

REMLA Project Report – Group 15

Marianna Louizidou, Mitali Patil, Shantnav Agarwal

1 INTRODUCTION

In this report, we explain the design and engineering of the URL phishing detection project.

2 RELEASE PIPELINE

In this section, we present a detailed overview of the release workflow for both the software packages (the 'lib-ml' and 'lib-version' Python libraries) and the container images for the Flask app and model service.

2.1 Releasing a Software Package

The 'lib-ml' and 'lib-version' repositories contain the code for two distinct Python packages. The 'lib-ml' package is designed for pre-processing both training data and incoming queries for the REST-API call, while the 'lib-version' package implements a version-aware library capable of returning the library version. Both packages are uploaded to the PyPI package index via a GitHub Actions workflow, triggered when a new tag starting with "v" is pushed, followed by the (semantic) version of the package to be published. The steps in the "Publish package to PyPI" workflow are as follows:

- **Checkout Code** : This step uses the 'actions/checkout@v2' to checkout the latest version of the repository.
- **Set Up Python** : This step uses 'actions/setup-python@v2' to ensure python version 3.10 is set up in the working environment.
- **Extract Version from Tag** : Extracts the version number from the tag name and sets it as an output variable VERSION. This is used for versioning the package.
- **Update version.py** : This step updates the version.py file with the extracted version number from the previous step if the current reference is a tag. This file contains the version of the package. The cat command outputs the contents of 'version.py' to the log for verification.
- **Install Poetry** : Since we used Poetry[8] for our dependency and package management across the project, we installed poetry for the required dependencies.
- **Update pyproject.toml Version** : Updates the version in .pyproject.toml using the extracted version number to be used by Poetry to install the dependencies in the .toml file.
- **Lint with Flake8** : Runs flake8 to perform linting on the code in the 'lib_version_URLPhishing' directory. Linting helps us identify and fix potential issues and code style violations if the pipeline breaks.
- **Commit and pull latest changes** : The updated versions are committed and the latest changes are pulled. This step is crucial to make sure the build uses the most recent updates.
- **Build and Publish to PyPI** : This step builds the package using poetry build, and publishes the built package to PyPI using the provided credentials (stored as secrets in the repository).

- **Install and generate documentation**: Pdoc is installed for generating documentation for this published package. It is also published on GitHub pages for viewing online.

The pypi published for libml[9] and lib-version [5]

2.2 Releasing Container Images

The app repository contains a Flask app which is the front-end service that exposes our phishing URL model to users. The 'app' relies on 'model-service' to preprocess and predict the queries using REST API call and 'lib-version' to make the version visible on the front-end. We release two container images for the app and the model-service respectively. For the model prediction, we do not include the model in the model-service's image, instead, it is fetched from the model-training repository into the docker container. Using a GitHub Actions workflow, the container images for these repositories are built and uploaded to the GitHub Container Registry, allowing us to easily start instances of the app and model-service using Docker or Kubernetes. Just like the workflows for the packages, the workflow is triggered when a new tag starting with v is pushed, which should be followed by the (semantic) version of the package that is to be published. The different steps to Build and Publish Docker Image workflow are as follows:

- **Checkout Code** : This step uses the 'actions/checkout@v2' to checkout the latest version of the repository.
- **Parse Version Info from Tag**: Extracts version information from the tag (e.g., v2.3.5), splitting it into major, minor, and patch components, and saves them as environment variables.
- **Login to GitHub Registry**: We use the GitHub token to log into the Container Registry(ghcr.io). This authentication allows us to push images to the container.
- **Build and Push Docker Image** : This step builds the Docker image and tags it with the semantic git tag and latest which is then pushed to the container registry.

The workflow also generates and deploys documentation to GitHub Pages. The respective images are available in the GitHub container registry of the Organisation under the packages section. The container registry for model-service [6] and app[4]

3 PROVISIONING AND OPERATING A WEB APP

3.1 Provisioning

A vagrant file is used to provision the virtual machines. Two types of machines have been defined - controller, worker#. The worker nodes are provisioned within a loop; therefore the number of worker nodes deployed can be indicated in the vagrant file. Two ansible playbooks are written, one for each type of machine. On the worker machines, K3S and Docker are installed. On the controller machine, helm charts and the required repos are also installed in addition to K3S and Docker.

