# Use of AI for Log Analysis in CI/CD Pipelines

Bachelor Thesis - Defence

MaidAlišić

July 23, 2025

FH Oberösterreich · Campus Hagenberg

Supervisor: FH-Prof. Dipl.-Ing. Dr. Stefan Wagner

## Road map

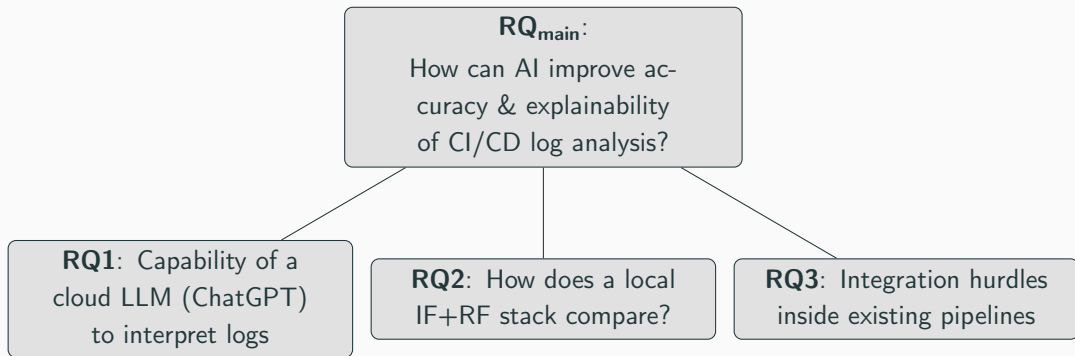Research questions

Problem context

Method

Architecture

Data & evaluation

Results

Impact

# Research questions

RQ_main:
How can AI improve accuracy & explainability of CI/CD log analysis?

RQ1: Capability of a cloud LLM (ChatGPT) to interpret logs

RQ2: How does a local IF+RF stack compare?

RQ3: Integration hurdles inside existing pipelines

# Problem context

- CI/CD emits $\approx$ 10-20 GB of build, test & deploy logs *per day*.

## Why do we care?

- CI/CD emits ≈ 10-20 GB of build, test & deploy logs *per day*.
- Manual `grep` slows the merge queue; critical faults slip through.

- CI/CD emits $\approx$ 10-20 GB of build, test & deploy logs *per day*.
- Manual `grep` slows the merge queue; critical faults slip through.
- Business **Service-Level Objective**: feedback within $\leq$ 200 ms per pipeline.

- CI/CD emits ≈ 10-20 GB of build, test & deploy logs *per day*.
- Manual `grep` slows the merge queue; critical faults slip through.
- Business **Service-Level Objective**: feedback within ≤ 200 ms per pipeline.
- Logs may expose customer IDs, therefore they must remain on-premises (no cloud export).

1. **Context-sensitivity** - identical tokens can be harmless or fatal.

## Operational pain points

1. **Context-sensitivity** - identical tokens can be harmless or fatal.
2. **Concept drift** - each merge may rename tests or switches.

## Operational pain points

1. **Context-sensitivity** - identical tokens can be harmless or fatal.
2. **Concept drift** - each merge may rename tests or switches.
3. **Latency pressure** - analysis must finish before runner teardown.

1. **Context-sensitivity** - identical tokens can be harmless or fatal.
2. **Concept drift** - each merge may rename tests or switches.
3. **Latency pressure** - analysis must finish before runner teardown.
4. **Alert fatigue** - regex rule sets grow without bound.

# Method

1. **Normalise** - strip timestamps, colours, IDs.

```
┌──────────┐        ┌───────────────┐
│ log line │  ───▶  │ sparse vector │
└──────────┘        └───────────────┘
```

1. **Normalise** - strip timestamps, colours, IDs.
2. **Tokenise** into uni- and bi-grams.

log line $\longrightarrow$ sparse vector

1. **Normalise** - strip timestamps, colours, IDs.
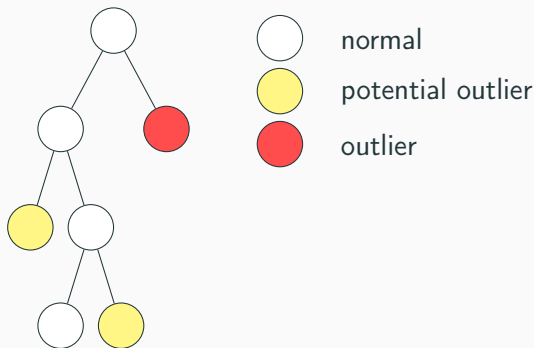2. **Tokenise** into uni- and bi-grams.
3. Weight with TF-IDF.

```
log line  ──────▶  sparse vector
```

