

GUing: A Mobile GUI Search Engine using a Vision-Language Model

JIALIANG WEI, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, France

ANNE-LISE COURBIS, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, France

THOMAS LAMBOLAIS, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, France

BINBIN XU, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, France

PIERRE LOUIS BERNARD, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, France

GÉRARD DRAY, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, France

WALID MAALEJ, University of Hamburg, Germany

App developers use the Graphical User Interface (GUI) of other apps as an important source of inspiration to design and improve their own apps. In recent years, research suggested various approaches to retrieve GUI designs that fit a certain text query from screenshot datasets acquired through automated GUI exploration. However, such text-to-GUI retrieval approaches only leverage the textual information of the GUI elements in the screenshots, neglecting visual information such as icons or background images. In addition, the retrieved screenshots are not steered by app developers and often lack important app features, e.g. whose UI pages require user authentication.

To overcome these limitations, this paper proposes GUing, a GUI search engine based on a vision-language model called UIClip, which we trained specifically for the app GUI domain. For this, we first collected app introduction images from Google Play, which usually display the most representative screenshots selected and often captioned (i.e. labeled) by app vendors. Then, we developed an automated pipeline to classify, crop, and extract the captions from these images. This finally results in a large dataset which we share with this paper: including 303k app screenshots, out of which 135k have captions. We used this dataset to train a novel vision-language model, which is, to the best of our knowledge, the first of its kind in GUI retrieval. We evaluated our approach on various datasets from related work and in manual experiment. The results demonstrate that our model outperforms previous approaches in text-to-GUI retrieval achieving a Recall@10 of up to 0.69 and a HIT@10 of 0.91. We also explored the performance of UIClip for other GUI tasks including GUI classification and Sketch-to-GUI retrieval with encouraging results.

1 INTRODUCTION

The Graphical User Interfaces (GUI) of mobile apps serve as the primary means of interaction between users and their devices. A well-designed GUI can markedly streamline the navigation process, facilitate the accomplishment of user tasks, and improve the overall user experience—contributing towards a higher user engagement and retention [12, 25]. Furthermore, a user-friendly, modern, and attractive GUI design can potentially differentiate an app from its market counterparts, thereby amplifying its likelihood of success in the highly competitive mobile app market [44].

To design a good GUI, app designers and developers may relay on multiple sources. In addition to the initial needs and requirements of the users, designers often initiate the ideation process with inspiration searches of other GUIs. Numerous GUI retrieval approaches have thus been suggested, enabling the retrieval of GUIs from existing apps. Researchers have proposed GUI retrieval approaches that accept sketches [26, 40, 41], wireframes [11, 15, 38], or screenshots [7, 35] as input to locate similar or fitting designs. While certainly useful, these approaches require a

Authors' addresses: [Jialiang Wei](mailto:jialiang.wei@mines-ales.fr), jialiang.wei@mines-ales.fr, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, Ales, France; [Anne-Lise Courbis](mailto:anne-lise.courbis@mines-ales.fr), anne-lise.courbis@mines-ales.fr, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, Ales, France; [Thomas Lambolais](mailto:thomas.lambolais@mines-ales.fr), thomas.lambolais@mines-ales.fr, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, Ales, France; [Binbin Xu](mailto:binbin.xu@mines-ales.fr), binbin.xu@mines-ales.fr, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, Ales, France; [Pierre Louis Bernard](mailto:pierre-louis.bernard@umontpellier.fr), pierre-louis.bernard@umontpellier.fr, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, Montpellier, France; [Gérard Dray](mailto:gerard.drays@mines-ales.fr), gerard.drays@mines-ales.fr, EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, Ales, France; [Walid Maalej](mailto:walid.maalej@uni-hamburg.de), walid.maalej@uni-hamburg.de, University of Hamburg, Hamburg, Germany.

preliminary graphical prototype, which might not be available or might be too restrictive for the ideation in early development phases. Therefore, text-to-GUI retrieval approaches are more helpful when only ideas and requirements in text form are available.

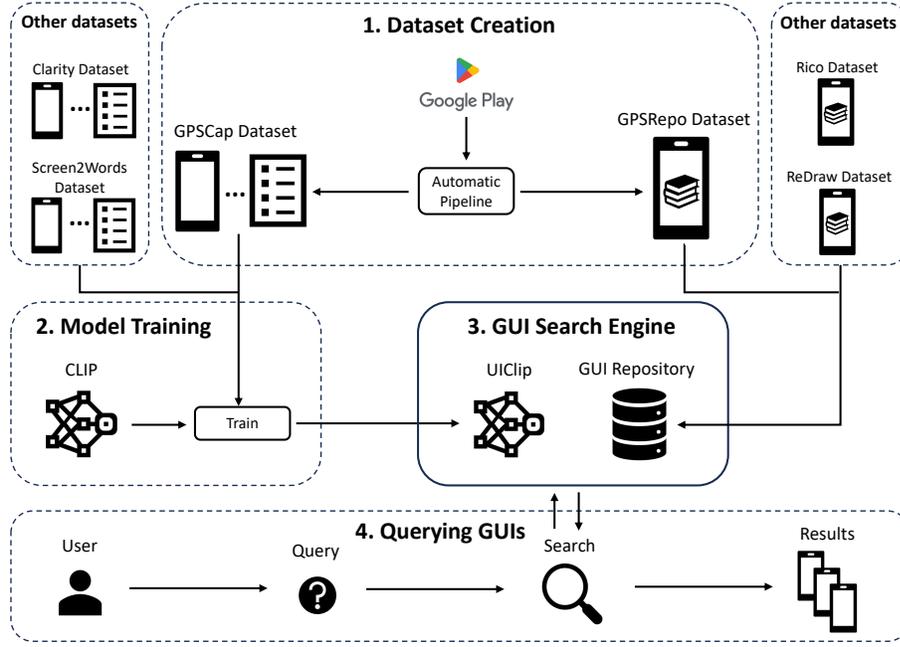


Fig. 1. Overview of our approach: including the creation of the datasets, vision-language model training, development of the GUI search engine, and the process of querying on the search engine.

Recent GUI search engines, such as GUIGLE [6] and RaWi [28], enable users to search GUI datasets using text queries. Both models use of the app metadata and GUI-related text for query matching. GUIGLE adopts basic keyword matching, while RaWi calculates semantic similarity through a BERT model. However, these purely text-based approaches solely use textual information available with the GUI, such as text blocks or button labels, neglecting key visual information such as images, layout, or background. Furthermore, the GUI datasets used by these search engines (namely Rico [15] and ReDraw [43]) are created by automatically exploring app screenshots at runtime. However, the access to certain pages may require app authorization or initial configurations, which are often lacking in automatic exploration. Thus, screenshots scrawled at runtime may fail to comprehensively encapsulate the essential features of the app.

We propose a novel GUI search engine based on vision-language machine learning models and the app introduction images from Google Play as shown on Figure 1. Recent Vision-Language Models (VLMs), such as CLIP [49], BLIP [34], and BLIP-2 [33], are trained in large-scale image-caption data by contrastive learning. These models have the ability to transform images and text into a multimodal embedding, ensuring that semantically similar images and texts are mapped closely in the embedding space. By computing the similarity between the images and the textual query, these models can be used in text-to-image retrieval tasks. For the effective application of the VLM in the app GUI domain, a large screenshot-caption dataset is a prerequisite. Currently available screenshot-caption datasets (Screen2Words [59]

and Clarity [45]) are inadequate for training a VLM as our evaluation shows in Section 5. Therefore, we have created a new large dataset.

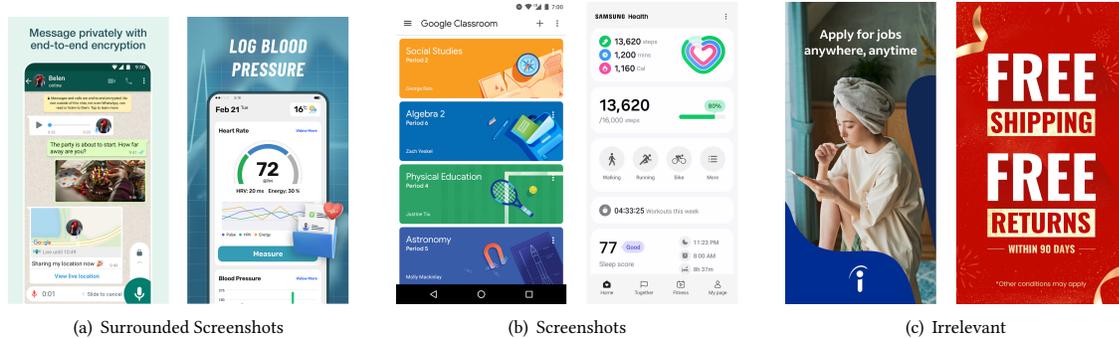


Fig. 2. Examples of app introduction images on Google Play.

Google Play is one of the most extensive mobile app stores. It offers thousands of apps from diverse domains, which makes it an abundant source of inspiration for requirements elicitation and app design [22]. The app introduction images on Google Play are a gold mine for GUI retrieval, as they are carefully selected by app developers to represent the important features of the apps. Figure 2 shows examples of these images. A large portion of the app introduction images can be categorized as *surrounded screenshots*, each of them contains a screenshot as well as a caption that succinctly describe the screenshot, as shown on Figure 2 (a). In order to extract the screenshots and respective captions, we developed an automated pipeline, optimized for these specific images. The pipeline includes image classification (to classify image into three categories as shown in Figure 2), screen cropping (for cropping the screenshots area from the *surrounded screenshots*), and caption recognition (to extract captions from the *surrounded screenshots*, like the “LOG BLOOD PRESSURE” in Figure 2 (a)). We demonstrate the efficacy of our pipeline in Section 2. By applying the pipeline on app introduction images from approximately 117k apps, we created a comprehensive dataset comprising 303k screenshots, of which 135k have a caption. We refer to the collection of 303k screenshots as the GPSRepo (Google Play Screenshot Repository), and the subset containing captions as GPSCap (Google Play Screenshot Caption). By leveraging the GPSCap, in conjunction with Screen2Words and Clarity datasets, we fine tuned the CLIP model to create a VLM specific to the GUI domain. We call the new model UIClip. Evaluation results demonstrated a promising performance of our UIClip model for text-to-image retrieval. It superseded both the text-only approach and the CLIP model. Drawing upon the UIClip model and the GUIRepo dataset, we devised our GUI search engine, Guing. The engine accepts textual queries as input to retrieve relevant GUI images from our GPSRepo dataset as well as the Rico and ReDraw datasets. Our engine can easily be extended with additional screen datasets. Our empirical evaluations show that the search engine deliver highly relevant GUIs. Additional experiments suggest that the UIClip model is also beneficial for other GUI-related tasks, such as sketch-to-GUI retrieval and GUI classification. The paper provides the following contributions, with source code and datasets publicly available <https://github.com/Jl-wei/guing>:

- A vision-language model named UIClip for a range of GUI-related tasks, including text-to-GUI retrieval, sketch-to-GUI retrieval, and GUI classification.
- A GUI search engine that can achieve high performance with textual query.

- Two large GUI datasets containing 303k screenshots, 135k of which include captions.
- An extensible pipeline for automatically extracting screenshots and captions from app introduction images.

2 DATASET CREATION

In this section, we provide a comprehensive description of our datasets:

- *Google Play Screenshot Caption (GPSCap)*: This is a screenshot-caption dataset containing 135k pairs of screenshots and captions which serve as the training data for our vision language model.
- *Google Play Screenshot Repository (GPSRepo)*: This screenshot repository is composed of 303k screenshots which form the source of our search engine.

As depicted in Figure 3, the initial step involved collecting an extensive number of app introduction images from Google Play. Subsequently, we developed an image classifier to categorize the images into *screenshots*, *surrounded screenshots* and *irrelevant*. The area of *surrounded screenshots* containing screenshots were precisely cropped from the surrounding imagery through the application of object detection. Additionally, the captions from the *surrounded screenshots* were extracted with the aid of OCR (Optical Character Recognition). The images classified as *screenshots* and the screenshots cropped from *surrounded screenshots* constitute the GPSRepo, while the caption and screenshot pairs extracted from *surrounded screenshots* constitute the GPSCap. In the subsequent sections, we will describe the construction details of our datasets.

