# Release Engineering for ML Applications

Team 14: Dani, Justin, Nadine, Yang Li
28 June 2024

TU Delft
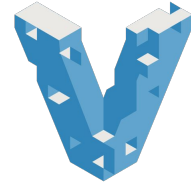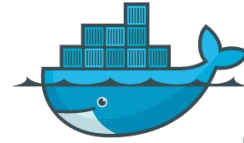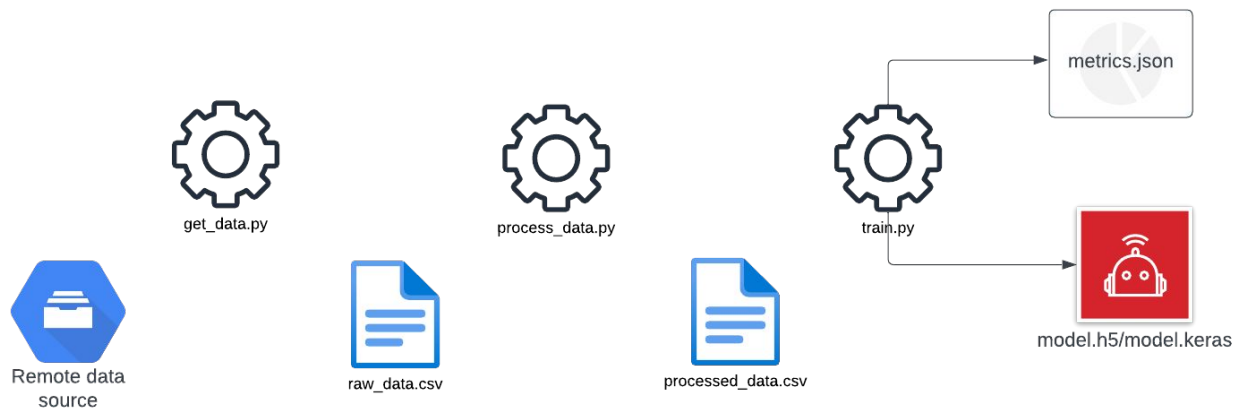
# Introduction

1. ML Training Pipeline
2. ML Testing Design
3. Release Pipeline
4. K8s Deployment
5. Experimental Setup: Istio & Prometheus
6. Additional Istio Use Case
7. Extension Proposal

**TU**Delft

# 1. ML Pipeline

- DVC
- AWS S3
- Poetry package/dependency manager



metrics.json

get_data.py

process_data.py

train.py

Remote data source

raw_data.csv

processed_data.csv

model.h5/model.keras

TUDelft

# 1. ML Pipeline: artifacts

- Trained model
- Raw and processed data
- metrics

| Metric | Value (%) |
|---|---|
| Accuracy | 0.00022 |
| Val Accuracy | 0.000002 |
| Precision | 44.6082 |
| Val Precision | 44.6080 |
| Recall | 99.9593 |
| Val Recall | 100.0000 |
| Loss | 0.35034 |
| Val Loss | 0.35619 |

# 1. ML Pipeline: linters

- Pylint & Flake8
- Allow common naming conventions
- Discourage 'bad' naming
- Handle more exceptions
- Display only warnings with high confidence levels
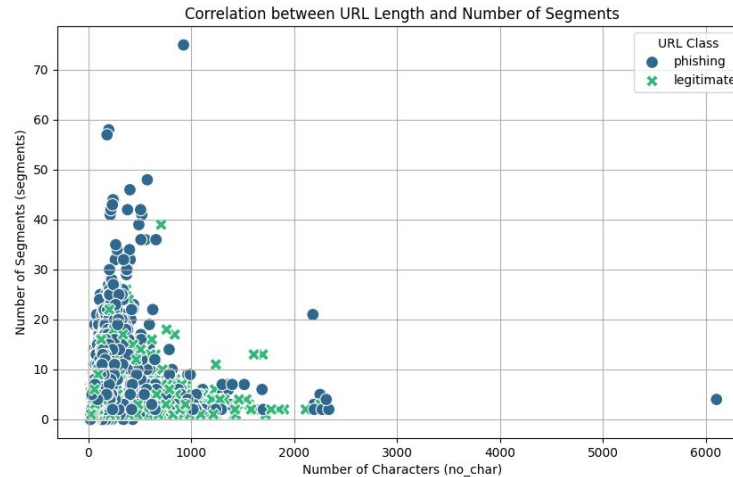
**T**U Delft

