

REFINE: Streamlining UI Mockup Iteration with Research Findings

Donghoon Shin

dhoon@uw.edu

University of Washington

Seattle, WA, USA

Bingcan Guo

bguoac@uw.edu

University of Washington

Seattle, WA, USA

Jaewook Lee

jaewook4@cs.washington.edu

University of Washington

Seattle, WA, USA

Lucy Lu Wang

lucylw@uw.edu

University of Washington,

Allen Institute for AI

Seattle, WA, USA

Gary Hsieh

garyhs@uw.edu

University of Washington

Seattle, WA, USA

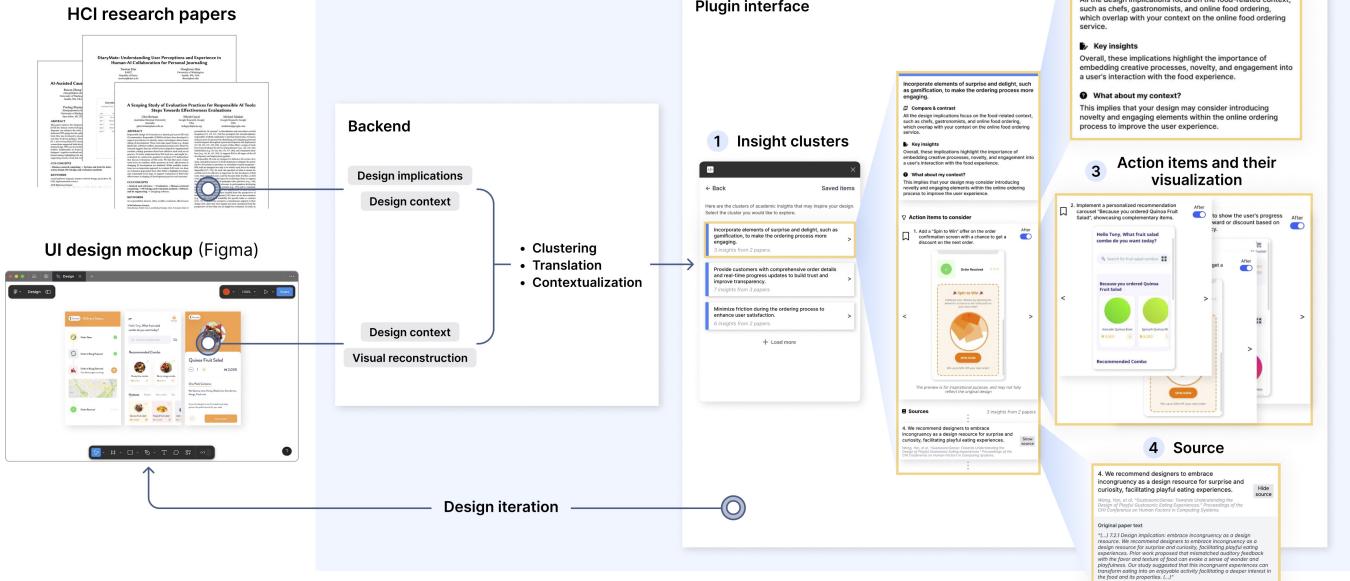


Figure 1: An overview of the REFINER system. When a user selects a UI design mockup on their Figma canvas, REFINER extracts the design context and retrieves relevant research papers from a scholarly repository. It then synthesizes design implications from these papers and organizes them into ① clusters. Each cluster includes ② an overview of the design implications tailored to the designer's specific context. Based on these implications, REFINER suggests ③ action items visualized on a reconstruction of the designer's mockup screens, and allows users to explore the ④ source papers that informed each cluster.

ABSTRACT

Although HCI research papers offer valuable design insights, designers often struggle to apply them in design workflows due to difficulties in finding relevant literature, challenges with understanding technical jargon, a lack of contextualization, and limited actionability. We present REFINER, a Figma plugin that supports real-time design iteration by surfacing contextualized insights from research papers. REFINER identifies and synthesizes relevant design implications from HCI literature, and tailors this research evidence to a specific design mockup by providing actionable visual guidance on how to update the mockup. To assess the system's effectiveness,

we conducted a technical evaluation and a user study. Results show that REFINER effectively synthesizes and contextualizes design implications, reducing cognitive load and improving designers' ability to integrate research evidence into UI mockups. This work contributes to bridging the gap between research and design practice by presenting a tool for embedding scholarly insights into the UI design process.

1 INTRODUCTION

Scholarly papers in HCI serve as an invaluable repository of emerging design implications. Guided by research findings, these implications are intended to support practical design decisions. Yet, design implications within these papers are often underutilized and fail to be effectively integrated into the workflows of design practitioners [13]. Prior research in translational science has identified several core challenges that impede the translation of scholarly knowledge into practical design action, including difficulties in finding relevant scholarly resources [8, 13]¹, adapting the communicated content and format to make the insights more understandable to designers [13, 38, 45], and contextualizing design implications within the designer’s action space [13, 38, 46]. These challenges are particularly pronounced during design prototyping, when designers who often work under time constraints iterate rapidly on emerging design ideas and mockups, making it impractical to continuously retrieve, translate, and contextualize scholarly insights in real-time [20, 48].

Recent work has explored making research papers more accessible to broader audiences. For instance, several research prototypes (e.g., personalized paper alerts [31], Q&A systems [7], recursive abstract exploration [21]) have been proposed to support the awareness and consumption of jargon-filled papers. However, these methods largely focus on addressing retrieval issues or summarizing paper insights, and do not directly address the needs of design practice. In the design domain, a handful of recent works have introduced tools to translate design implications into prescriptive formats (e.g., design cards [39, 46]), demonstrating potential to convey scholarly insights. Yet, these efforts have not addressed the retrieval problem or significantly improved the actionability of the research, which is especially critical during the prototyping stage, where identifying relevant insights and getting actionable recommendations are key to supporting iterations.

We introduce REFINE (REsearch FINDings for Evidence-informed UI design iteration), the first AI-powered design support system that facilitates real-time, in-situ UI mockup iteration by leveraging design implications from published research works. REFINE allows designers to draw inspiration from the latest research with minimal manual effort, by automatically retrieving, translating, and contextualizing design implications present in HCI research papers—directly within their existing design workspace (*i.e.*, Figma). To achieve this, our work addresses key challenges in translational science by (i) eliciting dimensions from each paper that represent what designers perceive as ‘relevant’ when consuming scholarly knowledge, and using these to retrieve papers, and identify and cluster design implications, (ii) generating summaries and translating scholarly insights by drawing analogies to bridge the gap between research findings and the mockup design, and (iii) providing ‘action items’ for their design and visualizing how these changes could be applied on a designer’s mockup reconstructed as an HTML.

REFINE leverages both large language models (LLMs) and vision-language models (VLMs) for retrieval, mockup understanding, paper understanding, and generating actionable insights.² To evaluate the effectiveness and accuracy of the system, we conducted a series of technical evaluations of intermediate model outputs on examples of mobile UI design iteration. Using a dataset constructed from research papers and mobile UI mockup design examples, we measured and reported the latency for generating each component in REFINE. Additionally, our comparative study evaluating REFINE components against alternative LLMs, input modalities, and techniques justifies our design choices. Finally, we assessed the reliability of VLM-driven mockup understanding and HTML generation³ within our system (*e.g.*, eliciting design dimensions, generating action items), finding that the outputs were largely faithful to their original inputs.

We further conducted a within-subjects user study with professional designers and design students in UI/UX ($N = 12$), where participants engaged in prototyping iterations both with and without REFINE. Our findings indicate that designers found the use of REFINE significantly less burdensome compared to when they tried to find and utilize research insights without it. Furthermore, REFINE improved the communication of scholarly insights across key metrics [44], *i.e.*, improving generativity, inspirability, actionability, validity, generalizability, and relevance, without undermining originality. Also, participants were able to create significantly more design edits and quickly reach design iterations on the canvas when using REFINE. Qualitative interviews revealed several factors contributing to these improvements, including reduced workload and a better understanding gained through visualizing action items. We also discuss potential future enhancements for supporting UI mockup iterations with research findings.

To summarize, we contribute the following:

- REFINE, a system that supports UI mockup iterations by retrieving relevant scholarly design implications, translating them into actionable insights, and contextualizing them within the current design mockup;⁴
- Results from an evaluation study, consisting of both technical evaluations and a user study, demonstrating the system’s ability to enhance the design iteration process by making scholarly insights more accessible and usable;
- Insights into how evidence-informed design can be further facilitated through leveraging research findings in the design workflow.

2 RELATED WORK

2.1 Translational Science for Design

In the field of human-computer interaction (HCI), much of the literature provides design implications (also known as guidelines) that could potentially be integrated into designers’ workflows to guide their design iterations and improve their design. Yet, these

¹While prior literature has used both *discovery* and *retrieval* to describe the process of finding relevant literature, for ease of reading, discovery and retrieval will be referred to for the rest of the paper as *retrieval*.

²In this work, we refer to models as LLMs when both the input and output are text, and as VLMs when an image is used as input.

³LLMs/VLMs have demonstrated a strong understanding of HTML [4, 11, 17, 26, 33]

⁴The codebases for both (i) the frontend plugin and (ii) the backend server are available here: [link anonymized for submission].

resources often remain underutilized by design practitioners, viewing them as impractical or disconnected from their needs [13, 38]. Previous research in translational science for design has pinpointed several obstacles that make it challenging to consume and leverage scholarly insights in design practice [13, 14, 46].

One major hurdle is finding and retrieving relevant resources [8, 13]. Designers often find it difficult to craft search queries that can identify papers relevant to their specific context [8, 13]. For example, a designer looking for UI patterns to boost accessibility might not know to search for terms like ‘inclusive design heuristics,’ leaving valuable design implications out of reach. This problem is worsened by existing paper search engines, which largely rely on basic metadata (e.g., titles, abstracts) for search (e.g., [6]), offering little depth or precision for design practitioners.

Another challenge is translating scholarly insights in a way that helps designers understand their use in design work [13, 38]. Despite the HCI community’s common practice to communicate design implications in published papers, designers often find these papers difficult to read [13]. Papers are often full of jargon, and the design implications may not be clearly communicated in a way that highlights insights that are generative, inspiring, and actionable [44].

Lastly, scholarly design implications are not sufficiently contextualized within a designer’s specific action space [13, 38, 46]. This problem is exacerbated by the fact that research papers themselves rarely match a designer’s project’s exact constraints or goals, such as specific target users or platforms [13]. While past work suggests even out-of-domain studies can offer valuable insights to support individual designers’ needs [45], designers often lack the know-how or tools to reinterpret these findings effectively for their unique situation, leading them to perceive the research as less useful or its implications as too generalized [13]. Additionally, even when designers dig into relevant papers and grasp the concepts, spelling out the practical and actual next actions to take within their own working design remains difficult. For instance, insights from a paper on online learning for older adults might hold value but fail to offer direct, actionable guidance for a designer working on a mobile design for seniors. In fact, prior research on auto-generated tools demonstrated that transforming scholarly insights into concise formats, while helpful, was insufficient to be fully supportive of the ‘what should I actually do in my design?’ question [13, 46].

