# AutoGameUI: Constructing High-Fidelity Game UIs via Multimodal Learning and Interactive Web-Based Tool

ZHONGLIANG TANG, Tencent TiMi L1 Studio, China

MENGCHEN TAN, Tencent TiMi L1 Studio, China

FEI XIA, Tencent TiMi L1 Studio, China

QINGRONG CHENG, Tencent TiMi L1 Studio, China

HAO JIANG, Tencent TiMi L1 Studio, China

YONGXIANG ZHANG, Tencent TiMi L1 Studio, China

We introduce an innovative system, AutoGameUI, for efficiently constructing cohesive user interfaces in game development. Our system is the first to address the coherence issue arising from integrating inconsistent UI and UX designs, typically leading to mismatches and inefficiencies. We propose a two-stage multimodal learning pipeline to obtain comprehensive representations of both UI and UX designs, and to establish their correspondences. Through the correspondences, a cohesive user interface is automatically constructed from pairwise designs. To achieve high-fidelity effects, we introduce a universal data protocol for precise design descriptions and cross-platform applications. We also develop an interactive web-based tool for game developers to facilitate the use of our system. We create a game UI dataset from actual game projects and combine it with a public dataset for training and evaluation. Our experimental results demonstrate the effectiveness of our system in maintaining coherence between the constructed interfaces and the original designs.

## 1 INTRODUCTION

With the widespread adoption of personal computers and smart mobile devices, the gaming industry has generated considerable economic value and cultural influence around the world. However, as a critical step in game development, game UI development still faces various challenges that need to be addressed. Compared to the numerous literature [10, 16, 21, 45, 47] in the computational UI domain for mobile apps and web pages, few academic research and industrial applications discovered and truly resolved the issues encountered in game UI development until now.

In conventional game UI development, two specialized teams collaborate in parallel to create a cohesive game user interface (GameUI) that can be played in the game engine. The first team, made up of UI designers or artists, focuses on

creating appealing and expressive visual elements to achieve diverse game styles. The second team, composed of UX designers or game engineers, pays more attention to emphasizing usability and functionality by designing effective logical structures. As illustrated in Fig .1, the distinct focus areas naturally lead to conflicting design elements, varying style patterns, and inconsistent layouts, making it difficult to construct the GameUI that requires a seamless integration of the aesthetic UI design and the functional UX design. Game project teams often employ a great deal of manpower to manually bridge the gap and construct the GameUI in the game engine, while the final quality cannot be guaranteed owing to the lack of industry-wide construction standards. Moreover, both UI and UX designs need to be frequently updated according to the latest game design script, which means that the previously completed GameUI must be revised. Regardless of developers' experience, manually handling such cases is extremely time-consuming and usually causes delays in the development timeline.
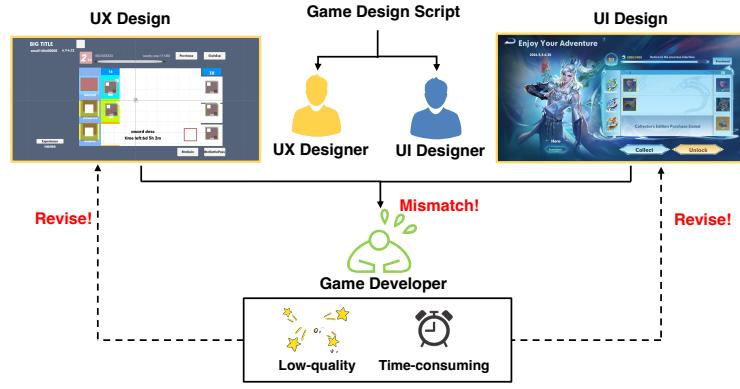


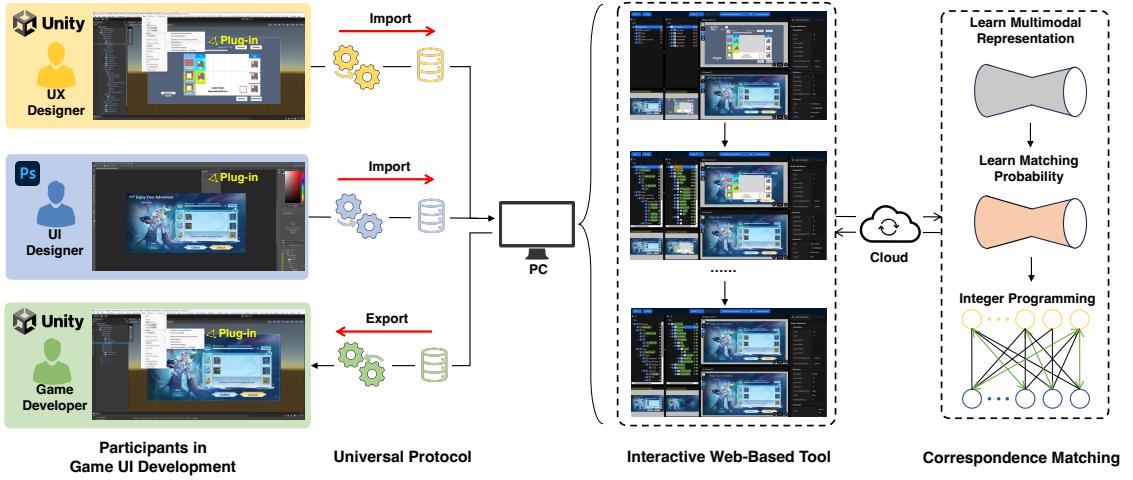Fig. 1. Example of the coherence issue in conventional game UI development.



Fig. 2. System framework.

To alleviate the bottleneck posed above, this paper introduces a novel automated system for GameUI construction, which is named as AutoGameUI. As depicted in Fig .2, our system consists of three components: automatic correspondence matching, universal data protocol, and interactive web-based tool. The automatic correspondence matching is composed of two multimodal models and a flexible integer programming algorithm. The multimodal models are based on the transformer model [42] and are learned in two separate stages, with the first stage focusing on representing UI and UX designs in a latent space, and the second stage aiming to estimate the matching probability between them. Previous studies [3, 16, 23, 51] typically examined the representation of graphical UIs from a limited number of perspectives and failed to integrate more comprehensive attributes. In contrast, our method inclusively considers various aspects such as spatial layout, functional semantics, textual content, hierarchy, and even rendering order, except image texture, which is useless in our task. Besides, our models adopt relative positional encoding to accommodate graphical UIs with varying lengths, which is different from some relevant attempts [13–15, 49, 50]. These works employed absolute positional encoding in the transformer model which can struggle to adjust the performance of original training data in extended test data. Furthermore, some works [30, 32, 46, 48] tried to estimate the correspondences between pairs of graphical UIs, relying on one-to-one element similarity or Hungarian algorithm [17] for bipartite matching. However, there exist complex correspondence rules present in real-world UI and UX designs, including but not limited to one-to-one, and many-to-one mapping, which conflicts with the existing methods. In our method, we first identify and summarize all correspondence rules between pairwise UI and UX designs. We then exploit grouped cross-attention modules to calculate one-to-group similarity with reduced dimensions. By integrating the similarity with the hierarchical and rendering constraints, we formulate an optimization problem that is effectively solved using integer programming to achieve optimal performance. Following prior works [2, 31, 38], we design a universal protocol to describe the UI and UX designs, enabling the data stream across different platforms. By utilizing the universal data protocol and the optimal correspondences, we can easily achieve high-fidelity effects in the GameUI construction. To streamline the usage of our system and provide users with intuitive feedback and guidance, we develop an interactive web-based tool to give a smooth experience during the construction process. The web-based tool also offers the ability of data annotation, enabling us to collect and build a specialized dataset from real-world game UIs, which aids in training and evaluating our approach. In summary, this paper makes the following contributions:

- We present a complete system named AutoGameUI, which innovatively utilizes artificial intelligence techniques for high-quality GameUI construction. This system breaks through the bottleneck limited by traditional manual construction and significantly accelerates the game UI development workflow.
- We propose an efficient multimodal learning pipeline that acquires the comprehensive representations of UI and UX designs and estimates their optimal correspondences. To enhance efficiency and accuracy, we incorporate grouped cross-attention modules with hierarchical and rendering constraints into the estimation process.
- We introduce a universal data protocol and an accompanying interactive web-based tool for our system. This not only supports cross-platform applications but also offers interactive features to assist users in rapid construction, thereby improving the overall user experience.
- We build a game UI dataset with annotated UI and UX designs for model training and evaluation.

## 2 RELATED WORK

### 2.1 Representation of UIs

How to extract and represent the data attributes of graphical UIs is crucial for downstream tasks such as UI recognition [21], UI completion [16], UI generation [15], UI structuring [47], etc. The essential attributes include spatial layout, semantics, visual appearance, textual content, view hierarchy, and rendering order. Most of the existing methods for UI representation exhibit equal treatment toward spatial layout and semantics but overlook others. For example, LayoutTrans [13] developed a transformer model for layout generation based on the 2D bounding box and the semantic label of each element. LayoutDM [15] also trained a transformer model as the backbone in the reverse diffusion process, to formulate the mapping between latent representation and original layout. Other efforts [16, 44, 51] attempted to enhance representation performance by integrating visual appearance and textual content, but they still fell short in comprehending the hierarchy and rendering order. Some studies [3, 24] delved into the modeling of the hierarchy, while the rendering order has not been explored yet. Few works extend adequate consideration to the rendering order for visual consistencies. This neglect is associated with the inherent limitations of public datasets like the RICO [10] and PubLayNet [56]. In these datasets, the layout of UI elements does not involve complex occlusion relationships. Moreover, the datasets only provide RGB screenshots, not the raw RGBA images used in the rendering pipeline.

Compared to previous works, we incorporate spatial layout, semantics, textual content, view hierarchy, and rendering order to comprehend the graphical UIs, which lay a robust foundation for subsequent computational tasks. Besides, we create a real-world game UI dataset based on actual game projects to measure and demonstrate the effectiveness of our method.

### 2.2 Correspondence between UIs

Identifying corresponding elements between pairs of UIs has been investigated for some time now. Earlier studies [18, 19] tried to establish the correspondences between hierarchical web pages with the probabilistic optimization, and automatically transfer design and contents from one to another side. Dayama *et al.* [8] presented an integer programming algorithm to optimize the correspondences and developed an interactive tool for layout transfer. LayoutBlending [48] recursively found the optimal correspondences between two hierarchical structures, which maintains the consistencies of hierarchy and avoids unsuitable rendering effects in layout blending applications. However, this work solved the correspondence problem using the Hungarian algorithm [17], which is restricted by the bipartite matching setting. Otani *et al.* [30] introduced a criterion to compute and measure the similarity between arbitrary pairs in layout generations. Although these works do not rely on extensive training data, the integrity of handcrafted features, heuristics, constraints, and optimization algorithms still affect their generalization performance in practical applications.

