

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

Bachelor Thesis – Defence

---

Maid Ališić

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

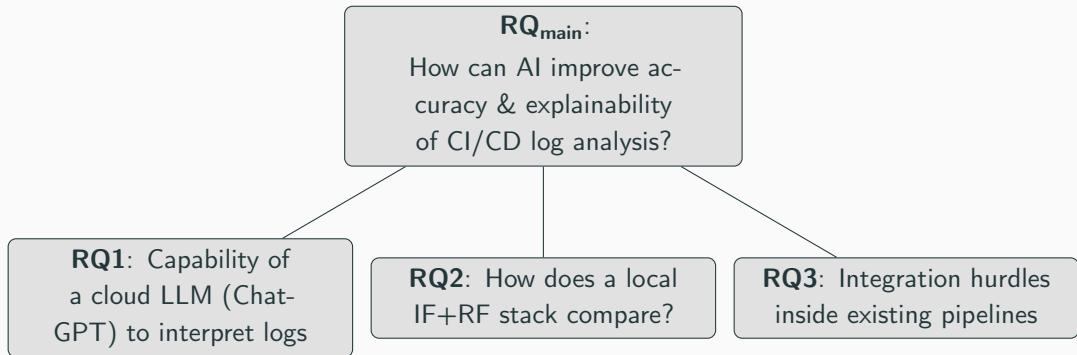
Results

Impact

## Research questions

---

# Research questions



## Problem context

---

## Why do we care?

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

## Why do we care?

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

## Why do we care?

- CI/CD emits  $\approx 10\text{--}20$  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.



## Why do we care?

- CI/CD emits  $\approx 10\text{--}20$  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.
- Logs may expose customer IDs, therefore they must remain on-premises (no cloud export).

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

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

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 the job completes.

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 the job completes.
4. **Alert fatigue** – regex rule sets grow without bound.

## Method

---

## Vectorisation pipeline ①

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



## Vectorisation pipeline ①

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





## Vectorisation pipeline ①

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



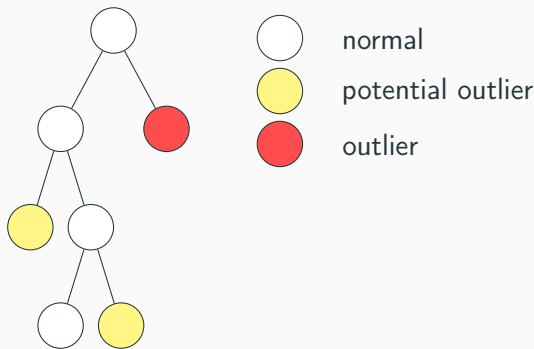
## Vectorisation pipeline ①

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



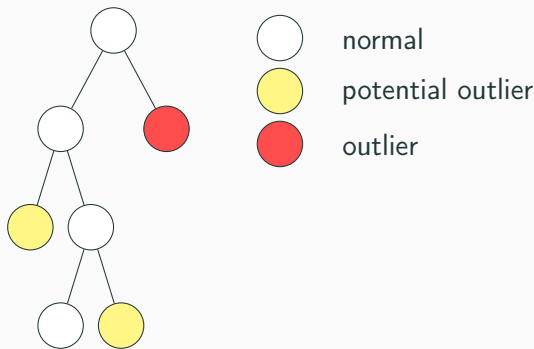
## Isolation Forest ② – intuition

- Random binary partitioning isolates unusual lines in fewer splits.