# 2. ML Testing Design: features and data

- View distribution of the data

# 2. ML Testing Design: Features and Data

- View distribution of the data
- Engineered features

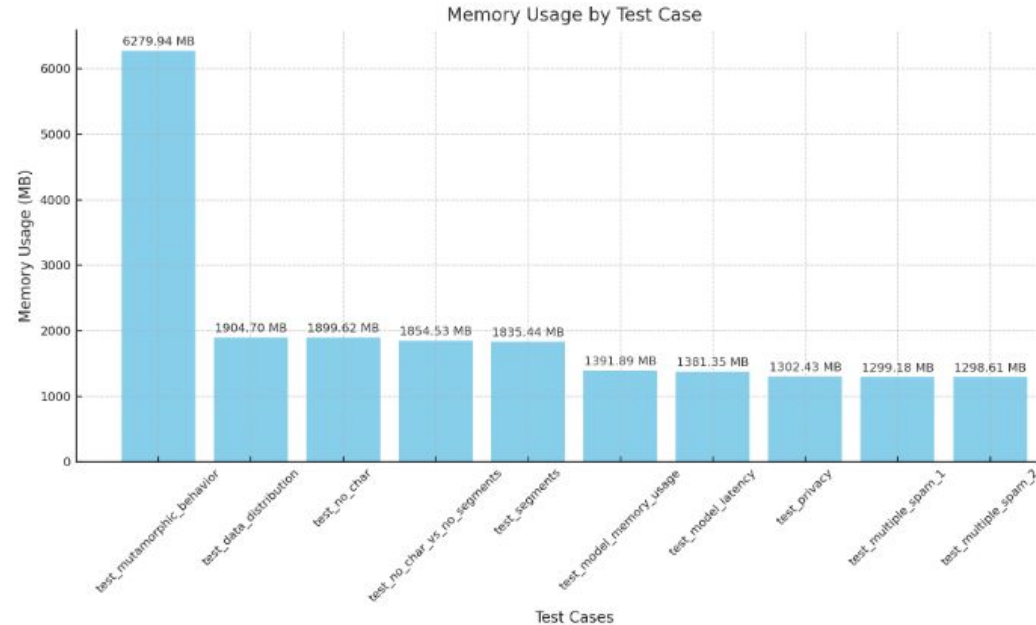Correlation between URL Length and Number of Segments

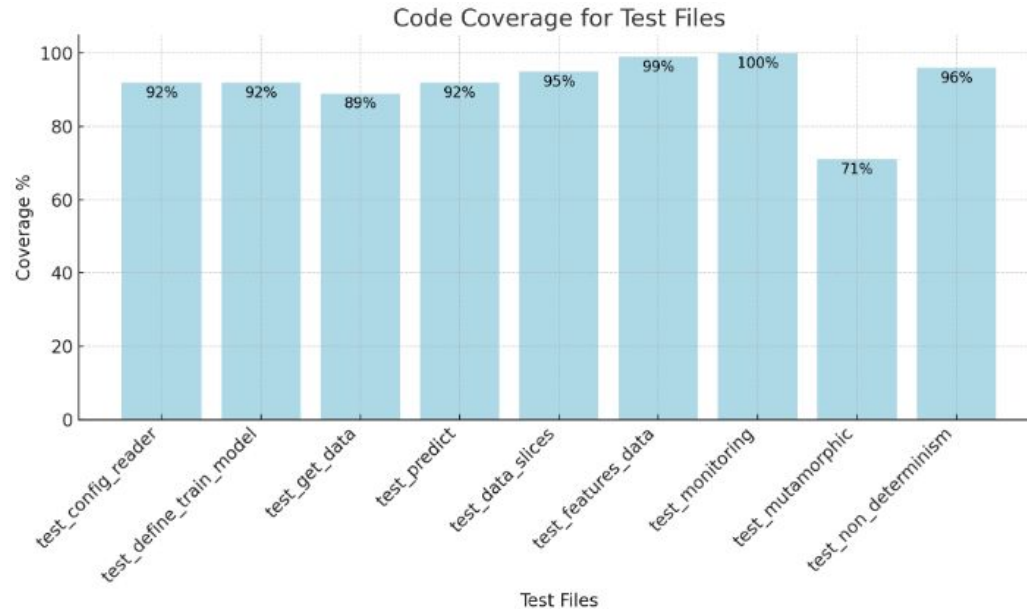# 2. ML Testing Design: Mutamorphic tests

