# Open Code Initiative
# Evaluation Package

**WG-DAISAM**

**Elora Schörverth, Fraunhofer HHI**
Steffen Vogler, Bayer AG
Pradeep Balachandran, Ehealth Consultant
Alixandro Werneck Leite, LAMFO UnB
Danny Xie Li, Tecnológico de Costa Rica
Kamran Ali, Hasso Plattner Institute
Luis Oala, Fraunhofer HHI

# WG-DAISAM, Evaluation and Reporting: From Paper to Practice - to Software



**Paper**

**Practice**

Oala, Luis, Jana Fehr, Luca Gilli, Pradeep Balachandran, Alixandro Werneck Leite, Saul Calderon-Ramirez, Danny Xie Li et al. "**ML4H Auditing: From Paper to Practice**." In *Machine Learning for Health*, pp. 280-317. PMLR, 2020.
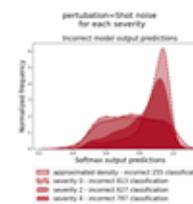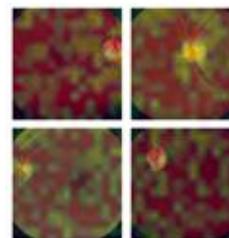
**Software**

Evaluation Package (EP)

Reporting Package (RP)

2019

2020

2021

# Evaluation Package: sub-streams

- Perturbed minds (research)
  - Plausible perturbations for robustness benchmarking in medical imaging
  - Fridays, 13.00 hrs Geneva time
  - Contact: Luis, Bruno
  - Outputs: Scientific publications, software
- Model reporting questionnaire (research)
  - A model reporting survey
  - Thursdays, 16.00 hrs Geneva time
  - Contact: Jana
  - More info today at 12.10 hrs
- Good practices (regulatory)
  - Tuesdays, 14.00 hrs Geneva time
  - Contact: Pradeep
  - Output: DEL 2.1/2.2

- "aiaudit.org"
  - Crowdsourced collection of SOTA AI quality assurance methods
  - Contact: Luis, Pat
  - Output: Web resource
- Regulatory ML4H "ICD-11"
  - An ontology to rule them all
  - Contact: Pradeep, Christian
  - Output: Software tool

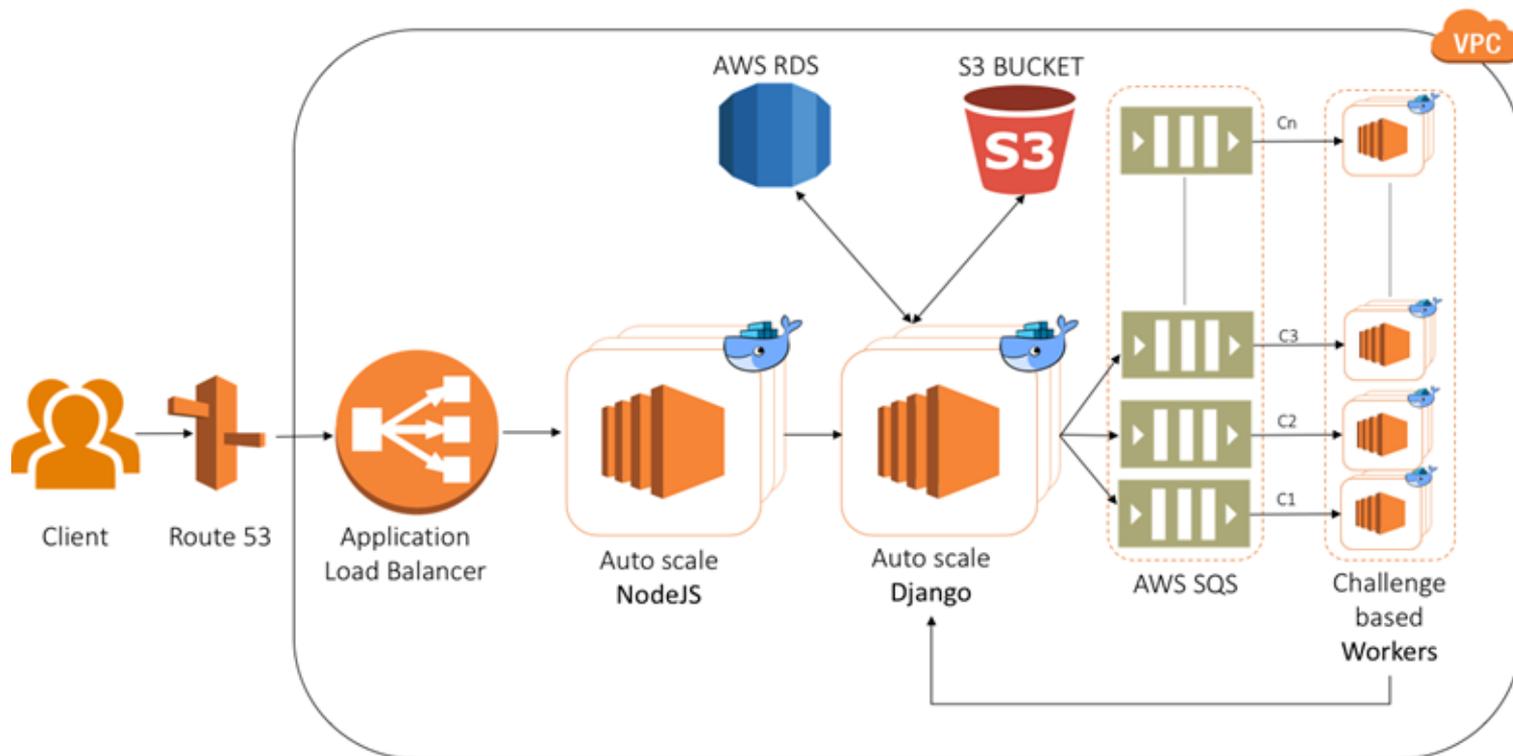| | |
|---|---|
| Luis | luis.oala@hhi.fraunhofer.de |
| Bruno | bruno.sanguinetti@dotphoton.com |
| Jana | Jana.Fehr@hpi.de |
| Pradeep | pbn.tvm@gmail.com |
| Pat | pat.baird@philips.com |
| Christian | christian.johner@johner-institut.de |

