# SummAct: Uncovering User Intentions Through Interactive Behaviour Summarisation

GUANHUA ZHANG, University of Stuttgart, Germany

MOHAMED AHMED, University of Stuttgart, Germany

ZHIMING HU, University of Stuttgart, Germany

ANDREAS BULLING, University of Stuttgart, Germany

Recent work has highlighted the potential of modelling interactive behaviour analogously to natural language. We propose *interactive behaviour summarisation* as a novel computational task and demonstrate its usefulness for automatically uncovering latent user intentions while interacting with graphical user interfaces. To tackle this task, we introduce SummAct – a novel hierarchical method to summarise low-level input actions into high-level intentions. SummAct first identifies sub-goals from user actions using a large language model and in-context learning. High-level intentions are then obtained by fine-tuning the model using a novel UI element attention to preserve detailed context information embedded within UI elements during summarisation. Through a series of evaluations, we demonstrate that SummAct significantly outperforms baselines across desktop and mobile interfaces as well as interactive tasks by up to 21.9%. We further show three exciting interactive applications benefted from SummAct: interactive behaviour forecasting, automatic behaviour synonym identification, and language-based behaviour retrieval.

Additional Key Words and Phrases: Interactive Behaviour, Intention recognition, Large language model, Next action prediction, Retrieval

## 1 INTRODUCTION

Recent work has demonstrated that interactive behaviour, e.g. when interacting with graphical user interfaces using the mouse or keyboard, shares similarities with the sequential and hierarchical nature of natural language [53]. In parallel, an increasing number of works have started to model interactive behaviour as natural language and process it using language models [12, 22, 30, 32, 50]. One key advantage of this language perspective is facilitating a more interpretable analysis and understanding of interactive behaviour, thus enabling novel paradigms for solving human-computer interaction (HCI) tasks.

Among these tasks, understanding users' intentions is key to intelligent interactive systems and anticipatory user interfaces [21, 53]. Recognising intentions based on the user's behaviour history has been widely studied and applied in HCI, including for unintentional error detection [1], next action prediction [4], or task automation [23, 50, 58]. Despite its potential for HCI and promising first results, predicting users' intentions from their interactive behaviour remains challenging, partly due to human behaviour's high variability and complexity. Existing works, therefore, typically assume a pre-defined and fixed set of intentions and treat intention recognition as a classification task. However, this approach neither captures the wide variety of user intentions in everyday scenarios nor can robustly adapt to unseen or context-dependent intentions [52]. It can also result in misinterpretations when users' needs do not align with predefined intention categories, which often happens in real-world applications [54].

In this work, we take inspiration from text and image summarisation tasks studied in natural language processing and computer vision. These tasks involve summarising long text or complex videos into a concise sentence description. Similarly, we formulate intention recognition as an interactive behaviour summarisation task: human interactive

behaviour is to be summarised into a sentence, i.e., a natural language description of users' underlying interactive intentions. In contrast to existing methods, interactive behaviour summarisation enables recognising an open-ended set of intentions. It allows for capturing more flexible and varied interaction intentions and handling unseen intentions.

To address this task, we propose SummAct – a novel large language model (LLM)-based method that uses a hierarchical summarisation process: the method initially summarises *low-level* actions into *mid-level* sub-goals, then uses them to augment the input, and finally summarise into a *high-level* intention. In this work, we focus on the interactive behaviour at the user interface (UI) element level, meaning that each input action sample consists of the interacted UI element and the operation the user conducts on this element (e.g., click or select). The UI element information includes its category (e.g., button or combo box), inherent (e.g., the name or the visible text on a button) and additional content (further values that users are interested and pick, e.g., a value selected from a combo box). On these input actions, SummAct first generates sub-goals using in-context learning via a pretrained, frozen LLM due to the lack of ground-truth annotations and then fine-tunes the LLM to produce the final summary. During fine-tuning, we further propose a UI element attention mechanism that assigns higher weights to the UI element contents, thereby preserving the detailed context information embedded within these elements. This is crucial for accurately interpreting intentions that exhibit subtle differences and for further applications like behaviour forecasting.

We evaluate SummAct on two datasets that cover a web (Mind2Web [12]) and a mobile (MoTIF [6]) interaction setting. We show that SummAct can accurately uncover the intentions underlying user actions, with a sentence embedding cosine similarity up to 0.842 compared to the ground-truth intentions. We also demonstrate the importance of our design choices with the full SummAct model significantly outperforming ablated versions by up to 21.9% in cosine similarity. We finally showcase three exciting applications enabled by interactive behaviour summarisation: providing contextual information of user intentions to enhance behaviour forecasting for proactive user interfaces; automatically identifying behaviour synonyms to understand user preferences, interaction strategies, system usability and common design patterns; and unlocking language-based behaviour retrieval that lays the foundation for building behaviour-related conversational agents.

In summary, the specific contributions of our work are three-fold:

- We formulate intention recognition as the novel open-ended task of summarising interactive behaviour into natural language descriptions. This formulation overcomes existing limitations associated with pre-defined intention sets and improves generalisability to unseen intentions. Towards this task, we propose an LLM-based method, SummAct[1] incorporating two distinct and novel designs – hierarchical summarisation and UI element attention.
- We show the effectiveness of these designs, and SummAct in general, for interactive behaviour summarisation, through a series of evaluations on two datasets covering both web and mobile interaction settings.
- We demonstrate the potential of interactive behaviour summarisation via three example applications: interactive behaviour forecasting, identifying behaviour synonyms, and language-based behaviour retrieval. These are widely relevant in HCI, particularly for developing intelligent interactive systems or UI optimisation.

## 2  RELATED WORK

We discuss related work on (1) understanding user intentions behind interactive behaviour, (2) large language models for interactive behaviour modelling, and (3) summarising non-language data.

---

[1] We will release our source code upon acceptance.

### 2.1 Understanding User Intentions behind Interactive Behaviour

Recognising users' intentions from their interactive behaviour is key to intelligent user interfaces, which can proactively support users by automatically adjusting the UIs or providing action recommendations [17, 54]. Therefore, an increasing number of works in HCI field have studied automatic intention prediction from interactive behaviour. For example, in virtual reality (VR) interactive environments, David-John et al. [11] recognised user intentions of selecting an item from gaze behaviour. Hu et al. [21] built a model based on convolutional neural networks and bidirectional gated recurrent units to recognise interactive tasks the users aimed to achieve from their eye and head movements in VR. In more pervasive, daily scenarios such as interacting with personal computers or mobile devices, researchers have also developed various methods to recognise the objectives users aim to achieve via their actions. These actions were captured through different modalities, including mouse movements and clicks, keyboard typing, eye tracking and touch interactions. For instance, Elbahi et al. [14] used hidden Markov model and conditional random field to recognise which e-learning task the users were performing from mouse movement. Koldijk et al. [26] developed classifiers based on different machine learning models including naive Bayes, KStar, decision tree and multilayer perceptron (MLP) to recognise which one out of 12 office tasks the user aimed to complete from their mouse and keyboard actions. Zhang et al. [54] proposed a multimodal random forest-based approach to recognise intentions from users' mouse, keyboard and gaze actions in a text editing task that included seven text formatting intentions pre-defined by the authors. In the mobile interaction settings, Xu et al. [51] identified intentional versus unintentional touches from gaze, head and screen touch behaviour. They built the model based on logistic regression, naïve Bayes, k-nearest neighbour, random forest, gradient boosting and MLP.