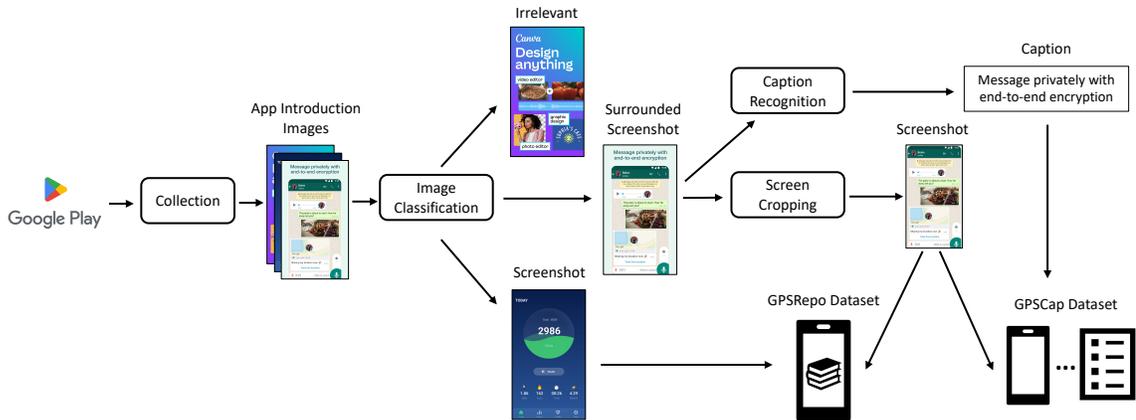


Fig. 3. Overview of dataset creation pipeline

2.1 Image Collection

Given the ID of an app, we can easily download the app introduction images via Google Play Scraper [46]. In order to obtain the app IDs, we initially gathered top ranked apps in each category from AppBrain [1], creating a seed list. An app on Google Play often suggests similar apps, and a developer associated with the app may also offer additional apps. This scenario can be conceptually treated as a graph, where apps represent nodes and the relationships (such as app similarity and common developer) represent edges. Starting from the seed list, after performing breadth-first searching this graph, we collected a total of 117,283 apps (gathered on October 30, 2023). All apps were collected from Google Play in the United States in English language.

The stylistic differences between the introduction images of games and non-game apps render the former unsuitable as references for general app design. Consequently, we excluded games from our analysis, leaving a dataset of 85,631 non-game apps. These include a total of 874,130 images. However, a portion of these images were specifically demonstrated for landscape, Wear OS or Android TV [23]. To exclude them from our dataset, we applied a filtering process based on aspect ratio, leading to 553,755 precisely filtered images. Additionally, a subset of images could be duplicates despite bearing different filenames. In order to eliminate duplications, we implemented a process where we measured the SHA-1 hash of each individual image, retaining only one image with the same hash. This process yielded a final count of 398,511 unique app introduction images.

2.2 Image Classification

The images obtained from the prior phase are classified into three categories for subsequent analysis, as depicted in Figure 3. A subset of these images encompasses the complete display of an app user interface: those are *screenshots*. For some images, only a segment of the image represents a captured screenshot, and are thus categorized as *surrounded screenshots*. The remaining images do not encapsulate screenshot or only contain partial or slanted screenshot are accordingly classified as *irrelevant*. Manual classification on around 400k images is laborious and time-consuming. Thus, we trained an image classifier to classify automatically these images into three categories: *screenshots*, *surrounded screenshot* and *irrelevant*.

2.2.1 Classification Model. The Vision Transformer (ViT) [17] was employed as the classifier for our dataset, based on its performance metrics, which will be discussed in subsequent sections. The ViT architecture essentially constitutes the encoder of transformer as described by Vaswani et al. [58]. It has demonstrated superior performance in image classification tasks. Within the realm of natural language processing, the dominant approach involves pre-training the model on a large text corpus followed by fine-tuning on a more targeted dataset [16]. Similarly, the ViT benefits from a pre-training phase using a comprehensive collection of images, thereby allowing the model to capture innate image features. Consequently, when deployed on novel tasks, the ViT requires a reduced amount of task-specific training due to the knowledge acquired during this pre-training process.

2.2.2 Training Data Creation. To train and fine tune the classification model, a labeled dataset is required. For this, we randomly sampled 3,615 images from the filtered images collected on Section 2.1. To annotate these images, we created an annotation guide to clarify the definition of *screenshots*, *surrounded screenshot* and *irrelevant*. Two of the authors then annotated these images independently, the conflicting labels were reviewed through discussion and consensus.

2.2.3 Classification Performance. The performance of the classifier was evaluated using the app introduction images from the dataset annotated in the previous step, which was split in an 80:20 ratio for training and testing respectively. The classification performance was measured by *precision*, *recall* and *F1*. We trained the classification model using mini-batch gradient descent, with a batch size of 64 and AdamW optimizer with an initial learning rate of $2e^{-5}$. It was trained on 5 epochs on a machine with a NVIDIA Tesla T4 GPU with 16 GB VRAM. In order to enhance the robustness of our results, we performed ten-fold cross-validation by repetitively splitting the dataset into distinct training and testing subsets ten times in a randomized manner. We then computed the mean performance across these runs. The experiment results presented in Table 1 show the great performance of our classifier, with an average F1 score 0.964.

2.2.4 Applying Classifier. The evaluation demonstrated the robust performance of our classifier. Subsequent to that evaluation, we trained a ViT model on all of the 3,615 images. This model was then employed to classify all of the filtered

Table 1. Image classification accuracy on evaluation set

	Precision	Recall	F1
Surrounded Screenshot	0.974	0.978	0.976
Screenshot	0.967	0.970	0.968
Irrelevant	0.937	0.926	0.931
Weighted Average	0.964	0.964	0.964

app introduction images. To augment the reliability and integrity of our GUI repository, we calibrated the classifier’s threshold at a high confidence level of 0.9, which means that only images that attain a classification probability exceeding the 0.9 are included in the respective categories. Our empirical investigations reveal that, at the 0.9 threshold, the precision scores for the categories *screenshots* and *surrounded screenshots* surpass 0.99. Following the classification process, our repository had 111,847 images classified as *screenshots* and 191,993 as *surrounded screenshots*.

2.3 Screen Cropping

The *surrounded screenshots* are the images that contain screenshots with additional visual frames. To automatically crop the screenshot areas from these *surrounded screenshots*, we trained an object detection model. This model is capable of localizing the screenshot within the larger image and subsequently generating a bounding box. A bounding box is defined as a rectangular delineation that encases the area of interest: in this case, the screenshot. Given the bounding boxes inside the *surrounded screenshots*, one can easily crop the screenshot area using an image processing library.

2.3.1 Object Detection Model. We use DETR (DEtection TRansformer) [8] as our screenshot detector. DETR is an end-to-end object detector, composed of a CNN (convolutional neural network) backbone and a encoder-decoder transformer [58]. Atop the transformer decoder, dual output heads are appended: a linear layer dedicated to categorize class labels, and a multi-layer perceptron (MLP) tasked with the generation of bounding boxes for the object location. The model uses so-called object queries to detect objects in an image. Each object query is designed to search specifically for one particular object of interest in the given image. In our work, the number of queries is set to 1, as we are interested in the largest screenshot within a *surrounded screenshot*.

2.3.2 Training Data Creation. To fine tune the detection model for screenshot cropping, we created a dataset. From the dataset created in Section 2.2.2, we collected all the *surrounded screenshots* and labeled the screenshot areas with a bounding box. The annotation was performed with Prodigy [20]. For each image, an author drew one bounding boxes that cover the screenshot, the results were reviewed by another author. The majority of the *surrounded screenshots* feature only a single screenshot. In instances where a single image contains multiple screenshots, annotation was selectively applied only to the most prominent (largest) screenshot. This resulted in a dataset comprising 1,768 annotated images, each including a delineated bounding box that specifies the screenshot area.

2.3.3 Detection Performance. Intersection over Union (IoU) is a standard metric in object detection tasks. It evaluates the extent of overlap between the predicted bounding box and the ground truth bounding box. If $IoU = 0$, it means that there is no overlap between the boxes, while $IoU = 1$ means that the overlap is perfect.

$$IoU = \frac{AreaOfOverlap}{AreaOfUnion}$$

Table 2. Screen localization accuracy

Precision	IoU=0.50:0.95	0.919
	IoU=0.50	0.981
	IoU=0.75	0.971
Recall	IoU=0.50:0.95	0.947

In this study, we followed the evaluation protocol established by the Common Objects in Context (COCO) dataset [37], which is a common benchmark in the field of object detection. As described in the protocol, we quantitatively evaluated the performance of our object detection by calculating precision and recall metrics at various IoU thresholds. Specifically, we computed the precision values at three distinct IoU levels: 0.50, 0.75, and an average over the range from 0.50 to 0.95. Similarly, recall was determined over the IoU range of 0.50 to 0.95.

We applied 80% of the data for training and the remaining 20% for evaluation. The detector was trained on the training set using mini-batch gradient descent, with a batch size of 16 and AdamW optimizer with an initial learning rate of $1e^{-5}$ and the training epochs is 10. Table 2 shows the results of ten-fold cross validation. The IoU at 0.75 exhibits a value of 0.971, indicating a high level of accuracy.

2.3.4 Applying the Detector. We trained a DETR model with all of the 1,768 images, and applied this model to the 191,993 images that were categorized as *surrounded screenshot* in Section 2.2.4. DETR predicted a bounding box that localize the screenshot for each of the images. We then cropped these images according to the bounding box. The screenshots cropped from the *surrounded screenshot* as well as the app introduction images which were previously classified as *screenshots* represent the GPSRepo dataset.

2.4 Caption Recognition

The majority of the images, classified as *surrounded screenshots*, include captions that provide a succinct descriptions of the screenshots. Typically, these captions are positioned either above or below the image. In order to accurately extract this textual information from the images, we employed Optical Character Recognition (OCR) technology. This technology is adept at detecting and recognizing text displayed within image files.

2.4.1 Applying OCR. In this study, we employed PaddleOCR [31], which is an industrially robust OCR model renowned for its accuracy and swift performance. Given that the textual content within the images under analysis is markedly legible, exhibiting a clarity not typically associated with hand-written text, the application of PaddleOCR is deemed adequately proficient for this task. PaddleOCR processes an input image and outputs a collection of identified elements. Each element is comprised of a bounding box delineating the text area, the recognized text within, and an associated confidence level reflecting the model’s recognition certainty.

We applied PaddleOCR on all of the 191,993 *surrounded screenshots*. For each *surrounded screenshot*, text bounding boxes that exhibit spatial overlap with the screenshot’s bounding box, as delineated in Section 2.3.4, were excluded. The remaining text bounding boxes were subsequently aggregated to form a coherent caption for the corresponding image.

2.4.2 Post-processing. Some of the *surrounded screenshots* may not contain captions. These images were removed from our analysis. The captions extracted are not exclusively in English; they encompass a variety of languages including French, German, and Arabic, among others. As we focus on English search in this work, captions in languages other

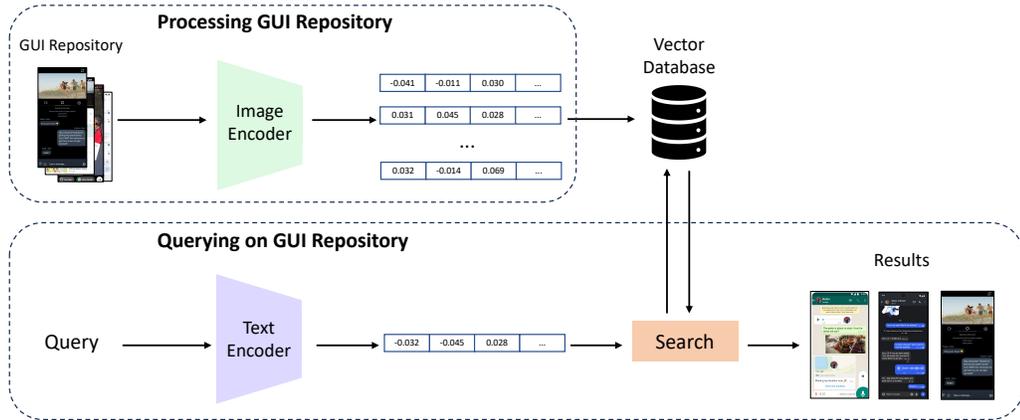


Fig. 4. Overview of our GUI search engine (GUing) for text-to-GUI retrieval.

than English were excluded from the analysis. To facilitate the identification of the language to which each caption pertains, we employed Lingua [55], an effective language detection tool. Furthermore, to correct any spelling error within the captions due to potential OCR errors, we utilized Autocorrect [54], a Python-based spelling corrector.

