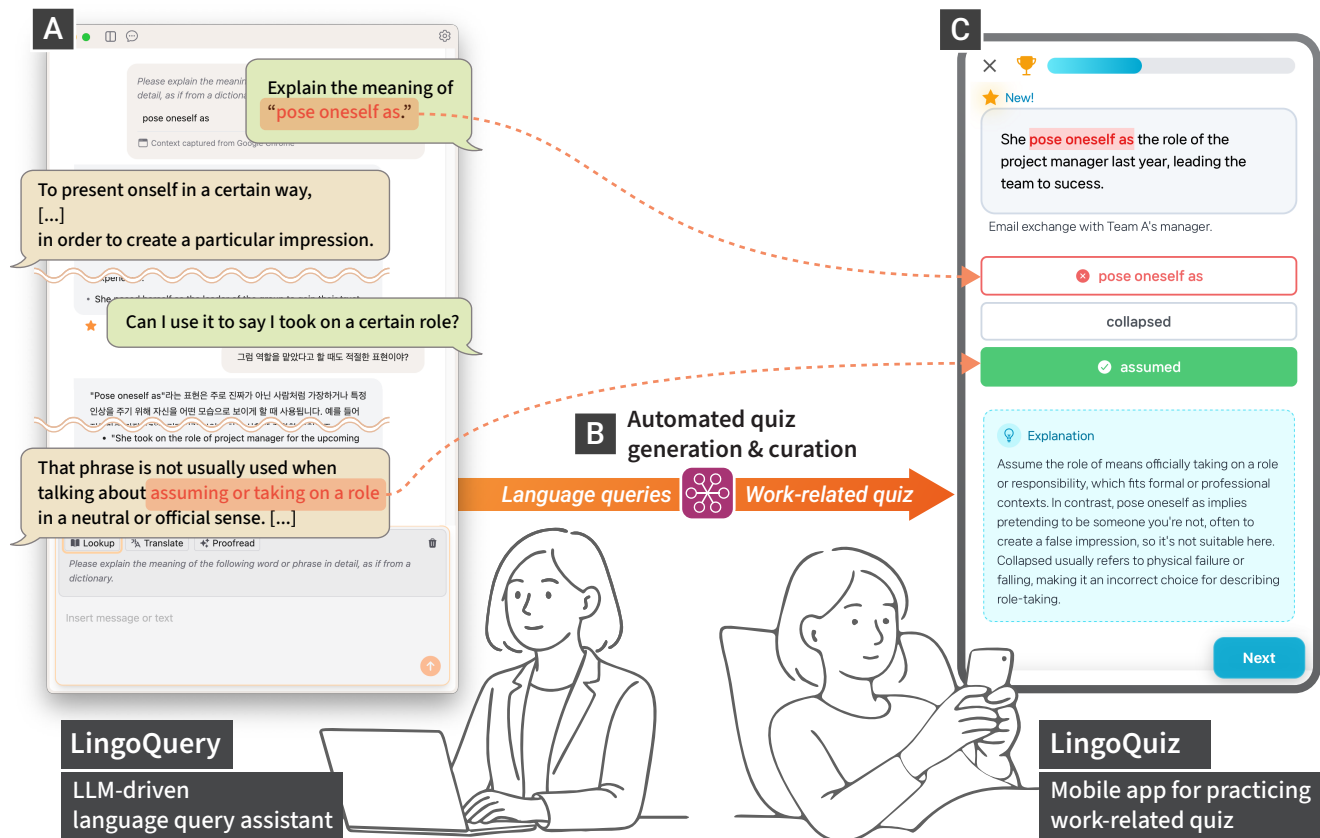


# LINGOQ: Bridging the Gap between ESL Learning and Work through AI-Generated Work-Related Quizzes

Yeonsun Yang\* DGIST Republic of Korea diddustjs98@dgist.ac.kr  
 Sang Won Lee† Virginia Tech Blacksburg, VA, USA sangwonlee@vt.edu  
 Jean Y. Song Yonsei University Republic of Korea jeansong@yonsei.ac.kr  
 Sangdoon Yun NAVER AI Lab Republic of Korea oodgnas@gmail.com  
 Young-Ho Kim NAVER AI Lab Republic of Korea yghokim@younghokim.net



**Figure 1: LINGOQ consists of three components. In LINGOQUERY for desktop **A**, information workers can interact with an LLM-based chatbot for English-related language queries. The automated quiz generation pipeline **B** produces and curates multiple choice English questions using the query interactions of LINGOQUERY as materials. In LINGOQUIZ **C** on a smartphone, workers can later review their language queries by completing the generated quizzes.**

## Abstract

Non-native English speakers performing English-related tasks at work struggle to sustain ESL learning, despite their motivation. Often, study materials are disconnected from their work context. Although workers rely on LLM assistants to address their immediate needs, these interactions may not directly contribute to their English skills. We present LINGOQ, an AI-mediated system that allows workers to practice English using quizzes generated from their LLM queries during work. LINGOQ leverages these queries using AI to generate personalized quizzes that workers can review

and practice on their smartphones. We conducted a three-week deployment study with 28 ESL workers to evaluate LINGOQ. Participants valued the relevance of quizzes that reflect their own context, constantly engaging with the app during the study. This active engagement improved self-efficacy and led to learning gains for beginners and, potentially, for intermediate learners. We discuss opportunities of leveraging users' reliance on LLMs to situate their learning in the user context for improved learning.

## CCS Concepts

• Human-centered computing → Natural language interfaces; Empirical studies in HCI.

\*Yeonsun Yang conducted this work as a research intern at NAVER AI Lab.

†Sang Won Lee conducted this work as a visiting scholar at NAVER AI Lab.

## Keywords

English as a second language, information workers, large language model, question generation, context awareness

## 1 Introduction

In the global economy, where online resources are predominantly in English, information workers who are not native English speakers often need English proficiency for their jobs. This requirement includes understanding English texts used in papers, articles, and reports, as well as the ability to communicate via email. To enhance their English skills, information workers<sup>1</sup> frequently engage in self-directed language learning using mobile applications like Duolingo [29], Babbel [8], Memrise [66], and RosettaStone [82] that offers access to learning anytime, anywhere through mobile phones [91]. These mobile apps offer access to learning anytime, anywhere through mobile phones [91], and make language learning engaging using gamification [85]. However, the mobile language learning apps often create study materials from generic situations like traveling or business meetings, limiting the depth and practical relevance of vocabulary and conversational skills. Such questions, based on ordinary scenarios, make it difficult to acquire the English skills needed for work-related tasks directly. For example, a programmer requiring English for API documentation gains little from fill-in-the-blank questions about family composition or reading comprehension exercises on threats to coral reefs in the Pacific Ocean.

To address the challenge of English studies being disconnected from workers' everyday tasks, we draw on task-based language teaching [70] and situated learning [59]. When study materials are embedded in workers' job contexts, they can improve task performance while sustaining engagement in language learning [4]. Prior research shows that grounding instruction in authentic tasks not only supports second language acquisition [1, 93], but also enhances motivation [46] and strengthens memory retention [18]. Building on these insights, we investigate whether generating study materials of mobile language learning apps directly from work tasks can further enhance language learning and foster long-term engagement.

To that end, we developed LINGOQ, an ensemble of intelligent systems designed to generate English quizzes from workers' LLM queries. LINGOQ consists of three components: (1) LINGOQUERY, an LLM-powered chatbot that answers workers' English-related queries, running on their computers similarly to ChatGPT (A in Figure 1); (2) LINGOQUIZ, a mobile application allowing workers to complete short, context-relevant quizzes at their convenience (C in Figure 1); and (3) the backend pipeline that processes queries from LINGOQUERY to generate and validate quizzes, considering the context captured from a worker's computer screen (B in Figure 1). Our system design was informed by findings from a formative study, which incorporated both an online survey and follow-up interviews, involving 49 non-native information workers in South Korea. This study identified several challenges that users face in sustaining their English studies, such as a lack of relevance to their work context

and difficulties in maintaining continuous engagement. In addition, we identified the workers' patterns of how they seek assistance at work: their heavy usage of LLM-based assistants and their frequent query types, such as look up, translate, and proofread. In designing LINGOQUERY, we incorporated these patterns to enhance query interaction and streamline the quiz generation pipeline, specifically focusing on areas where workers needed support.

LINGOQ leverages two commonly observed but separate modes of interaction: (1) using AI tools at work for English-related tasks and (2) learning English through smartphones. Our work integrates these activities into a single connected pipeline, making the established activities in English learning more directly relevant to participants' work and the domains in which they use English. While the mode of interaction in LINGOQ remains unchanged, we designed the system to make the learning content more relevant to workers' tasks and help workers benefit from personal big data collected from using LLMs at work.

To understand how LINGOQ effectively supported their English learning, we conducted a three-week field deployment with 28 information workers in South Korea. Our results showed that their self-efficacy in English skills increased significantly after actively using LINGOQ, and their engagement was consistent throughout the study period. Furthermore, we confirmed a notable learning gain from the beginner-level group. Overall, workers reported that quizzes generated by LINGOQ were more relevant to their work and were helpful in their tasks compared to their prior experience of studying English. It also encouraged them to consider the educational aspect of their work-related queries. Lastly, they perceived that studying with LINGOQ could be more sustainable than their previous experiences with other study methods.

The key contributions of this work are as follows:

- (1) Understanding the challenges non-native information workers face in maintaining consistent English study and how they leverage intelligent tools at work.
- (2) LLM-infused system architecture design pattern that links the use of LLM-based assistants to the generation of foreign language learning practices relevant to users' work contexts.
- (3) Empirical evidence that validates the effectiveness of generating context-relevant study materials in foreign language learning.

## 2 Related Work

### 2.1 ESL Learning for Information Workers

English plays a critical role in the workplace as the common language of global industries and disciplines [55, 90]. The proliferation of digital work environments has posed unique challenges for non-native English-speaking information workers, who navigate diverse English-language texts as part of their computer-based tasks. For example, Amano *et al.* found that researchers spend 46.6% more time reading English papers and 50.6% more time writing them compared to native speakers, while facing higher rejection rates due to writing quality [2]. Similarly, programmers using English as a second language struggle with technical documentation, professional communication, and code comprehension [42]. Consequently, English-as-a-second-language (ESL) learning has become critical

<sup>1</sup>In this work, we use the term *information worker* to refer to workers whose primary job role involves gathering, synthesizing, and producing new information [56]. In this paper, we will use the term 'workers' to specifically refer to information workers.

for information workers to develop practical English proficiency—the ability to achieve goals through the use of English in relation to the specific purposes or occasions [46].

In the field of HCI, previous research has explored approaches to directly support ESL workers’ use of English in the workplace, such as assisting with email composition, providing on-demand evaluations of workers’ own English sentences, simplifying complex texts, and paraphrasing with AI explanations, in order to reduce context-switching burdens and other disruptions during work [12, 14, 20, 45, 50, 54, 57]. However, effective support for ESL workers in the long run requires fostering practical English skills that directly enhance their ability to accomplish work tasks. In response, our work focuses on cultivating English proficiency within professional contexts.

Conventional English education, with its focus on general language competency, often fails to meet the specific demands of professional contexts. To address this gap, various education theories such as task-based language teaching [70] and situated learning [59] and English for Specific Purposes underscore the critical role of *authentic* materials—texts and resources originally produced for real-world application [37]. First, authentic materials provide learners with up-to-date information available in the field, in contrast to predefined materials that often fail to reflect recent developments [9]. Second, engaging with content directly related to one’s professional field fosters learner motivation and self-efficacy, as successfully comprehending this material builds confidence in handling real-world situations [10]. Lastly, using context-specific, task-based materials helps learners develop active knowledge, which is directly applicable to their daily work [94].

Building on this emphasis, our work investigates flexible ways to foster ESL learning by providing authentic, work-related materials generated from workers’ daily English tasks.

## 2.2 Creating Context-Aware Learning Materials

Context is fundamental to learning because knowledge is inseparable from the situations and activities in which it is acquired [16, 24]. Context-aware learning approaches [47], which situate learning in the learner’s personal context by selecting, adapting, or generating content, are particularly well-suited for language learning. Because contextually relevant materials allow learners to experience language in authentic settings [61] and foster situational interests [44], they thereby enhance motivation and engagement [44], prevent inert knowledge [64], and ultimately promote meaningful and effective learning [28].

The field of HCI has explored context-aware personalized learning materials that adapt to factors such as learners’ location [32, 43], surrounding elements [28, 48], social media content [15, 95], and other contextual information [72]. For example, MicroMandarin [32] suggested flashcards relevant to nearby venues, while Vocabura [43] generated L1-L2 word pairs from walking commute routes, both leveraging GPS coordinates to support vocabulary learning. Draxler *et al.* further explored an object-based approach that automatically generates exercises by detecting elements in learner-captured photos from their daily contexts [28]. In addition, Yamaoka *et al.* introduced a method that extracts keywords from Instagram posts to generate example sentences, helping learners

acquire new words aligned with their interests and improving retention [95]. Beyond academic research, recent ESL services have begun adopting generative AI to provide situated and contextualized learning experiences [30, 38].

With the growth of digital environments, research has increasingly focused on usage-based learning to help learners acquire practical language skills by leveraging learner–computer interactions as sources of contextual content, such as eye gaze [25], clicked hyperlinks [21, 84], translations [65], and other digital traces [7]. Ding *et al.* identified unknown words through gaze trajectories while learners read foreign language texts, offering real-time translations and explanations for just-in-time vocabulary acquisition [25]. Lungu *et al.* proposed a comprehension approach that generated mobile exercises from learners’ translated sentences during web reading, which could serve as potential learning cues [65]. This work demonstrated the feasibility of this ecosystem, highlighting engagement and learning benefits.

Our work extends this line of research by leveraging emerging conversational interactions between learners and LLM-based chatbots. We argue that queries to LLMs reflect learners’ immediate language difficulties and learning intentions, serving as valuable cues for situated and usage-based learning. In this work, we aim to generate work-related learning exercises from ESL workers’ queries collected in the course of their professional tasks.