Some works also applied deep learning methods for correspondence estimation. Patil *et al.* [32, 33] separately introduced a recursive neural network and a graph neural network to generate layouts and measure their structural similarity. The graph neural network, named LayoutGMN, utilizes the fully-connected graph to organize layout elements and employs metric learning to determine the similarity of both the entire layout and individual elements. Wu *et al.* [46] carried out a multimodal transformer model to generate the latent representations of layout elements and then infer correspondences across two UIs. This work pruned the dissimilar matches in the optimization process and avoided making sub-optimal decisions. Most of these works computed the similarity matrix between pairwise elements for correspondences. Once numerous elements are designed within the two homogeneous structures, incorporating
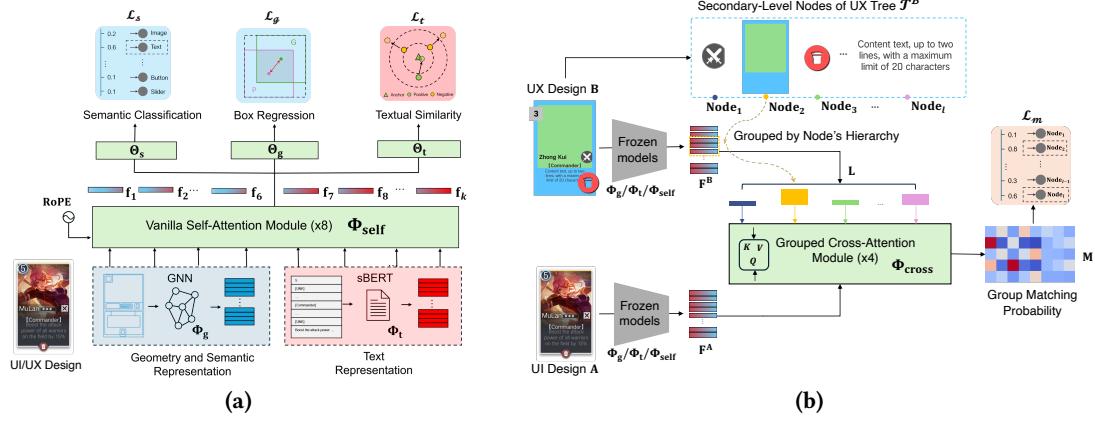
Fig. 3. Illustration of two-stage multimodal learning. (a) Given a UI or UX design with geometric layout, semantics, and textual contents, we combine a graph neural network, a pre-trained language model, and a transformer model equipped with vanilla self-attention modules, to acquire the multimodal representation of each element. (b) After the first stage of learning, we separately generate the multimodal representations of UI and UX designs. We group the representations of UX design based on its hierarchy and train the second transformer model equipped with grouped cross-attention modules, to calculate the node-to-group matching probability.

additional constraints into the optimization process may be extremely time-consuming. In contrast, our method exhibits advanced performance in generalization and time efficiency.

## 2.3 Artificial Intelligence in Gaming Industry

AI technology has been increasingly applied in the gaming industry, involving game-playing bots and procedural content generation. Prior works [11, 43, 52, 53] utilized reinforced AI to develop and manage agents based on player data, enabling collaborative playing or controlling Non-Playable Characters (NPCs) to enhance the player experience. Besides, more works employed generative AI to assist in the production of game assets, such as 2D graphics [39, 55], 3D character and scene models [26, 35], music and voices [1, 41], facial expression and animations [6, 7, 28].

Some efforts [12, 22, 54, 57] appeared to intersect with human-computer interaction and game development, but they chiefly centered around gameplay design and interaction experience, with limited relevance to game UI development itself. The applications of AI in game UI development remain a relatively unexplored field, with both academic research and industrial systems still in their early stages. Existing intelligent tools for UI development primarily focus on web pages or mobile apps such as Uizard[1] and Framer[2]. These tools can generate UI/UX prototypes based on prompts, screenshots, or wireframes, which are not as urgently required in game UI development. Moreover, these tools often perform poorly in real-world game UIs which have highly stylized and aesthetic design. In comparison, our work pioneers an automated system for GameUI construction, offering a real-time solution that directly addresses the core bottleneck encountered in game UI development.

---

[1]https://uizard.io
[2]https://www.framer.com

## 3 GAMEUI CONSTRUCTION

Constructing GameUI fundamentally involves estimating correspondences between UI and UX designs and integrating the visual appearance of UI design into UX design through these correspondences. This section provides an overview of the proposed system, as illustrated in Fig. 2 and Fig. 3. Given a pair of UI and UX designs, we first clarify the data attributes by factorizing them into multiple modalities with discrete parameters. Next, we employ a self-supervised learning approach to generate comprehensive multimodal representations from the discrete parameters. Utilizing the multimodal representations of UI and UX designs, we formulate the GameUI construction task as a constrained correspondence matching problem. To effectively address this problem, we introduce a novel cross-attention module and flexible integer programming. The novel cross-attention module learns the matching probability for identifying potential matches and the integer programming acquires optimal matching results. With the established correspondences, we design a universal protocol to facilitate data description and attribute integration across UI and UX designs. Lastly, we introduce a web-based tool to help users interactively construct GameUI.

### 3.1 Learning Multimodal Representation

*3.1.1 Data Attributes.* In game development, UI and UX designs are typically cached and managed in a tree structure, with some common attributes like spatial layout, image texture, textual content, hierarchy, rendering order, and semantics. Spatial layout is essential for describing the 2D geometry of leaf and non-leaf nodes. Image texture and textual content are the renderable attributes exclusively presented in leaf nodes. Unlike previous works [16, 21, 46], we exclude image texture to avoid uncontrollable noise, as the visual appearance of UI and UX designs are usually different. Hierarchy reveals the depth and arrangement of nodes that high-level nodes have a broader scope, while lower-level nodes have a narrower scope. Rendering order defines the display priority that nodes are rendered from top to bottom or conversely on the screen. Semantics indicate the interaction of the nodes in the tree structure such as TEXT, IMAGE, BUTTON, SLIDER, LIST, etc.

Mathematically, we can denote the tree structure as $\mathcal{T} = (\mathcal{V}, \mathcal{E})$ where $\mathcal{V}$ is the node set and $\mathcal{E}$ is the edge set. Let $\mathbf{n}$ be the number of nodes, the $k^{th}$ node can be defined as:

$$v_k = (g_k, s_k, t_k), \tag{1}$$

where $1 \leq k \leq \mathbf{n}$. $g_k = [\frac{x_k}{W}, \frac{y_k}{H}, \frac{w_k}{W}, \frac{h_k}{H}]$ is a four-dimensional parameter that details the 2D geometry of $v_k$. Each value of $g_k$ is normalized with the height and width of the entire layout, and helps improve scale stability. $s_k$ is an integer that indicates the semantic label. $t_k$ is the textual content and can be transformed into a fixed-length integer sequence with a particular vocabulary.

Inspired by [16, 21, 29, 33], we also explore the integration of the tree structure and the fully-connected graph structure into the edge. We can define the directed edge $e_{ij}$ from node $v_i$ to node $v_j$ as:

$$e_{ij} = (\Delta g_{ij}, \Delta h_{ij}, \Delta r_{ij}), \tag{2}$$

where $1 \leq i \leq \mathbf{n}, 1 \leq j \leq \mathbf{n}, i \neq j$. $\Delta g_{ij}$ indicates the spatial relationship between $v_i$ and $v_j$. We adopt the identical formula utilized in [29, 33], which measures the translation, IoU area, aspect ratio, and orientation. $\Delta h_{ij} \in \{0, 1, 2\}$ indicates the hierarchical relationship in $\mathcal{T}$ and the values 0, 1, 2 correspond to three conditions respectively: $v_i$ is the ancestor of $v_j$; $v_i$ is the descendant of $v_j$; $v_i$ does not have any hierarchical relationship with $v_j$. $\Delta r_{ij} \in \{0, 1, 2\}$ indicates

the rendering relationship, and similar to $\Delta h_{ij}$, the values 0, 1, 2 also correspond to three conditions: $v_i$ is overlapped with $v_j$ and rendered after that; $v_i$ is overlapped with $v_j$ but rendered before that; $v_i$ is not overlapped with $v_j$.

With the formulation of the tree structure, the UI and UX designs can be represented as $\mathcal{T}^A$ and $\mathcal{T}^B$. We should note that there exist differences in semantics and hierarchical relationships between $\mathcal{T}^A$ and $\mathcal{T}^B$. For semantics, $s_k^A$ is limited to two low-level semantic categories - IMAGE and TEXT, which are significantly fewer than $s_k^B$. This is because $v_k^A$ can only represent the leaf nodes of $\mathcal{T}^A$ where the non-leaf nodes serve as the containers of renderable resources and lack high-level semantics in interaction. On the contrary, $v_k^B$ can represent both leaf and non-leaf nodes of $\mathcal{T}^B$ where the non-leaf nodes are usually interactive widgets or components. For hierarchical relationships, $\Delta h_{ij}^B = 0$ can denote that $v_i^B$ is the ancestor of $v_j^B$ and likely to be enclosed by a larger rectangle box and have higher-level semantics. However, we cannot explain $\Delta h_{ij}^A = 0$ in the same way because $v_i^A$ and $v_j^A$ are limited to the leaf nodes. To overcome this limitation, we specifically simulate the hierarchical relationship $\Delta h_{ij}^A$ based on a simple geometric heuristic that once the bounding box of $v_j^A$ is included in $v_i^A$, $v_i^A$ is the ancestor of $v_j^A$.

### 3.1.2 Multimodal Representation.

Upon defining the tree representation with discrete parameters, we decide to seek a continuous representation with fixed-size length, which is more suitable for multimodal learning. To meet the requirements, we incorporate a graph neural network $\Phi_{\mathbf{g}}$ to embed the spatial parameter $g_k$ and the semantic parameter $s_k$. Besides, a pre-trained language model $\Phi_{\mathbf{t}}$ is employed to embed the content parameter $t_k$. Notably, the edge parameters $\Delta h_{ij}$ and $\Delta r_{ij}$ are not embedded into continuous feature vectors, which will be discussed later in Section 3.2.3.

The graph neural network $\Phi_{\mathbf{g}}$ is broadly similar to [29, 33] and has been proven effective in graph encoding tasks. We use $\Phi_{\mathbf{g}}$ to embed the node parameters $s_k$, $g_k$, and the adjacent edge parameters $\{\Delta g_{k1}, \ldots, \Delta g_{kj}, \ldots, \Delta g_{kn}\}$ for yielding a 128-dimensional feature vector $\mathbf{f}_k^g$ as:

$$\mathbf{f}_k^g = \Phi_{\mathbf{g}}(g_k, s_k, \{\Delta g_{k1} \ldots, \Delta g_{kj}, \ldots, \Delta g_{kn}\}), \tag{3}$$

The pre-trained language model $\Phi_{\mathbf{t}}$ is the sentence BERT model [36]. With the vocabulary of sentence BERT, every text sentence is segmented into multiple tokens and transformed into an integer sequence. After feeding the integer sequence into the encoder-only model, the original text sentence is embedded as a 768-dimensional feature vector. It is important to note that some nodes may be devoid of textual content, resulting in an empty $t_k$. To address this issue, we employ the special token [UNK] to handle such cases. The [UNK] token performs as a padding element, ensuring every node can be encoded as an equal-sized feature vector $\mathbf{f}_k^t$:

$$\mathbf{f}_k^t = \Phi_{\mathbf{t}}(t_k). \tag{4}$$

Leveraging the modules $\Phi_{\mathbf{g}}$ and $\Phi_{\mathbf{t}}$, we acquire two heterogeneous feature vectors $\mathbf{f}_k^g$ and $\mathbf{f}_k^t$, which offer different insights into multiple modalities. To fuse them into a unified representation, we combine principal component analysis (PCA) projection and dimensional concatenation, which are widely used in works [16, 21, 44], to initialize a 256-dimensional feature vector $\mathbf{f}_k^i$:

$$\mathbf{f}_k^i = [\mathbf{f}_k^g; \Phi_{\mathbf{p}}(\mathbf{f}_k^t)], \tag{5}$$

where $\Phi_{\mathbf{p}}$ is a trainable linear layer for PCA operation. We use $\Phi_{\mathbf{p}}$ to project the 768-dimensional feature vector $\mathbf{f}_k^t$ as the same as the 128-dimensional feature vector $\mathbf{f}_k^g$ and make them have equivalent significance in the fusion phase.

