

DebiAI: Open-Source Toolkit for Data Analysis, Visualisation and Evaluation in Machine Learning

Tom Mansion
IRT SystemX
91120 Palaiseau, France
email:tom.mansion@irt-systemx.fr

Raphaël Braud
IRT SystemX
91120 Palaiseau, France
0000-0002-8717-2512

Ahmed Amrani
IRT SystemX
91120 Palaiseau, France
0009-0002-5390-4321

Sabrina Chaouche
IRT SystemX
91120 Palaiseau, France
0009-0005-2638-4182

Faouzi Adjed
IRT SystemX
91120 Palaiseau, France
0000-0002-0100-9352

Loïc Cantat
IRT SystemX
91120 Palaiseau, France
email:loic.cantat@irt-systemx.fr

Abstract—DebiAI is an open-source tool designed for data analysis, visualization, as well as evaluation and comparison of machine learning (ML) models. It is intended to be used both at the stage of the project data preparation, and for the evaluation of the ML models performances. It has a rich and user-friendly graphical interface that allows to visualize, analyze, select, edit and annotate data and metadata, as well as for bias detection and contextual evaluation of ML models. The tool relies on a generic data model, making it applicable to any type of ML task: classification, regression, object detection in images and more. It is an open source code distributed under the Apache License, Version 2.0 . The code is publicly available¹ and further information along with guidelines for the users can be found on its dedicated website².

Index Terms—Data Analysis; Data Visualization; Bias Detection; Contextual Evaluation; Machine Learning; Trustworthy AI.

I. INTRODUCTION

Data analysis and visualization are fundamental tasks in Machine Learning (ML) projects. They are playing a crucial role in a typical ML process, and are contributing not only in the data preparation phase but also during and after the model building. This emerging research topic, which combines several interactive systems and domains for ML processes, focused on human interaction and collaboration [1] created a new field that is Human-Centred Machine Learning (HCML) interaction [2]. Thus, a typical HCML framework allows an interactive visual analytic and ML evaluation [3]. Consequently, efficient tools are needed to assist, during the whole ML process, in data preparation and quality inspection before and after using the data for model training.

An effective tool should help enhance the iterative ML process across various phases: data preparation, analysis, anomaly detection, and annotation during the preparation phase, as well as the analysis and evaluation of models results. This aids in identifying the model's weaknesses and detecting issues at

the data level. In real-world ML projects, like for example in industrial ones, the data is enriched by metadata including operational context, expert knowledge, ...etc; that can help better understanding the raw data and can be informative for the learning process, and hence improving the quality of the model's learning and predictions. Having such tools contribute also into enhancing the trustworthiness of the used ML algorithm.

In ML-based engineering systems, it is crucial to guarantee key properties like accuracy, robustness, explainability, fairness, privacy, among many other primary values of Trustworthiness. Current research and development challenges of deploying trustworthy ML solutions are covered by wide programs like Confiance.ai³, the French AI flagship program to industrialize trustworthy AI-based critical systems [4] [5] and the TAILOR⁴ network at the European level.

DebiAI has been developed by the IRT SystemX in the framework of Confiance.ai program to contribute in ensuring trustworthiness by data, and serves as the main interface to view, analyze, select, edit and/or annotate any type of data and metadata. The rest of the paper is organized as following: Section II presents a brief state of the art of developed HCML tools, Section III presents methodology, followed by implementation description and application in a real use case described in Sections IV and V, respectively. Finally, conclusions and perspectives are drawn in Section VI.

II. LITERATURE REVIEW

Data visualization is the practice of representing information using graphical representations, employing technology-driven tools and software. Its fundamental objective is to enhance pattern recognition, improve understanding of complex concepts and facilitate in-depth exploration, thereby fostering the generation of new insights. Well-designed data visualizations can help in understanding large datasets and establishing