## 2.3 Retrieval Practice for Second Language Acquisition

In second language acquisition, the *practice* of difficult linguistic features offers learners opportunities for meaningful language use, reinforces task performance, and fosters adaptive language proficiency [52, 88]. In particular, retrieval practice, which involves actively recalling knowledge from memory, typically through exercises such as quizzes or self-testing, is one of the most effective review strategies [83, 89]. It has been shown to be beneficial than simple ‘restudy’ in strengthening long-term memory and supporting the transfer of knowledge to new context [79, 83]. Accordingly, prior work has integrated retrieval-based exercises into second language learning systems to enhance vocabulary acquisition, reading comprehension, and communicative fluency [19, 27, 31].

When combined with microlearning [36]—an approach that delivers educational content in small, easily digestible units—retrieval practice becomes particularly effective for busy adults. Microlearning helps sustain motivation and engagement while minimizing the time burden, making it well-suited for learners balancing work and study [49, 53]. Many popular mobile-assisted language learning (MALL) systems, such as Duolingo [29], Anki [3], and Quizlet [78], leverage this principle by offering interactive exercises (e.g., flashcards, multiple-choice questions, and fill-in-the-blanks). In addition, the field of HCI has explored retrieval-based microlearning in various contexts, such as spaced practice for vocabulary acquisition, adaptive exercise scheduling, and bite-sized practice integrated into daily routines [27, 31].

Building on this, our work provides a mobile practice environment that leverages retrieval practice for information workers learning practical English skills. We focus on a particular exercise type—multiple choice fill-in-the-blank questions—which are widely used



in standardized proficiency tests [33]. Furthermore, recent research has demonstrated that AI-generated multiple-choice questions can reach expert-level quality [26, 34], highlighting their potential as a scalable and effective method for creating adaptive practice. This allows us to focus on the content’s quality, which we ensure through our systematic generation pipeline (see Section 4.3).

### 3 Formative Study

To inform the design of LINGOQ, we conducted an online survey with 49 information workers whose native language is Korean, followed by semi-structured interviews with ten volunteers. We aimed to understand the type of barriers they face during daily English-related tasks, limitations with the existing digital tools they use to handle these tasks, and effective ESL learning practices they have experience using to develop work-related English skills.

#### 3.1 Procedure and Analysis

**Online Surveys.** Through both closed and open questions, we asked participants about the challenges they face as non-native English speakers at work, the digital tools they use to support English-related tasks, the effectiveness of their ESL learning strategies, and their perceived need for continued learning. We also asked about the willingness to participate in a follow-up interview. The online survey was advertised to native Korean speakers on social media and our internal network, inviting information workers who use computers for their work and regularly perform tasks that require English. Forty-nine people (25 females; aged 22–49) completed the survey, which included 22 researchers, 11 engineers, and 16 professionals from various fields, including strategic planning, sales and marketing, design, general affairs, and healthcare. The survey took approximately 20 minutes to complete. We compensated 5,000 KRW (approx. 4 USD) for survey respondents.

**Interviews.** For in-depth analysis, we conducted follow-up interviews with ten survey respondents who indicated their willingness to attend as part of the online survey. Each interview lasted about 40 minutes and was conducted in person or remotely, depending on the participants’ availability. We revisited the interviewees’ survey responses and asked them to elaborate on their open-ended answers. Using screen sharing and think-aloud protocols [17], participants walked through recent scenarios involving English-related tasks, demonstrating queries they had made to generative AI (e.g., ChatGPT, Gemini) or other tools. They also described their review practices focused on work-specific English content. We compensated 20,000 KRW (approx. 14 USD) for interview participants.

**Analysis.** All interviews were audio-recorded and transcribed for analysis. We summarized the closed-ended survey questions using descriptive statistics. We used Thematic Analysis [11] to qualitatively analyze both the open-ended questions of the survey and interview transcripts. One researcher coded survey responses as well as interview transcripts simultaneously, grouping them into broader themes. The research team iterated through several rounds of discussion to refine these themes. In the following sections, we present findings from both the survey and the interviews, referring to each interview participant as I1 through I10.

#### 3.2 Finding 1: Understanding the Difficulties of ESL Workers in the Workplace

Participants worked with large amounts of information written in English as part of their daily tasks, ranging from (1) communication through emails or messengers, (2) accessing online resources, and (3) writing professional documents such as reports or papers.

One common linguistic difficulty that the majority of respondents (25/49) pointed out was **lexical disruption**, noting that unfamiliar domain-specific terminology often hindered their comprehension and prompted them to look up words frequently. I5, who works in governance administration at an international research lab, remarked that “*The official materials from the UN Headquarters are often overly formal and full of UN-specific terms, which slows me down as I have to look them up.*” Similarly, I1, an international business development manager, noted “*For example, I used to think ‘airway’ only meant a flight route, but later learned it also refers to a respiratory tract [a human body part].*” Relatedly, participants reported not only linguistic challenges but also affective challenges—stemming from a lack of confidence, which consequently hindered their workflow. Four participants mentioned in their surveys and interviews that they proofread emails for grammar, formality, and tone before sending, concerned that mistakes might appear impolite or give a negative impression of their professional competence.

Participants often felt that current ESL learning was **disconnected from work context**. More than 2/3 of the survey respondents (33/49) rated that they often (45.0%) or always (22.5%) feel the need for learning English for work on a 5-point Likert-type scale question. Yet, the majority of them (28/49) were not currently studying English, demonstrating the difficulty of constant engagement in ESL learning practices during work. For those who were currently studying English, all of them (21/21) reported that they used easily accessible and self-directed mobile apps, outside of work context (e.g., Duolingo [29], Speak app [86]). Other common practices included online tutoring (e.g., Ringle [80]; 8/21), reading English novels or articles (8/21), shadowing (4/21), and in-person courses (3/21). However, the majority of them (14/21) struggled to sustain their learning because irrelevant learning materials did not translate into practical support in their work contexts. In the follow-up interview, I10 noted, “*What I really want to learn right now is material I can use immediately in business meetings, but finding a suitable platform or tool has been very difficult.*” I1 also remarked, “*I’m often exposed to highly specialized medical terms, but when the material is from a learning app or an article outside my field, it tends to use more general vocabulary and expressions. While this is helpful for conversations, it’s not very useful when reading work-related articles or clinical papers.*”

Challenges occurred during the reviewing phase as well due to **lack of sustainable review routines**. To align ESL learning with their work contexts, six interviewees once tried to review unfamiliar words and expressions from work by compiling personal glossaries and organizing them with tools such as Notion, Google Docs, or the open-source flashcard app Anki [3]. However, participants failed to maintain engagement with such review routines, as manually collecting work-related vocabulary or expressions was time-consuming and burdensome. Moreover, reviewing these materials with explicit exercises further discouraged continued practice.

In the follow-up interviews, I9 noted, “After work, I don’t want to revisit the traces of what I did during the day. Reviewing would mean opening my daily logs in a workspace like Notion, finding the target words, gathering them on another page, and then asking GPT or searching Google for their meanings. Most days, it just feels too much.” Also, I10 remarked, “In one-on-one business English tutoring, my teacher listed my mistakes in Google Docs, but reviewing them felt like just reading meeting minutes and was neither fun nor motivating. I wish I could review them in more engaging ways, like quizzes or other formats for sustainable practice.” While interactive apps such as Anki, with flashcard and quiz features, were available, participants (I8, I9) found it overly complex and overwhelming to customize and manually upload word lists, especially after work.

### 3.3 Finding 2: Common Patterns of English Language Queries

To address language barriers, all participants used language assistance tools for lookup and double-checking, including dictionaries, web search engines, translators, AI-based writing assistants (e.g., Grammarly [40], DeepL [22]), and LLM-based chatbots (e.g., ChatGPT [73], Gemini [39], Claude [5]). In particular, most survey respondents (46/49) commonly used LLM-based chatbots, which offered convenient conversational support and context-aware explanations for a wide range of English-related difficulties. Interviewees entered queries primarily by copying and pasting text, ranging from single words to full passages, along with a prompt for linguistic support. From their usage scenarios, we identified three prominent query patterns to LLM-based chatbots: **look-up**, **translation**, and **proofreading**.

Most interviewees (7/10) often used chatbots just like dictionaries to **look up** definitions of unfamiliar words or description of confusing grammar during their tasks. I5 noted, “These days I just ask ChatGPT when I’m unsure about grammar, like ‘an MBA or a MBA?’” I6 and I8 found LLMs useful for clarifying domain-specific terms or subtle nuances, as they offered context-specific explanations, especially for words with multiple possible meanings.

All interviewees (10/10) used chatbots to **translate** text in English to Korean and vice versa. They translated text in English into Korean to ensure their comprehension and translated Korean into English to compose professional writing more efficiently. I6 noted, “I usually ask LLMs to align the original and the translation side by side, so I can double-check whether each part conveys the intended meaning,” highlighting the need for a dual-language view to enable rapid comparison under time pressure at work.

The majority of interviewees (6/10) also frequently asked chatbots to **proofread** their own draft, ranging from formal business documents to casual conversation with colleagues. Participants refined grammar, tone, and style to fit the context of the communication, often by providing additional information (e.g., their relationship with the interlocutor) to the assistant. For example, I5 copied entire email threads into the LLMs and asked, “Please proofread my reply,” to ensure their response was both grammatically sound and aligned with the ongoing exchange. I9 even checked short messages for online meetings, such as “Will you be joining soon?”, to understand the nuance of the message that they may be

implying in the message: “I worry it might sound like I’m pushing, so I ask ChatGPT to review even simple texts before sending.”

## 4 LINGOQ

Our formative study revealed that ESL workers suffer from *lexical disruption* while handling information work in English. In addition, our participants noted that most existing ESL learning systems rely on generic materials disconnected from their work contexts—despite research in ESL education showing that authentic, usage-based practice fosters learner engagement and improves both proficiency and self-efficacy [9, 10, 37]. To address these challenges, we designed and developed LINGOQ, a language querying and self-directed learning system that provides work-related quizzes generated from language queries. In this section, we discuss our design rationales from the formative study and literature. We then describe our system design and generative pipelines, along with implementation details.

### 4.1 Design Rationales

**DR1. Leverage AI-assisted language queries as a source of learning material.** In our formative study, nearly all participants relied on LLM-based AI assistants for work-related English tasks, such as looking up unfamiliar terminology, resolving confusing grammar, translating text, or proofreading their writing. We therefore treated users’ language queries with an AI chatbot as an authentic source from which we can learn the English assistance that they need and generate learning materials. Moreover, certain types of user queries, such as searching for a definition of a word or comparing input text with edited text, explicitly reveal their weakness in English proficiency.

**DR2. Optimize AI assistant interface for language querying.** Since our formative study participants frequently used generative AI assistants, such as ChatGPT, for language practice, we observed that their querying interactions were often inefficient and tedious. For example, participants had to repeatedly type boilerplate commands in the input (e.g., “Translate this into Korean:”) whenever they initiated a new query. In addition, participants often issued follow-up requests to format responses for their language tasks (e.g., requesting to display Korean and English text side by side or highlighting edited portions to track changes), which led to unnecessary back-and-forth dialogue turns.

Hence, we incorporated LINGOQUERY, an LLM-based assistant dedicated to language queries, with interactions and interfaces optimized for such usage. Given that participants in our formative study often copied text to query digital tools, we implemented a keyboard shortcut that directly copies and pastes selected text from the computer into a new chat message. We also introduced *query intents* that users can attach to an input message, which automatically insert predefined yet customizable prompts for three frequent request types—look up, translate, and proofread—thereby avoiding manual typing of boilerplate instructions. The AI responses to these query intents are rendered in language-relevant message displays (see Figure 2).

**DR3. Streamline reviewing work-related language activity.** Participants in our formative study attempted to review vocabulary or

expressions used in daily tasks at work, but the burden of manually collecting and revisiting materials without explicit exercises hindered sustained engagement, particularly after work. Inspired by literature suggesting that microlearning with short practice sessions embedded into daily routines can foster sustained engagement and improve proficiency [36, 49, 53], we designed our system to generate bite-sized interactive quizzes. These quizzes are generated directly from automatically collected queries, supporting continuous practice without extra burden outside work. We adapted multiple-choice fill-in-the-blank question formats from standardized tests (e.g., TOEFL, TOEIC, GRE) to support vocabulary and grammar practice. While the format of problems can be diversified (e.g., reading/listening comprehension, open-ended questions, short writing tasks, etc.), we limit our study material to the fill-in-the-blank multiple-choice question (MCQ) format for its simplicity to support easy access to English study anytime, anywhere on a mobile phone. The effectiveness of other question formats lies beyond the focus of this work, which is to generate study material relevant to work.

Based on these design rationales, we developed LINGOQ that consists of LINGOQUERY (4.2), LINGOQUIZ (4.4), and the backend pipeline (4.3). LINGOQUERY is a desktop-based AI assistant that users can share English-related queries or discuss freely. The backend pipeline manages user query data from LINGOQUERY and generates questions and curates them into quizzes. LINGOQUIZ is a mobile application that offers 10-question quizzes generated from the user’s dialogues with LINGOQUERY.

## 4.2 LINGOQUERY

**4.2.1 Interaction Components of LINGOQUERY.** LINGOQUERY adopts the typical interface design of desktop versions of LLM-based AI assistants, such as ChatGPT [73] and Claude [5], while incorporating bespoke interaction components tailored to the English-language query contexts. The sequence of chat messages is organized into chat threads, and users can either start a new thread or append messages to existing ones by selecting them from the sidebar (Ⓐ in Figure 2). By default, the AI response messages are rendered as a markdown-formatted view to accommodate any LLM-generated responses.

**Language Query Intent Selection and Prediction.** The system supports three predefined query intents: (1) *Look up*, (2) *Translate*, and (3) *Proofread*. When composing a new message, users can explicitly select a query intent in the chatbox by pressing the buttons at the top, which load a predefined prompt (Ⓒ in Figure 2). Each query intent applies a prompt template that is concatenated with the user’s message input; for example, “Please explain the meaning of the following word (or expression) in detail in dictionary format” for *Look up*, which the user can edit before sending. If no query intent is selected, the message in the chatbox is treated as a plain prompt in text, and the system automatically predicts its query intent when generating a response. If a user’s query corresponds to one of the three query intents, the system offers a customized view that highlights the structured information of the intent type. The *Look up* response type follows a typical dictionary format (Ⓑ in Figure 2); the *Translate* response

type provides a side-by-side view for comparing original and translated text (Ⓔ in Figure 2); and the *Proofread* response type displays a formatted container showing the proofread text with the rationale for edits underneath, along with an option to toggle track changes so users can quickly see where edits were made (Ⓕ in Figure 2). These views were informed by how interviewees in the formative study customized responses when using generic LLM assistants.