- Check for non-deterministic behaviour

# 2. ML Testing Design: Monitoring Tests
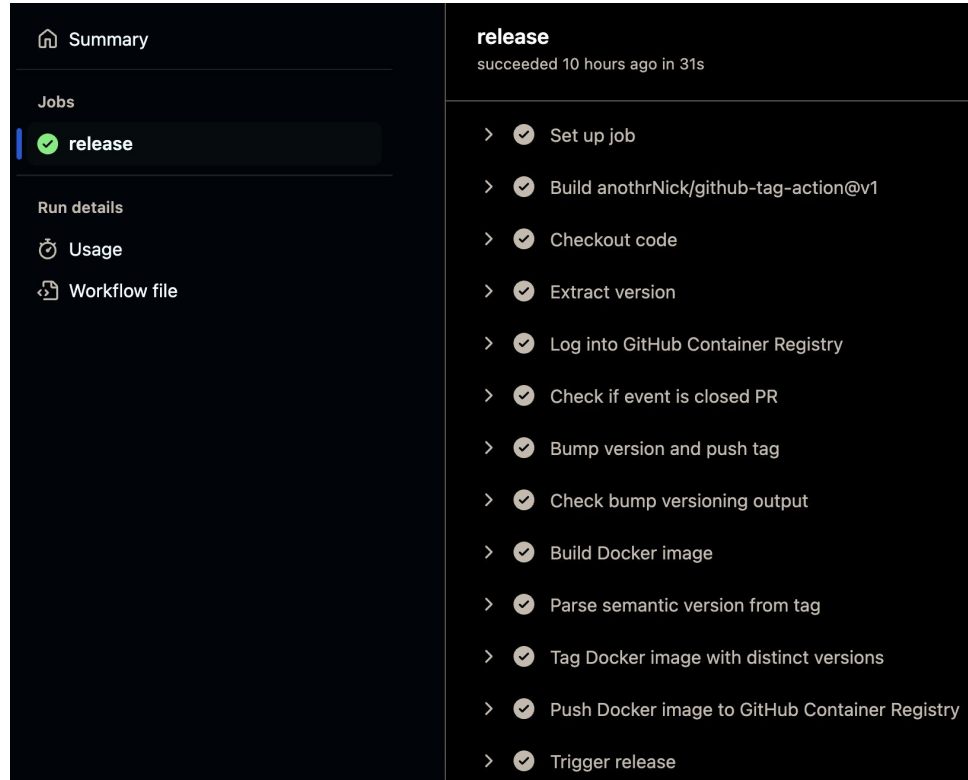
- Ram usage per method

# 2. ML Testing Design: Test Adequacy

- Statements, covered, missed



Code Coverage for Test Files

# 3. Release Pipeline: Container Image (app)

# 3. Release Pipeline: Container Image (app)

app / .github / workflows / **release.yml**

nadinekuo  Update CI

Code    Blame    108 lines (94 loc) · 4.45 KB

```yaml
1    name: Release to GitHub Container Registry
2
3    on:
4      pull_request:
5        branches:
6          - main
7        types:
8          - closed
9      push:
10       tags:
11         - 'v[0-9]+.[0-9]+.[0-9]+'         # Semantic versioning (e.g., v1.0.0, v2.3.4)
12
13   jobs:
14     version-and-release:
15       runs-on: ubuntu-latest
16       permissions:
17         contents: write          # Needed for automatic patch bumping
18
```

# 3. Release Pipeline: Container Image (app)

```
18
19      steps:
20        - name: Checkout code
21          if: github.event.pull_request.merged == true
22          uses: actions/checkout@v3
23          with:
24            ref: ${{ github.event.pull_request.merge_commit_sha }}   # Needed for automatic pre-releases
25            fetch-depth: '0'
26
27        - name: Checkout code
28          if: github.event.pull_request.merged != true
29          uses: actions/checkout@v3
30
31        - name: Extract version
32          id: version              # Set ID to refer back to in later steps
33          run: echo "::set-output name=version::${GITHUB_REF#refs/tags/}"
```