¹<https://github.com/debiai>

²<https://debiai.irt-systemx.fr>

³<https://www.confiance.ai>

⁴<https://tailor-network.eu/>

connections between ideas, concepts, and processing stages. Therefore, visual analysis can contribute to the optimization of AI approaches by actively participating in all aspects of the model building process [6] [7]. Similarly, Hohman et al. [8] highlights that successful ML applications often require iterations in data handling and continuous adjustments of the model. The authors introduced CHAMELEON, an interactive tool designed to attribute data iteration, thereby enhancing model performance, data validation, and the overall quality of ML projects. To improve data quality, Kandel et al. [9] presented Profiler, a tool using data mining to automatically detect issues and recommend coordinated visualizations for context-based assessment. Profiler offers methods for integrated statistical and visual analysis and view suggestions. ScrutinAI [10] is a Visual Analytics tool specifically tailored for enhancing the comprehension of deep neural network (DNN) predictions. Its primary objective is to identify and investigate potential weaknesses within models. To facilitate this, ScrutinAI provides interactive visualizations of input and output data, along with interactive plots and data filtering for comprehensive analysis of predictions. This tool is specifically designed for object detection and semantic segmentation, whereas DebiAI is applicable to a wide range of use cases. Zhang et al. [11] presented Manifold, a visual analytics platform designed for comparing and debugging ML models. The platform empowers users to categorize instances based on the model’s accuracy and confidence, identify symptomatic instances that generate incorrect results and continually help to enhance the model’s performance. As DebiAI, Manifold is created as a generic tool that operates independently of the internal logic of the ML model. It focuses on the input and the output. In order to improve the models’ performances and gain insights into their limitations, relying solely on the overall results from the test and training sets is insufficient. To address this limitation, the ModelTracker [12] tool offers instance-level visualization of results and enables users to individually inspect each instance. The tool has been applied to the binary classification task. In [13], the authors apply the approach for a multiclassification task. They employ Parallel Coordinates plots (PCP) to visually represent the multi-class predictions for a subset of instances. To complete these shortcomings and improve the understanding of model results and their comparison, the proposed DebiAI enables the analysis of model outcomes across various levels of granularity (instance, subset and dataset). This functionality has been applied to multiple tasks such as regression, classification, object detection in images, etc. In line with the approach outlined in [13], DebiAI leverages the use of PCP for model result analysis. However, the implementation of PCP within DebiAI is flexible, enabling its use not only for result analysis but also for assessing attributes.

III. METHODOLOGY

A. Contents

DebiAI is a web-based visual analytics application designed to support ML and data analysis. Its primary emphasis lies in

two crucial phases of the ML pipeline: pre-model construction and post-model building. It facilitates the development of ML models, aiding in data analysis during project preparation and enabling model performance comparison.

In the pre-model construction phase, DebiAI serves as an indispensable resource for data scientists and machine learning engineers during project preparation. It enables them to visually identify biases and errors in data inputs, detect anomalies and outliers throughout the data life cycle, assess data quality and domain coverage through relevant metrics, and select and analyze subsets of data to improve the quality of ML models.

In the post-model building phase, DebiAI serves as a visual analytics solution, simplifying the interpretation of ML model outputs. Its primary objective is to present model results in an intuitive and easily understandable manner, ultimately enhancing user confidence in the model’s predictions. Additionally, DebiAI offers features for identifying model weaknesses, comparing performance, and evaluating model effectiveness. This comprehensive approach fosters ongoing model refinement, tailored to the specific needs of the use case.

In both phases, DebiAI provides tools for creating and sharing statistical visualizations of project data and results with collaborators (team or clients).

B. Functional Description

DebiAI is an intuitive visualization tool designed to simplify the creation of interactive dashboards, empowering users without programming skills. It offers a diverse set of graphical widgets, including charts, tables, parallel categories, parallel coordinate plots, interval plots, night stars plots, and sample arrays. Moreover, DebiAI provides a user-friendly and flexible solution for interactive dashboard design, allowing users to effortlessly configure, adjust, resize, and position these widgets within their dashboards, ensuring the utmost customization of data presentation. This includes the ability to generate and share statistical visualizations of project data with team members or clients, fostering collaboration and informed decision-making by providing clear insights into the data. One of DebiAI’s standout features is its dynamic data selection and filtering capability, encouraging continuous exploration. Users can effortlessly create data subsets (selections) and apply filters based on various variables and contexts. This ensures that the dashboard consistently presents the identified subset of data. Furthermore, DebiAI assists users in identifying biases and inaccuracies in input, results, project data context, or ground truth, thereby improving data integrity. DebiAI facilitates the evaluation and comparison of ML model performance within the whole dataset or specific data subset. The tool enables the analysis of results across multiple levels of granularity. Indeed, the models’ performances are calculated at the level of each instance. Consequently, it is possible to identify the contexts, or combination of contexts in which the models encounter difficulties. The tool also simplifies the generation and organization of datasets, supporting in-depth analysis and potential retraining. DebiAI relies on a generic data model

