# Screen Correspondence: Mapping Interchangeable Elements between UIs

Jason Wu, Amanda Swearngin, Xiaoyi Zhang, Jeffrey Nichols, Jeffrey Bigham

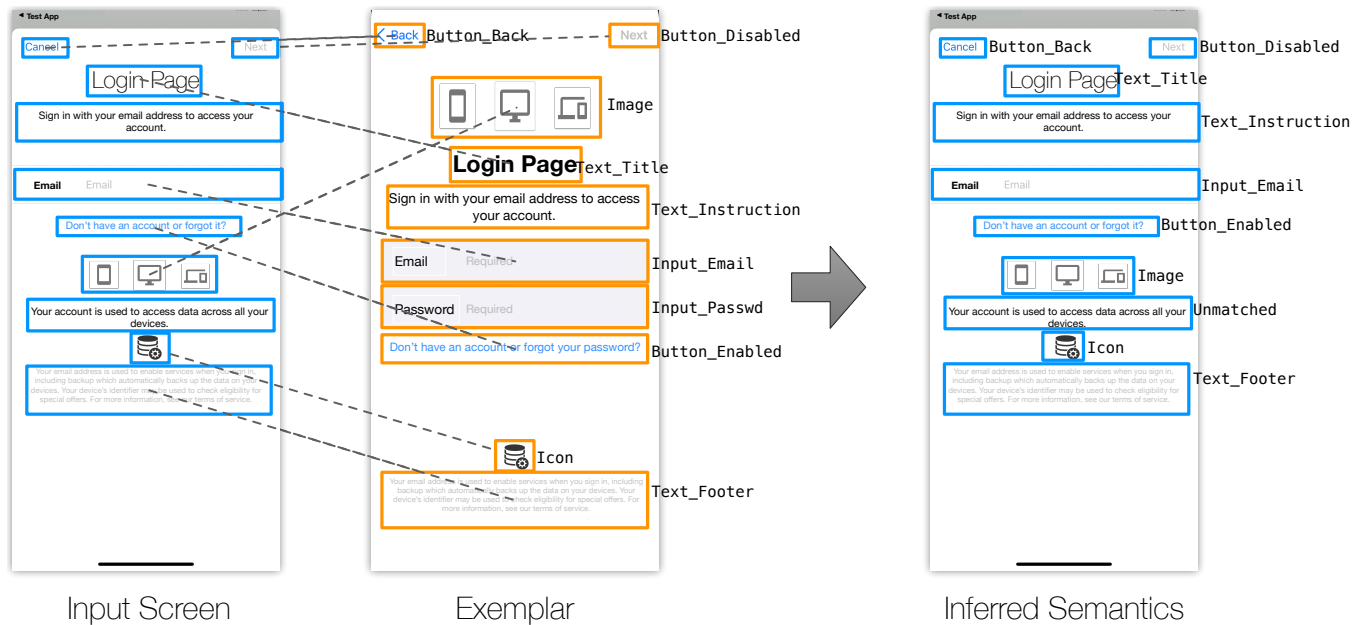{jsonwu,aswearngin,xiaoyiz,jwnichols,jbigham}@apple.com

**Figure 1: Screen correspondence produces a mapping of similar UI elements across two UIs that have related elements. Screenshots are encoded using a multi-modal model that segments and featurizes UI elements. Mappings are generated that link element pairs that have the same or similar functionality across UI screens.**

## ABSTRACT

Understanding user interface (UI) functionality is a useful yet challenging task for both machines and people. In this paper, we investigate a machine learning approach for *screen correspondence*, which allows reasoning about UIs by mapping their elements onto previously encountered examples with known functionality and properties. We describe and implement a model that incorporates element semantics, appearance, and text to support correspondence computation without requiring any labeled examples. Through a comprehensive performance evaluation, we show that our approach improves upon baselines by incorporating multi-modal properties of UIs. Finally, we show three example applications where screen correspondence facilitates better UI understanding for humans and machines: *(i)* instructional overlay generation, *(ii)* semantic UI element search, and *(iii)* automated interface testing.

## KEYWORDS

user interface modeling, ui semantics, element correspondence

## 1 INTRODUCTION

Understanding how user interfaces (UIs) can be operated to achieve some goal can be challenging for both machines and humans, especially those who are less tech-savvy. While automated systems in the right circumstances can provide useful assistance [62, 63] or automatically complete the task themselves [34, 36], people can be hindered or completely blocked by apps that do not provide necessary metadata, such as the view hierarchy. A promising approach involves inferring UI functionality solely from the pixels rendered

to the screen, but to date this method has primarily been useful for identifying the location and type of typical UI elements [65] and not higher-level semantics. For example, these algorithms can identify that a screen has a button that contains the text "Login," but are unaware of the higher-level concept of logging in to a service, and they cannot infer what the role of this button would be in that process.

There are many higher-level semantics in user interfaces (e.g., login, account registration, shopping carts), which would correspond to an enormous number of classes if we attempted to use a classifier to predict their occurrence. Instead of making class predictions, an alternate approach to data inference involves directly comparing inputs to previously encountered examples with known properties. Studies on human [55] and machine [56] learning suggest that direct comparison is a useful tool, especially when relevant examples are available in the form of analogies [4, 24] or templates [20, 22]. This concept can be highly effective for UIs, as many belong to the same app or are constructed to serve a similar purpose. For example, knowledge of how an app screen was previously interacted with by an app crawler or automated UI tester could aid in producing more robust and consistent results when visited again. Similar inferences can also be made for related screens from different apps, such as by determining that a button with the label "Login" in a new app is likely used to submit a login request because that is how a similar button is used in a known app.

In this paper, we pose the problem of *screen correspondence* to map interchangeable elements between two UI screenshots (Figure 1). We introduce a multi-modal transformer model for detecting, featurizing, and matching UI elements. Our approach is *unsupervised*, which allows it to work without a large dataset of labeled examples, which could be costly and time-consuming to collect. In a performance evaluation with strong baselines, we compare our approach to existing correspondence algorithms used in computer vision (CV) and heuristics such as schema-matching. Our results indicate that our multi-modal model outperforms all existing baselines.

We describe and implement three example applications that show the utility of screen correspondence for humans and machines to understand UI functionality. We create an application to **generate instructional overlays** by transferring high-quality human-authored coach marks (a type of instructional label) from one screen to another of the same category (*e.g.*, two registration screens). To support UI design search and exemplar-based exploration, we used our model to **index a large dataset of UI elements and screens**. Finally, we built a system to **aid an automated app crawler** by identifying mappings between the elements of screens from different runs.

To summarize, we make the following contributions:

- We introduce *screen correspondence* as a method of mapping interchangeable elements between UI screens from their screenshots.
- We describe a machine learning approach to generating correspondence between two UI screenshots, and we show it outperforms existing baselines.
- We show the utility of screen correspondence in three example applications that improve both human and machine understanding of UI functionality.

## 2 RELATED WORK

Our work is related to recent work in understanding user interfaces from their pixels, and also a variety of methods for understanding applications in terms of their many screens. We also overview machine learning solutions to correspondence problems in other domains, such as computer vision and natural language processing.

### 2.1 Predicting Screen Semantics

Computational representation of user interfaces are useful for many downstream tasks, such as design assistance [32, 39], accessibility improvement [65], and task-oriented systems. Screen Recognition [65] generates accessibility metadata of a UI from screenshots using an object detection model and heuristics. Screen Parsing [60] generates structured UI models from screenshots of UIs. Several models [6, 15, 41, 64] have also been trained to predict the semantics of unlabeled icons found in mobile apps. These models can be applied to improve the accessibility of mobile apps, either as a tool during design time or as an automated system that repairs existing apps at runtime. Most of these models map UI elements to a pre-defined set of classes (*e.g.* UI element and icon type), which may exclude less common components [7].

An alternative is to train models using self-supervision [9, 18, 35], which allows them to take advantage of larger unlabeled datasets. Screen2Vec [35] and other pixel-based autoencoders [9, 39] map UIs to fixed-length embedding vectors which can be used to represent semantic properties. The Pixel-words model [18] employs a transformer model architecture and masked training objective based on prior work in NLP [11]. Our work builds upon these approaches to train a model for identifying UI element correspondences between screens.

### 2.2 Multi-screen Understanding

While many automated UI systems can benefit from understanding the semantics of a single screen, screens are rarely used in isolation. Any task or interaction trace requires reasoning about multiple app screens and how they are related to each other.

