

---

# Papt: A Pairwise GUI Dataset between Android Phones and Tablets

---

## 1 Dataset Ethics and Responsible Usage

2 In our commitment to promote responsible research practices, we present a comprehensive set of  
3 ethical guidelines for the usage of the Papt dataset. Researchers and developers are urged to be fully  
4 aware of these directives and adhere to them in their entirety.

5 **Intended Use:** The Papt dataset is designed and released strictly for non-profit research purposes.  
6 Any commercial use or application is strictly prohibited.

7 **No Verbatim Publication:** While the dataset provides insights into the GUI layout of various  
8 applications, it may contain content from apps that researchers inadvertently downloaded. It is crucial  
9 to understand and respect that these contents may be protected under copyright laws. Therefore, users  
10 of the Papt dataset must not publish any verbatim content from the applications, such as articles from  
11 recognized outlets like BBC news.

12 **Respect for Copyrights:** The dataset might encompass copyrighted material inherent in some apps'  
13 GUI. Users should be aware of potential copyright infringements and must not reproduce, distribute,  
14 or create derivative works based on these copyrighted components without necessary permissions.

15 **Data Responsibility:** While leveraging the dataset for research, users should be cautious not to  
16 misuse any part of the data, ensuring that their research does not infringe upon the rights of original  
17 content creators or the privacy of users.

18 **Citing the Dataset:** Any research work that benefits from the Papt dataset should duly acknowledge  
19 it. Proper citation ensures that the creators and contributors receive appropriate credit, fostering a  
20 community of shared resources and collaborative research.

21 **Report Violations:** Users are encouraged to report any violations or potential misuse of the Papt  
22 dataset. This proactive step will help safeguard the ethical principles that underscore this dataset and  
23 the broader research community.

24 We trust that these guidelines will be meticulously followed, ensuring that the Papt dataset remains a  
25 valuable, responsible, and ethical resource for the research community.

## 26 2 Applications and Impact of the Pairwise GUI Dataset in Research and 27 Industry

28 A pairwise GUI dataset encompassing both phone and tablet platforms constitutes a pivotal asset in  
29 the burgeoning field of automated GUI development. Analyzing this dataset allows for the extraction  
30 of underlying patterns, verification of existing methodologies, and the potential training of models to  
31 recognize or generate analogous structures. Such an extensive dataset fosters innovation in several  
32 crucial domains, including but not limited to cross-platform design, user experience optimization,  
33 rigorous testing and validation protocols, adaptive design strategies, and the enhancement of accessi-

34 bility features. Beyond merely serving practical developmental needs, this dataset creates a fertile  
35 landscape for theoretical exploration, extending the current horizons of modern interface design  
36 research.

37 **Facilitating Cross-Platform Design:** A pairwise GUI dataset that includes representations for  
38 both phone and tablet interfaces provides a comprehensive understanding of how GUI elements  
39 adapt across different devices. This can enable researchers and developers to create algorithms that  
40 facilitate automatic resizing and rearranging of UI components to fit different screens, enhancing  
41 cross-platform compatibility.

42 **Improving User Experience:** By analyzing the relationships between phone and tablet GUI pairs,  
43 researchers can identify optimal design patterns that offer a seamless user experience across devices.  
44 These insights can guide automated GUI development tools to generate interfaces that maintain  
45 consistency and usability, whether viewed on a phone or a tablet.

46 **Enabling Multi-modal Development:** The rich dataset that includes various components and their  
47 metadata in GUI pairs allows for exploration into multi-modal interface development. Researchers  
48 can investigate how to intelligently adapt GUI components based on user interaction patterns on  
49 different devices. This could lead to automated development processes that personalize the interface  
50 according to the specific platform.

51 **Enhancing Testing and Validation:** With a pairwise dataset, automated testing tools can be de-  
52 veloped to compare and contrast how a GUI performs and appears on both phones and tablets.  
53 This allows for comprehensive validation processes that ensure consistency and functionality across  
54 platforms, thereby reducing development time and improving quality.

55 **Encouraging Adaptive Design Research:** The differences between phone and tablet GUI pairs  
56 highlight the need for adaptive design principles. Researchers can leverage this dataset to experiment  
57 with algorithms that automatically adjust GUI elements to different orientations, resolutions, and user  
58 preferences, fostering more adaptive and responsive design techniques.

59 **Facilitating Collaboration:** By establishing a common dataset for phone and tablet GUI pairs,  
60 collaboration among researchers is enhanced. It creates a standardized platform for development and  
61 evaluation, encouraging innovation in automated design tools, methodologies, and principles.