Since the node-level representation $\mathbf{f}_k^i$ is generated by a simple fusion method, we proceed to utilize a transformer network $\Phi_{\mathbf{self}}$ to enrich the initial representation. $\Phi_{\mathbf{self}}$ is stacked by eight vanilla self-attention modules [42] and each module lies in the multi-head attention mechanism. By attending to the relative importance between pairwise nodes, multi-head attention contributes to a better understanding of the context among all nodes and enables each node to focus on the most relevant nodes. Motivated by previous works [25, 27, 40], we apply a rotary positional encoding (RoPE) for the computation of attention score in each self-attention module. RoPE is a variety of relative positional encoding and capable of capturing long-range dependencies and achieving long context extrapolation. The ability to extrapolate long context is crucial for our transformer-based model $\Phi_{\mathbf{self}}$ because the number of nodes is not fixed during either the training or testing phases. Ultimately, a 512-dimensional multimodal representation $\mathbf{f}_k$ is acquired as follows:

$$\mathbf{f}_k = \Phi_{\mathbf{self}}(\mathbf{f}_k^i, \mathbf{RoPE}). \tag{6}$$

To ensure that the multimodal representation $\mathbf{f}_k$ is sufficiently expressive for multiple modalities, we introduce a self-supervised manner to train the transformer network $\Phi_{\mathbf{self}}$. The self-supervised manner relies on three task-specific MLP heads, namely semantic classification head $\Theta_{\mathbf{s}}$, bounding box regression head $\Theta_{\mathbf{g}}$ and textual similarity head $\Theta_{\mathbf{t}}$. These heads all perform as decoding models to transform the representation $\mathbf{f}_k$ back into separate modalities that consist of semantic category, spatial geometry, and textual content. To aid in the learning of $\Phi_{\mathbf{self}}$ and $\Theta_{\mathbf{s}}, \Theta_{\mathbf{g}}, \Theta_{\mathbf{t}}$, we define three training objectives: cross-entropy loss $\mathcal{L}_s$, spatial loss $\mathcal{L}_g$ and contrastive loss $\mathcal{L}_t$. These criteria guide the training process toward the optimal solution and are jointly measured as follows:

$$\mathcal{L} = \lambda_s \mathcal{L}_s + \lambda_g \mathcal{L}_g + \lambda_t \mathcal{L}_t, \tag{7}$$

where $\mathcal{L}$ is the total loss of multimodal representation learning. The cross-entropy loss $\mathcal{L}_s$ is computed using $\Theta_{\mathbf{s}}(\mathbf{f}_k)$ and semantic label $s_k$. The spatial loss $\mathcal{L}_g$ is a linear combination of $\ell_1$ loss and generalized IoU loss [5, 37], both of which are calculated with $\Theta_{\mathbf{g}}(\mathbf{f}_k)$ and spatial parameter $g_k$. The contrastive loss $\mathcal{L}_t$ is calculated with $\Theta_{\mathbf{t}}(\mathbf{f}_k)$, $\Phi_{\mathbf{t}}(t_k^+)$ and $\Phi_{\mathbf{t}}(t_k^-)$ where $t_k^+$ and $t_k^-$ are the positive and negative text sentences augmented by large language models. We randomly select these sentences to form contrastive pairs in training. $\lambda_s$, $\lambda_g$, and $\lambda_t$ are the hyper-parameters that balance the importance of multiple modalities. In the following experiments, we set them to 0.5, 1.0, and 0.1 respectively.

## 3.2 Learning Correspondence Matching

*3.2.1 Objective and Constraints.* The foundation of GameUI construction is to establish correspondences between UI and UX designs. After receiving the multimodal representation of UI design $\mathcal{T}^A$ and UX design $\mathcal{T}^B$, we can easily find the correspondence by evaluating the similarity between pairwise multimodal representation. However, the naive node-to-node similarity fails to satisfy the real-world correspondence rules and may lead to unreasonable matches. This arises from the fact that UI and UX designers often work concurrently with unaligned perspectives. Due to the lack of timely and clear communication, they may craft mismatches with missing or redundant elements on both ends. As depicted in the Fig. 4, we have identified several real-world rules as below:

*Rules 1.* From the viewpoint of tree structures, each node in $\mathcal{T}^A$ has at most one matching node in $\mathcal{T}^B$. This suggests that the leaf nodes in $\mathcal{T}^A$ (e.g text or image) may either have matching leaf nodes in $\mathcal{T}^B$, or be undefined children of some non-leaf nodes in $\mathcal{T}^B$, or be entirely redundant without matching nodes in $\mathcal{T}^B$.
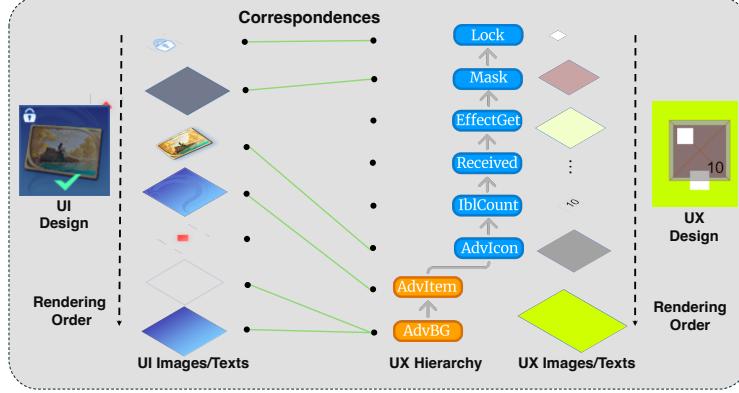
Fig. 4. Illustration of the correspondences between UI and UX designs. These correspondences consist of one-to-empty, one-to-one, and many-to-one mapping. Notably, the blue and yellow in the UX hierarchy respectively indicate the leaf nodes and the no-leaf nodes.

*Rules 2.* For each node in $\mathcal{T}^B$, the number of matching nodes in $\mathcal{T}^A$ is no greater than the count of the entire set. This suggests that both leaf and non-leaf nodes in $\mathcal{T}^B$ may either have multiple matching leaf nodes in $\mathcal{T}^A$, or be completely redundant without matching nodes in $\mathcal{T}^A$.

Following these rules, we can determine a constrained objective function as:

$$\underset{\mathbf{P} \in U(\mathbf{r}, \mathbf{c})}{\operatorname{argmin}} \sum \mathbf{P} \odot \mathbf{S} + \gamma \Omega(\mathbf{P}, \mathbf{S}), \tag{8}$$

with

$$U(\mathbf{r}, \mathbf{c}) = \{\mathbf{P} \in \mathbb{R}^{\mathbf{m} \times \mathbf{n}} | \mathbf{P} \cdot \mathbf{1_n} = \mathbf{r}, \mathbf{P}^{\mathbf{T}} \cdot \mathbf{1_m} = \mathbf{c}\}. \tag{9}$$

Eq. (8) and Eq. (9) are combined to formulate a classic optimal transport problem ($\mathcal{NP}$-hard). $\mathbf{m}$ and $\mathbf{n}$ is the number of available nodes in $\mathcal{T}^A$ and $\mathcal{T}^B$. $\mathbf{P}$ is the assignment matrix with integer values of either 0 or 1. When the $i^{th}$ node of $\mathcal{T}^A$ is matched with the $j^{th}$ node of $\mathcal{T}^B$, the value of $\mathbf{P}_{ij}$ is equal to 1, and 0 otherwise. $\mathbf{S} \in \mathbb{R}^{\mathbf{m} \times \mathbf{n}}$ is the cost matrix, which is critical to determine the optimal transport plan. As the similarity between two nodes across $\mathcal{T}^A$ and $\mathcal{T}^B$ increases, the values in matrix $\mathbf{S}$ gradually decrease from 1 to 0. $U(\mathbf{r}, \mathbf{c})$ is the transportation polytopes [4, 9] which are composed of two equalities. The two equalities separately constrain the feasible matches in matrix $\mathbf{P}$ along rows and columns. $\mathbf{r}$ and $\mathbf{c}$ are two one-dimensional integer vectors with lengths equal to $\mathbf{m}$ and $\mathbf{n}$, and the range of each value of $\mathbf{r}$ and $\mathbf{c}$ is constrained as $\forall i, 0 \le \mathbf{r}_i \le \mathbf{1}; \forall j, 0 \le \mathbf{c}_j \le \mathbf{m}$. $\sum \mathbf{P} \odot \mathbf{S}$ is the summation of the Hadamard product of $\mathbf{P}$ and $\mathbf{S}$. $\Omega$ is the regularization function used to prune the ambiguous matches, which are illustrated in Fig. 6. To enhance the flexibility of the optimization process, a hyper-parameter $\gamma = 1.0$ is introduced as the factor of the regularization function.

We introduce an efficient method to minimize the objective function Eq. (8) by combining a transformer-based model and an integer programming solver. The transformer-based model incorporates novel grouped cross-attention modules to determine the unknown matrices $\mathbf{S}, \mathbf{r}, \mathbf{c}$. Integer programming is a powerful technique to solve the optimal transport problem, aiming to find the optimal matrix $\mathbf{P}$ by exploring a vast solution space.

*3.2.2 Grouped Cross-Attention Module.* In contrast to prior studies [46, 48], we do not directly compute the node-to-node similarity matrix to solve the optimal transport problem. Given the trees $\mathcal{T}^A$ and $\mathcal{T}^B$ with hundreds or even

thousands of nodes, the number of penalty terms defined in regularization function $\Omega$ will exhibit quadratic growth as the nodes increase, potentially resulting in a time-consuming or unsolvable optimization process.

To overcome the limitation, we review the hierarchy of $\mathcal{T}^B$ and proceed to partition all nodes in $\mathcal{T}^B$ into several groups in terms of their association with secondary-level nodes in $\mathcal{T}^B$. As a child of the root node in $\mathcal{T}^B$, each secondary-level node corresponds to an individual group. If a node is the descendant of a secondary-level node, this node should be categorized into the corresponding group. According to the grouping criterion, we can resolve the optimal transport problem with node-to-group optimization instead of node-to-node optimization. The node-to-group optimization also enables us to recursively identify all matches between $\mathcal{T}^A$ and $\mathcal{T}^B$, which significantly reduces the complexity of the optimization.
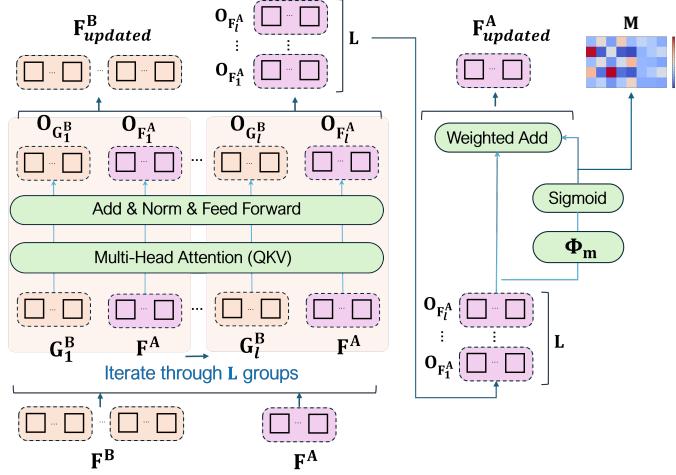


Fig. 5. Illustration of grouped cross-attention module. We iterate through the groups of UX representations and repeatedly compute the cross attention between each group and the entirety of UI representations.