It is noteworthy that in some instances, multiple *surrounded screenshots* of an app may feature identical captions. This repetition suggests that the caption is either overly generic or simply a logo of the app, and thus may not adequately convey the screenshot content. To address this redundancy, we conducted a filtering process whereby, for each app, only the first *surrounded screenshot* from the set containing identical captions was retained. We observed that the first *surrounded screenshot* displayed by an app on Google Play typically showcases the most representative app feature. After the post-processing stage, we obtained 135,357 screenshot-caption pairs, collectively referred to as the GPSCap dataset.

3 GUI SEARCH ENGINE

We introduce GUing, an advanced search engine based on vision-language model that leverages textual queries to retrieve relevant screenshots from a large GUI repository which consists of GPSRepo, Rico [15] and ReDraw [43]. The architecture of GUing is illustrated in Figure 4.

GUing is developed based on the UIClip model. This search engine utilizes the image encoder and text encoder modules of UIClip to embed screenshots and textual queries within a unified latent space. Consequently, it becomes possible to facilitate an efficient search process for the screenshots by calculating the cosine similarity between the text query embedding and screenshot embedding. In the following, we elaborate on the details of our GUI search engine.

3.1 Constructing UIClip Model

UIClip, our vision-language model, serves as the fundamental component of our GUI search engine, integrating the image and text modalities. Building upon the foundation of CLIP [49], UIClip offers enhanced capabilities in multimodal representation learning within GUI domain. CLIP (Contrastive Language-Image Pre-training) model based on vision language contrastive learning on large-scale image-caption data is a significant foundation model in multimodal

representation learning. In this work, we construct a GUI-specific CLIP model by training a vision-language model on a large-scale dataset of screenshot-caption pairs.

3.1.1 Training Data and Pre-processing. For the training of UIClip, we combined three datasets to learn the visual representation of mobile screenshots:

- Screen2Words dataset [59] is a subset of Rico. It contains 112,085 language summarization across 22,417 unique screenshots. The screenshot summaries are written by professional annotators. For each screenshot, five summaries are created by five different annotators.
- Clarity dataset [45] consists of 45,998 descriptions for 10,204 screenshots from popular Android apps. Crowd workers summarize each screenshot in different granularity, with one high-level caption to up to four low-level detailed descriptions.
- Google Play Screenshot Caption (GPSCap) dataset. Here the textual description is created by the respective developers and app vendors. In the Section 2, we created a dataset with 135k screenshot-caption pairs by mining Google Play. The screenshots are cropped from app introduction images, while the corresponding captions describe briefly the screenshots.

The captions were firstly tokenised with the CLIP tokeniser [47]. Captions longer than 77 tokens were truncated to 77, captions shorter were padded to 77. In terms of image processing for the screenshots, the resolution was resized to 224×224 . Furthermore, the pixel value of the images was proportionately rescaled from a range of 0-255 to 0-1. This new value was subsequently normalized using the parameters (means and standard deviations) featured in the official CLIP implementation documentation [47].

3.1.2 Model Architecture. UIClip consists of two components, an image encoder and a text encoder. The image encoder, which adopts a Vision Transformer (ViT) model [17], ingests image data and generates its corresponding image embedding. The text encoder is a transformer encoder [58] which accepts text as input to produce a associated text embedding. Employing contrastive learning techniques, the image and text embeddings are mapped into a multi-modal embedding space. In this shared space, images and texts presenting semantic similarity are mapped close to each other.

Contrastive learning on image-text pairs constructs a bridge between visual perception and natural language. Given a batch of N (image, text) pairs, the training objective is to identify the authentic pairings from $N \times N$ potential (image, text) combinations within the batch [49]. To achieve this objective, the model formulates a multi-modal embedding space, concurrently training an image encoder and text encoder to maximize the cosine similarity of the image and text embeddings derived from the N authentic pairs in the batch. Simultaneously, it minimizes the cosine similarity of the embeddings from the remaining $N^2 - N$ incorrect pairings.

3.1.3 Training Details. In order to adapt the CLIP model to UIClip, rather than training the CLIP model from scratch, where the weights of image encoder and text encoder are randomly initialized, our training is based on the checkpoint of "openai/clip-vit-base-patch32" [48]. This checkpoint has been pre-trained on a 400 million image-text pairs. This method allows us to harness the benefits inherent to pre-trained models and, as a result, substantially diminishes the time required for training.

The model was trained on all of the collected screenshot-caption pairs for 5 epochs. We applied mini-batch gradient descent with a batch size of 128 and AdamW optimizer with an initial learning rate of $5e^{-5}$.

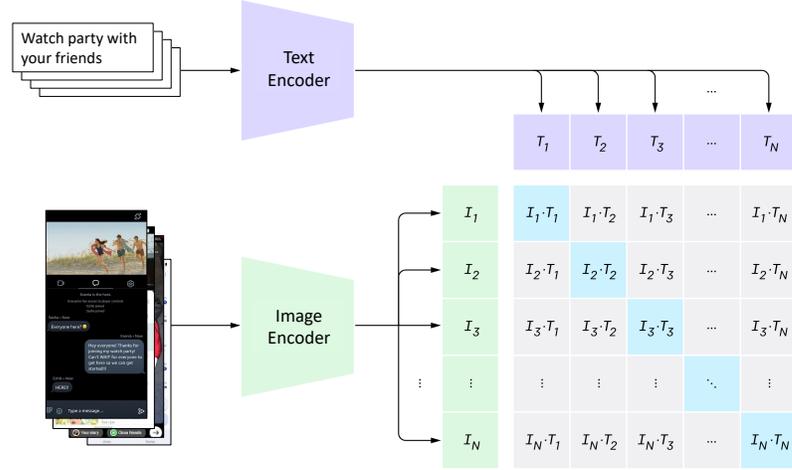


Fig. 5. Overview of contrastive learning to train our vision-language model UIClip (adapted from [49])

3.2 Searching the GUI Repository

3.2.1 Processing the GUI Repository. The GUI repository is composed of 383k screenshots from GPSRepo (303k), Rico (66k) and ReDraw (14k). These screenshots are on image format, thus cannot be used directly for querying. Therefore, we embed all the screenshots in the GUI repository and save them in a vector database.

The screenshots are firstly resized, rescaled and normalized with the same parameters as those in Section 3.1.1. These processed screenshots subsequently serve as the input for the image encoder of UIClip. The image encoder accepts an image as an input, subsequently outputting a 512-dimensional embedding. The entire GUI repository, comprising 383k screenshots, is processed via the image encoder, consequently generating 383k 512-dimensional embeddings. The mapping between the screenshot id and the resulted embeddings are saved in a vector database for further retrieval.

3.2.2 Querying on the GUI Repository. In the querying phase, the textual query is encoded by UIClip’s text encoder, resulting in a 512-dimensional embedding. Subsequently, this query embedding is utilized to identify the top-k most similar screenshots embeddings from the vector database. In our research, we employ cosine similarity as the measure of similarity. Consequently, the query results consist of the top-k screenshots that bear the most resemblance.

3.2.3 Accelerating the Searching. Given the large scale of the GUI repository, employing a brute force approach to compare the query embedding with all screenshot embeddings is inefficient. To tackle this problem, we applied approximate nearest neighbor search [3] to accelerate the querying as described following, as detailed below.

We established n Voronoi cells within the embedding space, with each screenshot embedding falling into one of these designated cells. During a search, the query embedding is initially compared with the centroid embeddings of each cell. This method results in a considerable reduction of necessary comparisons when contrasted with the total number of screenshot embeddings. The k cells with centroids closest to the query embedding are identified, whereupon the screenshot embeddings within these cells are thereafter compared with the query embedding in order to pinpoint similar screenshot embeddings. Our implementation is based on Faiss [18], a library for efficient similarity search. We set the number of cells n to 3000, and the number of cells that are visited to perform a search k to 1000.

Table 3. Evaluation dataset statistics

Split	Dataset	#Apps	#Screenshots	#Captions
Training	GPSCap	32720	133477	133477
	Screen2Words	5684	20379	101895
	Clarity	3346	8204	36964
Validation	GPSCap	246	1000	1000
	Screen2Words	281	1000	5000
	Clarity	413	1000	4475
Test	GPSCap	245	1000	1000
	Screen2Words	286	1000	5000
	Clarity	416	1000	4559

4 EVALUATION DESIGN

We present the design of the empirical evaluation for our GUI search engine. Specifically, we investigate the following research questions:

RQ₁: How does GUing perform in text-to-GUI retrieval tasks compared to state of the art approaches?

As the CLIP model [49] was trained on 400 millions of image-caption pairs from internet, it can also be applied on text-to-GUI retrieval task. An alternative method involves leveraging solely the textual information - the text displayed on the GUI images, computing its semantic similarity with the input textual query, and subsequently ranking this similarity for effective GUI retrieval. Thus, the question is whether the performance our proposed UIClip model surpasses the CLIP model and text-only approaches.

RQ₂: How relevant are the retrieved screenshots for queries representing app features? As a search engine, GUing processes textual queries as input and searches for relevant images within the GUI repository. It is thus important to evaluate how effectiveness GUing in delivering useful results that actually align with user expectations.

To answer these research questions, we designed two experiments: One (Expt₁) is based on benchmarking three screenshot-caption datasets, and the other (Expt₂) is based on a manual evaluation of the search engines.

4.1 Expt₁: Evaluation of Models Performance on Three Datasets

In the first experiment, the text-to-GUI retrieval performance of our UIClip model is evaluated together with baseline models on three distinct datasets. During the evaluation, every screenshot and its associated captions in test dataset are embedded using the models under evaluation. For every single caption, all screenshots within the test dataset are ranked by the cosine similarity of the caption embedding and screenshot embedding. A comparative analysis of the performance of UIClip against various baselines provide answers to RQ₁.

4.1.1 Experimental Data. From each of the screenshot datasets described in Section 3.1.1 (GPSCap, Screen2Words, and Clarity), we randomly selected 1000 screenshots for validation, 1000 screenshots for test, and the rest for model training, as shown in Table 3. Similar to the approach suggested by Wang et al. [59], we split the data to *not* share the screenshots from the same app across different splits. That is, all the apps and screenshots in the test and validation set were completely unseen during the training. This arrangement allows us to assess how well our model generalizes to previously unseen screenshots from unseen apps during the test phase.

4.1.2 UIClip Training Details. We trained the "openai/clip-vit-base-patch32" [48] checkpoint on the training set for a total of 5 epochs. During the training, a mini-batch gradient descent approach was employed with a batch size of 128. To optimize the learning process, the AdamW optimizer was utilized with an initial learning rate set to $5e^{-5}$.