**Shortcuts and Contexts.** To enable users to receive assistance within the context of work-related applications, LINGOQUERY provides an operating system-level shortcut to trigger a query. When the user highlights text anywhere on the computer and presses ‘Ctrl + Cmd + C’ on MacOS or Ctrl + Alt + C on Windows, the LINGOQUERY window opens with a new chat thread and the copied text pre-filled in the chatbox. At the same time, the system captures a screenshot of the active window and displays it alongside the text when the shortcut is pressed. Before submitting the query, the user can choose to include the screenshot with the message or remove it if it contains sensitive information. If the screenshot is included, the system runs image understanding on OpenAI’s GPT-4o model to extract the text surrounding the copied content and infer the nature of the tasks based on metadata of the visible application on a screen, enriching the content for question generation later. The inferred context will be displayed in the UI as well.

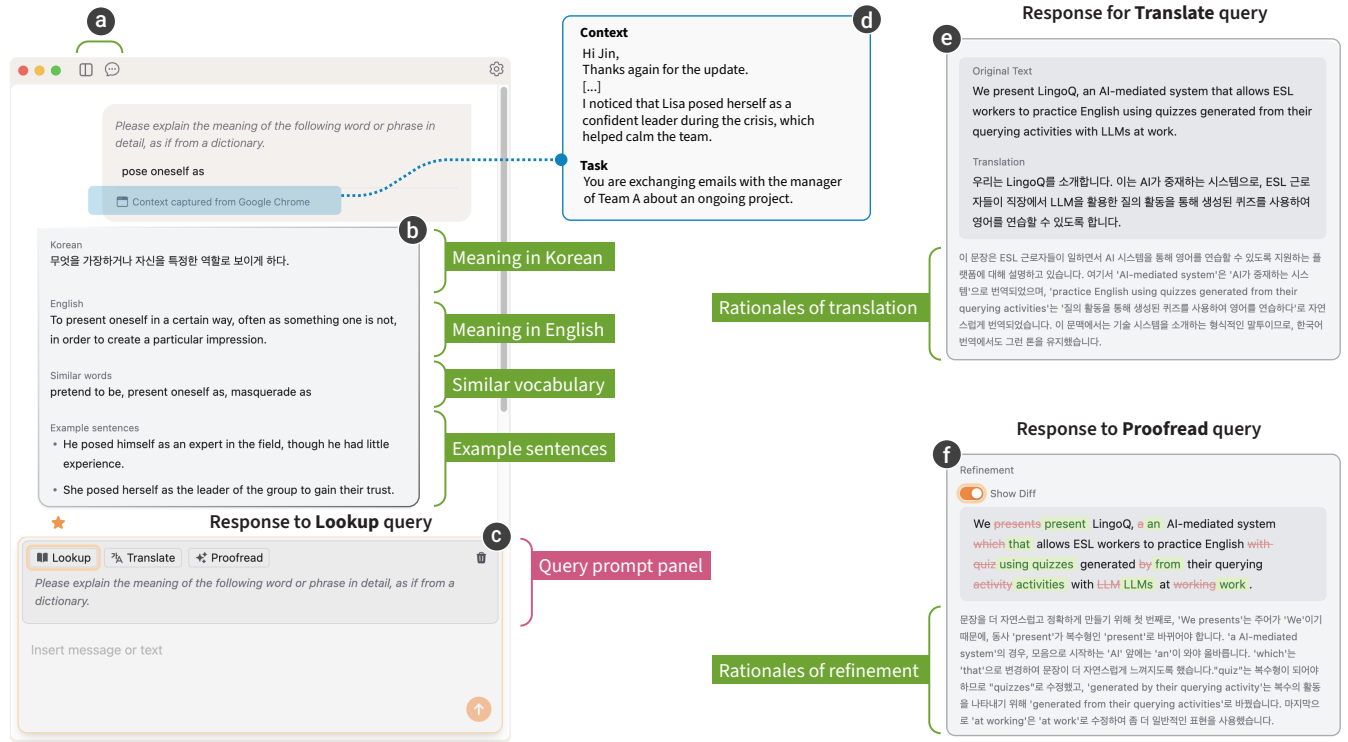
**Marking Messages for Prioritizing Question Generation.** Users can also mark noteworthy AI responses via the ★ star icon (right above Ⓒ in Figure 2). Marked messages are processed in LINGOQUIZ, increasing the likelihood that the corresponding question will be included in a quiz. This mark feature allows users to flag particular words or expressions they wish to review within the context of LINGOQUERY, removing the need to manually track what they want to study.

**4.2.2 LINGOQUERY Conversational Pipelines.** LINGOQUERY is a self-contained app that generates responses to user requests, similar to an LLM-based AI assistant with additional customization for specific query types. A user’s query input is processed before being sent to the LLM engine (OpenAI API in our case). Figure 3 illustrates the response generation pipeline of the LINGOQUERY conversational agent when it receives a new user message (Ⓐ in Figure 3). LLM-based **Intent Classifier** (Ⓑ in Figure 3; see Section A.1 for the instruction provided to the LLM) determines the corresponding query intent (Ⓒ in Figure 3). Both the chat history (Ⓓ in Figure 3) and the detected intent are then passed to the **Response Generator** (Ⓔ in Figure 3), which produces an AI response using an LLM (See Section A.2 for the instruction provided to the LLM). When the query intent does not correspond to a plain-text message but instead falls into one of the intents *Look up*, *Translate*, or *Proofread*, the LLM output is returned as a JSON object containing relevant attributes (e.g., the *original* input, *refined* text, and the *rationale* of refinement in the case of a *Proofread* intent). The structured format enables the application interface to render the output through a bespoke UI (Ⓔ in Figure 3).

## 4.3 Question Generation Pipelines

The question generation in LINGOQ follows three pipelines—question generation, question quality evaluation for filtering, and question





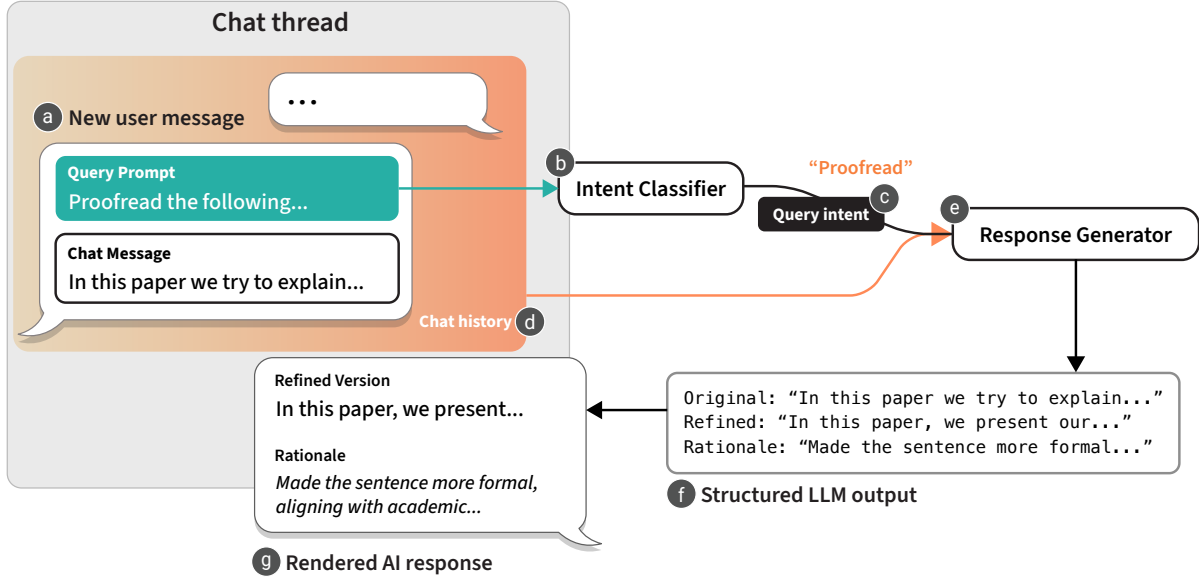
**Figure 2: Main window and the interface components of LINGOQUERY.** Users can open new chat threads of their choice in the thread list or via the New Chat button (a). The AI’s responses for three major types of query intents—*Look up* (b), *Translate* (e), and *Proofread* (f)—provide the UI components tailored to each query type. The new message panel (c) incorporates a query prompt panel at the top. Users can insert a template query prompt in their message via the three quick-access buttons. When a query is sent via a keyboard shortcut, the system analyzes a screenshot of the user’s active window, and the user can review both the surrounding context of the copied sentence and the task at the time of the query (d). By clicking the star icon ★ below AI responses (b), users can mark the message to increase the likelihood that the corresponding question will be included in LINGOQUIZ.

selection—that transform interactions collected through LINGOQUERY into validated quiz questions for LINGOQUIZ.

**4.3.1 Generating Questions from Conversation and Context.** Figure 4 illustrates the question generation pipeline, which periodically produces fill-in-the-blank multiple-choice questions from query-response pairs collected in LINGOQUERY and captured context data. Every five minutes, the system checks for new query-response pairs collected from LINGOQUERY (a in Figure 4). For each pair, an LLM-based module evaluates whether the query is English-related (b in Figure 4); if not, it is filtered out from the question generation pipeline. For eligible queries, the system takes the eligible pair and additional information—the conversation history with contextual data captured from the screenshot—to initiate generation (d in Figure 4). The pipeline applies a different system prompt with few-shot examples modeled after TOEFL, TOEIC, and GRE example questions (e in Figure 4). For each conversation, two distinct questions are generated to ensure variety (f in Figure 4). Each question will generate one structured output in JSON format that contains: stem with a blank, key, distractors, explanation, and rationale for question generation (g in Figure 4).

Finally, the contextual information extracted from a screenshot and/or conversation history is fed into the question generation pipeline to produce questions aligned with that context. The LLM-based question generation module ensures that the question stem is relevant to the worker’s task context. For example, a user query may simply involve searching for a word (e.g., “airway”), in which case the surrounding text (e.g., “the patient’s airway to ensure proper breathing”) from the screenshot can be used to extract the context and generate a relevant question stem. When a user submits an answer in LINGOQUIZ, the explanation provided will include this context to supplement the rationale for the correct answer.

**4.3.2 Quality Assurance of Generated Questions.** To ensure the quality of generated questions, they are evaluated by an LLM-based evaluation module informed by prior literature [26, 34]. The evaluation applies two binary criteria: answerability, that is, whether the question can be clearly and correctly answered, and proficiency, that is, whether the question requires an appropriate level of English skill and is not too easy to answer (h in Figure 4). Questions that fail one of the criteria are iteratively refined—up to three times in total—by feeding the evaluator’s rationale for failure, along with



**Figure 3: Conversational pipeline of LINGOQUERY.** When the user sends a new message (a), the intent classifier (b) identifies the query intent (c), which is then passed to the response generator (e) together with the chat history (d). The response generator produces an appropriate response (f) structured according to the query intent. Finally, LINGOQUERY renders this structured response accordingly (g).

the original input, back into the question generation module (i) in Figure 4). If a question passes within three iterations, it is added to the question pool (j) in Figure 4), otherwise discarded.

**4.3.3 Question Pool and Selection Logic.** After passing quality checks, questions undergo a final format validation to ensure that all required components—stem with blank, distractors, key, and explanation—are properly structured. Validated questions are then added to the question pool. When a user initiates a quiz session, ten questions are drawn from the pool. Each quiz contains 10 questions: 7 are selected from newly generated questions, and the remaining 3 are randomly drawn from the pool of questions that a user has solved previously using weighted probability. The system assigns higher weights to questions that have been repeated less frequently, answered incorrectly in the past, marked with the ★ star icon in LINGOQUERY, or not practiced recently. As a result, each quiz balances new questions (70%) with review questions that workers need to revisit. The examples of generated questions for each type are available in Figure 6.

## 4.4 LINGOQUIZ

Users can practice work-related English vocabulary and grammar by solving questions in LINGOQUIZ, generated from their dialogues with the LINGOQUERY chatbot. LINGOQUIZ provides a dashboard (A in Figure 5) that helps users track how many quizzes they have completed that day, how many they have completed in total since the beginning, and how many new questions are available in the question pool (i in Figure 5). Clicking ‘Start Quiz’ launches 10 multiple-choice questions, each requiring users to fill a blank with the best of three options (B in Figure 5). When the question is

generated from a marked AI response or when it is the first attempt, it is displayed with a ★ or ‘new!’ badge (2 in Figure 5). Questions generated from screenshots or sufficient thread context include the inferred task as a hint (3 in Figure 5). When users press the Submit button after selecting an option, the app provides immediate feedback on whether the answer is correct and the explanation (4 in Figure 5). Within the quiz, any incorrect questions reappear until users provide the correct answer. Once all ten questions are answered correctly, the progress bar completes (5 in Figure 5), and the quiz ends with a completion screen (C in Figure 5). Afterward, users may proceed to a new quiz or return to the dashboard.

## 4.5 Implementation

We implemented the core system in Python running on a FastAPI [35] server that provides REST APIs for both LINGOQUIZ and LINGOQUERY. The chat history, generated quizzes, and user interaction data are stored in a PostgreSQL [41] database on the server. The conversation and the question generation pipelines leverage OpenAI’s Chat Completion APIs [75] on top of the LangChain [58] framework to run the underlying LLM inferences. All LLM inference and image understanding tasks are performed using a gpt-4o model. To protect user queries that may contain sensitive information, we used OpenAI Enterprise, which neither uses our data for training nor retains them.

We built LINGOQUERY as a cross-platform desktop application using Electron [76], to support both Windows and MacOS desktop computers. The LINGOQUIZ app was implemented using React Native [67] as a cross-platform mobile application running on both iOS and Android phones. Both apps were written in TypeScript [68] and communicate with the server via REST API.