For each node-to-group optimization depicted in the right column of Fig. 3, $\mathcal{T}^A$ and $\mathcal{T}^B$ are respectively fed into the frozen models, generating the 512-dimensional multimodal representation for all their nodes as $\mathbf{F}^A \in \mathbb{R}^{\mathbf{m} \times \mathbf{512}}$ and $\mathbf{F}^B \in \mathbb{R}^{\mathbf{n} \times \mathbf{512}}$. Assume that $\mathcal{T}^B$ have $\mathbf{L}$ secondary-level nodes, $\mathbf{F}^B$ can be dynamically partitioned into $\mathbf{L}$ groups, one of which is denoted as $\mathbf{G}_l^B (1 \leq l \leq \mathbf{L})$. We adopt a transformer-based model $\mathbf{\Phi_{cross}}$ which is stacked by four novel grouped cross-attention (GCA) modules, to learn the cost matrices $\mathbf{S}$ and determine other unknown matrices $\mathbf{r}, \mathbf{c}$. As illustrated in Fig. 5, every individual group $\mathbf{G}_l^B$ and the entirety of $\mathbf{F}^A$ are iteratively processed by the GCA module. Accordingly, we can formulate the procedure for each grouped cross-attention as:

$$\mathbf{O}_{\mathbf{G}_l^B} = \mathbf{Attention}_{\mathbf{G}_l^B \to \mathbf{F}^A}(\mathbf{F}^A, \mathbf{G}_l^B), \tag{10}$$

$$\mathbf{O}_{\mathbf{F}_l^A} = \mathbf{Attention}_{\mathbf{F}^A \to \mathbf{G}_l^B}(\mathbf{F}^A, \mathbf{G}_l^B), \tag{11}$$

where $\mathbf{O}_{\mathbf{G}_l^B}$ and $\mathbf{O}_{\mathbf{F}_l^A}$ are the outputs in $l^{th}$ iteration. The calculation of multi-head attention is represented as **Attention** and its subscript arrow indicates the target where attention scores are applied. Upon completing the iterative procedure, we generate groups of outputs that need to be integrated to update the original multimodal representations for

subsequent GCA modules. Considering that $O_{G_l^B}$ is derived from a part of $F^B$ while $O_{F_l^A}$ is derived from the entirety of $F^A$, we can represent the groups of outputs as $O_{G^B} \in \mathbb{R}^{n \times 512}$ and $O_{F^A} \in \mathbb{R}^{L \times m \times 512}$.

By using the index of each group, $O_{G^B}$ can be directly added into the original $F^B$ for updating as follows:

$$F_{updated}^B = F^B + O_{G^B}. \tag{12}$$

In contrast, $O_{F^A}$ requires additional dimensional reduction to update the original $F^A$. We utilize a two-layers MLP to compute a lower-dimensional matrix $M \in \mathbb{R}^{L \times m}$ as:

$$M = \mathrm{Sigmoid}(\Phi_m(O_{F^A})), \tag{13}$$

where $M$ can implicitly represent the matching probability between $F^A$ and each $G_l^B$. This is attributed to cross attention that the more relevant $G_l^B$ carries more weights in $O_{F_l^A}$, leading to distinct significance in $O_{F^A}$ along the first dimension $L$, so that we can calculate the cost matrix $S$ in Eq. (8) by:

$$S = 1 - M^\top. \tag{14}$$

The updating of $F^A$ can be computed by using weighted add among the groups of output as follows:

$$\hat{M}_l = \frac{M_l}{\sum_{l=1}^{L} M_l}, \tag{15}$$

$$F_{updated}^A = F^A + \sum_{l=1}^{L} \hat{M}_l^\top \odot O_{F_l^A}, \tag{16}$$

where $\hat{M} \in \mathbb{R}^{L \times m}$ is proportional normalization of $M$ along the same first dimension $L$. As the matching probability increases, the corresponding weight becomes larger.

Finally, we exploit a binary cross-entropy loss $\mathcal{L}_m$ to train the transformer-based model $\Phi_{cross}$. The binary cross-entropy loss $\mathcal{L}_m$ measures the difference between the matching probability $M$ and the given group labels within a binary log-likelihood framework, which naturally supports unmatchable cases in both training and testing. For clarity, we represent unmatchable labels as a zero-filled vector during training and filter out unmatchable cases with a preset threshold $\sigma = 0.5$ during testing. Threshold $\sigma$ also allows us to ascertain the values of matrix $r$ and $c$ in Eq. (9), because the number of feasible matches in each row and column of $M$ can be approximated by comparing the matching probability with $\sigma$.

*3.2.3 Integer Programming.* Although the variables $P$, $S$, $r$, $c$ have been defined in Section 3.2.2, we notice that the regularization function $\Omega$ is not determined yet. The function $\Omega$ aims to eliminate suboptimal solutions, especially ambiguous matches depicted in Fig. 6, which often arise from the varying hierarchical relationship or converse rendering order between pairwise matches. Assume that two nodes $v_i^A$ and $v_j^A$ in the UI design have an ancestor-descendant relationship (i.e. $\Delta h_{ij}^A = 0$), but after matching, the corresponding nodes $v_{i'}^B$ and $v_{j'}^B$ in the UX design have a descendant-ancestor hierarchical relationship (i.e. $\Delta h_{i'j'}^B = 1$). In this case, the pairwise matches are recognized as ambiguous. Considering that the ambiguity in rendering order can be explained similarly, we determine the regularization function with penalty terms as follows:

$$\Omega(P, S) = \sum_{P_{ii'}=1, P_{jj'}=1} \frac{S_{ii'} + S_{jj'}}{\tau}, \tag{17}$$

Fig. 6. Two ambiguous matches for regularization. (a) The ambiguity in varying hierarchical relationships. After matching, the "ancestor-descendant" becomes the "descendant-ancestor". (b) The ambiguity in converse rendering order. After matching, the "rendered-after" becomes the "rendered-before". Dark gray indicates the higher priority.

where $\mathbf{P}_{ii'} = 1$ and $\mathbf{P}_{jj'} = 1$ together indicate the ambiguous matches. The penalty term, composed of the corresponding costs $\mathbf{S}_{ii'}$, $\mathbf{S}_{jj'}$, is adaptively integrated into regularization through a hyper-parameter $\tau$, which is set to group size $\mathbf{L}$ in our experiments. Because our method explicitly incorporates the edge parameters $\Delta h_{ij}$ and $\Delta r_{ij}$ into the regularization function, it yields greater improvements in maintaining the coherence of hierarchy and rendering. We leverage an open-source solver provided by Google's OR-Tools [34] to implement the objective function Eq. (8) and the transportation polytopes Eq. (9). To obtain the optimal assignment matrix $\mathbf{P}$ in real-time, we also apply some acceleration tricks such as adding heuristic hints.

## 3.3 Interactive Construction

*3.3.1 Universal Data Protocol.* Based on the estimated correspondences, we begin to construct GameUI by integrating the visual attributes of the elements in the UI design into the logic structure of the UX design. This process involves resizing the height and width of each node, updating text content, font, and color, and replacing image textures, etc. As depicted in Fig. 7, we introduce a universal protocol for data description and attribute integration.

The protocol is composed of two main components: language and compiler. The language defines a set of UI entities ranging from low-level to high-level semantics. Each entity is also characterized by a sequence of attributes that precisely describe common UI attributes such as position, rotation, scale, canvas anchor, opacity, image texture, text fonts, etc. The compiler is to create a parser that interprets the language into another programming language and enables execution across various platforms. We employ Antlr4[3], a powerful parser generator to support multiple programming languages, including C++, C#, Python, and TypeScript. Benefiting from Antlr4's capabilities, we can seamlessly adapt the protocol for mainstream game engines such as Unity and Unreal Engine (UE), as well as some digital content creation (DCC) software like Photoshop and Figma. To streamline the usage, we have developed accompanying plug-ins with embedded compilers to enable cross-platform applications.

*3.3.2 Interactive Web-based Tool.* We implemented a web-based interactive tool that follows a client-server architecture to remotely execute our proposed algorithm and display the matching results on the user's desktop PC. The input of the tool comprises two description files (ends with UIPROTO and UXPROTO) as well as the RGBA images that composite the UI design. The output is the updated UXPROTO, with the visual attributes integrated from the UI design.

---
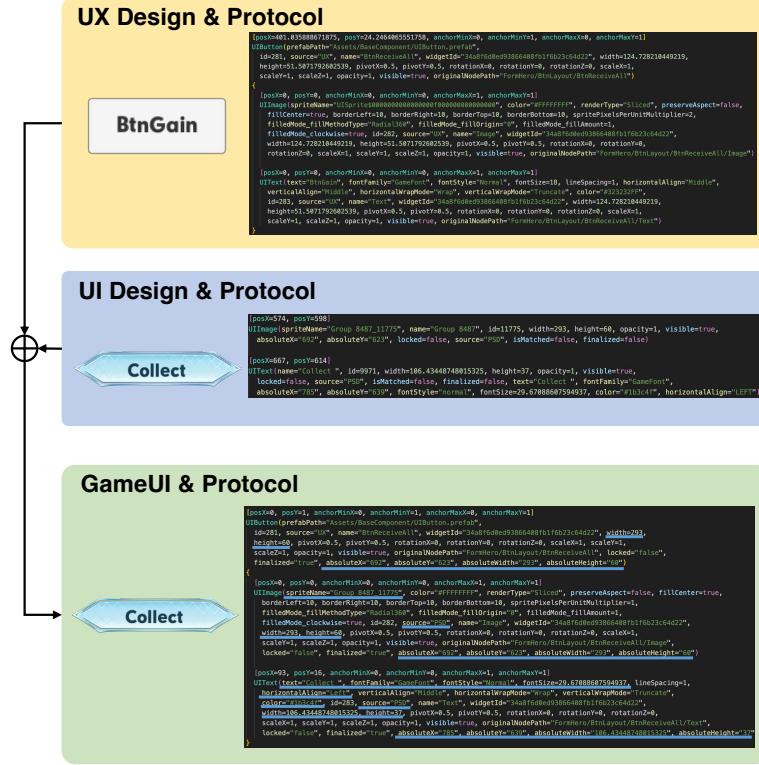
[3]https://github.com/antlr/antlr4

Fig. 7. Integration of UI and UX protocols. The blue lines in the GameUI protocol indicate the visual attributes transferred from the UI protocol.

As illustrated in Fig. 8, the interactive web-based tool provides users with several capabilities, including automatic and manual construction, previewing, editing, and annotating.

*Automatic Construction.* This is the core function of our system, which automatically invokes the backend algorithm to estimate the correspondence between the imported UI and UX designs, and integrates the visual attributes of the UI design into the original UX design based on the correspondence.

*Manual Construction.* This enables users to manually modify the unexpected matching results, because the automatic results may not always be accurate. Specifically, we use the matching probability $\mathbf{M}$ as confidence values to guide users by highlighting the potential matches. This function refines unsettled matching results and contributes to future operations.

*Data Preview.* This gives users various perspectives to understand and validate UI and UX designs, including options like multiple rendering modes. These perspectives also aid users in subsequent interactions.