However, all the above works studied a pre-defined, fixed, and closed set of user intentions, which inherently limits the adaptability and scalability of systems to new intentions. Inspired by the prior finding that interactive behaviour shares a similar sequential and hierarchical structure with natural language [53], in this paper, we formulate intention recognition as summarising interactive behaviour into a sentence. This attempt accommodates an open-ended set of user intentions and thus allows for more flexible and comprehensive interactive behaviour modelling.

### 2.2 Large Language Models for Interactive Behaviour Modelling

LLMs have recently achieved ground-breaking success in HCI research, bringing novel insights and methodology to model user actions for different applications. For example, Liu et al. [36] proposed HintDroid, an LLM-based method using in-context learning to generate hint-text in Android applications based on the user's input and its corresponding UI context. Wang et al. [46] used a pretrained LLM to investigate the conversational interactions with mobile user interfaces via prompt engineering and zero-shot learning. Their results demonstrated the potential of using LLMs for language-based mobile interactions. Huang et al. [22] applied pretrained LLMs and a chain-of-thought technique to extracting macros from mobile interaction traces in existing datasets. Other research focuses on building LLM-based automatic agents that navigate through interactive systems and complete pre-defined tasks [9, 32, 50, 59]. For instance, Deng et al. proposed MindAct [12] to automatically perform given tasks in complex web environments. MindAct first fine-tuned a language model to rank all the UI elements available on the web page based on the task description and action history. Selecting the top-ranked as the candidates, MindAct then formulated task automation as a multi-choice question-answering task and used in-context learning for task automation.

Despite the acknowledged potential of LLMs in interactive behaviour modelling, their application in intention recognition remains largely under-explored. In this paper, we approach intention recognition through an interactive behaviour summarisation task and propose an LLM-based method, SummAct, to address this task.

### 2.3 Summarising Non-Language Data

In language processing, summarisation has been widely studied and applied in condensing large amounts of information into concise sentences, enabling efficient content consumption across various domains, such as documents [13, 35], code [19], and speech [44]. Inspired by the success of these works, researchers have started to summarise non-language data also into coherent textual descriptions for quick and accessible understanding of complex data [46]. For instance, image captioning enhanced the comprehension of the images and meanwhile enables text-based image retrieval [45]. Kawamura et al. [25] proposed a multimodal method to summarise lecture videos using audio transcripts, on-screen images and texts enabling users to effectively obtain information from lengthy video content. Lin et al. [34] and Chen et al. [7] summarised human motion videos, which not only enhanced the understanding of the motion sequence, but also had the potential to allow controllable text-to-motion generation. Chen et al. [8] annotated natural-language explanations for fixations in scanpaths, providing insights into implicit gaze behaviour change and benefited explainable scanpath prediction. In HCI, researchers have studied the summarisation of graphical user interfaces. Wang et al. [47] summarised core information of mobile UI screens into natural language via the proposed multimodal method, Scree2Words, integrating the text, image, structures and UI semantics. They showcased that the summarisation could potentially be used for language-based UI retrieval, enhancing screen readers and screen indexing for conversational mobile interactions.

These applications highlight the transformative potential of applying summarisation techniques to diverse data modalities. Building upon this, our work introduces a novel method to summarise the complex interactive behaviour into human-interpretable natural language sentences, which reflect users' latent intentions. Additionally, we present three examples benefited from interactive behaviour summarisation, behaviour forecasting, behaviour synonym identification and language-based behaviour retrieval, which are widely relevant for intelligent interactive systems and user interface optimisation.

## 3 INTERACTIVE BEHAVIOUR SUMMARISATION USING SUMMACT

Building on recent advances demonstrating the potential of analysing interactive behaviour similarly to natural language [22, 32, 53], our method SummAct addresses the novel task of interactive behaviour summarisation. The input of SummAct is a sequence of UI element-level interaction actions caused by the user interacting with a graphical user interface. Every action consists of the UI element the user has interacted with and the operation performed on this element (e.g., click or select). The UI element contains the information of its category (e.g., button or combo box), the inherent (e.g., its name or the text on it) and additional content (the specific value the user is interested in, e.g., the value selected from a combo box). The output of our method is a natural language sentence that concisely summarises their interactive behaviour and, as we show here, their latent interaction intentions. In stark contrast to existing methods for the classification of intentions, which are limited to a closed and predefined set of possible intentions [14, 15, 55], interactive behaviour summarisation allows us to recognise an open-ended set of intentions, including intentions not seen during training. As such, SummAct can provide a more comprehensive and generalisable understanding of interactive behaviour.
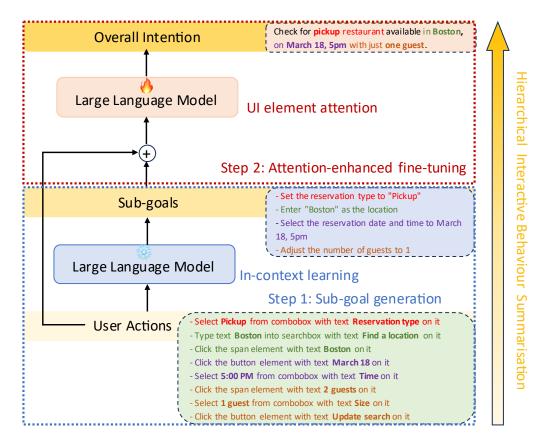
Fig. 1. Overview of SummAct for uncovering user intentions during user interface interactions through interactive behaviour summarisation. SummAct employs a hierarchical process that initially generates sub-goals and produces the overall intention in natural language. The input is a sequence of user actions, including the interacted UI element and the user's operation on this element. SummAct uses in-context learning to infer an arbitrary number of sub-goals using a pretrained, frozen LLM (Step 1) and then fine-tunes the LLM while introducing a UI element attention (Step 2) to keep detailed context embedded in UI element contents, as highlighted in **bold**. Actions in the same colour are summarised into the same sub-goal and then to a phrase in the overall intention. The summary of the output reflects the latent intentions that underlie these actions.