that facilitates seamless application across various datasets, data types and use cases while maintaining consistent data processing practices. This essential feature provides DebiAI with flexibility, allowing it to transition effortlessly between various datasets or model results. In addition to its visualization capabilities, DebiAI incorporates implementations of statistical measures to enhance graphs, such as correlation analysis using methods like Pearson or Spearman. To support these visualizations, DebiAI also integrates techniques for discretizing continuous variables. DebiAI enables the utilization of internal or external algorithms to compute metrics or indicators on the data. Consequently, these metrics can be calculated either before integrating the data into DebiAI or during the data analysis phase on the tool. Various types of calculations can potentially be carried out by these algorithms, including the computation of new features, the assessment of model result quality, as well as indicators of data quality and distribution.

IV. IMPLEMENTATION DESCRIPTION

A. DebiAI Technical architecture

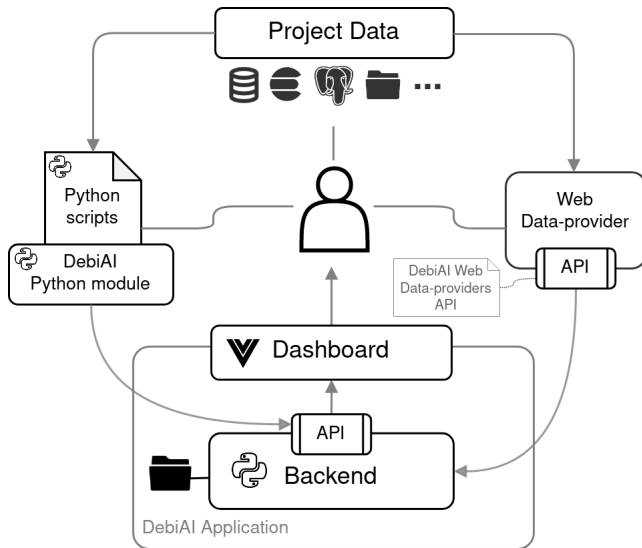


Fig. 1: DebiAI Architecture overview

The architecture of DebiAI is bifurcated into two main environments (Fig. 1):

The project environment :

- **Project Data:** This is the source of data that the user intends to analyze. It may originate from a variety of sources and formats, such as a CSV or JSON.
- **Python Scripts and DebiAI Python Module:** Using the DebiAI Python module, users can adapt their existing scripts and workflows to create selections and insert data and model results into DebiAI.
- **Web Data-Provider:** These are the services created by the user's project that enable DebiAI to fetch data directly from the project's data sources. A Web Data-Provider can be developed using any programming language, access

data from any type of database, and be hosted on any server. The only stipulation is that it should to implement and expose a specific REST API according to a defined contract.

DebiAI application environment:

- **DebiAI Web Dashboard:** This is the user interface of DebiAI, developed using VueJs. It provides users with an interactive platform to manage and view their data, and is hosted and served by the DebiAI backend.
- **Backend and API:** This is a Python-powered backend that not only provides an API but also serves the Web dashboard. This API is employed by the dashboard for data retrieval and by the Python module for data insertion. Additionally, it manages communications with the Web Data-providers and processes computational requests made by the dashboard.
- **Data storage:** DebiAI uses a folder-based data store that contains data in a JSON format. This data store supports the DebiAI backend by retaining projects created by the Python module and some specific dashboard elements, including layout configurations for project dashboards.

B. DebiAI Generic Data Model

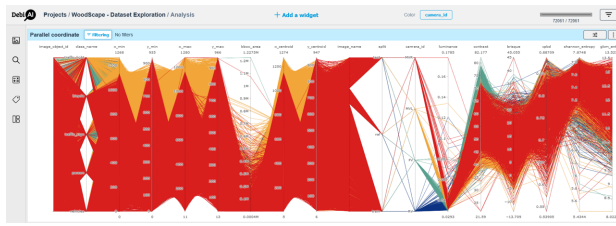
The data format required must follow a CSV structure with no missing values.