To tackle these key challenges in translational science for design, we present a systematic approach that weaves HCI research findings directly into the iterative UI design process. Our tool simplifies resource discovery with context-aware paper retrieval tailored to designers’ needs, supports the translation of insights for specific design scenarios, and serves up practical, tailored recommendations for immediate integration. As such, we aim to bridge the gap between scholarly knowledge and design practice, making HCI research a more useful asset for everyday design work.

2.2 AI-supported UI Design Iterations

The process of refining UI designs based on feedback, testing, or newly added insights [3, 37] is critical throughout the UI design

lifecycle—from early-stage wireframing to high-fidelity prototyping [20]. Designers continuously refine interfaces to enhance usability and creative expression, ensuring alignment with user needs and design principles, while validating design choices and incorporating new knowledge effectively. AI has been supporting this process by automating the mundane tasks to allow designers to concentrate on mentally demanding and creative work. Recent advancements in AI, especially in LLM/VLMs, are providing designers with useful methods to approach the iterative process, enabling innovative ways to generate ideas, assess designs, and incorporate insights.

One key area where AI provides support for UI design iterations is in exploring inspirational UI examples. While designers traditionally relied on manually searching through image repositories [52], AI supports are being increasingly used to enhance semantic search capabilities to facilitate fast and relevant queries. For instance, building on foundations laid by large-scale UI datasets (e.g., [15, 32]), systems like VINS [10] enabled visual search by leveraging embedded screen segments. Another work explored the use of multimodal LLMs to capture deeper semantic categories even without the screen-associated app metadata for more refined search results [42]. However, merely providing inspiration might be limited in that designers often seek supporting evidence behind AI recommendations to substantiate them [55]—a gap that these inspiration-focused approaches may not fully address. Furthermore, adapting retrieved examples effectively to specific design contexts remains a challenge [28].

Beyond inspirational support, AI is increasingly being applied to provide more contextualized assistance within the current design mockup. For instance, LLMs have begun to assist designers in understanding their mockups, such as by facilitating conversational interactions with UI designs [49]. Not limited to understanding UI, prior work has also proposed more direct forms of support, such as interfacing UI design iteration with code development loop [36] and enhancing UI iterations by blending example design components [35] using LLMs. While many of these systems operate in specialized, custom-built environments, a handful of efforts have extended such capabilities to widely adopted platforms like Figma—as seen in tools that automate usability evaluation using heuristics [19], though their application remains limited in providing textual feedback using a predefined set of design guidelines.

Building on this body of research, in this work, we propose a novel approach that orchestrates the rich knowledge contained within and retrieved from published research papers and AI to support UI design iterations. Specifically, we demonstrate using the generative and understanding capabilities of LLM/VLMs to retrieve relevant design insights from scholarly literature, translate these findings into actionable suggestions, and contextualize them within a designer’s specific ongoing design iteration. We posit that this method addresses the need for scholarly research-backed guidance—moving beyond generic heuristics or serendipitous inspiration—while significantly reducing the effort required to discover and apply these scholarly insights, integrating seamlessly into designers’ workspace.

3 THE REFINER SYSTEM

To help designers seamlessly consume scholarly knowledge relevant to their work and integrate it into their design iterations, we developed REFINER (see Figure 1), an in-situ design support system powered by design implications from HCI research papers.

REFINER is implemented as a Figma plugin. Users of REFINER initiate the process by selecting a design mockup within Figma. REFINER then analyzes the selected mockup, automatically extracting its design context across multiple dimensions. Users are prompted to verify this extracted context, ensuring that the system understood their design accurately. Behind the scenes, the system uses this context to search its repository and retrieve the most relevant scholarly papers from a preprocessed set of papers. REFINER identifies pertinent design implications from these papers and presents them back to the users, organized into clear design insight clusters.

Within each insight cluster, users of REFINER are provided with actionable insight on how the retrieved design implications can be applied to their own design mockup. They receive a summary of the cluster along with insights specifically tailored to the context present in their mockup. Additionally, the system provides users with a set of action items, illustrating potential design changes that could be made to the mockup inspired by the research implications. All transformed design insights are linked to their source research papers. Following, we detail the design and technical considerations of each REFINER system component.

3.1 Preprocessing

3.1.1 Building the paper index. Unlike traditional paper repositories designed for general use, REFINER aims to help designers retrieve papers relevant to their design work and to directly leverage the papers’ design implications. We therefore chose to focus on indexing these papers on key dimensions of their design context, which influence how designers evaluate the relevance of research papers. Prior research [45] found that designers consider several core dimensions when determining the relevance of research papers: (i) target user, (ii) domain (*i.e.*, the industry or field the design is intended for), (iii) modality (*i.e.*, the primary mode of interaction or medium used in the design), (iv) pain point (*i.e.*, the main problem or challenge the design addresses for the user), (v) client (*i.e.*, the entity or stakeholder commissioning or benefiting from the users interacting with the design), and (vi) metric (*i.e.*, the key performance indicator (KPI) used to measure the design’s success). We utilized these dimensions to index papers in REFINER, and use them as a retrieval mechanism.

Starting with the raw paper PDFs, REFINER converts each file into structured XML using GROBID [34]—a machine learning-based paper parsing model. From the XML, we instructed the LLM to extract the above six dimensions of each paper’s design context. Since not all papers contain all dimensions, the model is prompted to abstain from responding when a dimension is absent. The text for each extracted dimension was then converted to vector representations using a text embedding model (Google’s text-embedding-004 [24]). If a dimension is absent or not evident in the paper, an embedding of an empty string (*i.e.*, “”) is assigned instead; if the model is unable to identify any design context in the paper, it is instructed to return all entities as “”, and the paper is excluded from the retrieval index. REFINER stores all embeddings as vector arrays in our cloud storage.

For our system demonstration and user study, we used the exhaustive list of papers ($N = 1060$) from the ACM CHI ’24 proceedings. Since CHI is a broad HCI conference encompassing various subfields [2], we believed this selection would capture a diverse range of domains and ensure topical generalizability for our implementation. Similarly, a paper index could be constructed from alternate proceedings or curated lists of papers and easily substituted into our system.

3.1.2 Identifying design implications in papers. While HCI papers often include explicit design implications or guideline sections, many convey these insights implicitly, making it hard to adopt a rule-based approach. On such an account, Sas *et al.* [44] defined a ruleset characterizing design implications, defining their functions, taxonomy, sources, and heuristics. Motivated by prior work that has shown the effectiveness of using rulesets or heuristics for locating and identifying relevant entities within a document (*e.g.*, [21]), we provided this ruleset as instructions to an LLM to identify design implications from HCI papers.

From the XML representation of each paper, REFINER utilizes LLMs to identify a set of design implications contained therein, if any. The model is prompted with the definition of design implications from [44], and to output an array of design implications, each of which contains (i) the design implication text, (ii) the paragraph from the paper containing the design implication, and (iii) the rationale for the model’s derivation. If the design implication is presented in a complete and original form within the paper, the model uses it as-is; otherwise, the model may make slight adjustments to the design implication to ensure the implication is self-contained and clearly conveys its full meaning. Since not all HCI papers have design implications [44], the model is instructed to return an empty array when none are identified.

3.2 Retrieving Papers & Design Implications

3.2.1 Understanding the design mockup. To identify papers relevant to the user’s design mockup, REFINER extracts six key dimensions from the design mockup displayed on the user’s Figma canvas, using a method similar to our paper preprocessing approach. First, the user opens the plugin and selects the mockup screen(s) for iteration. REFINER converts these selected mockup screens into an image. Using this base64-encoded PNG image data as input, REFINER applies similar instructions and prompts as before to a VLM, analyzing the design’s visual representation to extract its dimensions. Again, if a certain dimension is missing, the model is prompted to return an empty string. The text of each extracted dimension is embedded using the same text embedding model used for the papers.

3.2.2 Retrieving papers relevant to the current design. We use the array of embeddings corresponding to the six dimensions extracted from the user’s mockup screens as our query vector: $m = [m_1, \dots, m_6]$, where each m_i corresponds to one of the six dimensions. Similarly, the paper vector index consists of embedding arrays representing the six dimensions of each paper in the repository:

$$P = [[p_{11}, \dots, p_{16}], [p_{21}, \dots, p_{26}], \dots]$$

The system prioritizes embeddings of the design dimensions that are not missing from the design mockup or the paper (*i.e.*, valid dimensions). We compute a summed embedding over all valid

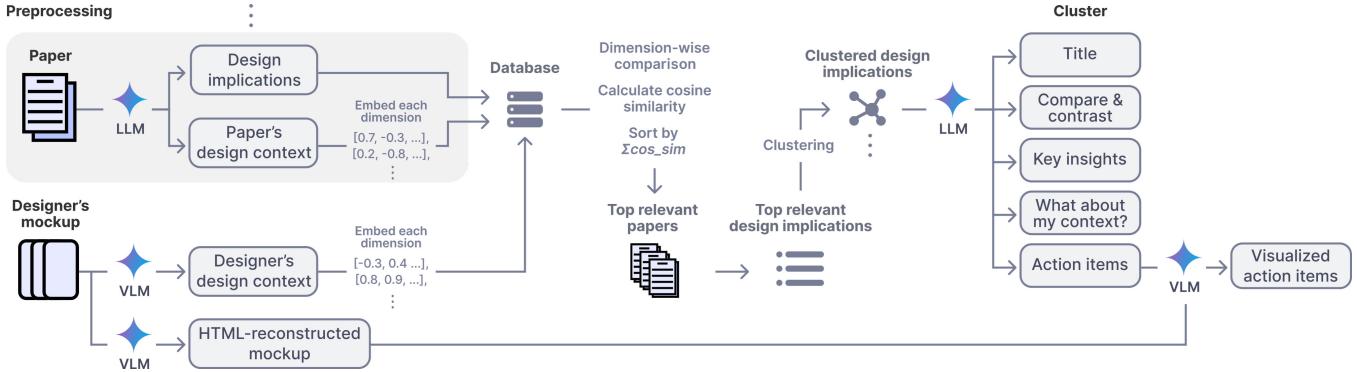


Figure 2: Pipeline for generating the components provided by REFINE

dimensions for the mockup representation as $S_m = \sum_{i \in \text{valid}} m_i$ and the corresponding dimensions of each paper in our retrieval index as $S_{p_i} = \sum_{j \in \text{valid}} p_{i,j}$. REFINE retrieves papers based on the cosine similarity between S_m and each S_{p_i} as:

$$\text{sim}(S_m, S_{p_i}) = \frac{S_m \cdot S_{p_i}}{\|S_m\| \|S_{p_i}\|} \quad \forall p_i \in P$$

To reduce information overload for downstream tasks, our algorithm retains the 8 most similar papers, though the user can always change this number in our system settings. All design implications from these papers are used in subsequent steps.

3.2.3 Clustering the design implications. REFINE clusters these design implications into themes, enabling the generation of coherent groups of design insights. This approach eliminates the need to analyze each implication individually while enhancing validity by synthesizing insights from multiple sources [44].

