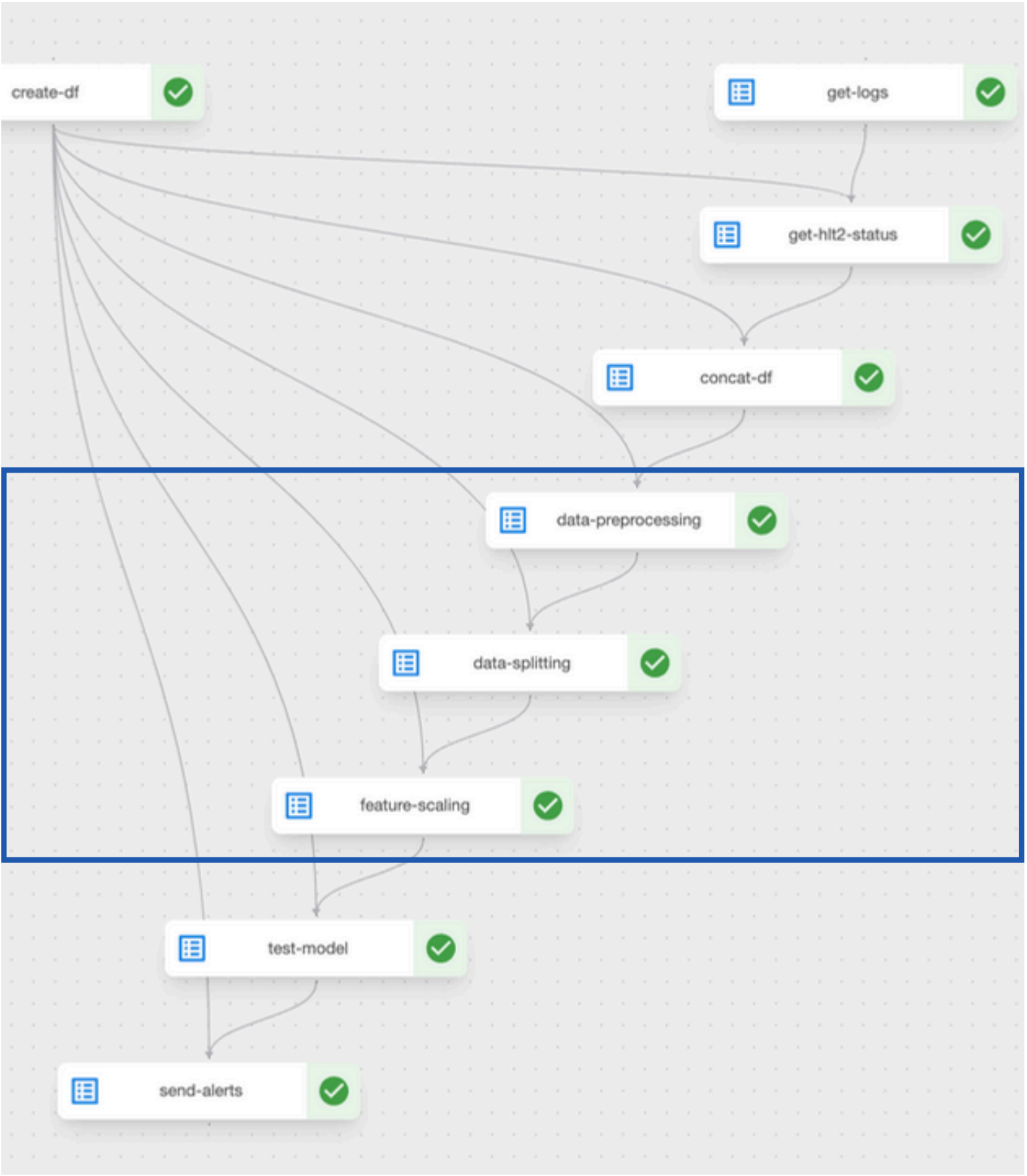
Split dataset into train and validation

Feature scaling



Correlation matrix

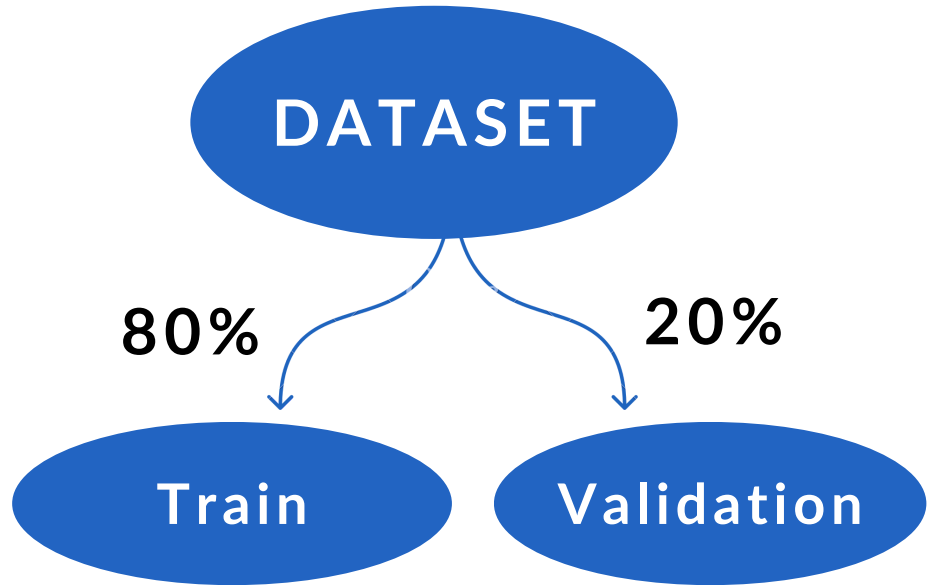
DATA PREPROCESSING



Fill missing values and choose important feature

Split dataset into train and validation

Feature scaling

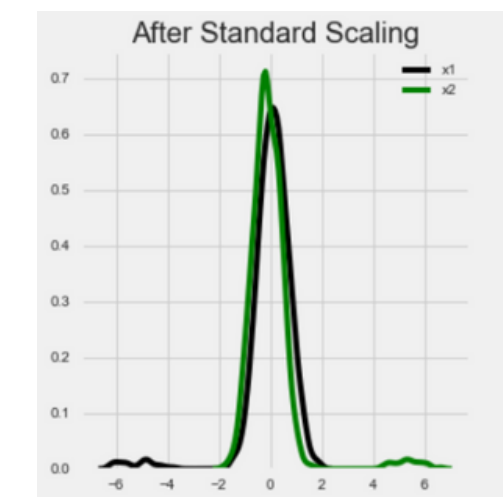
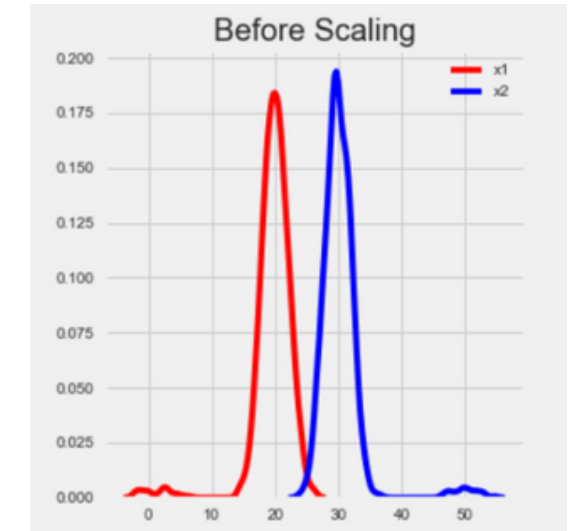