62 **Handling Complex Design Challenges:** The detailed information within the dataset, such as UI  
63 property and GUI hierarchies, allows for sophisticated modeling of complex design scenarios. This  
64 can lead to breakthroughs in handling intricate design challenges like non-correspondence between  
65 phone and tablet GUIs, leading to more robust and flexible design solutions.

66 **Boosting Accessibility Research:** The pairwise dataset may help researchers understand how  
67 accessibility features should be adjusted across different devices. Automated tools can then be  
68 developed to ensure that interfaces are inclusive and meet accessibility standards, whether accessed  
69 via a phone or a tablet.

### 70 **3 Descriptive Statistics of the Papt Dataset**

71 Using the methodologies described in the paper, we have gathered 10,035 valid phone-tablet GUI  
72 pairs. To provide a detailed insight into our dataset, we statistically analyze these pairs from two  
73 main aspects: the distribution of UI view types discussed in subsection 3.2, and the evaluation of  
74 GUI similarities between pairs as outlined in subsection 3.3. Subsection 3.4 details the format of  
75 the data in our dataset. Lastly, we discuss the merits of our dataset in comparison to existing GUI  
76 datasets in subsection ??.

#### 77 **3.1 Source App Pairs**

78 We first crawl 6,456 tablet apps from Google Play. Then we match their corresponding phone apps  
79 by their app names and app developers. Finally, we collect 5,593 valid phone-tablet app pairs from

80 22 app categories. Table 1 shows the top 15 categories of 5,593 app pairs. Due to the effect of the  
 81 data’s long tail disctribution, we only display the top 15 categories. The column *Category* represents  
 82 the category of these apps. The column *#Count* and *P(%)* denote the number of apps in this category  
 83 and their percentage of the overall number of apps, respectively. These 5,593 phone-tablet app  
 84 pairs are the data source for this dataset. The three most common categories of apps in the data  
 85 source are: *Entertainment* (8.87%), *Social* (7.04%) and *Communication* (5.83%). As shown in  
 86 Table 1, the categories of apps in our data source are scattered and balanced. Most of the categories  
 87 occupy between 4% and 6% of the total dataset. This balanced distribution ensures the dataset’s  
 88 generalizability and diversity.

Table 1: Top 15 categories of source apps.

Category	#Count	P (%)
Entertainment	496	8.87
Social	394	7.04
Communication	326	5.83
Lifestyle	318	5.69
Books & Reference	286	5.11
Education	279	4.98
News & Magazines	271	4.85
Shopping	270	4.83
Sports	267	4.78
Music & Audio	266	4.76
Weather	265	4.73
Finance	262	4.68
Bussiness	261	4.67
Travel & Local	255	4.57
Medical	254	4.54

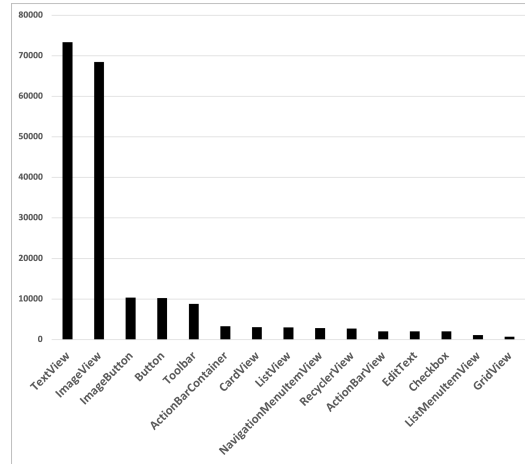


Figure 1: Distribution of top 15 UI view types in the dataset.

### 89 3.2 Distribution of UI View Types

90 In Android development, a UI view is a basic building block for creating user interfaces. Views are  
 91 responsible for drawing and handling user interactions for a portion of the screen [1]. For example, a  
 92 button, a text , an image, and a list are all a type of view.

93 Figure 1 illustrates the distribution of UI View types in the dataset. Considering the data’s long-tail  
 94 distribution, we only display the top 15 types. We can see that the *TextView* and *ImageView* types,  
 95 including all their derived categories such as *AppCompatActivity*, *AppCompatActivity*, and  
 96 *AppCompatActivity*, are the most common UI view types in the dataset. Their numbers (73,349  
 97 and 68,496) significantly outnumber all other view types. The GUI primarily presents information via  
 98 text and images, so text and image-related views are the most prevalent in the database. *ImageButton*  
 99 (10,366) and *Button* (10,235) are the third and fourth most UI views. Users interact with the GUI  
 100 mainly through clicks and click operations rely heavily on button views, so button-related views are  
 101 also common in GUI datasets.