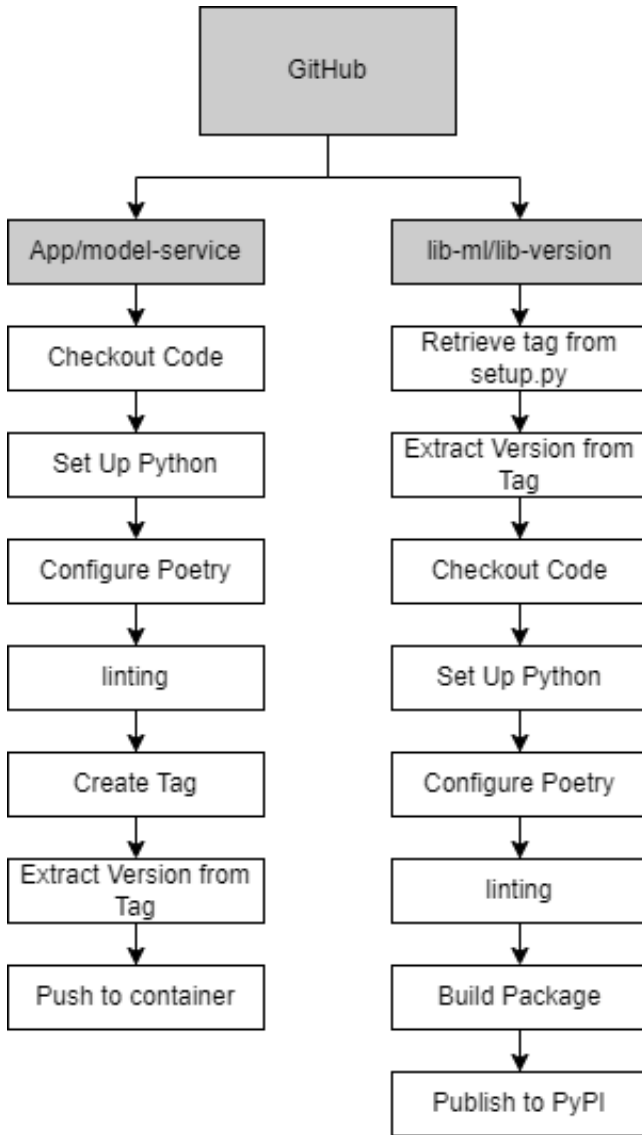


Figure 1: Release Pipeline Diagram

All machines are setup as part of a K3S cluster. The controller machine acts as the K3S server while the worker machines register as K3S agent nodes. The pods for Prometheus, Grafana and the URL App are deployed as helm charts in the controller machine/ node. The cluster automatically deploys these pods to any of the nodes.

3.2 Kubernetes

The URL App consists of two main deployments: the URL-app -> this serves the app front-end; and the URL-model -> this serves the ML model in the back end. An ingress *my-ingress* receives the request from the *load-balancer* and forwards it to the service *URL-app-serv*. This service forwards the request to the deployment. The *URL-app-depl* knows the address of the *URL-model-serv* from the *my-config* configuration. The query with the URL to be tested

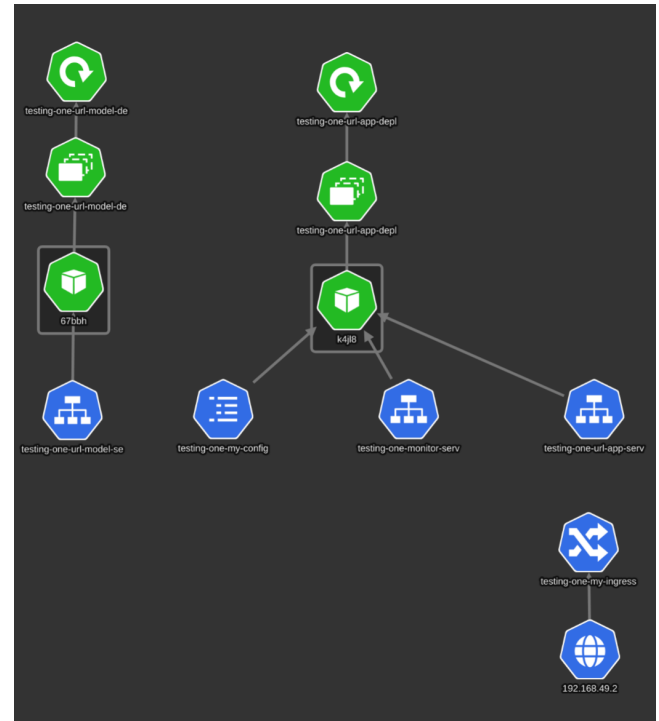


Figure 2: Nodes deployed for the URL-App. Shown here as a helm chart deployment with name *testing-one*. Visualised using kubeview

is forwarded to this address for inference from the ML model. A service monitor is defined with the same release label as the name of the prometheus-kube chart installation. This helps the URL app to be discovered by prometheus.