**T**U**Delft**

# 3. Release Pipeline: Container Image (app)

```
35          - name: Log into GitHub Container Registry
36            run: echo "${{ secrets.TOKEN }}" | docker login ghcr.io -u ${{ github.actor }} --password-stdin
37
38          # In the case of closing PRs, automatically bump current version (default: patch)
39          # Any commit message that includes #major, #minor, #patch, or #none will trigger the respective version bump.
40          # If no #major, #minor or #patch tag is contained in the merge commit message, it will bump whichever DEFAULT_BUMP is set to
41          - name: Bump version and push tag
42            if: github.event.pull_request.merged == true
43            id: bump
44            uses: anothrNick/github-tag-action@v1
45            env:
46              GITHUB_TOKEN: ${{ secrets.GITHUB_TOKEN }}
47              WITH_V: true                    # Tag with v character
48              DEFAULT_BUMP: patch             # Which type of bump to use when none explicitly provided in commit msg
49              PRERELEASE: true                # Define if workflow runs in prerelease mode
50              PRERELEASE_SUFFIX: beta         # Suffix for your prerelease versions. Note this will only be used if a prerelease branch.
```

TUDelft

# 3. Release Pipeline: Container Image (app)

```
60            - name: Build Docker image
61              run: docker build -t app .
62
63            - name: Parse semantic version from tag
64              id: semantic-version
65              run: |
66                echo "GITHUB_REF: ${GITHUB_REF}"
67                if [[ "${{ github.event.pull_request.merged }}" == 'true' ]]; then
68                  echo "Case PR closed: the automatically bumped version will be used"
69                  VERSION=${{ steps.bump.outputs.new_tag }}
70                else
71                  echo "Case pushed tag: the manually created tag will be used"
72                  VERSION=${GITHUB_REF:11}
73                fi
74                echo "$VERSION"
75                MAJOR=`echo "$VERSION" | cut -d . -f 1`
76                MINOR=`echo "$VERSION" | cut -d . -f 2`
77                PATCH=`echo "$VERSION" | cut -d . -f 3`
78                echo "::set-output name=version::$VERSION"
79                echo "::set-output name=version_major::$MAJOR"
80                echo "::set-output name=version_minor::$MINOR"
81                echo "::set-output name=version_patch::$PATCH"
82
```
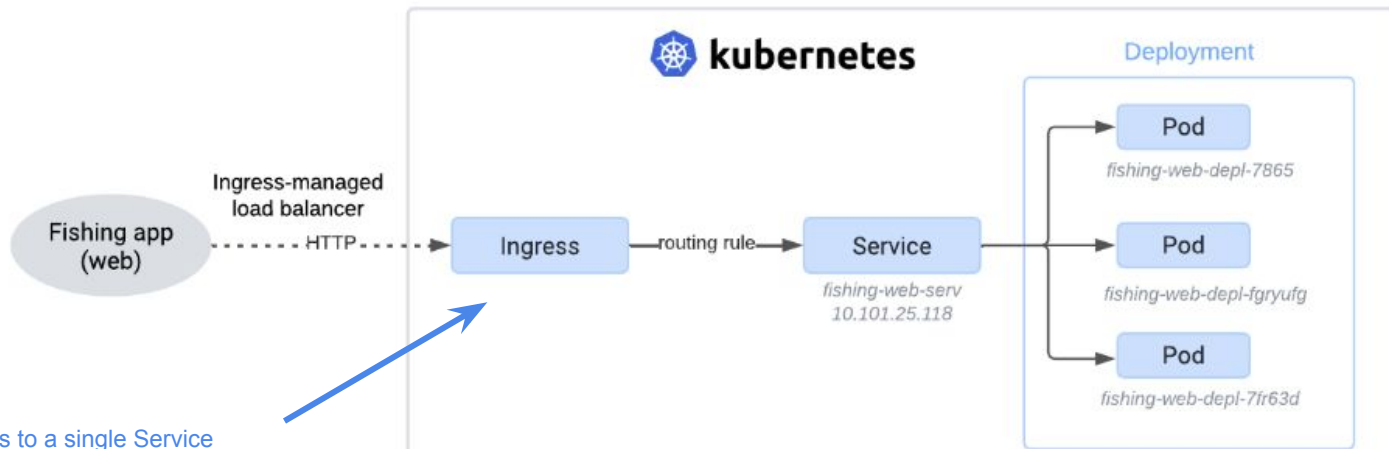
# 3. Release Pipeline: Container Image (app)

```
83          - name: Tag Docker image with distinct versions
84            run: |
85              IMG=ghcr.io/${{ github.repository }}
86              docker tag app $IMG:${{ steps.semantic-version.outputs.version }}
87              docker tag app $IMG:latest
88              docker tag app $IMG:${{ steps.semantic-version.outputs.version_major }}.${{ steps.semantic-version.outputs.version_minor }}.latest
89              docker tag app $IMG:${{ steps.semantic-version.outputs.version_major }}.latest
90
91          - name: Push Docker image to GitHub Container Registry
92            run: |
93              IMG=ghcr.io/${{ github.repository }}
94              docker push $IMG:${{ steps.semantic-version.outputs.version }}
95              docker push $IMG:latest
96              docker push $IMG:${{ steps.semantic-version.outputs.version_major }}.${{ steps.semantic-version.outputs.version_minor }}.latest
97              docker push $IMG:${{ steps.semantic-version.outputs.version_major }}.latest
98
99          - name: Trigger release
100           uses: actions/create-release@v1
101           env:
102             GITHUB_TOKEN: ${{ secrets.TOKEN }}
103           with:
104             tag_name: ${{ steps.version.outputs.version }}
105             release_name: Release ${{ steps.version.outputs.version }}
106             body: |
107               Release app version ${{ steps.version.outputs.version }}
```