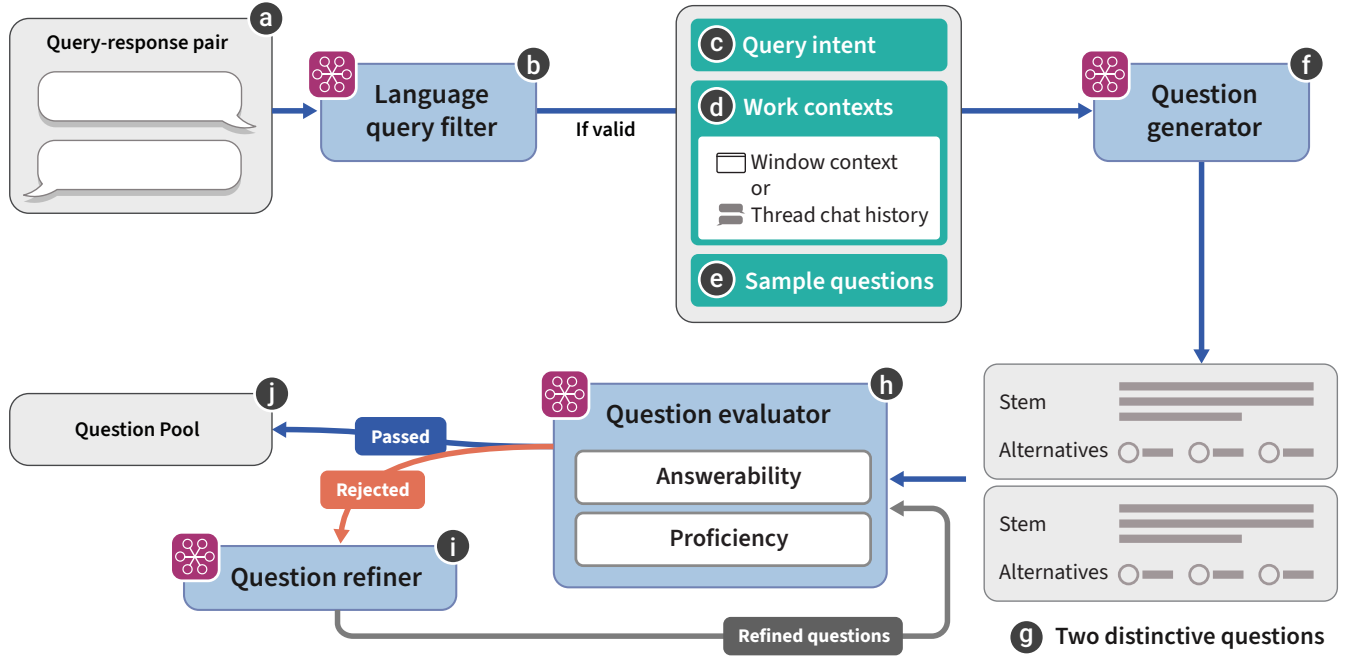


Figure 4: Question generation pipeline of LInGoQ. When a query–response pair arrives ①, the language query filter ② identifies the query intent ③, which is then passed to the Question generator ⑥ together with work contexts ④ and exam samples ⑤. The generator produces two candidate questions ⑦, which are evaluated by the Question evaluator ⑧ on two criteria: answerability and proficiency. The Question refiner ⑨ refines each question up to two iterations, and items that still fail are discarded. Accepted questions are stored in the question pool ⑩.

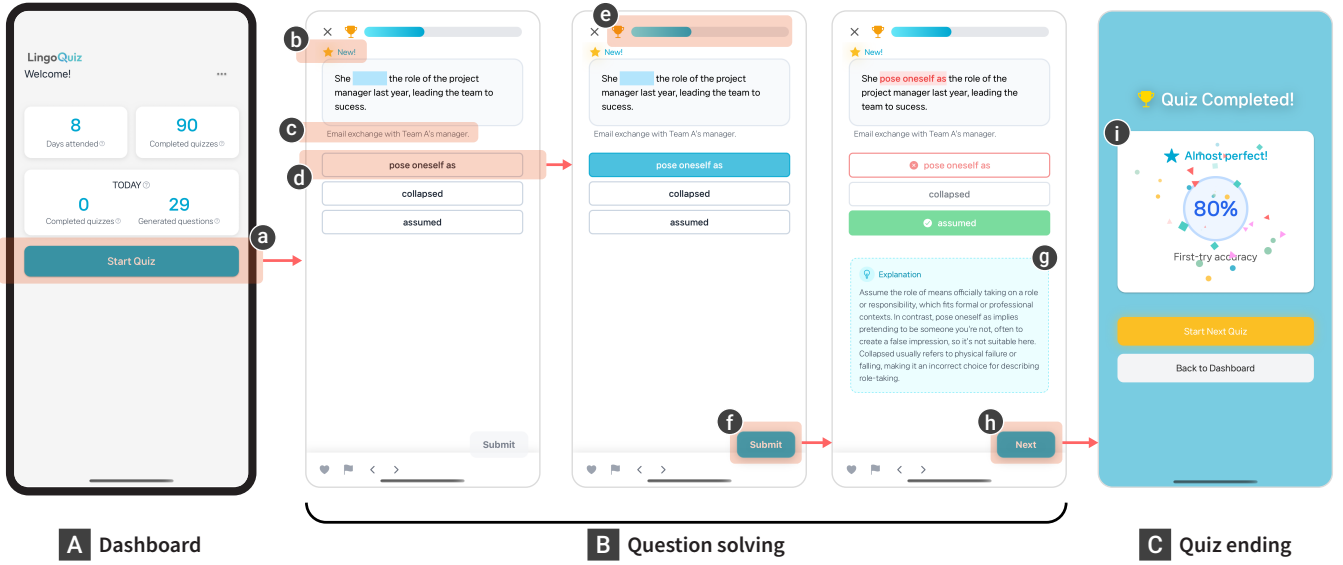
Table 1: Demographic information of the participants in our deployment study

Alias	Age	Gender	CEFR	Job title	Alias	Age	Gender	CEFR	Job title
P1	28	Female	C1	Global Business Developer	P15	29	Male	A2	Clinical Psychology Trainee
P2	28	Female	A2	Software Engineer	P16	41	Male	A1	IT Security Manager
P3	39	Female	B2	Hospital Operations Manager	P17	27	Female	B1	Graduate Student
P4	28	Male	A2	Medical Resident	P18	28	Male	B1	Machine Learning Engineer
P5	32	Female	B1	Product Designer	P19	25	Male	B2	Graduate Student
P6	36	Female	B2	Software Engineer	P20	29	Female	C1	Real Estate Professional
P7	47	Female	B1	Kindergarten Counselor	P21	32	Female	A1	Biotech Researcher
P8	45	Female	B2	Administrative Coordinator	P22	48	Female	C1	Logistics Specialist
P9	37	Female	B2	Office Manager	P23	25	Female	B1	Graduate Student
P10	29	Male	A2	International Patient Coordinator	P24	30	Female	B2	Marketing and Project Manager
P11	30	Male	B2	Sports Event Manager	P25	33	Male	A1	Network Engineer
P12	38	Female	B2	Export-Import Specialist	P26	31	Female	C1	Nuclear Policy Researcher
P13	42	Male	C1	Professor	P27	32	Female	C1	Sports Event Manager
P14	35	Female	B2	Apparel Export Manager	P28	30	Male	C1	Governance Administrator

## 5 Deployment Study

To evaluate the effectiveness of LInGoQ in engaging information workers in ESL learning, we conducted a three-week field deployment study with 28 ESL workers. We aimed to examine how ESL

workers engage with LInGoQ compared to their past learning experiences, and how it affects their English proficiency and self-efficacy. In addition, we conducted an expert evaluation to assess the quality of questions generated by LInGoQ. The study was conducted



**Figure 5: Main screens of LINGOQUIZ.** In the Dashboard screen **A**, users can check their records and stats, along with the number of new questions add to their question pool today. When starting a quiz by pressing the Start Quiz button **(a)**, a quiz with 10 unique questions are provided sequentially **B**. The new question that appears to the user for the first time is indicated by the star icon **(b)**. For questions generated from messages with context, the task description is provided **(c)**. To solve the question, the user can select an option **(d)** and press the Submit button **(f)** to submit an answer. Then the result the question is shown immediately, with explanation **(g)**, regardless of whether the user had selected a correct answer or not. After solving the ten questions, questions with wrong answers appears again, until all are answered correctly. The progress bar **(e)** indicates the current progress. In the Ending screen **C**, users can practice a new quiz or return to the Dashboard screen.

in South Korea with Korean native speakers and approved by the Institutional Review Board at the author’s university.

## 5.1 Participants

We recruited computer-based information workers by advertising our study on social media and a local social community platform. Our inclusion criteria was information workers who are: (1) working at least 30 hours per week, (2) using a computer as the primary work tool, (3) regularly performing tasks involving English, such as information access, communication, and document writing, (4) being a native Korean speaker, and (5) being an ESL learner As a minimum requirement for study completion, we instructed the participants to submit at least two queries in LINGOQUERY or complete at least one quiz in LINGOQUIZ, for at least ten days. Initially, 32 workers participated in the study but four were excluded during the analysis for not meeting the minimum requirement.

Finally, a total of 28 ESL workers (Table 1; P1–P28; 18 females and 10 males) completed the study. Participants were aged between 25 and 48 years old ( $M = 33.4$ ,  $SD = 6.5$ ) and represented diverse professional domains, including IT (7), healthcare (5), science (4), finance/business (4), international trade (3), education/welfare (2), sports management (2), and design (1). Regarding occupations, participants included office workers (14), engineers (5), graduate students (3), researchers (2), educators (2), and medical professionals

(2). Based on CEFR [71]<sup>2</sup> self-assessed English proficiency, 3 participants identified themselves as *A1* (beginner), 4 as *A2* (elementary), 14 as *B1* (intermediate), and 7 as *C1* (advanced). All participants reported daily use of LLM-based chatbots during their workdays. Additionally, they had prior ESL learning experience for work, including tutoring (17), vocabulary apps (13), and English media (12; TV shows, news, or novels). As compensation for their participation, we offered 200,000 KRW (approx. 144 USD) based on the required system usage over three weeks.

## 5.2 Procedure

**Pre-study Preparation.** Before the introductory session, we sent participants a link to a pre-study survey and a pre-study English proficiency test. The survey included three items on a 5-point Likert scale that generally assess the perceived relevance, effectiveness, and engagement of their past ESL learning methods, along with 16 items from the Questionnaire of English Self-Efficacy (QESE; eight on reading and eight on writing) on a 7-point scale, excluding speaking and listening to align with our research focus [92].

The English proficiency test consisted of 28 multiple-choice items selected from questions of TOEIC (Test of English for International Communication) [33], a standardized English proficiency test for

<sup>2</sup>The Common European Framework of Reference for Languages (CEFR) defines proficiency levels as basic (*A1*: beginner; *A2*: elementary), independent (*B1*: intermediate; *B2*: upper-intermediate), and proficient (*C1*: advanced; *C2*: proficient).

**(a)**

When these results are comprehensively considered with the brain MRI results conducted in 2022, a high likelihood of Alzheimer's disease dementia with [ ] is suggested.

Context: Reviewing a medical report

☒ ischemia

☐ anemia

☐ asthma

**(b)**

The first stall in the women's restroom frequently gets [ ] due to improper disposal of sanitary products, so feedback was sent to the building management.

Context: Reporting an issue to building manager

☒ clogged

☐ blocked

☐ jammed

**(c)**

Please find attached the [ ] plan for the 17th Marine Festival Regatta.

Context: Reviewing the plan for the Marine Festival Regatta

☒ transportation

☐ transmission

☐ transitory

**Figure 6: Selected questions generated by LINGOQ during the deployment study.** Each question consists of a stem with a blank, context reminding the task at the time of query, and three alternatives with the correct key marked. (a) P15, (b) P20, and (c) P11 show actual question examples.

general business. The test included two types of questions: 16 simple fill-in-the-blank items—each with a single sentence and one blank—and three sets four fill-in-the-blank items, each requiring participants to complete blanks within a single paragraph. Before the study, we finalized the items from 46 questions, by administering them to 29 information workers—who are not our deployment study participants—and selecting those whose percentage of correct answers fell between 40% and 80%, as suggested in Classical Test Theory [23]. (see Appendix B for details of the item validation.) Since the TOEIC questions are proprietary, we do not report the actual questions.

**Introductory Session.** A group of 2–3 participants attended an 1-hour, in-person introductory session, bringing their laptop to install LINGOQUERY. After explaining our study motivation and goals, we assisted participants with installing and setting up the system on their work computers and mobile devices. To ensure participants fully understood how to use the system, we ran a hands-on tutorial covering the main features of LINGOQ. After the walkthrough, participants practiced using LINGOQUERY and LINGOQUIZ on their devices. They were instructed to enter at least five English-related queries into LINGOQUERY, drawn from recent questions they had asked LLM assistants at work. They then tried LINGOQUIZ after five minutes when the system generated a quiz from the submitted queries. This step ensured that participants had quizzes available during the early stage of the deployment period. After this session, participants were asked to install LINGOQUERY on other computers they use, is possible.

**Deployment.** Immediately after the introductory session, participants began using LINGOQ for three weeks. During this period, participants were instructed to direct their English-related queries

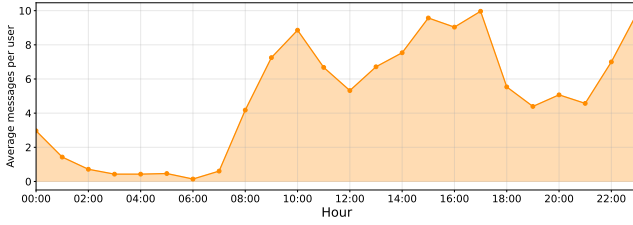
to LINGOQUERY during their study, instead of using ChatGPT or other tools. In parallel, LINGOQUIZ was accessible at any time of the day, and questions could be paused and resumed at their convenience.

After the first and second weeks, we sent participants a text message summarizing LINGOQ usage, including the number of submitted queries and completed quiz sessions to remind them of the minimum requirement for study completion. Additionally, when participants were inactive for more than three days, or had not attempted the questions generated from that day’s LINGOQUERY conversations, LINGOQUIZ sent an evening push notification.