### 102 3.3 Distribution of GUI Pair Similarity

103 The similarity analysis between phone-tablet GUI pairs is important for downstream tasks. Given a  
 104 GUI pair, there are a total of  $M$  and  $N$  GUI views in the GUIs of the phone and tablet, respectively.  
 105 Suppose there are  $L$  the same views in the GUIs of the phone and the tablet. The similarity of their  
 106 GUIs is calculated as

$$Sim(M, N) = \frac{2 * L}{M + N} \quad (1)$$

107 Figure 2 shows the frequency histogram of GUI similarities of our phone-tablet GUI pairs in the  
 108 dataset. The similarity between the GUIs of phones and tablets in most pairs is between 0.5 and 0.7.  
 109 Considering the difference in screen size between tablets and phones, the current phone GUI page

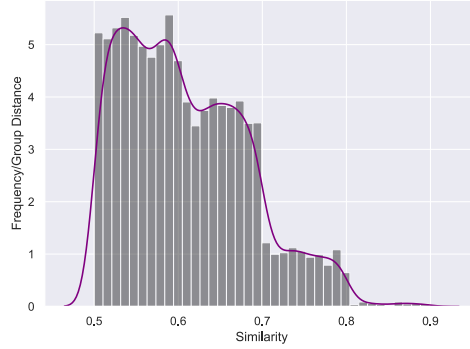


Figure 2: The frequency histogram of GUI similarity of collected pairs

110 can only contain part of the UI views in the corresponding tablet GUI page, and the current data  
 111 reminds us that when performing downstream tasks such as GUI layout generation, search, etc., we  
 112 should consider filling in the contents that are not available in the mobile phone GUI page.

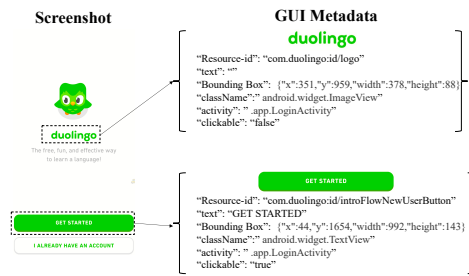


Figure 3: A screenshot of a GUI and its part UI metadata

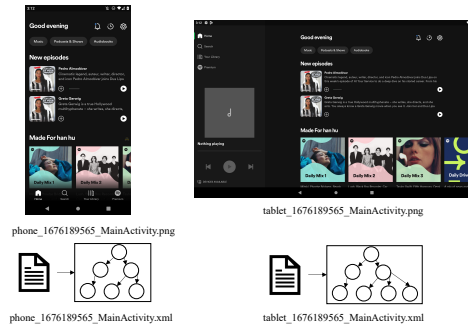


Figure 4: An example GUI pair in Spotify

### 113 3.4 Data Format

114 In this section, we first show an example screenshot of a GUI and its corresponding UI metadata.  
 115 Then, we introduce the format of each GUI pair in the dataset.

116 **GUI Screenshot and its Metadata** We install and run phone-tablet app pairs in Section ?? on the  
 117 Pixel6 and Samsung Galaxy tab S8, respectively. We use uiautomator2 [19] to collect screenshots  
 118 and GUI metadata of the dynamically running apps. Figure 3 shows an example of a collected GUI  
 119 screenshot and metadata of some UI components inside the GUI. This example is from the app  
 120 'Duolingo [9]. The metadata is a documentary object model (DOM) tree of current GUIs, which  
 121 includes the hierarchy and properties (e.g., class, bounding box, layout) of UI components. We can  
 122 infer the GUI hierarchy from the DOM tree hierarchy in metadata.

123 **GUI Pair** Figure 4 shows an example of pairwise GUI pages of the app 'Spotify' in  
 124 our dataset. All GUI pairs in one phone-tablet app pair are placed in the same direc-  
 125 tory. Each pair consists of four elements: a screenshot of the GUI running on the phone  
 126 (*phone\_1676189565\_MainActivity.png*), the metadata data corresponding to the GUI screenshot on  
 127 the phone (*phone\_1676189565\_MainActivity.xml*), a screenshot of the GUI running on the tablet  
 128 (*tablet\_1676189565\_MainActivity.png*), and the metadata data corresponding to the GUI screenshot  
 129 on the tablet (*tablet\_1676189565\_MainActivity.xml*). The naming format for all files in the dataset  
 130 is *Device\_Timestamp\_Activity Name*. As shown in Figure 4, The filename *tablet\_1676189565\_-*  
 131 *MainActivity.xml* indicates that this file was obtained by the tablet and was collected with the  
 132 timestamp *1676189565*, this GUI belongs to *MainActivity* and this file is a metadata file in XML  
 133 format. We use timestamps and activity names to distinguish phone-tablet GUI pairs.