StoryDroid is a system that extracts a storyboard of Android apps from APK files as an "App Transition Graph" [8]. ActionBert [23] models the relationship between two consecutive UI screens by predicting, among other things, which UI element was tapped on the first screen to reach the second (*i.e.*, link component prediction). Longer sequences of touch interactions have also been modeled to better understand user behavior and app usage [31, 69].

A particular problem that many multi-screen systems aim to address is identifying whether two screens are instances of the same UI, a problem which we refer to as "screen fingerprinting." NEAR [61] detects near-duplicate pages on the web using a combination of visual and DOM-based features. Prior work [14] used supervised learning to predict the relationship and transitions between screenshots by, among other things, classifying whether inputs were different instances of the same screen (*e.g.*, a news app with dynamic content).

Screen fingerprinting is useful for comparing screens to known examples; however, finer grain mappings (*e.g.*, element-level fingerprinting) can result in higher fidelity comparisons and additional benefits. Bricolage [30] is a system that renders the content of one

web page using the style and layout of another. It employs a supervised element matching model that featurizes web elements based on their DOM representation and was trained on a dataset of human-generated mappings. Interaction proxies [67] rely on a set of equivalency heuristics to identify UI components and structures found in Android view hierarchies to facilitate accessibility repair.

Our work is related to these approaches in multi-screen understanding and specifically element fingerprinting. While many previous examples relied heavily on the availability of a structured UI representation (*e.g.,* DOM, view hierarchy) and were trained on labeled data, our approach requires only screenshots of related apps with optional labels.

## 2.3 Machine Learning of Correspondence

Machine learning has been used to learn correspondences in other domains, such as computer vision (CV) and natural language processing (NLP), which are closely related to our approach.

A longstanding problem in CV is inferring accurate correspondence of objects from different images. Homography estimation [21] involves finding a mapping, either sparse or dense, or a transformation matrix that describes perspective changes in two images of the same object or scene. Optical flow [25] applies a similar concept to finding mappings between consecutively taken images. A common approach involves computing appearance descriptors (keypoints), then creating a mapping that optimizes the global correspondence [16, 38]. Recent work [1] has extended these approaches using learned semantic features to infer correspondence between images of inter-class or inter-domain objects.

Correspondence learning has also been useful for many tasks in NLP such as pronoun co-reference resolution and commonsense reasoning, both of which rely on modeling correspondences between words to resolve ambiguities [28, 33, 59]. Language translation and understanding, particularly for low-resource languages, benefit from learning word alignments to higher-resource languages [12, 47]. Finally, other types of conditional natural language generation have benefited from learning alignments between words with similar meanings [2, 48].

Screen correspondence is related to many machine-learning driven approaches to identifying correspondence. Our transformer model builds upon many of these approaches by combining visual information and word-alignment techniques to produce screen correspondence. In our evaluation, we compared our system to several baseline techniques from the CV and NLP literature. We show that by incorporating multiple sources of information, our model generates better representations for UI elements, which leads to more accurate correspondence predictions.

## 3 SCREEN CORRESPONDENCE

We define *screen correspondence* as the task of mapping interchangeable UI elements between two UI screens. While matching UI elements between screens may seem simple, it is a complex problem (especially from pixels alone) with many practical use-cases. Previous work relied on mappings to retarget UIs [30], provide help [62, 63], assist design [3], and test GUIs [5] and specifically called for more robust matching to improve performance.

We primarily consider cases where two UI screens are of the same category (*e.g.,* Login or Registration) but from different apps (*i.e.,* intra-class examples). This is challenging because UI element pairs across such screen type pairs may not share similar appearance, text, or position. Instead, the model must reason about the purpose of each element in the context of its screen. To give intuition why even a seemingly simple example is hard, consider two Login screens (Figure 1): one contains *Username* and *Password* fields, while another contains *Email* and *Password* fields. Position and element type information alone is unreliable for matching, since the text fields may have different sizes or appear at different locations on each screen. Appearance information alone is also noisy for matching, since the text fields may have different visual themes. State-of-the-art text encoders, even those trained on phrases, are unreliable (*e.g.,* most text models would produce a higher similarity score for *Username* and *Password* than *Email*).

To detect UI element correspondences between different UI screens, we built a system that *(i)* automatically detects UI elements and text from screenshots, *(ii)* generates multi-modal embeddings for each element, and *(iii)* establishes mappings between individual UI elements with high similarity. Figure 2 shows a high-level overview of our approach.

## 3.1 UI Element Detection

The first stage of our system identifies semantically relevant pieces of information from a UI screenshot, such as UI elements and text. The input is a bitmap and the output is a list of detected UI elements and pieces of text.

We use a pre-trained object detection model from previous work [60] that was trained to recognize UI elements in iOS app screens. The pre-trained model uses the Faster-RCNN architecture [45] and was trained on the AMP dataset [65], which is separate from the main dataset used in this paper. It achieved a class-weighted mAP score of 0.8. We use post-processing procedures, such as score-based thresholding and inter-class non-max suppression (NMS), to improve the quality of the output. Optical character recognition (OCR) is performed using Tesseract [53], an open-source, off-the-shelf OCR software package. We run OCR on regions of the screen that correspond to text elements as detected by our element detector.

## 3.2 UI Element Encoder

Using the elements segmented from the screenshot, we generate representations that encode properties useful for comparison. In our work, we consider relative positioning, element/icon category, visual appearance, and text, properties which we hypothesize to be relevant to element semantics. We used pre-trained models to predict these properties (element detection [60] and icon type ("common icon classifier" from previous work [7])) from screenshots. Note that the pre-trained models were trained on different datasets (*i.e.,* no sample overlap) than the ones used in our paper.

*3.2.1 Modality Representations.* We use off-the-shelf models to generate modality-specific features for each element, then feed their output into a screen transformer model, which combines and learns further associations between them.
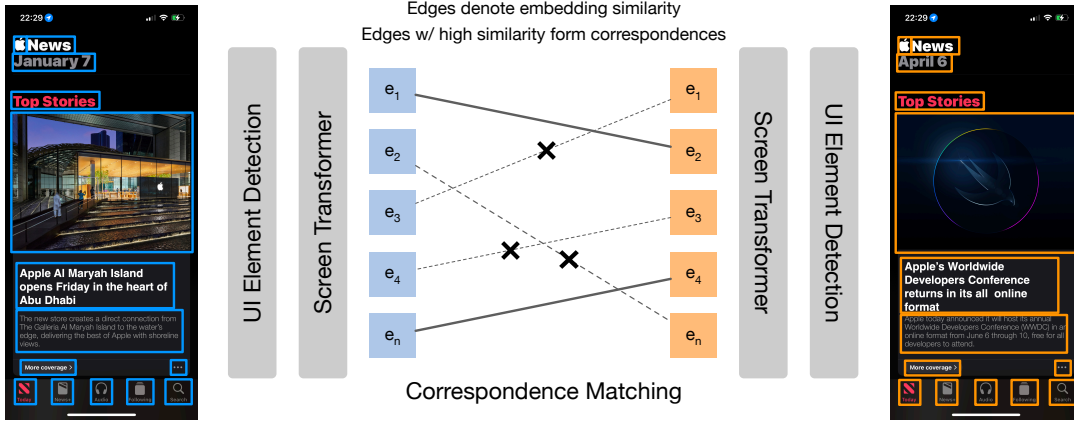
**Figure 2: Overview of our screen correspondence approach. Elements and text from two screenshots are first extracted using UI element detection then featurized using a screen transformer model. Finally, a correspondence between UI elements are generated from element pairs with highly similar embeddings relative to other candidates.**
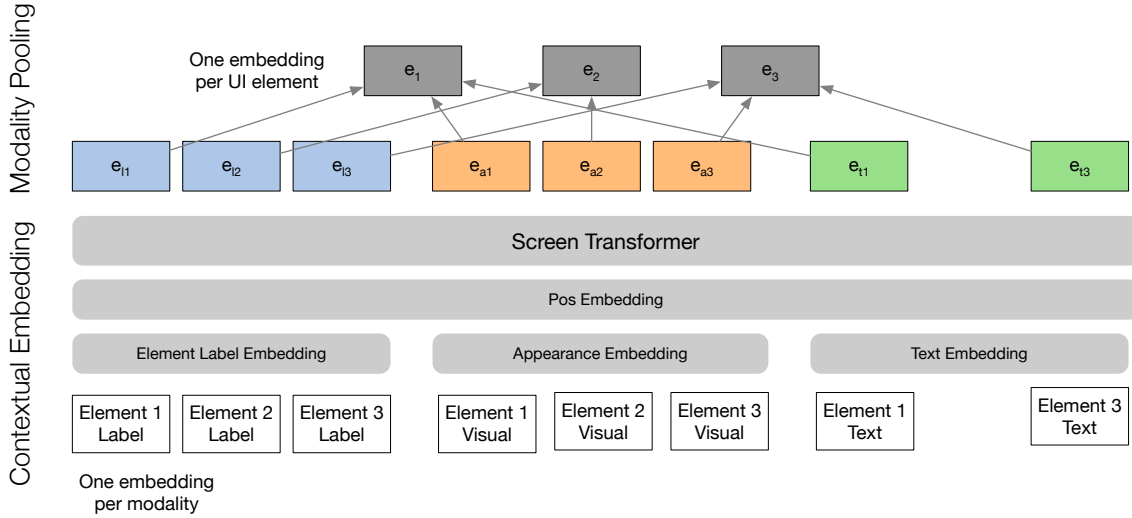