A helm chart of this deployment has been created. Using a value config, variables such as ML model URL, container image path, and resource limits (cpu, memory) can be defined. Multiple installations of this helm chart can be run at the same time.

3.3 Monitoring

The Prometheus & Grafana are launched using the prometheus-kube-stack helm chart. The values file has been modified to configure the ingress for Prometheus and Grafana. This value file also configures the alerts for Prometheus. The Grafana dashboard is defined as a *ConfigMap* and applied using kubectl.

The URL app publishes the metrics at the endpoint */metrics* about various things such as request duration, count and exceptions. These metrics are reported separately for each endpoint. This implementation has been done using the prometheus *flask* python library.

3.3.1 Prometheus. The prometheus service scrapes the metrics endpoint of the URL app and collects information about the number of requests to each endpoint, time taken to complete these requests etc. Based on information from the kubernetes cluster and the URL App, 3 alerts have been configured.

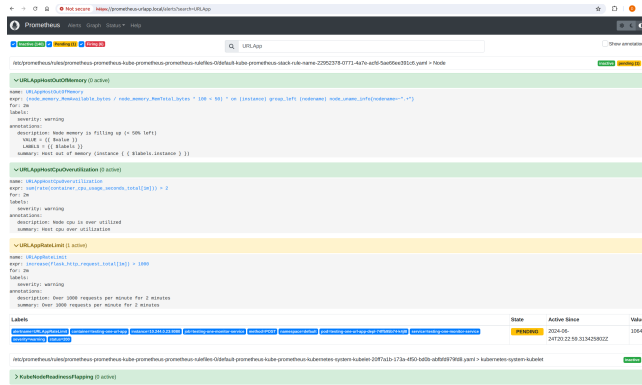


Figure 3: Alerts show up based on CPU usage, memory consumption and request rate (triggered).

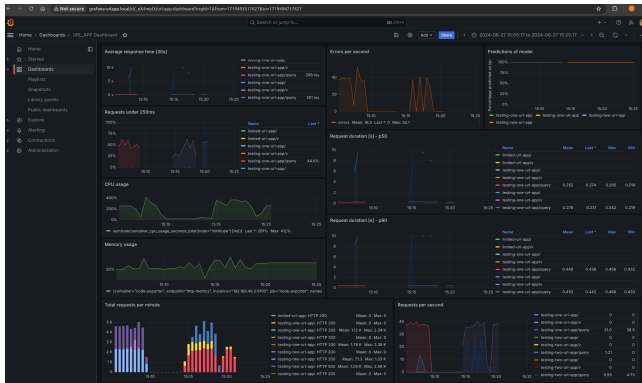


Figure 4: Visualisations shown using the grafana dashboard. This dashboard show the high rate of requests which triggered the alert in Figure 3.

3.3.2 *Grafana*. Grafana is used to continuously monitor the services using multiple visualisations. Where required, metrics such as rate aggregations, averages are also shown. Please refer to Figure 4.

4 ML PIPELINE

ML tasks involve various stages like data gathering, data preparation, training, and so on. We automate this workflow to be able to manage the code, reproduce, deploy at scale, maintain versions, and follow the best practices through pipeline orchestration. Each stage has its own code to process the data, and running the entire pipeline after every change is inefficient and resource-intensive. Hence, we used Data Version Control (DVC) [3] to configure an ML pipeline 5 for this project.

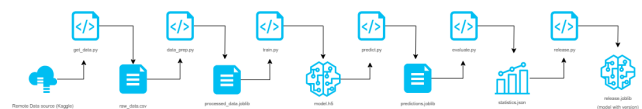


Figure 5: ML Config Pipeline.

