

UI screenshot classifier

by Yurii Voievidka

github repo: <https://github.com/yvoievic/cv-project#ui-screenshot-classifier>

data: <https://huggingface.co/datasets/yuriivoievic/classified-ui>

Introduction

The world of UI design is vast and diverse, especially when it comes to desktop interfaces. Each UI serves a distinct purpose, and having a well-structured dataset of UI images is crucial for analysis and model training. Manually capturing screenshots of various UIs is not only tedious but also inefficient. Fortunately, an automated approach can be employed—extracting UI images directly from App Store data manifests of different applications.

However, there is a challenge: **not all screenshots contain valuable information**. Additionally, many of these images are from mobile applications, which are irrelevant to our study.

Objective

This repository aims to develop and evaluate methods for **classifying Desktop UI images** into three categories:

- **clean-ui** – Images that accurately and clearly represent the UI without unnecessary elements.
- **ui-to-crop** – Images that include redundant information, such as advertisements, macOS desktop backgrounds, or toolbars, which need to be cropped.
- **unnecessary** – Images that are irrelevant, such as mobile UI screenshots, text-heavy images without UI elements, or visuals that cannot be cleaned for meaningful use.

Approach

High-resolution image classification presents unique challenges, and this report explores multiple techniques to address them. For this purpose, a manually labeled dataset of **4,000 images** has been created, and the trained models will be evaluated on a larger dataset of **20,000 UI images**.

The classification experiments will be conducted using **four different approaches**:

1. **Vision Transformer (ViT) Fine-Tuning**
2. **CLIP Encoder + XGBoost**
3. **BLIP2 Encoder + XGBoost**
4. **SigLIP2 Encoder + XGBoost**

By comparing these methods, we aim to determine the most effective and efficient approach for classifying desktop UI images.

Hypothesis

I hypothesize that **SigLIP2** has the most advanced encoder architecture, which will lead to the highest performance metrics among the tested approaches. Given its ability to capture complex visual patterns and semantic information effectively, I expect it to outperform other models in classification accuracy.

To evaluate model performance, the following metrics will be used:

- **Accuracy** – Measures overall correctness of predictions.
- **Precision** – Assesses how many of the predicted positive instances are actually correct.
- **Recall** – Determines how well the model identifies all relevant instances.
- **F1 Score** – Balances precision and recall, providing a comprehensive assessment of classification performance.

Through this experiment, I aim to validate whether **SigLIP2 truly provides superior classification performance** compared to other encoding methods.

Desktop UI images examples

Clean UI :

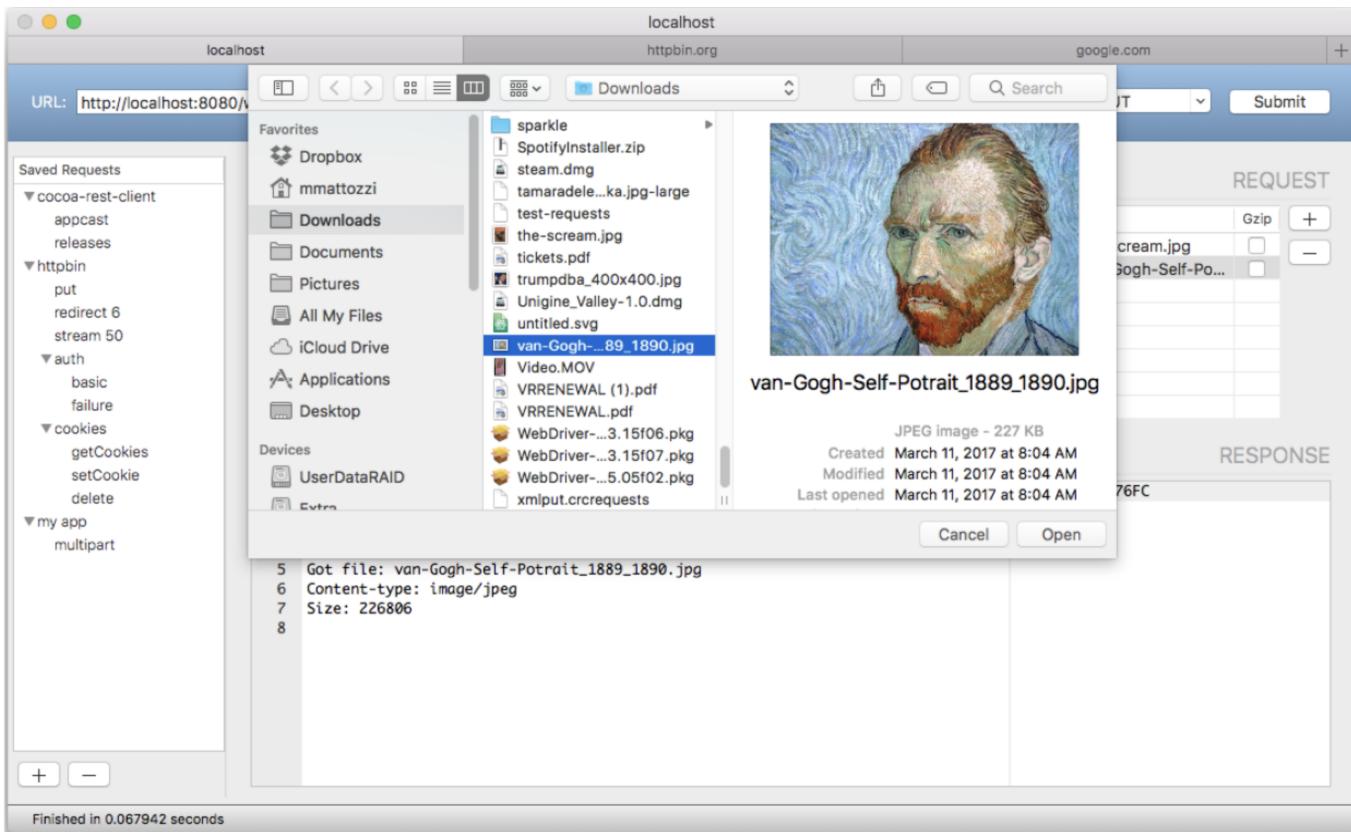
The screenshot shows a Trello board titled '#capturing-feedback'. The board is organized into several lists:

- Icebox:** Contains cards like "Better global button for capturing feedback" and "Once an idea is linked, paste the site2 URL of the idea either into the ticket text or make it easily accessible from the 'link to idea' section".
- Inbox:** Contains cards like "Allow agents to immediately search for newly created suggestions so they can link feedback to them".
- Next:** Contains cards like "[1] Make Edit button in legacy UI consistent with the rest of the interface" and "[2] Revert 'suggestion' to 'idea' nomenclature".
- Doing →:** Contains cards like "[1] Free plans: don't allow admin to select or create new categories in the capture feedback dialog" and "[2] Add unsubscribe links to suggestion status update and direct message notification emails".
- Testing:** Contains cards like "[1] Add time events" and "[2] Add additional Mixpanel events".

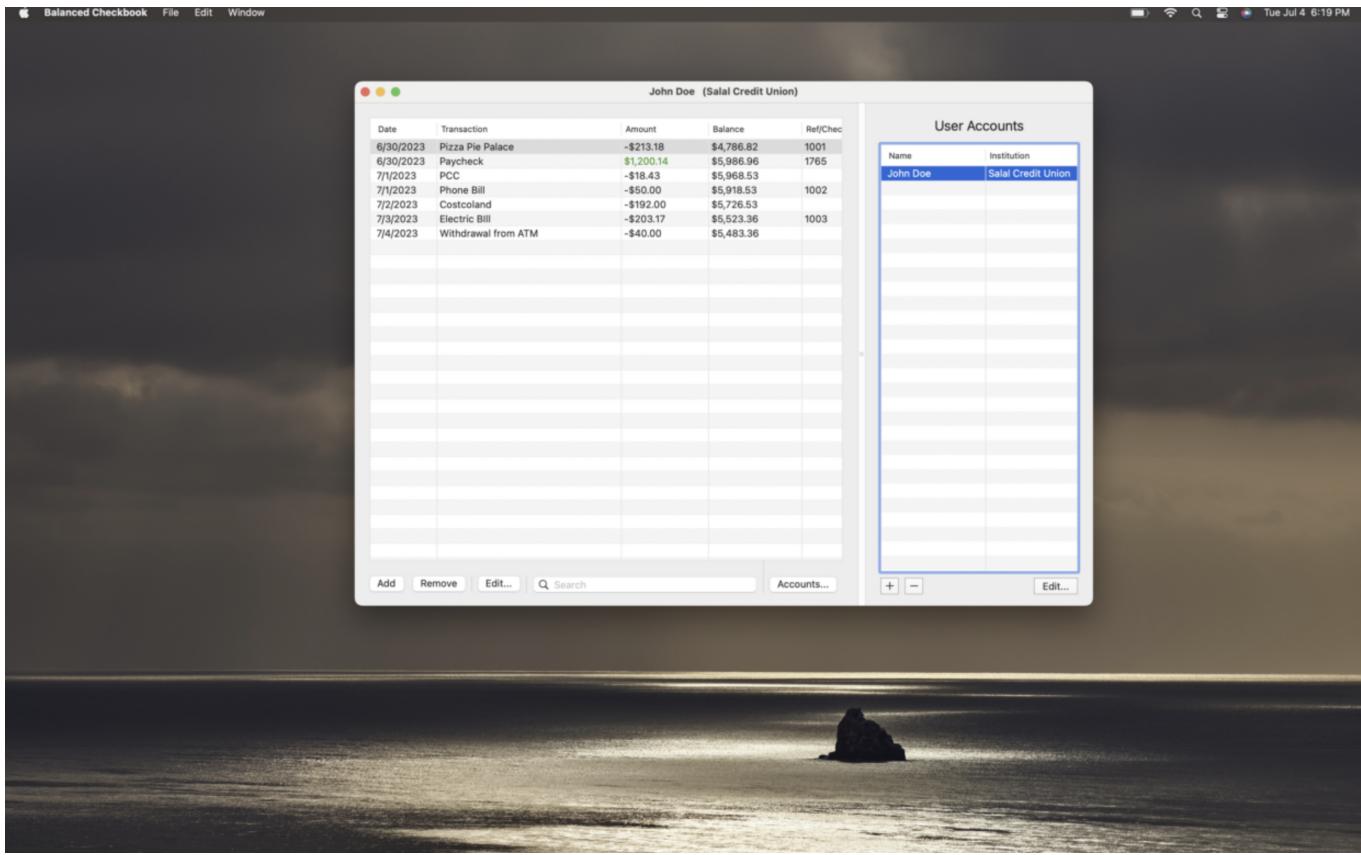
The interface is clean with a light blue background and white cards. Each card has a title, a description, and a small preview image. The Trello logo is visible at the top right.

The screenshot shows a medical search interface with the following components:

- Header:** Includes icons for Compendium, Favorites, and Interactions, and buttons for Print and CoMed.
- Search Bar:** A search bar with the placeholder "Mère". Below it are dropdown menus for "Préparation", "Titulaire", "Principe Actif / Code ATC", "No d'autorisation", "Thérapie", and "Plein Texte".
- Sidebar:** A sidebar on the left with a search bar and a list of search results:
 - Mères:** 7 Treffer
 - mère:** 1616 Treffer
 - mère-enfant:** 14 Treffer
 - mères:** 770 Treffer
 - mères/enfants:** 2 Treffer
- Content Area:** A list of pharmaceutical products:
 - Aciclovir Labatec® i.v.** | Labatec Pharma SA
Préclinique
 - Acide méfénamique Sandoz®** | Sandoz Pharmaceuticals AG
Préclinique
 - Acide zolédronique Onco Sandoz® 4 mg/100 mL, solution pour perfusion** | Sandoz Pharmaceuticals AG
Préclinique
 - Acide zolédronique Onco Sandoz® 4 mg/5 mL, concentré pour perfusion** | Sandoz Pharmaceuticals AG
Préclinique
 - Acivir® Crème** | Spirig Pharma AG
Préclinique
 - Acnatac® Gel** | MEDA Pharma GmbH
Préclinique
 - Actemra®** | Roche Pharma (Schweiz) AG
Préclinique
 - Acyclovir-Mepha 200/400/800, Comprimés** | Mepha Pharma AG
Préclinique
 - Acylovir-Mepha® i.v. 250** | Mepha Pharma AG
Préclinique
 - Adasuve** | Orpha Swiss GmbH
Préclinique
 - Adempas®** | Bayer (Schweiz) AG
Préclinique
 - Aggrastat** | Orpha Swiss GmbH
Préclinique
 - Alopxy® 2%** | Pierre Fabre (Suisse) S.A.
Préclinique
 - Amavita Lopéramide** | Amavita Health Care AG
Préclinique
 - Amitiza** | Takeda Pharma AG
Préclinique
 - Apydan® extent** | Desitin Pharma GmbH
Préclinique
 - Arnuity® Ellipta®** | GlaxoSmithKline AG
Préclinique
 - Arzerra®** | Novartis Pharma Schweiz AG
Préclinique
- Right Sidebar:** A sidebar on the right with a list of categories and their counts:
 - Grossesse / Allaitement (13)
 - Grossesse, allaitement (79)
 - Composition (4)
 - Grossesse / allaitement (2)
 - Surdosage (2)
 - Indications (8)
 - Cinétique (304)
 - Propriétés/Effets (23)
 - Conduite (1)
 - Grossesse / Allaitement (2)
 - Interactions (10)
 - Préclinique (410)
 - Contre-indications (10)
 - Effets indésirables (11)
 - Remarques (14)
 - Précautions (48)
 - Grossesse/Allaitement (1084)
 - Posologie (16)



UI To Crop:





Clean up hidden cache, trash & junk

IMPROVED!



613 MB
IN OUR TESTS

- ✓ Application caches
- ✓ iPhoto Library Trash
- ✓ Trash
- ✓ Mail Downloads
- ✓ Junk files
- ✓ Downloads

Unnecessary:

Open Multiple Account for WhatsApp

The image shows a WhatsApp interface on a Mac OS X desktop. It displays two open accounts: Lucy Elena and Zack Canson. Lucy Elena's account has several messages from contacts like Jessica Miller, Angelique Embry, Jack Theo, Navy Clover, and Thomas Jacob. Zack Canson's account also has messages from Jessica Miller. A prominent yellow banner at the bottom of the screen reads "2 Account 1 Device".



Visual Transformer Fine Tuning

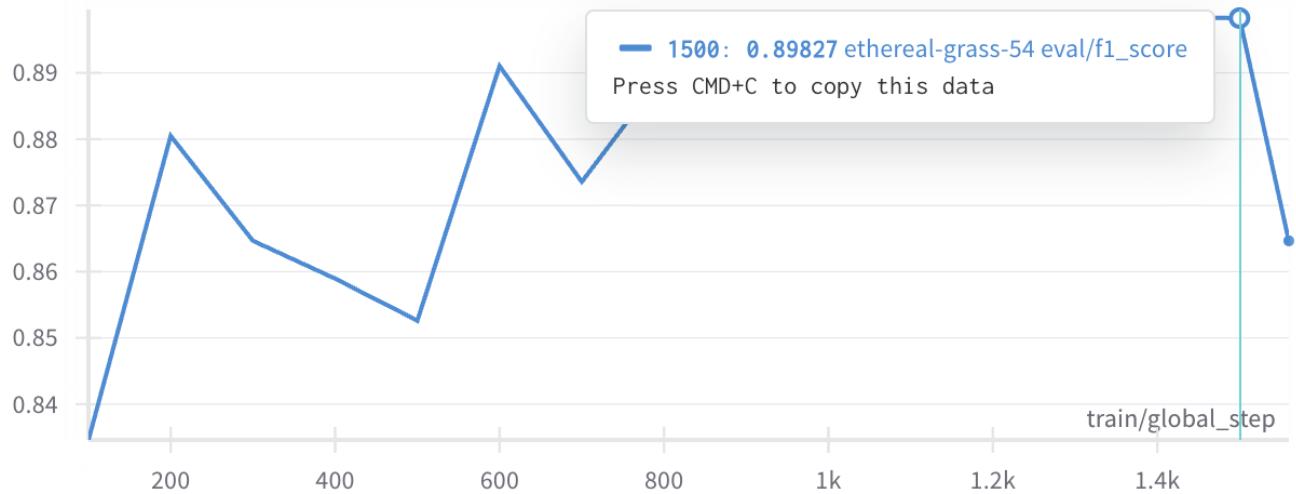
Start Date	2025-02-10
End Date	2025-02-11
dataset	DesktopUI

Motivation

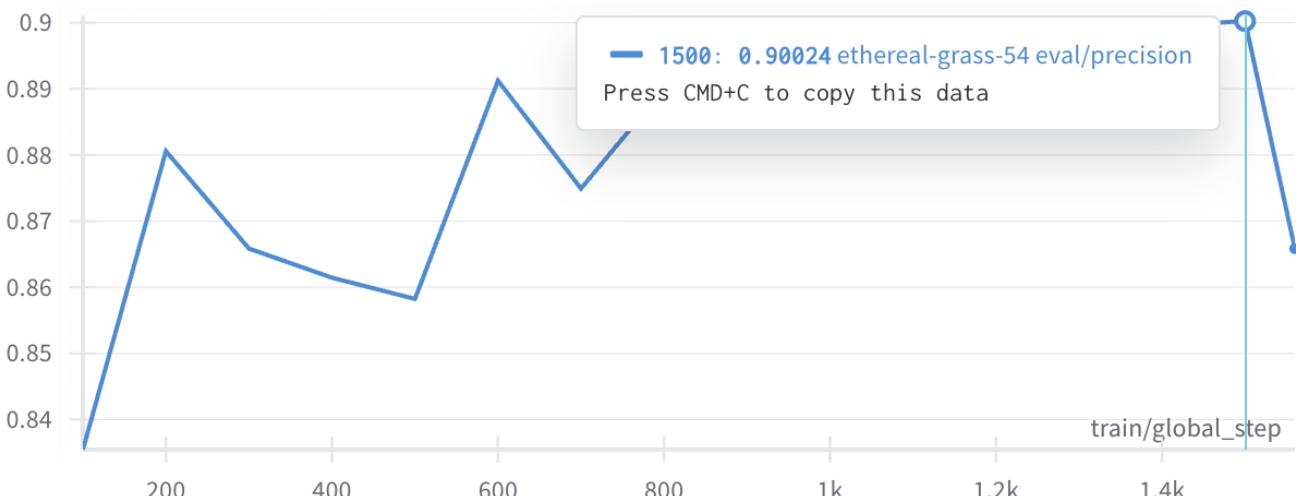
Try to fine tune Visual Transformer on custom dataset of desktop UI images.