**Figure 3: Architecture diagram of our screen transformer model. Each modality-specific input is treated as separate inputs to our transformer model, which implicitly aligns them based on their positional information. Note that elements may be missing modalities (Element 2 in this example). After the per-modality inputs are processed by our transformer, we generate element embeddings (*i.e.,* one per element) by pooling together outputs corresponding to the same original element.**

*Positional information:* Previous work [14, 18] encoded element position as a simple concatenation of bounding box coordinates. We hypothesized that relative position may be more effective, since UI interactions such as scrolling, text flow, and dynamic content could cause changes in absolute position but have less effect on relative ordering. We adopted a relative positional encoding scheme used to improve the performance of language models [49] that incorporates pairwise distance when calculating the attention score between two elements.

*Element Category:* We categorized elements based on their UI and icon type. Our pre-trained element detector classifies elements into 12 categories, as defined by previous work [65]. Three of these can

be further delineated into sub-categories. We separate the Toggle and Checkbox classes based on their selection state (*e.g.,* Toggle on and Toggle off). We also classified common icon types using a separate pre-trained CNN model [65]. In total, we consider 83 unique categories of elements and represent them as one-hot vectors.

*Visual Appearance:* We featurized regions using the intermediate representations of a proposal-based object detector. Similar approaches have also been used by visual question-answering models, which also need to take into account multiple visual information [? ] Since our UI element detector is based off of a similar proposal-based architecture, we retrieve the activations of the object proposals corresponding to detected elements using the fc6 layer [52].

This approach to featurizing appearance is beneficial since it results in a fixed-size representation for image regions without the need to explicitly resize or crop them.

*Text:* Numerous embedding methods have been developed for representing words, sentences, and documents. Sentence transformers are transformer-based models for encoding variable-length text into an embedding space representative of semantic meaning [44]. Since much of the text on UI screens is relatively short, we use a variant specifically trained to encode phrases [58].

*3.2.2 Transformer Model.* To further enrich and learn associations between the modality-specific element representations, we designed a model that generates a fixed-size embedding for each detected UI element. Our model is based on the transformer architecture (Figure 3), which has been used for UI representation learning [18, 23]. The modifications we described (*e.g.,* relative positioning and appearance features) are aimed at improving performance on the correspondence task.

Because not all elements have the same attributes, *e.g.,* not all UI elements have text, we rely on the transformer's attention mechanism to implicitly align information. Each modality-specific representation with the exception of position is first embedded with a separate linear layer to a common size. Instead of creating one input vector for each element by concatenating features from each modality, we create an input vector for each modality of each element. For example, a login button could result in three input vectors for element category, visual appearance, and text. All inputs are fed into a series of stacked self-attention blocks, which results in one output embedding for each original input vector. Finally, we use pooling to recover the output embeddings associated with each original UI element and compute their mean to incorporate information from all of the modalities.

*3.2.3 Unsupervised Training.* We did not have access to labeled data during the development of our model, so we used unsupervised training to learn its parameters. Masked element prediction is a training objective that requires the model to reconstruct an input that has been corrupted by randomized masking (*i.e.,* replacing a portion of the input with 0's). Previous work [54] has shown that this training objective encourages the model to learn semantically relevant representations since it must learn to associate masked information with other sources of information.

The reconstruction loss was measured separately for all modalities (element category, visual appearance, and text) then added together to obtain the model's total loss. L2-loss was used for reconstruction of visual and text features, and cross-entropy loss was used for reconstruction of element category.

## 3.3 Correspondence Matching

After we used our screen transformer model to featurize UI elements on two screens, we perform a matching procedure to predict correspondences between them. For a pair consisting of a source screen with $M$ elements and a target screen with $N$, we construct a $M \times N$ cost matrix $C \in \mathbb{R}^{MxN}$ to represent correspondence scores. The matching cost $C_{i,j}$ is computed using cosine similarity.

Several approaches have been used to generate correspondences from cost matrices [47]. A simple approach of matching based solely on highest cosine similarity may make suboptimal decisions when one element has more than one likely match. In our final implementation, we formulated correspondence mapping as an optimization problem that finds the assignment between two sets that results in the lowest overall cost [29]. We employ this approach for matching elements between screens, since elements are more likely to be dissimilar. To reduce false positives, we employ additional pre-processing and post-processing steps. Before running the best-cost optimization, we prune unlikely matches from the cost-matrix so that each element only considers its $k$ closest neighbors. Afterwards, we ignore matches where $C_{i,j} < c$. We tuned the values of $k = 5, c = 0.4$ based on manual examination of a small number of examples. This approach is similar to approaches in text decoding models that consider only the top $k$ most likely tokens, which have been shown to generate higher quality output by reducing the effect from low-probability outputs.

## 4 DATASET

We developed and trained our transfer model on two datasets of app screens that were generated by manual crawling of popular mobile apps: CRAWLS and RICO. The CRAWLS dataset, which was used by prior UI modeling work [7], consists of 750,000 iOS app screens from 6,000 apps and was collected by crowdworkers who were instructed to manually explore mobile applications through a remote interface that periodically captured screenshots and additional metadata of the current app screen. The RICO dataset [9] is a publicly available dataset of 72,000 Android screens from 9,700 apps that was also collected by crowdworkers remotely operating devices. We divided each dataset into training (70%), validation (15%), and testing (15%) splits by their crawl ID, which corresponds to which app was crawled.

## 4.1 Evaluation Dataset

While our training algorithm does not depend on labeled data (*i.e.,* unsupervised), we manually collected a small set of labeled examples (900 pairs across 90 screens) from each dataset to evaluate our system.

*4.1.1 Data Collection.* Our evaluation dataset consists of data from 9 types of screens that we hypothesized could have correspondences: Media Player, In-App Purchases, Login, Permission Request, Register, Pre-Login, Pop-up, Search, and Web Views. We initially asked crowdworkers to categorize a set of screenshots outside of the training split based on a criteria for each category. Unlike *app categories*, which might be used to categorize apps (*e.g.,* finance, health, social media), we focused on *screen categories*, since both a health app and a banking app might both contain a login screen that could contain correspondences. For each of our two datasets, we sampled a small number of screens from each category for correspondence labeling (9 categories x 10 screens = 90 screens total). The 9 categories that we chose do not cover all possibilities, but we believe they constitute a reasonable subset. More detailed descriptions and criteria of each category is available in the appendix of this paper.

*4.1.2 Data Labeling.* We created a labeling interface to annotate our small evaluation split. First, a randomly selected element was

shown on a screen, and the interface displayed a prompt asking if the element was likely to appear on other screens of the same type: *"Are elements of similar functionality likely to appear on other Login screens?"* If the user responded "Yes," the application displayed a prompt for a label: *"What is the role of this element in the current screen?"* We built our interface to include auto-complete functionality to encourage labelers to identify correspondence categories that could generalize across screens, *e.g.,* "login button" instead of "button to log into my credit card account." The autocomplete list was pre-populated with 5 choices for each category and was auto-updated with novel descriptions. If a label was provided on the first step, then the user was shown other screens from the same category and asked to select elements with a similar role, if they were present on the screens.

A drawback of this approach is that it is slow, since it requires providing a role description before elements are matched. However, we found the additional consideration of element role is useful for reasoning about correspondences.

## 5 EVALUATION

We evaluate our model against several baselines and ablated versions of our model. Our results show that compared to heuristic and traditional key-point methods used in CV, multi-modal transformer encodings lead to better correspondences. Furthermore, our ablation experiments show that the architectural improvements we made lead to modest performance gains.