4.1.3 Baselines. As highlighted by Bernal-Cardenas et al. [6] and by Kolthoff et al. [28], existing approaches perform their queries based on GUI metadata, which is absent in screenshots collected by processing app introduction images from Google Play. Consequently, we implemented a text-only retrieval approach by ourselves. In addition, we compared our model with (a) the original CLIP model without further tuning and (b) with the CLIP model fine-tuned *without* our GPSCap data. This results in three baselines:

- **OCR + BGE (Text Only).** As discussed in the Section 2.4.1, PaddleOCR [31] is an industrial grade OCR library. BGE [63] is a state-of-the-art text embedding model, trained in three steps: including pre-training with plain texts, contrastive learning on text pair dataset, and task-specific fine-tuning. First, the screenshots are embedded by using OCR + BGE in three steps: 1) the texts displayed on the GUI image are extracted with PaddleOCR; 2) the extracted text is then concatenated into a sentence, separated by semicolons; 3) the sentence is then embedded by BGE. Then, the captions are directly embedded using the BGE model.
- **CLIP.** The screenshots are embedded by the image encoder of CLIP and the captions are embedded by the text encoder of CLIP. The "openai/clip-vit-base-patch32" checkpoint is used for this evaluation.
- **UIClip-CS.** The UIClip-CS is a fine-tuned model from CLIP model using the Clarity and Screen2Words dataset. The training steps is exactly same as the training of UIClip, except that our GPSCap dataset is excluded during training process. The image encoder of UIClip-CS embeds the screenshots, while its text encoder embeds the captions.

4.1.4 Evaluation Metric. The recall@k is commonly used for evaluating text-to-image retrieval performance [33, 34, 49, 53]. Recall@k is defined as the percentage of captions whose corresponding screenshots fall into its top-k most similar screenshots. For our evaluation, we used recall@1, recall@3, recall@5, recall@10, recall@50 and recall@100.

4.2 Expt₂: Manual Evaluation of Search Engines

A GUI search engine accepts textual queries as input and identifies the corresponding images within the GUI repository as output. Expt₁ focuses on evaluating the performance of retrieval models on the test dataset. Thereby, the main limitation is the use of the screenshot-caption pairs as ground truth, that is the assumption of a direct one-to-one correspondence between a single caption (search query) and a relevant screenshot. However, such assumption contradicts real-world scenarios where one query can pertain to multiple screenshots.

To address this limitation, experiment 2 aims to evaluate different search engines, which are composed of various retrieval models and GUI repositories, through a manual evaluation. Specifically, we measure the proportion and rank of images retrieved by the GUI search engines that are deemed relevant to the query entered by the user. Through comparative analysis of the usefulness between GUIing and baseline approaches, we are able to address RQ₂.

4.2.1 Baselines. To answer RQ₂, we undertook a comparative analysis of our search engine with two baselines, RaWi and GUIing-CS.

- **RaWi.** RaWi [28] is a data-driven GUI prototyping approach that retrieves screenshots for reuse from Rico dataset based on natural language searches. The search engine of RaWi is a BERT-based Learning-To-Rank (LTR) model that is trained on GUI relevance data collected from a crowd-sourcing task. As Rico dataset provides

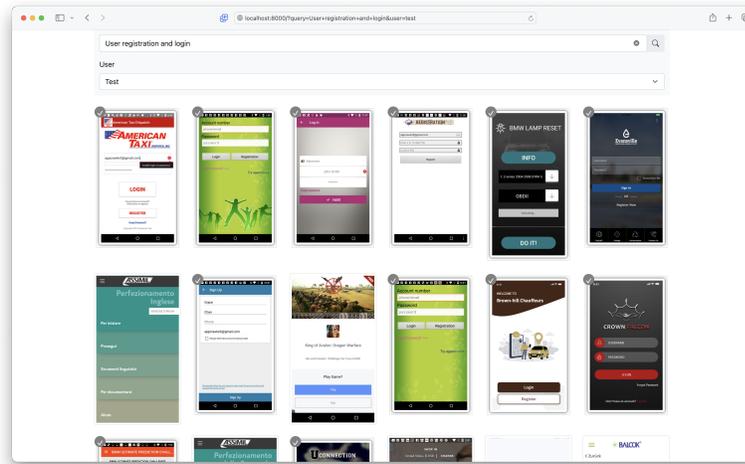


Fig. 6. Interface of the evaluation tool. It displays repeated GUIs due to their retrieval by more than one search engine.

the metadata of the screenshots, RaWi creates the textual representation of the screenshots by assimilating elements such as the activity names, UI component text, resource ids, and icon ids. This integrated textual representation is then compared with the user text query, forming the basis for GUI ranking.

- *GUing-CS*. The architecture of *GUing-CS* replicates precisely that of *GUing*. However, the retrieval model of *GUing-CS* is *UIClip-CS*, which only uses the *Clarity* and *Screen2Words* datasets for model training. That is we explicitly remove our *GPSCap* dataset from the training process. Furthermore, the GUI repository for *GUing-CS* is comprised exclusively of the *Rico* and *ReDraw* datasets, with no utilization of the *GPSRepo* dataset.

4.2.2 Evaluation Method. To compare between the various search engines, we developed a tool as shown on Figure 6. The tool accepts textual queries as input and produces the top-10 screenshots from the three search engines: *GUing*, *GUing-CS*, and *RaWi*. The tool mixes and randomly shuffles the resulting 30 screenshots to prevent the possible bias of the evaluators. That is the origin of each search results is unclear and cannot be guessed by the evaluator.

During the assessment, evaluators are asked to perform search with the provided queries, select all screenshots they consider relevant to the query and confirm their selections via the submission button. A screenshot is deemed relevant to a query if it either displays a UI implementation of the feature delineated in the query or if it can serve as a reference or inspiration for the development of the feature.

The evaluation queries were derived from the articles of *AttractGroup* and *BuildFire*, which are two companies specialized in mobile app development services and have substantial knowledge in this field. *AttractGroup* enumerated 138 mobile app features to be considered during mobile app development, it consists of 45 general features and 93 domain specific features spanning across 11 domain (such as E-learning, mHealth, FinTech, etc.) [57]. *BuildFire* proposed 50 innovative app ideas for 2024, forming a good source of inspiration for prospective entrepreneurs [27]. From these two resources, a total of 188 mobile app features were collected. These were then divided into two sets: one set consisting of the 45 general and 50 innovative features, and another set comprising 93 domain-specific features.

Table 4. Text-to-GUI retrieval performance on test set (Experiment 1)

Dataset	Model	Recall@1	Recall@3	Recall@5	Recall@10	Recall@50	Recall@100
	(Chance)	0.001	0.003	0.005	0.010	0.050	0.100
GPSCap	OCR + BGE	0.164	0.291	0.355	0.441	0.624	0.705
	CLIP	0.127	0.242	0.299	0.392	0.627	0.723
	UIClip-CS	0.066	0.132	0.164	0.244	0.479	0.603
	UIClip	0.377	0.526	0.593	0.687	0.857	0.908
Screen2Words	OCR + BGE	0.130	0.229	0.287	0.361	0.559	0.650
	CLIP	0.097	0.171	0.220	0.298	0.508	0.615
	UIClip-CS	0.195	0.349	0.429	0.532	0.800	0.887
	UIClip	0.216	0.371	0.449	0.566	0.825	0.903
Clarity	OCR + BGE	0.152	0.242	0.281	0.338	0.500	0.585
	CLIP	0.080	0.138	0.166	0.219	0.373	0.464
	UIClip-CS	0.073	0.146	0.198	0.278	0.515	0.648
	UIClip	0.075	0.150	0.207	0.287	0.533	0.667

A total of four evaluators took part in this assessment, each holding a master’s degree in computer science and possessing a minimum of five years’ experience in software development. Two of them performed the evaluation with the first set of queries, the other two used the second set of queries, which ensures that each query has been evaluated by two evaluators. The evaluators performed the searching, clicking and submitting for their respective queries, as described above.

4.2.3 Evaluation Metric. To measure the performance of the GUI search engine, we applied well-known information retrieval metrics (similar to Kolthoff et al. [28]), including Precision@k, Mean Reciprocal Rank, and HITS@k.

- *Precision at k (P@k)* quantifies the proportion of retrieved screenshots considered relevant, denoted as $|R_k|$, with respect to the top-k retrieved GUIs, denoted as k .

$$P@k = \frac{|R_k|}{k} \quad (1)$$

- *Mean Reciprocal Rank (MRR)* is a statistic measure for evaluating any process that produces an ordered list of possible responses to a sample of queries. $rank_i$ signifies the rank of the highest ranked relevant screenshot of the i -th query and Q stands for the total number of queries.

$$MRR = \frac{1}{Q} \sum_{i=1}^{|Q|} \frac{1}{rank_i} \quad (2)$$

- *HITS@k* is a binary measure that validates the presence of at least one relevant GUI among the top-k ranked screenshots. It has the value 1 if there is at least one relevant screenshot within the top-k, otherwise, it is 0.

$$HITS@k = \begin{cases} 1 & |R_k| > 0 \\ 0 & |R_k| = 0 \end{cases} \quad (3)$$

5 EVALUATION RESULTS

5.1 RQ₁ Performance of UIClip

Table 4 shows the evaluation results for the text-to-GUI retrieval task using UIClip (last row) and the baseline models on three distinct datasets.

Compared to the retrieval by chance, which has a recall@10 of 0.01, the CLIP model proves to be more effective in text-to-GUI retrieval, with a recall@10 exceeding 0.2 across the three datasets. This implies that, given a caption, the CLIP model has a 20 percent probability that the corresponding screenshot lies within the top 10 retrieved screenshots. Nevertheless, the UIClip model consistently outperforms the CLIP model in almost all metrics spanning the three datasets. This indicates that our fine-tuned CLIP model, UIClip, is more adept at bridging the semantic gap between the caption and the screenshot.

UIClip-CS were trained merely on the training set of Screen2Words and Clarity. Notably, UIClip delivers better performance than UIClip-CS on the Screen2Words and Clarity datasets, registering a performance improvement ranging between 0.1 to 0.3 across various metrics. This performance distinction is further accentuated in the GPSCap dataset, where UIClip-CS exhibits an inferior performance, even in comparison to the CLIP model. This anomalous outcome is hypothesized to derive from caption style variances. Specifically, the captions associated with the Screen2Words and Clarity datasets are characterized by their significant length and detail, which may possibly impede UIClip-CS’s capacity to generalize effectively on the GPSCap dataset. On a holistic level, the results of UIClip suggest that fine-tuning on our GPSCap dataset contributes substantially to the performance of CLIP models for the text-to-GUI retrieval task. This dataset enhances the performance of the CLIP model not only on the GPSCap dataset, but it also leads to performance gains on other screenshot-caption datasets.

As for the text-only approach, when analysing the GPSCap and Screen2Words datasets, it becomes evident that UIClip surpasses the performance of OCR+BGE. Interestingly, when the comparison is carried out on the Clarity datasets, UIClip exhibits a deterioration in comparison to text-only approach within the given metrics recall@1, Recall@3, recall@5 and recall@10. This discrepancy can presumably be attributed to the variability in caption styles. The majority of captions in the GPSCap dataset are simple features, such as “Track your health trends” or “Custom color themes”. And the captions found within the Screen2Words dataset generally encompass high-level descriptions of the screenshot, examples being “screen showing settings options” and a “display of news articles in a news app”. In contrast, the Clarity dataset captions offer substantially more details about the screenshots, like “at the bottom left there is a skip button for users to skip the introduction”, and “in the middle of the screen a popup is displayed with a label called set date”. These detailed descriptions can easily be matched with the text displayed on the screenshots, thereby improving the performance of text-only approaches. Nevertheless, UIClip continues to surpass OCR+BGE in both recall@50 and recall@100.

Overall, our findings underscore that UIClip demonstrates a superior performance than the baseline models in text-to-GUI retrieval, particularly when the goal is to search descriptions of the overall functionality and features implemented in the screens rather than verbose, documentation-like descriptions of the UI elements in the screens.