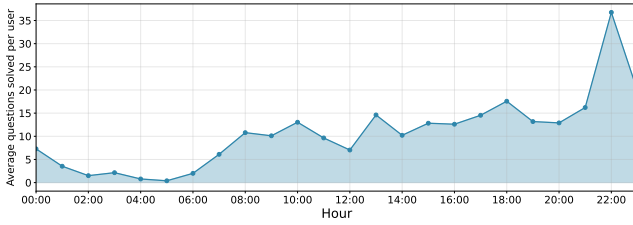
**Exit Survey.** The day after the 3-week deployment period, we sent an online post-study survey and a post-study English proficiency test. The survey reassessed QESE [92] and the three ratings of perceived learning experience used in the pre-study survey, but targeted for LINGOQ, and was supplemented with follow-up questions probing the reasons for participants’ ratings. It also asked about participants’ intention to use LINGOQ on a 5-point Likert scale and included open-ended questions regarding overall user experiences and suggestions for design improvements. In addition, the post-study English proficiency test contained the same questions as the pre-study test, but the order of questions and answer options was randomized.

### 5.3 Expert Evaluation of Questions

To assess the performance of question generation pipeline from the perspective of language education, we conducted an expert evaluation using a subset of questions generated during the deployment study. We recruited three English language teaching experts (E1–E3; all female) through social media advertisements. Their professional backgrounds included university and high-school teaching as well



**Figure 7: Hourly engagement patterns of LINGOQUERY during the three-week deployment. The orange line plot shows the average number of user messages per hour across a 24-hour day, with notable peaks around 10 a.m., 5 p.m., and 11 p.m.**



**Figure 8: Hourly engagement patterns of LINGOQUIZ during the three-week deployment. The blue line plot shows the average number of solved questions per hour across a 24-hour day, with a clear peak around 10 p.m.**

as the development of standardized English test materials. They were aged 39, 53, and 39, with 10, 20, and 15 years of experience in English education, respectively. The experts evaluated 30 questions from LINGOQ, randomly sampled in our problem generation pipeline (24 passed items and 6 failed items based on the Quality Assurance logic stated in 4.3.2). Evaluations were conducted remotely via Zoom, combining quantitative measures collected through an online survey with qualitative insights obtained from follow-up interviews that asked about the potential benefits of the approach that we take in LINGOQ, which generated ESL learning from their work context. The evaluation and interview took 90 minutes to complete. We compensated the experts with 100,000 KRW (approx. 72 USD).

Building on prior work on educational question quality evaluation [26, 34], we defined two quantitative criteria for each question: (1) Answerability, based on whether the key is actually a correct answer, and (2) Proficiency, based on the presence of unique choices and the absence of obviously incorrect options (see detailed rubrics in Appendix C). To enable comparison with our problem generation pipeline, we applied the same criteria but expanded the response format from true/false to four or five categorical options for more specific evaluations. For each question, we also included an open-ended field requesting experts’ rationale when they indicated uncertainty. In the follow-up interviews, we introduced LINGOQ and explored experts’ perspectives on question quality relative to

authentic English test items, as well as their views on the potential and concerns of AI-generated questions and suggestions for improvement.

## 5.4 Data Analysis

To explore participants’ engagement and usage patterns with LINGOQ, we first conducted a descriptive statistics analysis. We analyzed participants’ interactions with LINGOQUERY and LINGOQUIZ. For LINGOQUERY, this included the number of user messages–AI response pairs, usage days and patterns, and the distribution of query prompts. For LINGOQUIZ, we focused on the number of solved questions, usage days, and solving patterns, and quiz progress across repeated attempts.

To evaluate the generated questions, we examined pipeline performance and human evaluations. Pipeline performance was evaluated on 30 generated questions by comparing its binary judgments with expert labels. Ratings from three experts were aggregated by majority vote for answerability and proficiency, coded as true or false (with “unknown” mapped to false), and used as ground truth. We then calculated precision, recall, and F1-score, with precision measuring agreement on pipeline-accepted items, recall measuring agreement on expert-accepted items, and F1 as their harmonic mean. To gain a deeper qualitative understanding, we examined the questions generated by LINGOQ against fill-in-the-blank items from standard proficiency tests, and analyzed participant and expert feedback through open coding of surveys and interviews.

We analyzed pre- and post-study surveys and tests to assess the effects of LINGOQ on participants’ learning performance and experiences. To examine changes in learning performance over the study period, we analyzed the pre- and post-study English proficiency test using a mixed-effects model and the English self-efficacy questionnaires (QESE). To investigate learning experiences relative to prior ESL practices, we compared participants’ perceived effectiveness of LINGOQ with their past ESL learning using the Wilcoxon signed rank test across three criteria: relevance of materials to work, helpfulness for actual tasks, and sustainability of engagement. To gain deeper qualitative insights into participants’ interactions with LINGOQ and their feedback, we applied open coding of the survey responses.

## 6 Findings

This section presents the study’s findings, moving from how participants engaged with and used LINGOQ, to evaluations of the generated quiz questions, to the impact on ESL learning, and finally to feedback for future improvements.

### 6.1 Usage Patterns and Engagement

The interaction logs and usage data indicated that participants actively engaged with LINGOQ, frequently using both LINGOQUERY and LINGOQUIZ. Here we report the descriptive statistics regarding participants’ usage patterns and engagement of the two apps.

**Active Querying with LINGOQUERY.** Across three weeks, participants opened a total of 652 conversation threads, and submitted 3,325 messages ( $M = 118.8$  per participant) through LINGOQUERY. On average, participants used LINGOQUERY for 13.2 days ( $SD = 2.5$ ,  $min = 10$  [P15],  $max = 19$  [P5]), exceeding the required 10 days of



use. This indicates that participants engaged on most weekdays. [Figure 7](#) presents participants' hourly engagement patterns, showing peak usage during work hours, particularly around 17 o'clock.

Participants actively used pre-defined query prompts—i.e., *Look up*, *Translate*, and *Proofread*—or wrote their own when submitting a message. Out of the 3,325 query-response pairs, *Translate* responses were the most common (1,271 responses; 38.2%), followed by *Look up* (399 responses; 12.0%) and *Proofread* (287 responses; 8.6%). The rest of the responses (1,369 responses; 41.2%) were plain text messages, such as responses to queries asked in plain text or follow-ups.

Eighteen Participants regularly used the *marking* feature ★, which ensures that the particular message pairs would appear in future quizzes. They marked 13.4 AI responses per person on average ( $SD = 16.1$ ). Of the 241 marked messages, 91 messages (37.8%) were responses for *Look up* queries, indicating participants' desire to revisit vocabulary or expressions. Participants opened LINGOQUERY by directly capturing the selected text and surrounding context using keyboard shortcuts ⌘, for 6.9% of all queries. However, only three participants (P11, P19, P26) dominated the usage of this feature and accounted for 65.8% of all shortcut-triggered messages.

**Consistent Language Practice with LINGOQUIZ.** Of the 3,325 query-response pairs, our question generation pipeline classified 2708 (81.4%) pairs as English language queries and included to the question generation pool. The pipeline initially produced 5,682 questions, of which 997 (17.5%) were removed during the automated quality assessment. As a result, 3,290 questions were eventually exposed to participants (117.5 per participant). Including the reappeared cases, participants solved a total of 7,155 questions (255.5 per participant). These questions were curated in 604 quizzes and participants completed most ones, leaving only 10 quizzes incomplete (1.7%) throughout the study period. Over the three weeks, Participants completed at least one quiz in LINGOQUIZ for 13.4 days on average ( $SD = 2.8$ ,  $min = 10$  [P27],  $max = 19$  [P21]), indicating similar compliance with LINGOQUERY. Participants completed an average of 1.04 quizzes per day ( $min = 0.32$  [P9],  $max = 2.11$  [P8]), spending about 9.3 minutes ( $SD = 3.0$ ) per quiz. This result is more than twice the number of quizzes required to qualify for study completion; the minimum requirement was 10 days out of 21, at least one quiz per day, or roughly 0.5 quizzes per day on average. [Figure 8](#) presents participants' hourly engagement patterns, showing a more than twofold increase in usage around 10 p.m. This pattern aligns with trends observed in other popular mobile applications.

Of 3,290 unique questions, 927 questions (28.2%) appeared more than once, with the most frequently reappeared item occurring 15 times. On average, each question appeared across 2.86 quizzes ( $SD = 1.48$ ). As the questions were repeatedly presented, participants became more likely to answer them correctly. The average accuracy for questions presented for the first time was 82.6%, which gradually increased to 89.5% upon the second exposure in another quiz, and further to 92.4% upon the third exposure.

## 6.2 Utility and Quality of Questions in LINGOQUIZ

In this section, we cover the utility of the generated questions, with respect to participants' perceptions of their relevance and

helpfulness, as well as experts' evaluations and feedback on a subset of questions.

In particular, we compared participants' perceived relevance and helpfulness of learning materials of their prior ESL practices (measured in the pre-study survey) and LINGOQ (measured in the post-study survey), to their work tasks.

**Content Relevance.** The Wilcoxon signed-rank test revealed that participants rated the relevance of LINGOQ ( $M = 4.21$ ,  $SD = 0.57$ ) significantly higher than prior ESL practices ( $M = 2.61$ ,  $SD = 0.83$ ),  $z = -4.39$ ,  $p < 0.001$  (see [Figure 9\(a\)](#)). In the post-study survey, 25 participants valued the quizzes for reflecting practical English they actually encountered in their professional work, as opposed to general English. In particular, P15 emphasized the value of context-relevant words: “*I liked repeatedly practicing verbs specific to the medical field rather than casual spoken language.*” ([Figure 6\(a\)](#); see [Quiz](#)). Moreover, P11 and P20 found that domain-relevant distractors in quiz alternatives helped them contrast similar terms and deepen their understanding of subtle distinctions ([Figure 6\(b\)](#); see [Quiz](#)).

**Helpfulness for Work Tasks.** The Wilcoxon signed-rank test revealed that participants rated LINGOQ ( $M = 4.14$ ,  $SD = 0.65$ ) as significantly helpful for daily work tasks compared to their prior ESL practices ( $M = 2.89$ ,  $SD = 0.79$ ),  $z = -4.05$ ,  $p < 0.001$  (see [Figure 9](#)). In the post-study survey, participants reported that practicing work-related content with LINGOQUIZ not only improved retention but also made English-related work tasks smoother and more efficient. They remarked that LINGOQUIZ reinforced their learning by making them “*encounter the content again through quizzes*” (P17), which they “*had previously skimmed over while working*” (P2) and “*normally read quickly and move on from*” (P27). In particular, P18 and P19 highlighted that reviewing previously translated or looked-up content through quiz questions enhanced their reading fluency, especially when thoroughly re-reading research papers: “*Since I had the opportunity to revisit documents and papers I had read before, I found that when rereading, I was able to process them more quickly in English.*” (P18).

**Expert Evaluation.** We compared the expert's assessments of the quality of questions with the assessment from our automated quality assessment pipeline. We created the expert's assessment of acceptance/rejection of the criteria for each sample question by majority voting (e.g., if two or more experts marked the question as rejection for Answerability, we eventually labeled it as rejection). For this analysis, since experts were allowed to mark ‘unknown’ for uncertain cases, we conservatively treated them as labeled rejection. Our comparison yielded precision/recall/F1-scores of 0.91/0.81/0.86 for Answerability and 0.85/0.92/0.88 for Proficiency, suggesting that our automated filtering provided highly aligned decisions compared with the experts' judgment, with minor discrepancies. We identified two main reasons for the discrepancies. First, experts often marked domain-specific questions as “unknown” (5 items for Answerability, 10 items for Proficiency), making it difficult to assess their quality as even the experts were not familiar with domain-specific terms. Second, they applied a higher bar for proficiency, according to the follow-up interview, as they were accustomed to carefully adjusting difficulty to learners, to maintain discrimination of the test design.

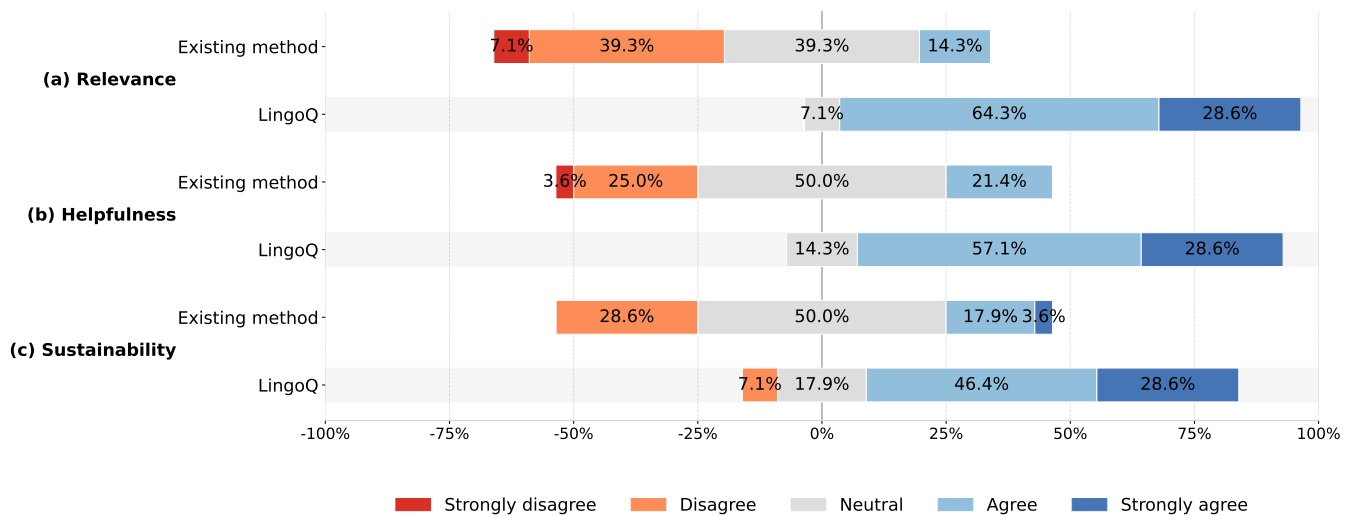


Figure 9: Stacked bar charts of five-point Likert ratings from participants ( $N = 28$ ) on (a) content relevance, (b) helpfulness for work tasks, and (c) sustainability in learning. Upper bars indicate pre-study evaluations of existing ESL methods, while lower bars indicate post-study evaluations of LINGOQ.