REFINE employs hierarchical clustering⁵ over the embeddings of the design implications to identify clusters. These embeddings are represented as $\mathcal{D} = \{d_1, d_2, \dots\}$, where each design implication d_i has an associated embedding vector $d_i \in \mathbb{R}^n$. To determine the optimal number of clusters, the system tests various values of $n_{\text{clusters}} \in \{2, 3, \dots, n_{\text{max}}\}$ and computes the silhouette score [43] for each clustering.

For each data point $d_i \in \mathcal{D}$, the silhouette score is defined as:

$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

where a_i is the average distance from d_i to other points in the same cluster, and b_i is the average distance from d_i to points in the nearest cluster that d_i is not part of. The silhouette score for a clustering C is then the mean over all points:

$$S(C) = \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} s_i$$

The optimal number of clusters, n_{best} , is the one that maximizes the mean silhouette score:

$$n_{\text{best}} = \underset{n_{\text{clusters}}}{\operatorname{argmax}} S(C_n)$$

⁵For hierarchical clustering, we use average linkage and cosine distance as the distance metric.



Figure 3: List view of the clusters generated by REFINE

To select n_{max} , we conducted an exploration of our clustering algorithm using 20 UI mockup examples. Our results revealed that the median number of clusters with the highest silhouette score was 6, and no example with more than 10 clusters achieved the highest score. Following this, we set n_{max} to 10 to optimize computational efficiency.

The implications are grouped into n_{best} clusters, and each cluster is assigned a label. Let C_k represent the set of design implications in cluster k . These groups are stored as:

$$\mathcal{G} = \{C_1, C_2, \dots, C_{n_{\text{best}}}\}$$

The final output for each cluster C_k can be represented as:

$$\mathcal{R}_k = \{(\mathcal{I}_k, \mathcal{S}_k) \mid C_k = \{d_1, d_2, \dots\}\} \quad \forall C_k \in \mathcal{G}$$

where \mathcal{I}_k is the list of original implications in cluster k , each element of which contains the paper ID, original implication text, and its source paragraph from the paper, and \mathcal{S}_k denotes the set of translated insights derived from each cluster, which we describe in Section 3.3.

In summary, REFINE retrieves the papers relevant to the user-provided design, and clusters design implications present in these papers to generate a set of synthesized design insights.

3.3 Translating Insights based on Designers' Design Context

3.3.1 Identifying similarities and differences between design implications and designer context. Although design implications with similar semantics are grouped together, they do not always share the same design context. Previous research has highlighted the importance of understanding both the similarities and differences between the context of the papers and the designer's design context to enhance comprehension of how the implications might be applied [45]. To support this, REFINE provides an overview of



Figure 4: Contents of the cluster and their formulation

similarities and differences (*i.e.*, *Compare & contrast*) between the designer’s context and the design implications of the cluster.

Initially, we calculated pairwise cosine similarities at the dimension level and presented the comparisons and contrasts accordingly. However, with six dimensions in play, this approach resulted in overly lengthy descriptions and computational complexities. Instead, drawing from prior work that demonstrated the effectiveness of LLMs in comparing and contrasting multiple entities [50, 53], we decided to leverage LLMs to identify and describe key similarities and differences. To achieve this, the LLM is given a list of design implications, their source paragraphs, and the design context of each paper, and instructed to analyze similarities and differences in relation to designers’ perspectives on design contexts through few-shot prompting. The generated summary of these similarities and differences is presented to the user.

3.3.2 Tailoring design insights through analogy generation. It is unlikely that the design insights from a paper’s design implications would perfectly align with a designer’s specific context, which necessitates the need for translating the insights tailored to an individual designer’s context. To bridge this gap, we leverage the concept of analogical ideation [23, 54], a widely used approach for adapting design insights across different contexts. Analogical ideation emphasizes the relations between surface-level information [23]—in this case, the relationship between a design implication and its original design context. This approach has been extensively applied in design to facilitate the transfer of insights from one context to another [12, 22, 23, 54].

First, the LLM is prompted to define the relation [45, 54] between each design implication within a cluster and its paper’s design context using Chain-of-Thought prompting [51]. The model is then prompted to generate a summary of these relations, which serves as the cluster summary (*i.e.*, *Key insight*). This relational summary is then applied to the designer’s specific context, producing a tailored design insight that the designer can reference (*i.e.*, *What about my context?*), which was derived from both the cluster of design implications and their design context. Lastly, the model is prompted to transform this tailored insight into a call-to-action format, producing a catchy title (*i.e.*, *Cluster title*) that the designer can view from the list of clusters.

3.4 Contextualizing Insights within the Designer’s Design Mockup

3.4.1 Generating action items. Based on insights tailored to the designer’s specific design context, REFINE generates up to three actionable design suggestions that users can implement in their actual designs (*i.e.*, action items). To achieve this, we provided the

model with (i) an image representation of the user’s mockup and (ii) the cluster’s implications, analogical ideas, summarized relationships, and translated insights. To prevent generating duplicate action items that are already accounted for in the mockup screens, we instructed the model to provide justification in the rationale, explaining why the action item has not been implemented. Additionally, the model is instructed to return fewer than three action items or none if there are insufficient action items to be recommended given the current screens.

Also, we recognized that some action items could not be seamlessly visualized within the Figma workspace, including those involving animations/effects. To address this, we prompted the model to return a boolean flag indicating whether the application of an action item is visually representable based on the following criteria, given the input image representation of the mockup: (i) the action item can be directly applied to any of the current mockup screens as visual elements, and (ii) the model does not need to generate an additional screen to implement the action. If these conditions are met, REFINE generates a preview illustrating how the design would change upon applying the action item as follows; otherwise, it presents the action item in text format and informs users that a preview was not generated.

3.4.2 Constructing a preview of the mockup with the action item applied. To support the contextualization of a design item within the designer’s design, REFINE visualizes the mockup with the action item applied. This is enabled by two technical modules: (i) mockup reconstruction and (ii) action item visualization.

(i) Reconstructing the mockups into HTML to visualize the action items. Most commercial design tools, including Figma, allow users to export designs as JSON, yet have limited support for directly importing them or rendering the design, making programmatic manipulation for visualization purposes impractical. To tackle this, we opted for an approach that converts designs into HTML—a format that LLM/VLMs have demonstrated a strong understanding of [4, 11, 17, 26, 33]—allowing for the reconstruction of user mockups. A handful of previous works have attempted to reconstruct the website design to automate frontend development (*e.g.*, [47]), yet their focus was not on replicating visual design but on enabling functionalities. Thus, we aimed to visually replicate the mockup into HTML to enable the visualization of text-based action items within the mockup. While this may not perfectly replicate the original visual rendering, our goal is to provide designers with inspirational sources within the context of their mockup design, making this level of fidelity sufficient. Moreover, this approach not only enables visual representation but also allows REFINE to modify designs through textual queries for applying action items, enabling us to visualize action items within the user’s mockup design to better inspire designers through visual digestion of the action item.

To reconstruct the Figma mockups, various export options from the Figma plugin⁶ can be leveraged to provide the model with sources for reconstruction, including visual formats (*e.g.*, image files) and structured text-based representations (*e.g.*, JSON) of the mockup screens. While both input modalities have demonstrated efficacy in UI understanding when coupled with LLMs (*e.g.*, [49]

⁶<https://www.figma.com/plugin-docs/api/ExportSettings/>



Figure 5: An example of mockup reconstruction in our pipeline. Although the reconstructed HTML may not perfectly match the original mockup, it faithfully mirrors the components in the original mockup and allows for visualization of action items.

vs. [19]), their effectiveness in *reconstructing* visually and semantically similar UI in comparison with the other modality remains unclear. Therefore, in Section 4.2.3, we conducted a comparative study to guide our choice of input modality.

While REFINEx is generating the clusters, the VLM concurrently takes the mockup as input, and is prompted to generate HTML code that visually replicates the mockup screen. Specifically, the model produces the screen in HTML with CSS styles embedded inline, streamlining the process of applying action items by screens in step (ii). Here, through a cursory exploration of prompting, we observed that the model tends to overfit to specific visual details (e.g., icons), increasing the latency, without a noticeable gain in accuracy. To address this, we guided the model to use emojis as stand-ins for visuals. Similarly, directly replicating or transferring original images from the mockup to the reconstructed mockup led to inconsistencies in positioning or distorted proportions with increased computational complexities. As such, we prompted the model to represent images using a container with a similar gradient background, instead of attempting direct replication.

To validate our ideas, we presented the original design and the reconstructed mockup to three UI/UX designers, all of whom unanimously indicated that these stand-ins did not compromise the semantic and visual understanding of the original design. Following this, we decided to adopt our design choice to improve efficiency while maintaining design clarity. Additionally, to reduce the latency of generating a mockup consisting of multiple screens at once, REFINEx is designed to run the VLM computation for reconstructing individual screens in parallel. Figure 5 illustrates an example of mockup reconstruction using our pipeline.

(ii) Applying an action item to the reconstructed mockup. Based on the reconstructed mockups from (i), REFINEx enables the visualization of action items on the mockups created in Section 3.4.1. Here, while we found the latency for task (i) to be manageable—since it runs with other tasks in parallel and user interactions and is generated only once for repeated use—the computation time for task (ii) is critical; given its role in visualizing action items for inspiration, and it needs to be generated multiple times per every cluster. To address this, instead of having the model generate an entirely new HTML output for each query, we designed an ‘edit-only generation’ approach that produces only the necessary edits to be made within the HTML code.

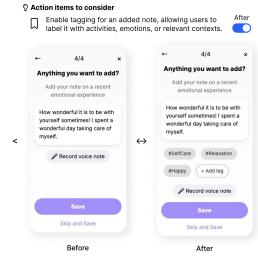


Figure 6: An example of an action item and its visualization. With the reconstructed HTML, REFINEx selects the appropriate screen(s) from the design mockup and visualizes how each action item can be applied. Users can toggle between the before and after visualization, assisting designers in grasping the differences. Users can also bookmark each action item.

Specifically, given (i) the user’s design rendering and (ii) the action item, we prompted the VLM to identify the indices of screens that need to be edited when generating action items. Then, when the user opens the cluster, it generates the required modifications to apply (ii) to the HTML mockup generated from (i). Specifically, the model outputs edits in three forms—*add*, *remove*, and *replace*—each linked to an ID of the HTML component, enabling comprehensive iterations in response to a wide range of plain language action items. For example, the edits for *replace* are generated through the following part of the system instruction:

```
...
[Rules (for 'replace')]
- When a certain element should be modified, target the smallest
  possible (lower-level) DOM elements to avoid redundant updates.
- Do not modify parent or sibling elements unless absolutely required
  for the change to work.
- {style-related instructions}

[Output format]
{
  "reference_element_id": String, // an id for the element that
  needs to be changed
  "edited_element": String, // an HTML code for the modified element
  "rationale": String
}
...
```

After generation, REFINEx uses an HTML parser to parse the HTML, locate the corresponding component from the HTML using its ID, and apply the changes. Figure 6 illustrates an example of applying an action item to the reconstructed mockup.