Figure 1 provides an overview of SummAct's hierarchical approach to interactive behaviour summarisation: given a sequence of user input actions encoded in natural language descriptions, SummAct first summarises these low-level actions into a set of sub-goals. Due to the absence of ground-truth data, we use expert annotations in combination with in-context learning to adapt a pretrained, frozen LLM to generate sub-goals. In the second step, these sub-goals are combined with the original actions and summarised into high-level intentions via fine-tuning the LLM. To preserve UI element content (indicated in **bold** in Figure 1) in the summary, our method uses a novel UI element attention during fine-tuning. Two previous findings inspire this hierarchical approach: hierarchical modelling of language data can robustly handle extensive and complex input, such as long documents [35]; and interactive behaviour has an inherent hierarchical nature similar to that observed in natural language [53]. In the following, we describe each of these steps in more detail.

## 3.1 Sub-Goal Generation

The first step involves generating sub-goals from the low-level, individual input actions. As shown in Figure 1, actions marked in the same colour are summarised into the same sub-goal, which later becomes a phrase integrated into the overall intention. Given the lack of HCI datasets offering annotations of interaction sub-goals, we used in-context learning. In-context learning involves giving an LLM a small set of examples presented within the context (the prompt) at inference time to guide its response [49]. This approach leverages LLM's ability to understand and adapt to patterns presented in the immediate context of the query without the need to fine-tune the model. To obtain these examples, we asked three HCI, GUI, and behaviour modelling experts to annotate the sub-goals on five samples from the training set [31]. These samples are five different action sequences completing five different tasks. In Appendix B.1, we provide the used prompt, including the example of sub-goal annotation.

## 3.2 Attention-Enhanced Fine-Tuning

In the second step, we fine-tune the LLM to summarise the overall intention from the generated mid-level sub-goals and the original low-level actions. In Appendix B.2, we provide sample prompts used in this fine-tuning step. LLMs are typically structured as sequence-to-sequence models, i.e. they are trained to generate output sequences based on input sequences, such as summarising an input. Therefore, LLMs are commonly trained using a next token prediction task in a teacher-forcing setup, where the model is guided by a ground-truth token rather than the previously predicted token to predict the next token [43]. This training strategy helps stabilise the training process and accelerates convergence by reducing the propagation of errors through the sequence [46]. Thus, based on the input prompt, LLMs iteratively predict the next token and continually update their predictions as each new token is added to the output sequence. Next token prediction is formulated as a classification task and thus uses a cross-entropy loss $L_{NextToken}$ [28]. For the $j$-th token in the input sequence, this loss is calculated as

$$L_{NextToken_j} = log(\mathbb{P}_\theta(Token_j | Token_1, ..., Token_{j-1})) \tag{1}$$

where $\theta$ are the model parameters.

In preliminary experiments, we found that fine-tuning the LLM only using next token prediction led to UI element content is getting excluded from the final interactive behaviour summary. This may be because the model tends to rely on frequent patterns of general natural language rather than focus on task-specific information embedded in the UI elements [2]. Figure 1 shows examples of such information related to the interactive behaviour summarisation task (highlighted in **bold**). Let us consider the first input action ("Select **Pickup** from combobox with text **Reservation type** on it") as an example: "Reservation type" is the name of the combo box, i.e., the inherent content of the combo box, representing what this combo box is about; while "Pickup" is an additional content of the combo box, namely one value the combo box provides that the user is interested in and ultimately selects. Retaining such UI element contents in the final summaries is particularly important for interactive behaviour summarisation: First, such content provides interactive context information necessary to distinguish between subtle intentions. For instance, actions include selecting "1" from a combo box named "guest number" on a booking site by some users versus selecting more guests by other users. However, these detailed contents are ignored, and the summarised intentions are both *finding a hotel room*; the system may further inaccurately suggest unsuitable accommodations, e.g., family rooms for solo travellers and vice versa, causing decreased usability and potential frustration. Second, they are important for downstream applications, such as behaviour forecasting, to ensure the prediction is relevant and consistent with the

current context and underlying goals [50]. For example, if a user clicks on a button named "gluten-free" but this content is overlooked, the interactive system may mistakenly predict the upcoming actions to involve browsing or purchasing products containing gluten, leading to a worse user experience.

To address this challenge we propose a UI element attention mechanism to enhance the fine-tuning process by guiding the model to focus on these contents. This is similar to ensuring that a text summary covers essential keywords in natural language processing [13]. Specifically, for the $i$-th training sample, we create an attention vector $\mathbf{K}_i$, where each component $K_{ij}$ denotes the amount of attention assigned to the $j$-th token in the ground-truth summary:

$$K_{ij} = \begin{cases} \lambda & \text{if} \quad Token_j \in Token_{Detail} \\ 1 & \text{otherwise} \end{cases} \tag{2}$$

As such, tokens that contain action details receive $\lambda$ times the attention compared to other tokens. We empirically set $\lambda = 2$ in our experiments. The overall fine-tuning loss integrating this attention mechanism $\mathcal{L}_{Enhanced}$ is then computed as a weighted version of the original cross-entropy loss:

$$\mathcal{L}_{Enhanced} = \mathbf{K} \circ \mathcal{L}_{NextToken} \tag{3}$$

### 3.3 Implementation

We opted for the lightweight open-source Mistral-7B model [24] as the LLM backbone, known for its efficiency and effectiveness in handling various NLP tasks. Mistral-7B incorporates advanced techniques such as grouped-query attention for fast inference and sliding window attention for managing long sequences and complex contexts. These features contribute to Mistral-7B's superior performance on various benchmarks compared to other state-of-the-art models while using fewer parameters, thus conserving computational resources [41]. We used a batch size of 16 and a maximum input sequence length of 1024, with an initial learning rate of 1e-6. We used the Adam optimiser with $\beta_1 = 0.9$ and $\beta_2 = 0.95$ [43], and a cosine annealing scheduler for a progressive reduction of the learning rate following a cosine curve, a strategy proven to stabilise the training phase [37]. We fine-tuned the model for 15 epochs using eight Tesla V100-SXM2-32GB GPUs, completing the training within ten hours.

## 4  EXPERIMENTS

We conducted experiments to evaluate the quality of interactive behaviour summaries generated by SummAct. Given the novelty of this task and the lack of existing baseline methods, we compare the full model with several ablated versions instead. More specifically, starting with using an off-the-shelf, pretrained LLM – the common practice in HCI research currently [5, 22] – we incrementally add fine-tuning, sub-goal generation, and the UI element attention mechanism. We report quantitative metrics that measure how similar the generated summaries are compared to the ground truth, as well as qualitative similarities and differences of the generated summaries.

### 4.1  Datasets

We conducted all evaluations using two prominent datasets that encompass both the web (Mind2Web [12]) and mobile (MoTIF [6]) interaction contexts. These datasets are extensively used to understand and model user interfaces and interactive behaviours [5, 59]. They encode interactive behaviour as user actions to achieve specified interaction objectives. Each user action is annotated with information related to the UI element (category and content) and the user operation (e.g., click or swipe) associated with this element.