### 5.1 Baselines

In this section, we describe the baselines used in our performance evaluation. We focus primarily on other *unsupervised* approaches, since our constraint was that we didn't have any labeled data available for training. Similar supervised approaches exist [30], but they depend on element-level annotations and access to underlying source code (*i.e.,* HTML).

For comparison, we chose a variety of baselines that include keypoint-based methods used for image matching and heuristics such as schema-matching. Our main constraint was that we did not have large quantities of labeled data for supervised machine learning methods, so we selected unsupervised techniques for comparison.

*ORB:* As a review, image correspondence relies on the computation of semantic features from regions of the image. Semantic features, usually invariant to surface-level changes such as translation and scale, are first calculated for small, localized regions of the image. When this process is repeated recursively, the receptive field increases, and globally-aware features can be learned. ORB [46] is a traditional CV approach to generating descriptor features. We first computed ORB descriptors for each screenshot which resulted in numerous keypoints at salient points of the image. Using brute-force matching, keypoints from one image were matched onto keypoints from another image based on descriptor similarity. Finally, to translate keypoint similarity to UI element similarity, we used an object detector to compute the boundaries of UI elements and matched elements based on the number of matching keypoints contained within them.

*Neural Best Buddies:* Neural Best Buddies (NBB) uses the internal representations of a deep CNN to featurize and match image regions. Like ORB, it also generates keypoint descriptors but uses activations from a convolutional neural network (CNN). One advantage that CNN features have over traditional methods is that learned features can better correspond to semantic properties that the network was trained on (e.g., image classification). To run our experiments, we used the code released by the authors of the paper [1]. The original paper focuses on finding correspondences between "natural images" and use a VGG-19 model [51] that was pretrained on ImageNet [10]. Since UI screenshots have different properties than images found in ImageNet, we initially tried to train a CNN model better representative of UIs using an unsupervised autoencoder objective due to the lack of labels in our training set. However, we found the autoencoder model did not produce good outputs, so we report results using the pre-trained ImageNet model.

*Schema-matching Heuristic:* One drawback of keypoint-based methods that we explored is that keypoints are generated using the entire image as input and without knowledge of UI element locations. Schema-matching is an approach that first considers each predicted element as a discrete object, then uses its attributes (*i.e.,* schema) to compare similarity to other candidates. We implemented a heuristic that uses schema matching through incorporating the predicted UI element/icon type by concatenating their one-hot class predictions into a single vector and applying the same best-cost matching algorithm [29]. More sophisticated schema-matching may incorporate additional information, such as UI hierarchy (*e.g.,* an element that belongs in a list should be matched to another element in a list). While possible to predict [60], we did not incorporate hierarchical information since it requires complex techniques for tree matching but expect it would perform similarly to [30], which uses hierarchical information.

*Screen Transformer Ablations:* Our performance evaluation includes ablated variations of our main transformer model. Transformer models allow learning more sophisticated representations of elements through data, which provides advantages over manually-defined schemas. We evaluate several ablated versions of our model to understand the performance impact of our architectural changes. Specifically, the ablated versions of our transformers removes certain components that we hypothesized to improve correspondence matching, such as relative positional embedding, visual appearance information, and text. In addition, we evaluated the Pixel-words transformer [18], which our model is based on, but we adjust the number of element classes, layers, attention heads, and hidden dimensions to be the same as our other models. The Pixel-words transformer also includes a "layout embedding" network which featurizes the layout of UI using a semantic map which is fed into an autoencoder. To summarize, the Pixel-words configuration *(i)* considers categorical and text information, *(ii)* uses absolute positional encodings, and *(iii)* includes an additional layout embedding component.

### 5.2 Results

*5.2.1 Baseline Comparison.* Our evaluation results (Table 1) shows the benefit of our multi-modal model over simpler baselines. We

---

[1]https://github.com/kfiraberman/neural_best_buddies

**Table 1: Performance of our approach and other baselines for screen correspondence. Our approach leads to the best performance, reaching an F1 score 0.61. We also included ablated versions of our model without relative positional embeddings, appearance features, and text features.**

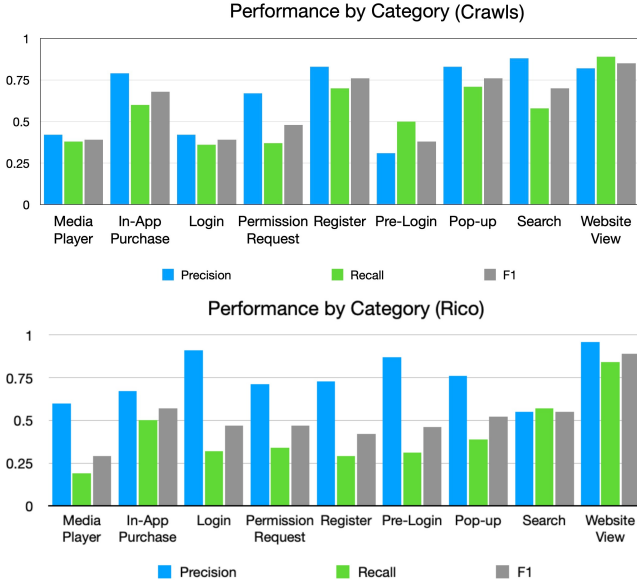|                                  | CRAWLS | | | RICO | | |
| -------------------------------- | ---- | ---- | ---- | ---- | ---- | ---- |
| Model Configuration              | P    | R    | F1   | P    | R    | F1   |
| Screen Trans.                    | 0.66 | 0.57 | 0.61 | 0.83 | 0.41 | 0.55 |
| Screen Tran. (w/o Relative)      | 0.58 | 0.53 | 0.56 | 0.74 | 0.37 | 0.49 |
| Screen Trans. (w/o Appearance)   | 0.74 | 0.44 | 0.55 | 0.77 | 0.38 | 0.51 |
| Screen Trans. (w/o Text)         | 0.66 | 0.63 | 0.59 | 0.77 | 0.37 | 0.50 |
| Screen Trans (Pixel-words)       | 0.70 | 0.49 | 0.58 | 0.83 | 0.22 | 0.35 |
| Heuristic                        | 0.48 | 0.59 | 0.53 | 0.80 | 0.32 | 0.45 |
| ORB                              | 0.25 | 0.17 | 0.20 | 0.63 | 0.21 | 0.31 |
| NBB                              | 0.22 | 0.15 | 0.18 | 0.58 | 0.17 | 0.26 |



**Figure 4: Performance across different categories in the Crawls (Top) and Rico (bottom) datasets using the full Screen Transformer model configuration. The average classification performance was F1=0.61 on Crawls and F1=0.55 on Rico.**

employ standard classification metrics to measure the accuracy of element-to-element correspondences generated by our model. Since elements in our evaluation dataset are labeled using their ground truth bounding boxes instead of our element detector's predictions, we first match predicted detections to ground truth elements using the best Intersection-over-Union (IoU) score. Due to our labeling procedure where one element is highlighted at a time, the examples in our evaluation set were only partially labeled, meaning that screens contained only a randomly sampled subset of all possible corresponding pairs. Our best model configuration reaches an F1 score of 0.61. Screens in our dataset contained an average of around 20 elements, so correct correspondence required

finding the best out of match out of all possible candidates. The *ORB* and *NBB* baselines are based on keypoint-based matching, which is commonly used in CV to compare images. Among them, ORB performs the best, achieving higher correspondence accuracy but performed poorly due to low recall. One possible reason is that keypoints are generated at visually salient locations of the image, such as edges and corners, and without any knowledge of where UI elements are. Thus, some UI elements may not contain many keypoints within them, reducing the quality of matches. The schema-matching heuristic performed substantially better than keypoint-based methods and reached high recall by directly using the outputs of pre-existing models (*i.e.,* element detection, icon classification). Precision was lower, possibly due to the difficulty of accurately matching ambiguous elements without knowledge of additional context.

Our ablation experiments revealed that our modifications to the base transformer architecture led to modest improvements in terms of F1 score but also had other consequences for precision and recall. For example, our model trained without appearance information was the lowest performing variation but reached the highest precision score. We attribute these variations to the information encoded in each modality and may warrant different configurations based on intended use-case.

*5.2.2  Performance across UI Categories.* Figure 4 provides a more in-depth breakdown correspondence by UI category. Our model achieved the best performance on the *Website View* and *In-App Purchase* categories and the worst performance on the *Media Player* and *Pre-Login* categories.