DATA PREPROCESSING

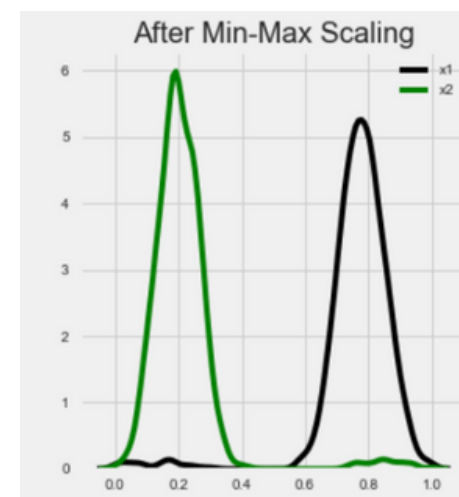
Fill missing values and
choose important
feature

Split dataset into train
and validation

Feature scaling

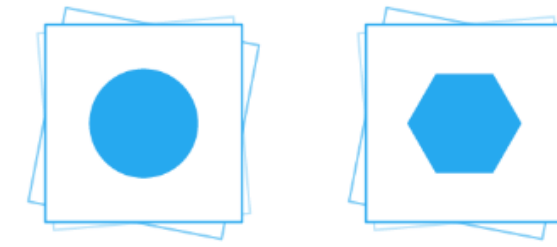


Standard scaler



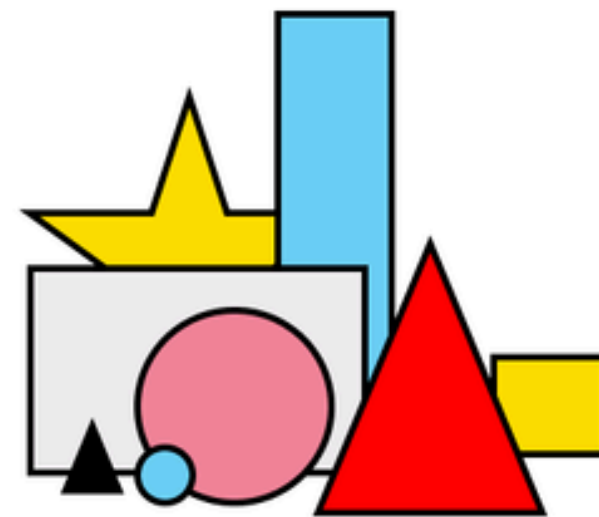
Min-Max scaler

MODEL TRAINING

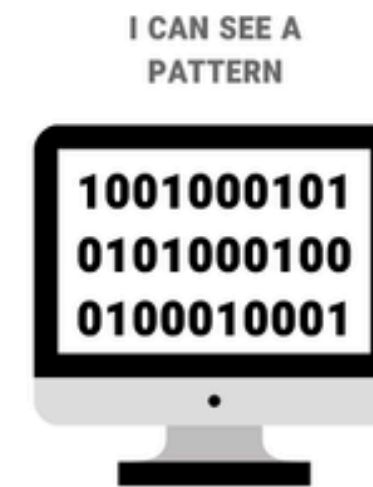


??
unlabeled data
??