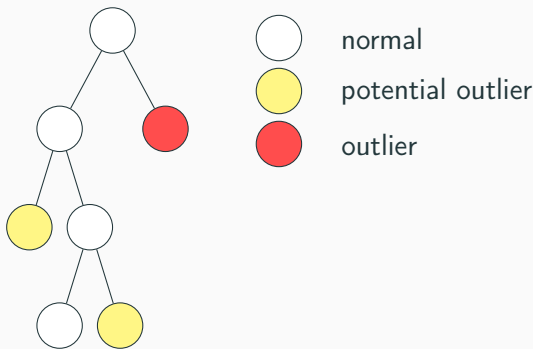
## Isolation Forest ② – intuition

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



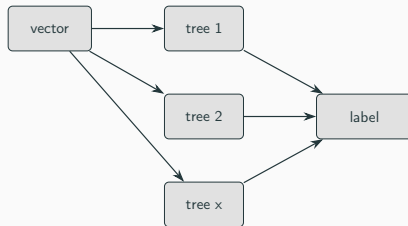
## Isolation Forest ② – intuition

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



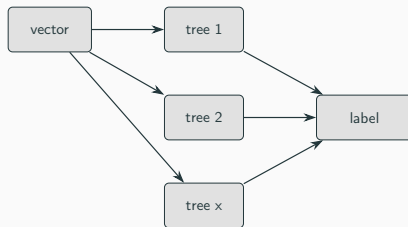
## Random Forest ③ – error labelling

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



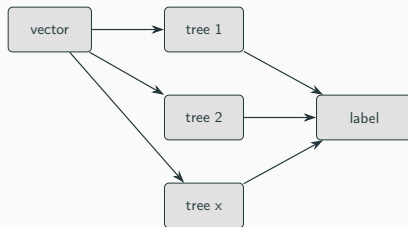
## Random Forest ③ – error labelling

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



## Random Forest ③ – error labelling

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

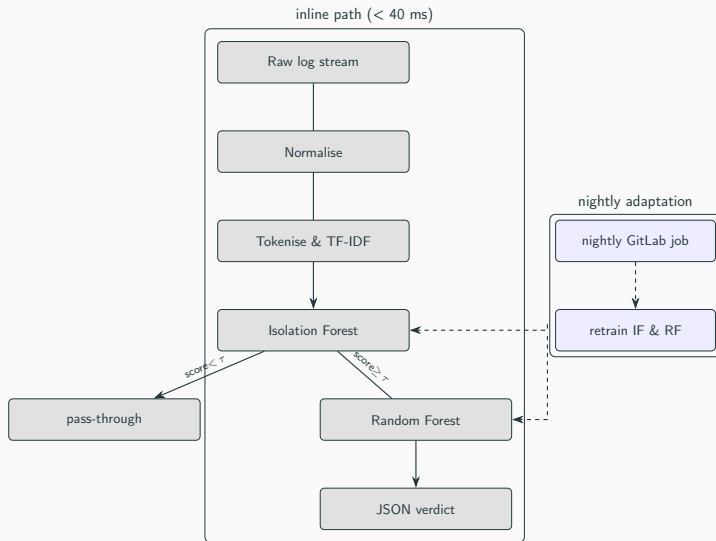




# Architecture

---

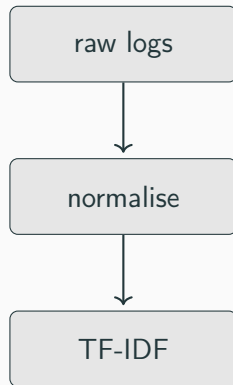
# End-to-end pipeline (< 40 ms inline)



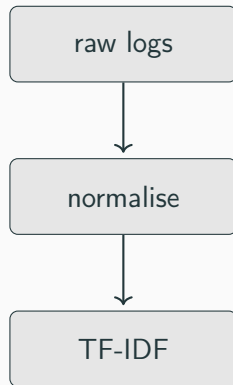
## Data & evaluation

---

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

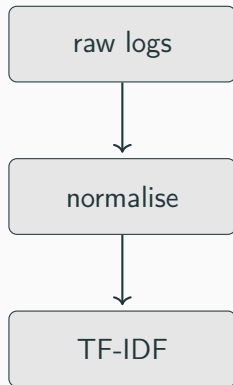


- 117 k *macOS* logs + 655 k *OpenSSH* logs
- 504 labelled anomalies (class imbalance  $\approx 1 : 200$ )



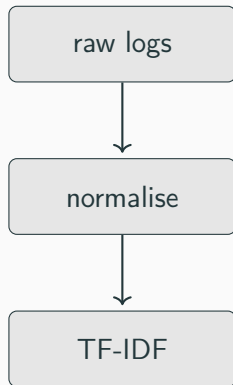
# Datasets & metrics

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



# 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

---



## Headline numbers

	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

---

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

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

- **Latency:** minutes → **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!