5.2 RQ₂ Relevance Search Results

Table 5 shows the results of the manual evaluation for GUing and the baseline search engines. The results show that GUing consistently outperforms RaWi across all evaluation metrics. P@1 and HIT@1 show that there is 0.372 probability that GUing will return a relevant screenshot as the first result, a considerable improvement over RaWi’s rate of about

Table 5. Relevance of the retrieved screenshots by three GUI search engines (Experiment 2 / Manual assessments)

Search Engine	MRR	P@k				HIT@k			
		P@1	P@3	P@5	P@10	HIT@1	HIT@3	HIT@5	HIT@10
RaWi	0.410	0.291	0.264	0.254	0.214	0.291	0.487	0.594	0.701
GUing-CS	0.255	0.134	0.138	0.139	0.152	0.134	0.291	0.401	0.615
GUing	0.510	0.372	0.352	0.339	0.343	0.372	0.666	0.773	0.914

0.291. The P@10 value, indicative of the proportion of relevance results of the top-10 results returned by GUing, stands at 0.343, compared with RaWi’s rate of 0.214. The HIT@10 evaluation metric shows that in more than 91% of cases, GUing can return at least one pertinent screenshot among the top-10 results, far exceeding the percentage for RaWi, which is recorded at 70%. This is particularly interesting for UI ideation tasks, as it seems reasonable for a designer to check 10 screen suggestions for identifying at least one relevant screen (similar to screening the first result page of a google search). Finally, the MRR metric reveals that the ranking of the first relevant result returned by GUing is superior to that of RaWi.

A comparative analysis of the evaluation outcomes for GUing and GUing-CS shows a significant outperformance of GUing. This can be quantified by notable gaps in the following metrics: a gap exceeding 0.2 in regards to P@k, approximately 0.3 for HIT@k, and 0.25 in terms of MRR. The architecture of the search engines GUing and GUing-CS is completely analogous, with their only difference being the vision-language models’ training set and the repository used for searches. GUing-CS does not incorporate our GPSCap and GPSRepo datasets for its training and search operations. The large gap between GUing and GUing-CS serves as further validation of the critical importance of our GPSCap and GPSRepo datasets.

We carried out additional analyses on GUing and RaWi using a few queries in which the original search engine of the screenshots was not concealed. The example of the returned screenshots of GUing and RaWi with the query "sleeping track" are shown on Figure 7 and 8. Upon examination, it became apparent that the stylistic differences between screenshots returned by two search engines were considerable. This difference lays on three aspects:

- The shape of the screenshots returned by GUing is less neat. A large portion of screenshots within the GPSRepo dataset is derived from the *surrounded screenshots*, in which the screenshot area varies in shape across diverse images. Even the shape is less neat, these screenshots can still be used for the inspiration of GUI design.
- GUing returns a wider variety of screenshots than RaWi. As shown on Figure 8, the presence of empty or settings pages among the screenshots returned by RaWi is thought to be an inherent limitation of the text-only retrieval approach, as they only offer minor contribution to sparking creativity in GUI design. Text-only retrieval approach compares the similarity of the query with text or metadata found on each screenshot. As setting pages typically contain an abundance of text, the likelihood of heightened query similarity is increased. Empty pages, on the other hand, are likely a consequence of the Rico data collection method, which involves the automatic exploration of apps. When performing automatic exploration, there is no initial data in the app, such as sleep events record.
- The modernity of the screenshots returned by GUing vastly exceeds those resulting from RaWi. This discrepancy can feasibly be attributed to the different collection periods of the underlying datasets. The GPSRepo dataset was collected recently, whereas the Rico dataset was created in 2017. The GPSRepo dataset can be easily updated by processing introduction images of new apps when added to the Google Play store.

Overall, our findings demonstrate that GUing consistently outperforms the baseline approaches in text-to-GUI retrieval task in term of relevance of the retrieved screens.

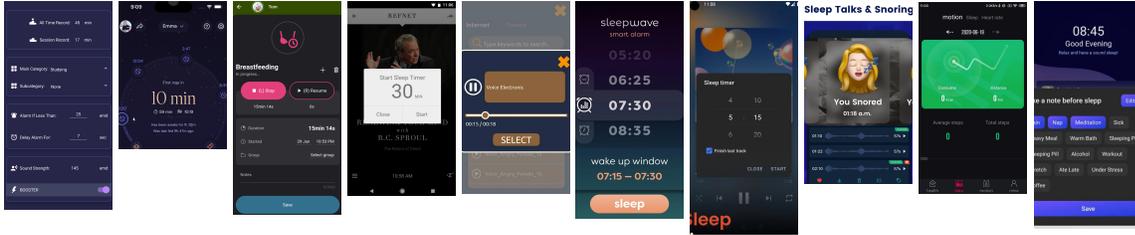


Fig. 7. Top-10 screenshots returned by GUing with query "sleep tracking"

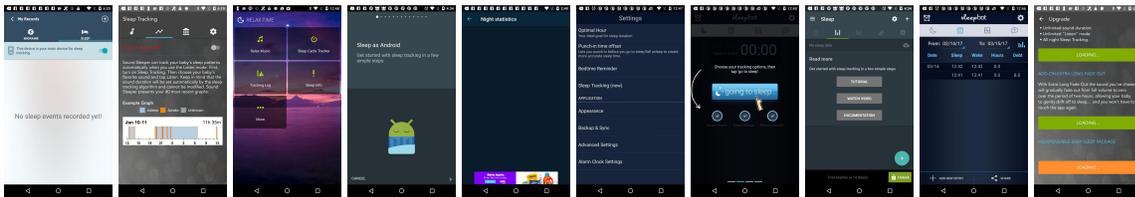


Fig. 8. Top-10 screenshots returned by RaWi with query "sleep tracking"

6 EXPLORING UICLIP AS FOUNDATION MODEL FOR GUI TASKS

In the previous section, we evaluated the performance of UIClip (as a main component of GUing) for the text-to-GUI retrieval task. However, as a tuned foundation model, UIClip can be adapted for various other tasks. In this section, we discuss the broader application of UIClip to GUI-related tasks, including GUI classification, sketch-to-GUI retrieval, GUI captioning, and GUI generation—along with a preliminary evaluation.

6.1 GUI Classification

GUI classification aims to categorize a screenshot as a whole under a specific label. It can be the fundamental task for other GUI-related tasks, like GUI captioning [30], or the association of user screenshots (e.g. submitted in support requests) to certain system requirements and components [24, 39, 51]. In this section, we explore the performance of UIClip on GUI classification under zero-shot and linear probe setting using the Enrico dataset [29].

6.1.1 Model Architecture. As performed by Radford et al. [49], we evaluated the performance of UIClip and baselines on GUI classification under zero-shot setting and linear probe setting as shown on Figure 9. Under both settings, the image encoder is frozen.

- **Zero-shot.** In the zero-shot process, all classification labels are embedded to text embedding with the text encoder. During the classification process, the screenshot image is encoded by the image encoder, and the resulting screenshot embedding is compared to all text label embeddings. The label exhibiting the highest similarity to the GUI image embedding is subsequently declared as the predicted outcome.
- **Linear Probe.** While in the linear probe setting, we added a linear layer of shape $m \times k$ to the output of UIClip image encoder. m refers to the output dimension of image encoder, while k is the number of classification labels.

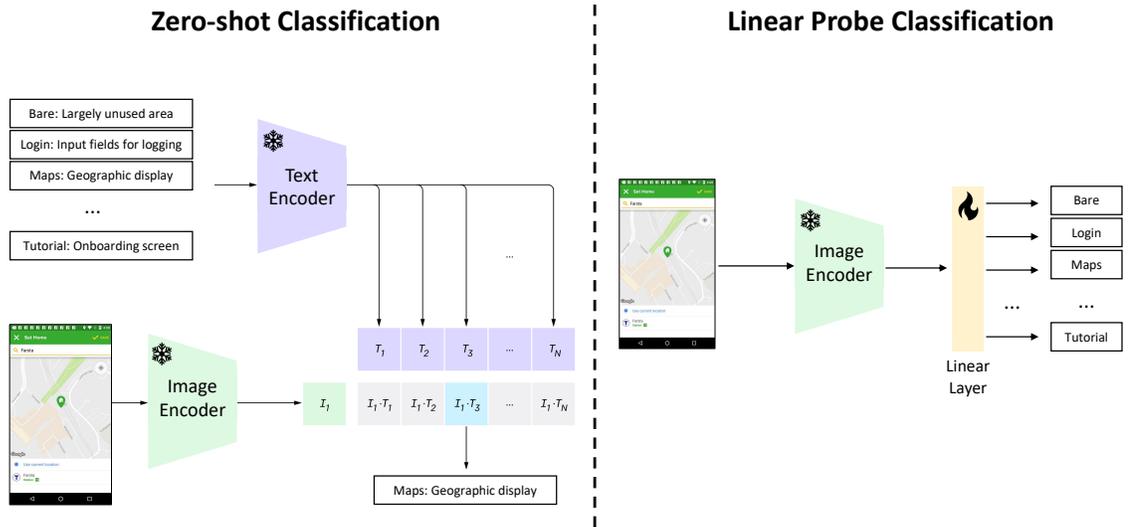


Fig. 9. Overview of zero-shot and linear probe classification on Enrico dataset with UIClip (left side is adapted from [49])

6.1.2 Dataset. We use the Enrico dataset to evaluate the performance of GUI classification. Leiva et al. [29] revised a random sample of 10k screenshots from Rico dataset to create Enrico (shorthand of Enhanced Rico): a human annotated dataset comprising 1460 screenshots and 20 design topics, including camera, chat, media player, etc.

In the linear probe setting, the Enrico dataset was divided into an 80:20 ratio to establish a training set and a test set, which were used respectively for the model training and testing. For the zero-shot setting, all available data was employed solely for testing purposes.

6.1.3 Baselines. We compared our model in GUI classification with the following baselines, performing both zero-shot classification and linear probe classification:

- *CLIP.* We also evaluated the performance of CLIP on GUI classification. The setting of zero-shot and linear probe evaluation of CLIP is identical to the setting of UIClip. Their only difference is the parameters' weights of the two models.
- *OCR + BGE.* Rather than encoding the screenshot image, this method focuses on the textual information presented on the screenshots. As described on the Section 4.1.3, the texts displayed on screenshots are extracted with PaddleOCR [31], then concatenated into sentence and embedded by BGE [63].

6.1.4 Training Details of Linear Probe. For the training and evaluation of linear probe classifier, we split the dataset into training and test set in a ratio of 80:20 in a stratified manner. During training, the image encoder is frozen, the output embedding is served as input of the final linear layer. We trained the model for 100 epochs using AdamW optimizer with an initial learning rate $1e^{-5}$, batch size 128. We conducted a ten-fold cross-validation, which involved randomly splitting the training and testing sets ten instances, and subsequently computing the average performance spanning across these runs.

6.1.5 Evaluation Metric. The evaluation of the classifiers is performed by utilizing precision, recall, and the F1 score. Given the relatively small size of the Enrico dataset, especially when compared to the 20 available classification labels,

Table 6. Classification accuracy on Enrico dataset

	Zero-shot			Linear Probe		
	Precision	Recall	F1	Precision	Recall	F1
OCR + BGE	0.210	0.094	0.081	0.038	0.186	0.059
CLIP	0.317	0.165	0.138	0.402	0.466	0.395
UIClip	0.415	0.347	0.334	0.613	0.640	0.600

we opt not to calculate the precision, recall, and F1 score for each individual label. In contrast, we calculate their weighted average.

6.1.6 Results and Discussion. As shown in Table 6, vision-language models, specifically CLIP and UIClip, clearly outperform OCR+BGE as expected. This reaffirms the assertion that mere analysis of text displayed on screenshots is insufficient for substantial GUI understanding, making the vision information crucial. Notably, UIClip achieves superior results compared to CLIP, and records an F1 score that exceeds that of CLIP by 20 percent in both settings. This highlights the advantage of our UIClip model for GUI classification tasks over the baseline models, given its enhanced knowledge about the GUI domain accumulated during the pre-training phase. Particularly, the Linear Probe led to encouraging precision and recall of over 60% (given the large number of classes in the classification task and the stright application of the model).

The efficacy of vision-language models, as measured under the linear probe setting, surpasses that determined under the zero-shot setting. This superior performance can be rationalized in terms of the additional knowledge imparted during the training phase about the Enrico dataset. However, a contrary pattern was found concerning OCR+BGE, which showed diminished performance following linear probe training. We hypothesize that this outcome is attributable to the text-only approach’s inadequate competency in extracting visual information from the screenshots, which renders the the training stage ineffective.

6.2 Sketch-to-GUI Retrieval

Rather than searching GUI by textual query, it is also possible to perform the retrieval using a UI sketch. Based on UIClip, we built a Dual UIClip ViT model for sketch-to-GUI retrieval. An evaluation on the Swire dataset [26] shows that our model outperform Swire and CLIP on sketch-to-GUI retrieval, with a recall@10 of 0.685.