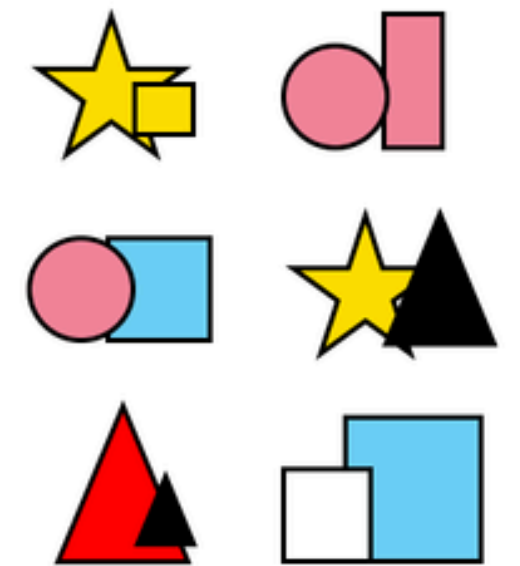
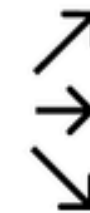
↓
Unsupervised learning



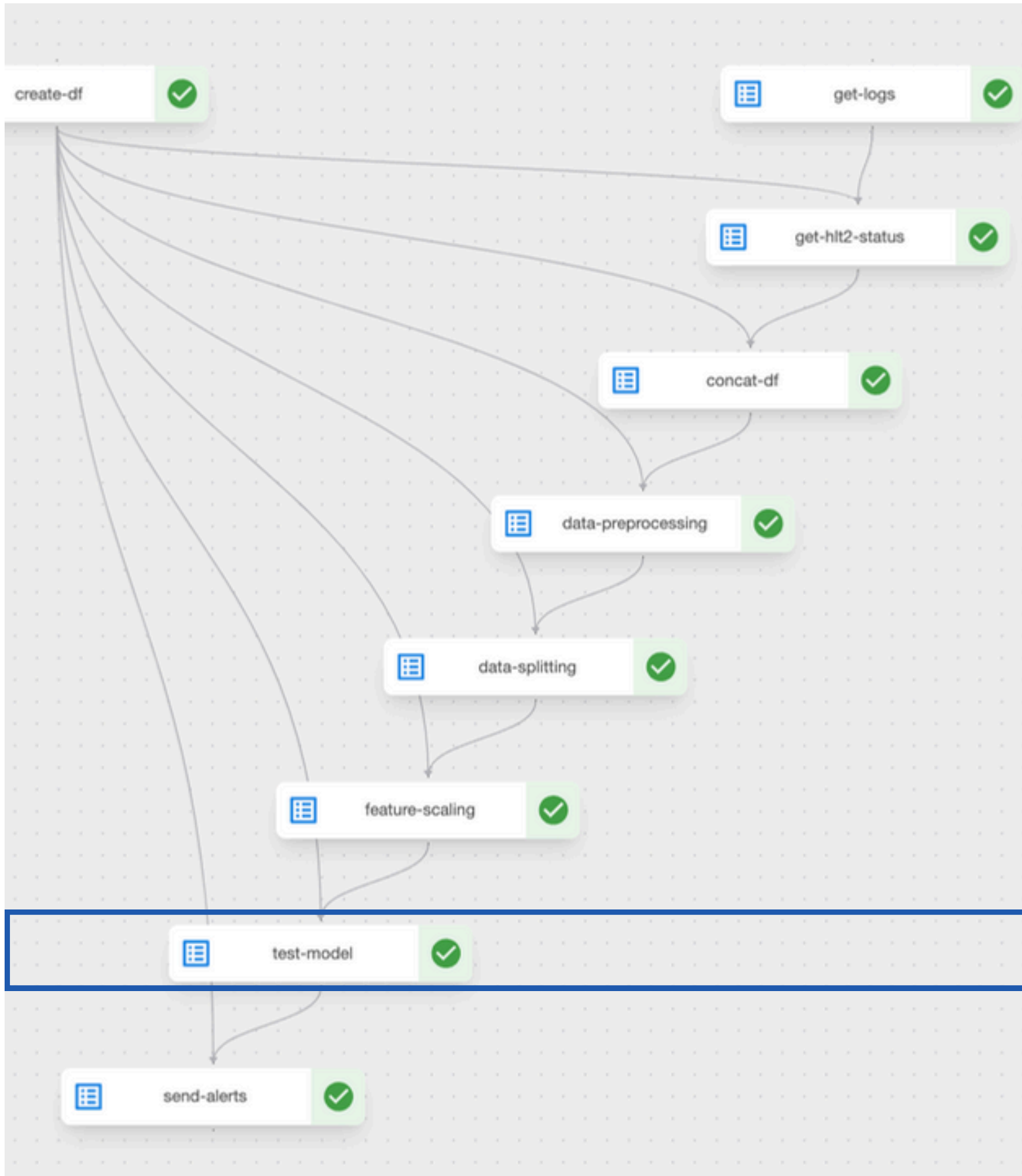
INPUT



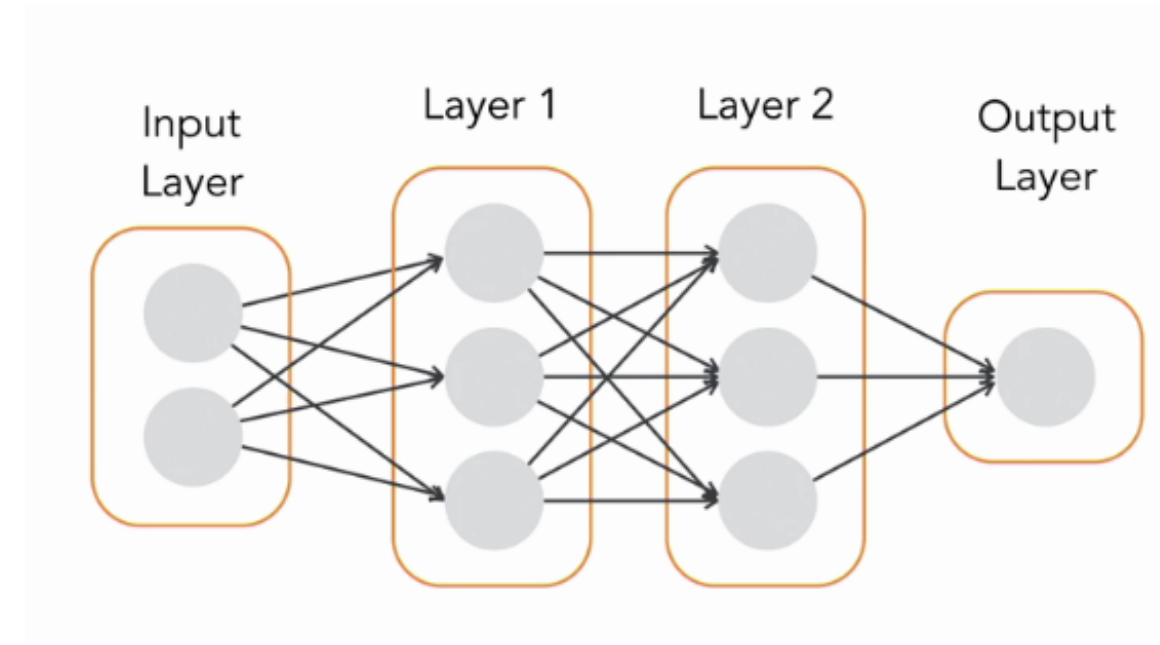
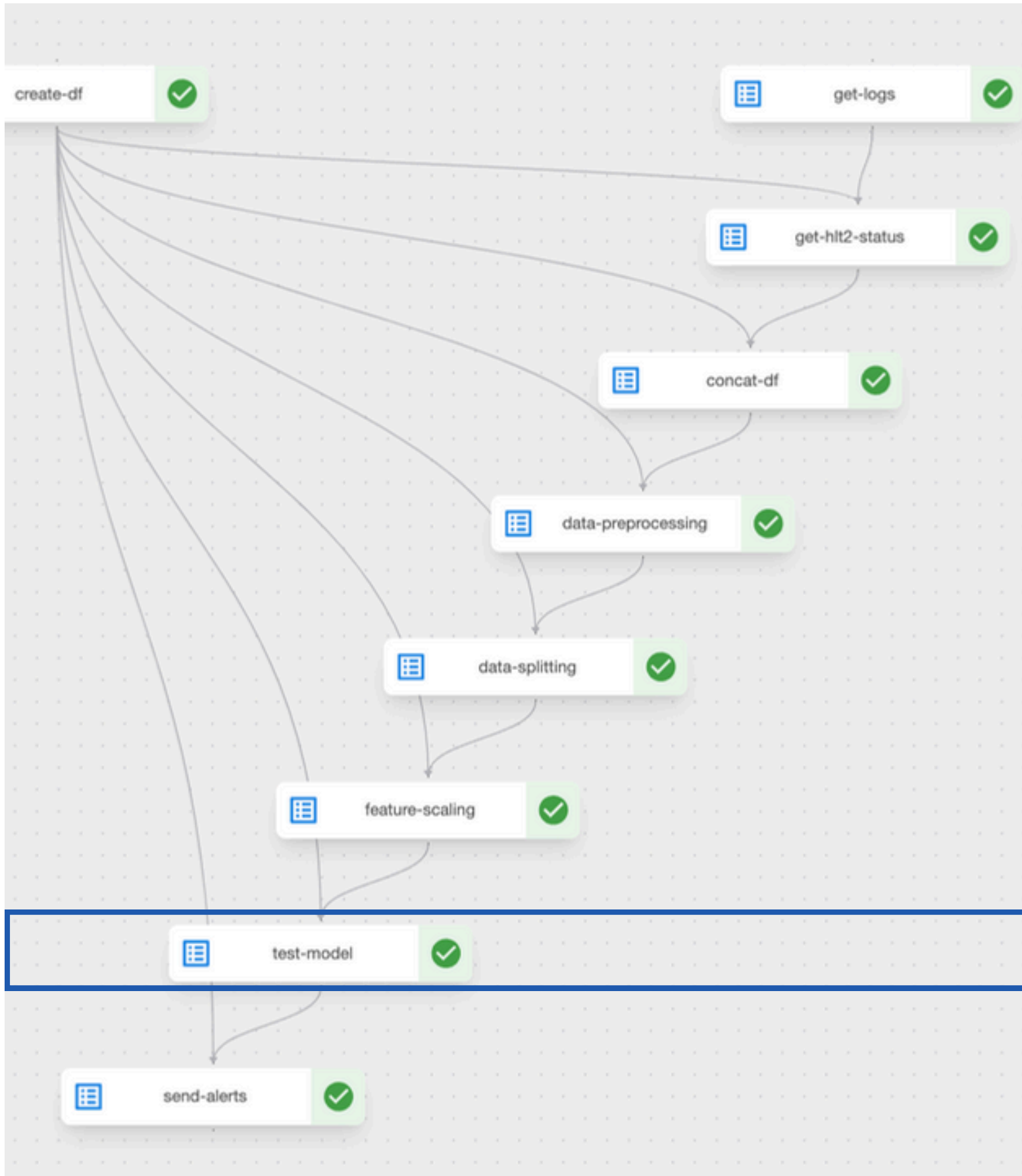
MODEL