One of the most important features of DebiAI is its data model. The main objective is to enable the determination of the format of instances and the relationship between instances, models, model outputs, and model evaluation metrics per instance. Syntactically, each instance is composed of attributes, contexts, and annotations. The instance is linked to multiple ML models, where each model produces an output. Evaluation metrics are also associated with the model outputs. This structure is applicable to all types of data and ML tasks (classification, regression, object detection, etc.).

C. DebiAI Data Integration Process

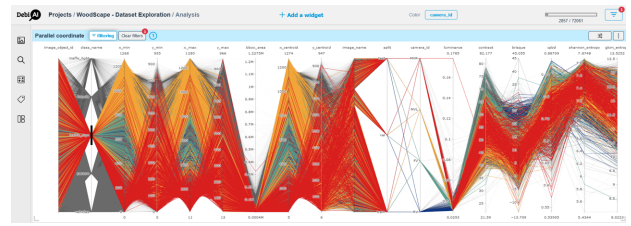
DebiAI offers two main ways to add data, each suited for different types of users and projects:

- 1) **Python Module:** This principal method enables seamless integration of project data into DebiAI via a dedicated Python module. Made for an integration within Python workflows, this approach, for example, facilitates the direct transfer of model results post-evaluation. This method is especially handy for those who primarily use Python.
- 2) **Data Providers:** Alternatively, DebiAI can interface with data through RESTful services, termed 'Data Providers'. This method is database-agnostic, allowing DebiAI to directly request project data, thereby making the data loading process faster and more efficient. Unlike the Python module, it doesn't require DebiAI to duplicate data within its integrated database. Although setting up a Data Provider is more time-consuming than using the Python module, it offers greater efficiency and



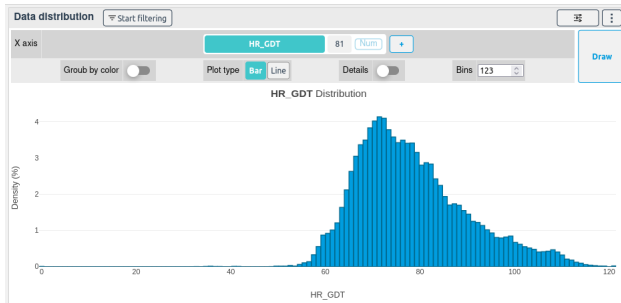
(a) original data

Filtering
↑



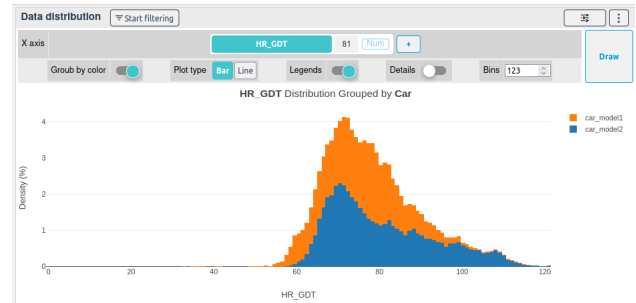
(b) Filtered data

Fig. 2: parallel coordinates widget. (a) represents the original data uploaded and (b) represents the same widget by selecting a subset of variables interactively.



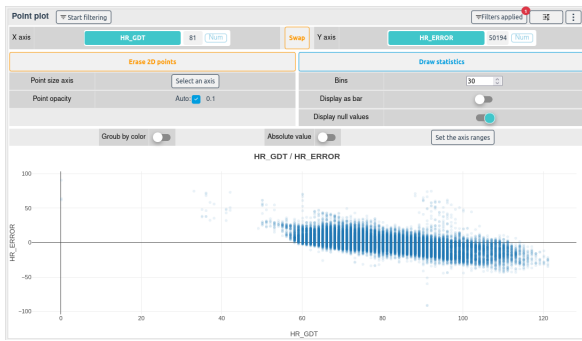
(a) data distribution of the data

Grouped by
↑



(b) data distribution + "grouped by" option

Fig. 3: illustration of data distribution by adding the option of "grouped by". (a) represents an example of data and (b) represents the same data grouped by another variable with two different colors.