3.5 Implementation

The REFINEx system consists of two main components: (i) a backend server and (ii) a frontend web app. The backend server operates on Python 3.11 using Google Cloud’s Ubuntu infrastructure and handles the LLM/VLM computations and clustering. The frontend web app is built with SvelteKit⁷, a JavaScript-based web framework, and is integrated with Figma’s Plugin API⁸ to function as a Figma

⁷<https://svelte.dev/>

⁸<https://www.figma.com/plugin-docs/>

plugin and enable interaction between the Figma workspace and the plugin. For parsing HTML code, we used BeautifulSoup⁹ library.

For natural language and vision language tasks, various LLMs and VLMs can be substituted for the modular components.¹⁰ In the following section, we further describe our model choices, as well as input modalities for reconstructing the mockups, based on our comparative evaluations.

4 TECHNICAL EVALUATION STUDY

REFINE leverages LLM/VLMs to perform various component tasks in our pipeline. We justify our design choices, including selection of the base model, input modality, input content, and technique for visual reconstruction, based on comparative performance evaluations. Here, we describe these technical evaluations—how we measure the performance accuracy and latency of each module and compare them with alternative approaches.

4.1 Datasets for Technical Evaluation

We constructed datasets of (i) UI mockups and (ii) HCI papers to support each component of our technical evaluation. We used the UI mockup dataset to evaluate visual reconstruction and mockup design context understanding, and the HCI paper dataset to evaluate paper design context understanding, retrieval/clustering, and generating translated insights. We describe these datasets below.

4.1.1 UI mockup dataset. In this study, we specifically focused on mobile UI, following prior work on UI design support tools (e.g., [42, 49]) that used mobile UI as a testbed. Although several existing datasets (e.g., Rico [15], Enrico [32], ERICA [16]) support data-driven mobile UI mockup evaluation, these datasets were not suitable for our evaluation because they lacked the structured JSON representations available through Figma. Therefore, we opted to manually curate a dataset by retrieving UI mockups designed on Figma from the Figma community library.

Two authors from our research team reviewed Figma’s publicly accessible mobile app mockup design templates.¹¹ Starting with the top designs, the team manually screened and excluded those that were (i) too thematically similar to earlier designs or (ii) primarily layout-focused without meaningful semantic content. We continued this process until we reached a total of 50 mockups for inclusion, representing reasonably diverse applications. Each mockup consisted of 3 to 5 screens.

4.1.2 HCI paper dataset. Our system retrieves from all papers from the CHI ’24 conference proceedings [1]. For evaluation, we randomly selected and reviewed outputs corresponding to 50 full papers from this collection.

4.2 Validating Mockup Visual Reconstruction & Action Item Visualization

Among various system components, reconstructing the mockup in HTML is both dependent on nuanced visual structure understanding and essential to incorporating actionable design feedback,

⁹<https://www.crummy.com/software/BeautifulSoup/>

¹⁰The prompts used for running the LLM/VLM components of REFINE are available in our codebase: [link anonymized for submission].

¹¹<https://www.figma.com/community/mobile-apps>

making it a useful probing point for evaluating different model choices. Figma does not export HTML, but can provide both image exports and a JSON representation of the mockup. To explore alternatives and guide our choices, we evaluated several flagship commercial VLMs (*i.e.*, Gemini 2.0 Flash [25], Claude 3.5 Sonnet [5], and GPT-4o [40]) for HTML reconstruction.¹² For comparing input modality candidates, we tested three variants of input modality for HTML reconstruction: providing (1) only the image representation of the mockup, (2) only the exported JSON of the mockup, and (3) both the image and JSON.

We evaluated the following components using our UI dataset. For each of the 50 mobile UI mockups in the dataset, the model was instructed to generate an HTML representation based on the provided input modalities using an instruction prompt, with minor modifications to denote the different types of input that can be provided. To assess reconstruction accuracy, we embedded both the original mockup image and an image of the rendered reconstructed HTML using a visual transformer (vit_base_patch16_224 [18]), and computed their visual similarity using cosine distance to understand how well each model reconstructed the mockup. We also computed the latency of each model and input setting to understand performance implications.

4.2.1 Comparing base models for HTML reconstruction. Visual similarity and latency for all three models, with the input modality controlled as image only, are provided in Table 1 ($N = 50$). Although Claude 3.5 Sonnet achieved the highest visual similarity, its latency was significantly higher than the other models. Compared to GPT-4o, Gemini 2.0 Flash demonstrated a slightly lower latency with similar visual similarity. We selected Gemini 2.0 Flash as the VLM in our system as it provided the best trade-off between accuracy and efficiency.

Table 1: Comparison of VLMs for HTML reconstruction using image inputs (metrics are computed per-screen \pm stdev)

| Model | Visual similarity | Latency (seconds) |
|-----------------------|---------------------|--------------------|
| Gemini 2.0 Flash [25] | 0.7875 ± 0.1437 | 10.794 ± 3.021 |
| Claude 3.5 Sonnet [5] | 0.7963 ± 0.1590 | 17.006 ± 4.914 |
| GPT-4o [40] | 0.7852 ± 0.1330 | 11.060 ± 3.382 |

4.2.2 Comparing input modalities for HTML reconstruction. Table 2 ($N = 50$) provides a comparison across the three input modality settings, controlling for the base model. Using images alone as input resulted in the highest visual similarity. Combining images and JSON as input led to a small reduction in visual similarity. Providing JSON alone performed the worst, with lower accuracy and significantly increased latency. We hypothesize that this is because, without visual input, the model may have overfitted to the available structural cues in the JSON, leading to less effective visual reconstruction. Based on these results, we decided to use only images as the input for mockup HTML reconstruction.

¹²State-of-the-art commercial models at the time of our study, serving to justify our base model choice and demonstrate the efficacy of our approach. The pipeline of REFINE is modular, thus can be readily adapted to newer models as they emerge.

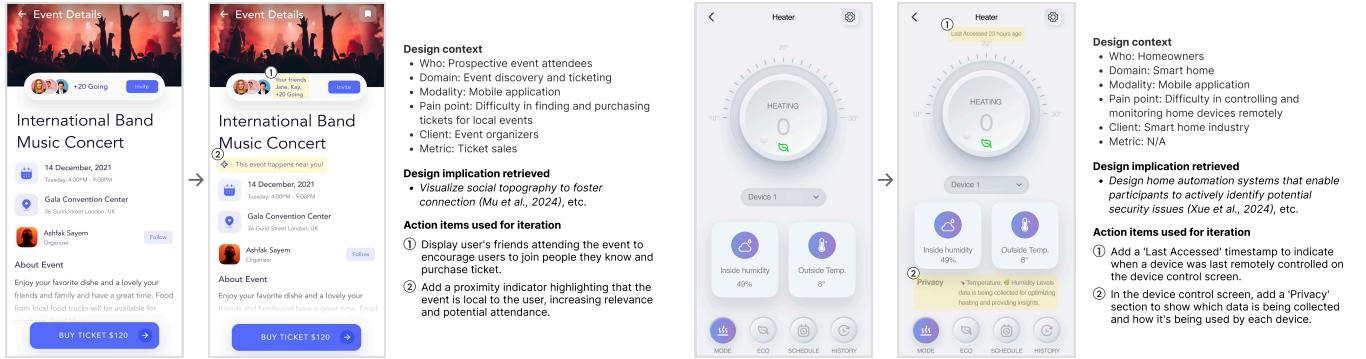


Figure 7: Design outputs from iterating two example mockups with REFINEx-generated action items. For each example, we present the extracted design context, an example of retrieved design implications, and action items used for the iteration. For conciseness, a single screen pair (before/after) is shown for each mockup, and edits are highlighted with a yellow background.

Table 2: Comparison of different input modalities on HTML reconstruction (base model is Gemini 2.0 Flash, metrics are computed per screen \pm stdev)

| Input modality | Visual similarity | Latency (seconds) |
|----------------|---------------------|--------------------|
| Image | 0.7875 ± 0.1437 | 10.794 ± 3.021 |
| JSON | 0.7661 ± 0.1314 | 15.432 ± 7.047 |
| Image + JSON | 0.7711 ± 0.1411 | 10.630 ± 2.427 |

4.2.3 Comparing techniques for visualizing action items. As detailed in Section 3.4.2, REFINEx prompts LLMs to generate action items as additions, deletions, or modifications on the HTML reconstruction. We do this instead of prompting the models to produce a complete HTML regeneration due to the significant latency improvement when generating the action item.

To validate this, we first ran REFINEx on our mockup dataset to identify evidence clusters and generate corresponding action items. We randomly selected 150 visually representable action items and generated visualizations of action items using both (i) REFINEx’s method, which generates the necessary edits and preprocesses them through an HTML parser, and (ii) a prompt for complete HTML regeneration, which takes the same input as (i) (*i.e.*, reconstructed HTML code and the action item text) and regenerates HTML from scratch. The research team manually evaluated whether each action item had been correctly applied to the screen as described in the action item text.

The results are shown in Table 3 ($N = 150$). Paired t-test revealed that our method significantly reduced processing time ($t = 26.34$, $p < .001$) while maintaining the accuracy of visualizing the action item, demonstrating the efficacy of our method for visualizing the action items.

4.3 Validating REFINEx’s Extractive and Generative Components

REFINE uses LLM/VLMs to interpret design contexts from both the designer’s mockup and relevant papers, helping match designers with inspiring research. Additionally, REFINEx relies on the LLM’s generative capability to generate action items based on the cluster

Table 3: Comparison of methods for generating visualizations of action items applied to designs (metrics are computed per action item \pm stdev)

| Method | Accuracy | Latency (seconds) |
|----------------------------|----------|--------------------|
| Complete HTML regeneration | 96.0% | 10.248 ± 3.488 |
| Edit-only generation | 95.3% | 2.913 ± 1.215 |

contents. Here, it is crucial to understand whether these generative components accurately reflect the provided inputs, as LLM/VLMs are known to be susceptible to hallucinations [30]. With our model and technique choices from Section 4.2.3, we ran a validation study to evaluate whether these generated contents are accurate and faithful to the provided inputs.

4.3.1 Assessing design context extraction from papers. We examined how REFINEx extracts design context dimensions from research papers. Applying the pipeline to the dataset of 50 HCI papers, the system generated six dimensions per paper, totaling 300 dimensions. To facilitate these, we developed a Streamlit-based annotation tool that allows annotators to view inputs and outputs side-by-side. Using our annotation tool, two annotators manually read both the paper and the dimensions, and assessed whether each dimension faithfully represented the design context described in the paper. When they determined that a dimension was not faithfully derived, a text field appeared for them to provide a justification.

Overall, 95.7% of the extracted dimensions faithfully represented the design context of the source paper. Nonetheless, we identified three key error patterns: misclassification ($N = 7$; *e.g.*, labeling ‘online survey’ as a modality rather than a research method), scope misalignment ($N = 4$; *e.g.*, specifying the target user as ‘users of LLMs’ when the system was designed for any individual), and focus misinterpretation ($N = 2$; *e.g.*, treating ‘incidental learning’ as a primary research domain when it was actually a study finding). The annotation demonstrated high inter-rater reliability, with a Cohen’s Kappa (κ) of 0.82.