One major source of error for our model was the presence of sub-categories within our dataset. For example, we manually examined examples from the *Pre-Login* and *Login* categories, which received relatively low performance. We discovered a considerable difference between apps that used different authentication providers, such as OAuth and Single-Sign-On (SSO). For example, a "traditional" login screen might include text fields for entering a username and password, but an app using a SSO provide (*e.g.,* Sign in with Apple) might only contain a button without any text fields. We found that there was also variance within media player screens – video players and music players had significant differences and some media players were full screen while others were not. Since our correspondence model uses contextual information (*i.e.,* information from other elements on the same screen) and relative positional encoding, this could significantly affect the computed representation. One strategy to address this is the formation of sub-categories with a more consistent set of elements *e.g.,* creating separate categories for traditional login screens and those with other types of authentication.

*5.2.3  Performance across Datasets.* We evaluated all models and baselines on both the Crawls and Rico dataset. Overall performance between the two datasets were similar, although the Rico models performed slightly worse (F1=0.55) than ones trained on Crawls (F1=0.61). One possible reason for the performance discrepancy is that Crawls is an order of magnitude larger and the model was exposed to more variation during training time, which is beneficial for unsupervised training techniques. While the full transformer model is the best-performing configuration for both

**Table 2: Performance of our approach and other baselines screen correspondence for *same-screen* pairs in the Crawls dataset. Many configurations, including our model, reach a maximum F1 score of 0.76. We attribute labeling noise and the IoU element matching process used to assign predicted element locations to ground-truth boxes.**

| Model Configuration | P | R | F1 |
|---|---|---|---|
| Screen Transformer | 0.85 | 0.68 | 0.76 |
| Screen Transformer (w/o Relative) | 0.86 | 0.68 | 0.76 |
| Screen Transformer (w/o Appearance) | 0.87 | 0.66 | 0.75 |
| Screen Transformer (w/o Text) | 0.85 | 0.68 | 0.76 |
| Screen Transformer (Pixel-words) | 0.88 | 0.66 | 0.76 |
| Heuristic | 0.87 | 0.66 | 0.75 |
| ORB | 0.78 | 0.48 | 0.59 |
| NBB | 0.53 | 0.25 | 0.34 |

datasets, the relative performance ablated models were affected differently. Notably, the Rico models without text and the Pixel-Words model performed much worse, suggesting that its evaluation set may have contained more text-heavy screens.

*5.2.4 Correspondence between Same-screen Pairs.* In addition to evaluating our models on screens from different related apps (*i.e., intra-class* pairs), we also investigated performance on *same-screen* pairs. Same-screen correspondence is useful for identifying the same UI element across multiple versions of the same screen. For example, an app's appearance may change following an update or from dynamically updated content (*e.g.,* a news page loads content from a remote source). Following prior work [14], we consider two screenshots to be the "same" if they represent different instances of the same underlying implementation, possibly with significantly different appearance. Correspondence mapping can help guide automated systems such as crawlers to behave more consistently in these situations. We randomly selected screen groups with the same app ID as those in the testing split of our Crawls dataset, then randomly sampled two screenshots from each group, resulting in 888 total pairs. Upon manual inspection, we found that some of the sampled pairs had only minimal visual changes. To filter out "easy pairs," we constructed a heuristic that attempted to match elements based only on bounding box location. If all elements in a pair were successfully matched, we discarded the example, since it meant that no significant dynamic change (*e.g.,* scrolling, dynamic content) occurred. After this process, the final dataset contains 607 examples. We did not repeat this for the Rico dataset because the authors applied a heuristic to filter out repeated views of the same screen [9].

Our observations and performance results (Table 2) show that *same-screen* correspondence is generally higher. Since *same-screen* pairs are usually more visually similar, the model can rely more heavily on surface-level features and in many cases perform direct matching, such as looking for recurring text. Many configurations, including our model, reached a maximum F1 score of 0.76. Errors from labeling noise and IoU element matching (*e.g.,* matching ground truth bounding boxes to predictions) may have established

an effective ceiling, since our element detection model introduced errors (has a class-weighted mAP score of 0.8).

## 6  EXAMPLE APPLICATIONS

We describe three example applications that show the utility of screen correspondence to human and machine understanding of UI functionality. Generating and transferring a type of instructional overlay called coach marks can help users navigate unfamiliar UIs by mapping them to previously encountered ones of the same class. UI search is useful for app designers to find how concepts are expressed across apps (*e.g.,* what are different ways of expressing a search intent?). Finally, automated GUI testing can be made more robust by accounting for variations in visual presentation between different app versions without requiring platform-specific APIs or metadata. These example applications are not meant to be novel, but we believe they show that accurate screen correspondence allows many existing systems to work under a wider range of conditions, *e.g.,* using pixel data alone or improved robustness to dynamic visual changes.

### 6.1  Instructional Overlays

We used our model to improve users' understanding of complex or newly installed apps by creating an infrastructure that could be used to crowdsource coach marks for apps. Coach marks are instructional overlays that are sometimes shown to provide assistance to users when an app is first launched, and can be helpful for exposing UI functionality. While it is possible to automatically generate natural language for describing screens [57] and widgets [37] using deep models, they are often affected by surface-level appearance and may be prone to producing generic outputs [37]. Building a model that produces natural language also introduces significant complexity that can be similar achieved with a correspondence mapping. A better approach might be to crowdsource users [43] or developers to write coach marks for screens in a subset of apps, and then apply our screen correspondence technology to map these coach marks onto a much larger set of screens with similar purposes. This idea builds on the template-based matching scheme of Yeh et al. [62] for generating contextual help, and expands their idea beyond *same-screen* applications to also *intra-class* usage.

We applied our model's intra-class correspondence capabilities to automatically transfer annotations from one screen to another related app of the same category (Figure 5). We first populated a small database of instructional text for elements from app screens in one of the categories from our evaluation data. In a real implementation of this system, an interface would be created to allow users to author new instructional text for screenshots that they upload. Each screen in the database was associated with its featurized elements as a key, and each instruction in the database was associated with its element. Our current prototype is a proof-of-concept implementation where the user can upload a screenshot image file through a web interface. On the uploaded screen, we perform a nearest-neighbor search to retrieve the screens in our database that are most similar. If the distance is sufficiently close, we run our screen correspondence matching, which also returns a "matching cost." If enough matches are discovered and the matching cost is
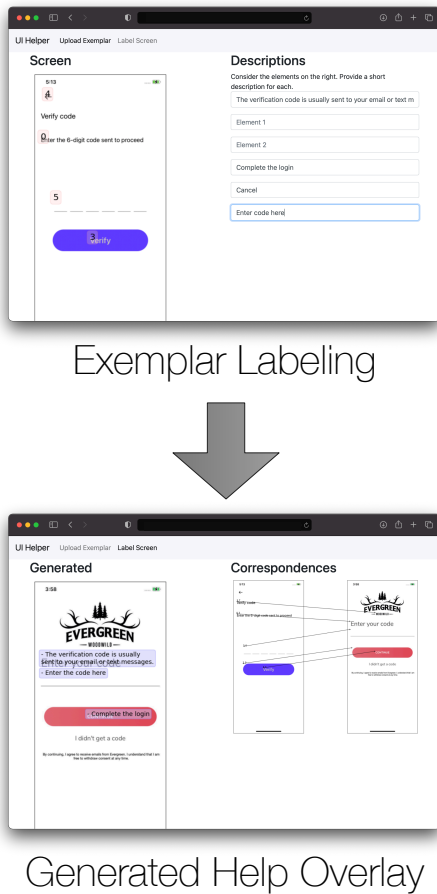
Exemplar Labeling

Generated Help Overlay

**Figure 5: Coach marks are useful for uncovering functionality in apps. High-quality natural language descriptions of UI components can be difficult to generate, so we curated a small number of labeled examples from different app categories. Element descriptions from this labeled set are transferred onto unseen app screens of the same type using the correspondence mapping. Depending on the use-case, the exemplar can be manually provided (*e.g.,* developer wishes to label many similar screens at once) or automatically retrieved (*e.g.,* a help-generation app uses a separate classifier to find an exemplar from a database of labeled screens.)**

below a heuristically set threshold, we directly render the annotations to the screenshot using image drawing APIs and display the annotated image.

In a complete implementation, the matching and rendering algorithms would be built into a mobile operating system and run on the user's mobile device so that it would not require the user to exit their current app to use our tool. In the future, we plan to improve the user experience and investigate ways that these overlays could be surfaced contextually.