# 3. Release Pipeline: Software Package (lib-ml)

```
53          - name: Set up Python
54            uses: actions/setup-python@v2
55            with:
56              python-version: '3.8'
57
58          - name: Install dependencies
59            run: |
60              python -m pip install --upgrade pip
61              pip install setuptools wheel twine
62
63          # Both source distribution (.tar.gz or .zip) and built distribution (.whl) are built
64          - name: Build Distribution
65            run: |
66              python setup.py sdist bdist_wheel
67
68          # Upload all files to PyPi
69          - name: Publish lib to PyPI
70            if: startsWith(github.ref, 'refs/tags/v')
71            run: |
72              twine upload dist/*
73            env:
74              TWINE_USERNAME: __token__
75              TWINE_PASSWORD: ${{ secrets.remla_a2 }}
```

# 4. K8s Deployment



- Routes requests to a single Service
- Ingress Controller determines relevant Service port

```
79    ---
80    apiVersion: networking.k8s.io/v1   # Create an Ingress for the application (to expose the ClusterIP service)
81    kind: Ingress
82    metadata:
83      name: fishing-gateway
84    spec:
85      defaultBackend:
86        service:
87          name: fishing-web-serv
88          port:
89            number: 8000
```

# 4. K8s Deployment

```
25      ---
26      apiVersion: v1     # Create service for image operation-model-service
27      kind: Service
28      metadata:
29        name: fishing-model-serv
30      spec:
31        selector:
32          app: fishing-model      # App defined in Deployment above
33        ports:
34          - port: 5001
35            targetPort: 5001
```

# 4. K8s Deployment

```yaml
36  ---
37  apiVersion: apps/v1      # Create deployment for image operation-app
38  kind: Deployment
39  metadata:
40    name: fishing-web-depl
41    labels:
42      app: fishing-web
43  spec:
44    replicas: 1
45    selector:
46      matchLabels:
47        app: fishing-web
48    template:
49      metadata:
50        labels:
51          app: fishing-web
52      spec:
53        containers:
54        - name: fishing-web
55          image: ghcr.io/remla2024-team14/app:latest
56          imagePullPolicy: Always
57          ports:
58          - containerPort: 8000
59          env:
60            - name: MODEL_SERVICE_URL
61              valueFrom:
62                configMapKeyRef:
63                  name: my-config    # See ConfigMap defined at bottom
64                  key: model.host
65        imagePullSecrets:
66        - name: ghcr-login-secret
```
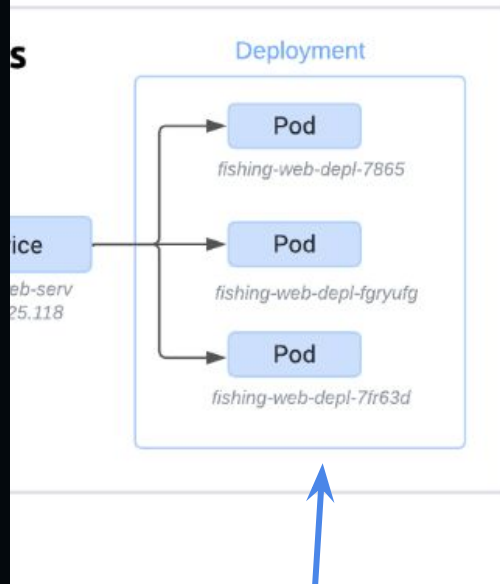
Fishing app
(web)