134 **Pairs in the Dataset** The dataset is made accessible to the public in accordance with the criteria  
135 outlined in the attached license agreements<sup>1</sup>. The pairs in the dataset are contained in separate folders  
136 according to the app. Most of the app folders are named after the package name of the app’s APK,  
137 for example, *air.com.myheritage.mobile* , and a few are named after the app’s name, for example,  
138 *Spotify*. In each app folder, as described in Section 3.4, each pair contains four elements: the phone  
139 GUI screenshot, the XML file of the phone GUI metadata, the corresponding tablet GUI screenshot,  
140 and the XML file of the tablet GUI metadata. We also shared the script for loading all GUI pairs in  
141 the open source repository.

### 142 3.5 Data Collection Tool

143 Based on the above-described two collection strategies, we develop two distinct collecting tools: the  
144 adjust resolution collector and the similarity matching collector.

145 The first tool dynamically adjusts the resolution of the current device using ADB instructions. When  
146 the running app detects a change in the screen’s resolution, it will call the layout file designed for the  
147 tablet and change the layout of the current GUI.

148 The second tool concurrently runs two apps of one app pair on a mobile phone and a tablet. The tool  
149 dynamically evaluates the similarity of the GUIs presented on two devices, and automatically collects  
150 the matched GUI page pair when the similarity exceeds a predetermined threshold.

151 These two data collection tools are also included in the repository of the publicly accessible dataset.  
152 With the installation instructions provided, more researchers can utilise our tools to collect more  
153 customised GUI datasets for future research.

### 154 3.6 Accessing the Dataset

155 The dataset is made accessible to the public in accordance with the criteria outlined in the attached  
156 license agreements<sup>2</sup>. The pairs in the dataset are contained in separate folders according to the  
157 app. Most of the app folders are named after the package name of the app’s APK, for example,  
158 *air.com.myheritage.mobile* , and a few are named after the app’s name, for example, *Spotify*. In each  
159 app folder, each pair contains four elements: the phone GUI screenshot, the XML file of the phone  
160 GUI metadata, the corresponding tablet GUI screenshot, and the XML file of the tablet GUI metadata.  
161 We also shared the script for loading all GUI pairs in the open source repository.

## 162 4 Preliminary Experiments

163 To demonstrate the usability of our dataset, we perform preliminary experiments on the dataset.  
164 These experiments contain two types of tasks: GUI conversion and GUI retrieval. We select current  
165 state-of-art approaches for these two tasks and use the automatic metrics for evaluation. We will  
166 discuss the qualitative and quantitative performance as well as the limitations of selected approaches  
167 on our proposed dataset.

### 168 4.1 Tasks

169 **GUI Conversion** This task is to automatically convert an existing GUI to a new layout of GUI [3].  
170 For example, given the present phone layout, generate a tablet layout automatically. It is also a GUI  
171 generation task. The goal of GUI conversion task is to generate a flexible interface that is applicable  
172 to different platform. By generating GUI layouts that fit the needs of developers, we can make  
173 developer more productive and decrease engineering effort.

---

<sup>1</sup><https://github.com/huhanGitHub/papt>

<sup>2</sup><https://github.com/huhanGitHub/papt>

174 **GUI Retrieval** This task is to retrieval the most relevant GUI from the database and recommend it  
175 to the most appropriate users [21]. GUI template search and recommendation is an essential direction  
176 for current automated GUI development to accelerate the development process.

## 177 4.2 Selected Approaches

### 178 4.2.1 GUI Conversion

179 As far as we know, no researcher has proposed a method for generating a tablet GUI from a phone  
180 GUI, so we will use some existing relevant GUI generation methods to apply to this task. Three  
181 approaches, which are all widely used in GUI generation, are selected in this task:

182 **LayoutTransformer** [10]: a simple yet powerful auto-regressive model based on Transformer  
183 framework that leverages self-attention to generate layouts by learning contextual relationships  
184 between different layout elements. It is able to generate a brand new layout either from an empty set  
185 or from an initial seed set of primitives, and can easily scale to support an arbitrary of primitives per  
186 layout.

187 **LayoutVAE** [13]: a stochastic model based on variational autoencoder architecture. It is composed  
188 of two modules: CountVAE which predicts the number of objects and BBoxVAE which predicts the  
189 bounding box of each object. It is capable of generating full image layouts given a label set, or per  
190 label layouts for an existing image given a new label. Besides, it is also capable of detecting unusual  
191 layouts, potentially providing a way to evaluate layout generation problem.

