

DeclarUI: Bridging Design and Development with Automated Declarative UI Code Generation

TING ZHOU*, Huazhong University of Science and Technology, China

YANJIE ZHAO*, Huazhong University of Science and Technology, China

XINYI HOU, Huazhong University of Science and Technology, China

XIAOYU SUN, Australian National University, Australia

KAI CHEN[†], Huazhong University of Science and Technology, China

HAOYU WANG[†], Huazhong University of Science and Technology, China

Declarative UI frameworks have gained widespread adoption in mobile app development, offering benefits such as improved code readability and easier maintenance. Despite these advantages, the process of translating UI designs into functional code remains challenging and time-consuming. Recent advancements in multimodal large language models (MLLMs) have shown promise in directly generating mobile app code from user interface (UI) designs. However, the direct application of MLLMs to this task is limited by challenges in accurately recognizing UI components and comprehensively capturing interaction logic.

To address these challenges, we propose DECLARUI, an automated approach that synergizes computer vision (CV), MLLMs, and iterative compiler-driven optimization to generate and refine declarative UI code from designs. DECLARUI enhances visual fidelity, functional completeness, and code quality through precise component segmentation, Page Transition Graphs (PTGs) for modeling complex inter-page relationships, and iterative optimization. In our evaluation, DECLARUI outperforms baselines on React Native, a widely adopted declarative UI framework, achieving a 96.8% PTG coverage rate and a 98% compilation success rate. Notably, DECLARUI demonstrates significant improvements over state-of-the-art MLLMs, with a 123% increase in PTG coverage rate, up to 55% enhancement in visual similarity scores, and a 29% boost in compilation success rate. We further demonstrate DECLARUI's generalizability through successful applications to Flutter and ArkUI frameworks. User studies with professional developers confirm that DECLARUI's generated code meets industrial-grade standards in code availability, modification time, readability, and maintainability. By streamlining app development, improving efficiency, and fostering designer-developer collaboration, DECLARUI offers a practical solution to the persistent challenges in mobile UI development.

*Ting Zhou and Yanjie Zhao are the co-first authors.

[†]Kai Chen (kchen@hust.edu.cn) and Haoyu Wang (haoyuwang@hust.edu.cn) are the corresponding authors.

Authors' Contact Information: **Ting Zhou**, tingzhou27@hust.edu.cn, Hubei Key Laboratory of Distributed System Security, Hubei Engineering Research Center on Big Data Security, School of Cyber Science and Engineering, Huazhong University of Science and Technology, Wuhan, China; **Yanjie Zhao**, Yanjie_Zhao@hust.edu.cn, Hubei Key Laboratory of Distributed System Security, Hubei Engineering Research Center on Big Data Security, School of Cyber Science and Engineering, Huazhong University of Science and Technology, Wuhan, China; **Xinyi Hou**, xinyihou@hust.edu.cn, Hubei Key Laboratory of Distributed System Security, Hubei Engineering Research Center on Big Data Security, School of Cyber Science and Engineering, Huazhong University of Science and Technology, Wuhan, China; **Xiaoyu Sun**, Xiaoyu.Sun1@anu.edu.au, Australian National University, Canberra, Australia; **Kai Chen**, kchen@hust.edu.cn, Hubei Key Laboratory of Distributed System Security, Hubei Engineering Research Center on Big Data Security, School of Cyber Science and Engineering, Huazhong University of Science and Technology, Wuhan, China; **Haoyu Wang**, haoyuwang@hust.edu.cn, Hubei Key Laboratory of Distributed System Security, Hubei Engineering Research Center on Big Data Security, School of Cyber Science and Engineering, Huazhong University of Science and Technology, Wuhan, China.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2025 Copyright held by the owner/author(s).

ACM 2994-970X/2025/7-ARTFSE011

<https://doi.org/10.1145/3715726>

CCS Concepts: • **Software and its engineering** → **Software development techniques; Development frameworks and environments.**

Additional Key Words and Phrases: Declarative UI, Mobile App Development, Multimodal Large Language Models, Code Generation, React Native, Flutter, ArkUI

ACM Reference Format:

Ting Zhou, Yanjie Zhao, Xinyi Hou, Xiaoyu Sun, Kai Chen, and Haoyu Wang. 2025. DeclarUI: Bridging Design and Development with Automated Declarative UI Code Generation. *Proc. ACM Softw. Eng.* 2, FSE, Article FSE011 (July 2025), 23 pages. <https://doi.org/10.1145/3715726>

1 Introduction

As the mobile app ecosystem continues to expand, with a constant influx of new apps entering the market, the development of user interfaces (UIs) has evolved to meet the increasing demands for better user experiences and more efficient development processes. Traditional imperative UIs, which specify step-by-step instructions for rendering UI elements, have given way to declarative UIs, which describe the desired state of the UI rather than the sequence of operations to achieve it [30, 54]. Declarative UIs offer several advantages over their imperative counterparts, such as improved maintainability, testability, and separation of concerns by decoupling the UI logic from the underlying implementation details.

The adoption of declarative UI frameworks in mobile app development has become widespread due to these benefits [21, 30]. However, while these frameworks have simplified many aspects of UI development, they have not significantly reduced the amount of manual coding required. In real-world scenarios, the development process typically begins with the design of the UI [42], followed by the complex and time-consuming task of translating these visual designs into functional declarative UI code. This process requires developers to manually map visual elements to their corresponding implementation details, which remains prone to errors and can be particularly challenging despite the advantages of declarative frameworks. Consequently, **there is a need for automated tools to streamline declarative UI code generation directly from UI designs.**

Despite this evident need, current research efforts in this area are still inadequate. Some methods rely on heuristics or machine learning models to extract UI components and layout information from design images [7, 10], but they often struggle with complex designs and lack ability to handle interactive logic. Other approaches use program synthesis techniques to generate code from natural language descriptions [56], but they require detailed specifications and cannot directly interpret visual designs. Researchers have also explored generating UI code across different frameworks [4, 14], but these approaches are carried out on traditional imperative UI frameworks and cannot be extended to declarative UI frameworks. In summary, **existing research falls short in addressing specific challenges of automatically generating declarative UI code for real-world development scenarios and lacks ability to generalize across different UI frameworks.**

Recent advancements in multimodal large language models (MLLMs) have opened up new possibilities for generating code from visual inputs. MLLMs, which integrate image processing capabilities into large language models (LLMs), have demonstrated superior performance in understanding and interpreting visual information compared to traditional computer vision (CV) models based on convolutional neural networks (CNNs) [6, 55, 61]. Furthermore, LLMs have exhibited remarkable proficiency in various code intelligence tasks, such as code generation [12, 13, 32, 40], translation [17, 33], summarization [2, 8, 23, 27], and repair [16, 39, 59, 60]. The synergy between visual comprehension and code generation abilities within MLLMs presents a promising avenue for translating UI designs into functional code, bridging the gap between visual representations and programmatic implementations.

Despite the powerful image understanding and code generation abilities of MLLMs, directly applying them to UI code generation faces several significant challenges. **Imprecise component recognition** often occurs due to the complexity of UI designs and the lack of explicit annotations. MLLMs also struggle with a **limited understanding of interactive logic**, failing to accurately infer navigation flow and event handling from static images. Furthermore, they often generate **inconsistent behavior across pages**, producing code that may work for individual screens but fails to maintain coherence in a multi-page app. Lastly, the generated code frequently suffers from **reliability and compilation issues**, containing syntax errors or violating framework-specific conventions. These challenges underscore the need for a specialized approach that can bridge the gap between MLLMs' capabilities and the specific requirements of declarative UI code generation, addressing accuracy, consistency, and reliability in the process.

To fill this gap, we propose DECLARUI, a novel approach to enhance MLLM-based declarative UI code generation from UI designs. DECLARUI integrates advanced CV techniques with MLLMs to decompose complex UI designs into structured component information. We introduce the Page Transition Graph (PTG) to represent the app's navigation logic, which, combined with the structured component information, forms a comprehensive prompt for the MLLM. After initial code generation, DECLARUI employs an iterative refinement process, performing navigational integrity and compilation checks to identify and rectify common errors. This multi-stage approach ensures the generated UI code is both visually accurate and functionally robust across multiple pages.

In summary, our key contributions are as follows:

- We propose DECLARUI, an approach that enhances MLLM-based declarative UI code generation from UI designs. DECLARUI combines CV techniques with MLLMs to decompose UI designs into structured component information, introduces the PTG to represent navigation logic, and employs an iterative refinement process.
- Our study compiled a dataset of 50 top-ranked UI design sets (50% mockups, 50% app screenshots) from design websites and app stores across 10 categories. For React Native [47], DECLARUI significantly outperformed state-of-the-art MLLMs like Claude-3.5 [1]¹ and GPT-4o [44] in generating UI code. DECLARUI achieved a 123% improvement in PTG coverage rate, up to 55% increase in visual similarity scores, and a 29% higher compilation success rate. Furthermore, DECLARUI demonstrated strong generalization to Flutter [22] and ArkUI [29], effectively bridging UI designs and functional declarative UI code across frameworks.
- We conducted a user study with professional developers to evaluate the performance of DECLARUI. The results demonstrate that DECLARUI significantly outperforms the baseline in terms of code availability (4.97 vs. 4.15), modification time (4.03 vs. 2.45), readability (4.62 vs. 3.75), and maintainability (4.55 vs. 3.37). Developers reported that the UI code generated by DECLARUI meets industrial-grade standards in multiple dimensions, highlighting its potential to streamline app development processes and improve designer-developer collaboration.

2 Background and Motivations

The rapid advancement of MLLMs has ushered in a new era of possibilities for automating mobile app development. Models such as OpenAI's DALL-E [46] and Google's Imagen [48] have showcased remarkable abilities in generating images and text from natural language prompts, inspiring researchers to explore their potential in generating app code directly from UI design mockups or app screenshots [28, 52, 58]. This innovative approach holds the promise of streamlining the app development process, potentially allowing designers to transform their visual concepts into functional UI code with unprecedented ease and speed.