Our DVC pipeline has 6 stages each automating specific tasks of the workflow:

- **Get Data** : This stage involves fetching and storing the raw text data from a remote repository. In our case we are fetching the train, test and evaluation text data from Kaggle [10].
- **Prepare Data** : This stage pre-processes the raw data. This involves cleaning the data, tokenizing, and padding to generate training, testing, and validation datasets, along with a character index file used for feature encoding. These are stored as .joblib files.
- **Train** : In this stage, the machine learning model is trained using the processed training data. The training script utilizes the specified parameters (like loss function, optimizer, batch size, epochs) from the params.yaml file to configure the training process. The resulting trained model is saved as an .h5 file.
- **Predict** : This stage uses the trained model to make predictions on the test dataset. The predictions are then saved as a .joblib file for further analysis.
- **Evaluate** : The evaluation stage assesses the performance of the model using the predictions made in the previous stage. Metrics such as accuracy, precision, recall, and F1-score are calculated and saved as a JSON file.
- **Release** : The release stage involves packaging the final model and any other necessary artifacts for deployment. This includes the trained model, character index, and the model version which helps in versioning the model releases, ensuring that the deployment process can track and manage different versions of the model.

DVC makes collaboration for ML tasks better and allows to work with remote files which is a significant advantage over utilities like Makefile. Moreover, it makes reproducibility, tracking, and saving of changes easier to handle allowing for better comparisons among various versioned experiments. A close alternative to DVC was considered, the Git LFS but it has drawbacks as to how much data can be uploaded per user per day. Considering our project data size and experiments, it was ruled out for data versioning.

We also used Poetry [8] to be able to manage our dependencies. Poetry has a significant advantage over using a plain 'requirements.txt' file to install packages. We do not have to maintain the dependency file every time we add a new package, unlike requirements.txt. Through poetry and a conda environment with the selected 3.10 Python version, we ensure that the dependencies are consistent across various systems and the builds are repeatable.

To ensure and maintain high quality code throughout our project, we have integrated the use of two different linters, pylint and flake8 into our CI/CD pipeline. Pylint provides a detailed analysis of Python source code, checking for errors, enforcing a coding standard, and looking for code smells. The configurations have been customized to fit our project needs, an example of this is the addition of the dslinter plugin. On the other hand, flake8 focuses more on PEP 8 compliance and common programming errors. It complements Pylint by providing an additional layer of checks.

The code quality checks are implemented in our tag creation workflow. This ensures that if the code does not pass the linter

checks, the workflow will break and no new tag will be created. Consequently, the release process is not triggered and no model is released unless the code quality meets our standards.

Automating the workflow reduced manual intervention and potential errors. Once the pipeline was set up, the build failures and linter issues could be easily handled to ensure code quality. Tags provide a clear reference to specific versions of the working code ensuring that experiments and results can be precisely reproduced. This practice enhanced the organization of our project by simplifying debugging and aiding the progress and performance of different iterations of changes every time.

5 ML TESTING

Our testing framework incorporates a variety of tests to ensure comprehensive coverage of the model's functionality and performance.



Figure 6: Testing Pipeline

Data Slices Test: The purpose of the data slices test is to evaluate the performance of the model on a representative subset of the data. This test was chosen to ensure that the model maintains consistent performance and generalizes well across different inputs that reflect the entire dataset. By creating a representative slice and comparing its performance to the entire dataset, we can confirm that the model does not favor specific segments and performs reliably on all data. We generate a random sample that preserves the original distribution of features and labels, and test whether the model's accuracy on this representative slice closely matches its accuracy on the full test set. This approach helps identify potential discrepancies and ensures that the model's performance is robust and unbiased.

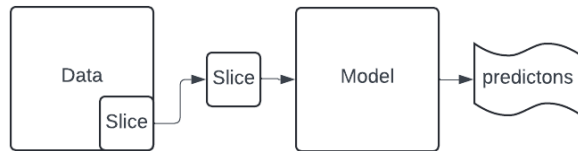


Figure 7: Data slices test

Features test: The feature test is designed to verify the integrity and correctness of the input features. This test was chosen to ensure that the data preprocessing steps were implemented correctly and that no sensitive information was included in the dataset. By confirming the correct tokenization and encoding of features, we prevent data leakage or incorrect feature extraction that could affect model performance. We test the tokenizer and encoder to ensure that they work as expected, and validate that sensitive patterns such as passwords or session tokens are appropriately masked.