Deployment

Pod

*fishing-web-depl-7865*

Pod

*fishing-web-depl-fgryufg*

Pod

*fishing-web-depl-7fr63d*

- Specification of desired state of Pods
- K8s ensures this state is retained + allows for scaling up or down replicas of Pods

# 4. K8s Deployment

- Prevent internal Pods from being accessed publicly
- Listen on specified port (fixed)
- Forward requests to containers in Pods

Deployment

kubernetes

Pod
*fishing-web-depl-7865*

Ingress-managed
load balancer

Fishing app
(web) ····· HTTP ····· ▶ Ingress ──routing rule──▶ Service

Pod
*fishing-web-depl-fgryufg*

*fishing-web-serv*
*10.101.25.118*

Pod
*fishing-web-depl-7fr63d*

- Routes requests to a single Service
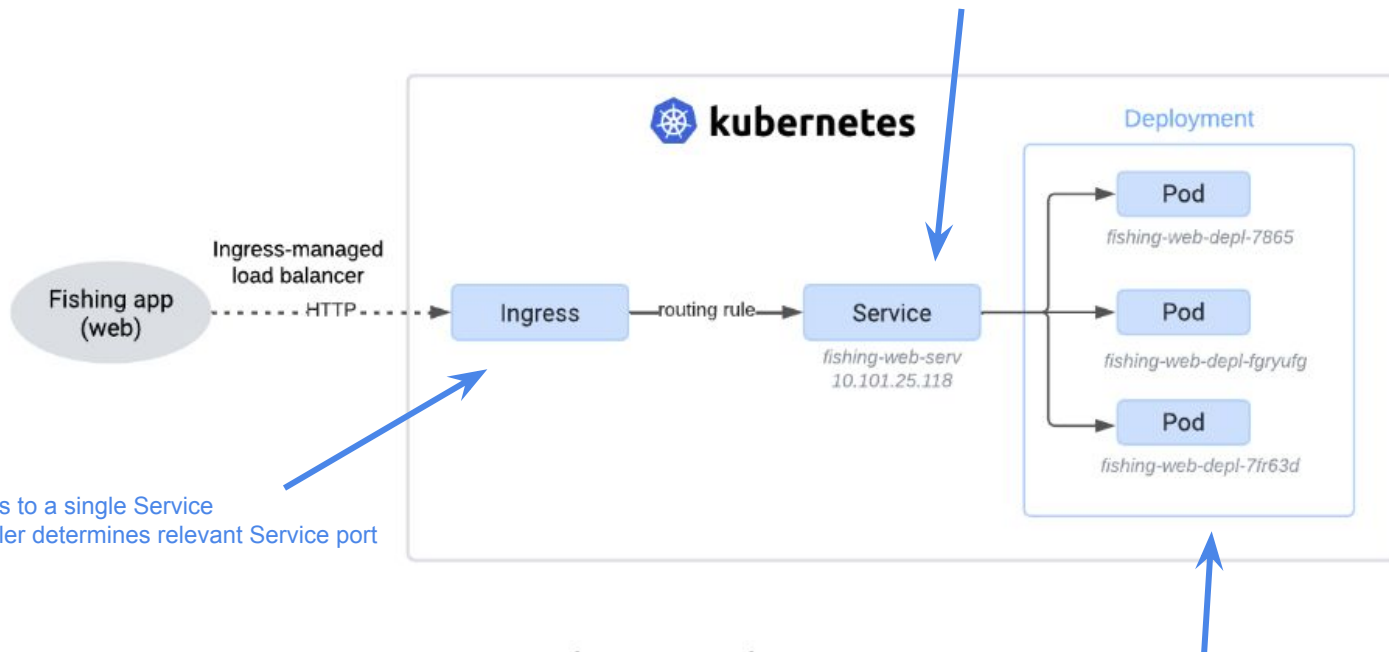- Ingress Controller determines relevant Service port