6.2.1 Model Architecture. In order to adapt UIClip for sketch-to-GUI retrieval, we employed its image encoder. The left side of the Figure 10 illustrates the architecture of the model, which encompasses a screenshot encoder and a sketch encoder. Both of these components are initialized from UIClip’s image encoder. During the training phase, the weights of the screenshot encoder are held constant, whilst training is solely conducted on the sketch encoder. While during the inference phase, which is illustrated at the right side of the Figure 10, both of the two encoders are frozen. The embedding of the query sketch is compared with the embeddings of each screenshot, the top-k screenshots whose embedding is the most similar with the sketch embedding are retrieved.

6.2.2 Dataset. The dataset employed in this study was procured from Swire [26], which encompasses 3802 sketches corresponding to 2201 screenshots from 167 apps within the Rico dataset. This dataset was compiled by four designers, who contributed to sketching 505, 1017, 1272, and 1008 screenshots respectively. However, it was discovered that certain images from this dataset were not present within the Rico dataset, hence after appropriate filtration, a refined dataset

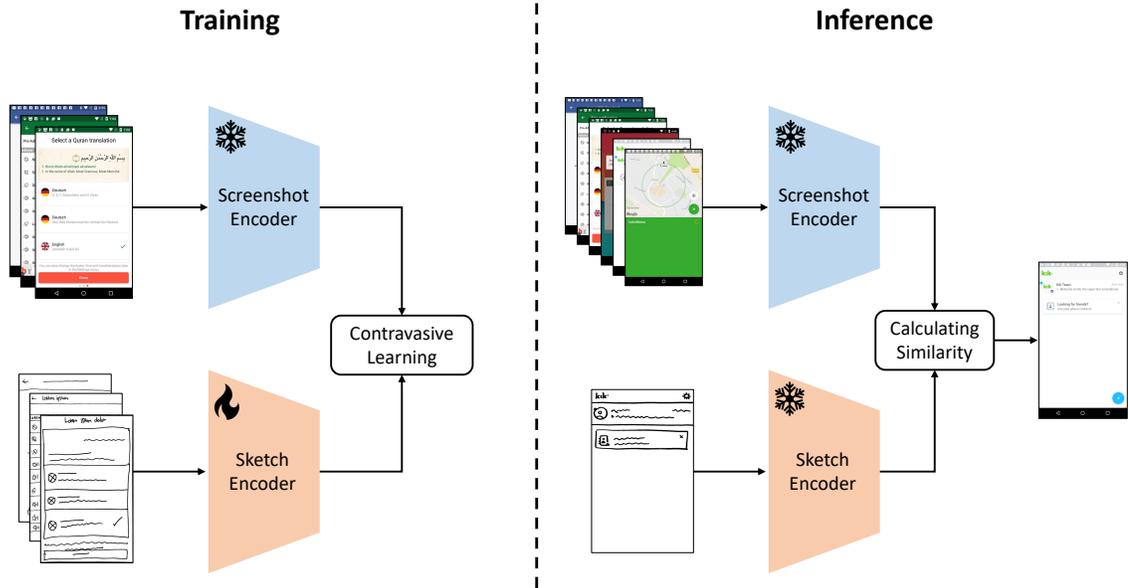


Fig. 10. Overview of the training and inference of Dual UIClip ViT model for sketch-to-GUI retrieval

featuring a total of 3551 sketches was obtained. The split between the training and test sets within Swire’s dataset was not explicitly delineated in their work. Consequently, for the purposes of our study, we adopted 486 sketches from one designer as the test set, while the sketches from three other designers constituted the training set. This is to ensure that the model exhibits generalization across sketches developed by various designers.

6.2.3 Baselines. To investigate the benefit of using UIClip in sketch-to-GUI retrieval, we compared our model with the following baselines.

- *Swire*. Swire establishes itself as the pioneer sketch-to-GUI retrieval model [26]. This model comprises two distinct VGG-A networks [52] that independently compute the embeddings for screenshots and sketches. The retrieval process is subsequently effectuated by measuring the similarity between these embeddings. Given that Swire’s implementation is not publicly accessible, we took it upon ourselves to recreate it according to the details laid out in the original paper.
- *Dual CLIP ViT*. The architecture of the Dual CLIP ViT is identical to that of UIClip, with their main distinguishing factor being their respective parameters’ weights. The weights of the image encoder of Dual CLIP ViT are derived directly from the CLIP model.

6.2.4 Training Details of Fine-tuning. Under zero-shot setting, the Dual UIClip ViT and Dual CLIP ViT have not been trained. While under fine-tuning setting, the Dual UIClip ViT and Dual CLIP ViT model has been trained in contrastive learning manner on the training set for 10 epochs. We used AdamW optimizer with an initial learning rate $1e^{-5}$, batch size 256.

For the training of Swire, we tried to use the same configurations as specified in the original paper, including a learning rate of 0.01 and a batch size of 32. However, we noticed a significant shortfall in performance with these

Table 7. Retrieval accuracy on the test set of Swire dataset

	Zero-shot				Fine-tuned			
	Recall@1	Recall@3	Recall@5	Recall@10	Recall@1	Recall@3	Recall@5	Recall@10
Swire	-	-	-	-	0.027	0.062	0.117	0.214
Dual CLIP ViT	0.019	0.029	0.037	0.053	0.177	0.296	0.422	0.533
Dual UIClip ViT	0.062	0.117	0.156	0.224	0.368	0.580	0.685	0.772

settings, wherein the top-10 recall stood at a mere 0.02. To address this, we made the decision to utilize a batch size of 64 and a learning rate of $1e^{-5}$. The original article does not specify the number of training epochs employed; thus, we proceeded to train it for 100 epochs utilising the training set, until we observe convergence in the training loss.

6.2.5 Evaluation Metric. Similar as the evaluation of text-to-GUI retrieval presented on the Section 4.1.4, we use the recall@1, recall@3, recall@5 and recall@10 for evaluation.

6.2.6 Results and Discussion. Table 7 presents the evaluation outcomes of UIClip and various baseline models on the Swire dataset within both zero-shot and fine-tuning settings. Our observation revealed that even under zero-shot setting, the Dual UIClip ViT has a recall@10 score of 0.224, which indicates its valuable application in sketch-to-GUI retrieval. In contrast, the Dual CLIP ViT was not as effective, as evidenced by its low recall@10 score of 0.053.

Fine-tuning can significantly improve the models’ performance. After fine-tuning, the recall@10 score of the Dual CLIP ViT increased to 0.553, while the Dual UIClip ViT exhibited a great improvement, reaching the score of 0.772. These results substantiate the superior efficacy of UIClip in comparison to CLIP for sketch-to-GUI retrieval.

The performance of our implementation of Swire model, with a recall@1 of 0.027 and a recall@10 of 0.214, is not congruent with the metrics reported in the originating study by Huang et al. [26]. The cited research demonstrated significantly higher recall values of 0.159 and 0.609, at ranks 1 and 10 respectively. This discrepancy could potentially be attributed to two factors: 1) Our evaluation used 486 pairings for testing, in contrast to the 276 pairings used in the original research. The size of the test set can dramatically influence retrieval performance outcomes; 2) The divergence in implementation methods could also account for differing results. It’s important to note, however, that the original Swire model’s implementation is not open-source for replication. Nonetheless, even if we take into consideration the results from Huang et al. [26], our Dual UIClip ViT model still demonstrates superior performance, with a score of 0.368 for recall@1 and 0.772 for recall@10.

6.3 Additional Tasks and Applications

Besides GUI classification and sketch-to-GUI retrieval, our model has the potential to be adapted for other GUI-related tasks.

6.3.1 GUI Captioning. GUI captioning is a fundamental task in GUI understanding, where the model predicts a textual informative caption to a given GUI [30, 32, 59]. ClipCap [42] is an image captioning model. It applies CLIP as image encoder, followed by a simple mapping network to generate the caption prefix. Subsequently, this caption prefix is utilized for the production of the caption through a language model. Given the superior performance of the UIClip model in understanding GUI information, it is possible to enhance ClipCap on GUI Captioning task by replacing the CLIP image encoder with the image encoder of UIClip.

6.3.2 GUI generation. GUI generation is the task that generates GUIs that aligns to the given specifications. Wei et al. [60] introduced UI-Diffuser that leverage stable diffusion model to generate mobile GUIs through simple textual descriptions and UI components. The text encoder of CLIP is the core component of the stable diffusion model, it converts the textual description to embedding that can be processed by the subsequent networks. However, given that the CLIP model is trained on general image-caption pairs, it may fall short on understanding the GUI-related descriptions. A potential improvement to the stable diffusion model for GUI generation tasks could be the replacement of the text encoder used in CLIP with that of UIClip’s.

7 RELATED WORK

7.1 Mobile GUI Retrieval

GUI retrieval refers to the search of existing GUI designs with queries of various format.

By wireframe: Wireframe is an image that represents the skeletal layout of a screenshot. Given a wireframe image, one can retrieve the corresponding screenshots. Deka et al. [15] and Liu et al. [38] trained autoencoder [5] by reconstructing the wireframe to generate the wireframe embedding for GUI retrieval by embedding similarity. Chen et al. [11] performed wireframe-based GUI retrieval by training an autoencoder for encoding the visual semantics of UI designs using a large database of UI design wireframes.

By sketch image: UI sketch is often used at the early stage of design, it can also be used as the query for the GUI retrieval. Huang et al. introduced Swire [26], which comprises two separate VGG-A networks [52]. These networks calculate the embeddings for screenshots and sketches, respectively. The retrieval process is subsequently performed by measuring the similarity of these embeddings. Unlike Swire that treats the sketch image as a unit, Mohian et al. [40, 41] performed object detection on sketch to identify the UI components on a sketch image, and then find the screenshots in Rico dataset whose type and location of UI components match best the sketch.

By screenshot image: Using a screenshot from an app, it is possible to retrieve similar screenshots from the GUI repository. Screen2vec [35] extracts the UI components, layouts and app description to generate a screen embedding, querying similar screenshots is done by comparing the similarity of screen embedding. On the other hand, VINS [7] creates screen embedding by combining both image embedding and UI components embedding for screen searching. Instead of querying screenshots, GUIFetch [4] searches for apps in public repositories that are as similar as possible to the provided app’s screens and transitions between them.

By textual query: There are also other works that use textual queries for GUI retrieval. GUIGLE [6] has indexed a collection of GUIs along with their metadata from the ReDraw dataset [43]. GUIGLE leverages various information sources such as text displayed on the screen, name of UI components, and app name to retrieve relevant screens based on the keyword matching to the given query. On the other hand, RaWi [28] uses a BERT-based ranking model to retrieve GUIs from the Rico dataset [15]. This retrieval process involves matching the GUI text (text displayed on the screen, GUI activity name and GUI component identifier) with the textual query in order to identify the most relevant images. Unlike the preceding studies that focus on retrieving entire screenshot, Gallery D.C. [9, 21] allow users to search UI components, like button and checkbox. It provide a GUI gallery that allow users to search UI components with different filters, including component type, size, color, app category and displayed text. The aforementioned methods used the textual query to search the text presented on the GUIs, while the non-textual information have not been taken into account.

7.2 Vision Language Models for Mobile GUI

Vision-language models combines both the vision and language modalities, it has been trained on image-text pairs and can be applied to a range of tasks [19, 65].

Up until now, the main use case of VLMs in mobile GUI domain is the GUI captioning. GUI captioning refers to the task of generating a high-level summary of a screenshot to describe its contents and functionalities. The Screen2Words model [59], built on the foundations of the Transformer Encoder-Decoder architecture [58], uses multi-modal data, inclusive of the screenshot image, view hierarchy, and app description, to yield the screen caption. Contrarily, Spotlight, a vision-only approach employed for multiple tasks including widget captioning and screen captioning, possess architecture rooted in Vision Transformer [17] and the T5 model [50].