Results

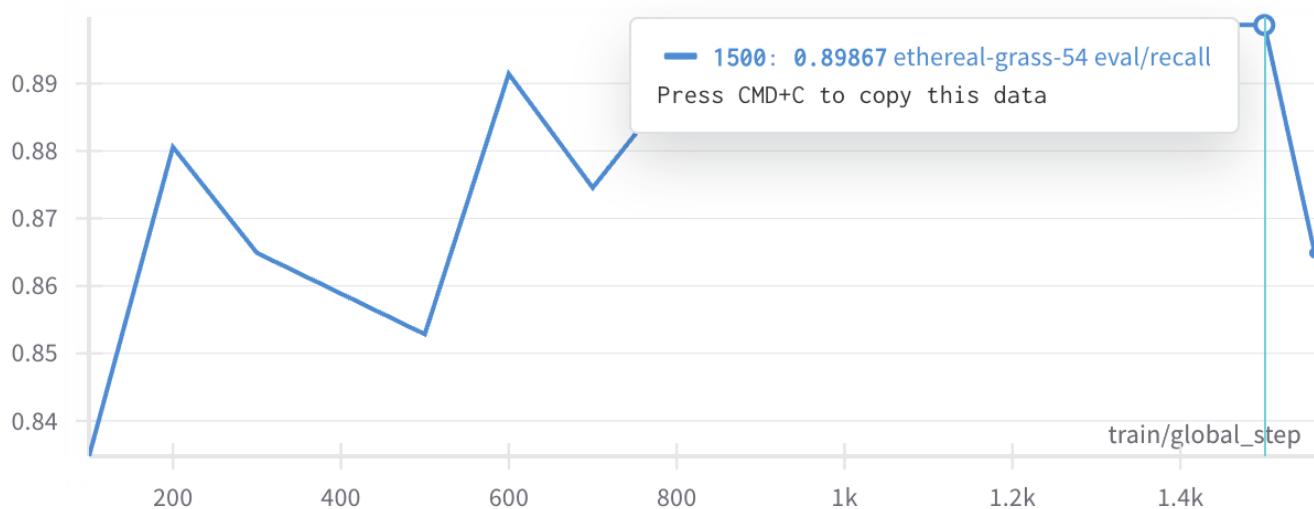
eval/f1_score



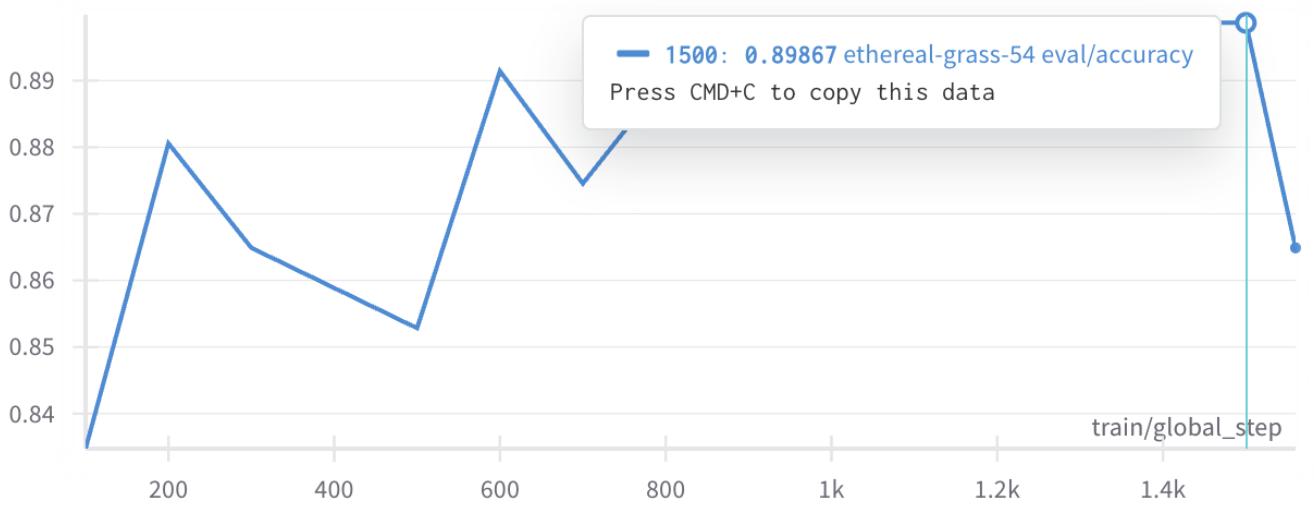
eval/precision



eval/recall



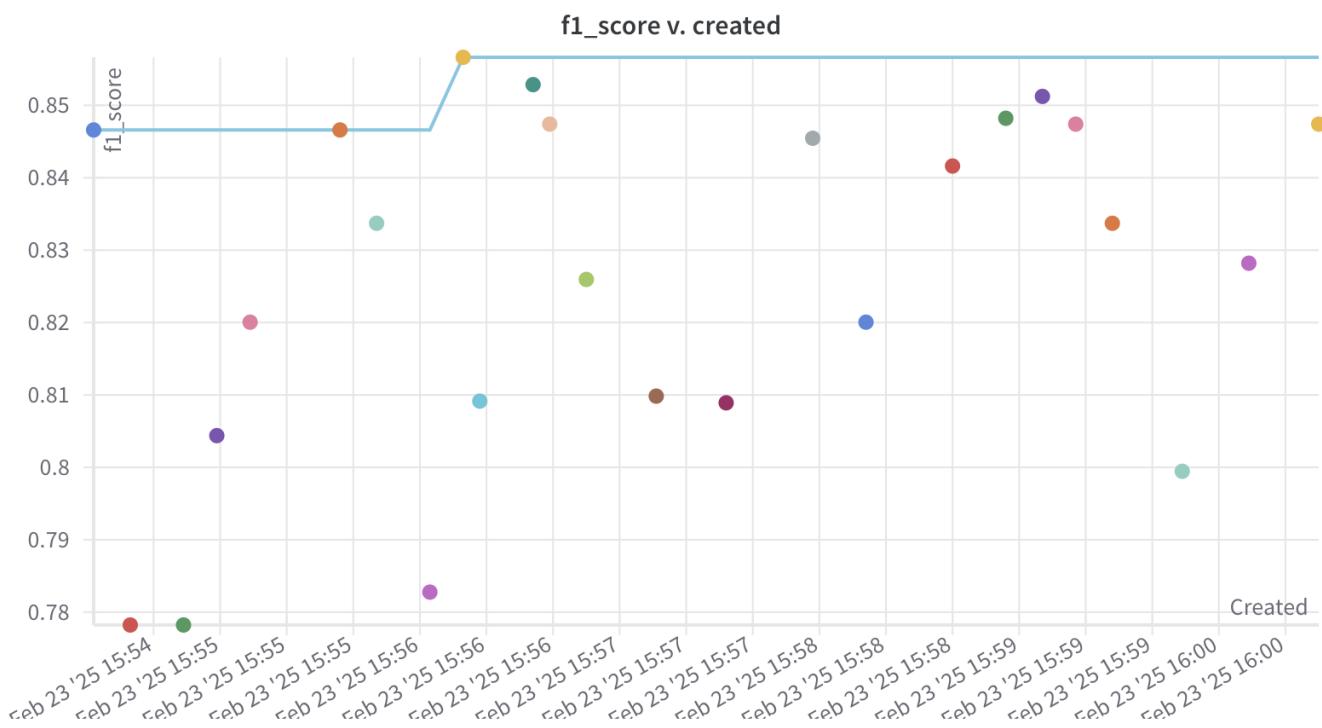
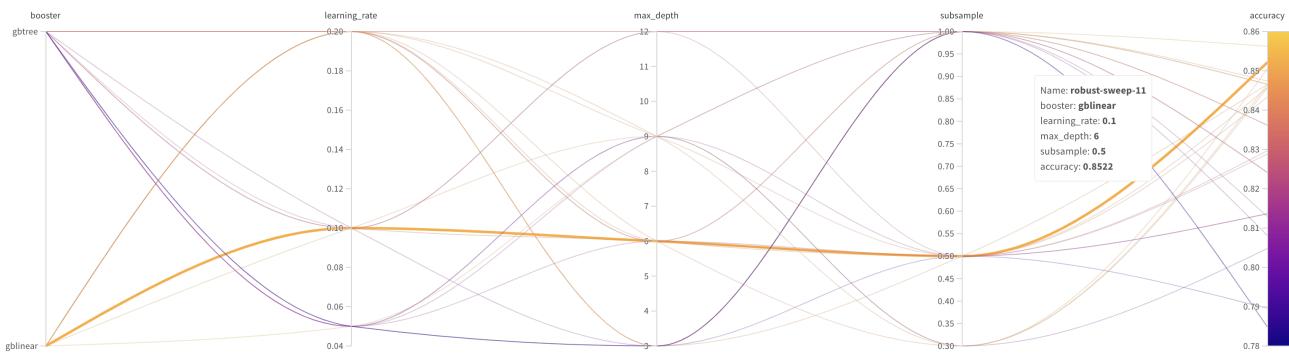
eval/accuracy