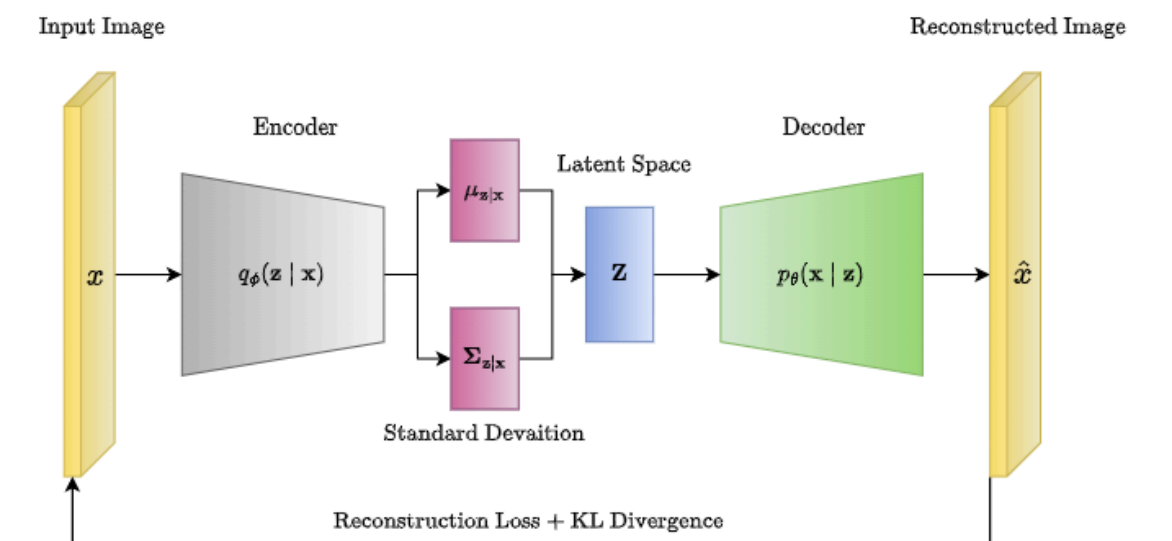
OUTPUT



MODEL TRAINING



Sequential model



VAE

TESTING

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

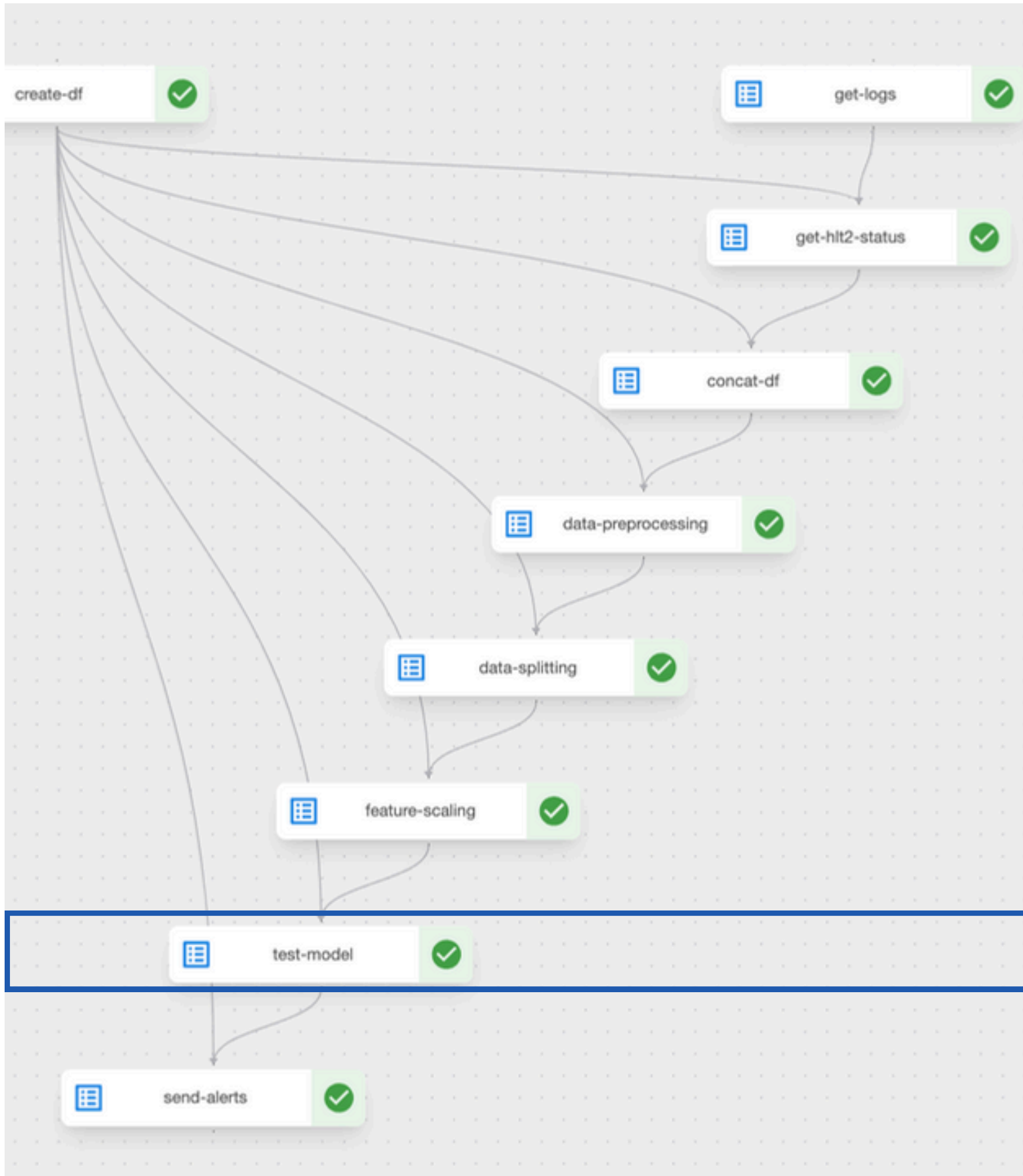
$$MSE = \frac{\sum (y_i - \hat{y}_i)^2}{n}$$