4.3.2 Assessing design context extraction from mockups. To analyze the design context of mobile UI mockups, we applied REFINEx’s

extraction pipeline to 50 mockup sets. For each mockup, the system identified six distinct dimensions, yielding a total of 300 dimensions. Using the same annotation procedure as with Section 4.3.1, two annotators reviewed these mockup screens and extracted dimensions to assess extraction accuracy.

Our analysis showed that 94.3% of the extracted dimensions accurately captured the intended design context of the mockup. However, the annotations revealed a few classes of errors, including overgeneralization ($N = 5$; e.g., categorizing as ‘publisher’ without specifying the type of media), misclassification ($N = 4$; e.g., identifying ‘touch’ as a modality instead of ‘mobile app’), and incorrect metric identification ($N = 8$; e.g., extracting metrics that were not explicitly signaled in the design). The results again showed strong agreement between annotators ($\kappa=0.88$).

4.3.3 Assessing the relevance of generated action items to the source design implications. To evaluate how well the action items generated by ReFINE align with their corresponding source design implications from papers, we conducted a manual evaluation. From the dataset described in Section 4.3.2, we randomly sampled 50 (action item, design implication) pairs. Then, two annotators used our annotation tool to independently reviewed each pair, reading the original design implication and rating the relevance of the corresponding action item on a 5-point Likert scale (1: not relevant at all = the action item does not reflect or implement the core idea of the design implication in any meaningful way; 5: highly relevant = the action item clearly and accurately operationalizes the core idea of the design implication within a similar domain). Along with each rating, annotators were also asked to provide a brief written explanation of their judgment.

The evaluation revealed that, even after the ReFINE’s translational processes, the generated action items maintained a high relevance to their source design implications, with an average rating of 3.82/5. Only 14.0% of the pairs were rated 2 (slightly relevant) or below. Given the ordinal nature of the data, we used weighted Cohen’s Kappa to assess inter-rater reliability, which indicated strong agreement ($\kappa = 0.81$). The examples of annotators’ evaluation for each score, along with their rationale, are included in Table 5.

4.3.4 Assessing the relevance of action items to their assigned clusters. Additionally, we evaluated the alignment between action items and their corresponding cluster titles. Using the same set of action items from Section 4.2.3, two annotators were provided with corresponding sets of mockup images, cluster titles, and assigned action items. They assessed if each action item was faithful to its corresponding cluster title and the mockup image.

Overall, 96.0% of the action items were deemed well-aligned with their respective titles. Identified misalignments include: returning an action item that was already realized in the screens ($N = 1$), generating overly generic action items ($N = 3$), and producing action items with an irrelevant focus ($N = 2$). The results showed strong inter-rater reliability ($\kappa = 0.79$).

Furthermore, annotators assessed whether the action items were well-situated within the mockup ($\kappa = 0.74$). The model correctly identified the appropriate screens for 98.0% of generated action items. However, in three instances, it suggested updates to screens that already included the requested changes.

4.4 Overall Latency Analysis

Several ReFINE components operate on LLM/VLM computations, which may introduce latency. Given our model (*i.e.*, Gemini 2.0 Flash) and input modality (*i.e.*, image for mockup representation) choices based on our comparison in Section 4.2, we report the evaluation of latency for every component of ReFINE using our dataset, as outlined in Table 4. For components involving multiple computations in parallel or outputs from distinct units (*e.g.*, mockup reconstruction, generating cluster contents), we measured the time spent on each individual computation. The latency analysis for visualizing action items is based on the dataset used in Section 4.2.3, with its latency value referenced from that section.

Table 4: Overall system latency analysis (metrics are computed per task \pm stdev, *: per unit; otherwise per design)

| System step | Latency (seconds) |
|---------------------------------|--------------------|
| Eliciting design contexts | 3.167 ± 1.019 |
| Mockup reconstruction* | 10.794 ± 3.021 |
| Retrieval & clustering | 2.004 ± 0.400 |
| Generating translated insights* | 2.982 ± 0.326 |
| Generating action items | 4.435 ± 0.480 |
| Visualizing action item | 2.913 ± 1.215 |

5 USER STUDY

Our technical evaluation presented an overview of the performance of ReFINE components, comparison of their effectiveness to alternatives, and accuracy of the generated outputs. To further understand how ReFINE is perceived and utilized in actual design iteration, we conducted a within-subjects user study involving 12 designers.

5.1 Recruitment

We began by recruiting designers from three design-focused university communities and one online community. Our recruitment post specified that participants must either (i) be currently working as professional UI/UX designers or (ii) be students pursuing a professional UI/UX design degree, who have experience using Figma in design work. As a result, we recruited 12 participants; of all, 8 identified as female, 4 male, and the average age was 26.0 ($SD = 4.5$). The average length of their design experience was 4.2 years ($SD = 2.3$).

5.2 Procedure

Each study was conducted remotely via Zoom. Before starting, each participant was asked to install the Figma desktop app. Also, we retrieved two publicly available mobile UI mockups from the Figma mobile app library—one for a financial banking app and the other for a travel app (see Appendix A), each containing four key screens.

Once the user joined the study, we introduced the objectives of our study. Then, each participant went through two design conditions, with order randomized to minimize any ordering effects: (i) engaging with scholarly insights and iterating on the mockup using ReFINE-supported design (*i.e.*, ReFINE condition), and (ii) manually searching for papers in ACM Digital Library and iterating on the mockup by consuming knowledge by reading the papers (*i.e.*, baseline condition). Each condition was randomly assigned one of

the two UI mockups we prepared. To ensure a fair comparison with our current implementation that leverages the CHI '24 proceedings, we restricted the papers available for retrieval in both conditions to those published in CHI '24. Prior to each condition, we instructed participants on how to import the REFINER plugin into Figma, or access the ACM DL to search and retrieve full-text papers, based on the condition.

Each condition lasted 30 minutes. During the first 25 minutes, designers were asked to inform their designs with insights from research papers, and we encouraged them to iterate on the design without focusing heavily on the trivial visual aesthetics of their edits. During each condition, we notified them every 5 minutes to help them manage their time. Once the participants completed the design condition, they proceeded to a survey and filled out questionnaires for 5 minutes.

After completing both conditions, participants took part in a qualitative interview session, where they answered questions about their preferences and discussed the strengths, weaknesses, and potential enhancements for both components of REFINER and the overall system. The interview lasted approximately 20 minutes.

In total, the entire process lasted approximately 90 minutes. Each participant was compensated with a 50 USD gift card for their participation. The study protocol was reviewed and approved by the university's Institutional Review Board.

5.3 Measure & Analysis

5.3.1 Survey measures. First, by utilizing REFINER's design support, we hypothesized that designers could significantly lessen the effort required to extract scholarly insights from scholarly repositories. To evaluate the cognitive workload of participants, we employed the NASA-TLX workload index [27] on a 7-point scale.

Additionally, previous research in design has identified six core dimensions [44] for assessing the communication quality of design implications: validity, generalizability, originality, generativity, inspirability, and actionability. To determine if the quality of communicated design insights differed when using REFINER compared to not using it, we included these six dimensions as survey questions, along with the original definitions from the literature to ensure a shared understanding. Lastly, to measure the effect of translational processes, we added relevance as an additional measure, measuring the relevance of communicated insights with their design task.

5.3.2 Analysis. To analyze the quantitative results (*i.e.*, survey responses, behavioral data from usage logs), we initially checked and confirmed that every measure met the normality assumption using the Shapiro-Wilk test. After confirming normality, we conducted a paired t-test to assess the differences in these measures between the two conditions (REFINE versus baseline).

For the qualitative results, we conducted a thematic analysis [9] with a bottom-up approach. First, two authors manually read through the transcripts and identified the initial set of themes individually. Then, they regularly met to discuss and refine the themes, which repeated for three rounds. As a result, the authors identified the themes and corresponding quotes as detailed in Section 5.4.

5.4 Results

5.4.1 Overall perception of REFINER. Participants thought positively about using REFINER to support their UI mockup iterations. When comparing between REFINER and baseline conditions, 11 participants preferred using the plugin to support their design iteration. Only one participant mentioned that their preference depends on the context, but indicated that they would only prefer to manually gather design insights from research papers if they had unlimited time to conduct iterative paper searches, admitting that this timeframe is practically impossible.

As shown in Figure 8, participants found the REFINER-driven iteration to be significantly less burdensome than the baseline condition. They found that interacting with REFINER to guide their design iterations was significantly less mentally ($M = 2.92, SD = 1.51$ vs. $M = 6.33, SD = 0.78$; $t = 6.84, p < .001$), temporally ($M = 2.92, SD = 1.73$ vs. $M = 6.00, SD = 0.85$; $t = 5.41, p < .001$), and physically ($M = 2.42, SD = 1.38$ vs. $M = 4.58, SD = 1.98$; $t = 3.86, p < .01$) demanding, compared to the baseline condition. They also perceived the REFINER support as leading to better performance ($M = 5.08, SD = 1.62$ vs. $M = 2.75, SD = 1.66$; $t = 4.84, p < .001$), while requiring less effort ($M = 3.58, SD = 1.38$ vs. $M = 5.58, SD = 1.31$; $t = 3.13, p < .01$) and resulting in lower frustration ($M = 2.17, SD = 1.40$ vs. $M = 4.92, SD = 1.08$; $t = 5.40, p < .001$).

Our findings also indicate that the design insights communicated by REFINER exhibited enhanced communication qualities. Participants rated REFINER as providing insights that were more relevant ($M = 5.58, SD = 0.99$ vs. $M = 3.58, SD = 2.11$; $t = 2.97, p < .01$), generative ($M = 4.92, SD = 1.51$ vs. $M = 2.75, SD = 1.54$; $t = 3.03, p < .01$), inspirational ($M = 5.17, SD = 1.11$ vs. $M = 3.00, SD = 1.41$; $t = 3.68, p < .01$), actionable ($M = 6.00, SD = 1.13$ vs. $M = 3.17, SD = 1.75$; $t = 4.53, p < .001$), valid ($M = 5.58, SD = 1.08$ vs. $M = 3.42, SD = 1.68$; $t = 4.73, p < .001$), and generalizable ($M = 4.83, SD = 1.40$ vs. $M = 3.08, SD = 1.56$; $t = 2.96, p < .01$) compared to the baseline. Our translational process did not compromise perceived originality ($M = 3.83, SD = 1.40$ vs. $M = 3.92, SD = 1.38$; $t = 0.15, p = 0.44$).

Supporting these perceived improvements in workload and communication qualities, participants were able to make significantly more edits with REFINER within the same timeframe. On average, participants made 5.5 design edits ($SD = 1.4$) using REFINER, which was significantly higher than without the support, where they made 2.4 edits on average ($SD = 1.4$; $t = 6.59, p < .001$).