¹The specific reference to "Claude-3.5" in this paper refers to **Claude 3.5 Sonnet**.

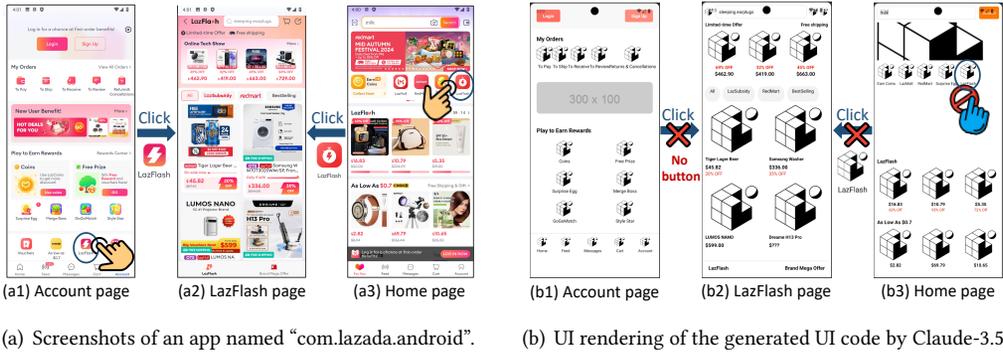


Fig. 1. Motivating example comparing the original app screenshots with the rendering of code generated by an MLLM. Please note that this paper focuses on code generation evaluation and does not consider the impact of image assets on the generation process.

Unfortunately, despite the exciting potential of MLLMs for UI code generation, current approaches face significant challenges when processing raw UI design mockups or screenshots. In [Figure 1](#), we provide an example of a shopping app named “com.lazada.android” [26], sourced from Google Play, to illustrate some of the issues with generating UI code using an MLLM, i.e., Claude-3.5:

Imprecise Component Recognition. The MLLM struggles with accurate identification and rendering of UI components. As evident in [Figure 1\(b1\)](#), the generated *Account page* lacks several crucial elements present in the original design ([Figure 1\(a1\)](#)). Most notably, the *LazFlash button* is missing, which is a critical component for navigation within the app. This omission demonstrates that MLLMs may not fully capture all the necessary visual elements, potentially leading to incomplete or inaccurate UI implementations.

Limited Understanding of Interactive Logic. UI design encompasses more than just visual elements; it includes complex interaction patterns and user experience considerations. The MLLM faces difficulties in accurately interpreting the designer’s intentions and the relationships between different elements. This is clearly illustrated by the lack of navigation logic in the generated UI code. The arrows with crosses in [Figure 1\(b\)](#) indicate that the expected transitions between pages (such as from the *Account page* to the *LazFlash page*) are not implemented in the generated code, despite being clearly present in the original designs shown in [Figure 1\(a\)](#).

Inconsistent Behavior Across Pages. The MLLM struggles to recognize and implement consistent behaviors across different pages of an app. This challenge is exemplified in [Figure 1](#), where both the *Account page* (a1) and the *Home page* (a3) should navigate to the *LazFlash page* (a2) when the *LazFlash button* is clicked. However, the generated UI code fails to maintain this consistency, as indicated by the crossed-out arrows in [Figure 1\(b\)](#). This suggests that MLLMs have difficulty in analyzing and implementing complex attention patterns across multiple screenshots, potentially leading to inconsistent user experiences in the generated app.

Code Reliability and Compilation Issues. The output from the MLLMs can be unpredictable, often resulting in code that contains compilation errors or runtime issues. While not directly visible in [Figure 1](#), this challenge is implicit in the discrepancies between the original app screenshots and the rendered UI from the generated code. These issues necessitate significant human intervention to correct errors and ensure the code is functional, potentially offsetting the time-saving benefits of using MLLMs for UI code generation.

These challenges underscore the current limitations of MLLM-based approaches in UI code generation. While the potential for automating app development is promising, significant improvements in component recognition, understanding of interactive logic, cross-page consistency, and code reliability are necessary before MLLMs can reliably transform visual designs into fully functional UI code. Addressing these challenges will be crucial for realizing the vision of streamlined app development through AI-assisted code generation.

Our Solutions. To address the issue of imprecise component recognition, we leverage **advanced CV techniques**. By integrating the Grounding DINO object detection model [38] with the Segment Anything image segmentation model [34], we are able to precisely isolate individual components within UI designs. This preprocessing step significantly enhances the MLLM’s ability to recognize UI elements, ensuring that the generated code more accurately reflects the design intent. To tackle the limitations in understanding interactive logic and maintaining consistent behavior across pages, we introduce the concept of **Page Transition Graph (PTG)**. The PTG serves as a structured representation that clearly delineates the navigation logic between different pages. By providing this high-level abstraction to the MLLM, we compensate for the model’s weaknesses in comprehending complex inter-page relationships, thereby generating interaction logic that more closely aligns with the designer’s intentions. Finally, to improve code reliability and reduce compilation errors, we implement **an iterative refinement process**. This mechanism not only checks the navigational integrity of the MLLM-generated code but also performs compilation checks, ensuring consistent code quality. This approach significantly reduces the need for manual intervention, lowering time costs while simultaneously improving the overall quality of the generated code.

3 Approach

The workflow of DECLARUI, as illustrated in Figure 2, starts with UI designs as input. These designs include high-quality Figma [20] design mockups and screenshots of Google Play [25] apps. Design mockups typically include both visual elements and text-based specifications, providing comprehensive design information. While text-based design specifications are readily interpretable by LLMs, DECLARUI aims to address more challenging scenarios where **only visual designs are available**. Therefore, DECLARUI processes both sources uniformly as image data, focusing exclusively on visual content. This unified image-based analysis allows DECLARUI to process diverse inputs consistently, whether they are design mockups or app screenshots.

DECLARUI preprocesses the input through two key steps: **PTG Construction** (§ 3.1), which captures navigation logic, and **UI Component Extraction and Representation** (§ 3.2), which utilizes CV techniques. Subsequently, the **Prompt Synthesis** step (§ 3.3) integrates the preprocessed data—including the generated PTG, component analysis results, and the complete UI screenshot—into a comprehensive prompt. This prompt serves as the foundation for the MLLM to generate complete UI code. Lastly, an **Iterative Code Refinement** process (§ 3.4) checks for errors and navigation absence, ensuring the reliability of the final code output.

3.1 PTG Construction

To address challenges rooted in the complexities of modern UI designs and MLLM limitations, we propose the utilization of PTGs as a key tool for capturing and representing UI interaction logic. We formally define a PTG where nodes represent individual app pages, and edges represent transitions between pages along with their triggering conditions. The *PTG* is defined as a tuple:

$$PTG = (N, E)$$

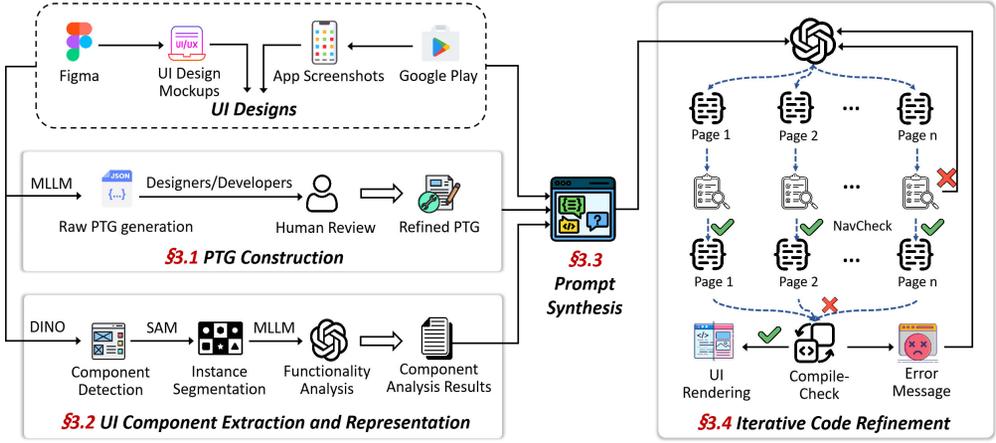


Fig. 2. The workflow of DECLARUI.

where:

- N is a finite set of nodes representing app pages. Each node includes an *id* as a unique identifier, a *name* for the page name, a *type* indicating the page type, and *property* containing additional properties. N is defined as:

$$N = \{n_i = (id_i, name_i, type_i, property_i) \mid i = 1, 2, \dots, k\}$$

- E is a set of directed edges representing transitions between pages. Each edge has specific attributes: *id* is a unique identifier for the edge, $ns, nt \in N$ denote the source and target nodes, while *condition* is the condition that triggers the transition, and *action* is the operation performed during the transition. E is defined as:

$$E = \{e_i = (id_i, ns_i, nt_i, condition_i, action_i) \mid i = 1, 2, \dots, m\}$$

Our introduction of PTGs seeks to offer MLLMs a clear, structured representation of interaction logic, effectively bridging the gap in expressing dynamic interactions in static UI designs. The PTG generation process leverages the robust inference capabilities of MLLMs combined with the analysis of UI designs. Initially, a set of UI designs is provided to an MLLM, i.e., Claude-3.5, accompanied by a formal definition of the PTG structure. The MLLM analyzes the UI designs and generates a comprehensive PTG based on the provided prompt, as exemplified in Figure 3. It is guided to output its analysis in a structured JSON format to facilitate further processing and integration. Figure 4 illustrates an example PTG for a simplified e-commerce app scenario, designed to clearly demonstrate the PTG concept. Finally, designers or developers review the generated PTG, making necessary corrections or additions to ensure the accuracy and completeness of the captured interaction logic.

PTGs offer a structured representation of UI elements, encompassing page hierarchies, transition logic, and navigation rules. PTGs play a crucial role in addressing key challenges associated with utilizing MLLMs for UI code generation from designs, including limited understanding of interactive logic and inconsistent behavior across pages.