*Data Editing.* This allows users to fine-tune the UX design to reduce disparities between the UI and UX designs, which may involve deleting or creating UX nodes, repositioning, and updating the semantic type of UX nodes.
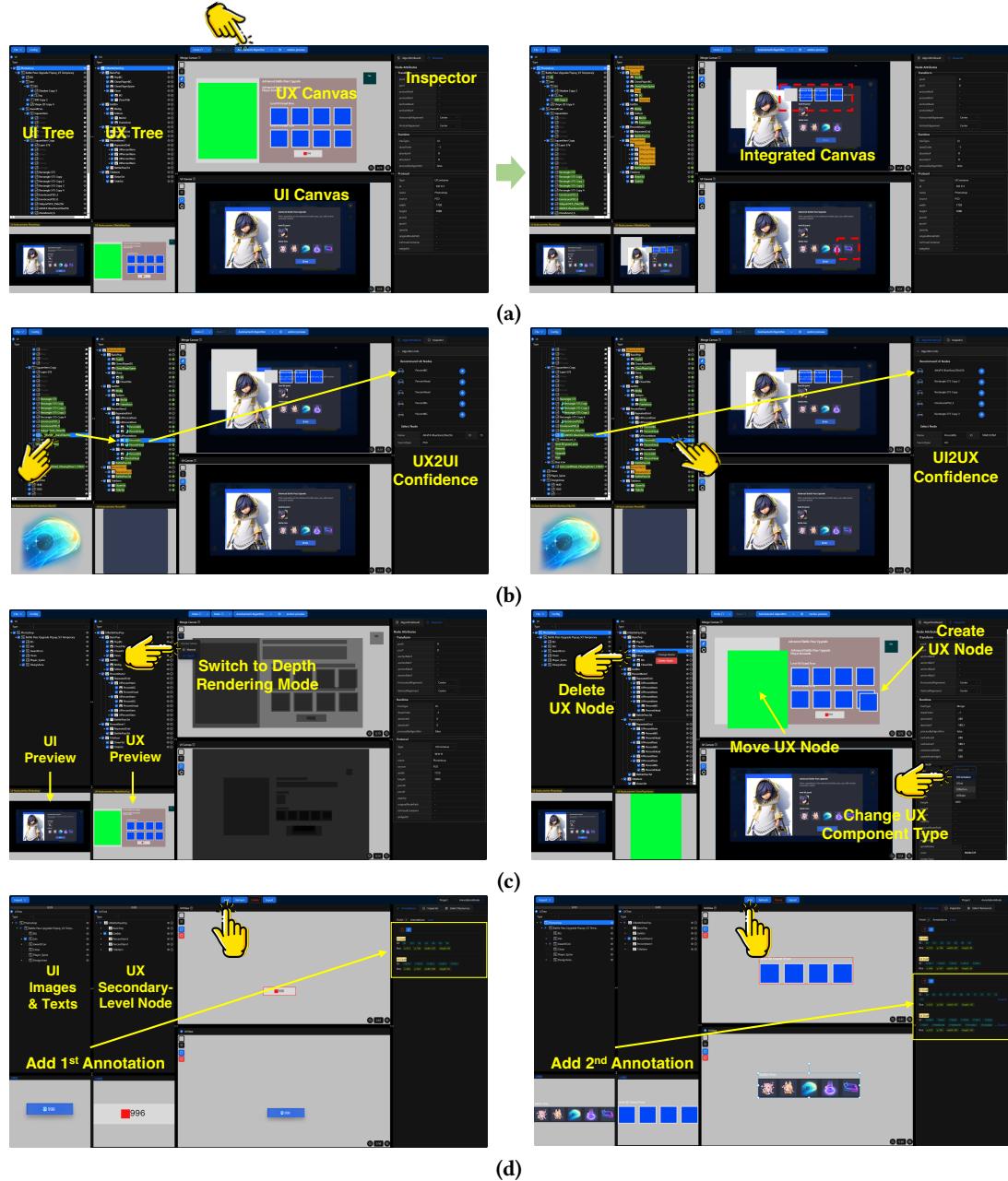
Fig. 8. Illustration of the interactive web-based tool. (a) Automatic GameUI construction from UI and UX designs. The red dashed box indicates the mismatches between UI and UX designs. (b) Manual construction with bidirectional confidence values. (c) Data preview and editing. (d) Data annotation for recording the correspondences between the UI leaf nodes (images and texts) and the UX secondary-level nodes.

*Data Annotation.* This is to create a specialized dataset from actual game projects. The internal users can also utilize this function to submit their data to the cloud server, augmenting the dataset.

With these interactive features, the tool offers a complete solution for GameUI construction, adequately addressing the limitations of automatic algorithms. By combining human supervision with machine feedback, the tool ensures a balance between efficiency and accuracy, consequently improving the experience of the overall process.

## 4  EXPERIMENTS

In this section, we primarily evaluate the performance of our proposed approach on the GameUI construction task. Besides, we use the UI retrieval task as a supplementary evaluation to strengthen our assessment. To conduct these experiments, we build a specialized dataset and combine it with the public RICO dataset, to compare our method against the state-of-the-art methods. Finally, We demonstrate the effectiveness of our approach in the GameUI construction through both quantitative and qualitative analyses.

### 4.1  Datasets

To the best of our knowledge, there is no existing work in academia and industry that truly formulates and tackles the GameUI construction task. Most relevant research, such as those on UI generation or UI retrieval, primarily focuses on the RICO dataset [10] or ENRICO dataset [20]. These public datasets are largely crawled from mobile Apps in categories such as communication, medicine, and finance, which are quite simple in aesthetics without complex visual elements. Compared to real-world game UIs, these datasets exhibit differences in data attributes and distribution. Therefore, we additionally collected pairs of UI and UX designs from actual game projects to create a specialized dataset, which is called GAMEUI, aiming to better evaluate the effectiveness of our approach.

The GAMEUI dataset comprises 42 sets of game UIs provided by different game projects. Each set includes a pair of UI and UX designs. Some of these sets have already been launched in operation, while others remain in the development phase. We use the interactive web-based tool mentioned earlier for data annotation and cleansing. As shown in Fig. 8(d), we can switch to the annotation mode and manually establish the correspondences between UI and UX designs, starting from the root node of the UX design and iteratively recording the correspondence between the image or text resources in the left UI tree and the secondary-level node in the right UX tree. All annotations within the same UX secondary level are exported as a single data sample. After removing low-quality samples from the 42 interfaces, we gather 381 samples in total. These samples are split into 282/33/66 for model training, validation, and testing. Especially, the samples in different splits cannot belong to the same interface. To increase the validity and credibility of our results, we also conduct experiments on the RICO dataset, which is divided into 35851/2109/4218 samples for training, validation, and testing.

### 4.2  Implementation Details

In the experiments with the GAMEUI dataset, we train our models in two stages. The first stage concentrates on learning the multimodal representations of UI and UX designs, while the second stage aims to learn their correspondences. In the first stage, we freeze the weights of the pre-trained language model $\Phi_t$ due to the lack of sufficient text data. Meanwhile, we adopt several data augmentation schemes to avoid overfitting. These data augmentation schemes include using a large language model to generate positive and negative samples for each text sentence in advance, randomly masking a few words in the text sentences, and randomly translating or scaling the 2D bounding box within a limited range. We set the batch size to 4 and train the first-stage model for 300 epochs. The data in each batch is padded and processed to

a fixed length of 200, trimming any excess. The initial learning rate is set to 1e-5, with a factor of 0.2 multiplied every 50 epochs. In the second stage, we freeze the weights of the first-stage model and train the second-stage model for 200 epochs using identical learning configurations.

| Method | Train Set | Similarity | Precision↑ | Recall↑ | F1-score↑ | RMSE↓ | PER↓ |
|---|---|---|---|---|---|---|---|
| **LayoutGMN** [33] | RICO | N2N | 91.4 | 60.4 | 62.5 | 1.66 | 1.52 |
| | GameUI | N2N | 92.7 | 61.8 | 65.2 | 1.56 | 0.81 |
| **LayoutTrans** [13] | RICO | N2N | 76.0 | 54.3 | 42.1 | 3.84 | 1.75 |
| | GameUI | N2N | 78.6 | 54.8 | 45.1 | 5.48 | 2.48 |
| **LayoutDM** [15] | RICO | N2N | 80.4 | 55.2 | 47.2 | 3.24 | 4.35 |
| | GameUI | N2N | 81.6 | 55.5 | 48.6 | 4.00 | 4.46 |
| **Ours** (w/o Reg) | RICO | N2N | 93.4 | 63.7 | 66.8 | 1.63 | 0.94 |
| | GameUI | N2N | 94.2 | 64.1 | 68.9 | 1.65 | 1.05 |
| | GameUI | N2G | 96.0 | 68.0 | 74.4 | 1.16 | 0.58 |
| **Ours** (w/i Reg) | GameUI | N2N | - | - | - | - | - |
| | GameUI | N2G | **96.2** | **68.7** | **75.2** | **0.78** | **0.36** |

Table 1. Evaluation of construction results on the GAMEUI dataset's test set. We separately report the accuracy of the correspondence matching and the visual consistency of the constructed image, with different training datasets and similarity computation. N2N refers to the node-to-node similarity obtained from the multimodal representations, while N2G specifically refers to the node-to-group similarity obtained from the GCA modules. Reg denotes the hierarchical and rendering constraints. Note that if we apply the hierarchical and rendering constraints on the node-to-node similarity, we cannot get the optimal correspondences within a limited time budget.

In the experiments with the RICO dataset, we cannot follow the aforementioned two-stage training strategy due to the absence of correspondence annotations. To ensure a fair comparison with other baselines, we solely train our first-stage model to learn the multimodal representation of each element in the RICO dataset. During the training procedure, the batch size is set to 64, the initial learning rate is set to 1e-4, and the number of training epochs is 50.

### 4.3 Evaluation

*4.3.1 Setup.* In the following studies, we select LayoutGMN [33], LayoutTrans [13], LayoutDM [15], and LayoutBlending [48] as baselines in the GameUI construction. We measure these methods on the GAMEUI dataset's test set and estimate the correspondences between pairwise UI and UX designs. Using these correspondences, we can integrate the raw RGBA images for image re-rendering and compute the differences between the rendered image and the original UI screenshot. It's worth noting that the raw RGBA images are exported from the UI design, and the original UI screenshot serves as the ground truth image.

Among the baselines, LayoutBlending is the unique heuristic approach that does not require training. It exploits spatial layout, semantics, and hierarchical structure to compute node-to-node similarity and adopts the Hungarian algorithm [17] to determine correspondences. LayoutGMN, LayoutTrans, and LayoutDM are based on deep neural networks and exclusively combine spatial layout and semantics for multimodal representation learning. The DNNs-based methods are individually trained in contrastive, autoregressive, and generative ways. We adopt the same learning configurations as our first-stage model to train these baselines on the RICO and GAMEUI datasets. Then, we use the trained models to extract the multimodal representation of each element in UI and UX designs, and compute the node-to-node cosine similarity as the matching probability matrix. Particularly, we extend the maximum length of the absolute positional embeddings in LayoutDM and LayoutTrans from 25 to 200, to better accommodate the GAMEUI

dataset. Besides, since LayoutDM is trained through a conditional generation process, we set the condition type of LayoutDM as the ground truth in testing, which means that no additional constraints are imposed on the model, generating an unbiased multimodal representation.

After obtaining the node-to-node similarity matrix, we can identify the UX node with the highest similarity (in the columns) as the matching object for each UI leaf node (in the rows). Since the ground truth samples in the GAMEUI dataset only record the correspondences from the UI leaf nodes to the UX secondary-level node, we must verify whether each matching UX node belongs to the correct UX secondary-level node in terms of the UX hierarchy. In other words, any node-to-node similarity matrix should be converted into a node-to-group matrix before it can be evaluated on the GAMEUI dataset.

We prioritize two aspects for the quantitative benchmarks in the GameUI construction: the accuracy of the correspondence matching, and the visual consistency of the constructed image. To quantify the matching accuracy, we compute precision, recall, and the F1-score. These metrics are defined using positive and negative samples, which enable us to identify various types of matching scenarios, including matchable, unmatchable, correctly matched, and incorrectly matched. To mitigate the effects of imbalanced categories, we use a macro-average setting to balance the scores. To measure the visual consistency of the constructed image, we adopt the Root Mean Squared Error (RMSE) and Pixel-wise Error Ratio (PER) as the metrics. The RMSE and PER can be determined as:

$$\text{RMSE} = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(\mathbf{I}_i^T - \mathbf{I}_i^C)^2}, \tag{18}$$

$$\text{PER} = \frac{1}{n}\sum_{i=1}^{n}\mathbb{I}(\mathbf{I}_i^T \neq \mathbf{I}_i^C), \tag{19}$$

where $\mathbf{I}_i^T$ and $\mathbf{I}_i^C$ respectively denote the $i^{th}$ pixel in the ground truth image and the constructed image. $n$ is the total number of pixels in each image. $\mathbb{I}$ is the indicator function that takes the value of 1 when the inside inequality is satisfied, and 0 otherwise. The unit of RMSE is the pixel intensity ranging from 0 to 255, while the unit of PER is the percentage. For the GAMEUI dataset's test set, both PER and RMSE values are computed as weighted averages, with each item weighted by the number of pixels in every image.