Ongoing    Upcoming

# Scope

- Compare performance ML algorithms

- Customizable benchmarking

- Upload own datasets

# EvalAi

- Open Source platform for evaluating ML algorithms

- Django-based

- Customizable benchmarking/evaluation as challenges

- Docker-based submissions by participants

- Computation in AWS

# Architecture

# Current State

- Adapted Frontend

- Implemented questionnaire for qualitative evaluation

- Demo version running on AWS

- Submission of dockerized diabetic retinopathy model

AI for Health

An ITU Focus Group
In collaboration with WHO

# Evaluating AI models in health

This platform is part of the AI4H assessment framework, developed by ITU in partnership with WHO. It enables you to test and compare the performance of your AI models.

Host Challenge    Participate

## Developed by

ITU

| 20 | 10.000 | 4 | 500 |
|---|---|---|---|
| Evaluation methods | Users | Organizations | Submissions |

AI for Health

Dashboard

**All Benchmarking Tasks**

Hosted Benchmarking

Contestant Teams

Results

⬛ Fraunhofer
Heinrich-Hertz-Institut

# Retinopathy Model Evaluation

Organized by: HHI_Fraunhofer
Starts on: Jan 1, 2019 1:00:00 AM
Ends on: Jun 1, 2099 1:59:59 AM

★ 0

ℹ Overview    📊 Evaluation    ⚡ Phases    ⬆ **Submit**    👁 My Submissions    📈 Leaderboard