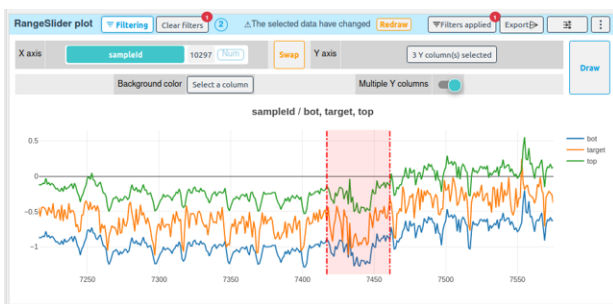
(a) original cloud data

Treatment
↑



(b) cloud data with statistical measurements

Fig. 4: Statistical treatment for an example of cloud data. (a) represents the original data and (b) represents the data by adding a set of statistical measures. Here illustrated measures are: mean, standard deviation, min and max and deciles



(a) original time series

Analysis
↑



(b) times series selected and analyzed

Fig. 5: Statistical analysis applied on times series example. (a) represents the original time series (b) represents the times series filtered and analyzed by adding a set of statistical measures

flexibility. This is particularly beneficial for long-term projects that regularly update their data.

Each method offers distinct benefits, and the choice depends on the specific requirements and scale of the project.

V. DEBIAI APPLICATION

A. Introduction

DebiAI is built upon a generic data model, and does not depend on data type (for instance images and time series), making it applicable to various use cases across a multitude of datasets. This intrinsic adaptability allows it to be valuable in a wide range of scenarios. It demonstrates its utility in the analysis of time series data, simplifying essential tasks such as regression. Furthermore, its functionality seamlessly extends to computer vision applications. Indeed, DebiAI provides tailored visual support for each stage of the process, enhancing models in tasks such as objects detection and image classification.

A starting example is provided⁵. In the following two sections, we will present the use of multiple screens in DebiAI for various use cases and provide an overview of a use case related to 2D objects detection in images.

B. DebiAI Visual Functionalities

As described in section III-B, DebiAI gives the ability to visualize and create interactive dashboards. Moreover, it can visualize various data types such as time series, point clouds, and tabular data, displaying computed attributes derived from images. However, for image viewing, it can establish links with external tools. In this section, we review a set of graphics implemented on different datasets with different data types. We also illustrate the main filtering features proposed by DebiAI. Five graphical visualizations are presented by exploring the parallel coordinates, the data distribution, the points plots, and time series widgets enhanced with interactive options. Fig. 2 illustrates a visualization of a dataset by using a parallel coordinates and the possibility to filter⁶ directly a set of variables. Another graphic visualization to analyse data distributions variables with the possibility of grouping by other variables is shown in Fig. 3. The third visualization selected from DebiAI is the possibility to apply statistical measures. Fig. 4, captures a data cloud visualized with its primary statistical measures; an envelope of min and max of the data, the average, the confidence interval $\pm\sigma$, where σ represents the standard deviation and also two deciles of the data. Fig.5 visually encapsulates the two distinct stages in the statistical analysis of a time series. Initially, Fig.5a displays the original data over an extended period. Subsequently, in the second stage, Fig.5b illustrates the time series after a more detailed analysis and filtering, focusing on a shorter timeframe. Among the statistical measures integrated into DebiAI, a noteworthy one to explore when analyzing dataset's variables is examining their correlation, a task effortlessly

⁵<https://debiai.irt-systemx.fr/dataInsertion/pythonModule/quickStart.html#requirements>

⁶The filtering widget is visualized in https://debiai.irt-systemx.fr/dashboard/widgets/1_parallel_coordinates/

accomplished using the DebiAI's correlation matrix widget. Fig. 6 displays the Pearson and Spearman correlation matrices for a set of a dataset's variables, showcasing a significant correlation between Shannon and glm entropies.