**Figure 2: Kubernetes Deployment Overview**

- Specification of desired state of Pods
- K8s ensures this state is retained + allows for scaling up or down replicas of Pods

**T**U Delft

# 5. Experimental Setup

Requirements:

- Istio Service Mesh that deploys two versions of app-frontend and app-service.
    - Prometheus annotations to scrape.
- Verifiable hypothesis with a reject/validate criterion.
- Prometheus scrapes the relevant metrics.
    - Summary, counter, gauge, histogram.
    - Metrics about inference time, user interaction metrics, etc…
- Grafana displays results
- AlertManager sends an email for extreme cases (example: too many calls)

Hypothesis: *"Given a version of the web-app that returns binary predictions (valid/phishing) to the user and a version that returns probabilities instead, the user will prefer the version of the web-app that returns binary predictions"*

TUDelft

# URL Phishing Detection

Library Version: v5.0.1

App version: v1

Enter Text:

google.com

Choose a model:

model.h5

Submit

This link is valid

Is this prediction accurate?

Yes

No

# URL Phishing Detection

Library Version: v5.0.1

App version: v2

Enter Text:

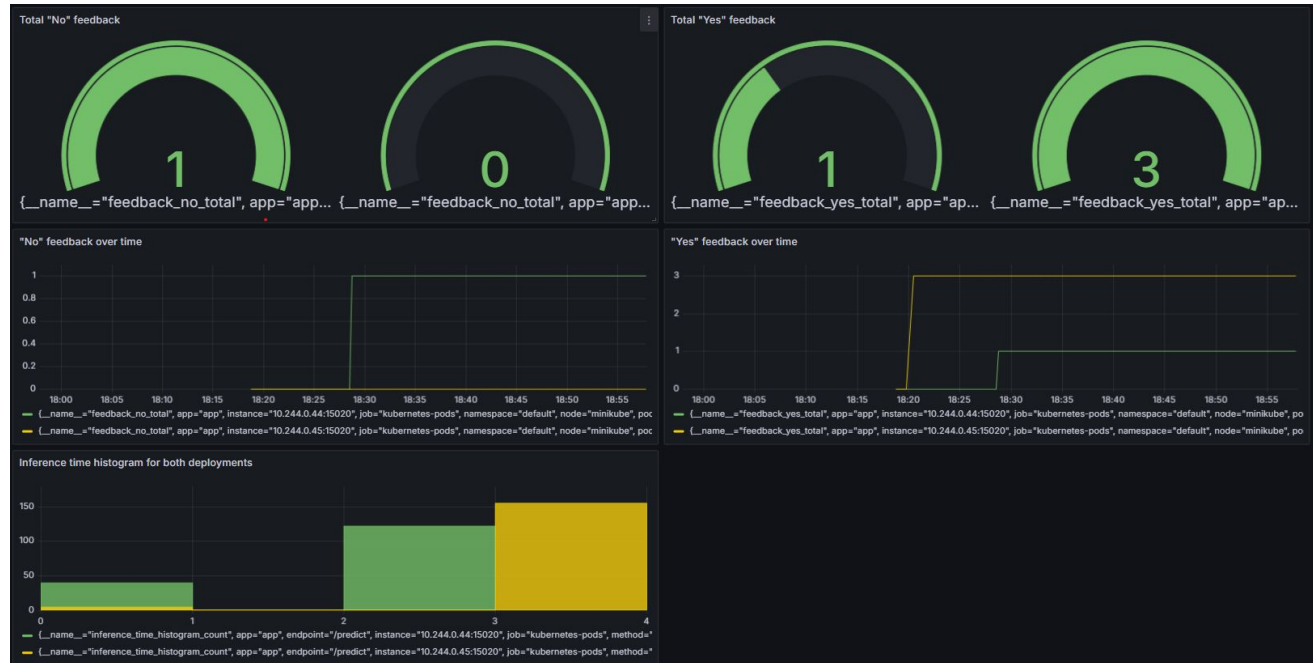google.com

Choose a model:

model.h5

Submit

0.5110589

Is this prediction accurate?

Yes

No

**T̃U**Delft

# 5. Experimental Setup: Prometheus and Grafana

# 5. Experimental Setup: Istio

# 6. Additional Istio Use Case-Rate Limiting

**Egress Extension**

In Istio, the Egress extension is used to manage and control traffic to external services in the service grid.