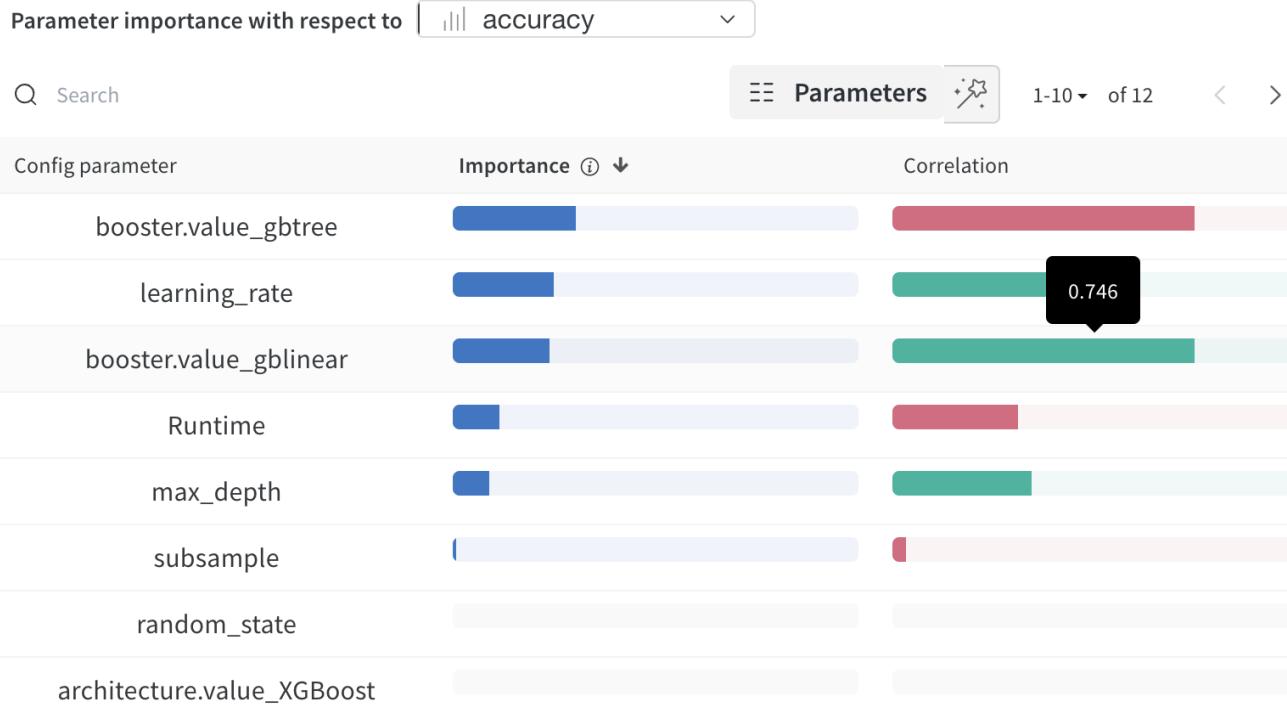


Metric	Value
step	1500
Eval Accuracy	0.898673
Eval F1 Score	0.898269
Eval Loss	0.541631
Eval Precision	0.900237
Eval Recall	0.898673

CLIP + XGBoost

Start Date	2025-02-13
End Date	2025-02-15
dataset	DesktopUI





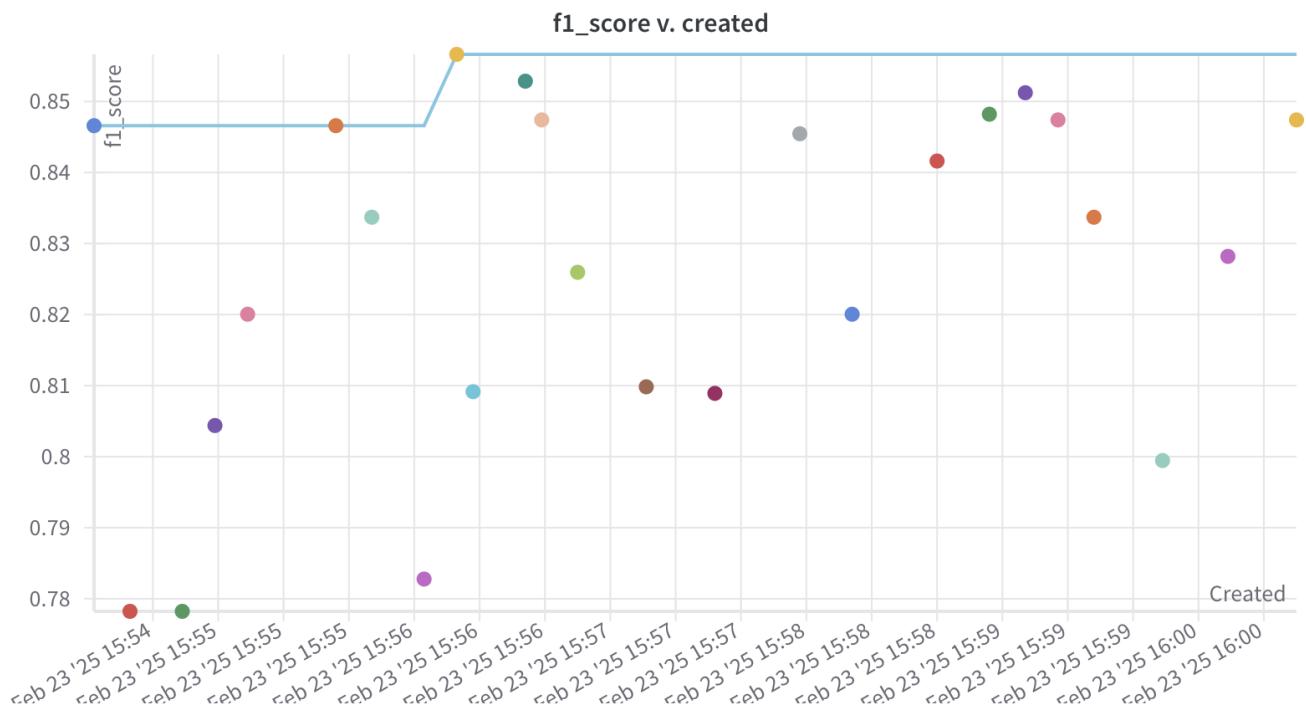
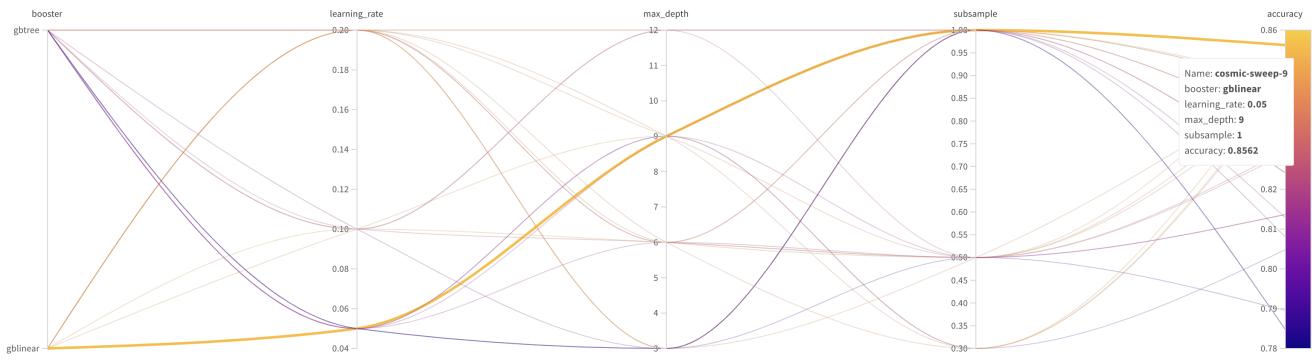
Results:

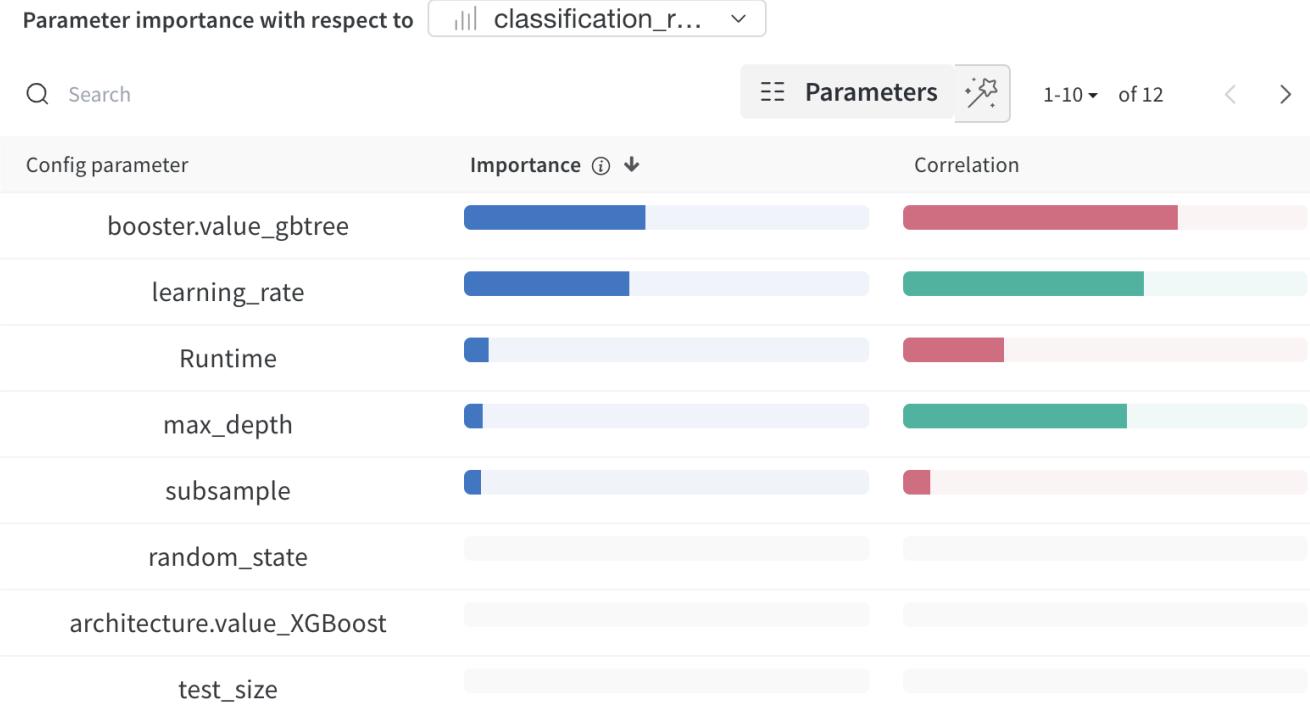
"Unfortunately, Weights & Biases did not save the history of metric progression. However, you can check them in the repository."