Given the following set of UI design images:
[Image Set]
 And the formal definition of a Page Transition Graph (PTG):
[PTG Formal Definition]
 Please analyze the UI design images and generate a complete Page Transition Graph (PTG) that captures all pages, their elements, and the possible transitions between them. Ensure that you include both page-specific transitions and global transitions where applicable.

Fig. 3. Example of prompt used to construct PTG.

```

{
  N = [{id: 'login', name: 'Login Page', type: 'login'},
    {id: 'home', name: 'Home Page', type: 'list', props: {layout: 'grid', itemType: 'product'}},
    {id: 'search', name: 'Search Page', type: 'search'},
    {id: 'product', name: 'Product Detail Page', type: 'detail', props: {layout: 'vertical'}},
    {id: 'cart', name: 'Shopping Cart Page', type: 'list', props: {layout: 'vertical', itemType: 'cartItem'}}],

  E = [{id: 'login->home', from: 'login', to: 'home', conditions: 'isLoggedIn'},
    {id: 'home->search', from: 'home', to: 'search', actions: 'passSearchParams'},
    {id: 'home->product', from: 'home', to: 'product', actions: 'passProductId'},
    {id: 'search->product', from: 'search', to: 'product', actions: 'passProductId'},
    {id: 'product->cart', from: 'product', to: 'cart', actions: 'addToCart'},
    {id: 'cart->product', from: 'cart', to: 'product', actions: 'passProductId'}]
}
    
```

Fig. 4. Example of the PTG for a simplified e-commerce app scenario.

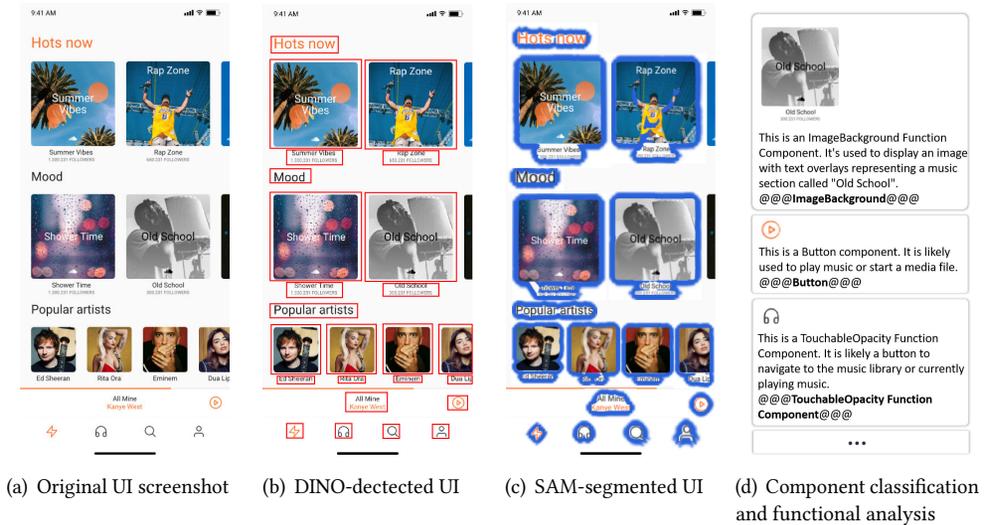


Fig. 5. Example of the UI component extraction process.

3.2 UI Component Extraction and Representation

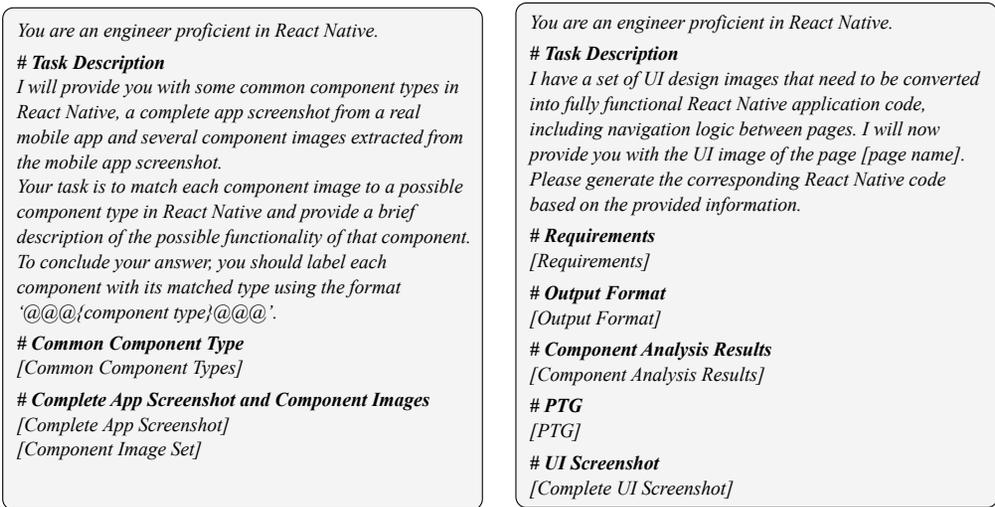
This process aims to transform the complex UI design images into a structured representation of components, which is more easily understood and processed by MLLMs for generating high-quality code, as illustrated in Figure 5. It includes three stages as follows:

Automated Component Detection. Instead of relying on manual annotation, DECLARUI leverages the capabilities of Grounding DINO [38], a state-of-the-art object detection model. Grounding DINO excels in zero-shot object detection, making it ideal for identifying potential UI components without requiring extensive training on UI-specific datasets. The model generates bounding boxes around detected objects, which serve as initial component localization.

Automated Instance Segmentation. To achieve precise instance segmentation of the UI components, DECLARUI employs the Segment Anything Model (SAM) [34]. The bounding boxes produced by Grounding DINO act as prompts for SAM, guiding the segmentation process. This approach combines the strengths of both models: Grounding DINO’s ability to detect objects without prior training, and SAM’s capability to produce highly accurate segmentation masks.

Component Classification and Functional Analysis. We then leverage an MLLM (i.e., Claude-3.5) to process the segmented component images. As shown in Figure 6(a), a carefully designed prompt is used to extract the component type. The prompt includes the segmented images and instructs the MLLM to classify the component type (e.g., button, text field, dropdown menu, etc.) and infer its potential functionality within the UI. This step utilizes the MLLM’s understanding of UI design patterns and its ability to reason about the purpose of visual elements.

The output of this three-stage process is a comprehensive, structured representation of the UI design images. This representation includes precise spatial information and instance segmentation masks for each component, type classifications, and inferred functional descriptions.



(a) Example of the prompt used to extract component types.

(b) Example of a synthesized prompt.

Fig. 6. Examples of the prompts used in § 3.2 and § 3.3.

3.3 Prompt Synthesis

Prompt synthesis integrates the preprocessed data into a comprehensive prompt for the MLLM, enhancing its ability to generate UI code that accurately reflects both the visual elements and the underlying interactive logic of the original design. This process combines three key elements: the **PTG**, which provides a structured representation of the navigation logic and inter-page relationships; **component analysis results**, detailing UI element types and their functionalities; and the

complete UI images, providing visual context. Together, these components give the MLLM a holistic understanding of both the structure and visual design, enabling accurate and functionally complete UI code generation. [Figure 6\(b\)](#) illustrates an example of the synthesized prompt.

The resulting prompt is fed into the MLLM for the initial code generation, setting the stage for the subsequent iterative code refinement process, where the generated code undergoes navigation and compilation checks for further improvement.

3.4 Iterative Code Refinement

The initial code generated by MLLMs may fail to implement the navigation logic defined in the PTG fully, contain syntax errors, or exhibit logical inconsistencies. To address these potential issues and ensure high-quality code output, DECLARUI employs a two-pronged strategy for code refinement: **navigation consistency validation** and **compile-time error correction**. This iterative process systematically identifies and rectifies discrepancies between the generated code and the intended app structure and behavior.

Navigation Consistency Validation. To ensure fidelity to the PTG, DECLARUI implements a comprehensive navigation consistency validation process using static code analysis techniques. This process begins by parsing the generated code to extract all implemented navigation logic, specifically identifying framework-specific navigation statements without executing the code. The extracted navigation paths are then meticulously compared against the PTG definition, with a primary focus on identifying any missing transitions that are present in the PTG but absent in the generated code. This static analysis approach allows for additional navigation options in the code that may not be explicitly defined in the PTG, providing flexibility for potential enhancements or alternative user flows. Upon completion of the analysis, a detailed report is generated, highlighting any navigation inconsistencies, specifically focusing on missing transitions. This report serves as a crucial input for the MLLM in subsequent refinement iterations, offering targeted guidance to address specific navigation-related omissions and ensuring comprehensive implementation of the intended user journey as defined in the PTG.

Compile-Time Error Correction. An integral part of DECLARUI's iterative refinement process is leveraging compiler feedback to detect and correct syntactic errors in the generated code. This is achieved through an automated compilation checking script that systematically processes the code using the appropriate development tools. The script captures and analyzes compiler error messages generated during this process. To bridge the gap between the compiler's technical output and the MLLM's natural language processing capabilities, these error messages are processed and restructured into clear, actionable instructions. This process involves extracting key information from error messages, including file names, line numbers, and error descriptions. The extracted information is then combined with relevant code snippets to form a structured prompt. This prompt is fed back to the MLLM, enabling it to make targeted corrections to the code.

Iterative Refinement Process. The code refinement process in DECLARUI is iterative, designed to continuously enhance the quality of the generated code. Initially, DECLARUI generates code for each app page based on the PTG. Then, it undergoes iterative cycles to check and correct navigation inconsistencies and compile-time errors. After generating each page's code, a navigation consistency check is performed, producing a report of any missing navigation paths. This report is fed back into the MLLM to address the gaps, with up to three iterations per page to ensure completeness. Once all pages are generated and compiled into a project, DECLARUI conducts a compile-time error analysis. It identifies the relevant files and error messages, which are passed to the MLLM for targeted corrections. The MLLM refines the code based on both navigation and compile-time issues. This process is repeated, with compile-time error correction limited to three

iterations for the entire project to maintain efficiency. This iterative approach ensures continuous improvement, addressing specific issues while keeping the process efficient through iteration limits.