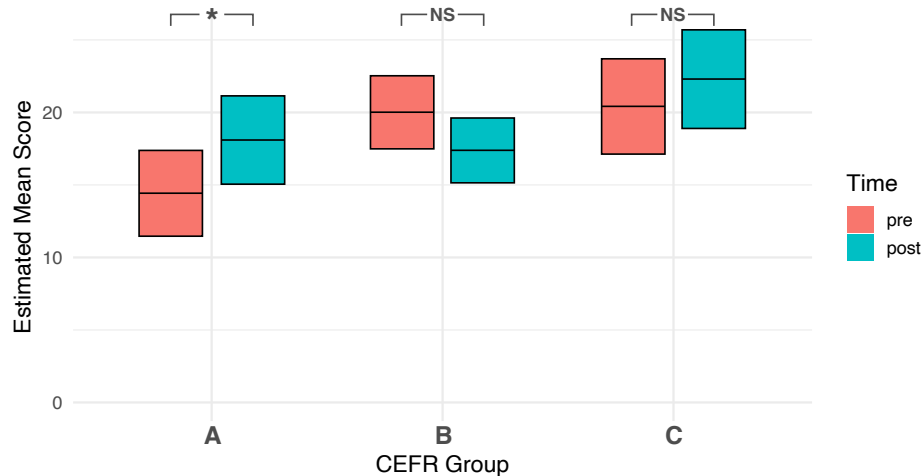


Figure 10: Mixed effect model estimates of pre- and post-test English proficiency scores by CEFR group. The plot shows estimated group means (with 95% CIs) for A (basic,  $N = 7$ ), B (independent,  $N = 14$ ), and C (proficient,  $N = 7$ ) on a 0–28 scale.

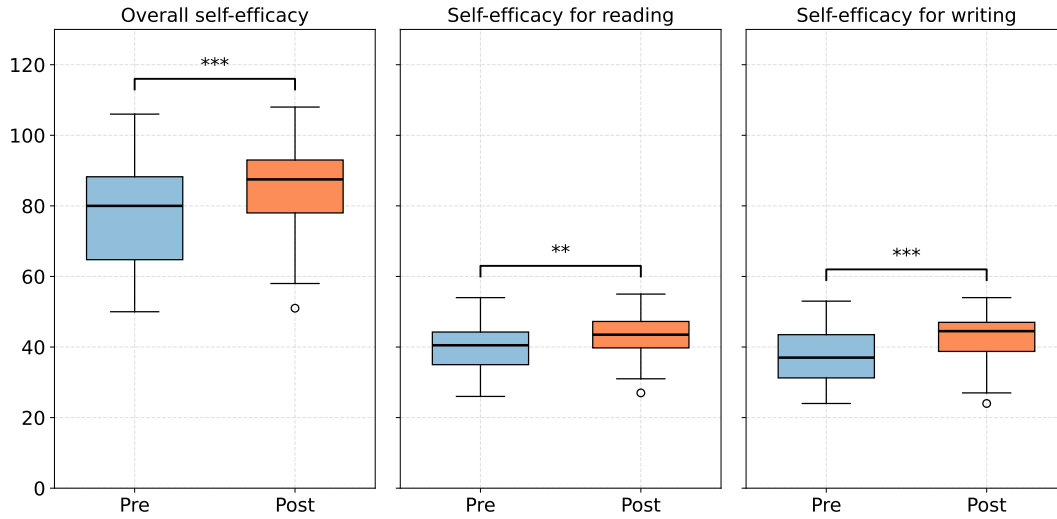
In the follow-up interviews, all experts emphasized that the key difference between LINGOQ questions and standardized English tests is in the contexts used in the question stems as E1 noted: “TOEIC usually covers general business contexts, but some of the questions from LINGOQ required knowledge confined to highly specific domains.” Two experts (E2, E3) valued the domain-specific stems, noting that the learners’ specialty in the domain could foster learning engagement and motivation, which are crucial factors for effective self-directed learning. Meanwhile, E3 remarked that the overall quality of LINGOQ questions was comparable to text-completion items in TOEIC or TOEFL, noting that some items (Figure 6(c); see

Quiz) resembled high-quality questions that could plausibly appear on actual tests.

### 6.3 Impact on Learning Performance

Over the three-week study, the use of LINGOQ significantly enhanced participants’ self-efficacy in using English at work. Proficiency gains varied by self-reported CEFR levels, with greater benefits for lower-level learners. These improvements were positively associated with LINGOQUERY usage.

**English Proficiency.** A mixed-effects model analysis revealed a significant main effect of time on English proficiency scores, with an average increase of 1 point across all participants ( $p = 0.01$ )



**Figure 11: Box plots of English self-efficacy (QESE) scores on a 7-point scale. The left plot shows overall self-efficacy (16 items), while the middle and right plots show the subscales of reading (7 items) and writing (7 items). Significant pre-post differences are observed in both the overall scale and the subscales. Significance is marked as  $p < 0.05$  (\*),  $p < 0.01$  (\*\*), or  $p < 0.001$  (\*\*\*).**

(see Figure 10). Post-hoc pairwise comparisons revealed that only basic (CEFR A) learners showed a significant improvement, gaining an average of 4 points (Total 30 points) from pre-to post-test ( $p = 0.01$ ), whereas independent (CEFR B) and proficient (CEFR C) participants showed no remarkable change. However, for the independent group, the interaction between time and the number of LINGOQUERY messages was significant ( $p = .01$ ). The result indicates that more frequent use of the LINGOQUERY was associated with greater learning gains among these ESL workers. We further explain the learning effects of querying activity based on participants' general reactions in Section 6.4.

**English Self-Efficacy.** We compared participants' perceived self-efficacy of English assessed using the Questionnaire of English Self-Efficacy (QESE) scores, measured twice in the pre- and post-study surveys. The paired  $t$ -test revealed significant improvement in QESE score from pre- ( $M = 77.43$ ,  $SD = 14.31$ ) to post-study ( $M = 84.75$ ,  $SD = 13.21$ ) measurements ( $t(27) = -4.30$ ,  $p < 0.001$  (see Figure 11). In addition, both the reading and writing subscales of QESE also demonstrated significant gains ( $t(27) = -3.67$ ,  $p < 0.01$  for reading, and  $t(27) = -4.29$ ,  $p < 0.001$  for writing), respectively. In the post-study survey, P15 highlighted the enhanced self-efficacy as the most notable benefit of LINGOQ, noting “*What improved the most was my confidence. It was really satisfying to go over the mistakes I often made, and over time, I found myself reading tough sentence structures much more easily.*” Additionally, an expert from the expert evaluation reinforced this point: “*Confidence is a key factor in conversational ability, as it often translates into greater written and spoken outputs by reducing fear and hesitation. Thus, fostering confidence is essential for advancing from intermediate to higher proficiency levels.*” (E1).

## 6.4 General Reactions to LINGOQ

In the post-study survey, we gathered participants' feedback on LINGOQ across multiple aspects. The Wilcoxon signed-rank test showed that the participants rated the sustainability of learning with LINGOQ ( $M = 3.96$ ,  $SD = 0.88$ ) significantly higher than prior ESL learning experiences ( $M = 2.96$ ,  $SD = 0.79$ ),  $z = -3.47$ ,  $p < 0.001$  (see Figure 9(c)). To gauge the utility of LINGOQ beyond the study context, we asked participants how much they would be willing to use LINGOQ in their real life (LINGOQ on their phones and computers would continue to work). 24 out of 28 participants (86%) expressed willingness (*i.e.*, “agree” or above) to adopt our system, including 15 (54%) who selected “agree” and 9 (32%) who selected “strongly agree”. In the following, we summarize participants' feedback on the strengths and drawbacks we identified, which inform potential design improvements.

**Expanding Types of Questions and Language Learning Disciplines.** Most participants (85.7%) found the quiz design effective for sustaining learning beyond work hours. P3 noted “*The quizzes weren't burdensome and fit easily into my routine, like during commutes or before bed.*” Similarly, P20 remarked “*I could learn just by doing the quizzes without the extra step of studying beforehand,*” contrasting LINGOQ to other mobile apps. Participants suggested diversifying the quiz formats beyond the current fill-in-the-blank design. They proposed that exercises could be more closely aligned with the types of queries submitted. For instance, for translation queries, quizzes could present multiple sentence options and ask learners to select the correct translation. Such alignment would make the exercises feel more relevant and effective for their specific learning needs. In addition, five participants (P7, P14, P19, P27, P28) suggested expanding the material modalities to include speaking and listening practice, aiming to better support verbal communication tasks such as video meetings.

**Perception Change of Queries as Learning Opportunities.** While we provided LINGOQUERY as a dedicated channel for language

querying, three participants mentioned that using such a scoped interface provided not only language support, but also meaningful learning opportunities by raising awareness of knowledge gaps and encouraging reflection on their English use.

They contrasted this experience with prior experience with AI assistants. P25 remarked, “Using LINGOQUERY instead of ChatGPT helped me develop the habit of looking more carefully at words in sentences I would have otherwise translated without much thought. I found myself learning by checking whether I already knew the words and whether the meaning was correct.” It seems that knowing their queries would generate learning materials to complete at the end of the day reinforced the connection between querying and language learning, reminding them of English learning even when they asked questions to LINGOQUERY. These reflections suggest the potential to reorient LLM reliance from passive consumption toward active learning.

**Backfire of Authentic Materials: Workplace Detachment.** While participants valued the activity of solving work-related quizzes, two also noted that practicing quizzes containing work-related materials outside of work hours sometimes discouraged them from practicing further. P8 noted, “Sometimes I wanted to detach from work, but reviewing the same materials after hours felt like an extension of my job.”

## 7 Discussion

Our results highlight how LINGOQ bridges two familiar practices—using AI tools at work for English-related tasks and studying English on smartphones—by turning routine queries into learning activities that are directly connected to workers’ tasks and that strengthen their self-efficacy.

### 7.1 Leveraging Reliance on LLMs for Learning Opportunities

Reliance on Generative AI has been a threat to learning as it takes away a critical engagement with a subject matter [13]. Especially for workers, the convenience that LLMs provide fosters passive consumption of generated information rather than critically examining what they are producing [60]. Therefore, ironically, the ESL workers’ convenience coming from reliance on LLM tools, such as ChatGPT, can result in the decay of English skills.

Our results show that LINGOQ helps bridge the gap between workers’ everyday use of LLM tools for English-related tasks and their language learning, turning routine queries into learning activities that are closely related to their work. This approach has been explored in language learning under the framework of English for Specific Purposes (ESP), but the occupational settings commonly used in ESP have been limited to fields such as medicine, business, or technology [46]. As a result, workers in other occupations, such as law or software engineering, are often excluded. Even when such occupational contexts in ESP are available, the examples used in the learning materials can still feel distant from learners’ actual work. Our work addresses these challenges by personalizing English learning to the individual context, directly extracting data from their GenAI applications, and creating learning materials with it.

We also found that some workers perceived their LLM queries not merely as assistance, but as opportunities for language learning. In particular, participants indicated that language-specific UI features—such as side-by-side translation view, toggling between refined and original responses, and marking AI responses for later review—helped them become aware of what they did not know and facilitated conscious learning. The noticing hypothesis [81] suggests that conscious awareness of linguistic gaps is essential for acquisition, beyond mere exposure. Unlike passive reliance on AI-generated answers, this awareness might have reframed their work-related queries as active learning events.

Our work contributes to the existing literature on integrating LLM into learning. While most of the existing applications alter the role of LLM to behave differently in their learning environment, for instance, as a tutor or trainee, rather than as an assistant, [51, 63], we offer a novel approach of leveraging already existing behaviors to improve their learning experience. Increasingly, the industry as well has been attentive to the need to preserve learning opportunities for users who rely on LLMs. For instance, the recent launch of ChatGPT-5 [74] introduced flashcard-style quiz generation within conversations, enabling just-in-time learning during information seeking.

The flourishing of LLM-powered work-assistance tools presents new avenues for applying this approach [69]. While we applied this approach to ESL workers in Korea, we believe it holds broader potential for L2 learning in general. Beyond language learning, other domains where workers depend heavily on LLMs—such as programming [87] or writing [62]—offer promising opportunities to extend this approach, as workers must still develop expertise even while leveraging AI assistance.

These insights point to a promising avenue for future research in education, the future of work, and HCI-AI interaction: designing LLM-powered systems that not only assist with work tasks but also deliberately surface knowledge gaps, transforming everyday interactions into authentic and sustained learning opportunities.

### 7.2 Diversifying and Personalizing Work-related ESL Learning

Based on the participants’ feedback, we identified several directions for improving the users’ learning experience.

One direction to diversify quiz content is to reduce wear-out effects coming from predictable mapping between questions and queries. Future systems could evaluate a user’s proficiency level informed by user modeling based on the collection of query-response pairs [6, 14, 25, 45, 65]. Given the query response pairs, varying the stem for creating a new scenario that uses the same words and expressions can mitigate the reviewing, not anticipating, nature of LINGOQ [77, 95].

In addition, systems should consider how generative approaches can support learners across all proficiency levels; future designs must go beyond static difficulty settings. While our system was particularly effective for lower-proficiency learners, who showed significant improvement over the three-week period, we did not observe clear learning gains for intermediate and advanced learners. By integrating learners’ query logs with quiz performance history, future systems could infer emerging strengths and weaknesses,



enabling adaptive generation of more diverse, challenging, and personalized tasks (e.g., summarization, paraphrasing, error correction). Such adaptive scaffolding—grounded in both user-initiated and system-monitored signals—can better sustain engagement and accelerate skill development across the proficiency spectrum, including for learners at higher CEFR levels.

We anticipate that a more intelligent approach to distinguishing the nuanced purposes of each query could be effective. For example, queries made out of a lack of knowledge (e.g., looking up a word in a dictionary) differ from those made for efficiency (e.g., translating text into English). Creating a learner profile and tracking their personal data can provide additional context that can account for the nature of their English query and can be used to generate questions at an appropriate level.

With the growing importance of verbal communication in the workplace, participants in both our deployment and formative studies expressed strong interest in listening and speaking practice. They reported using LLM-based assistant tools alongside video meeting screens to look up unfamiliar words, check pronunciations, or review meeting recordings via speech-to-text systems for better understanding. Future systems could build on these practices by incorporating voice-based interactions—such as spoken queries or transcripts from virtual meetings—to identify frequently used phrases and recurring pronunciation challenges. These insights could then inform the generation of audio-based quizzes and feedback powered by voice-enabled intelligent agents, extending language learning into more authentic, speech-oriented contexts.