Non-determinism robustness test: To ensure that model training is stable and produces consistent results, the non-determinism robustness test involves training the model multiple times with the same data and hyperparameters, and then comparing the weights and performance metrics. This test was chosen to identify any sources of randomness that could affect the reproducibility of the model, thereby ensuring consistent and reliable model training results. We trained our model multiple times with the same dataset and hyperparameters, and verified that the resulting models had nearly identical accuracy.

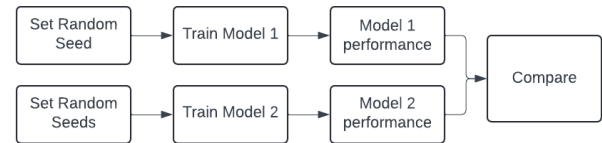


Figure 8: Non-determinism robustness test

Metamorphic Testing with Automatic Inconsistency Repair: Metamorphic testing detects inconsistencies in model predictions when inputs are slightly changed. This test was chosen to increase the robustness of the model by ensuring that it responds predictably to similar inputs. By generating mutated inputs and testing whether the model's predictions remain consistent with expectations, we can identify and repair any inconsistencies in the model's behavior. Specifically, we created new test inputs by replacing words in URLs with synonyms and observed whether the model's predictions for these altered inputs were consistent with the original predictions. When inconsistencies were detected, the system attempted to automatically repair the inconsistencies by adjusting the model's decision boundaries.

Cost of Features Test: The cost of features test evaluates the impact of different features on the model's performance. This test was chosen to determine which features are most important to the model, which helps with feature selection and dimensionality reduction. By focusing on the most important features, we can improve the efficiency of the model and potentially reduce the computational overhead. We conducted experiments to assess how including or excluding certain features, such as certain keywords or URL patterns, affected the model's performance metrics, such as accuracy and F1 score.

To ensure the reliability and completeness of our ML model, we integrated test coverage reporting into our CI pipeline using the pytest-cov plugin. Test coverage measures the proportion of our code base that is executed during testing, providing insight into which parts of the code are being tested and which are not. This metric is critical because it helps identify untested areas that are more prone to bugs. The coverage report is generated and stored in the reports directory, ensuring that we can continuously monitor and improve our coverage.

The process of implementing and running comprehensive tests, especially those involving large datasets and complex models, is time-consuming. This has highlighted the need for efficient testing

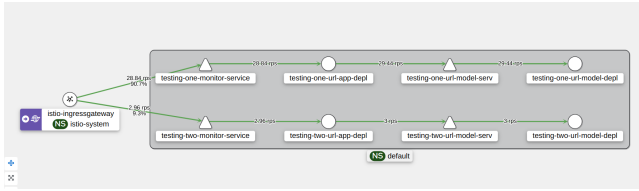


Figure 9: Istio Canary release with 10% traffic routed to testing-two

strategies and the importance of planning for adequate testing time in the project timeline.

6 ISTIO SERVICE MESH

6.1 Experiment

The original ML model, version 0.1.0, has *tanh* non-linearity. We train a new model, version 0.2.1, with *ReLU* non-linearity. Because *ReLU* is a more performant and mathematically simpler nonlinearity, According to our hypothesis, we expect the new model to have better prediction performance and inference speed.

The two ML models are deployed by specifying the appropriate URLs for the ML model via the Helm deployment for URL-App. The *testing-one* deployment has the ML models with version 0.2.1, *testing-two* has 0.1.0.

In our Grafana dashboard, we track the average response time for all requests. This metric, at the */query* endpoint is primarily affected by the inference time of the model and can be monitored to check if the inference speed improves.

We also compare the ratio of positive/ negative predictions of all models being deployed to see whether, on average, their predictions are the same. Overall, this metric highlights if a particular model is more biased towards positive or negative predictions. This metric can be further expanded to check if two different ML models are giving different predictions for the same inputs for further investigation as required.

6.2 Traffic Management

We define a Gateway and a Virtual service to route network to our URL App deployments. The Virtual service has multiple endpoints to show different possibilities:

- */testing-one*: Routes to testing-one deployment
- */testing-two*: Routes to testing-two deployment
- */testing-url-app*: Routes to testing-one with 90% weight, and testing-two with 10% weight. This endpoint shows our canary release with the new model being made available to only a fraction of users.
- */mirroring-url-app*: Routes to testing-one with mirroring to testing-two. This serves as our shadow release. The new model can be tested on real user data without being released. The metrics can be tracked on the Grafana dashboard.