192 **VTN** [2]: a VAE-based framework advanced by Transformer model, which is able to learn margins,  
193 alignments and other elements without explicit supervision. Specifically, the encoder and decoder  
194 are parameterized by attention-based neural network. During the variational process, VTN sample  
195 latent representations from the prior distributions and transform those into layouts using self-attention  
196 based decoder.

### 197 4.2.2 GUI Retrieval

198 Researchers have not yet attempted to get the comparable tablet GUI design from a phone GUI design,  
199 but numerous methods for locating relevant GUI designs have been proposed. These comparable  
200 GUI design retrieval approaches will be selected and adapted to fit this task. For GUI retrieval task,  
201 we employ three semantic matching and learning-based models:

202 **Rico** [8]: a neural-based training framework that utilizing content-agnostic similarity heuristic  
203 method for UI comparing and matching. To facilitate query-by-example search, Rico exposes a vector  
204 representation for each UI that encodes layout. Rico provides search engines with several visual  
205 representations that can be served up as results: UI screenshots, flows, and animations.

206 **GUIFetch** [4]: a method that takes an input the sketch for an app and leverages the growing number  
207 of open source apps in public repositories and then devise a component-matching model to rank the  
208 identified apps using a combination of static and dynamic analysis and computes a similarity metric  
209 between the models and the provided sketch.

210 **WAE** [7]: a wireframe-based UI searching model using image autoencoder architecture. Specifically,  
211 it is a neural based approach using convolutional neural network (CNN) in an unsupervised manner  
212 for building a UI design search engine that is flexible and robust in face of the great variations in UI  
213 designs. The enhancement of wireframe will facilitate the layout generation process.

Table 2: Automatic evaluation results on the test set.

Model	mIoU ↓	Overlap ↓	W class ↓	W bbox ↓	# Unique matches ↑	Matched rate ↓
LayoutVAE	0.10	0.23	0.29	0.012	356	0.13
LayoutTransformer	0.12	0.32	0.31	0.024	445	0.15
VTN	0.13	0.35	0.37	0.026	541	0.19

214 **4.3 Evaluation Metrics**

215 We consider automatic evaluation on both these two tasks.

216 **4.3.1 GUI Conversion**

217 For GUI conversion task, it’s important to evaluate layouts in terms of two perspectives: perceptual  
218 quality and diversity. We follow with the similar evaluation protocol in [2, 10, 5] and utilize a set of  
219 metrics to evaluate the quality and diversity aspects. Specifically, we use the following metrics:

220 **Mean Intersection over Union (mIoU) [18]**: also known as the Jaccard index [11], is a method to  
221 quantify the percent overlap between the ground-truth and generated output. The IoU is calculated by  
222 dividing the overlap area between predicted class positions and ground truth by the area of union  
223 between predicted position and ground truth. So, it is computed by

$$mIoU = \frac{1}{k} \sum_{i=0}^k \frac{TP(i)}{TP(i) + FP(i) + FN(i)} \tag{2}$$

224 where  $k$  means  $k$  classes in both images,  $TP(i)$ ,  $FP(i)$  and  $FN(i)$  represent the distribution of true  
225 positive, false positive and false negative of  $i_{th}$  class between two compared images.

226 **Overlap [20, 10]**: measures the overlap ratio between the ground-truth and our generated output.  
227 The overlap metric use the total overlapping area among any two bounding boxes inside the whole  
228 page and the average IoU between elements.

229 **Wasserstein (W) distance [15]**: is the distance of the classes and bounding boxes to the real data. It  
230 contains  $W$  class and  $W$  bbox metrics. Wasserstein distance is to evaluate the diversity between the  
231 real and generated data distributions.

232 **Unique matches [16]**: is the number of unique matchies according to the DocSim [16]. It measures  
233 the matching overlap between real sets of layouts and generated layouts. It is designed for diversity  
234 evaluation.

235 **Matched rate**: To more directly show how many UI components in the ground truth’s tablet GUI  
236 are successfully and automatically converted by the model, we select the metric: the matched rate.  
237 Suppose there are a total of  $m$  UI components in the ground truth, and there are  $n$  components in  
238 the generated tablet GUI that match the components in the ground truth, then the matched rate is  
239 calculated as  $n/m$ .

240 **4.3.2 GUI Search**

241 For GUI retrieval task, we evaluate the performance of a UI-design search method by two metrics:  
242 Precision and Mean Reciprocal Rank (MRR), which have been widely-used in GUI search [7, 8, 21,  
243 6, 12].