In this example, we show that accurate intra-class screen correspondence can facilitate transferring coach marks, which can help

users discover new app functionality and documentation. Other possible applications of screen correspondence to improving end-user usage include transferring more types of accessiblity meta-data, example-based re-targeting of UIs [30] and using input redirection techniques to improve the accessibility of UI components [66].

## 6.2 UI Element Search Engine

UI search can help app designers find how concepts are expressed across apps and provide example starting points when designing a new app. Previous work indexed databases of UI screens using visual properties [3], structural properties [60], sketches [26]. We focus specifically on returning relevant UI elements instead of screens, and leverage our model's *intraclass* matching abilities to improve the search process.

We integrated our screen correspondence model into a UI search engine to support tag-based search and exemplar-based refinement. The implementation of our UI search engine is a web app that indexed elements from 130,000 UI screens using a variety of metadata, including detected element classes, icon types, and text, which are stored in a database. Our app features a search page, where users can first perform an initial search by entering text or tags in a search bar. Results are returned based on matching attributes found in the property database. Matching elements are shown in the context of their app screen and highlighted with a bounding box. When a result is selected, users are brought to the element inspector page, where users can examine the properties of all elements on the screen.

One limitation of tag-based search is that it is difficult to specify target properties that do not belong to the pre-defined set of tags. For example, a "plus" icon displayed on the top or bottom of a list may indicate adding *to* the list while a "plus" icon displayed next to a list item is more likely to representing adding the item *from* the list. It would be difficult to disambiguate between these cases as they share the same tag. Thus, we used our correspondence model to enable exemplar-based search refinement, which allows users to "narrow in" on more specific results. To enable this functionality, we computed embeddings for UI elements in our database and stored this information into a vector data store which supports fast approximate nearest-neighbor search. We added a "search for similar items" button on the element inspector page, which finds results with a high similarity to the target element according to the cosine similarity metric. Figure 6 shows an example flow of our UI element search engine.

## 6.3 Automated GUI Testing

Finally, we used our model to improve the robustness of automated GUI tests using our model's *same-screen* matching capability. Automated testing is useful for ensuring the quality of GUIs. Specifically, visual-based methods can be employed in these systems to search for targets based on their rendered appearance, which allows for easier authoring of testing scripts and reduces the dependency of testing frameworks on specific UI toolkits [5]. However, strong reliance on visual similarity may lead to failures caused by change in visual style, such as updated application theme or icons [5].

In such applications, the quality of UI element matching is important for automated GUI testing because poor matching capability
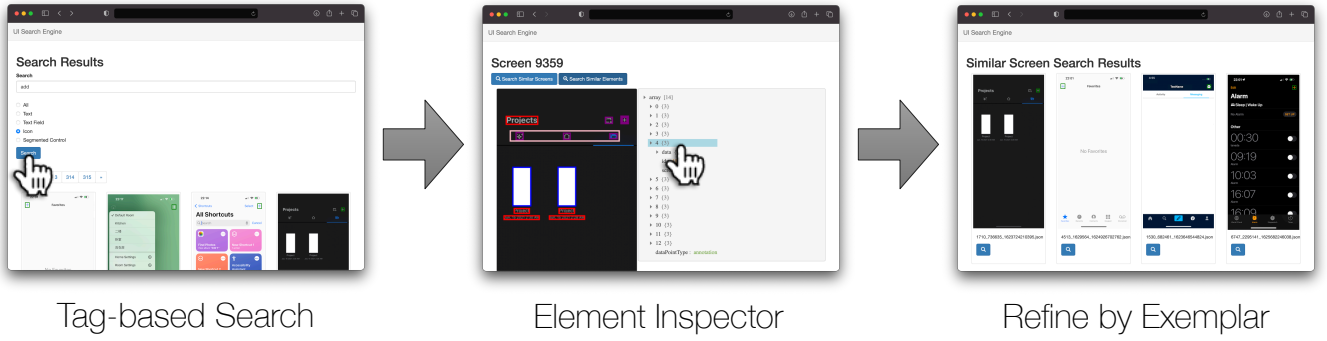
Tag-based Search             Element Inspector             Refine by Exemplar

**Figure 6: An example usage flow of our UI element search engine. The user first searches for icon elements that contain the "add" tag. The results page shows UI screens with a matching element highlighted (Left). The user selects a result screen where an add button is placed on the top right of the screen. The inspection page provides details about element info and allows searching for similar elements (Center). Another search query is run using the embedded element of interest. The new results are similar to the query in that they are all located at the top right of the screen and they appear to be used for adding items to a gallery (Right). This example shows how designers can start a search using natural language or tag-based queries then refine the results based on exemplars.**



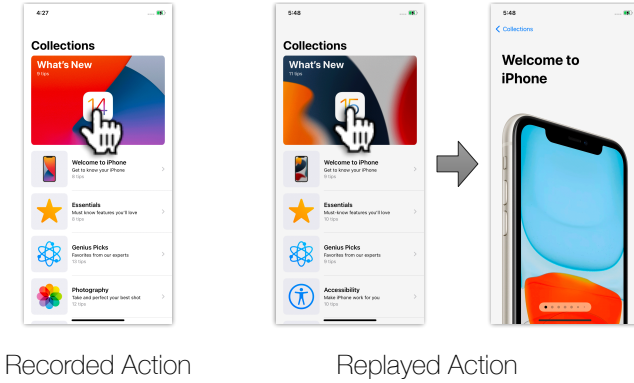Recorded Action          Replayed Action

**Figure 7: Automated UI testing techniques execute an interaction trace (either manually pre-defined or automatically generated) to detect functional regressions, visual regressions, and other unexpected behavior. Updated versions of apps may lead to small changes in layout and visual appearance and knowledge of same-screen correspondence can improve the consistency and robustness of tests. This example shows a automated application performing a previously recorded action, despite the target's appearance change.**

can lead to a failure to replicate recorded interaction traces in a scripted testing scenario, and repeated visits to the same screens in a random crawler stress test example. As shown by previous work on screen similarity [14], methods that rely heavily on surface-level appearance may have high precision but low recall due to possible variations between UIs. We applied our screen correspondence model to improve the robustness of these matches.

We built a prototype system that interacts with remotely connected smartphone devices through a VNC interface. Our software sends commands through this interface to simulate interactions, such as clicking and swiping. We also include a "recording" mode

that allows a tester to record an interaction trace, during which all of the screenshots and interactions are saved. When replaying the interaction trace, the saved screenshots and interacted elements are used to match the current state of the VNC output. Specifically, for each step in the saved trace, we identify the UI element with which the tester interacted, such as the button that was pressed. Then, on the live VNC view, we find the corresponding element and apply the recorded interaction to it, similar to previous work on tutorial consumption [68]. Figure 7 illustrates how our automated tester navigates an app where the appearance of a target element has changed. Used in conjunction with traditional template-matching techniques, which offer high precision but low recall, correspondence matching can help improve the overall performance of automated testers.

## 7 LIMITATIONS & FUTURE WORK

Our evaluation shows that correspondences can be automatically identified through machine learning and matching approaches. Some types of screens are more likely to have correspondences detectable by our system (*e.g.,* Website Views and In-App purchases) than others (*e.g.,* media players). The required accuracy level depends largely on the final application, since different use-case since different performance attributes. For example, using correspondences to generate contextual help (instructional overlays) may result in a better experience if only very confident matches are used, as incorrect instructions can lead to confusion and frustration from the user. GUI testing and crawling is less tolerant to mistakes, since an incorrect action can make it impossible to access the rest of an application. On the other hand, UI design search is more forgiving, since it can provide value if most of the returned elements are correct (does not need to be the top choice). Our current evaluation does not account for the requirements of down-stream applications, although based on the example applications we implemented, we found them to provide acceptable performance. We plan to further evaluate our system in down-stream applications.

A limitation of our current experiments is that they focus only on mobile UIs that belong to a set of 9 categories that we identified. These 9 categories do not cover all possibilities of app screens, but they cover a considerable subset. Our model is likely to perform better for complex app screens if given a small amount of annotations to fine-tune on. Moreover, since we only use pixel information as input to our model, we believe that our approach is likely to generalize well to other types of graphical UIs that also represent their output as pixels. In the future, we aim to replicate our experiments on other types of graphical UIs with varying screen sizes and shapes.