While some studies pertaining to GUI-related tasks may appear to employ VLMs, they do not in reality. For instance, XUI [30] model does not address the language modality, it generates the screen caption with template-based natural language generation (NLG) engine. As presented on the previous section, GUIGLE [6] and RaWi [28] conducted text-based GUI retrieval by utilizing the textual content exhibited in the screenshots or the metadata, instead of employing the vision information for searching purposes. Such an approach may result in the omission of crucial vision information of the UI screen. To the best of our knowledge, no existing research employs VLMs for text-to-GUI retrieval.

7.3 Mobile GUI Datasets

7.3.1 Large Scale Dataset. The large-scale mobile GUI dataset comprises a vast collection of screenshots, thereby serving as a valuable resource for GUI retrieval.

Rico [15, 38] is one of the largest GUI datasets in the literature. Deka et al. combined crowd sourcing and automation to mine design and interaction data from Android apps at runtime. It exposes visual, textual, structural, and interactive design properties of more than 66k unique UI screens from more than 9.3k Android apps. Several other dataset are built based on Rico [29, 56, 59, 64]. However, as the Rico dataset was created seven years ago, some of its UIs may be outdated for contemporary apps.

Chen et al. [11] developed an extensive dataset that comprises 54,987 screenshots derived from 7,748 Android apps, spanning 25 distinct app categories. These screenshots are obtained with automated GUI exploration [10].

Moran et al. [43] collected the ReDraw dataset in a fully automated manner by mining and executing the top-250 Android apps in each category on Google Play, excluding games. This yielded a total of 14,382 unique screenshots and 191,300 labeled GUI components after data cleaning.

WebUI [62] is a large dataset of 400,000 rendered web pages associated with extracted metadata using automated web crawling. Although this dataset is not designed specifically for mobile apps, it can be useful for understanding mobile UI through transfer learning.

More and more automated GUI exploration methods are introduced to explore app UI automatically from iOS apps [61] or Android apps [13, 14, 66]. Nonetheless, GUI datasets gathered through automated exploration may inadvertently omit important app features. Because the access to certain pages may necessitate app authorization or initial configurations, which are often absent in automatic exploration methods [13]. Generally automated UI explorer are rather useful for UI testing than for getting design insights.

7.3.2 GUI Caption Dataset. A GUI caption dataset is essential for the training of vision language models [33, 34, 49, 53].

Screen2Words dataset [59] is a subset of Rico, it contains 112,085 manually created descriptions across 22,417 unique screenshots. Moran et al. created Clarity dataset [45]. They selected a subset of screenshots from the ReDraw dataset

[43], and created a annotated dataset consisting of 45,998 descriptions for 10,204 screenshots from popular Android app. Instead of focusing on entire screens, Li et al. [36] released a dataset containing 162,859 language phrases created by human workers for annotating 61,285 UI components across 21,750 unique UI screens.

8 THREATS TO VALIDITY

This section discusses potential threats to the validity of our study.

Query Selection for Search Engine Evaluation. The selection of queries is crucial for evaluation, as the performance of a GUI search engine may vary depends on the query. For example, a GUI search engine that exhibits satisfactory results in the health and fitness domain may not necessarily produce the satisfactory outcomes in finance or education domains. During the evaluation process, it is imperative to consider the varying domains of applications. The study of Kolthoff et al. [28] provided a gold standard of queries and corresponding GUIs. However, this dataset has already been exhausted for the training of RaWi, thereby rendering it unsuitable for testing within our study. Given the lack of prior queries, we instead sourced from articles published by companies that specialize in mobile app development services. The app features referenced in [57] span ten distinct domains, while [27] details innovative features from even more domains. Based on this, we posit that these selected sets are representative enough for our evaluation.

Subjectivity of manual evaluation. As for every study that includes manual evaluation, the process of evaluation a search engine may introduce inherent bias. In an attempt to mitigate bias and subjectivity within the evaluation, we incorporated a number of strategic steps. The evaluation outcomes may be influenced by the evaluators' preference towards specific search engines. To tackle this issue, we created an evaluation tool which mixed and shuffled the resulting screenshots originating from the three distinct search engines. As a result, the evaluator cannot discern the origin of a specific screenshot, thereby curbing potential bias among diverse search engines. Moreover, the evaluation process was carried out by four evaluators, each with a minimum of five years' experience in software development. They conducted a careful evaluation based on a evaluation guide (e.g. stating what is or not a relevant GUI) and using a uniform evaluation tool. And each GUI was assessed by at least two evaluators to minimize bias attributed to individual preference.

9 CONCLUSION AND FUTURE WORK

In this paper, we presented GUiing, a GUI search engine, built upon UIClip, a novel vision-language foundation model for the mobile GUI domain. The UIClip model was trained using the GPSCap dataset, which is a comprehensive dataset comprised of 135k screenshot-caption pairs. In addition, we created the GPSRepo dataset to serve as the search source for GUiing, encompassing a total of 303k screenshots. GPSCap and GPSRepo, were generated through a fully automatic pipeline, which enables to regularly update the datasets, e.g. with new screen designs. The pipeline enables the collection, classification, cropping, and caption recognition of app introduction images as found in Google Play as well as in other stores too. Our evaluation, including a benchmarking on various datasets and a manual assessment of search results, indicate that our approach outperforms the state-of-the-art methods for GUI retrieval. Moreover, the evaluation highlights the critical role that the GPSCap and GPSRepo datasets play in achieving the top performance of the vision-language models within GUI retrieval. There are several future directions to follow up on our work:

- *Increasing and diversifying the screenshot dataset.* The performance of a vision-language model is profoundly correlated with the size of its training dataset. According to AppBrain [2], Google Play includes approximately 2.36 million apps as of 2024. Our training dataset includes merely 117k apps. Mining additional app introduction

images will thus likely improve the performance of GUing and relevance of the results. Moreover, other platforms like Apple iOS, Samsung Galaxy Watch, or Microsoft HoloLens have their stores with similar app introduction images. Learning and searching those screenshots can similarly inspire the design of corresponding apps.

- *Hybrid search methods.* The inherent limitations of natural language can pose difficulties in describing the desired GUIs exhaustively through text. To enhance the retrieval process, it might be beneficial to support a more diverse range of search methods. For instance, sketch-to-GUI retrieval would enable users to find GUI designs (or related follow up screens) based on their sketched designs. Additionally, GUI-to-GUI retrieval would enable users to retrieve other GUIs that exhibit a resemblance to the original interfaces.
- *Evaluations with practitioners and integration into the app development workflows.* Although our empirical evaluation show the efficacy of GUing, the reality of software development can be much more complex and nuanced. Thus, additional empirical studies with app developers, UI designers, and requirements engineers would likely lead to comprehensive feedback on real usage, scalability, further search parameters, as well as how GUing and UIClip can be ideally integrated into existing app development workflows and tools.

REFERENCES

- [1] AppBrain. 2024. Google Play Ranking in the United States (Oct 2023). <https://www.appbrain.com/stats/google-play-rankings>. Accessed: 2023-10-01.
- [2] AppBrain. 2024. Number of Android apps on Google Play (Mar 2024). <https://www.appbrain.com/stats/number-of-android-apps>. Accessed: 2024-3-10.
- [3] Sunil Arya, David M. Mount, Nathan S. Netanyahu, Ruth Silverman, and Angela Y. Wu. 1998. An optimal algorithm for approximate nearest neighbor searching in fixed dimensions. *J. ACM* 45, 6 (1998), 891–923. <https://doi.org/10.1145/293347.293348>
- [4] Farnaz Behrang, Steven P. Reiss, and Alessandro Orso. 2018. GUIfetch: Supporting app design and development through GUI search. *Proceedings - International Conference on Software Engineering* (2018), 236–246. <https://doi.org/10.1145/3197231.3197244>
- [5] Yoshua Bengio. 2009. *Learning deep architectures for AI* Vol. 2. 1–27 pages. <https://doi.org/10.1561/2200000006>
- [6] Carlos Bernal-Cardenas, Kevin Moran, Michele Tufano, Zichang Liu, Linyong Nan, Zhehan Shi, and Denys Poshyvanyk. 2019. Guile: A GUI search engine for android apps. *Proceedings - 2019 IEEE/ACM 41st International Conference on Software Engineering: Companion, ICSE-Companion 2019* (2019), 71–74. <https://doi.org/10.1109/ICSE-Companion.2019.00041> arXiv:1901.00891
- [7] Sara Bunian, Kai Li, Chaima Jemmali, Casper Hartevelde, Yun Fu, and Magy Seif Seif El-Nasr. 2021. VINS: Visual Search for Mobile User Interface Design. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (CHI '21)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3411764.3445762>
- [8] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-End Object Detection with Transformers. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 12346 LNCS (2020), 213–229. https://doi.org/10.1007/978-3-030-58452-8_13 arXiv:2005.12872
- [9] Chunyang Chen, Sidong Feng, Zhenchang Xing, Linda Liu, Shengdong Zhao, and Jinshui Wang. 2019. Gallery D.C.: Design search and knowledge discovery through auto-created GUI component gallery. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019). <https://doi.org/10.1145/3359282>
- [10] Chunyang Chen, Ting Su, Guozhu Meng, Zhenchang Xing, and Yang Liu. 2018. From UI design image to GUI skeleton: A neural machine translator to bootstrap mobile GUI implementation. In *Proceedings - International Conference on Software Engineering*, Vol. 6. 665–676. <https://doi.org/10.1145/3180155.3180240>
- [11] Jieshan Chen, Chunyang Chen, Zhenchang Xing, Xin Xia, Liming Zhu, John Grundy, and Jinshui Wang. 2020. Wireframe-based UI Design Search through Image Autoencoder. *ACM Transactions on Software Engineering and Methodology* 29, 3 (2020). <https://doi.org/10.1145/3391613> arXiv:2103.07085
- [12] Qiuyuan Chen, Chunyang Chen, Safwat Hassan, Zhengchang Xing, Xin Xia, and Ahmed E. Hassan. 2021. How Should I Improve the UI of My App?: A Study of User Reviews of Popular Apps in the Google Play. *ACM Transactions on Software Engineering and Methodology* 30, 3 (2021), 1–37. <https://doi.org/10.1145/3447808>
- [13] Sen Chen, Lingling Fan, Chunyang Chen, and Yang Liu. 2023. Automatically Distilling Storyboard With Rich Features for Android Apps. *IEEE Transactions on Software Engineering* 49, 2 (2023), 667–683. <https://doi.org/10.1109/TSE.2022.3159548> arXiv:2203.06420
- [14] Sen Chen, Lingling Fan, Chunyang Chen, Ting Su, Wenhe Li, Yang Liu, and Lihua Xu. 2019. StoryDroid: Automated Generation of Storyboard for Android Apps. In *Proceedings - International Conference on Software Engineering*, Vol. 2019-May. IEEE, 596–607. <https://doi.org/10.1109/ICSE.2019.00070> arXiv:1902.00476