Fig. 6: Correlation Matrices of WoodScape dataset - left: Pearson; right: Spearman

C. WoodScape Dataset: 2D Objects Detection Case Study

The WoodScape dataset [14] is a public dataset containing more than 100.000 images of urban scenes captured using fish-eye cameras for automotive driving tasks. The images are provided by 4 different cameras with different angles of view (front, rear, middle right and middle left) with 360 coverage and have annotations for a diverse set of computer vision tasks. In our study, we focus on the 2D bounding boxes detection task with five classes: vehicle, person, bicycle, traffic light and traffic sign. We applied two versions of YOLO-based architectures, specifically YOLOv5 and YOLOv8. A YOLOv5 and a YOLOv8 models trained on WoodScape dataset and a YOLOv8 model pre-trained on COCO2017 [15].

The first step is to obtain a comprehensive overview of the data distribution, understand its scope and how it can be effectively used in a ML process. This comprehension is essential, as it helps in the formulation of effective training strategies. The following figures illustrate this process. The Fig. 7 shows the distribution of the dataset's train set composition of each class of the five cited above grouped by the camera id. By applying the same configuration to display the distribution of each of the three sets (train, validation, and test), we observe a similar distribution among the three of them. This confirms the appropriate split of the data. Nevertheless, the figure highlights the imbalance among the distributions of the five classes, suggesting the need for adaptive training techniques, for example using a weight sampler. It is essential to consider this imbalance when interpreting models outputs to avoid biased and skewed conclusions.

In the second stage, we used DebiAI to analyze the results of our models applied on the WoodScape test set and put them back into the context of the dataset and its features as shown in Fig. 8. This approach ensures an accurate interpretation of the models' outputs and provides potential improvements directions. The Fig. 8a illustrates the relationship between the precision and the recall of each model using the *night stars plot* widget which helps to navigate the trade-off between the two depending on the context of the task, for example, are we

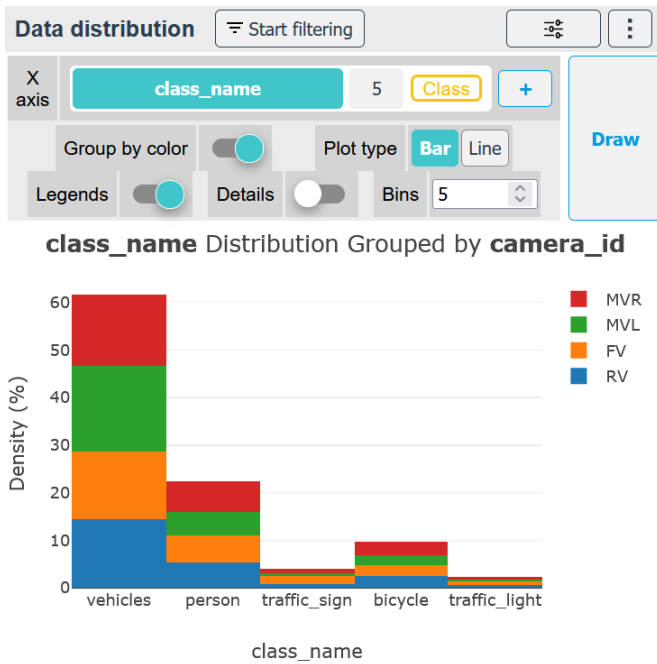


Fig. 7: Dataset’s train set split by objects class distribution grouped by Camera ID

prioritizing the detection quality over the quantity and vice-versa. The Fig. 8b shows the distribution of the f1-score of each model by the camera id, where we can easily spot the gap in performance between the three models: having the two models trained on WoodScape dataset showing higher scores compared to the one pre-trained on COCO 2017 dataset, which is expected giving the discrepancy between the two datasets. We can also notice that the Yolov5 trained on the WoodScape train set has better score on the data coming from the front and rear view cameras (FV & RV) of the vehicle while the Yolov8 also trained on WoodScape shows a greater score on the middle view cameras (left & right) data. This first observations led to further investigations using DebiAI in an attempt to understand the models’ outputs; you can check our tutorial on this case study on our website for more details.

VI. CONCLUSIONS AND PERSPECTIVES

In this paper, we introduce DebiAI, a versatile web-based visual analytics tool that enhances data preparation process, quality assessment, and model results analysis and comparison in ML projects. Its adaptability to various use cases and user-friendliness make it a valuable asset contributing to the trustworthiness in AI. We have illustrated its application in various use cases, primarily in a 2D objects detection task for driving assistance. As machine learning evolves, DebiAI can play a pivotal role in ensuring reliable and interpretable ML outcomes, solidifying its relevance in the field. In DebiAI’s outlook, the first priority is to simplify the data integration process, enhance interoperability with the learning process to retrieve and analyze data from each cycle. The concepts of robustness and explainability are also tied to model quality.