We additionally conduct experiments on the larger RICO dataset for the UI retrieval task. Within the RICO dataset's test set, LayoutTrans and Screen2Vec [24] are chosen to compare with our first-stage model. This comparison reveals slight differences in querying the nearest neighbor retrievals. Screen2Vec extracts the latent representation of the entire screen layout, thereby it can directly obtain the scalar similarity between two UI pages. However, our first-stage model and layoutTrans yield the node-to-node similarity matrix, which requires further processing. This involves separately calculating the maximum values in the rows and columns and averaging them to obtain the scalar similarity between two UI pages. By analyzing the retrieved results, we can assess the quality of the multimodal representations learned by the models and determine how effectively they comprehend the spatial layout and textual content of diverse UI pages.

We use the official github repositories [4] [5] [6] [7] to implement the above evaluation. Besides, we prioritize using the official pre-trained models for comparisons. Note that our experimental devices are composed of a desktop PC and a remote machine.

---

[4]https://github.com/agp-ka32/LayoutGMN-pytorch
[5]https://github.com/lyf7115/LayoutBlending
[6]https://github.com/CyberAgentAILab/layout-dm
[7]https://github.com/tobyli/Screen2Vec

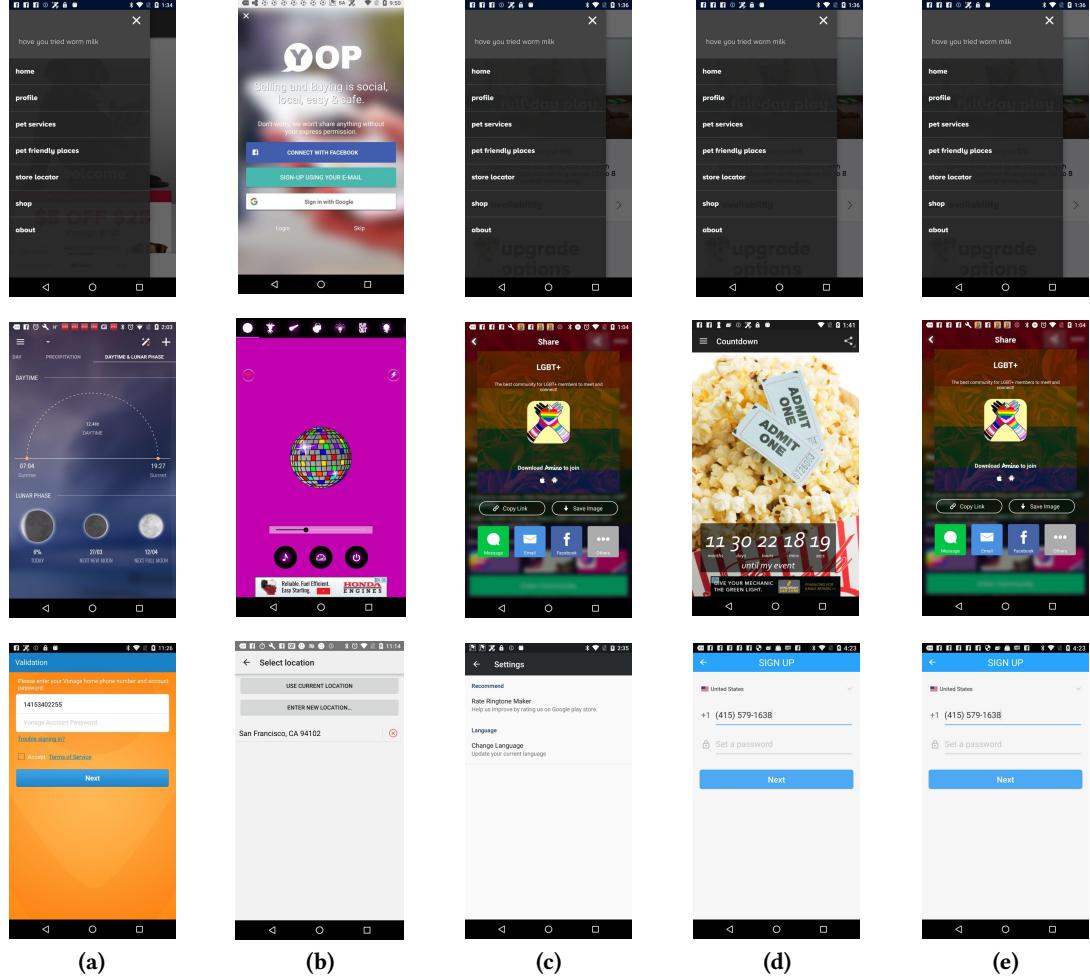|                 |                 |                 |                 |                 |
| :-------------: | :-------------: | :-------------: | :-------------: | :-------------: |
| (a)             | (b)             | (c)             | (d)             | (e)             |

Fig. 9. Retrieval results of our approach, Screen2Vec [24], and LayoutTrans[13] on the RICO dataset's test set. (a) Query. (b) Results of Screen2Vec trained on the RICO dataset. (c) Results of LayoutTrans trained on the RICO dataset. (d) Results of our first-stage model trained on the RICO dataset. (e) Results of our first-stage model trained on the GAMEUI dataset.

*4.3.2 Results.* We qualitatively and quantitatively compare the results of our first-stage model against LayoutGMN [33], LayoutTrans [13], LayoutDM [15] and Screen2Vec [24]. Table 1 illustrates that our first-stage model performs better than other baselines in the GameUI construction, even if only using a node-to-node similarity matrix based on the multimodal representations for matching. Given the comparison of the identical models trained on the RICO and GAMEUI datasets, we can observe that this advanced performance is not limited to a specific dataset. Moreover, the supplementary validations in Fig. 9 demonstrate that our first-stage model is still competitive in the UI retrieval task, covering not only spatial layout but also textual content. Especially, the middle retrieval result in Fig. 9(d) is the only one with high similarity to the query content, both of which are related to time and date. Fig. 9(e) shows that even when

trained on the smaller GAMEUI dataset, our first-stage model can achieve the second-best retrievals. This reconfirms the generalization performance of our first-stage model on different datasets.

Next, we perform ablation studies for the GCA modules and regularization terms in Section 3.2.2 and Section 3.2.3. Compared to the naive node-to-node results in Table 1, the GCA modules and regularization terms successively exhibit an improvement in both correspondence matching and visual consistency. As depicted in Fig. 10(a) and (b) with more details, we can see that the transparent background in the UI design shrinks the canvas and causes a shift to the right compared to the UX design. The three dragon patterns on the left side of the UI design are misaligned with their corresponding secondary-level node (marked in blue) in the UX design. The similarity between misaligned nodes can easily lead to unexpected correspondences in Fig. 10(c). In contrast, Fig. 10(d) shows that the GCA modules can achieve higher overall similarities between the secondary-level node in the UX design and the misaligned patterns in the UI design. As the regularization terms come into play in Fig. 10(e), we penalize ambiguous matches with hierarchical and rendering constraints, thoroughly resolving the trivial mismatches. Furthermore, by comparing the last two rows in Table 1, we can see that the GCA modules reduce the dimensions of the similarity matrix, effectively decreasing the computational load brought by the regularization terms. Conversely, if we apply the hierarchical and rendering constraints on the node-to-node similarity matrix, it would be impossible to get the optimal solution within a limited time budget (set to 30 seconds in our experiments). Table 2 presents the time statistics of our approach under different configurations. Although the inclusion of regularization terms increases the run time, the resolution time for a single sample remains below one second, ensuring that our method introduces minimal delays.

| Method | Similarity | Per-sample | Per-interface |
|---|---|---|---|
| **Ours** (w/o reg) | N2N | 0.190 | 0.786 |
| | N2G | 0.206 | 0.854 |
| **Ours** (w/i reg) | N2G | 0.472 | 1.957 |

Table 2. Average time consumption (in seconds) of our proposals on the GAMEUI dataset's test set. The Per-sample indicates the time consumed to match a single annotated sample, while the Per-interface indicates the time required for all samples within a pair of UI and UX designs. The time statistics are both given on the local desktop PC.

Ultimately, we employ four sets of real-world game UIs to illustrate the GameUI construction. In this process, all hierarchical nodes in the UX designs are recursively matched and constructed with the UI designs. As shown in Fig. 11 and Fig. 12, we display the overall matching results of each set from top to bottom, including the correspondences, the constructed image, and the error maps rendered in heat and binary mode. Compared to other results, our method achieves leading performance in matching accuracy and visual consistency. We can also analyze the shortcomings of the baseline methods. For example, LayoutTrans and LayoutDM perform poorly in the cases "Adventure Experience" and "Autumn Event" with more elements, but perform well in the cases "Settle Up" and "Super Reward" with fewer elements. This demonstrates the limitation of absolute positional encoding, namely the poor generalization when training data is insufficient. The mismatches caused by LayoutBlending are related to the use of the Hungarian algorithm [17] for bipartite graph matching, which does not satisfy the matching rules between UI and UX designs. It is worth noting that visual consistency doesn't always align with matching accuracy. It can be observed in Fig. 11 and Fig. 12 that LayoutGMN generates fewer visual errors. One reason is that mismatches do not disrupt visual consistency when the rendering order of overlapping images is not reversed, as shown in Fig. 10(c). Another reason is that the rendering constraints penalize the ambiguous instances, instead of the pixels. This contradicts the intent of error maps, which aim to evaluate the visual consistency with pixel-wise difference.