- [15] Biplab Deka, Zifeng Huang, Chad Franzen, Joshua Hibschan, Daniel Afergan, Yang Li, Jeffrey Nichols, and Ranjitha Kumar. 2017. Rico: A mobile app dataset for building data-driven design applications. *UIST 2017 - Proceedings of the 30th Annual ACM Symposium on User Interface Software and Technology* (2017), 845–854. <https://doi.org/10.1145/3126594.3126651>
- [16] Jacob Devlin, Ming Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL HLT 2019 - 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference*, Vol. 1. 4171–4186. arXiv:1810.04805 <https://arxiv.org/abs/1810.04805>
- [17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image Is Worth 16X16 Words: Transformers for Image Recognition At Scale. In *ICLR 2021 - 9th International Conference on Learning Representations*. arXiv:2010.11929
- [18] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The Faiss library. (2024). arXiv:2401.08281 [cs.LG]
- [19] Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. 2022. A Survey of Vision-Language Pre-Trained Models. *IJCAI International Joint Conference on Artificial Intelligence* (2022), 5436–5443. <https://doi.org/10.24963/ijcai.2022/762> arXiv:2202.10936
- [20] Explosion. 2024. Prodigy - An annotation tool for AI, Machine Learning and NLP. <https://prodigy.ai/>. Accessed: 2024-3-10.
- [21] Sidong Feng, Chunyang Chen, and Zhenchang Xing. 2022. Gallery D.C.: Auto-created GUI Component Gallery for Design Search and Knowledge Discovery. In *Proceedings - International Conference on Software Engineering*, Vol. 1. Association for Computing Machinery, 80–84. <https://doi.org/10.1109/ICSE-Companion55297.2022.9793764> arXiv:2204.06700
- [22] Alessio Ferrari and Paola Spoletini. 2023. Strategies, Benefits and Challenges of App Store-inspired Requirements Elicitation. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 1290–1302. <https://doi.org/10.1109/ICSE48619.2023.00114>
- [23] Google. 2024. Add preview assets to showcase your app - Play Console Help. <https://support.google.com/googleplay/android-developer/answer/9866151?hl=en&sjid=206438066775745925-EU>. Accessed: 2024-3-10.
- [24] Marlo Haering, Christoph Stanik, and Walid Maalej. 2021. Automatically matching bug reports with related app reviews. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE, 970–981.
- [25] Safwat Hassan, Cor Paul Bezemer, and Ahmed E. Hassan. 2020. Studying Bad Updates of Top Free-to-Download Apps in the Google Play Store. *IEEE Transactions on Software Engineering* 46, 7 (2020), 773–793. <https://doi.org/10.1109/TSE.2018.2869395>
- [26] Forrest Huang, John F. Canny, and Jeffrey Nichols. 2019. Swire: Sketch-based User Interface Retrieval. *Conference on Human Factors in Computing Systems - Proceedings* (2019), 1–10. <https://doi.org/10.1145/3290605.3300334>
- [27] Nic Hughart. 2023. 50 Best App Ideas For 2024. <https://buildfire.com/best-app-ideas>. Accessed: 2024-3-10.
- [28] Kristian Kolthoff, Christian Bartelt, and Simone Paolo Ponzetto. 2023. Data-driven prototyping via natural-language-based GUI retrieval. *Automated Software Engineering* 30, 1 (2023), 13. <https://doi.org/10.1007/s10515-023-00377-x>
- [29] Luis A. Leiva, Asutosh Hota, and Antti Oulasvirta. 2020. Enrico: A Dataset for Topic Modeling of Mobile UI Designs. *Extended Abstracts - 22nd International Conference on Human-Computer Interaction with Mobile Devices and Services: Expanding the Horizon of Mobile Interaction, MobileHCI 2020* (2020). <https://doi.org/10.1145/3406324.3410710>
- [30] Luis A. Leiva, Asutosh Hota, and Antti Oulasvirta. 2023. Describing UI Screenshots in Natural Language. *ACM Transactions on Intelligent Systems and Technology* 14, 1 (2023), 1–28. <https://doi.org/10.1145/3564702>
- [31] Chenxia Li, Weiwei Liu, Ruoyu Guo, Xiaoting Yin, Kaitao Jiang, Yongkun Du, Yuning Du, Lingfeng Zhu, Baohua Lai, Xiaoguang Hu, Dianhai Yu, and Yanjun Ma. 2022. PP-OCRv3: More Attempts for the Improvement of Ultra Lightweight OCR System. (2022). arXiv:2206.03001 <http://arxiv.org/abs/2206.03001>
- [32] Gang Li and Yang Li. 2022. Spotlight: Mobile UI Understanding using Vision-Language Models with a Focus. 2022 (2022), 1–16. arXiv:2209.14927 <http://arxiv.org/abs/2209.14927>
- [33] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models. (2023). arXiv:2301.12597 <http://arxiv.org/abs/2301.12597>
- [34] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation. *Proceedings of Machine Learning Research* 162, 2 (2022), 12888–12900. arXiv:2201.12086
- [35] Toby Jia Jun Li, Lindsay Popowski, Tom M. Mitchell, and Brad A. Myers. 2021. Screen2vec: Semantic embedding of GUI screens and GUI components. *Conference on Human Factors in Computing Systems - Proceedings* (2021). <https://doi.org/10.1145/3411764.3445049> arXiv:2101.11103
- [36] Yang Li, Gang Li, Luheng He, Jingjie Zheng, Hong Li, and Zhiwei Guan. 2020. Widget captioning: Generating natural language description for mobile user interface elements. *EMNLP 2020 - 2020 Conference on Empirical Methods in Natural Language Processing, Proceedings of the Conference* 2015 (2020), 5495–5510. <https://doi.org/10.18653/v1/2020.emnlp-main.443> arXiv:2010.04295
- [37] Tsung Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft COCO: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference*. https://doi.org/10.1007/978-3-319-10602-1_48 arXiv:1405.0312
- [38] Thomas F. Liu, Mark Craft, Jason Situ, Ersin Yumer, Radomir Mech, and Ranjitha Kumar. 2018. Learning design semantics for mobile apps. *UIST 2018 - Proceedings of the 31st Annual ACM Symposium on User Interface Software and Technology* (2018), 569–579. <https://doi.org/10.1145/3242587.3242650>
- [39] Walid Maalej, Hans-Jörg Happel, and Asarnusch Rashid. 2009. When users become collaborators: towards continuous and context-aware user input. In *Proceedings of the 24th ACM SIGPLAN conference companion on Object oriented programming systems languages and applications (OOPSLA)*.

- 981–990.
- [40] Soumik Mohian and Christoph Csallner. 2022. PSDoodle: Fast App Screen Search via Partial Screen Doodle. *Proceedings - 9th IEEE/ACM International Conference on Mobile Software Engineering and Systems, MOBILESoft 2022* January (2022), 89–99. <https://doi.org/10.1145/3524613.3527816>
- [41] Soumik Mohian and Christoph Csallner. 2023. Searching Mobile App Screens via Text + Doodle. (2023). arXiv:2305.06165 <http://arxiv.org/abs/2305.06165>
- [42] Ron Mokady, Amir Hertz, and Amit H Bermano. 2021. ClipCap : CLIP Prefix for Image Captioning. (2021). arXiv:2111.09734 <https://arxiv.org/abs/2111.09734>
- [43] Kevin Moran, Carlos Bernal-Cardenas, Michael Curcio, Richard Bonett, and Denys Poshyvanyk. 2020. Machine Learning-Based Prototyping of Graphical User Interfaces for Mobile Apps. *IEEE Transactions on Software Engineering* 46, 2 (2020), 196–221. <https://doi.org/10.1109/TSE.2018.2844788> arXiv:1802.02312
- [44] Kevin Moran, Boyang Li, Carlos Bernal-Cárdenas, Dan Jelf, and Denys Poshyvanyk. 2018. Automated reporting of GUI design violations for mobile apps. In *40th International Conference on Software Engineering*. 165–175. <https://doi.org/10.1145/3180155.3180246> arXiv:1802.04732
- [45] Kevin Moran, Ali Yachnes, George Purnell, Junayed Mahmud, Michele Tufano, Carlos Bernal Cardenas, Denys Poshyvanyk, and Zach H'Doubler. 2022. An Empirical Investigation into the Use of Image Captioning for Automated Software Documentation. *Proceedings - 2022 IEEE International Conference on Software Analysis, Evolution and Reengineering, SANER 2022* (2022), 514–525. <https://doi.org/10.1109/SANER53432.2022.00069> arXiv:2301.01224
- [46] Facundo Olano. 2015. Google Play Scraper. <https://github.com/facundoolano/google-play-scraper>. Accessed: 2024-3-10.
- [47] OpenAI. 2021. CLIP. <https://github.com/openai/CLIP/blob/main/clip/clip.py>. Accessed: 2024-3-10.
- [48] OpenAI. 2021. Model Card openai/clip-vit-base-patch32 on HuggingFace. <https://huggingface.co/openai/clip-vit-base-patch32>. Accessed: 2024-3-10.
- [49] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proceedings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 8748–8763. arXiv:2103.00020 <http://arxiv.org/abs/2103.00020><https://proceedings.mlr.press/v139/radford21a.html>
- [50] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research* 21 (2020), 1–67. arXiv:1910.10683
- [51] Tobias Roehm, Nigar Gurbanova, Bernd Bruegge, Christophe Joubert, and Walid Maalej. 2013. Monitoring user interactions for supporting failure reproduction. In *2013 21st International Conference on Program Comprehension (ICPC)*. IEEE, 73–82.
- [52] Karen Simonyan and Andrew Zisserman. 2015. Very deep convolutional networks for large-scale image recognition. In *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*. 1–14. arXiv:1409.1556
- [53] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. 2022. FLAVA: A Foundational Language And Vision Alignment Model. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition 2022-June* (2022), 15617–15629. <https://doi.org/10.1109/CVPR52688.2022.01519> arXiv:2112.04482
- [54] Filip Sondej. 2020. Autocorrect. <https://github.com/filyp/autocorrect>. Accessed: 2024-3-10.
- [55] Peter M. Stahl. 2022. Lingua. <https://github.com/pemistahl/lingua-py>. Accessed: 2024-3-10.
- [56] Srinivas Sunkara, Maria Wang, Lijuan Liu, Gilles Baechler, Yu Chung Hsiao, Jindong Chen, Abhanshu Sharma, and James Stout. 2022. Towards Better Semantic Understanding of Mobile Interfaces. *Proceedings - International Conference on Computational Linguistics, COLING 29, 1* (2022), 5636–5650. arXiv:2210.02663
- [57] Vladimir Terekhov. 2023. 138 Features to Consider While Developing a Mobile App. <https://attractgroup.com/blog/most-comprehensive-list-of-mobile-app-features-while-developing-a-mobile-application>. Accessed: 2024-3-10.
- [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in Neural Information Processing Systems 2017-Decem*, Nips (2017), 5999–6009. arXiv:1706.03762
- [59] Bryan Wang, Gang Li, Xin Zhou, Zhourong Chen, Tovi Grossman, and Yang Li. 2021. Screen2Words: Automatic Mobile UI Summarization with Multimodal Learning. In *UIST 2021 - Proceedings of the 34th Annual ACM Symposium on User Interface Software and Technology*. 498–510. <https://doi.org/10.1145/3472749.3474765> arXiv:2108.03353
- [60] Jialiang Wei, Anne-Lise Courbis, Thomas Lambolais, Binbin Xu, Pierre Louis Bernard, and Gérard Dray. 2023. Boosting GUI Prototyping with Diffusion Models. In *2023 IEEE 31st International Requirements Engineering Conference (RE)*. 275–280. <https://doi.org/10.1109/RE57278.2023.00035> arXiv:2306.06233
- [61] Jason Wu, Rebecca Krosnick, Eldon Schoop, Amanda Swearngin, Jeffrey P. Bigham, and Jeffrey Nichols. 2023. Never-ending Learning of User Interfaces. (2023). <https://doi.org/10.1145/3586183.3606824> arXiv:2308.08726
- [62] Jason Wu, Siyan Wang, Siman Shen, Yi Hao Peng, Jeffrey Nichols, and Jeffrey P. Bigham. 2023. WebUI: A Dataset for Enhancing Visual UI Understanding with Web Semantics. *Conference on Human Factors in Computing Systems - Proceedings* (2023). <https://doi.org/10.1145/3544548.3581158> arXiv:2301.13280
- [63] Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. 2023. C-Pack: Packaged Resources To Advance General Chinese Embedding. (2023). arXiv:2309.07597 <http://arxiv.org/abs/2309.07597>

- [64] Alaa Zaki and Mohamed Abdallah. 2023. MASC : A Dataset for the Development and Classification of Mobile Applications Screens. (2023), 1–15. <https://doi.org/10.21203/rs.3.rs-3786876/v1>
- [65] Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. 2023. Vision-Language Models for Vision Tasks: A Survey. March (2023), 1–23. [arXiv:2304.00685](https://arxiv.org/abs/2304.00685) <http://arxiv.org/abs/2304.00685>
- [66] Xiangyu Zhang, Lingling Fan, Sen Chen, Yucheng Su, and Boyuan Li. 2023. Scene-Driven Exploration and GUI Modeling for Android Apps. (2023). [arXiv:2308.10228](https://arxiv.org/abs/2308.10228) <http://arxiv.org/abs/2308.10228>