4 Evaluation

To comprehensively evaluate the performance of DECLARUI, we designed a series of experiments to address three key research questions (RQs):

RQ1: Effectiveness. How effective is DECLARUI? How does DECLARUI perform in comparison to baseline methods when operating autonomously without human intervention?

RQ2: Ablation Study. How do individual key modules contribute to the overall performance of DECLARUI?

RQ3: User Study and Manual Repair. How do professional developers evaluate DECLARUI's generated UI code compared to the baseline in real-world scenarios? To what extent is manual effort required to repair compilation failures in DECLARUI's output?

4.1 Dataset

We curated a dataset of 50 mobile app UI designs (totaling 250 app pages), including 25 UI design mockups and 25 app screenshots, across 10 functional categories. This diversity allows us to assess the generalizability and robustness of our approach across various UI design paradigms and app types. We now provide details on the data sources, selection criteria, and dataset composition.

Data Sources. Our dataset is drawn from two main sources: Figma [20] and Google Play [25]. Figma, a leading collaborative design tool, provides a selection of contemporary, high-quality mobile UI design mockups. Google Play, the official app distribution platform for Android devices, serves as our source for downloading apks, which we subsequently install on an emulator to manually capture screenshots. By incorporating cutting-edge designs from Figma and real-world apps from Google Play, we ensure that our approach is evaluated against a dataset that reflects both ideal design practices and practical implementation challenges in the mobile app development landscape.

Selection Criteria. To ensure diversity in our dataset, we employed strict selection criteria. We classified mobile apps into 10 functional categories, including *Social Networking*, *Entertainment*, *Productivity*, *Education*, *Wellness*, *Finance*, *Shopping & E-commerce*, *Travel & Navigation*, *News & Information*, and *Utility & Practical Tools*. For Figma designs, we selected 194 mobile UI design mockups with over 300 followers, indicating their recognition within the design community. We then excluded designs that didn't meet our page count requirements (i.e., less than five pages) or weren't specifically designed for mobile apps. For Google Play apps, we focused on those ranking in the top 10 by download count within their respective categories, ensuring the representation of popular and widely used apps. We excluded paid apps and those requiring complex registration processes. From each source, we selected 3-5 designs or apps per category that met our criteria. This range was established to accommodate potential scarcity in certain categories while maintaining a minimum threshold for diversity. For each category, we ultimately collected no fewer than five designs; for those with more than five, we randomly selected five for experimentation to control experimental costs while ensuring research quality and statistical significance.

To ensure the quality and reliability of our dataset, we implemented a comprehensive validation process for PTG construction. Two authors independently analyzed each selected app and design to create ground-truth PTGs. The validation process was conducted through direct interaction with real-world applications, allowing authors to map actual navigation patterns and understand the complete user journey. For Figma designs, authors performed semantic validation of UI components

and their relationships, ensuring an accurate representation of design intentions. Through cross-validation between the two independent annotations, we achieved 100% accuracy in our ground-truth PTG construction. This rigorous validation process establishes a reliable foundation for our evaluation metrics, particularly in assessing PTG coverage rates, and ensures that our results accurately reflect the real-world performance of DECLARUI.

Dataset Composition. The composition of our dataset is structured to provide a balanced and comprehensive representation of mobile UI designs. We selected 50 apps in total, equally divided between Figma designs and Google Play apps. These 50 apps are distributed across 10 functional categories, with five apps per category. The categorization was designed to encompass the broad spectrum of mobile app functionalities, aiming to provide a comprehensive representation of the app ecosystem. For each app, we captured five logically connected pages, allowing us to analyze the flow and interaction between related interfaces.

4.2 Experimental Setup

4.2.1 Evaluation Metrics. Our evaluation employs a set of metrics to assess the performance of DECLARUI, encompassing visual fidelity, navigation logic, code quality, and practical usability.

To evaluate visual fidelity, we utilize two metrics:

- **CLIP score** [50]: This metric quantifies the semantic similarity between the generated and original UIs, providing an indicator of how well the generated UI captures the intended visual elements and overall design concept.
- **SSIM (Structural Similarity Index Measure) score** [53]: It assesses the layout and compositional accuracy, focusing on the spatial arrangement and structural similarities between the generated and original UIs.

The completeness of navigation logic implementation is measured by **PTG Coverage Rate (PCR)**. This metric, manually calculated by human evaluators, assesses the alignment between the implemented app structure and the intended design. It compares the predefined PTG (e.g., Figure 4) with the UI Transition Graph (UTG) of the apps generated by DECLARUI. The UTG is obtained using DroidBot [36], a widely adopted automated testing tool. The PCR is calculated as:

$$PCR = \frac{|PTG_{edges} \cap UTG_{edges}|}{|PTG_{edges}|} \times 100\%$$

As described in § 3.4, we apply a compile-time error correction process, allowing up to three compilation attempts for each app. This iterative process aims to maximize the number of successfully compiled apps while minimizing manual intervention. We then employ different metrics based on the compilation outcome.

For successfully compiled apps (within three iterations), we calculate:

- **Compilation Success Rate (CSR)** [56]: This measures the percentage of generated app code that is compiled successfully without errors after our automated iterative refinement process, without any human intervention:

$$CSR = \frac{\text{Number of successfully compiled apps}}{\text{Total number of apps}} \times 100\%$$

- **Average Compilation Iteration Count (ACIC)**: This metric indicates code generation efficiency and is computed as:

$$ACIC = \frac{\sum_{i=1}^n \text{iterations}_i}{n}$$

where n is the total number of successfully compiled apps and $iterations_i$ is the number of compilation attempts for the i -th app (maximum 3). A lower ACIC suggests more efficient initial code generation.

For apps that fail to compile after three iterations (termed *compilation-failed apps*), we measure:

- **Average Manual Correction Time (AMCT):** This quantifies the average time required for human developers to manually correct remaining errors in the generated code:

$$AMCT = \frac{\text{Total manual correction Time}}{\text{Number of compilation-failed apps}}$$

AMCT provides insight into the practical efficiency gains of our method and the potential reduction in developer workload.

To provide a comprehensive assessment of the generated code's practical usability and efficiency, we employ **user study metrics** [11]. These metrics evaluate:

- **Code Availability:** Quantifies the proportion of usable code in the generated output.
- **Modification Time:** Measures time to adapt generated code to production standards.
- **Readability:** Assesses the structural understandability of the code.
- **Maintainability:** Measures the ease of adapting the code to new requirements.

Scores for these metrics are assigned using predefined scales and intervals, allowing for a nuanced evaluation of the code's quality from a developer's perspective.

4.2.2 Baselines. We compare DECLARUI against two baselines: GPT-4o and Claude-3.5 (specifically, Claude-3.5-Sonnet-20240620). Both baselines use unprocessed UI design images without our custom prompting techniques or additional processing. To ensure a fair comparison and avoid underestimating the capabilities of MLLMs, we meticulously designed the prompts for the baselines. These prompts adhere to established practices from existing image-to-code research [52, 56, 58] and widely-adopted techniques for LLM utilization [18, 51]. Both baseline approaches are evaluated within the React Native framework.

4.2.3 Implementation Details. Our experiments cover three popular mobile development frameworks that employ a declarative UI paradigm, i.e., React Native, Flutter, and ArkUI. For each framework, we generate UI code using our DECLARUI approach and evaluate its performance. All experiments were conducted on a server equipped with dual NVIDIA GeForce RTX 3090 GPUs, 64-core AMD EPYC processors, and 512GB of RAM. The code generation process was automated with custom Python scripts, totaling 2,935 lines, to minimize human intervention and bias.

4.3 Effectiveness (RQ1)

To address RQ1, we conducted experiments on a carefully curated dataset comprising 50 app UI designs, each containing five pages, as stated in § 4.1. Given their status as state-of-the-art MLLMs, we selected GPT-4o and Claude-3.5 as base models to power DECLARUI. The experiment's input consisted of these 50 app UI designs, while the output was the code generated by DECLARUI for each app's five pages. We integrated the generated code directly into an initial project in the IDE (i.e., Android Studio for React Native) without any manual intervention. The project was then packaged into an APK, which we installed and ran on the Android Studio built-in emulator to render the UI. We then evaluated the results according to the metrics outlined in § 4.2.1.

As shown in Table 1, the evaluation results reveal that DECLARUI (Claude-3.5) outperforms DECLARUI (GPT-4o) in several key areas. For example, DECLARUI (Claude-3.5) achieved a remarkable 98% Compilation Success Rate (CSR), meaning that only one out of the 50 apps failed to compile automatically. In contrast, DECLARUI (GPT-4o) had six apps that didn't compile under automated

conditions. Upon analyzing the single app that failed to compile with DECLARUI (Claude-3.5), we found that the issue was due to the app’s use of non-existent file resources, which caused it to fail multiple compilation checks. This issue primarily stems from occasional gaps in the MLLM’s domain knowledge about specific framework requirements. Notably, these technical issues can be systematically addressed through enhanced prompting strategies and framework-specific knowledge integration, as discussed in § 5.1.

While DECLARUI (GPT-4o) showed a marginally lower Average Compilation Iteration Count (ACIC) – 10 iterations for 44 compiled apps versus 11 for 49 apps with DECLARUI (Claude-3.5) – this advantage was minimal, differing by only 0.01 iterations per app. Despite this slight edge in iteration efficiency, DECLARUI (GPT-4o) lagged behind in other crucial metrics. The SSIM score showed a particularly significant lead of 17.2% (0.68 vs. 0.58), indicating that DECLARUI (Claude-3.5) excels in accurately reproducing the structural elements of the original UI. Moreover, it exhibited a better PTG Coverage Rate (PCR), suggesting a more comprehensive understanding of navigation logic. These results suggest that while both models show promise, **DECLARUI (Claude-3.5) offers a more robust and accurate solution for automating UI code generation from designs**. Its higher CSR, coupled with better visual fidelity and navigation logic understanding, positions it as the more reliable choice for practical applications. The marginally higher ACIC for DECLARUI (Claude-3.5) is a small trade-off for the significant improvements in other areas, particularly given the importance of visual accuracy and functional correctness in UI development.