Participants also found REFINER's generation speed sufficiently fast for their usage scenario. Out of all participants, ten responded that latency in REFINER's content generation did not affect their design process. Two others noted that they 'at least did notice' the delay, although they felt it did not significantly impact their design iteration. The system displayed a circular progress bar while loading; two participants suggested that providing detailed text to explain what is being generated, along with indicating the actual progress, could help mitigate the perceived effects of latency even further: "*I think it (latency) is okay (...) Just a suggestion, but maybe you can put something like 'oh, we're loading these' or 'we're retrieving these insights.'*" (P12)

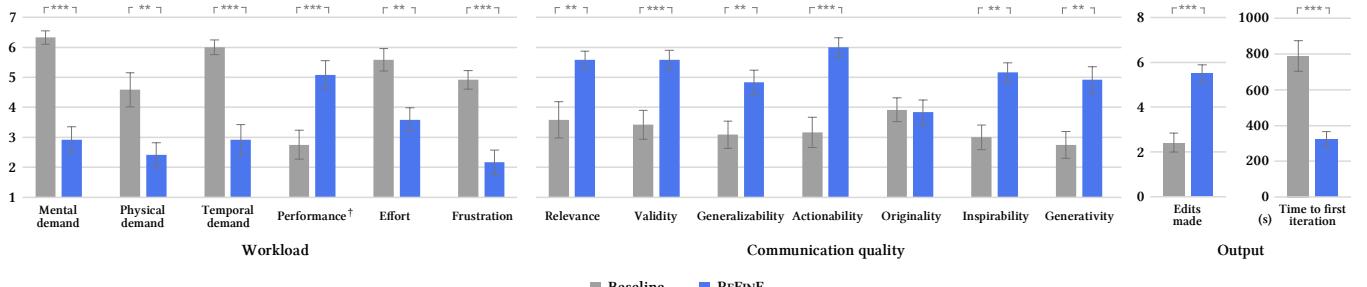


Figure 8: Perceived workload, communication quality, and the design output metrics from our user study. Participants in the REFINEx condition reported significantly reduced workload compared to the baseline condition, while finding the communicated design insights to be significantly improved. Additionally, participants made significantly more edits interacting with REFINEx, while reaching iteration more quickly. Error bars indicate standard errors. (: $p < .01$, ***: $p < .001$; †: reverse item)**

To this end, participants expected the REFINEx system to be highly beneficial for their future design work, as it provided scholarly evidence to support their iterations—resources they value but often struggle to utilize [13]. They illustrated the excessive iterations they typically face during the prototyping stage and the difficulty of looking up multiple sources to inform their edits, highlighting the potential of REFINEx in their workflows. Consequently, 10 participants expressed willingness to use the plugin in their future design process, and three interviewees directly inquired about its public release: “*When are you gonna publish this plugin? I would really love to use this in my designing.*” (P1)

Specifically, participants discussed several real-world scenarios where REFINEx could be particularly valuable. Four participants envisioned using it in corporate design teams, where justifying design iterations with scholarly-backed rationale is highly appreciated, viewing REFINEx as a powerful support for enhancing design iterations with well-founded insights: “*I would love to use that at work because sometimes being a designer is about presenting the design to all the stakeholders, especially the PMs. So it would be nice to have some of the evidence to back the design decisions up.*” (P3) Similarly, three participants saw its potential in classroom design projects, where students must justify their decisions with evidence. They noted that REFINEx could serve as an educational tool, fostering learning through hands-on engagement with insight-evidence pairs, ultimately enhancing both the quality of their work and their ability to articulate design rationale: “*When I need to make a digital product, and if I need to make a report, it would be really helpful, because every school assignment requires reference and asks me to think through design rationale about my design iterations.*” (P11)

5.4.2 Designer-centered retrieval facilitated the intentional discovery of scholarly insights. Of all, 9 responded that the most painful process during the baseline condition was searching for the papers that could potentially benefit their design iteration. More specifically, they responded that they had to try out multiple queries regarding the designs, as it was difficult to find the right query for retrieving papers, and even then, they could not consistently find inspiring papers. Supporting this, participants during the baseline condition spent the first 13m 8s on average ($SD = 4m 58s$) out of the 25-minute task time navigating scholarly resources before beginning their initial design iteration, frequently refining their

queries to gather relevant insights. This was significantly longer than in the REFINEx condition, where participants reached design iteration on the canvas more quickly, spending an average of 5m 23s ($SD = 2m 28s$) before starting their initial design iteration ($t = 5.88$, $p < .001$). During this time, while they viewed 5.8 papers ($SD = 4.8$) on average during the baseline condition, they only referred back to 1.4 papers ($SD = 1.0$) when designing. This repetitive and time-consuming process was extremely burdensome and the overall process less effective: “*The most frustrating thing (in the baseline condition) was that, I didn’t know the right keywords for finding the research paper. I felt like it was not working. I needed to try very hard to search for the relevant resources (...) which failed.*” (P5) This increased difficulty in retrieving papers likely contributed to the higher perceived workloads during the baseline condition.

On the other hand, participants noted that REFINEx’s design-centric indexing and retrieval enabled them to discover more relevant and applicable insights with less effort, in contrast to manually coming up with queries through a trial-and-error search process. First, every participant agreed that the design contexts extracted by REFINEx were accurately elicited for retrieval: “*The most challenging part (in the baseline condition) is that, I didn’t know how to use the right keyword. But I found they (design dimensions elicited by REFINEx) were really perfect ones.*” (P12) With its accuracy, REFINEx’s automatic retrieval was reported to facilitate more intentional exploration, while reducing frustration and cognitive load compared to conventional retrieval methods: “*Having a plugin that can automatically collect all the papers that are related to what you’re doing is really efficient.*” (P9)

Envisioning the use of REFINEx in their own projects, two participants expressed a desire to expand the design dimensions by incorporating their own design goals—ideas that are not yet visible in their mockups but exist in their minds or elsewhere in their workflow. Currently, REFINEx lets designers refine the design dimensions by iterating on the detected dimensions after the system elicits them. However, since these dimensions are derived based on what is already visible in the design, participants wished for more control over the outputs by integrating what remains invisible, such as undocumented thoughts or goals recorded elsewhere: “*(what if) this is the direction that I’m pushing for now and I’m just trying to*

find the evidence (...) possibly adding a way to provide my own goal to that plugin would be nice.” (P3)

5.4.3 Clustering helped them to avoid information overload while improving validity. Participants in our study found the design implications to be well-clustered. Among the seven participants who commented on clustering, six stated that the design implications were effectively grouped and accurately reflected the cluster content, and the other participant mentioned that a more thorough evaluation would be possible if they had access to the original papers: “*The clusters were accurate and related to the contents (implications) they had (...) the titles clusters had were really relevant to the sources.*” (P2)

Our survey indicated that participants considered the design insights communicated by REFINEx to be significantly more valid than those they identified in the baseline condition, while also requiring significantly lower temporal demand. Participants mentioned that multiple implications provided within the cluster reinforced the cluster’s core message, while enhancing its validity and eliminating the need for manually reading and verifying multiple papers: “*Cluster is very straightforward to me (...) The suggestion is very trusted, because, it’s all following its research paper resources and also very easy to read.*” (P5)

At the same time, our interview highlighted the potential need to allow designers to organize action items by each screen in their design workspace. Participants emphasized that screens serve as their primary focus, suggesting that action items derived from cluster contents should be grouped by screen first. This approach would allow designers to assess validity while maintaining alignment with their original workflow of interacting with screens: “*So, for example, click a page, and then I see the action items on that page, and then maybe it has guidelines like ‘add a search bar here because that will (...)’ to simplify the user flow.*” (P10)

5.4.4 Step-by-step translation helped them to broaden the scope of relevant insights without losing the originality of sources. Our survey suggests that, with REFINEx, participants found the communicated design insights from scholarly papers more relevant for their work. In Section 5.4.2, we discussed how that is partly due to participants being able to find more relevant papers (*i.e.*, supporting retrieval). However, at the same time, our interview results also suggest that the difference in relevance rating is partly due to the research translation that occurred. Of all, 10 noted that the translated insights from REFINEx broadened their understanding of what constitutes ‘applicable design insights’ from research papers. Participants reported a tendency during the baseline condition to fixate on finding a perfectly matching paper, which constrained their perception of applicable insights: “*I was kind of being oriented by the paper directions, not being led by my own thoughts. I don’t like that process, and I would never do that again.*” (P11) However, with REFINEx’s translation support, they were able to explore a wider range of relevant insights that they would have overlooked in traditional paper querying contexts: “*The use cases for me would be that it inspires me to consider aspects I might not think about otherwise if I weren’t using the tool.*” (P4) This reinforces previous findings that translating design insights from research papers expands the scope of relevant scholarly insights [45].

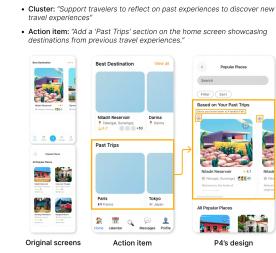


Figure 9: P4’s example of implementing an action item. After recognizing its potential usefulness, P4 creatively adapted the action item to a different screen, redesigning it with alternative UI components.

In this process, the step-by-step guidance provided by REFINEx played a crucial role in helping participants assess the validity of the translated insights without losing the originality of the sources. Beginning with a comparison and contrast, they could gain an overview of these implications before seeing how they could be applied to their specific contexts. This structured process enabled them to quickly and effectively evaluate the relevance and applicability of the translated insights: “*It told me that some of the papers were relevant, then they would explain that it’s because those papers also talk about travel (as a domain), and that allowed me to trust it. And then also with the plugin they had a summary. And they’re like, oh, this is how you can inspire your design. That was something that the papers I reviewed (in the baseline condition) didn’t have.*” (P10)

5.4.5 Visualizing action items improved designers’ ability to make informed design decisions quickly and enhanced learnability. Describing themselves as ‘visual learners,’ participants emphasized the significance of visual illustrations in understanding design insights. Before the study began, they expressed that the text-heavy nature of papers made it challenging to derive inspiration from written content, which diminished the papers’ utility—which aligns with previous findings in translational science for design that noted the difficulties associated with consuming text-heavy papers [13]. Similarly, participants during the study reported challenges in quickly grasping the design implications presented in the text of the papers: “*I feel like that’s just too much text (in the paper) and there are no highlights relating to the design suggestions. So I personally don’t want to read it.*” (P3)

Our survey showed that participants rated the actionability and validity of design insights significantly higher when using REFINEx compared to the baseline, with lower temporal demand and more design edits made within the same timeframe. Our interviews showed that visualizing action items allowed participants to quickly understand the required actions, streamlining their process of validating the insights and minimizing the time needed to integrate them into their designs. All participants reported the visual representations of the action items to be extremely useful, as they not only enabled the quick application of insights to their designs but also facilitated a better understanding of the messages conveyed. By rapidly grasping the insights through the visuals, participants could assess their usefulness based on their design experience and the relevance of surrounding content (*e.g.*, cluster contents, sources),

which streamlined the process while maintaining their validity: *“I think that (visualization) is the really intuitive and direct way of understanding of explanation on the rationale. I had a pretty good understanding when I just saw the visuals.”* (P11)

With a deeper understanding, we observed that the visualization also facilitated externalization, allowing participants to use it as a learning tool to translate insights into their own understanding. As shown in Figure 9, the enhanced comprehension of action items through visualization enabled participants to externalize these items and proactively design in ways they deemed improved. This demonstrated that visualization not only reinforced understanding but also empowered participants to take an active role in refining and applying design concepts.