244 **Precision@k (Pre@k)**: Precision@k is the proportion of the top-k results for a query UI that are  
245 relevant UI designs. Specifically, the calculation of ranking Precision@k is defined as follows:

$$\text{Precision@k} = \frac{\#relevantUIDesign}{k}, \tag{3}$$

246 As we consider the original UI as the only relevant UI for a treated UI in this study, we use the  
247 strictest metric Pre@1: Pre@1=1 if the first returned UI is the original UI, otherwise Pre@1=0.

248 **Mean Reciprocal Rank (MRR)**: MRR is another method to evaluate systems that return a ranked  
249 list. It computes the mean of the reciprocal rank (i.e.,  $1/\text{rank}$ ) of the first relevant UI design in the  
250 search results over all query UIs. Specifically, the calculation of MRR is defined as follows:

$$MRR = \frac{1}{Q} \sum_{i=1}^Q \frac{1}{rank_i}, \tag{4}$$

251 where  $Q$  refers to the number of all the query UIs. The higher value a metric is, the better a search  
252 method performs.

253 **4.4 GUI Conversion Task**

254 **4.4.1 Experimental Setup**

255 Since we are comparing the structure of the rendered GUI as opposed to its pixels, we do not consider  
256 the contents of images and texts inside the GUI. Therefore, we convert GUI pages to wireframes by  
257 converting each category of UI component to a box of the specified color, which has been widely  
258 adopted in recently related works [10, 14, 13]. In our GUI conversion work, unlike typical GUI  
259 generation tasks, we ask the model to learn to generate a tablet GUI comparable to the input phone  
260 GUI, rather than another phone GUI. As input to the training model, we encode all UI components of  
261 a phone GUI as a component sequence from top to bottom and left to right. Each UI component in  
262 the GUI is encoded as a quaternion  $(x, y, w, h)$ . Where  $x$  and  $y$  denote the x and y coordinates of the  
263 upper left corner of this UI component, and  $w$  and  $h$  denote the length and width of the component.  
264 Similarly, the tablet GUIs in the pair are encoded into a sequence that serves as the ground truth for  
265 training and validation.

266 Following the setup of experiments in LayoutTransformer and LayoutVAE, we randomly select  
267 1000 pairs ( 10%) in the total dataset as the test set and the rest of the data as the training set. All  
268 the results of the metrics are calculated on the test set. All of the results are based on 5-fold cross  
269 validation on the test set.

270 **4.4.2 Quantitative Evaluation**

271 We present the quantitative evaluation results on the Table 2. As expected, we observe that VTN  
272 model achieves the best performances in terms of all metrics. It yields large improvement regarding  
273 perceptual quality, such as *W class*, and *Unique matches*. Besides, it obviously outperforms the  
274 LayoutVAE model across all metrics and is slightly better than LayoutTransformer model. The reason  
275 heavily relies on the mutual enhancement between self-attention and VAE mechanisms. However,  
276 according to the final metric *Matched rate*, only less than 20% of the UI components can be matched  
277 between the generated tablet GUI layouts and ground truth. We demonstrate that most phone-tablet  
278 GUI pairs in the dataset have a similarity between 50% and 70%, so the current model’s capacity to  
279 learn the relationship between phone-tablet GUI pairs still requires improvement. The results suggest  
280 that further work is needed to design and train a more effective model in the GUI conversion task.  
281 We hope that our open source dataset will enable more researchers to participate in automated UI  
282 development and contribute more effective methods.

283 **4.4.3 Qualitative Evaluation**

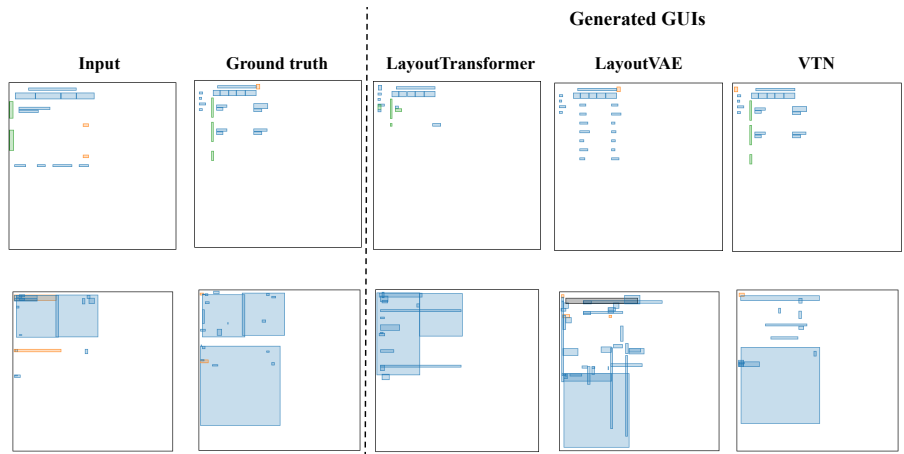


Figure 5: Examples of generated GUIs by selected three approaches

284 To better understand the performance of different models on our task, we present qualitative results  
285 of selected models and their generated outputs in Figure 5. In alignment with the quantitative results