We see several opportunities to improve the performance of our system. Since our system relies on several individual components, it may be useful to quantify the performance of each separately. We used a pre-trained element detector model that produced noisy output for the correspondence matching. Previous work [60] has shown that element detectors perform poorly on more complex screens due to the increased number of elements and sometimes miss smaller elements. Future work could investigate a screen correspondence system that uses a more accurate element detector model or accepts manual annotations as input. More advanced matching techniques can also be employed, such those that consider multi-scale correspondence, which first process smaller sub-regions before merging their predictions globally. Separately, prior work on image correspondence [27] has shown improved performance by scaling images during training and inference. A similar idea could be applied to UIs by first predicting their UI hierarchy [60] and generating mappings for groups of elements. Our model could also use different unsupervised pre-training objectives to help it build better representations of UI elements for our matching task [28, 50].

Our work focuses on mapping interchangeable elements with similar functionality between UI screens, however there are other relationships that can be modeled. Categorization of different relations in language analogies [19, 40] show that antonym, categorical, and functional connections can enrich the expressiveness of language and rhetoric. We plan to focus future modeling efforts on identifying and inferring a wider range of similar relationships that exist in UIs.

Finally, our work explores inferring UI functionality from a single previously encountered example, yet we believe our approach may extend to multiple examples [17]. For example, non-parametric machine learning methods such as the k-nearest neighbors algorithm often benefit from considering more than one example at a time.

## 8 CONCLUSION

In this paper, we explore *screen correspondence* as a machine learning technique for inferring UI functionality by directly leveraging previously encountered examples. We describe our model architecture and training procedure that incorporates information about UI semantics, appearance, and text when generating correspondence mappings between screenshots. In a comprehensive evaluation with strong baselines, we show that our approach outperforms correspondence algorithms by leveraging multiple information sources found in UIs. Finally, we show how three example applications of

screen correspondence: *(i)* transferring coach marks from related apps, *(ii)* UI element search, and *(iii)* automated GUI testing. Broadly, our work demonstrates the feasibility of learning UI semantics by mapping to prior examples.

## A MODEL HYPERPARAMETERS

| Model | Hyperparameter | Value |
|---|---|---|
| Screen Transformer | optimizer | Adam |
| | learning rate | 1e-4 |
| | weight decay | 1e-5 |
| | dropout | 0.25 |
| | hidden size | 256 |
| | num layers | 4 |
| | num heads | 4 |

We trained our models with early stopping and stopped training when validation loss did not improve for 10 epochs. We implemented our model using the PyTorch [42] and PyTorch Lightning [13] frameworks.

## B UI CATEGORY CRITERIA

We collected a small dataset of 9 screen categories for evaluation of our model's *intra-class* correspondence capabilities. We used the following guidelines to categorize apps.

- *Media Player* - A screen that allows users to play media content such as music or video. Usually contains controls for adjusting playback, volume, and sharing.
- *In-App Purchase* - A screen that asks users to make a purchase for a subscription or to access some part of an app. Usually contains buttons for making the purchase, dismissing the screen, or signing up for a trial.
- *Login* - A screen that asks users to log into an app or service. It may contain fields for entering username and password or buttons for third party authentication services.
- *Permission Request* - A screen that asks users to enable some permission, which are usually associated with security settings such as location or camera access.
- *Register* - A screen that asks the user to create an account. May contain a form to register or buttons for third party authentication providers.
- *Pre-Login* - A screen that contains controls to access other parts of the app either by logging in or registering for an account. This usually comes before the login page.
- *Pop-up* - A screen with a pop-up or dialog model that is displayed over other app content. Pop-ups may contain controls for accepting or dismissing it. For pop-ups that ask for permission or purchases, see other categories.
- *Search* - A screen for entering and submitting a search query. May include a search bar and filtering controls.
- *Website View* - A screen where an app opens an external website. May contain a URL bar and forward/backward controls.

# REFERENCES

[1] Kfir Aberman, Jing Liao, Mingyi Shi, Dani Lischinski, Baoquan Chen, and Daniel Cohen-Or. 2018. Neural best-buddies: Sparse cross-domain correspondence. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–14.

[2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).

[3] Sara Bunian, Kai Li, Chaima Jemmali, Casper Harteveld, 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*. 1–14.

[4] Joel Chan, Joseph Chee Chang, Tom Hope, Dafna Shahaf, and Aniket Kittur. 2018. Solvent: A mixed initiative system for finding analogies between research papers. *Proceedings of the ACM on Human-Computer Interaction* 2, CSCW (2018), 1–21.

[5] Tsung-Hsiang Chang, Tom Yeh, and Robert C Miller. 2010. GUI testing using computer vision. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 1535–1544.

[6] Jieshan Chen, Chunyang Chen, Zhenchang Xing, Xiwei Xu, Liming Zhu, Guoqiang Li, and Jinshui Wang. 2020. Unblind your apps: Predicting natural-language labels for mobile gui components by deep learning. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 322–334.

[7] Jieshan Chen, Amanda Swearngin, Jason Wu, Titus Barik, Jeffrey Nichols, and Xiaoyi Zhang. 2022. Towards Complete Icon Labeling in Mobile Applications. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–14.

[8] 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 *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 596–607.

[9] Biplab Deka, Zifeng Huang, Chad Franzen, Joshua Hibschman, Daniel Afergan, Yang Li, Jeffrey Nichols, and Ranjitha Kumar. 2017. Rico: A mobile app dataset for building data-driven design applications. In *Proceedings of the 30th Annual ACM Symposium on User Interface Software and Technology*. 845–854.

[10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 248–255.

[11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

[12] Chris Dyer, Victor Chahuneau, and Noah A Smith. 2013. A simple, fast, and effective reparameterization of ibm model 2. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 644–648.

[13] WA Falcon and .al. 2019. PyTorch Lightning. *GitHub. Note: https://github.com/PyTorchLightning/pytorch-lightning* 3 (2019).

[14] Shirin Feiz, Jason Wu, Xiaoyi Zhang, Amanda Swearngin, Titus Barik, and Jeffrey Nichols. 2022. Understanding Screen Relationships from Screenshots of Smartphone Applications. In *Proceedings of the 27th Annual Conference on Intelligent User Interfaces*. 1–12.

[15] Sidong Feng, Suyu Ma, Jinzhong Yu, Chunyang Chen, TingTing Zhou, and Yankun Zhen. 2021. Auto-icon: An automated code generation tool for icon designs assisting in ui development. In *26th International Conference on Intelligent User Interfaces*. 59–69.

[16] Martin A Fischler and Robert C Bolles. 1981. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM* 24, 6 (1981), 381–395.

[17] Huazhu Fu, Xiaochun Cao, and Zhuowen Tu. 2013. Cluster-based co-saliency detection. *IEEE Transactions on Image Processing* 22, 10 (2013), 3766–3778.

[18] Jingwen Fu, Xiaoyi Zhang, Yuwang Wang, Wenjun Zeng, Sam Yang, and Grayson Hilliard. 2021. Understanding Mobile GUI: from Pixel-Words to Screen-Sentences. *arXiv preprint arXiv:2105.11941* (2021).

[19] Anna Gladkova, Aleksandr Drozd, and Satoshi Matsuoka. 2016. Analogy-based detection of morphological and semantic relations with word embeddings: what works and what doesn't.. In *Proceedings of the NAACL Student Research Workshop*. 8–15.

[20] Kelvin Guu, Tatsunori B Hashimoto, Yonatan Oren, and Percy Liang. 2018. Generating sentences by editing prototypes. *Transactions of the Association for Computational Linguistics* 6 (2018), 437–450.

[21] Richard Hartley and Andrew Zisserman. 2003. *Multiple view geometry in computer vision*. Cambridge university press.

[22] Junxian He, Taylor Berg-Kirkpatrick, and Graham Neubig. 2020. Learning sparse prototypes for text generation. *Advances in Neural Information Processing Systems* 33 (2020), 14724–14735.

[23] Zecheng He, Srinivas Sunkara, Xiaoxue Zang, Ying Xu, Lijuan Liu, Nevan Wichers, Gabriel Schubiner, Ruby Lee, Jindong Chen, and Blaise Aguera y Arcas. 2020. ActionBert: Leveraging User Actions for Semantic Understanding of User Interfaces. *arXiv preprint arXiv:2012.12350* (2020).

[24] Aaron Hertzmann, Charles E Jacobs, Nuria Oliver, Brian Curless, and David H Salesin. 2001. Image analogies. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*. 327–340.