*4.1.1 Mind2Web.* This dataset provides crowdsourced actions across 2,350 tasks performed on 137 real-world websites (e.g., Booking, Uniqlo, IMDB) spanning 31 domains (e.g., travel, shopping, entertainment). As such, this dataset offers a wide variety of user actions and intentions and allows us to evaluate the performance of SummAct in real-world scenarios. Mind2Web includes three different data subsets: 1) *cross-domain* includes data instances from different domains, e.g., shopping vs travel; 2) *cross-website* includes instances from unseen websites, e.g., Booking vs Airbnb; and 3) *cross-task* includes unseen tasks, e.g., booking a flight vs buying a shirt. We preprocessed the dataset by generating natural language descriptions for the provided raw actions using a transformation template [39] (see Appendix A for more details).

*4.1.2 MoTIF.* This dataset targets a mobile interaction setting comprising screen touch data collected on 756 different tasks across 125 Android applications. The dataset directly provides synthetic natural language sentences describing each low-level action including the interacted UI element and the user's operation (click, type or swipe) on it. MoTIF includes both feasible and infeasible tasks such as tasks that are too unclear or cannot be completed in the given App. We only used the feasible tasks for our evaluations to ensure that user actions reliably reflect the corresponding intentions.

## 4.2 Ablations

We compared the full SummAct model **LLM+FT+SubGoal+Attn** (fine-tuned LLM using both the input actions and sub-goals with $\mathcal{L}_{Enhanced}$ for summarisation) with several ablated versions to evaluate the impact of the different modifications. To ensure a fair comparison, all methods used the same LLM and the same prompt templates.

- **LLM**: pretrained LLM using only the input actions.
- **LLM+SubGoal**: pretrained LLM using both the input actions and the sub-goals.
- **LLM+FT**: fine-tuned LLM using only the input actions with $\mathcal{L}_{NextToken}$ for summarisation.
- **LLM+FT+Attn**: fine-tuned LLM using only the input actions with $\mathcal{L}_{Enhanced}$ for summarisation.
- **LLM+FT+SubGoal**: fine-tuned LLM using both input actions and sub-goals with $\mathcal{L}_{NextToken}$ for summarisation.

Since the UI element attention mechanism is specifically incorporated into the loss function of fine-tuning, we do not have a standalone version of the pretrained LLM enhanced solely by the attention, i.e., LLM+Attn.

## 4.3 Quantitative Evaluations

We first quantify the similarity between ground truth and summarised intentions for all methods with four widely used NLP metrics [46, 57]. Specifically, we report Recall-Oriented Understudy for Gisting Evaluation (ROUGE) [33], Metric for Evaluation of Translation with Explicit ORdering (METEOR) [3], and Bilingual Evaluation Understudy (BLEU) [40] that count the overlapping n-grams between texts to assess their similarity, thus providing a measure that reflects lexical precision and recall. We further report an embedding-based metric using a state-of-the-art sentence encoder, Sentence Transformer[2], to obtain the sentence embeddings. Embedding-based metrics evaluate the cosine similarity between sentences and capture deeper and more robust semantic meanings that go beyond mere lexical matches [42]. All of these metrics indicate better results as their values increase.

Table 1 provides an overview of the results of this comparison. As can be seen from the table, our proposed full model consistently outperforms the ablated versions on all test sets and all metrics, obtaining a cosine similarity of up to 0.842 with the ground-truth user intentions. Among the three generalisation test sets from Mind2Web, SummAct performed better in *cross-website* and worse in the *cross-domain* setting. The former is likely because different websites within

---

[2]https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2

| Method | Metric | Mind2Web | | | MoTIF |
| | | Cross Domain | Cross Task | Cross Website | |
|---|---|---|---|---|---|
| LLM | CosSim | .357 | .348 | .380 | .203 |
| | BLEU | .004 | .004 | .004 | .012 |
| | ROUGE | .050 | .056 | .060 | .077 |
| | METEOR | .126 | .132 | .146 | .093 |
| LLM+SubGoal | CosSim | .381 | .374 | .404 | .230 |
| | BLEU | .006 | .004 | .006 | .009 |
| | ROUGE | .051 | .055 | .059 | .063 |
| | METEOR | .150 | .157 | .167 | .098 |
| LLM+FT | CosSim | .631 | .661 | .673 | .691 |
| | BLEU | .201 | .217 | .204 | .293 |
| | ROUGE | .313 | .301 | .312 | .511 |
| | METEOR | .303 | .306 | .322 | .510 |
| LLM+FT+Attn | CosSim | .705 | .740 | .756 | .785 |
| | BLEU | .349 | .296 | .305 | .407 |
| | ROUGE | .345 | .372 | .383 | .730 |
| | METEOR | .293 | .341 | .374 | .709 |
| LLM+FT+SubGoal | CosSim | .718 | .753 | .763 | .756 |
| | BLEU | .291 | .301 | .301 | .359 |
| | ROUGE | .381 | .392 | .391 | .689 |
| | METEOR | .323 | .332 | .363 | .662 |
| LLM+FT+SubGoal+Attn (SummAct) | CosSim | **.755** | **.796** | **.802** | **.842** |
| | BLEU | **.406** | **.453** | **.445** | **.453** |
| | ROUGE | **.390** | **.432** | **.447** | **.799** |
| | METEOR | **.376** | **.430** | **.454** | **.758** |

FT = Fine-tuning; SubGoal = Sub-goal generation; Attn = UI element attention

Table 1. Interactive behaviour summarisation results achieved by our proposed SummAct and its ablated versions. The evaluation is conducted on a web dataset Mind2Web (including three test subsets for generalisability assessment across domains, tasks and websites) and a mobile dataset MoTIF. We measure the summarisation quality with four metrics, cosine similarity between sentence embeddings, and n-gram based BLEU, ROUGE and METEOR. The best results are shown in **bold**.

the same domain share similar UI designs and thus require similar navigation patterns [12], which can be efficiently captured and integrated by our SummAct. These results also show, however, that generalisation across domains remains challenging due to the variations in context and user interactions.

We can also see from Table 1 that directly using a pretrained LLM performs the worst while adding fine-tuning, sub-goals, or the UI element attention mechanism improved performance notably. Although adding sub-goals (LLM+SubGoal) increased cosine similarity by up to 13.3% (0.203 vs 0.230 on MoTIF), the largest performance increase was achieved when adding fine-tuning (LLM vs LLM+FT), where the cosine similarity improved by up to 89.9% (0.348 vs 0.661, cross-task setting) on Mind2Web, and 240.4% (0.203 vs 0.691) on MoTIF.