286 in Table 2, we observe that the VTN model outperform other selected models. It demonstrates the  
 287 efficiency of self-attention mechanism on the layout generation task. However, we can observe that  
 288 the best selected model VTN still fall short of generating precise margins and positions towards the  
 289 ground-truth. There is still much room for improvement in learning the relationship between phone  
 290 and tablet GUIs, and the current results are far from helpful to GUI developers. The GUI conversion  
 291 in the crossing platform poses huge challenges for the existing state-of-the-art model. Therefore,  
 292 simply utilizing previous layout generation model cannot tackle the challenges of GUI conversion in  
 293 both Android phones and tablets.

## 294 4.5 GUI Retrieval

### 295 4.5.1 Experimental Setup

296 Following related GUI search works [7, 17], we convert the GUI pages into wireframes in the same  
 297 way as the GUI conversion task in Section 4.4.1. The input in this experiment is the phone GUI  
 298 wireframe and the ground truth is the corresponding tablet GUI wireframe. We train a GUI encoder  
 299 to encode GUI into numerous vectors and compare the similarity of these numerous vectors. We  
 300 randomly select 1000 pairs ( 10%) in the total dataset as the test set and the rest of the data as the  
 301 training set. The results of this test set is automatically checked.

302 Due to the existence of comparable GUI designs on the tablet, sometimes the search result is not the  
 303 ground truth, but the GUI design is reasonable and very close to the ground truth. We would also  
 304 consider this result to be an appropriate design. Therefore, we randomly selected 100 cases from the  
 305 test set and manually verify the results. We invited three industry developers with at least one year of  
 306 experience in Android development to manual evaluate search results. Each participant evaluate the  
 307 top 10 search results to determine if they are a reasonable tablet design. After the initial evaluating,  
 308 three volunteers have a discussion and merge conflicts. They clarify their findings, scope boundaries  
 309 among categories and misunderstanding in this step. Finally, they iterate to revise evaluation results  
 310 and discuss with each other until consensus is reached. To differentiate, we call these results manual  
 311 checked results.

Table 3: Qualitative results for GUI retrieval and recommendation task.

Model	Pre@1	Pre@5	Pre@10	MRR
Rico [8]	0.73	0.65	0.58	0.69
GUIFetch [4]	0.63	0.54	0.52	0.62
WAE [7]	0.80	0.77	0.75	0.83

### 312 4.5.2 Results

313 Table 3 shows the performance metrics of the three selected methods in our dataset. We report four  
 314 metrics in terms of precision following with previous UI search work: *Pre1*, *Pre5*, *Pre10* and *MRR*.  
 315 The results in column *Pre@1* and *MRR* are automatic checked results. The results in column *Pre@5*  
 316 and *Pre@10* are manual checked results. We can observe that the *WAE* model outperforms other  
 317 two approaches in terms of automated and manual results. Compared with GUI conversion task,  
 318 GUI search task is more developed and applicable. Learning-based approaches *Rico* and *WAE* both  
 319 achieved high search accuracy. It demonstrates the advantages of neural network models with huge  
 320 training parameters in extracting features and patterns of GUI. We hope that more researchers in the  
 321 future will design more search-related tasks based on our dataset, such as semantic-based search,  
 322 GUI template recommendation, etc.

## 323 References

- 324 [1] Android View, 2023. <https://data-flair.training/blogs/android-layout-and-views/>.
- 325 [2] Diego Martin Arroyo, Janis Postels, and Federico Tombari. Variational transformer networks  
 326 for layout generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
 327 *Pattern Recognition*, pages 13642–13652, 2021.