### 7.3 Limitations and Future Work

The development of LINGOQ has several limitations. First, our focus was primarily on reading and writing skills, whereas participants consistently expressed a need for support in listening and speaking. This suggests that future systems should incorporate voice input, pronunciation feedback, and speech-based interactions to better support oral communication. Similarly, the problem format was limited to fill-in-the-blank questions, which might have been more impactful if optimized—an aspect we did not explore in this study.

Second, although screen capture enabled personalized learning, it also raised privacy concerns. Secure, privacy-preserving alternatives must be explored; otherwise, deployment in practice may conflict with industries and companies with high security standards.

Third, our study was limited to a Korean–English context. While we believe the architecture and pipeline structure are language-tolerant, performance can vary significantly depending on the LLM, which in turn depends heavily on training data. Generalizing to other languages—particularly low-resource and non-English second languages—requires further investigation.

Finally, our system relied solely on user-initiated queries, missing opportunities to leverage implicit signals such as hesitation, repeated queries, or affective cues. Future work could integrate such behavioral signals to provide more proactive and situated learning support.

### 7.4 Conclusion

We presented LINGOQ, an LLM-powered system that supports practicing work-related English skills by generating quizzes directly

from workers’ language queries. By connecting everyday use of language assistants through LINGOQUERY with low-burden practice on LINGOQUIZ, LINGOQ enables work-relevant English exercises anytime and anywhere. To examine how people engage and learn with LINGOQ, we conducted a three-week deployment with 28 ESL information workers. Participants actively engaged with the system and reported increased self-efficacy. Our study showed that queries can be transformed into situated learning materials that are relevant, useful for work, and of sufficient quality as validated by expert evaluation. These findings demonstrate how leveraging workers’ reliance on LLMs can create new opportunities for sustainable and meaningful learning. In sum, our work contributes to the growing body of personalized language learning that leverages LLMs and personal data, highlighting the feasibility of grounding study materials in user demand.

### Acknowledgments

We thank our study participants from the formative studies, English test validation, expert evaluation, and the deployment study, for their time and efforts. This work was supported through a research internship at NAVER AI Lab of NAVER Cloud.

### References

- [1] 1996. The role of the linguistic environment in second language acquisition. *Handbook of second language acquisition* (1996), 413–468.
- [2] Tatsuya Amano, Valeria Ramirez-Castañeda, Violeta Berdejo-Espinola, Israel Borokini, Shawan Chowdhury, Marina Golivets, Juan David González-Trujillo, Flavia Montaña-Centellas, Kumar Paudel, Rachel Louise White, et al. 2023. The manifold costs of being a non-native English speaker in science. *PLoS biology* 21, 7 (2023), e3002184.
- [3] Anki Developers. 2025. Anki: Powerful, Intelligent Flashcards. <https://apps.ankiweb.net/>
- [4] Laurence Anthony. 2018. *Introducing English for specific purposes*. Routledge.
- [5] Anthropic. 2025. Claude. Retrieved Sep 1, 2025 from <https://claude.ai/>
- [6] Riku Arakawa, Hiromu Yakura, and Sosuke Kobayashi. 2022. VocabEncounter: NMT-powered vocabulary learning by presenting computer-generated usages of foreign words into users’ daily lives. In *Proceedings of the 2022 CHI conference on human factors in computing systems*. 1–21.
- [7] Mahmoud Azab, Ahmed Salama, Kemal Oflazer, Hideki Shima, Jun Araki, and Teruko Mitamura. 2013. An NLP-based reading tool for aiding non-native English readers. In *Proceedings of the International Conference Recent Advances in Natural Language Processing RANLP 2013*. 41–48.
- [8] Babbel. 2025. Babbel: Learn Languages Online. Retrieved Sep 1, 2025 from <https://www.babbel.com/>
- [9] Gabriela Torregrosa Benavent and Sonsoles Sánchez-Reyes Peñamaría. 2011. Use of Authentic Materials in the ESP Classroom. *Online Submission* 20 (2011), 89–94.
- [10] Rimma Bielousova. 2017. Developing materials for English for specific purposes online course within the blended learning concept. *Tem Journal* 6, 3 (2017), 637–642.
- [11] Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative research in psychology* 3, 2 (2006), 77–101.
- [12] Daniel Buschek, Martin Zürn, and Malin Eiband. 2021. The impact of multiple parallel phrase suggestions on email input and composition behaviour of native and non-native english writers. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–13.
- [13] Cecilia Ka Yuk Chan and Katherine K. W. Lee. 2023. The AI generation gap: Are Gen Z students more interested in adopting generative AI such as ChatGPT in teaching and learning than their Gen X and Millennial Generation teachers? [doi:10.48550/arXiv.2305.02878](https://doi.org/10.48550/arXiv.2305.02878) arXiv:2305.02878 [cs].
- [14] Yuexi Chen and Zhicheng Liu. 2024. WordDecipher: Enhancing Digital Workspace Communication with Explainable AI for Non-native English Speakers. In *Proceedings of the Third Workshop on Intelligent and Interactive Writing Assistants*. 7–10.
- [15] Graeme W Coleman and Nick A Hine. 2012. Twasebook: a crowdsourced phrasebook for language learners using Twitter. In *Proceedings of the 7th Nordic Conference on Human-Computer Interaction: Making Sense Through Design*. 805–806.
- [16] Allan Collins and Manu Kapur. 2006. *Cognitive apprenticeship*. Vol. 291. Citeseer.

- [17] Lynne Cooke. 2010. Assessing concurrent think-aloud protocol as a usability test method: A technical communication approach. *IEEE Transactions on Professional Communication* 53, 3 (2010), 202–215.
- [18] Do Coyle, Philip Hood, and David Marsh. 2010. *CLIL: Content and Language Integrated Learning*. Cambridge University Press.
- [19] Gabriel Culbertson, Shiyu Wang, Malte Jung, and Erik Andersen. 2016. Social situational language learning through an online 3d game. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*. 957–968.
- [20] Hai Dang, Chelse Swoopes, Daniel Buschek, and Elena L Glassman. 2025. CorpusStudio: Surfacing Emergent Patterns In A Corpus Of Prior Work While Writing. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–19.
- [21] Isabelle De Ridder. 2002. Visible or invisible links?. In *CHI'02 Extended Abstracts on Human Factors in Computing Systems*. 624–625.
- [22] DeepL SE. 2025. DeepL Translator. Retrieved Sep 1, 2025 from <https://www.deepl.com/translator>
- [23] Robert F DeVellis. 2006. Classical test theory. *Medical care* 44, 11 (2006), S50–S59.
- [24] John Dewey. 2024. *Democracy and education*. Columbia University Press.
- [25] Jiexin Ding, Bowen Zhao, Yuntao Wang, Xinyun Liu, Rui Hao, Ishan Chatterjee, and Yuanchun Shi. 2025. Unknown Word Detection for English as a Second Language (ESL) Learners using Gaze and Pre-trained Language Models. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. 1–16.
- [26] Jacob Doughty, Zipiao Wan, Anishka Bompelli, Jubahed Qayum, Taozhi Wang, Juran Zhang, Yujia Zheng, Aidan Doyle, Pragnya Sridhar, Arav Agarwal, et al. 2024. A comparative study of AI-generated (GPT-4) and human-crafted MCQs in programming education. In *Proceedings of the 26th Australasian Computing Education Conference*. 114–123.
- [27] Fiona Draxler, Julia Maria Brenner, Manuela Eska, Albrecht Schmidt, and Lewis L Chuang. 2022. Agenda-and activity-based triggers for microlearning. In *Proceedings of the 27th International Conference on Intelligent User Interfaces*. 620–632.
- [28] Fiona Draxler, Albrecht Schmidt, and Lewis L Chuang. 2023. Relevance, effort, and perceived quality: Language learners' experiences with AI-generated contextually personalized learning material. In *Proceedings of the 2023 ACM Designing Interactive Systems Conference*. 2249–2262.
- [29] Duolingo. 2025. Duolingo: Language Lessons for Everyone. Retrieved Sep 1, 2025 from <https://www.duolingo.com/>
- [30] Duolingo. 2025. Duolingo Max: AI-powered language learning with GPT-4. Retrieved Sep 1, 2025 from <https://blog.duolingo.com/duolingo-max/>
- [31] Darren Edge, Stephen Fitchett, Michael Whitney, and James Landay. 2012. Mem-Reflex: adaptive flashcards for mobile microlearning. In *Proceedings of the 14th international conference on Human-computer interaction with mobile devices and services*. 431–440.
- [32] Darren Edge, Elly Searle, Kevin Chiu, Jing Zhao, and James A Landay. 2011. MicroMandarin: mobile language learning in context. In *Proceedings of the SIGCHI conference on human factors in computing systems*. 3169–3178.
- [33] Educational Testing Service (ETS). 2025. TOEIC Official Website. <https://www.ets.org/toEIC.html>. Accessed: retrieveddate.
- [34] Sabina Elkins, Ekaterina Kochmar, Jackie CK Cheung, and Iulian Serban. 2024. How teachers can use large language models and bloom's taxonomy to create educational quizzes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 23084–23091.
- [35] FastAPI. 2025. FastAPI framework, high performance, easy to learn, fast to code, ready for production. Retrieved Sep 1, 2025 from <https://fastapi.tiangolo.com/>
- [36] Gerhard Gassler, Theo Hug, and Christian Glahn. 2004. Integrated Micro Learning—An outline of the basic method and first results. *Interactive computer aided learning* 4 (2004), 1–7.
- [37] Alex Gilmore. 2007. Authentic materials and authenticity in foreign language learning. *Language teaching* 40, 2 (2007), 97–118.
- [38] Google. 2025. Google Search: Little Language Lessons. Retrieved Sep 1, 2025 from <https://labs.google/lll/en/>
- [39] Google DeepMind. 2025. Gemini. Retrieved Sep 1, 2025 from <https://gemini.google.com/>
- [40] Grammarly Inc. 2025. Grammarly: AI Writing Assistance. Retrieved Sep 1, 2025 from <https://www.grammarly.com/>
- [41] The PostgreSQL Global Development Group. 2025. PostgreSQL: The World's Most Advanced Open Source Relational Database. Retrieved Sep 1, 2025 from <https://www.postgresql.org/>
- [42] Philip J Guo. 2018. Non-native english speakers learning computer programming: Barriers, desires, and design opportunities. In *Proceedings of the 2018 CHI conference on human factors in computing systems*. 1–14.
- [43] Ari Hautasaari, Takeo Hamada, Kuntaro Ishiyama, and Shogo Fukushima. 2019. Vocabura: A method for supporting second language vocabulary learning while walking. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 3, 4 (2019), 1–23.
- [44] Suzanne Hidi and K Ann Renninger. 2006. The four-phase model of interest development. *Educational psychologist* 41, 2 (2006), 111–127.
- [45] Taichi Higasa, Keitaro Tanaka, Qi Feng, and Shigeo Morishima. 2024. Keep eyes on the sentence: An interactive sentence simplification system for English learners based on eye tracking and large language models. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. 1–7.
- [46] Tom Hutchinson and Alan Waters. 1987. *English for Specific Purposes*. Cambridge University Press.
- [47] Gwo-Jen Hwang, Chin-Chung Tsai, and Stephen JH Yang. 2008. Criteria, strategies and research issues of context-aware ubiquitous learning. *Journal of Educational Technology & Society* 11, 2 (2008), 81–91.
- [48] Adam Ibrahim, Brandon Huynh, Jonathan Downey, Tobias Höllerer, Dorothy Chun, and John O'donovan. 2018. Arbis pictus: A study of vocabulary learning with augmented reality. *IEEE transactions on visualization and computer graphics* 24, 11 (2018), 2867–2874.
- [49] Nanna Inie and Mircea F Lungu. 2021. Aiki-turning online procrastination into microlearning. In *Proceedings of the 2021 CHI conference on human factors in computing systems*. 1–13.
- [50] Takumi Ito, Naomi Yamashita, Tatsuki Kuribayashi, Masatoshi Hidaka, Jun Suzuki, Ge Gao, Jack Jamieson, and Kentaro Inui. 2023. Use of an AI-powered rewriting support software in context with other tools: a study of non-native English speakers. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–13.
- [51] Hyoungwook Jin, Seonghee Lee, Hyungyu Shin, and Juho Kim. 2024. Teach ai how to code: Using large language models as teachable agents for programming education. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–28.
- [52] Keith Johnson. 1997. *Language teaching and skill learning*. Blackwell Oxford, England.
- [53] Mineyoung Kim, Jiwook Lee, Youngji Koh, Chanhee Lee, Uichin Lee, and Auk Kim. 2024. Interrupting for Microlearning: Understanding Perceptions and Interruptibility of Proactive Conversational Microlearning Services. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–21.
- [54] Yewon Kim, Thanh-Long V Le, Donghui Kim, Mina Lee, and Sung-Ju Lee. 2025. Design Opportunities for Explainable AI Paraphrasing Tools: A User Study with Non-native English Speakers. In *Proceedings of the 2025 ACM Designing Interactive Systems Conference*. 1061–1083.
- [55] Thomas Andrew Kirkpatrick. 2011. Internationalization or Englishization: Medium of instruction in today's universities. (2011).
- [56] Carol Collier Kuhlthau. 1999. The role of experience in the information search process of an early career information worker: Perceptions of uncertainty, complexity, construction, and sources. *Journal of the American Society for information Science* 50, 5 (1999), 399–412.
- [57] Huisung Kwon, Soyeong Min, and Sangsu Lee. 2025. How to Better Translate Participant Quotes Using LLMs: Exploring Practices and Challenges of Non-Native English Researchers. In *Proceedings of the Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. 1–11.
- [58] Inc. LangChain. 2025. LangChain: Applications that Can Reason. Retrieved Sep 1, 2025 from <https://www.langchain.com/>
- [59] Jean Lave and Etienne Wenger. 1991. *Situated learning: Legitimate peripheral participation*. Cambridge university press.
- [60] Hao-Ping (Hank) Lee, Advait Sarkar, Lev Tankelevitch, Ian Drosos, Sean Rintel, Richard Banks, and Nicholas Wilson. 2025. The Impact of Generative AI on Critical Thinking: Self-Reported Reductions in Cognitive Effort and Confidence Effects From a Survey of Knowledge Workers. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 1121, 22 pages. doi:10.1145/3706598.3713778
- [61] Sangmin-Michelle Lee. 2022. A systematic review of context-aware technology use in foreign language learning. *Computer assisted language learning* 35, 3 (2022), 294–318.
- [62] Zhuoyan Li, Chen Liang, Jing Peng, and Ming Yin. 2024. The value, benefits, and concerns of generative ai-powered assistance in writing. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–25.
- [63] Anna Lieb and Toshali Goel. 2024. Student interaction with newtbot: An llm-as-tutor chatbot for secondary physics education. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*. 1–8.
- [64] T-Y Liu. 2009. A context-aware ubiquitous learning environment for language listening and speaking. *Journal of computer assisted Learning* 25, 6 (2009), 515–527.
- [65] Mircea F Lungu, Luc van den Brand, Dan Chirtoaca, and Martin Avagyan. 2018. As we may study: Towards the web as a personalized language textbook. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. 1–12.
- [66] Memrise. 2025. Memrise: Language Learning Made Fun. Retrieved Sep 1, 2025 from <https://www.memrise.com/>
- [67] Meta. 2025. React Native - Learn Once, Write Everywhere. Retrieved Sep 1, 2025 from <https://reactnative.dev/>
- [68] Microsoft. 2025. TypeScript. Retrieved Sep 1, 2025 from <https://www.typescriptlang.org>