Also, adding our two novel designs of sub-goals and UI element attention contributes to the effectiveness of SummAct (LLM+FT vs LLM+FT+SubGoal+Attn), together leading to an up to 21.9% improvement on the cosine similarity (0.691 vs 0.842 on MoTIF). Comparing our method with LLM+FT+Attn, the UI element attention mechanism increased the cosine

similarity on Mind2Web by 7.1% (0.705 vs 0.755) in the cross-domain setting, 7.6% (0.740 vs 0.796) in the cross-task setting, 6.1% (0.756 vs 0.802) in the cross-website setting, and by 7.3% (0.785 vs 0.842) on MoTIF. On n-gram-based metrics, SummAct obtained improvements of up to 53.0% (0.296 vs 0.453 on Mind2Web cross-task setting) on BLEU, 16.7% (0.383 vs 0.447 on Mind2Web cross-website setting) on ROUGE, and 28.3% (0.293 vs 0.376 on Mind2Web cross-domain setting) on METEOR. In Section 5.1, we further show that a lack of the UI contents harms performance for behaviour forecasting. Similarly, SummAct outperformed its ablated version that removed the sub-goals (LLM+FT+Attn), where the cosine similarity increased by 5.2% (0.718 vs 0.755) in Mind2Web cross-domain setting, 5.7% (0.753 vs 0.796) in the cross-task setting, 5.1% (0.763 vs 0.802) in the cross-website setting and 11.4% (0.756 vs 0.842) on MoTIF. Moreover, BLEU improved by up to 50.5% (0.301 vs. 0.453), achieved on Mind2Web cross-task setting; the maximum enhancement of ROUGE 15.97% (0.689 vs. 0.799) on MoTIF; while the largest improvement of METEOR reached 29.5% (0.332 vs. 0.430), obtained in Mind2Web cross-task setting. Taken together, these evaluations show the effectiveness of the proposed components for interactive behaviour summarisation.

### 4.4 Qualitative Analysis

We further examined the summaries generated by SummAct and its ablations qualitatively to understand the impact of the different designs.

*4.4.1 Impact of UI element attention.* Compared to the summaries generated by the full SummAct implementation, the ablated version without UI element attention lacks detailed context information embedded in UI element contents. For example, the ground-truth intention to *Find a **campground** in Orlando for two adults to check in on **Mar 29** and check out on **Mar 30*** was correctly summarised by SummAct as *Find a **campground** in Orlando for two adults from **March 29** to **March 30***. On the contrary, the ablated version (LLM+FT+SubGoal) produced a less accurate summary, *Find **hotels** in Orlando for two adults in **March***, missing the information of the precise dates and the specific type of accommodation. This issue arose because during the summarisation, the ablation ignored the detailed content in a clicking action on the button of "Find a KOA", which specified the accommodation type as a campground instead of any hotel.

Another example is the intention *Find a highest rated dealer for Cadillac with a rating above **4 stars** within **20 miles** of **zip 60606***. SummAct effectively summarised this into *Find a highest rated Cadillac dealer above **4 star** within **20 miles** of **60606***. However, the ablation's prediction was simply *Find the highest rated dealer for Cadillacs*, missing the specific criteria of rating and proximity.

This analysis shows that without the UI element attention mechanism, although the summaries retain the overall logic, they lack crucial specific information provided by the UI elements.

*4.4.2 Impact of sub-goals.* We then examined the summaries generated by the other ablation (LLM+FT+Attn) in which the sub-goals were removed from SummAct. We found that the summaries retained specific information but often failed to capture the overarching logic or coherence, especially when handling complex, multi-step interactive behaviour. Figure 2 shows two examples of this phenomenon, each with the user's input actions and their underlying intentions, the ground-truth intention, the summary generated by the full version of SummAct, and the summary generated by the ablation. For example, at the top, the user browsed through top-trending content within a community about work, then picked one post with the heading "...woman-dominated work..." and saved it. Based on this interactive behaviour, SummAct's summarisation was the same as the ground truth, i.e. *save a rising post on a community about work*. However, the ablation produced *find a post about women-dominated workplace*, focusing disproportionately on specific content cues, such as the particular post's heading, rather than the overall context of these actions. This demonstrates that

User actions | Intentions

| | | |
|---|---|---|
| 1. Type text "work" into search box with text "Search all of Reddit" on it<br>2. Click the search box element with text "Search all of Reddit" on it<br>3. Click the button element with text "Communities" on it<br>4. Click the heading element with text "r/work" on it<br>5. Click the button element with text "Top" on it<br>6. Click the menu item element with text "Rising" on it<br>7. Click the heading element with text "It's interesting working in a woman-dominated work..." on it<br>8. Click the button element with text "Save" on it | Ground-truth → | Save a rising post on a community about work |
| | SummAct → | Save a rising post on a community about work |
| | w/o Sub-goals → | Find a post about women dominated workplace |

User actions | Intentions

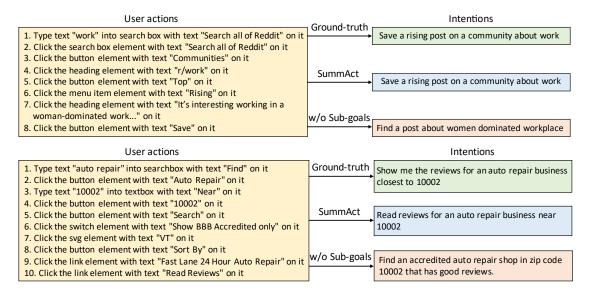| | | |
|---|---|---|
| 1. Type text "auto repair" into searchbox with text "Find" on it<br>2. Click the button element with text "Auto Repair" on it<br>3. Type text "10002" into textbox with text "Near" on it<br>4. Click the button element with text "10002" on it<br>5. Click the button element with text "Search" on it<br>6. Click the switch element with text "Show BBB Accredited only" on it<br>7. Click the svg element with text "VT" on it<br>8. Click the button element with text "Sort By" on it<br>9. Click the link element with text "Fast Lane 24 Hour Auto Repair" on it<br>10. Click the link element with text "Read Reviews" on it | Ground-truth → | Show me the reviews for an auto repair business closest to 10002 |
| | SummAct → | Read reviews for an auto repair business near 10002 |
| | w/o Sub-goals → | Find an accredited auto repair shop in zip code 10002 that has good reviews. |

Fig. 2. Two examples showing the input user actions, their underlying ground-truth intentions and those summarised by the full version of SummAct and its ablation removing sub-goals.

without sub-goals, summaries may lack the essence and broader context of user interactions and instead focus on specific keywords or aspects, resulting in a summary that does not reflect the overall intention.

In the second example, the user searched for *auto repair* with filtering conditions, selected one particular item, and then read its reviews. This was summarised correctly by SummAct: the summary includes that the user first looks for a business and then verifies its quality according to its reviews. However, the ablation summarised the intention as *find an accredited auto repair shop in zip code 10002 that has good reviews*, mistakenly understanding the goal of filtering the auto repair business based on their reviews. This occurs because, without sub-goals, the model processes all behaviour information indiscriminately and struggles to dissect the intricate dependencies and hierarchy among input actions. As a result, the ablated model erroneously swapped the sequence priorities between *finding auto repair* and *reading reviews*.

This analysis underscores the importance of our hierarchical approach in handling complex interactive behaviour. By generating intermediate sub-goals, SummAct not only distils key information from different stages of the user actions but also maintains a coherent understanding throughout each interaction stage, ensuring that the final summary encapsulates the overall context.