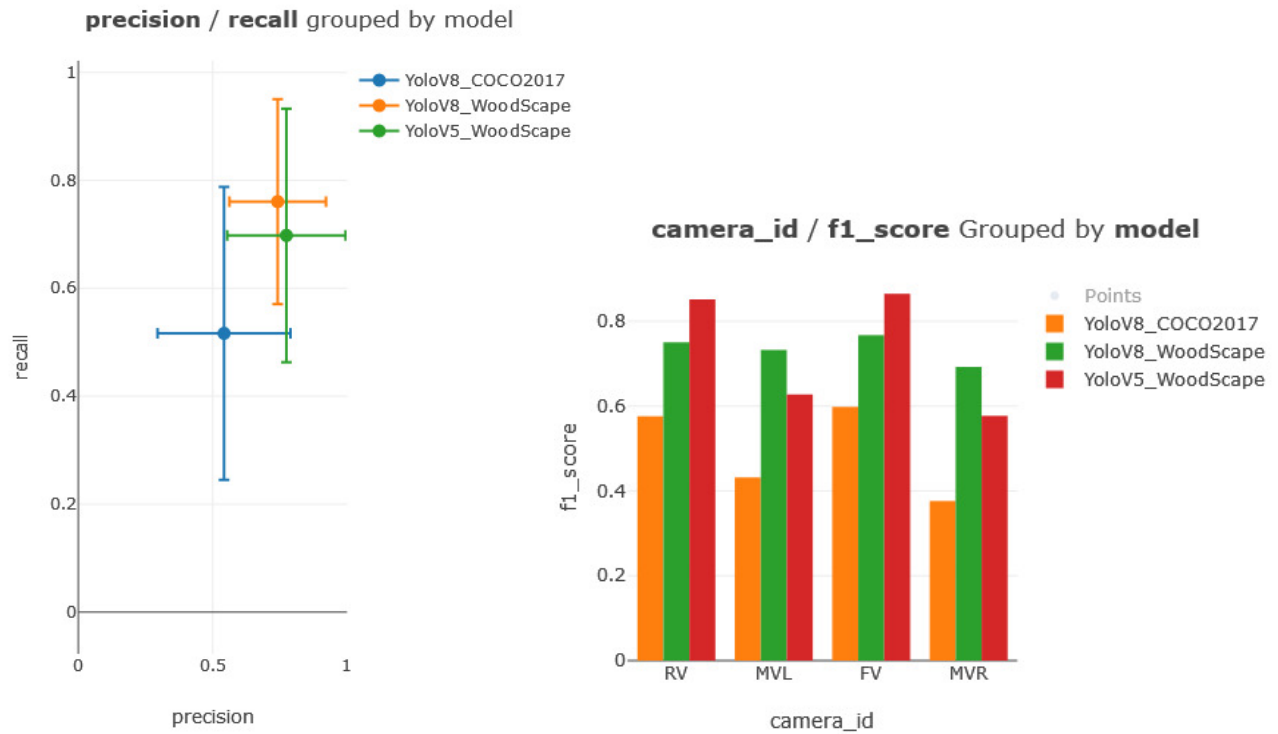
Therefore, incorporating these metrics into the process is critical for overall trust. Furthermore, we are considering coupling DebiAI with a data version control system to ensure traceability of dataset and model results evolution.

ACKNOWLEDGMENT

This work has been supported by the French government under the “France 2030” program, as part of the SystemX Technological Research Institute. This work was conducted as part of the Confiance.ai program, which aims to develop innovative solutions for enhancing the reliability and trustworthiness of AI-based systems.