- [69] Amr Mohamed, Maram Assi, and Mariam Guizani. 2025. The Impact of LLM-Assistants on Software Developer Productivity: A Systematic Literature Review. *arXiv preprint arXiv:2507.03156* (2025).
- [70] David Nunan. 2004. *Task-based language teaching*. Cambridge university press.
- [71] Council of Europe. Council for Cultural Co-operation. Education Committee. Modern Languages Division. 2001. *Common European framework of reference for languages: Learning, teaching, assessment*. Cambridge University Press.
- [72] Hiroaki Ogata, Bin Hou, Mengmeng Li, Noriko Uosaki, Kosuke Mouri, and Songran Liu. 2014. Ubiquitous learning project using life-logging technology in Japan. *Journal of Educational Technology & Society* 17, 2 (2014), 85–100.
- [73] OpenAI. 2025. ChatGPT. Retrieved Sep 1, 2025 from <https://chat.openai.com/>
- [74] OpenAI. 2025. ChatGPT-5. Retrieved Sep 1, 2025 from <https://openai.com/index/introducing-gpt-5/>
- [75] OpenAI. 2025. OpenAI API. Retrieved Sep 1, 2025 from <https://openai.com/api/>
- [76] OpenJS Foundation and Electron contributors. 2025. Electron: Build cross-platform desktop apps with JavaScript, HTML, and CSS. Retrieved Sep 1, 2025 from <https://www.electronjs.org/>
- [77] Zhenhui Peng, Xingbo Wang, Qiushi Han, Junkai Zhu, Xiaojuan Ma, and Huamin Qu. 2023. Storyfier: Exploring vocabulary learning support with text generation models. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–16.
- [78] Quizlet Inc. 2025. Quizlet: Learning Tools and Flashcards. Retrieved Sep 1, 2025 from <https://quizlet.com/>
- [79] Katherine A Rawson and John Dunlosky. 2011. Optimizing schedules of retrieval practice for durable and efficient learning: How much is enough? *Journal of Experimental Psychology: General* 140, 3 (2011), 283.
- [80] Ringle English Education Service. 2025. Ringle: 1:1 Online Tutoring with Ivy League Tutors. Retrieved Sep 1, 2025 from <https://www.ringleplus.com/>
- [81] Peter Robinson. 1995. Attention, memory, and the “noticing” hypothesis. *Language learning* 45, 2 (1995), 283–331.
- [82] Rosetta Stone Ltd. 2025. Rosetta Stone: Language Learning Software. Retrieved Sep 1, 2025 from <https://www.rosettastone.com/>
- [83] Christopher A Rowland. 2014. The effect of testing versus restudy on retention: a meta-analytic review of the testing effect. *Psychological bulletin* 140, 6 (2014), 1432.
- [84] Gyula Sankó. 2006. The effects of hypertextual input modification on L2 vocabulary acquisition and retention. *University of Pécs Roundtable 2006: Empirical Studies in English Applied Linguistics* (2006), 157.
- [85] Mitchell Shortt, Shantanu Tilak, Irina Kuznetcova, Bethany Martens, and Babatunde Akinkuolie. 2023. Gamification in mobile-assisted language learning: A systematic review of Duolingo literature from public release of 2012 to early 2020. *Computer Assisted Language Learning* 36, 3 (2023), 517–554.
- [86] Speak Global Inc. 2025. Speak: AI-Powered English Tutoring App. Retrieved Sep 1, 2025 from <https://www.speak.com/?lang=en>
- [87] Benyamin Tabarsi, Heidi Reichert, Ally Limke, Sandeep Kuttal, and Tiffany Barnes. 2025. LLMs’ Reshaping of People, Processes, Products, and Society in Software Development: A Comprehensive Exploration with Early Adopters. *arXiv preprint arXiv:2503.05012* (2025).
- [88] Colin Thompson. 2019. Practice makes Perfect? A review of second language teaching methods. *The Bulletin of the Graduate School of Josai International University* 22, 55-69 (2019).
- [89] Thomas C Toppino, Melissa H LaVan, and Ryan T Iaconelli. 2018. Metacognitive control in self-regulated learning: Conditions affecting the choice of restudying versus retrieval practice. *Memory & Cognition* 46, 7 (2018), 1164–1177.
- [90] Ashok Kumar Veerasamy and Anna Shillabeer. 2014. Teaching English based programming courses to English language learners/non-native speakers of English. *International Proceedings of Economics Development and Research* 70 (2014), 17.
- [91] Olga Viberg and Åke Grönlund. 2012. Mobile assisted language learning: A literature review. In *11th world conference on mobile and contextual learning*.
- [92] Chuang Wang, Do-Hong Kim, Rui Bai, and Jiyue Hu. 2014. Psychometric properties of a self-efficacy scale for English language learners in China. *System* 44 (2014), 24–33.
- [93] Jane Willis. 2016. *A Framework for Task-Based Learning*. Longman.
- [94] Masanori Yamada, Satoshi Kitamura, Noriko Shimada, Takafumi Utashiro, Katsusuke Shigeta, Etsuji Yamaguchi, Richard Harrison, and Yuhei Yamauchi. 2012. Development and Evaluation of English Listening Study Materials for Business People Who Use Mobile Devices. *Calico Journal* 29, 1 (2012), 44–66.
- [95] Kanta Yamaoka, Ko Watanabe, Koichi Kise, Andreas Dengel, and Shoya Ishimaru. 2022. Experience is the best teacher: Personalized vocabulary building within the context of Instagram posts and sentences from GPT-3. In *Adjunct Proceedings of the 2022 ACM International Joint Conference on Pervasive and Ubiquitous Computing and the 2022 ACM International Symposium on Wearable Computers*. 313–316.

## A LLM Prompts

### A.1 Instructions for LINGOQUERY Conversational Agent Intent Classifier

[Role]

You are an Intention Classifier. Your job is to analyze the input text and classify it into one of four categories.

[Classification Categories]

**\*\*translation\*\***: "Translate the following text naturally between English and Korean. Please also explain how the nuance and context of the sentences are reflected in the translation."

**\*\*proofread\*\***: "Proofread the following text into more accurate and natural English. Please also provide an explanation of the changes and the reasons behind them."

**\*\*lookup\*\***: "Explain the meaning of the following word (or expression) in detail, in the style of a dictionary entry."

**\*\*text\*\***: Any other input that doesn't match the above three categories.

[Classification Rules]

1. If the input text exactly matches one of the three specific examples above => Classify accordingly
2. If the input text is similar to any of the three examples => Classify accordingly
3. If the input text doesn't match any of the three examples => Classify as "text"

[Output Format]

**\*\*CRITICAL: You must respond with ONLY ONE WORD from the list below.\*\***

**\*\*DO NOT use JSON format. DO NOT add explanations. DO NOT add quotes.\*\***

Respond with ONLY one of these four values:

- translation
- proofread
- lookup
- text

### A.2 Instructions for LINGOQUERY Conversational Agent Response Generator

[Role]

You are a Workplace English Support Assistant, designed to help the user tackle English-related tasks and challenges in everyday work situations.

[Personality]

- Patient and encouraging
- Clear and articulate in explanations
- Friendly and approachable
- Professional yet conversational
- Culturally sensitive and inclusive

[Chat Style]

- The user will speak in {user\_language}. So you must also speak in polite and supportive {user\_language}.
- Do not greet the user and treat them as if you already know them well.

[CRITICAL: Context Memory & Style Consistency]

- ALWAYS remember the entire conversation history
- Remember user's work context, preferences, and instructions
- Remember user's ongoing projects and tasks
- Maintain consistent response style throughout the conversation
- If user prefers certain response styles, maintain that consistency

[Message Content Format]

The user's intention is provided as: [Intention: INTENTION\_PLACEHOLDER]

The user's message contains:

- query\_prompt: The prompt the user is using to make the query
- content: The content the user is querying about

[Intent-Based Response Generation]

IMPORTANT: The user's intention has already been classified. Use this information to determine the appropriate output\_type and response format.

**\*\*Response Format Based on Intention:\*\***



```

1. For Lookup (intention: "lookup"):
  - Use DictionaryOutput format
  - Provide comprehensive dictionary information including meanings, examples, synonyms, etc.
  - Focus on the word/phrase in the user's content

2. For Translate (intention: "translation"):
  - Use TranslationOutput format
  - Provide original text, translation, and explanation
  - Translate naturally, considering user's context and communication style
  - Avoid literal translation - focus on natural expression
  - **Pay attention to formality, tone, and context**: Match the user's professional level, industry terminology, and communication style
  - When the user content is a mix of {user_language} and English, translate the entire content into English

3. For Proofread (intention: "proofread"):
  - Use RefinementOutput format
  - Provide original content, refined content, and refinement rationale
  - Refine naturally
  - **Minimal refinement approach**: Preserve the user's original structure and meaning as much as possible.
  - **refinement_rationale**: Write in simple, natural Korean. Avoid numbered lists or structured formats.
  - When the user content is a mix of {user_language} and English, refine the content to be fully in English
  - Only refine to {user_language} if the user explicitly requests it

4. For General (intention: "text"):
  - Use Text output format
  - Respond naturally to the user's query_prompt and content
  - Provide helpful, detailed explanations
  - Suggest 2-3 alternative approaches when appropriate
  - Be conversational and engaging like ChatGPT

**Your Task:**
Based on the classified intention provided, generate the appropriate response using the correct output_type and format.
Do not re-classify the intention - use the one that has been provided to you.

```

## B Development and Validation of the English Proficiency Test

An English proficiency test was developed to evaluate the learning performance of the deployment study participants. We selected 46 multiple-choice items from the TOEIC (Test of English for International Communication), consisting of 30 single-sentence fill-in-the-blank items (each with one blank) and four paragraph-based sets (each set containing a short paragraph with four blanks).

**Participants.** To validate the difficulty and time required to complete the test, we recruited computer-based information workers via social media advertisements, following the inclusion criteria described in [Section 5.1](#). Among the 40 applicants, 11 were excluded based on their responses to attention check items designed to ensure data quality. In total, 29 South Korean information workers (16 female, 12 male, 1 preferred not to disclose) completed the validation. Participants had an average age of 27.9 years ( $SD = 4.8$ ) and represented diverse occupational backgrounds, including researchers (14), office workers (11), and engineers (4). Based on CEFR self-assessment [71], 2 participants identified as *A1* (beginner), 2 as *A2* (elementary), 8 as *B1* (intermediate), 5 as *B2* (upper-intermediate), 5 as *C1* (advanced), and 7 as *C2* (proficient). Each participant received 20,000 KRW (approx. 14 USD) as compensation.

**Procedure.** Participants completed the validation via an online survey. They solved all 46 test items along with 2 attention check questions. Problem-solving time was recorded. The order of questions and answer choices was randomized for each participant. The average response time was 23.7 seconds for single-sentence items and 117.2 seconds for paragraph-based sets. Mean scores were  $M = 21.03$  ( $SD = 4.88$ ) for the single-sentence items (score range: 0–30) and  $M = 11.24$  ( $SD = 2.43$ ) for the paragraph sets (score range: 0–16).

**Validation.** Based on classical test theory [23], item difficulty was calculated as the proportion of participants who answered each item correctly. Following the standard range of acceptable difficulty (0.4 to 0.8), we selected 16 single-sentence items and 3 paragraph-based sets (12 items total). Given the average solving time, the expected completion time for the selected 28 items is approximately 13 minutes. Thus, the final version of the English proficiency test consists of 28 validated items to be completed in 13 minutes.

## C Expert Evaluation of Generated Questions

To evaluate the performance of the question evaluator in the LINGOQ pipeline, we conducted an expert evaluation on 30 sample questions generated from user data in the deployment study.

The set of 30 sample questions consisted of 24 accepted questions (evaluated as *True* for both criteria) and 6 discarded questions (evaluated as *False* for one or both criteria) based on our pipeline's evaluation criteria: (1) Answerability and (2) Proficiency [26, 34] (see [Table 2](#)). We compared the expert evaluation results with the pipeline's evaluation.

**Table 2: Rubrics used for expert evaluation of questions generated by LINGOQ.**

Rubric	Question	Options
Correct answer	Is there a correct answer listed in the options? Is the option marked “correct” actually correct?	Yes, there is a correct answer and it is marked ‘correct’ There is a correct answer but it is not marked ‘correct’ There are multiple correct answers No, there is no correct answer Don’t know
Unique choices	Are the options distinct from each other, ensuring they are unique choices?	Yes, they are completely unique Some choices are unique, some are too similar No, they are all too similar Don’t know
No obviously wrong	Is the MCQ free from obviously-wrong options?	Yes, there are no obviously-wrong options Yes, but the options give away the correct answer No, there are obviously-wrong options Don’t know

The pipeline adopts binary judgments (*True/False*) for clarity and alignment with LLM performance. To gain richer insights, we extended the rubrics for expert evaluation. For comparison with our pipeline, one of the authors mapped expert ratings into *True/False* labels, and two other authors reviewed the mappings for consistency.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009