Metric	Value
Accuracy	0.43293
F1 Score	0.41927
Precision	0.41268
Recall	0.43293

BLIP2 + XGBoost

Start Date	2025-02-16
End Date	2025-02-19
dataset	DesktopUI



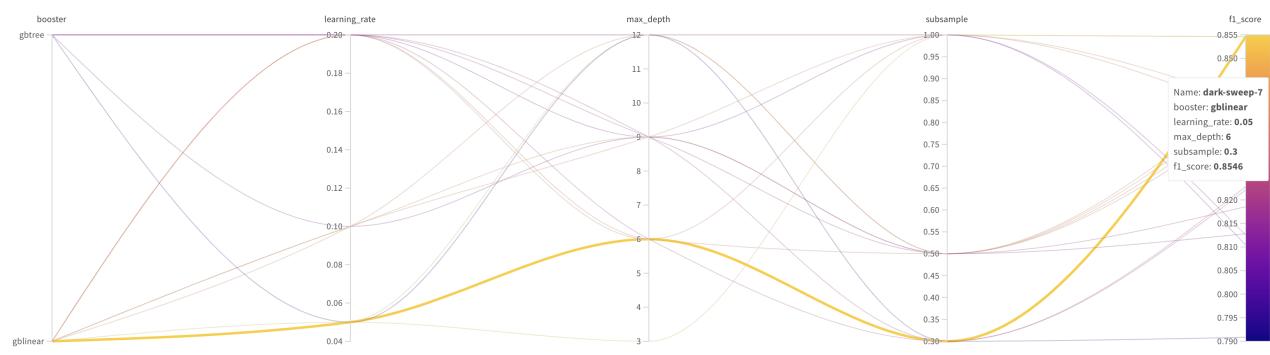


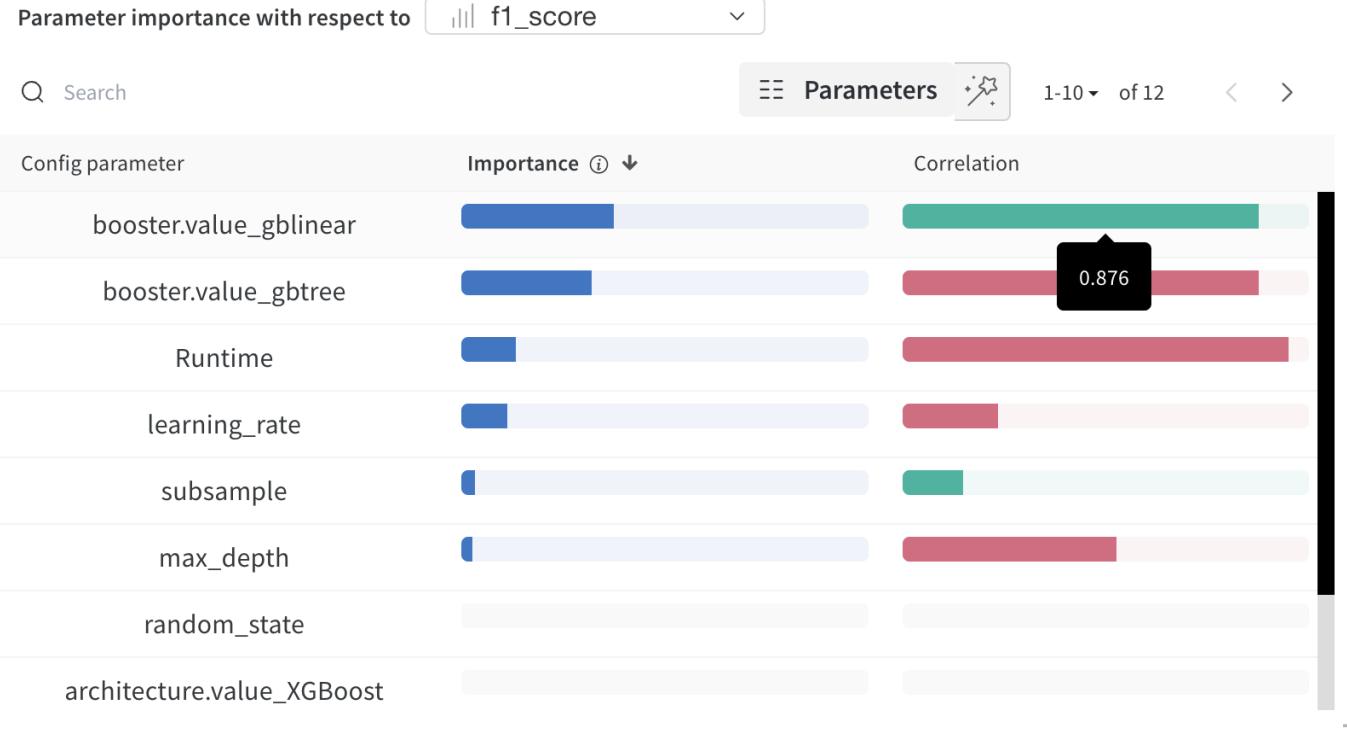
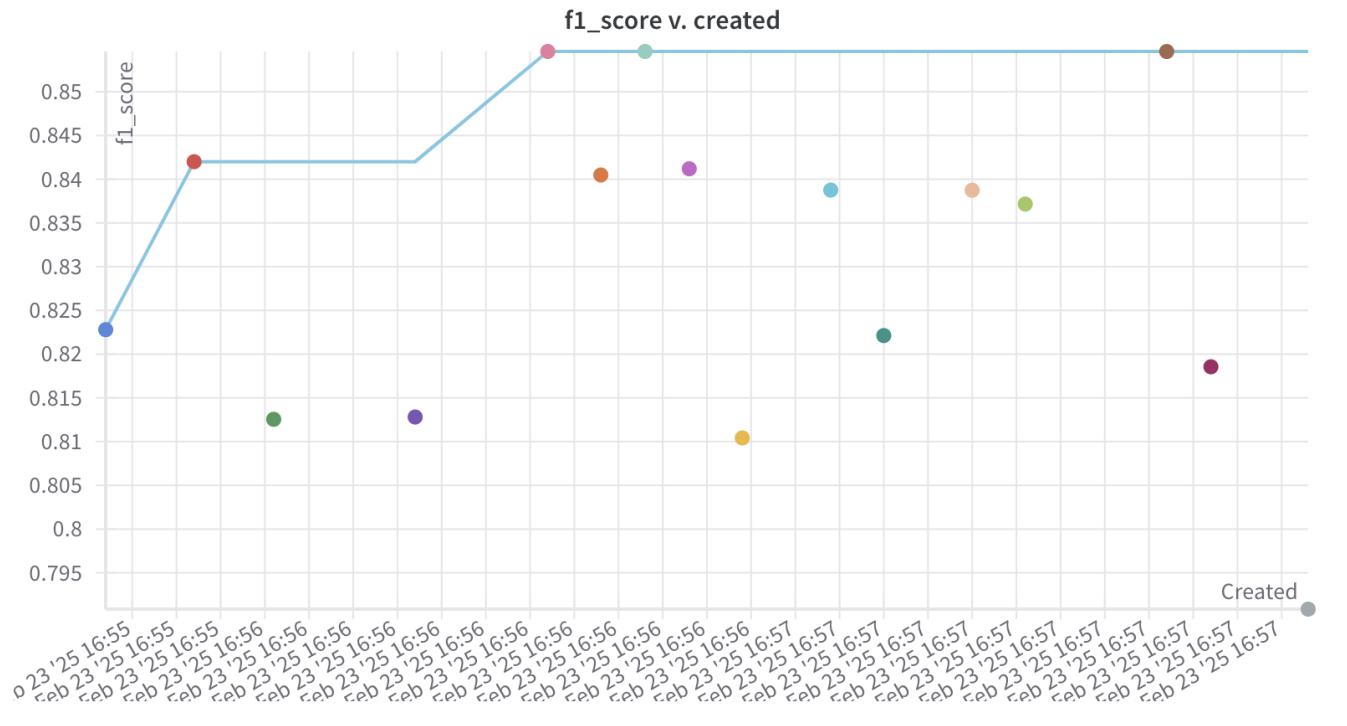
Results

Metric	Value
Accuracy	0.84659
F1 Score	0.8474
Precision	0.84967
Recall	0.84659

SigLIP2 + XGBoost

Start Date	2025-02-20
End Date	2025-02-23
dataset	DesktopUI





Metric	Value
Accuracy	0.85462
F1 Score	0.8546
Precision	0.85477
Recall	0.85462

Summary

Among the four approaches evaluated, the Vision Transformer (ViT) demonstrated the best performance, achieving an F1 score of up to 0.89 when trained for 10 epochs using FP16 precision. However, this improvement came at a significant computational cost—training ViT took approximately two hours on a GPU. In contrast, SigLIP2 embeddings were extracted within just 15 minutes on a CPU, making it a far more efficient alternative.

This observation leads to the conclusion that fine-tuning Visual Transformers can be an effective strategy for maximizing performance metrics. However, the process is both computationally expensive and time-consuming, making it less practical for scenarios with resource constraints. For future work, I would recommend exploring alternative approaches that strike a balance between efficiency and accuracy, such as lightweight transformer models, optimized embedding techniques, or self-supervised learning methods.