## Please fill out this form before submitting a challenge.

1. For which purpose was the ML model developed for?

Diabetic retinopahy

2. Where was the dataset collected from?

3. Who created the dataset?

4. For what purpose was the dataset created?

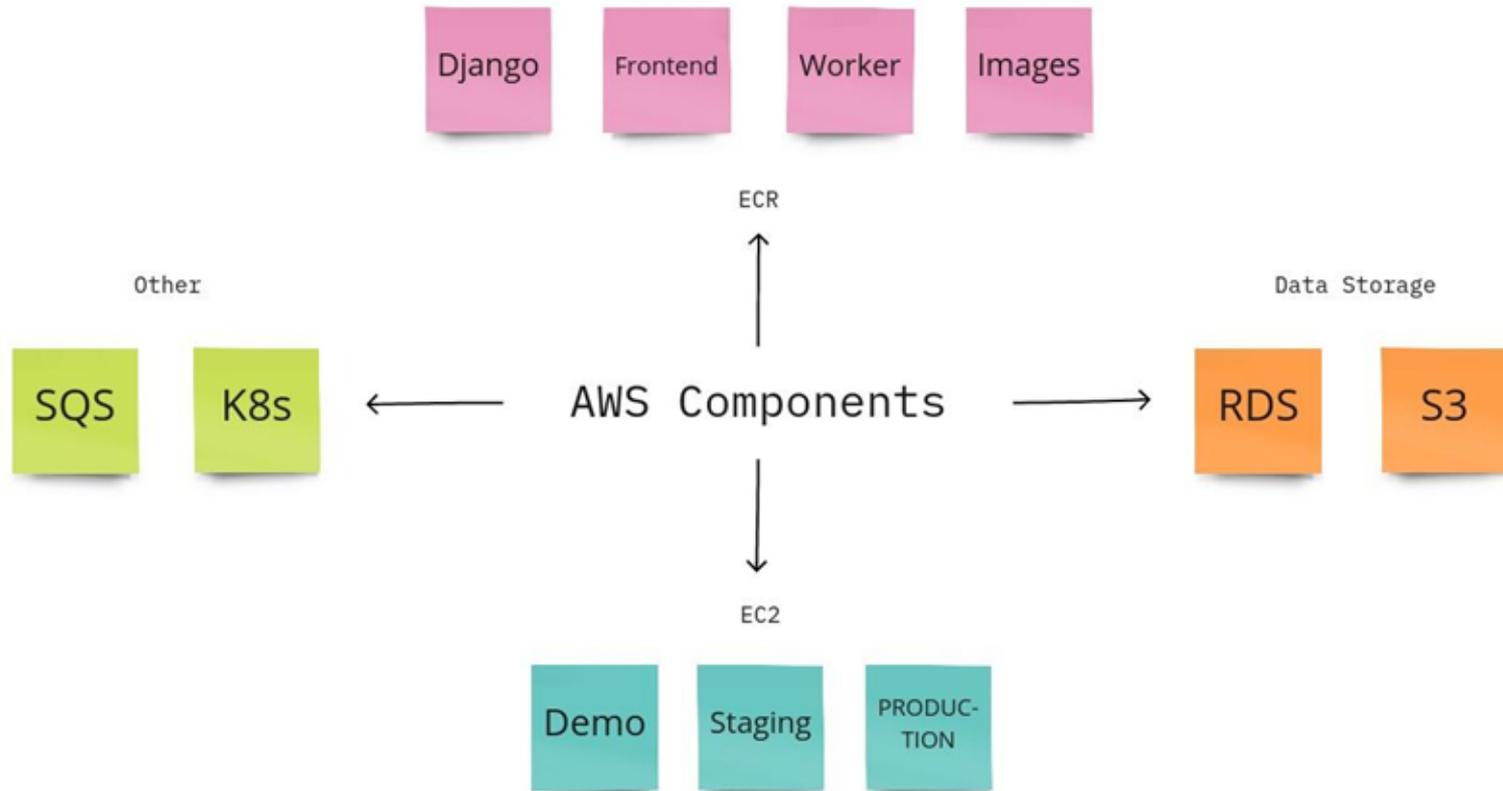5. When was the dataset collected?

mm/dd/yyyy

6. Are there archived versions of the raw/original dataset available?

7. What were the inclusion and exclusion criteria for the training dataset?

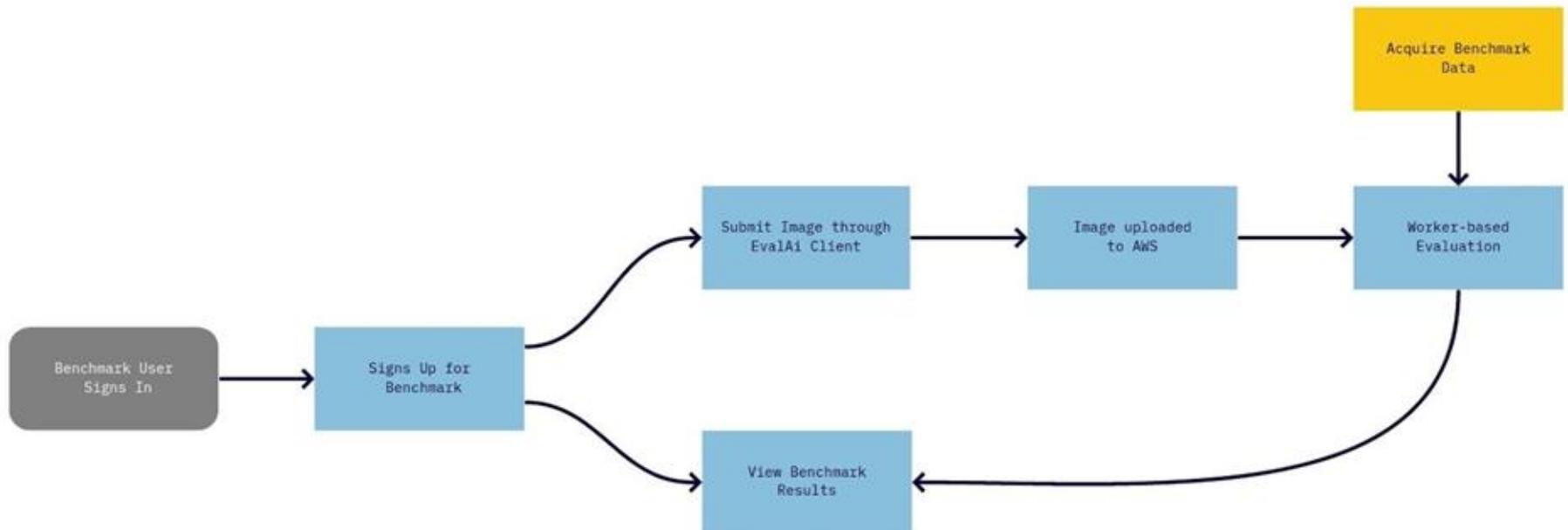8. What were the inclusion and exclusion criteria for the test dataset?