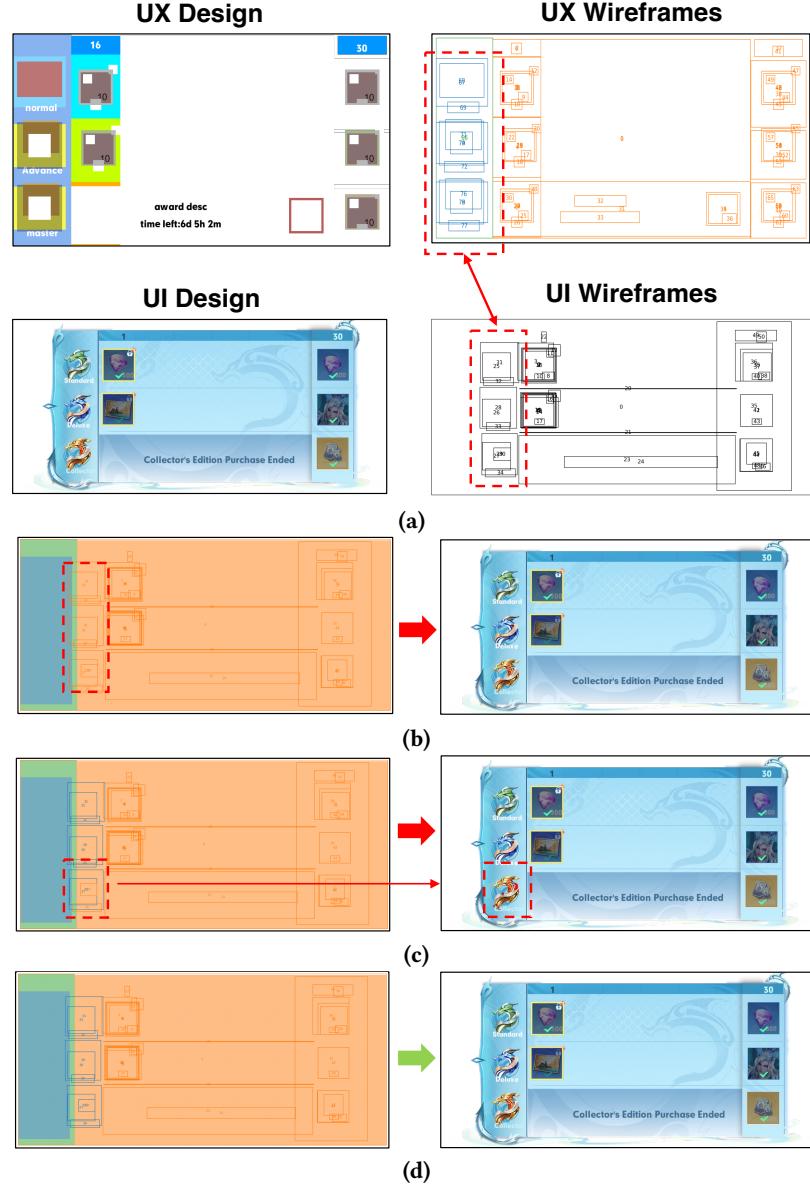**(a)**



**(b)**



**(c)**



**(d)**

Fig. 10. Matching results of our approach on a pair of UI and UX designs "CenterPart". (a) Designs and wireframes. The red dashed box and bidirectional arrow indicate the misalignment caused by the transparent background in UI wireframes. The three colors (yellow, blue, and green) in UX wireframes represent that these wireframes belong to three UX secondary-level nodes. Scalars within wireframes indicate rendering order, with larger numbers denoting higher priority. (b) Matching results estimated from node-to-node similarity matrix. The left side shows that UI wireframes and corresponding UX secondary-level nodes are filled with the same colors to display their correspondences. The red dashed box emphasizes incorrect correspondences. The right side shows the GameUI constructed through these correspondences. (c) Matching results estimated from node-to-group similarity matrix. The red dashed box and red mask on the right side emphasize the pixel-wise difference between the constructed GameUI and the original UI design. (d) Matching results based on (c) and further optimized with hierarchical and rendering constraints.
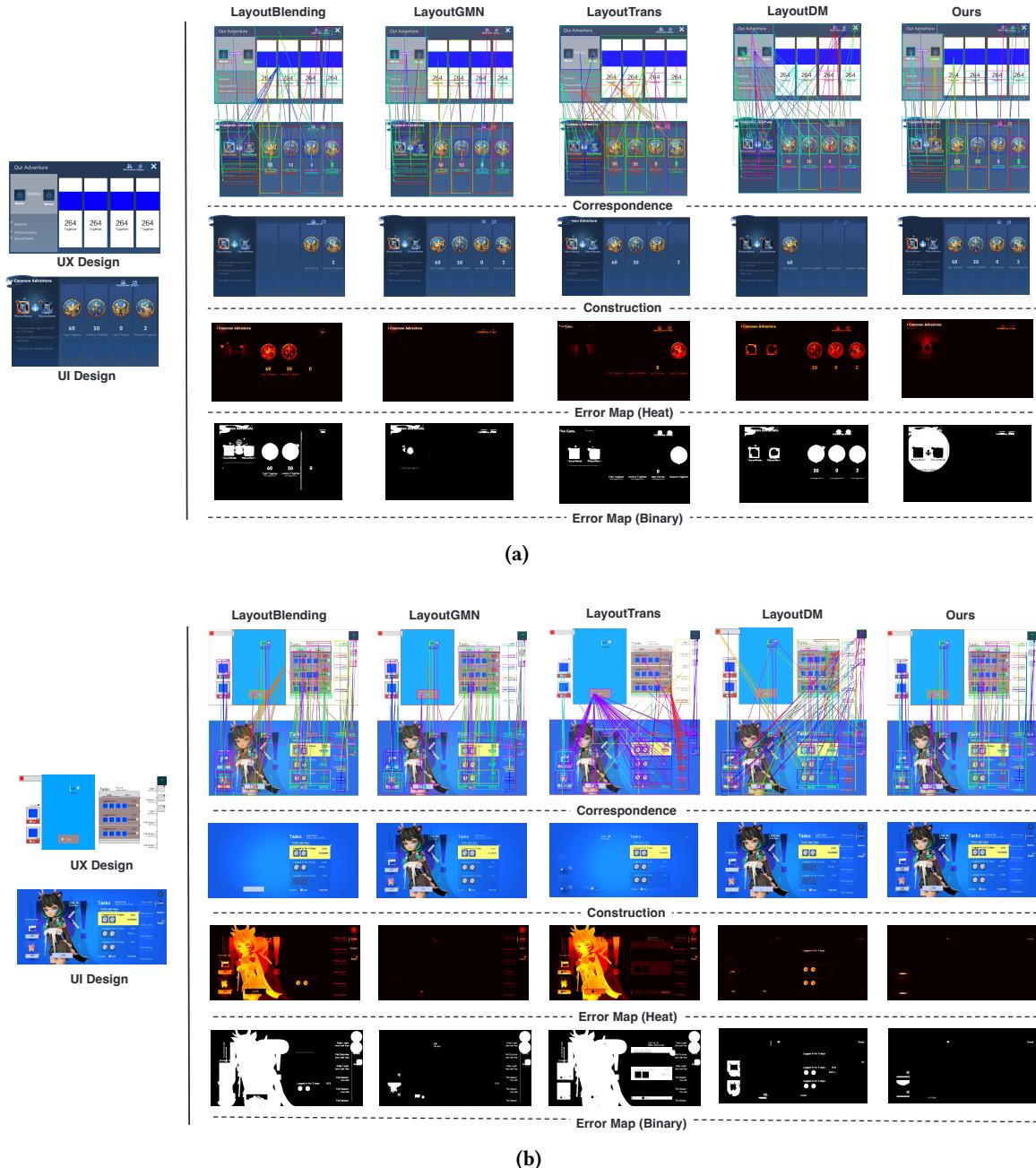
(a)



(b)

Fig. 11.  Construction results of two real-world cases. (a) Case "Adventure Experience". (b) Case "Autumn Event". The two cases come from the GAMEUI dataset's test set.

Zhongliang Tang, Mengchen Tan, Fei Xia, Qingrong Cheng, Hao Jiang, and Yongxiang Zhang
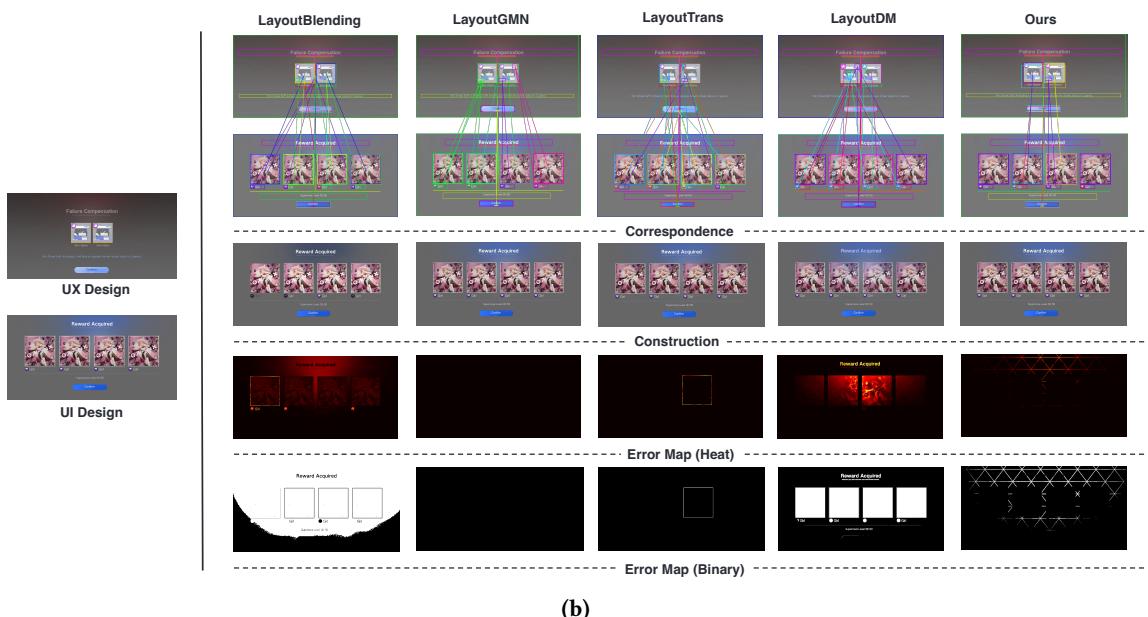


**(a)**



**(b)**

Fig. 12. Construction results of two real-world cases. (a) Case "Settle Up". (b) Case "Super Reward". The two cases come from the internal user data recorded by our system.

### 4.4 User Study

*4.4.1 Participants.* We invited three game UI developers from different game project teams to test our system. Each participant had 2-5 years of experience in game UI development, with an average age of under 30.

*4.4.2 Environments.* We provided participants with the same local desktop PC and identical remote server configurations. Before testing, we completed the installation of all necessary plug-ins used in the procedure.

*4.4.3 Materials.* We prepared five sets of game UIs for testing. The first four sets were gathered from previously completed game UIs, each with a pair of UI and UX designs. The fifth set was brand-new and contained 12 pairs of UI and UX designs. Notably, 10 of these pairs were designed for simple pop-up windows that required little time to construct. Within the five sets, the UI designs were Photoshop files and the UX designs were Unity prefabs.

*4.4.4 Criteria.* Our evaluation criterion was based on man-days, which refers to the working hours required for a single person to construct a single GameUI. A full man-day is equivalent to eight hours.

*4.4.5 Procedure.* Two developers were assigned to build the first four sets, and the third engineer was responsible for handling the brand-new set. They were asked to complete the GameUI construction under two conditions: the first one being the traditional approach, and the second one using our innovative system. The traditional approach refers to manually exporting the protocol and resources from Photoshop and building the cohesive game interface in Unity according to the prefab's hierarchy. In contrast, the innovative approach involves using plug-ins to automatically export the protocols and resources from Photoshop and Unity, which are then assembled in the interactive tool. Once construction is finished, the integrated protocol and resources are re-imported into the Unity engine for visual and functional testing. The three developers used our system before proceeding with the traditional manual method, which ensures that the time statistics do not favor our system, even if they are already familiar with the sets.
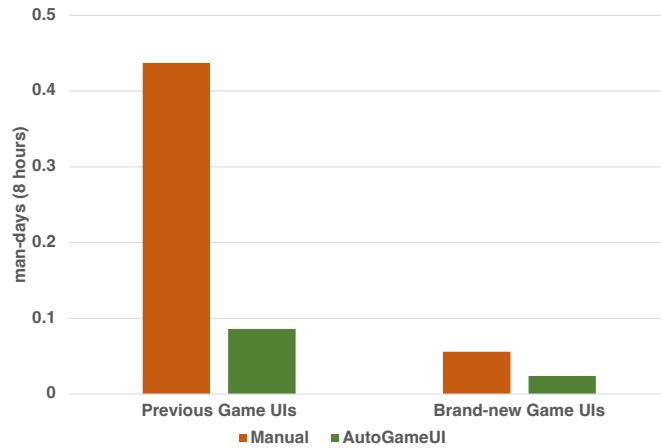


Fig. 13. Comparison of the average man-day cost between the AutoGameUI system and the traditional manual construction.

*4.4.6 Findings.* As shown in Fig. 13, the manual construction took an average of 0.437 man-days for the four previous sets, while using our system only consumed 0.086 man-days. On the other hand, the manual construction took an

average of 0.056 man-days for the brand-new set, while using our system took 0.024 man-days. The reason for the shorter average man-days on the brand-new set is that the 10 pairs of simple pop-up windows cost less time. Approximately, our tool achieved a 2-5 times speed improvement.

## 5　LIMITATIONS & FUTURE WORK

Until now, our system is in the alpha release stage, and we are constantly refining the interactive web-based tool based on users' feedback. For example, although we have introduced convenient interactive features like bidirectional confidence values to guide and support users in manual construction, users still need to frequently expand the hierarchical structure of the UX design to verify the matching results. The verification process slows down the construction speed in real-world workflows. This is because the hierarchical tree and 2.5D main canvas (i.e. RGBA and depth rendering mode) cannot fully represent the 3D UI or UX design. Users have expressed their desire for the interactive tool to provide more intuitive visualization capabilities to aid manual reconstruction, such as a 3D hierarchical preview on the main canvas.