**Adjustments to the Base Deployment**

- Increase ENVOYFILTER configuration

- Enable and verify the flow limit function

- Monitoring and log analysis

- Adjustment and optimization strategy



```
PS C:\Users\maart\Desktop\REMLA-A5\ratelimit> kubectl exec -it sleep-0 -n rl -- sh
~ $ for i in $(seq 1 15); do curl -s -o /dev/null -w "%{http_code}\n" http://httpbin.rl:80/get; done
200
200
200
200
200
200
200
200
200
200
429
429
429
429
429
```

# 7. Extension Proposal- Rollback Strategy

**Benefits:**

- Improve System Stability

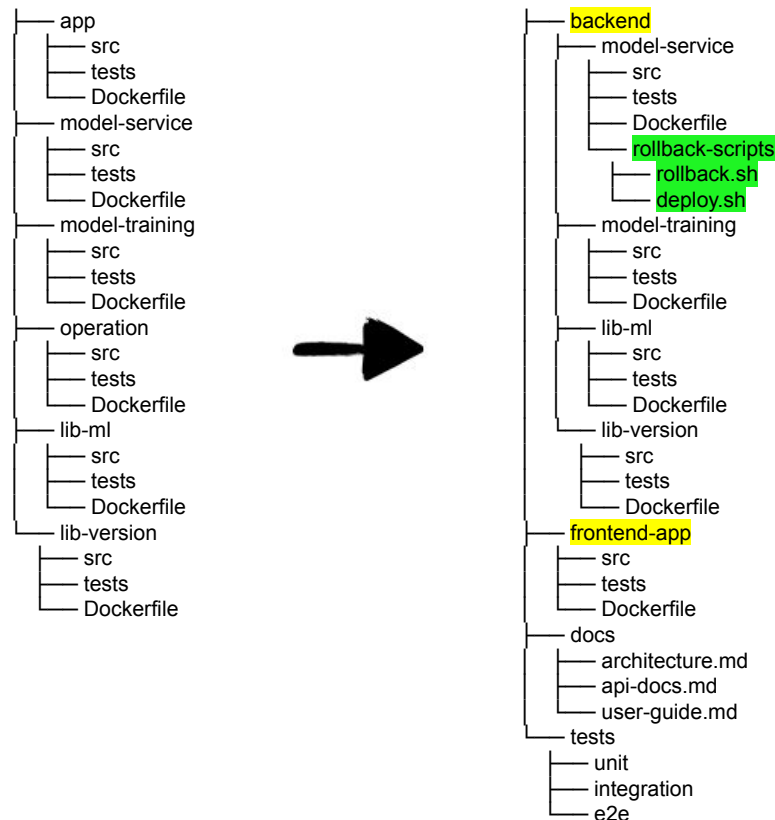- Reduce deployment risk

**Suggestions:**

- Configure Automated Rollback with CI/CD Tools

- Blue-Green Deployment

```
├── app
│   ├── src
│   ├── tests
│   └── Dockerfile
├── model-service
│   ├── src
│   ├── tests
│   └── Dockerfile
├── model-training
│   ├── src
│   ├── tests
│   └── Dockerfile
├── operation
│   ├── src
│   ├── tests
│   └── Dockerfile
├── lib-ml
│   ├── src
│   ├── tests
│   └── Dockerfile
└── lib-version
    ├── src
    ├── tests
    └── Dockerfile
```

→

```
├── backend
│   ├── model-service
│   │   ├── src
│   │   ├── tests
│   │   ├── Dockerfile
│   │   └── rollback-scripts
│   │       ├── rollback.sh
│   │       └── deploy.sh
│   ├── model-training
│   │   ├── src
│   │   ├── tests
│   │   └── Dockerfile
│   ├── lib-ml
│   │   ├── src
│   │   ├── tests
│   │   └── Dockerfile
│   └── lib-version
│       ├── src
│       ├── tests
│       └── Dockerfile
├── frontend-app
│   ├── src
│   ├── tests
│   └── Dockerfile
├── docs
│   ├── architecture.md
│   ├── api-docs.md
│   └── user-guide.md
└── tests
    ├── unit
    ├── integration
    └── e2e
```

**T̃U**Delft

# Q&A

**Thanks for listening! Questions?**