where,

n : number of observation

y_i : the actual value of the i^{th} observation

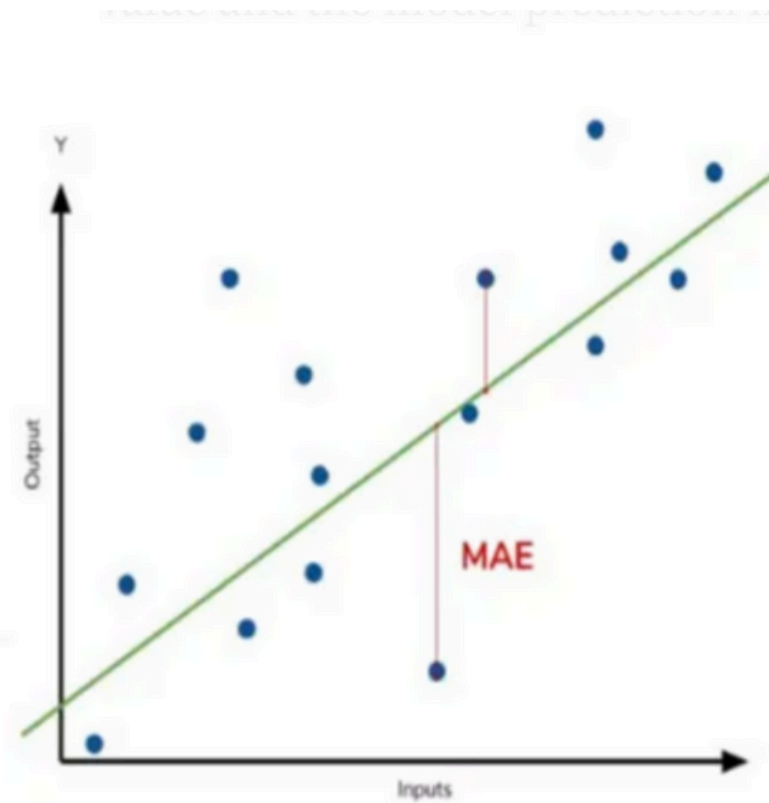
\hat{y}_i : the predicted value of the i^{th} observation



TESTING

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$MSE = \frac{\sum (y_i - \hat{y}_i)^2}{n}$$

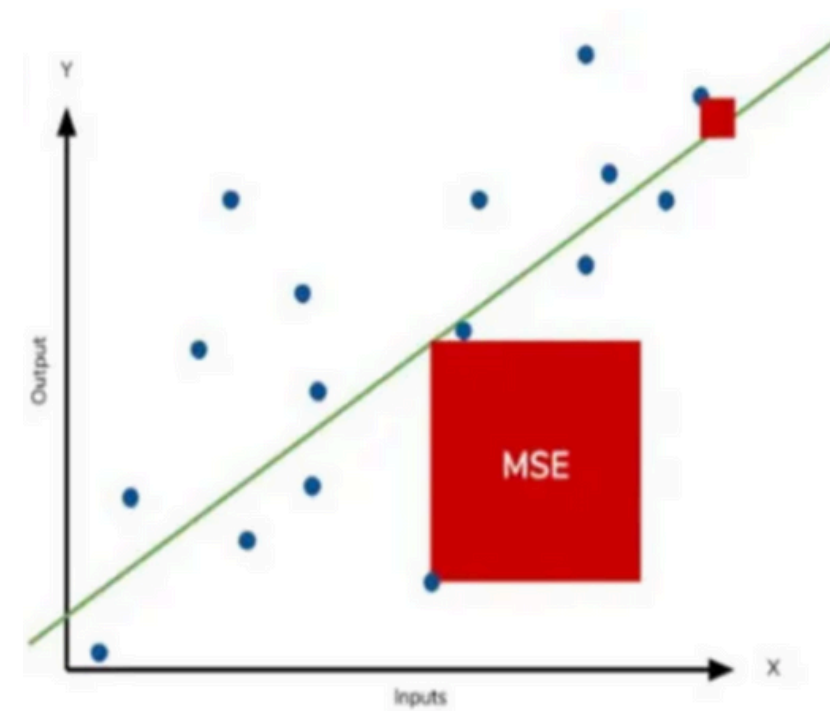


difference between the actual value and the model prediction over the entire data set

TESTING

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$MSE = \frac{\sum (y_i - \hat{y}_i)^2}{n}$$



the average *squared* difference between the estimated values and the actual value