1. **Normalise** - strip timestamps, colours, IDs.
2. **Tokenise** into uni- and bi-grams.
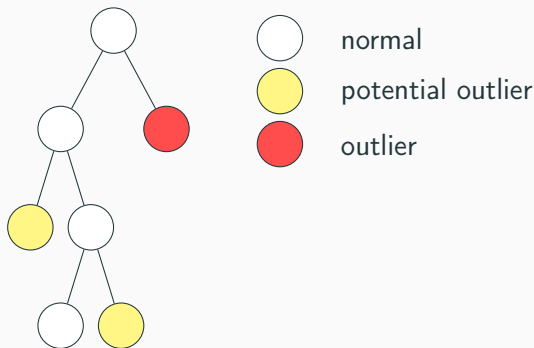3. Weight with TF-IDF.
4. Produce sparse vector.

log line → sparse vector

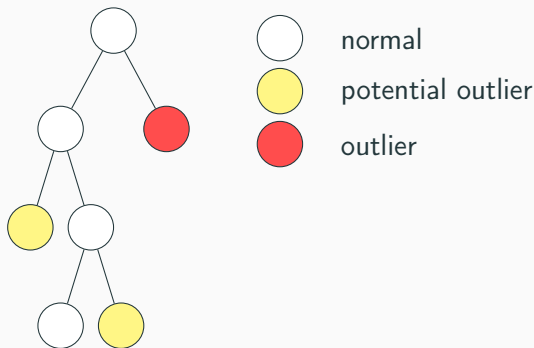- Random binary partitioning isolates unusual lines in fewer splits.



normal

potential outlier

outlier

- Random binary partitioning isolates unusual lines in fewer splits.
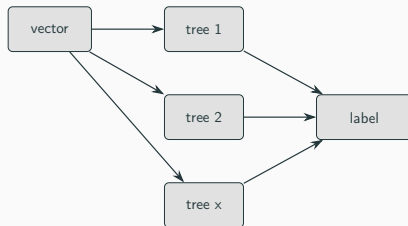- Score $s(x) = 2^{-h(x)/c(n)} \in [0,1]$ if high $\rightarrow$ outlier.



normal

potential outlier

outlier

- Random binary partitioning isolates unusual lines in fewer splits.
- Score $s(x) = 2^{-h(x)/c(n)} \in [0,1]$ if high $\rightarrow$ outlier.
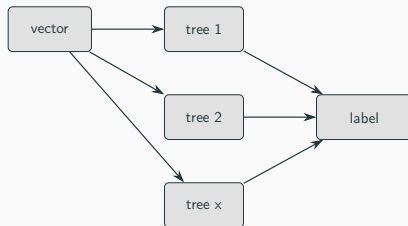- CPU-only: $\approx 30\,\mu s$ per line.



normal

potential outlier

outlier

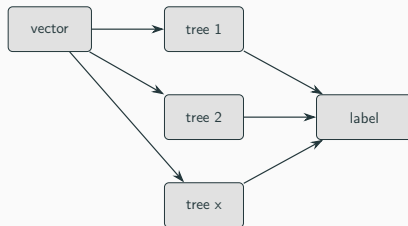- Maps each flagged line to a domain-specific error category.