In this process, every participant unanimously reported that the minor discrepancies in the reconstructed mockups did not affect their comprehension of the action items. They indicated that they were more focused on the overall ‘semantics’ conveyed by the visualization rather than minor visual details. Here, the presence of a toggle to turn the visual change on or off was reported to significantly assist in their prompt understanding of the modifications: *“It had all the same components, like same features, buttons, and all that. It was just different and kind of different in design. For me, I see all of these as components, and the design was up to me.”* (P9)

5.4.6 Enhanced scannability and optional presentation of contents to better help the navigation and validation of design insights. Although our quantitative results indicate that participants found the use of ReFINE to be significantly less effortful compared to the baseline, we identified opportunities to further reduce effort, such as incorporating a more catchy summary of clusters to streamline the process of discovering insights. For instance, upon encountering a cluster about the importance of improving financial literacy for mobile banking security, P9 suggested adding hashtags like ‘financial literacy’ to facilitate prompt understanding of what the cluster is talking about: *“I think maybe it could be better to have keywords, like ‘financial literacy’ (...) just showing those keywords would make it easier to navigate.”* (P9) Similarly, to further improve the scannability of the contents, participants highlighted the potential need for highlighting the subset of contents, such as boldfacing the important keywords: *“If only the core keywords can be highlighted, or bold, maybe it’s helpful for quickly scanning the clusters.”* (P4)

Additionally, participants acknowledged that supporting evidence (e.g., sources) contributed to originality and validity; yet, after reviewing these and finding the content trustworthy, they noted that sources were not their primary focus when reading in their design scenario. As a result, they felt these elements occupied excessive space and preferred to hide them by default, opting to review them only when insights contradicted their expectations and required validation. Consequently, they suggested making these sections hidden by default and expandable upon request: *“So maybe we can hide that, and then you can open it only when you are interested in learning more about the sources.”* (P12)

6 DISCUSSION & FUTURE WORK

In this paper, we introduced ReFINE, a novel system powered by LLM/VLMs that streamlines the consumption of design implications from research papers during UI design iterations in Figma.

ReFINE tackles key challenges across retrieval, translation, and contextualization—common barriers to applying design insights, as identified in prior research in translational science for design. Through our technical evaluations and user study, we demonstrated the reliability of ReFINE’s components and explored how designers perceive its impact during real design iterations.

One key implication of our study is that supporting the visualization of scholarly insights into their designs not only facilitates their quick integration into workflows but also significantly enhances the understanding and learnability of these insights. Designers, who are heavily trained to engage with visuals, often find text-heavy content overwhelming [13, 29, 41], as supported by our user study revealing their challenges in consuming full-text papers and lengthy summaries. With design implications presented in a visual format, they could quickly grasp the core messages of each implication cluster and more easily incorporate them into their design iteration. Still, our system currently limits visualization to static design components that can be integrated without adding extra screens or interactivity. Having demonstrated the efficacy of our approach, future research could explore ways to support dynamic elements or those requiring additional screens, expanding the applicability of our approach.

Another finding from our study is that, when designers were asked to search for papers directly in the baseline condition, most designers relied on narrow, domain-specific keyword searches in the paper repository, limiting their exposure to relevant papers outside those exact terms. Our system operates within the Figma canvas to help match users with papers along dimensions critical to design relevance and uses that structured understanding, augmented by translation, to retrieve papers that may not match the original keywords but still offer valuable insights. This helped designers move beyond keyword constraints and discover findings more aligned with their design needs. In this process, the informational content within each insight cluster served as an explainability layer, guiding designers on how to adapt insights from imperfectly matching papers. These results highlight the potential for a more expansive notion of relevance when retrieving from paper repositories. At the same time, our findings identify a need for more concise cluster summaries; future work could explore presentation methods for highlighting key takeaways to enhance navigation and scannability.

While our system demonstrated its effectiveness in supporting design iteration using a limited set of papers (i.e., CHI ‘24 proceedings), we see potential in increasing the scope—both extending (i.e., covering multiple-year proceedings) and expanding (i.e., diversifying the type of proceedings) the sources of literature. First, extending the literature to multiple proceedings could increase the granularity of topical coverage, improving the chances of finding more relevant literature. Second, unlike CHI—which serves as an umbrella venue—the focus of more domain- and topically-specialized HCI conferences (e.g., UIST, DIS, ASSETS, MobileHCI) could provide designers with a filtering mechanism to retrieve papers more specialized to their designs. As such, incorporating the unique focus of each venue as an additional dimension could enhance ReFINE’s ability to recommend design insights tailored to specific designer needs, yet it would be necessary to understand if the papers retrieved from these added conferences have sufficient

diversity and relevance. Beyond HCI, future work could explore incorporating publications from other applied disciplines (e.g., psychology, communication), which often produce research relevant to design practices. REFINEx will likely need to be modified to support the translation of papers from these scholarly communities with different practices for communicating practical implications.

In this paper, we presented an evaluation of REFINEx in a mobile UI design setting, building on prior works that focused on mobile UIs (e.g., [19, 49]). However, our system should be readily adaptable to other types of interfaces. For example, Figma is widely used by designers to create UIs on other modalities (e.g., web¹³, smart watch¹⁴), and REFINEx can dynamically detect the design modality through the elicitation of design context, enabling the retrieval of insights tailored to the corresponding modality. Beyond Figma, REFINEx can also be integrated with commercial web editors that support WYSIWYG and HTML exports, enabling faster and more accurate mockup visualizations by directly utilizing their native HTML exports, thereby eliminating the need for mockup reconstruction.

Lastly, participants envisioned various potential use cases within their real-world design workflows, such as applying REFINEx in corporate settings or class projects, while at the same time revealing the need to extend the design contexts extracted by REFINEx beyond visible elements to include aspects embedded in their thoughts or workflows. To better support these scenarios, we propose that future translational systems for UI design iterations could leverage documents from existing workflows and extract implicit design goals, as these contexts are often guided by well-documented guidelines or objectives (e.g., marketing briefs, syllabi) that may not be immediately apparent in a design mockup. By incorporating these additional dimensions for characterizing a design context, the system could further enrich results by aligning them with organizational goals.

7 CONCLUSION

In this work, we introduce REFINEx, a system designed to seamlessly integrate scholarly design implications into the UI mockup iteration process. REFINEx automates the process of retrieving, translating, and contextualizing research insights from HCI papers, enabling designers to incorporate evidence-informed knowledge directly within their existing workspace (i.e., Figma). Our comprehensive evaluations, including technical evaluations and a user study with designers, demonstrate the reliability of REFINEx components and the system's ability to reduce the burden on designers, improve the communication of scholarly insights, and streamline design iteration. These results highlight REFINEx's potential to bridge the gap between academic research and practical design, providing insights for developing tools that support evidence-informed design workflows.

REFERENCES

- [1] 2024. *CHI '24: Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA). Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3613904>
- [2] CHI 2024. *Papers*. <https://chi2024.acm.org/for-authors/papers/>
- [3] Robin S Adams. 2002. Understanding Design Iteration: Representations from an Empirical Study. In *Common Ground - DRS International Conference 200*. <https://dl.designresearchsociety.org/drs-conference-papers/drs2002/researchpapers/2>
- [4] Armen Aghajanyan, Dmytro Okhonko, Mike Lewis, Mandar Joshi, Hu Xu, Gargi Ghosh, and Luke Zettlemoyer. 2021. HTLM: Hyper-Text Pre-Training and Prompting of Language Models. *arXiv preprint arXiv:2107.06955* (2021). <https://arxiv.org/abs/2107.06955>
- [5] Anthropic. 2024. *Claude 3.5 Sonnet*. <https://www.anthropic.com/news/clause-3-5-sonnet>
- [6] arXiv. 2018. *New Release: arXiv Search v0.1*. <https://blog.arxiv.org/2018/04/17/new-release-arxiv-search-v0-1/>
- [7] Tal August, Lucy Lu Wang, Jonathan Bragg, Marti A Hearst, Andrew Head, and Kyle Lo. 2023. Paper Plain: Making Medical Research Papers Approachable to Healthcare Consumers with Natural Language Processing. *ACM Transactions on Computer-Human Interaction* 30, 5 (2023), 1–38. <https://doi.org/10.1145/3589955>
- [8] Peter Bailis, Simon Peter, and Justine Sherry. 2016. Introducing Research for Practice. *Commun. ACM* 59, 9 (2016), 38–41. <https://doi.org/10.1145/2909474>
- [9] Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative Research in Psychology* 3, 2 (2006), 77–101. <https://doi.org/10.1191/1478088706QP063OA>
- [10] Sara Bunian, Kai Li, Chaima Jemmali, Casper Harteveld, Yun Fu, and Magy Seif El-Nasr. 2021. VINS: Visual Search for Mobile User Interface Design. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–14. <https://doi.org/10.1145/3411764.3445762>
- [11] Andrea Burns, Krishna Srinivasan, Joshua Ainslie, Geoff Brown, Bryan A Plummer, Kate Saenker, Jianmo Ni, and Mandy Guo. 2023. A Suite of Generative Tasks for Multi-Level Multimodal Webpage Understanding. *arXiv preprint arXiv:2305.03668* (2023). <https://arxiv.org/abs/2305.03668>
- [12] Joel Chan, Katherine Fu, Christian Schunn, Jonathan Cagan, Kristin Wood, and Kenneth Kotovsky. 2011. On the Benefits and Pitfalls of Analogies for Innovative Design: Ideation Performance Based on Analogical Distance, Commonness, and Modality of Examples. *Journal of Mechanical Design* 133, 8 (08 2011), 081004. <https://doi.org/10.1115/1.4004396>
- [13] Lucas Colusso, Cynthia L Bennett, Gary Hsieh, and Sean A Munson. 2017. Translational Resources: Reducing the Gap Between Academic Research and HCI Practice. In *Proceedings of the 2017 Conference on Designing Interactive Systems*. 957–968. <https://doi.org/10.1145/3064663.3064667>
- [14] Lucas Colusso, Ridley Jones, Sean A Munson, and Gary Hsieh. 2019. A translational science model for HCI. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–13. <https://doi.org/10.1145/3290605.3300231>
- [15] 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. <https://doi.org/10.1145/3126594.3126651>
- [16] Biplab Deka, Zifeng Huang, and Ranjitha Kumar. 2016. ERICA: Interaction Mining Mobile Apps. In *Proceedings of the 29th Annual Symposium on User Interface Software and Technology*. 767–776. <https://doi.org/10.1145/2984511.2984581>
- [17] Xiang Deng, Prashant Shiralkar, Colin Lockard, Binxuan Huang, and Huan Sun. 2022. DOM-LM: Learning Generalizable Representations for HTML Documents. *arXiv preprint arXiv:2201.10608* (2022). <https://arxiv.org/abs/2201.10608>
- [18] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2020. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *arXiv preprint arXiv:2010.11929* (2020). <https://arxiv.org/abs/2010.11929>
- [19] Peitong Duan, Jeremy Warner, Yang Li, and Bjoern Hartmann. 2024. Generating Automatic Feedback on UI Mockups with Large Language Models. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–20. <https://doi.org/10.1145/3613904.3642782>
- [20] Jennifer Ferreira, James Noble, and Robert Biddle. 2007. Agile Development Iterations and UI Design. In *Agile 2007*. IEEE, 50–58. <https://doi.org/10.1109/AGILE.2007.8>
- [21] Raymond Fok, Joseph Chee Chang, Tal August, Amy X Zhang, and Daniel S Weld. 2024. Qlarify: Recursively Expandable Abstracts for Dynamic Information Retrieval over Scientific Papers. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*. 1–21. <https://doi.org/10.1145/3654777.3676397>
- [22] Dedre Gentner and Arthur B Markman. 1997. Structure Mapping in Analogy and Similarity. *American Psychologist* 52, 1 (1997), 45. <https://doi.org/10.1037/0003-066X.52.1.45>
- [23] Karni Gilon, Joel Chan, Felicia Y Ng, Hila Liifshitz-Assaf, Aniket Kittur, and Dafna Shahaf. 2018. Analogy Mining for Specific Design Needs. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. 1–11. <https://doi.org/10.1145/3173574.3173695>
- [24] Google. 2024. *Embeddings*. <https://ai.google.dev/gemini-api/docs/embeddings>