- 328 [3] Miroslav Behan and Ondrej Krejcar. Adaptive graphical user interface solution for modern user  
329 devices. In *Intelligent Information and Database Systems: 4th Asian Conference, ACIIDS 2012,*  
330 *Kaohsiung, Taiwan, March 19-21, 2012, Proceedings, Part II 4*, pages 411–420. Springer, 2012.
- 331 [4] Farnaz Behrang, Steven P Reiss, and Alessandro Orso. Guifetch: supporting app design and  
332 development through gui search. In *Proceedings of the 5th International Conference on Mobile*  
333 *Software Engineering and Systems*, pages 236–246, 2018.
- 334 [5] Sara Bunian, Kai Li, Chaima Jemmali, Casper Hartevelde, Yun Fu, and Magy Seif Seif El-Nasr.  
335 Vins: Visual search for mobile user interface design. In *Proceedings of the 2021 CHI Conference*  
336 *on Human Factors in Computing Systems*, pages 1–14, 2021.
- 337 [6] Chunyang Chen, Sidong Feng, Zhenchang Xing, Linda Liu, Shengdong Zhao, and Jinshui  
338 Wang. Gallery dc: Design search and knowledge discovery through auto-created gui component  
339 gallery. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–22, 2019.
- 340 [7] Jieshan Chen, Chunyang Chen, Zhenchang Xing, Xin Xia, Liming Zhu, John Grundy, and  
341 Jinshui Wang. Wireframe-based ui design search through image autoencoder. *ACM Transactions*  
342 *on Software Engineering and Methodology (TOSEM)*, 29(3):1–31, 2020.
- 343 [8] Biplab Deka, Zifeng Huang, Chad Franzen, Joshua Hibschan, Daniel Afergan, Yang Li,  
344 Jeffrey Nichols, and Ranjitha Kumar. Rico: A mobile app dataset for building data-driven  
345 design applications. In *Proceedings of the 30th annual ACM symposium on user interface*  
346 *software and technology*, pages 845–854, 2017.
- 347 [9] Duolingo, 2023. <https://www.duolingo.com/>.
- 348 [10] Kamal Gupta, Justin Lazarow, Alessandro Achille, Larry S Davis, Vijay Mahadevan, and  
349 Abhinav Shrivastava. Layouttransformer: Layout generation and completion with self-attention.  
350 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1004–  
351 1014, 2021.
- 352 [11] Lieve Hamers et al. Similarity measures in scientometric research: The jaccard index versus  
353 salton’s cosine formula. *Information Processing and Management*, 25(3):315–18, 1989.
- 354 [12] Forrest Huang, John F Canny, and Jeffrey Nichols. Swire: Sketch-based user interface retrieval.  
355 In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages  
356 1–10, 2019.
- 357 [13] Akash Abdu Jyothi, Thibaut Durand, Jiawei He, Leonid Sigal, and Greg Mori. Layoutvae:  
358 Stochastic scene layout generation from a label set. In *Proceedings of the IEEE/CVF Interna-*  
359 *tional Conference on Computer Vision*, pages 9895–9904, 2019.
- 360 [14] Jianan Li, Jimei Yang, Aaron Hertzmann, Jianming Zhang, and Tingfa Xu. Layoutgan: Gener-  
361 ating graphic layouts with wireframe discriminators. *arXiv preprint arXiv:1901.06767*, 2019.
- 362 [15] Victor M Panaretos and Yoav Zemel. Statistical aspects of wasserstein distances. *Annual review*  
363 *of statistics and its application*, 6:405–431, 2019.
- 364 [16] Akshay Gadi Patil, Omri Ben-Eliezer, Or Perel, and Hadar Averbuch-Elor. Read: Recursive  
365 autoencoders for document layout generation. In *Proceedings of the IEEE/CVF Conference on*  
366 *Computer Vision and Pattern Recognition Workshops*, pages 544–545, 2020.
- 367 [17] Jerome Revaud, Jon Almazán, Rafael S Rezende, and Cesar Roberto de Souza. Learning with  
368 average precision: Training image retrieval with a listwise loss. In *Proceedings of the IEEE/CVF*  
369 *International Conference on Computer Vision*, pages 5107–5116, 2019.
- 370 [18] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio  
371 Savarese. Generalized intersection over union: A metric and a loss for bounding box regression.  
372 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages  
373 658–666, 2019.

- 374 [19] Uiautomator2, 2023. <https://play.google.com/store>.
- 375 [20] Varduhi Yeghiazaryan and Irina Voiculescu. Family of boundary overlap metrics for the  
376 evaluation of medical image segmentation. *Journal of Medical Imaging*, 5(1):015006–015006,  
377 2018.
- 378 [21] Tom Yeh, Tsung-Hsiang Chang, and Robert C Miller. Sikuli: using gui screenshots for search  
379 and automation. In *Proceedings of the 22nd annual ACM symposium on User interface software*  
380 *and technology*, pages 183–192, 2009.