## 5 APPLICATIONS OF INTERACTIVE BEHAVIOUR SUMMARISATION

### 5.1 Interactive Behaviour Forecasting

Interactive behaviour forecasting is core to anticipatory and proactive interactive systems [54]. Specifically, we conducted the next action prediction following [55]. Most current next action prediction methods are based solely on the historical information without explicitly understanding overall intentions [26, 27, 56]. The summarisation of action history can potentially enhance the next action prediction by providing contextual information on users' goal trajectories.

We approached next action prediction as a multi-choice question-answering task, i.e., the model selects from a list of candidate UI elements that users may interact with, following [10, 18]. The pipeline includes three steps [12, 58]: First,

| Input | Cross Domain | | Cross Task | | Cross Website | |
|---|---|---|---|---|---|---|
| | Element | Operation | Element | Operation | Element | Operation |
| History | 31.2 | 42.1 | 34.2 | 46.2 | 30.6 | 40.4 |
| History+Summary (*w/o UI element attention*) | 37.8 | 45.9 | 43.8 | 47.7 | 34.9 | 43.9 |
| History+Summary (*Full*) | **46.8** | **50.5** | **47.8** | **54.5** | **40.1** | **45.0** |

Table 2. Element accuracy and operation F1 score of next action prediction achieved using 1) only behaviour history, 2) behaviour history plus its summary generated by SummAct *w/o UI element attention*, and 3) behaviour history and its summary generated by SummAct full version. Both metrics are in percentage. The best results are shown in bold.

we summarise the actions a user has performed using our SummAct. Then, we employ a candidate extraction method proposed by [12] to filter and rank UI elements on the current web page and retain the top-$k$ elements as the candidate targets for the next action. Retaining only top-$k$ elements is because the raw HTML contains a large amount of noisy UI data that can distract LLMs and cause hallucinations [38] and exceeds the maximum length of allowed input tokens. In our experiments, we set $k$ to 50 [12]. Finally, we fine-tune an LLM to select the next target UI element out of the 50 candidates and predict its corresponding operation out of three classes (click, select or type). We used the same fine-tuning set-up as the summarisation, i.e., Mistral-7B as the backbone LLM and the same learning rate, optimiser and scheduler. We fine-tuned the model for only three epochs, given its fast convergence. We required the behaviour history to include at least five past actions to offer adequate context, consistent with prior next action prediction works [27, 55]. The prompt for fine-tuning the LLM included these past actions, the intention summarised by SummAct, and the list of candidate UI elements (see Appendix B.3).

We compared our results with two baselines, as shown in Table 2. The first is when only using the action history to compare with and examine the effectiveness of summaries in next action prediction. The other is when using the history plus the summary generated by the ablated version of SummAct excluding the UI element attention, to check the specific contribution of keeping UI content as discussed in Section 3.2. Following [12, 59], we calculated the accuracy of UI element prediction, and F1 score of user operation prediction to measure the imbalanced operation classes. We showcased the performance of next action prediction on Mind2Web given that this dataset has more variety of intentions and interfaces than MoTIF, as shown in Section 4.1.

As presented in Table 2, integrating summarised intentions consistently enhanced the performance across domains, websites and tasks, with an average 12.9% higher element accuracy and 7.1% higher operation F1 score. Moreover, we observed adding the UI element attention improved the element accuracy and the operation F1 score by 6.1% and 4.2% on average, respectively, verifying that the UI content preserved by our method was helpful for next action prediction.

Enhancing next action prediction offers practical benefits for various interactive scenarios. For instance, the adaptive user interfaces can dynamically adjust the layout and functionality or directly recommend future actions to users to reduce required cognitive and physical demands and enhance usability [30]. Additionally, this approach allows automation agents to operate more efficiently by intuitively responding to user preferences without requiring explicit task instructions. This potentially leads to smoother, more personalised interactions adapting to evolving user behaviour.

## 5.2 Automatic Identification of Behaviour Synonyms

Besides directly capturing user intentions, our summarisation method can also identify "interactive behaviour synonyms", which, in our examples, are alternative action sequences reflecting the same underlying user intention.
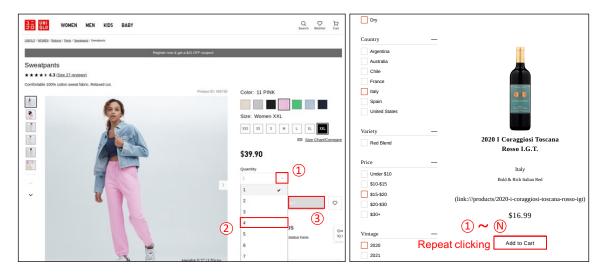
Fig. 3. An example of using synonyms to compare UI usability for the task of *adding N items into the shopping cart*. The Uniqlo website (left) allows users to add multiple items with just three clicks, while the Macy's website (right) requires one click per item, leading to more effort and less usability as *N* increases.

Practically, we used multiple windows of various lengths from two to the maximum sequence length to respectively segment each input action sequence into sub-sequences. Our SummAct processed each sub-sequence to summarise its underlying intention. We then combined all the sub-sequences and calculated the cosine similarity between their sentence embeddings. Two action sub-sequences are considered as synonyms if their cosine similarity is higher than a threshold.
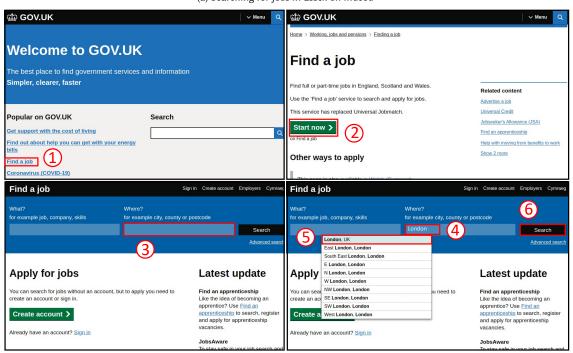
We showed three example types of synonyms that can provide interesting insights into interaction strategies and system usability: 1) when two synonym behaviours are different but generated from the same UI, the synonyms reflect different strategies or preferences users can take towards an intention; 2) when two synonym behaviours are different and generated from different UIs, the shorter action path indicates better usability; 3) when two synonym behaviours are generated from different UIs but have same actions, this presents that there are common behaviour patterns or UI designs.

*5.2.1 Different behaviours from the same UI.* Based on the same intention on the same UI, users can still generate various behaviours, showing different user preferences, interaction habits and strategies.