Additional limitations pertain to our algorithm and dataset. Compared to previous works, our method considers two crucial attributes: hierarchy and rendering order. However, we do not explicitly embed them into the multimodal representation. Instead, they serve as the objective constraints in correspondence matching. Future work can integrate the hierarchical relationship and the rendering order into the transformer-based models using positional encoding, potentially enhancing the performance of the multimodal representation. Furthermore, when we attempt to tackle the cohesive issue in game UI development, the game project teams cannot provide sufficient UI and UX designs for our GAMEUI dataset. The reason is that earlier UI and UX designs have been lost as game design scripts are updated, which inevitably affects the building of the GAMEUI dataset. To augment the GAMEUI dataset, we have employed the interactive tool to record the submitted data of internal users.

In game UI development, there remain many topics that deserve to be explored, such as how to automatically optimize the memory usage of game UI resources and how to automatically generate UX prototypes from UI designs. Some of these topics are closely related to classic computational UIs but need more in-depth investigation.

## 6　CONCLUSION

We present a novel system to integrate UI and UX designs into cohesive user interfaces in game UI development. Our system allows users to construct cohesive user interfaces in a web-based tool, which supports both automatic and manual operations. For automatic construction, our system exploits a two-stage multimodal learning pipeline to capture the comprehensive representations of UI and UX designs and determine their correspondences by using grouped cross-attention modules and flexible integer programming. Based on the correspondence, the visual attributes of the UI design are seamlessly integrated into the logic structure of the UX design. For manual construction, our system offers multiple perspectives and interactive features to gain a user-friendly experience. To guarantee the visual and logical consistency of the cohesive user interfaces, we introduce a universal data protocol, which also facilitates the application of our system across different platforms. Utilizing our newly built GAMEUI dataset and RICO dataset, we demonstrate the effectiveness of our system through extensive experimental results. These results show that our system performs well in maintaining the coherence of the constructed interfaces with the original UI and UX designs.

## REFERENCES

[1] Philip Anastassiou, Jiawei Chen, Jitong Chen, Yuanzhe Chen, Zhuo Chen, Ziyi Chen, Jian Cong, Lelai Deng, Chuang Ding, Lu Gao, et al. 2024. Seed-TTS: A Family of High-Quality Versatile Speech Generation Models. *arXiv preprint arXiv:2406.02430* (2024).

[2] Michaela Bacikova, Jaroslav Poruban, and Dominik Lakatos. 2013. Defining domain language of graphical user interfaces. In *2nd Symposium on Languages, Applications and Technologies (2013)*. Schloss-Dagstuhl-Leibniz Zentrum für Informatik.

[3] Yue Bai, Dipu Manandhar, Zhaowen Wang, John Collomosse, and Yun Fu. 2023. Layout representation learning with spatial and structural hierarchies. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 206–214.

[4] Ethan D Bolker. 1972. Transportation polytopes. *Journal of Combinatorial Theory, Series B* 13, 3 (1972), 251–262.

[5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *European conference on computer vision*. Springer, 213–229.

[6] Kang Chen, Zhipeng Tan, Jin Lei, Song-Hai Zhang, Yuan-Chen Guo, Weidong Zhang, and Shi-Min Hu. 2021. Choreomaster: choreography-oriented music-driven dance synthesis. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–13.

[7] Rui Chen, Mingyi Shi, Shaoli Huang, Ping Tan, Taku Komura, and Xuelin Chen. 2024. Taming Diffusion Probabilistic Models for Character Control. In *ACM SIGGRAPH 2024 Conference Papers*. 1–10.

[8] Niraj Ramesh Dayama, Simo Santala, Lukas Brückner, Kashyap Todi, Jingzhou Du, and Antti Oulasvirta. 2021. Interactive layout transfer. In *Proceedings of the 26th International Conference on Intelligent User Interfaces*. 70–80.

[9] Jesús A De Loera and Edward D Kim. 2013. Combinatorics and geometry of transportation polytopes: An update. *Discrete geometry and algebraic combinatorics* 625 (2013), 37–76.

[10] 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.

[11] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. 2022. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems* 35 (2022), 18343–18362.

[12] David Gonçalves, Manuel Piçarra, Pedro Pais, João Guerreiro, and André Rodrigues. 2023. " My Zelda Cane": Strategies Used by Blind Players to Play Visual-Centric Digital Games. In *Proceedings of the 2023 CHI conference on human factors in computing systems*. 1–15.

[13] Kamal Gupta, Justin Lazarow, Alessandro Achille, Larry S Davis, Vijay Mahadevan, and Abhinav Shrivastava. 2021. Layouttransformer: Layout generation and completion with self-attention. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1004–1014.

[14] Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. 2022. Layoutlmv3: Pre-training for document ai with unified text and image masking. In *Proceedings of the 30th ACM International Conference on Multimedia*. 4083–4091.

[15] Naoto Inoue, Kotaro Kikuchi, Edgar Simo-Serra, Mayu Otani, and Kota Yamaguchi. 2023. Layoutdm: Discrete diffusion model for controllable layout generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10167–10176.

[16] Yue Jiang, Changkong Zhou, Vikas Garg, and Antti Oulasvirta. 2024. Graph4GUI: Graph Neural Networks for Representing Graphical User Interfaces. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–18.

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

[18] 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.

[19] Ranjitha Kumar, Jerry O Talton, Salman Ahmad, Tim Roughgarden, and Scott R Klemmer. 2011. Flexible tree matching. In *Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence-Volume Volume Three*. 2674–2679.

[20] 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.

[21] Gang Li, Gilles Baechler, Manuel Tragut, and Yang Li. 2022. Learning to denoise raw mobile UI layouts for improving datasets at scale. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–13.

[22] Quan Li, Haipeng Zeng, Zhenhui Peng, and Xiaojuan Ma. 2021. Understanding players' interaction patterns with mobile game app UI via visualizations. In *Proceedings of the Ninth International Symposium of Chinese CHI*. 9–21.

[23] 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.

[24] 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.

[25] Yang Li, Si Si, Gang Li, Cho-Jui Hsieh, and Samy Bengio. 2021. Learnable fourier features for multi-dimensional spatial positional encoding. *Advances in Neural Information Processing Systems* 34 (2021), 15816–15829.

[26] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. 2023. Magic3d: High-resolution text-to-3d content creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 300–309.

[27] Philipp Lindenberger, Paul-Edouard Sarlin, and Marc Pollefeys. 2023. Lightglue: Local feature matching at light speed. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 17627–17638.

[28] Haiyang Liu, Zihao Zhu, Giorgio Becherini, Yichen Peng, Mingyang Su, You Zhou, Naoya Iwamoto, Bo Zheng, and Michael J Black. 2023. Emage: Towards unified holistic co-speech gesture generation via masked audio gesture modeling. *arXiv preprint arXiv:2401.00374* (2023).

[29] Dipu Manandhar, Dan Ruta, and John Collomosse. 2020. Learning structural similarity of user interface layouts using graph networks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII 16*. Springer, 730–746.

[30] Mayu Otani, Naoto Inoue, Kotaro Kikuchi, and Riku Togashi. 2024. LTSim: Layout Transportation-based Similarity Measure for Evaluating Layout Generation. *arXiv preprint arXiv:2407.12356* (2024).

[31] Konstantinos Panayiotou, Constantine Doumanidis, Emmanouil Tsardoulias, and Andreas L Symeonidis. 2024. SmAuto: A domain-specific-language for application development in smart environments. *Pervasive and Mobile Computing* 101 (2024), 101931.

[32] Akshay Gadi Patil, Omri Ben-Eliezer, Or Perel, and Hadar Averbuch-Elor. 2020. Read: Recursive autoencoders for document layout generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 544–545.

[33] Akshay Gadi Patil, Manyi Li, Matthew Fisher, Manolis Savva, and Hao Zhang. 2021. Layoutgmn: Neural graph matching for structural layout similarity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11048–11057.

[34] Laurent Perron and Vincent Furnon. [n. d.]. *OR-Tools*. Google. https://developers.google.com/optimization/

[35] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. 2022. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988* (2022).

[36] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics. https://arxiv.org/abs/1908.10084

[37] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. 2019. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 658–666.

[38] André Ribeiro and Alberto Rodrigues da Silva. 2014. Xis-mobile: A dsl for mobile applications. In *Proceedings of the 29th Annual ACM Symposium on Applied Computing*. 1316–1323.

[39] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.

[40] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing* 568 (2024), 127063.

[41] Atsushi Takada, Daichi Yamazaki, Yudai Yoshida, Nyamkhuu Ganbat, Takayuki Shimotomai, Naoki Hamada, Likun Liu, Taiga Yamamoto, and Daisuke Sakurai. 2023. Genélive! generating rhythm actions in love live!. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 5266–5275.

[42] Ashish Vaswani. 2017. Attention is all you need. *arXiv preprint arXiv:1706.03762* (2017).

[43] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *nature* 575, 7782 (2019), 350–354.

[44] Zilong Wang, Mingjie Zhan, Xuebo Liu, and Ding Liang. 2020. Docstruct: A multimodal method to extract hierarchy structure in document for general form understanding. *arXiv preprint arXiv:2010.11685* (2020).

[45] Jeremy Warner, Kyu Won Kim, and Bjoern Hartmann. 2023. Interactive Flexible Style Transfer for Vector Graphics. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–14.

[46] Jason Wu, Amanda Swearngin, Xiaoyi Zhang, Jeffrey Nichols, and Jeffrey P Bigham. 2023. Screen correspondence: Mapping interchangeable elements between uis. *arXiv preprint arXiv:2301.08372* (2023).

[47] 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.

[48] Pengfei Xu, Yifan Li, Zhijin Yang, Weiran Shi, Hongbo Fu, and Hui Huang. 2022. Hierarchical layout blending with recursive optimal correspondence. *ACM Transactions on Graphics (TOG)* 41, 6 (2022), 1–15.

[49] Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. 2020. Layoutlm: Pre-training of text and layout for document image understanding. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*. 1192–1200.

[50] Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, Wanxiang Che, et al. 2020. Layoutlmv2: Multi-modal pre-training for visually-rich document understanding. *arXiv preprint arXiv:2012.14740* (2020).

[51] Kota Yamaguchi. 2021. Canvasvae: Learning to generate vector graphic documents. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5481–5489.

[52] Deheng Ye, Guibin Chen, Peilin Zhao, Fuhao Qiu, Bo Yuan, Wen Zhang, Sheng Chen, Mingfei Sun, Xiaoqian Li, Siqin Li, et al. 2020. Supervised learning achieves human-level performance in moba games: A case study of honor of kings. *IEEE Transactions on Neural Networks and Learning Systems* 33, 3 (2020), 908–918.

[53] Deheng Ye, Zhao Liu, Mingfei Sun, Bei Shi, Peilin Zhao, Hao Wu, Hongsheng Yu, Shaojie Yang, Xipeng Wu, Qingwei Guo, et al. 2020. Mastering complex control in moba games with deep reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 6672–6679.

[54] Michael Yin and Robert Xiao. 2022. The Reward for Luck: Understanding the Effect of Random Reward Mechanisms in Video Games on Player Experience. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–14.

[55] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. 2023. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3836–3847.

[56] Xu Zhong, Jianbin Tang, and Antonio Jimeno Yepes. 2019. Publaynet: largest dataset ever for document layout analysis. In *2019 International conference on document analysis and recognition (ICDAR)*. IEEE, 1015–1022.

[57] Jichen Zhu, Jennifer Villareale, Nithesh Javvaji, Sebastian Risi, Mathias Löwe, Rush Weigelt, and Casper Harteveld. 2021. Player-AI interaction: What neural network games reveal about AI as play. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–17.