DETECT ANOMALIES

Reconstruct loss on
predicted data

Find threshold as 99.9th
percentile

Identify anomalies

DETECT ANOMALIES

Reconstruct loss on predicted data

$$L = \frac{1}{n} \sum_{i=1}^n |y_i - \tilde{y}_i|$$

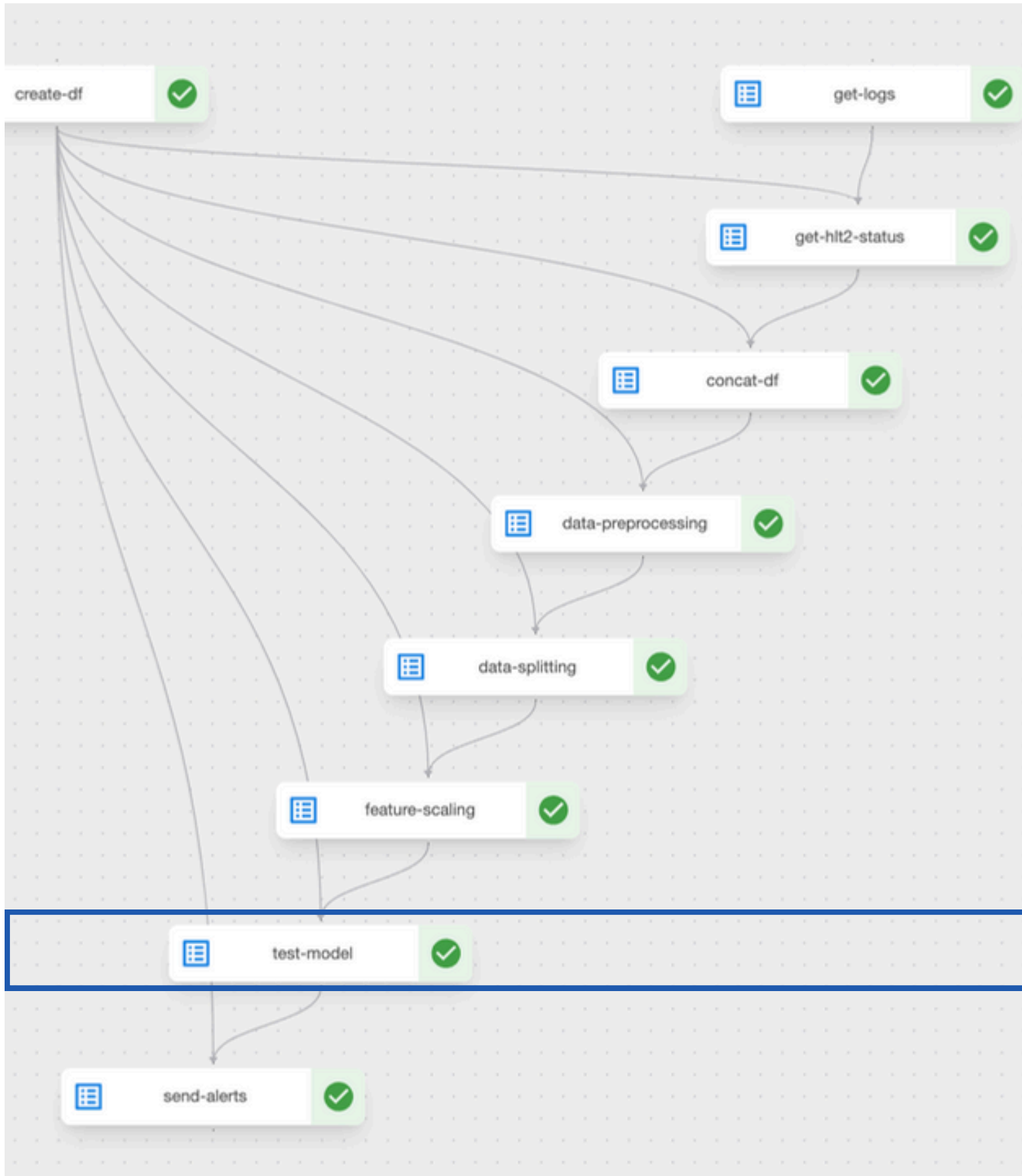
y_i is the original data point.

\tilde{y}_i is the reconstructed data point.