¹³<https://www.figma.com/community/design-tutorials/web-design>

¹⁴https://www.figma.com/community/search?resource_type=mixed&query=watch

- [25] Google. 2024. *Gemini 2.0 Flash*. <https://ai.google.dev/gemini-api/docs/models/gemini>
- [26] Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Saifdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. 2023. A Real-World WebAgent with Planning, Long Context Understanding, and Program Synthesis. *arXiv preprint arXiv:2307.12856* (2023). <https://arxiv.org/abs/2307.12856>
- [27] Sandra G Hart and Lowell E Staveland. 1988. Development of NASA-TLX (Task Load Index): Results of Empirical and Theoretical Research. In *Advances in Psychology*. Vol. 52. Elsevier, 139–183. [https://doi.org/10.1016/S0166-4115\(08\)62386-9](https://doi.org/10.1016/S0166-4115(08)62386-9)
- [28] Scarlett R Herring, Chia-Chen Chang, Jess Krantzler, and Brian P Bailey. 2009. Getting Inspired! Understanding How and Why Examples are Used in Creative Design Practice. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 87–96. <https://doi.org/10.1145/1518701.1518717>
- [29] Bryan Howell, Asa Jackson, Henry Lee, Julianne DeVita, and Rebekah Rawlings. 2021. Exploring the Experiential Reading Differences between Visual and Written Research Papers. In *Learn X Design 2021: Engaging with Challenges in Design Education*. https://doi.org/10.21606/drs_lxd2021.03.247
- [30] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2025. A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *ACM Transactions on Information Systems* 43, 2 (2025), 1–55. <https://doi.org/10.1145/3703155>
- [31] Yoonjoo Lee, Heyeonsu B Kang, Matt Latzke, Juho Kim, Jonathan Bragg, Joseph Chee Chang, and Pao Siangliu. 2024. PaperWeaver: Enriching Topical Paper Alerts by Contextualizing Recommended Papers with User-collected Papers. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–19. <https://doi.org/10.1145/3613904.3642196>
- [32] Luis A Leiva, Asutosh Hota, and Antti Oulasvirta. 2020. Enrico: A Dataset for Topic Modeling of Mobile UI Designs. In *22nd International Conference on Human-Computer Interaction with Mobile Devices and Services*. 1–4. <https://doi.org/10.1145/3406324.3410710>
- [33] Junlong Li, Yiheng Xu, Lei Cui, and Furu Wei. 2021. MarkupLM: Pre-training of Text and Markup Language for Visually-rich Document Understanding. *arXiv preprint arXiv:2110.08518* (2021). <https://arxiv.org/abs/2110.08518>
- [34] Patrice Lopez. 2009. GROBID: Combining Automatic Bibliographic Data Recognition and Term Extraction for Scholarship Publications. In *Research and Advanced Technology for Digital Libraries*. Springer Berlin Heidelberg, 473–474. https://doi.org/10.1007/978-3-642-04346-8_62
- [35] Yuwen Lu, Alan Leung, Amanda Szwarcin, Jeffrey Nichols, and Titus Barik. 2025. Misty: UI Prototyping Through Interactive Conceptual Blending. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–17. <https://doi.org/10.1145/3706598.3713924>
- [36] Jenny Ma, Karthik Sreedhar, Vivian Liu, Sitong Wang, Pedro Alejandro Perez, and Lydia B Chilton. 2024. DIDUP: Dynamic Iterative Development for UI Prototyping. *arXiv preprint arXiv:2407.08474* (2024). <https://arxiv.org/abs/2407.08474>
- [37] Jakob Nielsen. 2002. Iterative User-Interface Design. *Computer* 26, 11 (2002), 32–41. <https://doi.org/10.1109/2.241424>
- [38] Donald A Norman. 2010. The Research-Practice Gap: The Need for Translational Developers. *Interactions* 17, 4 (2010), 9–12. <https://doi.org/10.1145/1806491.1806494>
- [39] Novia Nurain, Chia-Fang Chung, Clara Caldeira, and Kay Connelly. 2024. Designing a Card-Based Design Tool to Bridge Academic Research & Design Practice For Societal Resilience. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–20.
- [40] OpenAI. 2024. GPT-4o. <https://openai.com/index/hello-gpt-4o/>
- [41] Hyerim Park, Joscha Eirich, Andre Luckow, and Michael Sedlmair. 2024. "We Are Visual Thinkers, Not Verbal Thinkers!": A Thematic Analysis of How Professional Designers Use Generative AI Image Generation Tools. In *Proceedings of the 13th Nordic Conference on Human-Computer Interaction*. 1–14. <https://doi.org/10.1145/3679318.3685370>
- [42] Seokhyeon Park, Yumin Song, Soohyun Lee, Jaeyoung Kim, and Jinwook Seo. 2025. Leveraging Multimodal LLM for Inspirational User Interface Search. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–22. <https://doi.org/10.1145/3706598.3714213>
- [43] Peter J Rousseeuw. 1987. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *J. Comput. Appl. Math.* 20 (1987), 53–65. [https://doi.org/10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7)
- [44] Corina Sas, Steve Whittaker, Steven Dow, Jodi Forlizzi, and John Zimmerman. 2014. Generating Implications for Design through Design Research. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 1971–1980. <https://doi.org/10.1145/2556288.2557357>
- [45] Donghoon Shin, Tze-Yu Chen, Gary Hsieh, and Lucy Lu Wang. 2025. What About My Design Context?: Exploring the Use of Generative AI to Support Customization of Translational Research Artifacts. In *Proceedings of the 2025 ACM Designing Interactive Systems Conference*. 1210–1227. <https://doi.org/10.1145/3715336.3735686>
- [46] Donghoon Shin, Lucy Lu Wang, and Gary Hsieh. 2024. From Paper to Card: Transforming Design Implications with Generative AI. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–15. <https://doi.org/10.1145/3613904.3642266>
- [47] Chenglei Si, Yanzhe Zhang, Ryan Li, Zhengyuan Yang, Ruibo Liu, and Diyi Yang. 2024. Design2Code: Benchmarking Multimodal Code Generation for Automated Front-End Engineering. *arXiv preprint arXiv:2403.03163* (2024). <https://arxiv.org/abs/2403.03163>
- [48] Sarah Sulteri, Nilda Kipi, Linh Chi Tran, and Matthias Jarke. 2019. UI Design Pattern-driven Rapid Prototyping for Agile Development of Mobile Applications. In *Proceedings of the 21st International Conference on Human-Computer Interaction with Mobile Devices and Services*. 1–6. <https://doi.org/10.1145/3338286.3344399>
- [49] Bryan Wang, Gang Li, and Yang Li. 2023. Enabling Conversational Interaction with Mobile UI using Large Language Models. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–17. <https://doi.org/10.1145/3544548.3580895>
- [50] Taylor Webb, Keith J Holyoak, and Hongjing Lu. 2023. Emergent Analogical Reasoning in Large Language Models. *Nature Human Behaviour* 7, 9 (2023), 1526–1541. <https://doi.org/10.1038/s41562-023-01659-w>
- [51] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, F. Xia, Quoc Le, and Denny Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. *Advances in Neural Information Processing Systems* 35 (2022), 24824–24837.
- [52] Ziming Wu, Qianyao Xu, Yiding Liu, Zhenhui Peng, Yingqing Xu, and Xiaojuan Ma. 2021. Exploring Designers' Practice of Online Example Management for Supporting Mobile UI Design. In *Proceedings of the 23rd International Conference on Mobile Human-Computer Interaction*. 1–12. <https://doi.org/10.1145/3447526.3472048>
- [53] Junchi Yu, Ran He, and Zhitao Ying. 2024. Thought Propagation: An Analogical Approach to Complex Reasoning with Large Language Models. In *The Twelfth International Conference on Learning Representations*. <https://doi.org/10.48550/arXiv.2310.03965>
- [54] Lixiu Yu, Aniket Kittur, and Robert E Kraut. 2014. Searching for Analogical Ideas with Crowds. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 1225–1234. <https://doi.org/10.1145/2556288.2557378>
- [55] Ruican Zhong, Donghoon Shin, Rosemary Meza, Predrag Klasnja, Lucas Colusso, and Gary Hsieh. 2024. AI-Assisted Causal Pathway Diagram for Human-Centered Design. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–19. <https://doi.org/10.1145/3613904.3642179>

A STUDY DETAILS

A.1 Technical Evaluation

Table 5: Examples of relevance ratings assigned by annotators when evaluating the connection between design implications and action items

| Relevance rating | Design implication | Action item | Annotator's explanation |
|----------------------------|---|---|--|
| 5 - Highly relevant | Reduce user burden by minimizing text input in form-heavy tasks | Implement a one-click autofill feature that uses stored preferences for all checkout fields | Directly addresses and solves the burden |
| 4 - Substantially relevant | Make it easier for users to navigate dense information in reports | Add a 'Back to Top' button on the settings screen | Helps with navigation, but not specific to navigating documents |
| 3 - Moderately relevant | Striking the balance between insufficient and overwhelming transparency will enable users to trust the system | Add a 'Report' button to the comment section to enhance user feedback and platform moderation | May support trust, but unclear link to transparency or balance |
| 2 - Slightly relevant | Allow users to control the visibility of their personal data in shared platforms | Let users toggle switches in the profile settings to show/hide recent activity | Some connection to visibility, but lacks clarity on shared context |
| 1 - Not relevant | Help users build trust in financial transactions through transparency | Display an autoplay video on the app with product promotions | Unrelated to trust or transparency |

A.2 User Study

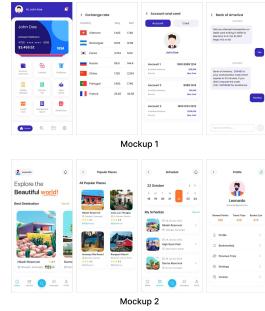


Figure 10: Prototype mockups used in our user study