Table 1. Comprehensive performance metrics across different methods and frameworks.

Framework	Method	CLIP Score	SSIM Score	PCR (%)	ACIC	CSR (%)
React Native	Baseline (GPT-4o)	0.66	0.44	1.4	-	76.0
	Baseline (Claude-3.5)	0.65	0.41	43.3	-	74.0
	DECLARUI (GPT-4o)	0.87	0.58	88.6	0.23	88.0
	DECLARUI (Claude-3.5)	0.88	0.68	96.8	0.22	98.0
Flutter	DECLARUI (Claude-3.5)	0.85	0.68	99.4	0.82	92.0
ArkUI	DECLARUI (Claude-3.5)	0.85	0.66	95.9	1.62	86.0

DECLARUI vs. Baselines. As detailed in § 4.2.2, we carefully selected and configured our baselines for comparison. The results, presented in Table 1, reveal significant limitations in using GPT-4o and Claude-3.5 directly for declarative UI code generation from mobile app UI designs. Under fully automated conditions, the baseline methods achieved a CSR of 76% (GPT-4o) and 74% (Claude-3.5). Compilation failures primarily stemmed from issues such as using non-existent packages or components, inadequate special character handling, and component utilization without proper package imports. While Baseline (GPT-4o) showed a slight edge in visual fidelity compared to Baseline (Claude-3.5), the advantage was marginal, with both falling short of our DECLARUI’s performance. Notably, there was a stark contrast in PCR between the baselines, with Baseline (GPT-4o) achieving 1.4% and Baseline (Claude-3.5) reaching 43.4%. This disparity highlights Claude-3.5’s substantial advantage in logical comprehension over GPT-4o in the baseline scenario.

Comparing DECLARUI to the baselines reveals significant improvements across key metrics. The PCR, a crucial indicator of an app’s navigational structure implementation, demonstrates DECLARUI’s substantial advantage. In the React Native framework, DECLARUI achieves 88.6% and 96.8% coverage with GPT-4o and Claude-3.5 respectively, far surpassing the corresponding baseline methods (1.4% and 43.3%). This marked improvement suggests DECLARUI’s enhanced ability to capture and implement complex app structures accurately. The compilation statistics further underscore DECLARUI’s superiority in code quality. The 98% CSR of DECLARUI (Claude-3.5)

compared to Baseline (Claude-3.5)'s 74%, and DECLARUI (GPT-4o)'s 88% versus Baseline (GPT-4o)'s 76%, clearly demonstrate DECLARUI's significant outperformance of baseline methods. This high CSR, coupled with low iteration requirements (only 11 iterations needed for 49 successfully compiled apps), indicates that DECLARUI generates more robust, deployment-ready code. Based on the comprehensive performance of Claude-3.5 and GPT-4o in both baseline scenarios and with DECLARUI, we have selected Claude-3.5 as the base model for our subsequent experiments.

Generalization Capabilities. The results for Flutter and ArkUI presented in Table 1 demonstrate DECLARUI's robust performance across diverse UI frameworks, highlighting its impressive generalization capabilities and adaptability. Consistent with our evaluation approach for React Native, we compiled Flutter projects in Android Studio and ArkUI projects in DevEco Studio, integrating the generated code directly into initial projects without manual intervention.

Flutter emerges as particularly noteworthy, with DECLARUI achieving an exceptional PCR of 99.4%. This near-perfect score indicates DECLARUI's profound understanding and accurate implementation of navigation logic within the Flutter ecosystem. React Native and ArkUI also display impressive results, with PCRs of 96.8% and 95.9% respectively, further reinforcing DECLARUI's versatility across different framework paradigms. The CSRs across frameworks, ranging from 86% to 98%, reveal both DECLARUI's strengths and areas for potential improvement. Specifically, the slightly lower CSR for ArkUI presents an intriguing avenue for future investigation. We explore potential strategies for improving DECLARUI's performance with ArkUI in § 5.1, where we discuss the implementation of advanced techniques to potentially bridge this performance gap.

Answer to RQ1: DECLARUI demonstrates high effectiveness in generating functional UI code across multiple frameworks, significantly outperforming baselines. It achieves up to 98% CSRs, superior visual fidelity, and navigation logic implementation (PCR up to 99.4%). With strong generalization across React Native, Flutter, and ArkUI, and Claude-3.5 as the preferred base model, DECLARUI proves to be a robust solution for automated UI development, despite some room for improvement in ArkUI, which will be discussed in § 5.1.

4.4 Ablation Study (RQ2)

To assess the contribution of each key component in DECLARUI, we conducted two ablation studies on a subset of 13 apps from our dataset. This selection was based on CLIP score distributions across app categories, ensuring diverse representation. Our sample includes: three apps from categories with the lowest median CLIP scores, five from categories with median performance, three from categories showing high CLIP score variability, one from a category with low variability, and one from a category exhibiting multiple outliers in CLIP scores. This diverse sample spans various complexities and visual effects, providing a robust basis for evaluating DECLARUI's performance across different scenarios. We focused on two key aspects: the impact of CV techniques, and the role of PTG and the refinement process. It's important to note that we couldn't isolate PTG construction alone due to its inherent coupling with the iterative refinement process, as the navigation checks during refinement rely on the PTG. These studies help isolate the effects of these components on the overall performance of DECLARUI.

Table 2. Comprehensive performance metrics across different DECLARUI versions.

Version	CLIP Score	SSIM Score	PCR (%)	ACIC	CSR (%)
DECLARUI (Full)	0.90	0.72	96.8	0.22	98.0
DECLARUI (without CV)	0.88	0.57	95.1	0	100.0
DECLARUI (without PTG+Refinement)	0.87	0.68	23.1	-	80.0

Impact of CV techniques on UI Generation. In this ablation study, we removed the UI component extraction and representation module while retaining the PTG and refinement processes to investigate the impact of CV techniques on UI generation. The results in Table 2 illustrate that the absence of the UI component extraction and representation module led to a significant decrease in visual similarity scores, with the SSIM score dropping from 0.72 to 0.57, and the CLIP score slightly decreasing from 0.90 to 0.88. These reductions in structural and semantic similarity metrics directly reflect CV techniques' impact on the MLLM's component recognition capabilities. The decrease in visual similarity scores without the combination of CV techniques indicates that MLLMs struggle more with accurately identifying and reproducing individual UI elements, as evidenced by the drop in PCR from 96.8% to 95.1%.

Impact of PTG and Refinement Process. In the second ablation study, we removed both the PTG and refinement components, relying solely on the UI component extraction and representation for UI generation. Although the visual similarity scores remained relatively high based on Table 2, the PCR dropped drastically to 23.1% and the CSR decreased to 80.0%. These results emphasize the critical role of PTG and refinement in ensuring accurate navigation logic, code compilability, and overall UI quality.

Answer to RQ2: The combination of CV techniques contributes to the visual accuracy of UI generation, while the PTG and refinement processes ensure functional correctness and high code quality. The synergistic integration of these three components forms the core technology of DECLARUI, enabling the generation of high-fidelity and functionally-correct UIs.

4.5 User Study and Manual Repair (RQ3)

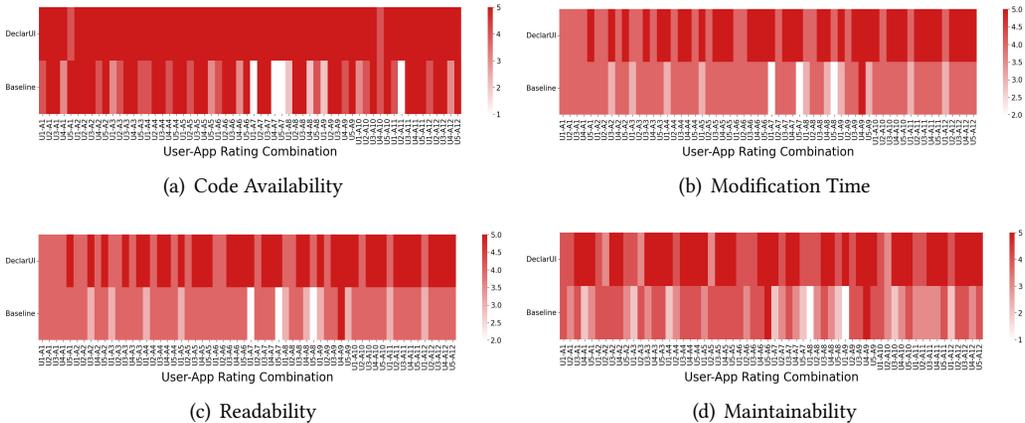


Fig. 7. Heatmaps of user study metrics.

4.5.1 User Study. To comprehensively evaluate the performance of our automatic generation method, we focused not only on the visual quality of the rendered UI but also on the overall quality of the generated code. In real-world application scenarios, particularly in industry settings, code readability, availability, and maintainability are three crucial evaluation dimensions [11, 43]. To assess these aspects, we conducted a user study with experienced developers.

Procedures. We recruited five React Native developers with an average of over two years of experience. All participants were compensated \$50 for their time and expertise. Each participant

reviewed 12 apps with UI code generated by DECLARUI and Baseline (Claude-3.5). These apps were carefully selected based on their CLIP scores: three well-performing apps (CLIP score > 0.85), six average-performing apps ($0.8 < \text{CLIP score} < 0.85$), and three poor-performing apps (CLIP score < 0.8), all randomly chosen within their respective groups. For code availability evaluation, participants modified each set of UI page code to meet basic business requirements. We tracked modifications using git, following the method of Chen et al. [11], mapping scores from 0.8 with each 5% interval corresponding to levels 1-5. After modifying each set, participants rated its readability and maintainability on a five-point Likert scale. Code modification time measures the time required to adjust the code to production standards. We scored this on 5-minute intervals, with 5 points for completion within 5 minutes and 1 point for over 20 minutes. To avoid tool bias, participants used Android Studio, with a 20-minute time limit per app. We recorded modification time and collected ratings for each set of UI page code. This study design allowed us to evaluate the generated code's quality through both objective measures (modification time and extent) and subjective assessments (readability and maintainability ratings) from experienced professionals.