[25] Berthold KP Horn and Brian G Schunck. 1981. Determining optical flow. *Artificial intelligence* 17, 1-3 (1981), 185–203.

[26] Forrest Huang, John F Canny, and Jeffrey Nichols. 2019. Swire: Sketch-based user interface retrieval. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–10.

[27] Wei Jiang, Eduard Trulls, Jan Hosang, Andrea Tagliasacchi, and Kwang Moo Yi. 2021. Cotr: Correspondence transformer for matching across images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 6207–6217.

[28] Vid Kocijan, Ana-Maria Cretu, Oana-Maria Camburu, Yordan Yordanov, and Thomas Lukasiewicz. 2019. A Surprisingly Robust Trick for the Winograd Schema Challenge. In *ACL*.

[29] Harold W Kuhn. 1955. The Hungarian method for the assignment problem. *Naval research logistics quarterly* 2, 1-2 (1955), 83–97.

[30] Ranjitha Kumar, Jerry O Talton, Salman Ahmad, and Scott R Klemmer. 2011. Bricolage: example-based retargeting for web design. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 2197–2206.

[31] Seokjun Lee, Rhan Ha, and Hojung Cha. 2018. Click sequence prediction in Android mobile applications. *IEEE Transactions on Human-Machine Systems* 49, 3 (2018), 278–289.

[32] Luis A Leiva, Asutosh Hota, and Antti Oulasvirta. 2020. Enrico: A dataset for topic modeling of mobile ui designs. In *22nd International Conference on Human-Computer Interaction with Mobile Devices and Services*. 1–4.

[33] Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In *Thirteenth international conference on the principles of knowledge representation and reasoning*.

[34] Toby Jia-Jun Li, Amos Azaria, and Brad A Myers. 2017. SUGILITE: creating multimodal smartphone automation by demonstration. In *Proceedings of the 2017 CHI conference on human factors in computing systems*. 6038–6049.

[35] Toby Jia-Jun Li, Lindsay Popowski, Tom Mitchell, and Brad A Myers. 2021. Screen2Vec: Semantic Embedding of GUI Screens and GUI Components. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–15.

[36] Yang Li, Jiacong He, Xin Zhou, Yuan Zhang, and Jason Baldridge. 2020. Mapping natural language instructions to mobile UI action sequences. *arXiv preprint arXiv:2005.03776* (2020).

[37] 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. *arXiv preprint arXiv:2010.04295* (2020).

[38] Ce Liu, Jenny Yuen, and Antonio Torralba. 2010. Sift flow: Dense correspondence across scenes and its applications. *IEEE transactions on pattern analysis and machine intelligence* 33, 5 (2010), 978–994.

[39] Thomas F Liu, Mark Craft, Jason Situ, Ersin Yumer, Radomir Mech, and Ranjitha Kumar. 2018. Learning design semantics for mobile apps. In *Proceedings of the 31st Annual ACM Symposium on User Interface Software and Technology*. 569–579.

[40] Hongjing Lu, Ying Nian Wu, and Keith J Holyoak. 2019. Emergence of analogy from relation learning. *Proceedings of the National Academy of Sciences* 116, 10 (2019), 4176–4181.

[41] Forough Mehralian, Navid Salehnamadi, and Sam Malek. 2021. Data-driven accessibility repair revisited: on the effectiveness of generating labels for icons in Android apps. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 107–118.

[42] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019).

[43] Vidya Ramesh, Charlie Hsu, Maneesh Agrawala, and Björn Hartmann. 2011. ShowMeHow: translating user interface instructions between applications. In *Proceedings of the 24th annual ACM symposium on User interface software and technology*. 127–134.

[44] Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084* (2019).

[45] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems* 28 (2015), 91–99.

[46] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. 2011. ORB: An efficient alternative to SIFT or SURF. In *2011 International conference on computer vision*. Ieee, 2564–2571.

[47] Masoud Jalili Sabet, Philipp Dufter, François Yvon, and Hinrich Schütze. 2020. SimAlign: High quality word alignments without parallel training data using static and contextualized embeddings. *arXiv preprint arXiv:2004.08728* (2020).

[48] Abigail See, Peter J Liu, and Christopher D Manning. 2017. Get to the point: Summarization with pointer-generator networks. *arXiv preprint arXiv:1704.04368* (2017).

[49] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. 2018. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155* (2018).

[50] Ming Shen, Pratyay Banerjee, and Chitta Baral. 2021. Unsupervised Pronoun Resolution via Masked Noun-Phrase Prediction. *arXiv preprint arXiv:2105.12392* (2021).

[51] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).

[52] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards vqa models that can read. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8317–8326.

[53] Ray Smith. 2007. An overview of the Tesseract OCR engine. In *Ninth international conference on document analysis and recognition (ICDAR 2007)*, Vol. 2. IEEE, 629–633.

[54] Hao Tan and Mohit Bansal. 2019. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490* (2019).

[55] Michael Tomasello. 2005. *Constructing a language: A usage-based theory of language acquisition.* Harvard university press.

[56] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. 2016. Matching networks for one shot learning. *Advances in neural information processing systems* 29 (2016).

[57] Bryan Wang, Gang Li, Xin Zhou, Zhourong Chen, Tovi Grossman, and Yang Li. 2021. Screen2words: Automatic mobile UI summarization with multimodal learning. In *The 34th Annual ACM Symposium on User Interface Software and Technology*. 498–510.

[58] Shufan Wang, Laure Thompson, and Mohit Iyyer. 2021. Phrase-bert: Improved phrase embeddings from bert with an application to corpus exploration. *arXiv preprint arXiv:2109.06304* (2021).

[59] Jason Wu, Karan Ahuja, Richard Li, Victor Chen, and Jeffrey Bigham. 2019. ScratchThat: Supporting Command-Agnostic Speech Repair in Voice-Driven Assistants. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 3, 2 (2019), 1–17.

[60] Jason Wu, Xiaoyi Zhang, Jeff Nichols, and Jeffrey P Bigham. 2021. Screen Parsing: Towards Reverse Engineering of UI Models from Screenshots. In *The 34th Annual ACM Symposium on User Interface Software and Technology*. 470–483.

[61] Rahulkrishna Yandrapally, Andrea Stocco, and Ali Mesbah. 2020. Near-duplicate detection in web app model inference. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 186–197.

[62] Tom Yeh, Tsung-Hsiang Chang, Bo Xie, Greg Walsh, Ivan Watkins, Krist Wong-suphasawat, Man Huang, Larry S Davis, and Benjamin B Bederson. 2011. Creating contextual help for GUIs using screenshots. In *Proceedings of the 24th annual ACM symposium on User interface software and technology*. 145–154.

[63] Ja Eun Yu and Debaleena Chattopadhyay. 2020. "Maps are hard for me": Identifying How Older Adults Struggle with Mobile Maps. In *the 22nd international ACM SIGACCESS conference on computers and accessibility*. 1–8.

[64] Xiaoxue Zang, Ying Xu, and Jindong Chen. 2021. Multimodal Icon Annotation For Mobile Applications. In *Proceedings of the 23rd International Conference on Mobile Human-Computer Interaction*. 1–11.

[65] Xiaoyi Zhang, Lilian de Greef, Amanda Swearngin, Samuel White, Kyle Murray, Lisa Yu, Qi Shan, Jeffrey Nichols, Jason Wu, Chris Fleizach, et al. 2021. Screen Recognition: Creating Accessibility Metadata for Mobile Applications from Pixels. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–15.

[66] Xiaoyi Zhang, Anne Spencer Ross, Anat Caspi, James Fogarty, and Jacob O Wobbrock. 2017. Interaction proxies for runtime repair and enhancement of mobile application accessibility. In *Proceedings of the 2017 CHI conference on human factors in computing systems*. 6024–6037.

[67] Xiaoyi Zhang, Anne Spencer Ross, and James Fogarty. 2018. Robust annotation of mobile application interfaces in methods for accessibility repair and enhancement. In *Proceedings of the 31st Annual ACM Symposium on User Interface Software and Technology*. 609–621.

[68] Mingyuan Zhong, Gang Li, Peggy Chi, and Yang Li. 2021. HelpViz: Automatic Generation of Contextual Visual Mobile Tutorials from Text-Based Instructions. In *The 34th Annual ACM Symposium on User Interface Software and Technology*. 1144–1153.

[69] Xin Zhou and Yang Li. 2021. Large-Scale Modeling of Mobile User Click Behaviors Using Deep Learning. In *Fifteenth ACM Conference on Recommender Systems*. 473–483.