9. Were instances from the original dataset excluded for the ML model training?

10. Does the dataset contain confidential/personal information?

# Next Steps

- Writing evaluation script for diabetic retinopathy

- Execute full benchmarking cycle

- Adapt the questionnaire

- Interested to join and work with an interesting tech stack (Django, AWS, Docker...)? Contact
  - Elora

    [elora.schoerverth@hhi.fraunhofer.de](mailto:elora.schoerverth@hhi.fraunhofer.de)

  - Alixandro

    alixandrowerneck@outlook.com

# Open Code Initiative
# Reporting Package

**WG-DAISAM**

**Pradeep Balachandran**, Alixandro Werneck, Andrea Garcia, Danny Xie Li, Dominik Schneider, Elora Schörverth, Joachim Krois, Kamran Ali, Marc Lecoultre, Shobha Iyer, Shruti Choudhary, Steffen Vogler, Luis Oala

# Scope

- Prepares and presents AI4H model evaluation results generated by the Evaluation Package (EP)

- Provides a customizable reporting interface to support ease of comparison, classification and reproducibility of different types of AI4H model evaluation results

# Reporting service-workflow

# Architecture



| Module | Implementation Technology |
|---|---|
| Backend module | Django-VIEW (Python) |
| Database module | Django-MODEL (SQLIte DB) |
| Frontend module | Django-TEMPLATE file (HTML, CSS) |

# Current Status

Prototype implementation done using Django framework

**Backend**

✓ Web API to fetch the Diabetic Retinopathy model evaluation result data from the Evaluation Service

**Frontend**

✓ Web based GUI for model report ID submission

✓ Web based rendering template to display Diabetic Retinopathy model report

✓ Utility to enable downloading the model report in PDF format

# Binary Diabetic Retinopathy Model-Summary

| | |
|---|---|
| **Model Name** | Binary Diabetic Retinopathy Model |
| **Model Developer** | Xtend.AI |
| **Model Task** | Image Classification |
| **Model Algo** | CNN (Resnet 101) |
| **Model Output** | Disease Class Probability (Normal Vs DR) |
| **Accuracy** | 0.90 |
| **Sensitivity** | 0.90 |
| **Specificity** | 0.90 |
| **F-Score** | |
| **AU-ROC** | 0.96 |
| **Clinical Implications** | 1. Model serves as a tool for early detection of Diabetic Retinopathy( DR) in clinical / primary care setting 2. Model can be used to reject non-gradable and this reduces sampling errors and frees the clinician from looking at non-gradable images 3. Model can be used to prioritize the cases at higher-risk and refer them to a clinician 4. Model performance is comparable to the performance scores or the level of competence of the clinician/specialist/user in the clinical setting |
| **Safety Implications** | 1. Stored on secure servers. 2. Used SSL for all web access |
| **Efficiency** | 1. Model can be used to reject non-gradable images – which typically represent 10 – 20% of the input dataset. 2. This can increase efficiency by reducing sampling errors and freeing the clinician from looking at non-gradable images |

# Prototype Demo ⟶

# Next Steps

Module integration with Evaluation Package

- Report data schema validation with Evaluation Service

- Data-READ Web API validation with Evaluation Service
- Interested to join? Contact
    - Pradeep                              [pbn.tvm@gmail.com](mailto:pbn.tvm@gmail.com)
    - Alixandro
        alixandrowerneck@outlook.com