Table 3. Participant ratings of DECLARUI vs. Baseline across different metrics.

Metric	Code Availability	Modification Time	Readability	Maintainability
Baseline Score	4.15	2.45	3.75	3.37
DECLARUI Score	4.97	4.03	4.62	4.55
P-value	9.23E-08	2.20E-10	2.57E-13	1.33E-12

Results. We present the average scores and p-values using the Mann-Whitney U test in Table 3, while more detailed statistical analyses are visualized using a set of heatmaps in Figure 7. Professional React Native engineers consistently rated DECLARUI-generated code significantly higher in terms of availability, modifiability, readability, and maintainability. The substantial improvements observed in all these metrics were statistically significant ($p < 0.05$), underscoring the robustness of our findings. Notably, DECLARUI achieved a near-perfect score in code availability (4.97 out of 5), demonstrating its exceptional performance in generating immediately usable code. The marked reduction in modification time (4.03 for DECLARUI vs. 2.45 for the baseline) indicates that DECLARUI-generated code requires considerably less effort to adapt to production standards. This is further evidenced by the fact that users were able to complete modifications on 73.3% of DECLARUI-generated code within 10 minutes. In contrast, 80% of the baseline code required more than 10 minutes to modify, and the actual modification time for the baseline is likely underestimated, as **27.1% of this 80% was not completed even within 20 minutes**. Furthermore, the superior readability (4.62 vs. 3.75) and maintainability (4.55 vs. 3.37) scores highlight the enhanced long-term value of the code produced by DECLARUI.

Case Study. To demonstrate DECLARUI's superiority, we present concrete examples in Figure 8 and Figure 9. Our analysis reveals three primary areas where DECLARUI excels over baseline methods: component fidelity, navigation completeness, and interaction consistency. We illustrate these using the "com.amazon.mShop.android.shopping" app [24]. DECLARUI achieves superior **component fidelity** by accurately replicating all elements from the original design (Figure 8(c)), including the bottom navigation bar omitted by the baseline (Figure 8(b)). This completeness extends to **navigation completeness** (Figure 9). DECLARUI maintains comprehensive navigation pathways (Figure 9(a)), addressing gaps present in the baseline's structure (Figure 9(b)). Consequently, our approach demonstrates improved **interaction consistency**, particularly in navigation between related UI components. This consistency is evident in the uniform navigation patterns to "Orders" from both "Profile" and "Account" pages (Figure 9(a)), ensuring a more coherent user experience.

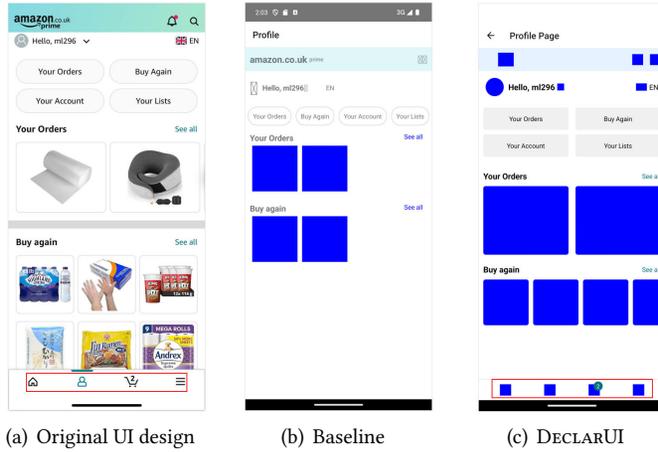


Fig. 8. The original UI design image, output from Baseline (Claude-3.5), and DECLARUI’s generated result.

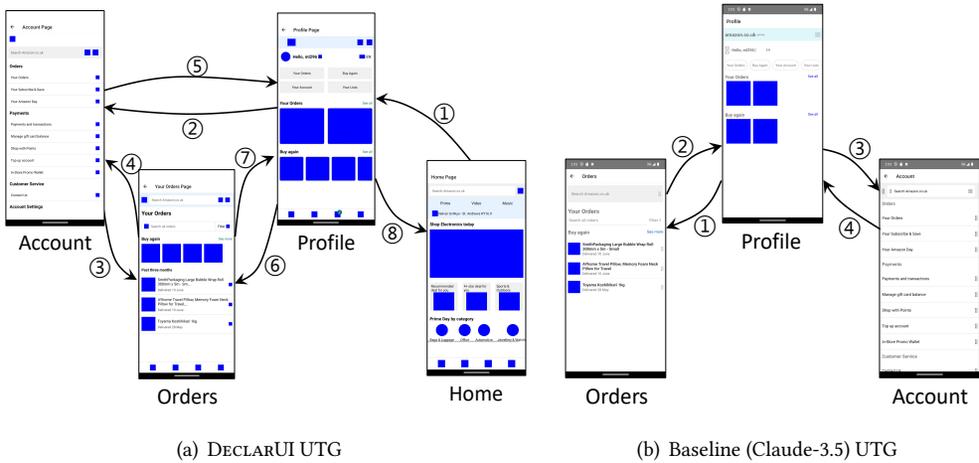


Fig. 9. Navigation completeness & interaction consistency.

4.5.2 *Manual Repair.* To comprehensively assess the human effort required to repair compilation failures of DECLARUI, we used the Average Manual Correction Time (AMCT) metric. Our analysis revealed that for the React Native framework, DECLARUI (Claude-3.5) required an AMCT of 48s, while DECLARUI (GPT-4o) needed 43.2s. In comparison, Flutter required 62s, and ArkUI took 183s. These results demonstrate DECLARUI’s efficiency in generating code that requires minimal manual intervention, highlighting its superiority across different frameworks.

Answer to RQ3: DECLARUI significantly outperforms the baseline in code quality, usability, and efficiency. The user study and case study demonstrate DECLARUI’s superiority in generating high-quality, readily usable UI code that closely matches original designs. Notably, DECLARUI requires minimal manual intervention, a characteristic observed consistently across multiple UI frameworks.

5 Discussion

5.1 Retrieval-Augmented Generation (RAG): Impact and Implications

The performance of DECLARUI across the three frameworks is uneven, with ArkUI showing relatively poor performance compared to the others (Table 1). This suboptimal performance likely results from the MLLM's limited knowledge of ArkUI, a newer and less mainstream framework, compared to its more comprehensive understanding of established frameworks.

To address this challenge, we explored integrating a knowledge graph (KG) [45] to supplement the MLLM's understanding of ArkUI. This approach follows a Retrieval-Augmented Generation (RAG) [19] method to compensate for the MLLM's lack of knowledge about less mainstream frameworks. When generating UI code, we first query the KG and then incorporate the retrieved information into the prompt. We tested this approach with 20 apps with UI designs from our dataset, totaling 100 app pages. The results showed promising improvements: ACIC is relatively reduced by 0.52 (1.1 vs. 1.62), CSR increased by 4% (90% vs. 86%), and AMCT reduced by 139.5s (43.5s vs. 183s). While our preliminary results demonstrated potential, the improvements were not as significant as anticipated, primarily due to time constraints limiting the sophistication of our KG implementation.

Future research should focus on refining the RAG-based approach, particularly for less common frameworks like ArkUI. This could involve developing a more comprehensive and sophisticated KG that captures the intricacies of ArkUI syntax and design patterns. Additionally, exploring techniques to dynamically update the KG with new framework features and best practices could ensure the system remains current and effective. Further investigation into optimizing the integration of KG-retrieved information with the MLLM's existing knowledge could potentially yield more substantial performance improvements. Moreover, extending this approach to other emerging UI frameworks could enhance the versatility and adaptability of DECLARUI.

5.2 Integration into Development Workflow

The practical adoption of DECLARUI in real-world development environments necessitates careful consideration of its integration into existing workflows. We envision DECLARUI functioning as an integral part of the development ecosystem through seamless IDE integration, enabling developers to generate UI code directly from design files within their familiar development environment. This integration could extend to continuous integration pipelines, where DECLARUI automatically processes design updates and generates corresponding code changes, maintaining synchronization between design and implementation throughout the development cycle. Furthermore, direct integration with design tools like Figma could facilitate real-time code generation during the design process, providing immediate feedback on implementation feasibility. The specific implementation of these integrations would naturally adapt to various team sizes and development practices, offering flexibility while maintaining efficiency. Future research should examine the effectiveness of these integration approaches through practical deployment studies, ensuring that DECLARUI not only generates high-quality code but also enhances existing development processes.

5.3 Threat to Validity

Internal Validity. The primary threats stem from our evaluation metrics, CV technique selection, and method implementation across different frameworks. CLIP and SSIM scores, used for UI similarity assessment, may not fully capture human-perceived UI similarity, potentially overlooking subtle differences crucial to user experience. While our combination of Grounding DINO and SAM achieves a 92.3% accuracy rate in UI component recognition (based on the evaluation of 10 randomly sampled pages), alternative CV approaches might yield different results. This particular choice of CV techniques, though effective for our purposes, represents one possible solution among

many, and the impact of different CV approaches on overall system performance remains to be fully explored. PTG Coverage, employed for navigation logic accuracy, ensures implementation of expected page transitions but may not reflect the actual user experience of navigation fluidity and intuitiveness. Fortunately, our user study confirmed the validity of our evaluations, demonstrating that DECLARUI indeed outperformed the baselines.

External Validity. Our study faces several threats related to the performance of MLLM influenced by prompt design and the experience levels of user study participants. Different prompts could lead to varying results, making it challenging to isolate the true impact of our approach. Although we selected participants with similar backgrounds, inherent biases may still affect the results. Fortunately, Figure 7 showed that participants' opinions on the same app were relatively consistent. Another threat involves the scope and size of our UI design sample. Despite efforts to include diverse UI designs within the mobile app domain, our sample may not fully represent all mobile app UI types. The limited sample size of 50 UI design sets, due to MLLMs' token cost constraints, may not capture the full variability of mobile app UIs. Future research should expand the size and diversity of the UI sample to ensure comprehensive coverage of mobile app UI types.