$|y_i - \tilde{y}_i|$ is the absolute difference

Find threshold as 99.9th percentile

Identify anomalies

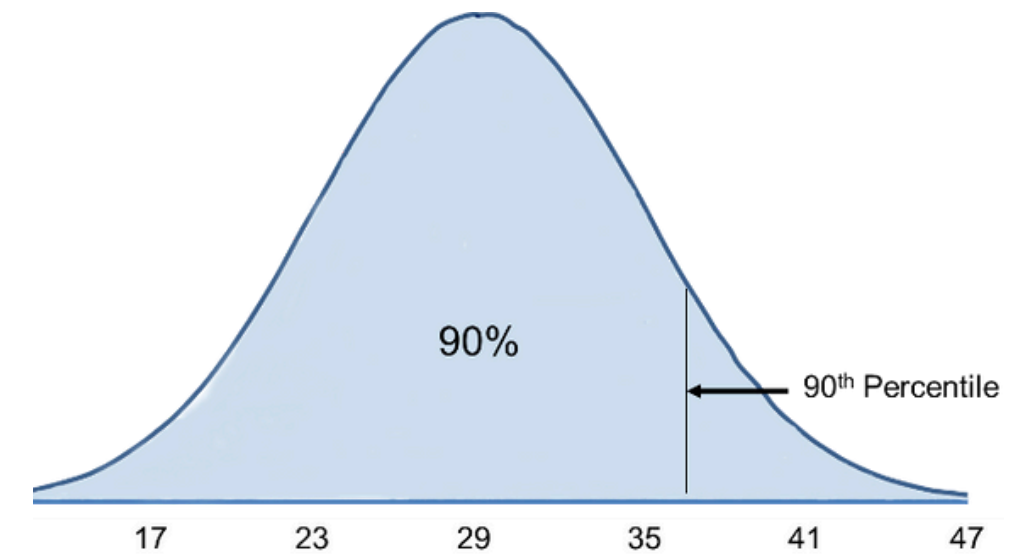


DETECT ANOMALIES

Reconstruct loss on predicted data

Find threshold as 99.9th percentile

Identify anomalies

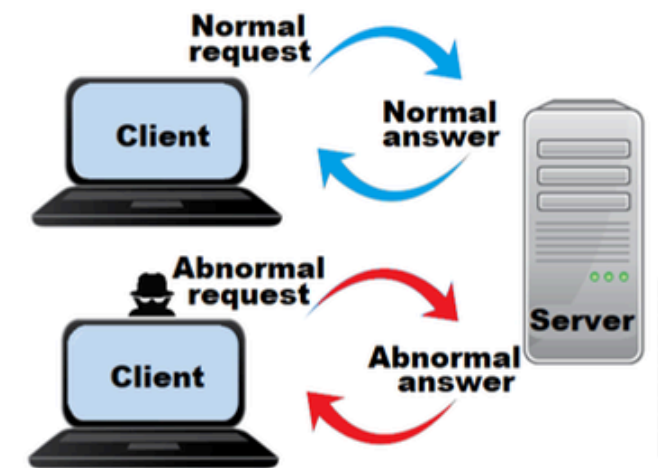


DETECT ANOMALIES

Reconstruct loss on predicted data

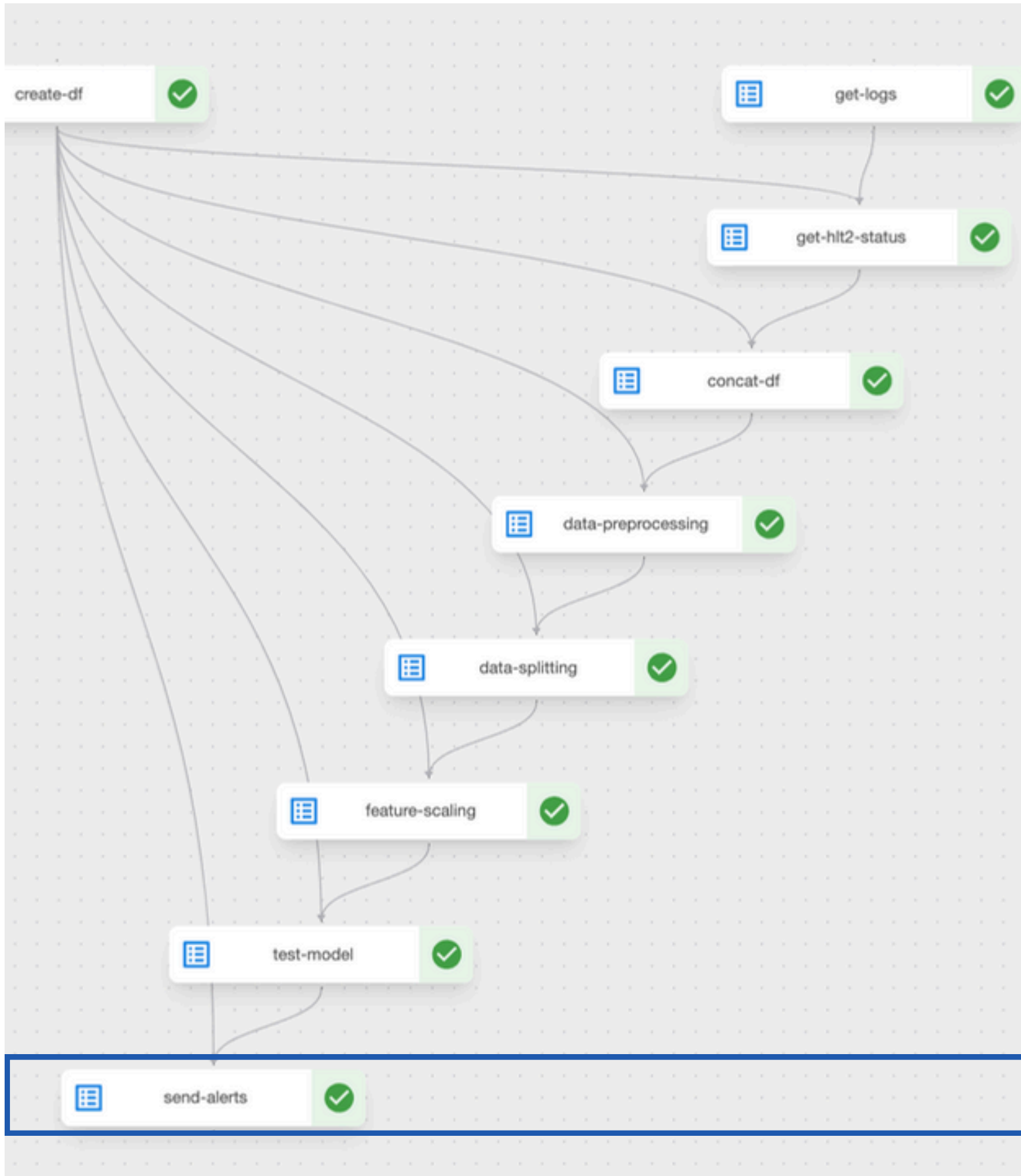
Find threshold as 99.9th percentile

Identify anomalies



when the reconstruction loss exceeds the threshold

SEND ALERTS



Anomaly Detection Report

1 message

anastasiia.petrovych@cern.ch <anastasiia.petrovych@cern.ch>
To: anastasiyapetrovych25@gmail.com

Dear Colleague,

This is the detected anomaly report for the past day. Please find the details below:

Anomalies detected:

- Node n2011704 has 1 windows with high anomaly scores.
- Node n2012503 has 1 windows with high anomaly scores.
- Node n2020501 has 1 windows with high anomaly scores.
- Node n2020504 has 1 windows with high anomaly scores.
- Node n2020901 has 1 windows with high anomaly scores.
- Node n2022501 has 2 windows with high anomaly scores.
- Node n2022504 has 1 windows with high anomaly scores.
- Node n2022901 has 1 windows with high anomaly scores.
- Node n2024101 has 1 windows with high anomaly scores.
- Node n2024104 has 1 windows with high anomaly scores.
- Node n2024304 has 1 windows with high anomaly scores.
- Node n2024503 has 2 windows with high anomaly scores.
- Node n2040102 has 1 windows with high anomaly scores.
- Node n2040104 has 1 windows with high anomaly scores.
- Node n2040304 has 1 windows with high anomaly scores.
- Node n2040703 has 1 windows with high anomaly scores.
- Node n2041501 has 1 windows with high anomaly scores.
- Node n2041503 has 1 windows with high anomaly scores.
- Node n2044302 has 4 windows with high anomaly scores.
- Node n2052504 has 1 windows with high anomaly scores.
- Node n2053901 has 3 windows with high anomaly scores.
- Node n2054301 has 5 windows with high anomaly scores.
- Node n2062502 has 1 windows with high anomaly scores.
- Node n2062702 has 3 windows with high anomaly scores.

Regards,
LHCb team.



Recurring Runs

Status	Trigger
ENABLED	Every 1 days



CONCLUSION

- Automated the anomaly detection process and increased efficiency
- Improved the accuracy and reliability of the whole pipeline
- Implemented daily monitoring for continuous detection and timely alerts



THANK YOU!

Anastasiia Petrovych