REFERENCES

- [1] F. Sperrle, M. El-Assady, G. Guo, R. Borgo, D. H. Chau, A. Endert, and D. Keim, “A Survey of Human-Centered Evaluations in Human-Centered Machine Learning,” *Computer Graphics Forum*, vol. 40, no. 3, pp. 543–568, Jun. 2021, publisher: John Wiley & Sons, Ltd. [Online]. Available: <https://doi.org/10.1111/cgf.14329>
- [2] T. Kaluarachchi, A. Reis, and S. Nanayakkara, “A review of recent deep learning approaches in human-centered machine learning,” *Sensors*, vol. 21, no. 7, p. 2514, 2021.
- [3] D. Sacha, M. Sedlmair, L. Zhang, J. A. Lee, J. Peltonen, D. Weiskopf, S. C. North, and D. A. Keim, “What you see is what you can change: Human-centered machine learning by interactive visualization,” *Neurocomputing*, vol. 268, pp. 164–175, 2017.
- [4] M. Adedjouma, C. Alix, L. Cantat, E. Jenn, J. Mattioli, B. Robert, F. Tschirhart, and J.-L. Voirin, “Engineering Dependable AI Systems,” in *17th Annual System of Systems Engineering Conference (SOSE)*. Rochester, United States: IEEE, Jun. 2022. [Online]. Available: <https://hal.science/hal-03700300>
- [5] M. Juliette, L. R. Xavier, B. Bertrand, C. Loic, T. Fabien, R. Boris, G. Rodolphe, and N. Yves, “AI engineering to deploy reliable ai in industry,” *AI4I*, vol. 25, no. 1, 2023.
- [6] J. Yuan, C. Chen, W. Yang, M. Liu, J. Xia, and S. Liu, “A survey of visual analytics techniques for machine learning,” *Computational Visual Media*, vol. 7, no. 1, pp. 3–36, Mar. 2021. [Online]. Available: <https://doi.org/10.1007/s41095-020-0191-7>
- [7] J. Wang, S. Liu, and W. Zhang, “Visual analytics for machine learning: A data perspective survey,” *CoRR*, vol. abs/2307.07712, 2023. [Online]. Available: <https://doi.org/10.48550/arXiv.2307.07712>
- [8] F. Hohman, K. Wongsuphasawat, M. B. Kery, and K. Patel, “Understanding and visualizing data iteration in machine learning,” in *Proceedings of the 2020 CHI conference on human factors in computing systems*, 2020, pp. 1–13.
- [9] S. Kandel, R. Parikh, A. Paepcke, J. M. Hellerstein, and J. Heer, “Profiler: Integrated statistical analysis and visualization for data quality assessment,” in *Proceedings of the International Working Conference on Advanced Visual Interfaces*, 2012, pp. 547–554.
- [10] E. Haedecke, M. Mock, and M. Akila, “ScrutinAI: A visual analytics tool supporting semantic assessments of object detection models,” *Computers & Graphics*, vol. 114, pp. 265–275, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S009784932300105X>
- [11] J. Zhang, Y. Wang, P. Molino, L. Li, and D. S. Ebert, “Manifold: A model-agnostic framework for interpretation and diagnosis of machine learning models,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 1, pp. 364–373, 2019.
- [12] S. Amershi, M. Chickering, S. Drucker, B. Lee, P. Simard, and J. Suh, “ModelTracker: Redesigning Performance Analysis Tools for Machine Learning,” in *Proceedings of the Conference on Human Factors in Computing Systems (CHI 2015)*. ACM - Association for Computing Machinery, Apr. 2015, edition: Proceedings of the Conference on Human Factors in Computing Systems (CHI 2015). [Online]. Available: <https://www.microsoft.com/en-us/research/publication/modeltracker-redesigning-performance-analysis-tools-for-machine-learning/>
- [13] D. Ren, S. Amershi, B. Lee, J. Suh, and J. D. Williams, “Squares: Supporting interactive performance analysis for multiclass classifiers,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 23, no. 01, pp. 61–70, jan 2017.



(a) A Night Stars Plot of the Precision / Recall relationship

(b) f1-score results by Camera ID grouped by Models

Fig. 8: Some plots of the models results on the WoodScape dataset

- [14] S. Yogamani, C. Hughes, J. Horgan, G. Sistu, P. Varley, D. O’Dea, M. Uricár, S. Milz, M. Simon, K. Amende *et al.*, “Woodscape: A multi-task, multi-camera fisheye dataset for autonomous driving,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9308–9318.
- [15] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollár, “Microsoft coco: Common objects in context,” 2015.