6 Related Work

Mobile App UI-to-Code Generation. Mobile app UI-to-code translation has seen various approaches in recent years. [4, 14] were limited by traditional UI frameworks and lacked support for multi-page apps. [52] explored code generation from natural language descriptions but struggled with complex designs and fine-grained component recognition. [7, 9, 10, 57] introduced intermediate representations to aid UI code generation, yet still required significant manual effort to produce the final code. Unlike previous works, DECLARUI is the first to focus on serving declarative UI frameworks, combining advanced CV techniques for precise component recognition with PTG-based inter-page logic capture to enable fully automated, multi-page mobile app UI code generation.

Image-to-Code Generation. Image-to-code generation research, particularly in UI generation, has advanced significantly through three main approaches: Deep Learning-based methods [3–5], CV-based methods [15, 31, 49], and Multimodal Learning methods [35, 37, 41]. [4] pioneered deep learning techniques to generate code from GUI screenshots. CV-based methods, exemplified by [31], converted hand-drawn sketches into HTML code, while [49] advanced this by leveraging GPT-4V and DALL·E 3 to transform screenshots into functional code. Addressing limitations of single-modal approaches, Multimodal Learning methods emerged. [37] presented a multimodal autoregressive model for various vision and language tasks, including code generation from images. [35] demonstrated how combining different data types can enhance model robustness and versatility, potentially improving accuracy and contextual understanding in code generation tasks.

7 Conclusion

In this paper, we presented DECLARUI, an approach that combines CV techniques, MLLMs, and iterative compiler-driven optimization to generate high-quality declarative UI code from design mockups and screenshots. Our evaluation demonstrates that DECLARUI outperforms state-of-the-art methods across multiple metrics, including PTG coverage rates, visual similarity, average compilation iteration counts, and compilation success rates. DECLARUI's multi-framework compatibility is evidenced by its ability to generalize across React Native, Flutter, and ArkUI. The positive feedback from professional developers in our user study underscores DECLARUI's practical value. By addressing challenges in component recognition, interactive logic understanding, and code reliability, DECLARUI represents a significant advancement in bridging UI design and implementation.

Data Availability and Ethics

Open Science. The artifact is available at <https://github.com/security-pride/DeclarUI>.

Ethics Considerations. Our study adhered to ethical research practices. Participants provided informed consent and their data was anonymized. The \$50 compensation was deemed fair, and the 20-minute time limit per app was reasonable. The study design was reviewed and approved by our institution prior to implementation.

Acknowledgments

This work was supported in part by the Key R&D Program of Hubei Province (2023BAB017, 2023BAB079) and CCF-Huawei Populus Grove Fund.

References

- [1] Anthropic Claude 3.5 [n. d.]. Introducing Claude 3.5 Sonnet | Anthropic. <https://www.anthropic.com/news/claude-3-5-sonnet>
- [2] Shushan Arakelyan, Rocktim Jyoti Das, Yi Mao, and Xiang Ren. 2023. Exploring Distributional Shifts in Large Language Models for Code Analysis. *CoRR* abs/2303.09128 (2023). doi:10.48550/ARXIV.2303.09128 arXiv:2303.09128
- [3] R. Archana and P. S. Eliahim Jeevaraj. 2024. Deep learning models for digital image processing: a review. *Artif. Intell. Rev.* 57, 1 (2024), 11. doi:10.1007/S10462-023-10631-Z
- [4] Tony Beltramelli. 2017. pix2code: Generating Code from a Graphical User Interface Screenshot. *CoRR* abs/1705.07962 (2017). arXiv:1705.07962 <http://arxiv.org/abs/1705.07962>
- [5] Bo Cai, Jian Luo, and Zhen Feng. 2023. A novel code generator for graphical user interfaces. *Scientific Reports* 13 (2023). <https://api.semanticscholar.org/CorpusID:265349539>
- [6] Jinyuan Chang, Jing He, Jian Kang, and Mingcong Wu. 2023. Statistical inferences for complex dependence of multimodal imaging data. arXiv:2303.03582 [stat.ME] <https://arxiv.org/abs/2303.03582>
- [7] Chunyang Chen, Ting Su, Guozhu Meng, Zhenchang Xing, and Yang Liu. 2018. From UI design image to GUI skeleton: a neural machine translator to bootstrap mobile GUI implementation. In *Proceedings of the 40th International Conference on Software Engineering, ICSE 2018, Gothenburg, Sweden, May 27 - June 03, 2018*, Michel Chaudron, Ivica Crnkovic, Marsha Chechik, and Mark Harman (Eds.). ACM, 665–676. doi:10.1145/3180155.3180240
- [8] Fuxiang Chen, Fatemeh H. Fard, David Lo, and Timofey Bryksin. 2022. On the Transferability of Pre-trained Language Models for Low-Resource Programming Languages. *CoRR* abs/2204.09653 (2022). doi:10.48550/ARXIV.2204.09653 arXiv:2204.09653
- [9] Jieshan Chen, Amanda Swearngin, Jason Wu, Titus Barik, Jeffrey Nichols, and Xiaoyi Zhang. 2022. Towards Complete Icon Labeling in Mobile Applications. In *CHI '22: CHI Conference on Human Factors in Computing Systems, New Orleans, LA, USA, 29 April 2022 - 5 May 2022*, Simone D. J. Barbosa, Cliff Lampe, Caroline Appert, David A. Shamma, Steven Mark Drucker, Julie R. Williamson, and Koji Yatani (Eds.). ACM, 387:1–387:14. doi:10.1145/3491102.3502073
- [10] Jieshan Chen, Mulong Xie, Zhenchang Xing, Chunyang Chen, Xiwei Xu, Liming Zhu, and Guoqiang Li. 2020. Object Detection for Graphical User Interface: Old Fashioned or Deep Learning or a Combination? *CoRR* abs/2008.05132 (2020). arXiv:2008.05132 <https://arxiv.org/abs/2008.05132>
- [11] Liuqing Chen, Yunnong Chen, Shuhong Xiao, Yaxuan Song, Lingyun Sun, Yankun Zhen, Tingting Zhou, and Yan-fang Chang. 2023. EGFE: End-to-end Grouping of Fragmented Elements in UI Designs with Multimodal Learning. *CoRR* abs/2309.09867 (2023). doi:10.48550/ARXIV.2309.09867 arXiv:2309.09867
- [12] Liguang Chen, Qi Guo, Hongrui Jia, Zhengran Zeng, Xin Wang, Yijiang Xu, Jian Wu, Yidong Wang, Qing Gao, Jindong Wang, Wei Ye, and Shikun Zhang. 2024. A Survey on Evaluating Large Language Models in Code Generation Tasks. *CoRR* abs/2408.16498 (2024). doi:10.48550/ARXIV.2408.16498 arXiv:2408.16498
- [13] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating Large Language Models Trained on Code. *CoRR* abs/2107.03374 (2021). arXiv:2107.03374 <https://arxiv.org/abs/2107.03374>