For instance, to *add an item to a new shopping list*, users could choose a shorter action path, i.e., creating a new list when adding an item, or a longer trajectory first navigating to the page showing all the existing lists, adding a list there, then returning to the item page, and finally adding the item. In another example, when the intention was to *find a top-rated restaurant in Miami*, one behaviour directly navigated to a page listing all restaurants and then selected the desired city. In contrast, another user path first identified the city, browsed a broad range of "things to do", and then narrowed down to restaurants. These variations may reflect different user preferences, browsing habits, or the clearness of user intention: the former path shows that the user may have a straightforward motivation to look for a *restaurant*. At the same time, the latter shows that the user may just look for a *place to go* in the city, not necessarily for

(a) Searching for jobs in Essex on Indeed



(b) Searching for jobs in London on Gov.UK

Fig. 4. An example of using synonyms to compare UI usability for the intention of *Searching for jobs in city A*. On the upper interface, the users can finish the task with only three actions; on the lower interface, the users must perform six actions, indicating worse usability.

a restaurant. These synonyms can help designers understand user preferences and habits, find the optimal interaction strategies, and design user-tailored interfaces.

*5.2.2   Different behaviours from different UIs.* UI designs impact the efficiency of achieving interactive intentions, shown by user behaviours. Therefore, the length of the synonyms found through our summarisation can be used as a metric of

UI efficiency. Unlike classical metrics like the keystroke-level model (KLM) that measure interactive system's usability via task completion time [16, 20], this metric will compare UIs via the number of actions required for the same intention. As shown in Figure 3, to *add N items into the shopping cart*: the Uniqlo website (left) enables users to directly choose the total quantity and add all of them to the cart at once, i.e., finishing with three clicks; while Macy's (right) only allows adding one item each time, i.e., with *N* clicks. When *N* has a large value, users on the latter interface will have to perform many more actions, harming efficiency and usability. Another example is shown in Figure 4, where the intention was to *search for jobs in a city* (Essex in the upper and London in the lower example). In the upper example, users only needed three clicks to navigate from the main page to view all jobs from the city, i.e., *Countries →England →Jobs in Essex*. On the contrary, in the lower example, users had to perform more actions to see the list of jobs, i.e., *click on Find a job →click on Start now →click on the text box under Where →type London →click on London →click on Search.* The reason is that the former website, Indeed, is specifically designed for job searching, optimising its UI to streamline this function. In contrast, the latter website, Gov.UK, serves multiple functions, not focusing primarily on job hunting, consequently leading to relatively lower efficiency. As such, these types of synonyms can help UX designers optimise interface design to facilitate quicker and more intuitive access, such as creating shortcuts for the interactive intentions that are dominant among users.

*5.2.3 Same behaviours from different UIs.* We found cases where the synonym behaviours followed the same patterns, although in different interactive systems. For example, when the intention was *to book a flight ticket from A to B*, user behaviours on different interfaces, including Kayak, Trip.com, American Airlines, and Expedia, typically followed a uniform process: users clicked flights, typed and selected the departure city or airport, and then typed and selected the destination. When *making an appointment with a doctor* on various medical platforms such as Zocdoc, Mayoclinic.org and Healthgrades, users also employed the same procedure: typed and clicked on the type of specialist, browsed and selected an available doctor, and finally picked the appointment time.

Such common patterns in interactive behaviour are why our model can generalise across different websites and tasks (as shown in Table 1). Understanding these patterns also gives UX designers intuitive starting points to create a new interface that aligns with established user behaviour and expectations. This can reduce the learning curve while ensuring a consistent, friendly user experience.

## 5.3 Language-based action retrieval

Through summarisation, SummAct lowers the barrier of understanding and interpreting interactive behaviour. Instead of analysing the low-level actions and complex UI contexts, users can directly read their semantic meaning in natural language and further utilise language queries to access specific behaviour trajectories. This capability opens up various practical application scenarios. In **technical troubleshooting**, the system streamlines the diagnostic process by providing quick access to relevant troubleshooting steps based on user-provided problem descriptions and their intended functions. In **educational contexts**, when novice users struggle to complete an interactive goal, the system can quickly pull up the most relevant teaching actions or tutorials from experts for them to watch and learn. For instance, when users are learning to use complex softwares that require professional skills, such as editing an image with Adobe Photoshop or drawing with AutoCAD, they can type a wanted function or effect, then retrieve, learn and apply other users' action paths implementing it. Moreover, such retrieval lays the foundation for the development of **question-answering agents** and **conversational user interfaces**. These systems observe and summarise a user behaviour, then retrieve similar behaviours or query related features, such as behaviour frequency from databases

to provide more context-aware and personalised responses. For example, based on a task, if we can retrieve a large amount of behaviours, this task may be a common or routine task so that the system can offer automating this task; furthermore, if the most frequent behaviour has a long action trajectory, the system can adjust the menu and dashboard to reduce users' time and effort.

## 6 DISCUSSION

### 6.1 Interactive Behaviour Summarisation

Although recent works have begun to model interactive behaviour from a natural language perspective [12, 50, 53], how to use this approach to understand underlying user intentions during interactions remains under-explored. Our work fills this gap by formulating intention recognition as an interactive behaviour summarisation task, i.e., summarising input action sequences into sentences reflecting latent user intentions. This formulation represents a paradigm shift in intention modelling. Traditionally, intention recognition tasks have been constrained by a closed set of predefined intentions, limiting the system's ability to adapt to the varied and evolving needs of users dynamically. By employing natural language to encapsulate user actions, we offer a more flexible and scalable framework that can intuitively interpret intentions, eliminate the exhaustive definition of every possible user intention, and cater to an open-ended array of these intentions. The open-ended nature of these summaries allows for recognising unseen, out-of-distribution user intentions, verified by the SummAct's robust performance in the cross-task setting on Mind2Web (see Section 4.3). The open-ended set also enables continuous learning and system refinement [48], which can integrate emergent user behaviours and preferences without requiring extensive manual updates or reconfigurations. Furthermore, leveraging a *natural language* representation can improve explainability by grounding the understanding of complex actions within a familiar linguistic framework, lowering the barrier of interpreting and analysing user behaviour.

In addition to its inherent utility, we also demonstrate various example applications enabled by the summary, including interactive behaviour forecasting, automatic behaviour synonym identification, and language-based behaviour retrieval. Understanding users' intentions from past actions is important to forecast future actions along the intention. Therefore, we enhanced current computational next action prediction methods based on action history by adding the summary as an additional input. SummAct's summarisation led to better prediction performance, as shown in Table 1. These improvements underscore the effectiveness of adding semantic information through interactive behaviour summarisation. Moreover, if two behaviours have the similar summary, they can be considered synonyms, i.e., alternatives achieving the same goal. Synonyms can provide insights into understanding diverse user preferences and interaction strategies, optimising interactive systems, and finding common patterns in behaviour and UI designs across systems. The linguistic nature of the generated summaries also unlocks language-based action retrieval, meaning that users can query and retrieve specific action sequences using simple language requests. This will facilitate learning from expert users in navigating complex interfaces or completing difficult tasks and enable question-answering agents and conversational user interfaces.