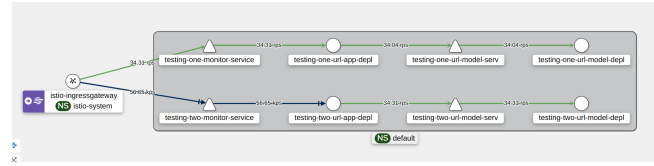


Figure 10: Istio Shadow release with all traffic mirrored to testing-two

7 EXTENSION PROPOSAL

There are a number of things we can improve in this project. Many of these things could not be completed because of the lack of time and resources.

7.1 Testing

Although most of the project is automated to speed up the deployment process, the testing setup takes a long time to run. The testing design involves training the model to verify the features of various classes involved in the experimentation. Right now the project, the testing is triggered ‘on-push’ which is expensive for every minor change made to the repository.

Proposed solution. The coverage metric could be used to limit the number of tests that need to be run each time the code is updated. Specifically, the *pytest-cov* plugin [7] could record coverage data for each test individually. This would provide the ability to map the files that are executed by each test, and then whenever there is an update in the code, identify which tests cover the updated files and run only those tests. Another solution could be to exclude expensive tests during minor updates and patches and include them in the workflow during major version bump. GitHub runner can run the workflow based on the semantic versioning and either skip or add the ‘run tests’ step in the workflow[2]

7.2 Model Training

The training pipeline is not configured to use GPUs right now. This translates to extremely long training times (and also tests). This limitation is partly because the Jupyter file was not programmed to switch between devices (CPUs, GPUs, MPS etc.) easily. Building a ML training pipeline that can take advantage of hardware specific accelerations in all environments is valuable!

Proposed solution. We need to refactor the model training pipeline to query available devices and select the appropriate one. This will introduce many difficulties with *Poetry* due to the additional dependencies on *CUDA* which may not be available on all hardware configurations. Difficulties and possible workarounds/ hacks have been discussed on this issue. This issue still remains open though.

7.3 Inference optimisation

During inference, we are running the model without any advanced optimisations. We ensure that the model is loaded from the disk only once - this is quite basic. To improve performance multiple things can be done:

- **Batching requests:** Combining multiple inference requests into a single batch.

- Quantization: Optimise precision of the model as some cost of accuracy to improve inference speed and request capacity
- Graph Optimisation: Use frameworks like TensorFlow and ONNX to simplify the computation graph and reduce the number of operations.

Proposed solution. Many of these optimisations are being researched upon and are improving at a tremendous rate! We found a well cited research paper [1] that shows how batching can be performed on stateless serverless deployments to improve performance. Quantization is used extensively to allow LLM models to be deployed on edge mobile devices regularly as also explained in this article. Finally, graph optimisation libraries also reduce inference time as well as memory consumption by reducing the number of redundant operations and better utilise available hardware resources [11].

REFERENCES

- [1] Ahsan Ali, Riccardo Pincioli, Feng Yan, and Evgenia Smirni. 2020. BATCH: Machine Learning Inference Serving on Serverless Platforms with Adaptive Batching. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*. 1–15. <https://doi.org/10.1109/SC41405.2020.00073>
- [2] <https://docs.github.com/en/actions/using-workflows/workflow-syntax-for-github-actions#pushbranchestagsbranches-ignoretags-ignore> [n.d.]. github runner for tags.
- [3] <https://dvc.org/doc> [n.d.]. DVC.
- [4] <https://github.com/REMLA24-TEAM-15/app/pkgs/container/app> [n.d.]. app.
- [5] <https://github.com/REMLA24-TEAM-15/lib-version> [n.d.]. lib-version.
- [6] <https://github.com/REMLA24-TEAM-15/model-service/pkgs/container/model-service> [n.d.]. model-service.
- [7] <https://pytest-cov.readthedocs.io/en/latest/plugins.html> [n.d.]. pytest-cov.
- [8] <https://python-poetry.org/> [n.d.]. Poetry.
- [9] <https://test.pypi.org/project/libml-URLPhishing/> [n.d.]. libml.
- [10] <https://www.kaggle.com/datasets/aravindhannamalai/dl-dataset> [n.d.]. Kaggle_data.
- [11] Rasmus Munk Larsen and Tatiana Shpeisman. 2019. TensorFlow Graph Optimizations.