# Random Forest ③ - error labelling

- Maps each flagged line to a domain-specific error category.
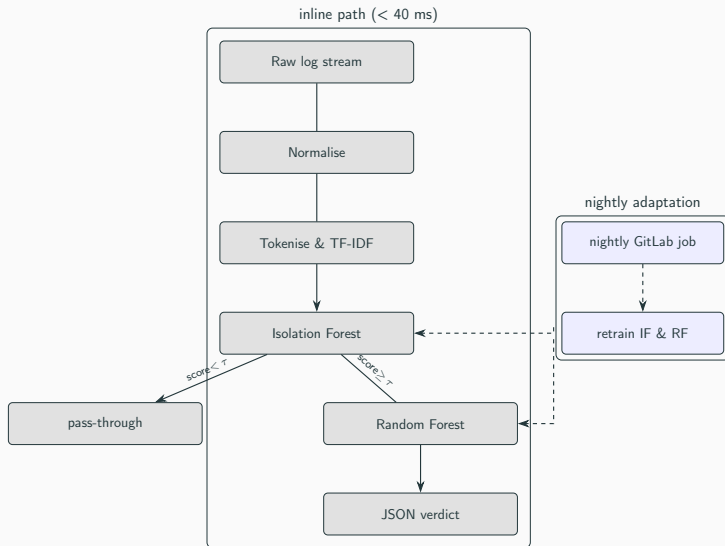- Majority vote = deterministic, auditable output.

# Random Forest ③ - error labelling

- Maps each flagged line to a domain-specific error category.
- Majority vote = deterministic, auditable output.
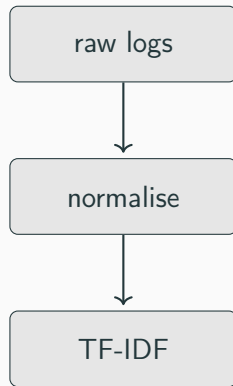- Nightly retrain < 90 s; warm-start handles drift.

# Architecture

# Data & evaluation

## Datasets & metrics

- 117 k *macOS* logs + 655 k *OpenSSH* logs

```
┌──────────────┐
│   raw logs   │
└──────────────┘
        │
        ▼
┌──────────────┐
│  normalise   │
└──────────────┘
        │
        ▼
┌──────────────┐
│    TF-IDF    │
└──────────────┘
```

## Datasets & metrics

- 117 k *macOS* logs + 655 k *OpenSSH* logs
- 504 labelled anomalies (class imbalance $\approx$ 1 : 200)

```
raw logs
   │
   ▼
normalise
   │
   ▼
TF-IDF
```

## Datasets & metrics

- 117 k *macOS* logs + 655 k *OpenSSH* logs
- 504 labelled anomalies (class imbalance $\approx$ 1 : 200)
- Split 70 / 15 / 15 % (training / validation / testing)
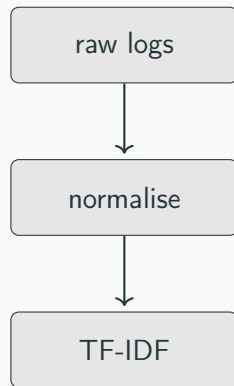
```
raw logs
   |
   v
normalise
   |
   v
TF-IDF
```
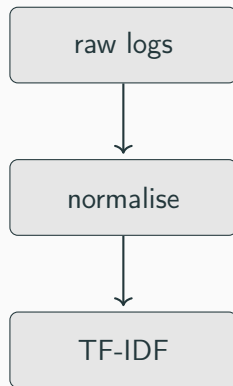
## Datasets & metrics

- 117 k *macOS* logs + 655 k *OpenSSH* logs
- 504 labelled anomalies (class imbalance $\approx$ 1 : 200)
- Split 70 / 15 / 15 % (training / validation / testing)
- Metrics: Precision, Recall and $F_1$

# Results

|                                | Precision | Recall | $F_1$ |
|--------------------------------|-----------|--------|-------|
| Detection (Isolation Forest)   | 0.91      | 0.88   | 0.89  |
| Classification (Random Forest) | 0.99      | 0.99   | 0.99  |
| Regex Baseline                 | 0.286     | 0.286  | 0.286 |

$$F_1 = 2 \cdot \frac{P \cdot R}{P + R}$$

# Impact

## Operational impact

- **Latency**: minutes → **milliseconds** (inline verdict).

## Operational impact

- **Latency**: minutes $\rightarrow$ **milliseconds** (inline verdict).
- **Cost-free**: 2.3 k lines of code, CPU-only, no token fees.

## Operational impact

- **Latency**: minutes $\rightarrow$ **milliseconds** (inline verdict).
- **Cost-free**: 2.3 k lines of code, CPU-only, no token fees.
- **GDPR compliant**: logs never leave the VPN.

# Wrap-up

## Light-weight on-prem ML matches AIOps SaaS

without latency, cost or privacy pain.

Questions welcome - thank you!