- [14] Wen-Yin Chen, Pavol Podstreleny, Wen-Huang Cheng, Yung-Yao Chen, and Kai-Lung Hua. 2022. Code generation from a graphical user interface via attention-based encoder-decoder model. *Multim. Syst.* 28, 1 (2022), 121–130. doi:10.1007/S00530-021-00804-7
- [15] Daniel de Souza Baulé, Christiane Gresse von Wangenheim, Aldo von Wangenheim, and Jean C. R. Hauck. 2020. Recent Progress in Automated Code Generation from GUI Images Using Machine Learning Techniques. *J. Univers. Comput. Sci.* 26, 9 (2020), 1095–1127. http://www.jucs.org/jucs_26_9/recent_progress_in_automated
- [16] Meghdad Dehghan, Jie JW Wu, Fatemeh H. Fard, and Ali Ouni. 2024. MergeRepair: An Exploratory Study on Merging Task-Specific Adapters in Code LLMs for Automated Program Repair. *CoRR* abs/2408.09568 (2024). doi:10.48550/ARXIV.2408.09568 arXiv:2408.09568
- [17] Jiangyi Deng, Xinfeng Li, Yanjiao Chen, Yijie Bai, Haiqin Weng, Yan Liu, Tao Wei, and Wenyuan Xu. 2024. RACONTEUR: A Knowledgeable, Insightful, and Portable LLM-Powered Shell Command Explainer. *CoRR* abs/2409.02074 (2024). doi:10.48550/ARXIV.2409.02074 arXiv:2409.02074
- [18] Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. 2023. Self-collaboration Code Generation via ChatGPT. *CoRR* abs/2304.07590 (2023). doi:10.48550/ARXIV.2304.07590 arXiv:2304.07590
- [19] Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A Survey on RAG Meeting LLMs: Towards Retrieval-Augmented Large Language Models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2024, Barcelona, Spain, August 25-29, 2024*, Ricardo Baeza-Yates and Francesco Bonchi (Eds.). ACM, 6491–6501. doi:10.1145/3637528.3671470
- [20] Figma. 2024. Figma: Collaborative Interface Design Tool. <https://www.figma.com/> Accessed: 2024-09-05.
- [21] Flatirons. 2024. Popularity of Flutter vs. React Native in 2024. <https://flatirons.com/blog/popularity-of-flutter-vs-react-native-2024/> Accessed: 2024-09-05.
- [22] Flutter. 2024. User Interface (UI) in Flutter. <https://docs.flutter.dev/ui> Accessed: 2024-09-05.
- [23] Shuzheng Gao, Xin-Cheng Wen, Cuiyun Gao, Wenxuan Wang, and Michael R. Lyu. 2023. Constructing Effective In-Context Demonstration for Code Intelligence Tasks: An Empirical Study. *CoRR* abs/2304.07575 (2023). doi:10.48550/ARXIV.2304.07575 arXiv:2304.07575
- [24] Google Play Store. 2024. Amazon Shopping. <https://play.google.com/store/apps/details?id=com.amazon.mShop.android.shopping&hl=en> Accessed: 2024-09-05.
- [25] Google Play Store. 2024. Google Play Store: Apps Section. <https://play.google.com/store/apps?hl=en> Accessed: 2024-09-05.
- [26] Google Play Store. 2024. Lazada - Online Shopping App. <https://play.google.com/store/apps/details?id=com.lazada.android&hl=en> Accessed: 2024-09-05.
- [27] Jian Gu, Pasquale Salza, and Harald C. Gall. 2022. Assemble Foundation Models for Automatic Code Summarization. *CoRR* abs/2201.05222 (2022). arXiv:2201.05222 <https://arxiv.org/abs/2201.05222>
- [28] Yi Gui, Zhen Li, Yao Wan, Yemin Shi, Hongyu Zhang, Yi Su, Shaoling Dong, Xing Zhou, and Wenbin Jiang. 2024. VISION2UI: A Real-World Dataset with Layout for Code Generation from UI Designs. *CoRR* abs/2404.06369 (2024). doi:10.48550/ARXIV.2404.06369 arXiv:2404.06369
- [29] Huawei. 2024. ArkUI Overview: HarmonyOS. <https://developer.huawei.com/consumer/cn/doc/harmonyos-guides-V5/arkui-overview-V5> Accessed: 2024-09-05.
- [30] Increment. 2024. The Shift to Declarative UI. <https://increment.com/mobile/the-shift-to-declarative-ui/> Accessed: 2024-09-05.
- [31] Vanita Jain, Piyyush Agrawal, Subham Banga, Rishabh Kapoor, and Shashwat Gulyani. 2019. Sketch2Code: Transformation of Sketches to UI in Real-time Using Deep Neural Network. *CoRR* abs/1910.08930 (2019). arXiv:1910.08930 <http://arxiv.org/abs/1910.08930>
- [32] Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and Sunghun Kim. 2024. A Survey on Large Language Models for Code Generation. *CoRR* abs/2406.00515 (2024). doi:10.48550/ARXIV.2406.00515 arXiv:2406.00515
- [33] Dharma KC and Clayton T. Morrison. 2023. Neural Machine Translation for Code Generation. arXiv:2305.13504 [cs.CL] <https://arxiv.org/abs/2305.13504>
- [34] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloé Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross B. Girshick. 2023. Segment Anything. *CoRR* abs/2304.02643 (2023). doi:10.48550/ARXIV.2304.02643 arXiv:2304.02643
- [35] Changwon Kwak, Pilsu Jung, and Seonah Lee. 2023. A Multimodal Deep Learning Model Using Text, Image, and Code Data for Improving Issue Classification Tasks. *Applied Sciences* (2023). <https://api.semanticscholar.org/CorpusID:261157795>
- [36] Yuanchun Li, Ziyue Yang, Yao Guo, and Xiangqun Chen. 2017. DroidBot: a lightweight UI-guided test input generator for Android. In *Proceedings of the 39th International Conference on Software Engineering, ICSE 2017, Buenos Aires, Argentina, May 20-28, 2017 - Companion Volume*, Sebastián Uchitel, Alessandro Orso, and Martin P. Robillard (Eds.). IEEE Computer Society, 23–26. doi:10.1109/ICSE-C.2017.8

- [37] Dongyang Liu, Shitian Zhao, Le Zhuo, Weifeng Lin, Yu Qiao, Hongsheng Li, and Peng Gao. 2024. Lumina-mGPT: Illuminate Flexible Photorealistic Text-to-Image Generation with Multimodal Generative Pretraining. *CoRR abs/2408.02657* (2024). doi:10.48550/ARXIV.2408.02657 arXiv:2408.02657
- [38] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, and Lei Zhang. 2023. Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection. *CoRR abs/2303.05499* (2023). doi:10.48550/ARXIV.2303.05499 arXiv:2303.05499
- [39] Yizhou Liu, Pengfei Gao, Xinchun Wang, Jie Liu, Yexuan Shi, Zhao Zhang, and Chao Peng. 2024. MarsCode Agent: AI-native Automated Bug Fixing. *CoRR abs/2409.00899* (2024). doi:10.48550/ARXIV.2409.00899 arXiv:2409.00899
- [40] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin B. Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, Ming Gong, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie Liu. 2021. CodeXGLUE: A Machine Learning Benchmark Dataset for Code Understanding and Generation. *CoRR abs/2102.04664* (2021). arXiv:2102.04664 <https://arxiv.org/abs/2102.04664>
- [41] Sanbi Luo. 2021. A Survey on Multimodal Deep Learning for Image Synthesis: Applications, methods, datasets, evaluation metrics, and results comparison. In *ICIAI 2021: 2021 the 5th International Conference on Innovation in Artificial Intelligence, Xiamen, China, March 5-8, 2021*. ACM, 108–120. doi:10.1145/3461353.3461388
- [42] MobileAppDaily. 2024. Mobile App Development Process: Step-by-Step Guide. <https://www.mobileappdaily.com/knowledge-hub/mobile-app-development-process> Accessed: 2024-09-05.
- [43] Yun nong Chen, Yan kun Zhen, Chu ning Shi, Jia zhi Li, Liu qing Chen, Ze jian Li, Ling yun Sun, Ting ting Zhou, and Yan fang Chang. 2022. UI Layers Merger: Merging UI layers via Visual Learning and Boundary Prior. arXiv:2206.13389 [cs.CV] <https://arxiv.org/abs/2206.13389>
- [44] OpenAI Hello GPT-4o [n. d.]. Hello GPT-4o | OpenAI. <https://openai.com/index/hello-gpt-4o/>
- [45] Ciyuan Peng, Feng Xia, Mehdi Naseriparsa, and Francesco Osborne. 2023. Knowledge Graphs: Opportunities and Challenges. *CoRR abs/2303.13948* (2023). doi:10.48550/ARXIV.2303.13948 arXiv:2303.13948
- [46] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-Shot Text-to-Image Generation. *CoRR abs/2102.12092* (2021). arXiv:2102.12092 <https://arxiv.org/abs/2102.12092>
- [47] React Native. 2024. Getting Started with React Native. <https://reactnative.dev/docs/getting-started> Accessed: 2024-09-05.
- [48] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Raphael Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. 2022. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. *CoRR abs/2205.11487* (2022). doi:10.48550/ARXIV.2205.11487 arXiv:2205.11487
- [49] ScreenshotToCode2024 [n. d.]. Screenshot to Code: Transform Screenshots into Code! <https://supertools.therundown.ai/content/screenshot-to-code>
- [50] Chenglei Si, Yanzhe Zhang, Zhengyuan Yang, Ruibo Liu, and Diyi Yang. 2024. Design2Code: How Far Are We From Automating Front-End Engineering? *CoRR abs/2403.03163* (2024). doi:10.48550/ARXIV.2403.03163 arXiv:2403.03163
- [51] The Decoder. 2024. ChatGPT Guide: Effective Prompt Strategies. <https://the-decoder.com/chatgpt-guide-prompt-strategies/> Accessed: 2024-09-05.
- [52] Yuxuan Wan, Chaozheng Wang, Yi Dong, Wenxuan Wang, Shuqing Li, Yintong Huo, and Michael R. Lyu. 2024. Automatically Generating UI Code from Screenshot: A Divide-and-Conquer-Based Approach. *CoRR abs/2406.16386* (2024). doi:10.48550/ARXIV.2406.16386 arXiv:2406.16386
- [53] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE Trans. Image Process.* 13, 4 (2004), 600–612. doi:10.1109/TIP.2003.819861
- [54] Wolfpack Digital. 2024. Declarative UI: A New Way of Building Mobile Apps. <https://www.wolfpack-digital.com/blogposts/declarative-ui-new-way-of-building-mobile-apps> Accessed: 2024-09-05.
- [55] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. 2023. Visual ChatGPT: Talking, Drawing and Editing with Visual Foundation Models. *CoRR abs/2303.04671* (2023). doi:10.48550/ARXIV.2303.04671 arXiv:2303.04671
- [56] Jason Wu, Eldon Schoop, Alan Leung, Titus Barik, Jeffrey P. Bigham, and Jeffrey Nichols. 2024. UICoder: Finetuning Large Language Models to Generate User Interface Code through Automated Feedback. *CoRR abs/2406.07739* (2024). doi:10.48550/ARXIV.2406.07739 arXiv:2406.07739
- [57] Jason Wu, Xiaoyi Zhang, Jeffrey Nichols, and Jeffrey P. Bigham. 2021. Screen Parsing: Towards Reverse Engineering of UI Models from Screenshots. *CoRR abs/2109.08763* (2021). arXiv:2109.08763 <https://arxiv.org/abs/2109.08763>
- [58] Shuhong Xiao, Yunnong Chen, Jiazhi Li, Liuqing Chen, Lingyun Sun, and Tingting Zhou. 2024. Prototype2Code: End-to-end Front-end Code Generation from UI Design Prototypes. *CoRR abs/2405.04975* (2024). doi:10.48550/ARXIV.2405.04975 arXiv:2405.04975

- [59] Kangwei Xu, Grace Li Zhang, Xunzhao Yin, Cheng Zhuo, Ulf Schlichtmann, and Bing Li. 2024. Automated C/C++ Program Repair for High-Level Synthesis via Large Language Models. *CoRR* abs/2407.03889 (2024). doi:10.48550/ARXIV.2407.03889 arXiv:2407.03889
- [60] Boyang Yang, Haoye Tian, Weiguo Pian, Haoran Yu, Haitao Wang, Jacques Klein, Tegawendé F. Bissyandé, and Shunfu Jin. 2024. CREF: An LLM-based Conversational Software Repair Framework for Programming Tutors. *CoRR* abs/2406.13972 (2024). doi:10.48550/ARXIV.2406.13972 arXiv:2406.13972
- [61] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023. The Dawn of LLMs: Preliminary Explorations with GPT-4V(ision). *CoRR* abs/2309.17421 (2023). doi:10.48550/ARXIV.2309.17421 arXiv:2309.17421

Received 2024-09-12; accepted 2025-01-14