### 6.2 Design of SummAct

We propose a novel LLM-based approach, SummAct, that includes two methodological contributions and design choices to generate natural language summarisation from user actions and their UI context. The first design choice is to employ a hierarchical process adept at handling the complexity and variability of user actions. The method first enriches the semantic context available for the final summary generation by producing intermediate sub-goals from low-level

actions. The second design choice is an UI element attention mechanism applied on the fine-tuning loss to prevent over-abstraction by preserving essential context information embedded in the UI elements, such as the meaning of the clicked button. Extensive quantitative comparisons with ablations showed that both design choices contribute to the effectiveness of SummAct (see Table 1). Additionally, in Section 4.4, the qualitative evaluations examined the generated text and verified that both designs were necessary and complementary to each other: without the sub-goals, the summaries can have logical errors; while without the detail attention, the summaries lacked action details (see Figure 2).

Moreover, as shown in Table 1, our method has significantly improved upon the results obtained from using pretrained LLMs (up to three times better, on MoTIF), which is for now a common practice in HCI research that use LLMs [22, 36, 46]. The substantial enhancement brought by SummAct reveals the considerable potential of further enhancements in these HCI works and positions our approach as a pioneering example in the field. Our findings suggest that the HCI community can gain immensely by further adapting advanced natural language processing techniques to specific HCI applications.

In implementing SummAct, we utilised Mistral-7B as the backbone LLM due to its lightweight and robust performance across various NLP tasks. However, our SummAct framework allows replacing Mistral-7B with other LLM models if computing resources are available or more powerful models appear.

### 6.3  Limitations and Future Work

In our work, we investigated summarising the interactive behaviour at the UI-element level, i.e., the users operate on UI elements such as buttons or combo boxes. In the future, we will dive deeper into the raw, pixel-level interactive behaviours, e.g., integrating the specific on-screen locations of each move, click or tap into our model, which will carry more information about users and their intentions [55, 56]. Currently, we integrated UI information via HTML and DOM elements. Future enhancements can adopt vision language models to process GUI screenshots [29, 59]. This will allow for capturing visual cues that HTML alone cannot provide, such as iconography, layout spatial arrangements, and thematic designs that influence user interactions. The proposed UI element attention implemented a crude match between tokens by simply considering whether they are identical, but moving forward, we could consider matching their semantics to increase the model's flexibility. Moreover, our exploration of interactive behaviour summarisation is based on the assumption that the observed UI trajectories accurately reflect user intentions without error. While essential for developing our method, in real-world interactions, it is possible that actions may not always convey true intentions due to errors in user operation, which will be interesting for developing more robust models in the future.

### 7  CONCLUSION

In this work, we model interactive behaviour from a natural language viewpoint and investigate a novel interactive behaviour summarisation task, summarising input actions into natural language descriptions. These descriptions reflect user intentions underlying their interactive behaviour. Towards this task, we propose SummAct – an LLM-based method with two specific designs, hierarchical summarisation and UI element attention. We evaluated our method on two datasets, covering both the web and mobile interactive settings, from both the quantitative and qualitative perspectives. Results demonstrated the effectiveness of our method in summarising intentions and the complementary contributions of our two designs. We then showcased example applications of interactive behaviour summarisation, including behaviour forecasting, automatic identification of behaviour synonyms, and language-based behaviour retrieval. The natural language representation of interactive behaviour can boost the explainability of computational

behaviour modelling and contribute to developing more intuitive and responsive interactive systems. Furthermore, our significant improvement over the common practice of directly using LLMs in HCI suggests large potential benefits from further adapting advanced NLP techniques to HCI tasks.

## REFERENCES

[1] Abdulaziz Almehmadi. 2021. Micro-Behavioral Accidental Click Detection System for Preventing Slip-Based Human Error. *Sensors* 21, 24 (2021), 8209.

[2] Gregor Bachmann and Vaishnavh Nagarajan. 2024. The pitfalls of next-token prediction. *arXiv preprint arXiv:2403.06963* 1, 1 (2024), 1–1.

[3] Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*. 1, 1, 65–72.

[4] Roman Bednarik, Hana Vrzakova, and Michal Hradis. 2012. What do you want to do next: a novel approach for intent prediction in gaze-based interaction. In *Proceedings of the symposium on eye tracking research and applications*. 1, 1, 83–90.

[5] Omri Berkovitch, Sapir Caduri, Noam Kahlon, Anatoly Efros, Avi Caciularu, and Ido Dagan. 2024. Identifying User Goals from UI Trajectories. *arXiv preprint arXiv:2406.14314* 1, 1 (2024), 1–1.

[6] Andrea Burns, Deniz Arsan, Sanjna Agrawal, Ranjitha Kumar, Kate Saenko, and Bryan A. Plummer. 2022. A Dataset for Interactive Vision Language Navigation with Unknown Command Feasibility. In *European Conference on Computer Vision (ECCV)*, Vol. 1. 1, 1, 1–1.

[7] Ling-Hao Chen, Shunlin Lu, Ailing Zeng, Hao Zhang, Benyou Wang, Ruimao Zhang, and Lei Zhang. 2024. MotionLLM: Understanding Human Behaviors from Human Motions and Videos. *arxiv:2405.20340* 1, 1 (2024), 1–1.

[8] Xianyu Chen, Ming Jiang, and Qi Zhao. 2024. GazeXplain: Learning to Predict Natural Language Explanations of Visual Scanpaths. *arXiv preprint arXiv:2408.02788* 1, 1 (2024), 1–1.

[9] Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. 2024. Seeclick: Harnessing gui grounding for advanced visual gui agents. *arXiv preprint arXiv:2401.10935* 1, 1 (2024), 1–1.

[10] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research* 25, 70 (2024), 1–53.

[11] Brendan David-John, Candace Peacock, Ting Zhang, T Scott Murdison, Hrvoje Benko, and Tanya R Jonker. 2021. Towards gaze-based prediction of the intent to interact in virtual reality. In *ACM symposium on eye tracking research and applications*. ACM, 1, 1–7.

[12] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. Mind2Web: Towards a Generalist Agent for the Web. arXiv:2306.06070 [cs.CL]

[13] Wafaa S El-Kassas, Cherif R Salama, Ahmed A Rafea, and Hoda K Mohamed. 2021. Automatic text summarization: A comprehensive survey. *Expert systems with applications* 165 (2021), 113679.

[14] Anis Elbahi, Mohamed Ali Mahjoub, and Mohamed Nazih Omri. 2013. Hidden markov model for inferring user task using mouse movement. In *Fourth International Conference on Information and Communication Technology and Accessibility (ICTA)*. IEEE, IEEE, 1, 1–7.

[15] Anis Elbahi and Mohamed Nazih Omri. 2015. Web user interact task recognition based on conditional random fields. In *Computer Analysis of Images and Patterns: 16th International Conference, CAIP 2015, Valletta, Malta, September 2-4, 2015 Proceedings, Part I 16*. Springer, Valletta, 740–751.

[16] Erik Frøkjær, Morten Hertzum, and Kasper Hornbæk. 2000. Measuring usability: are effectiveness, efficiency, and satisfaction really correlated?. In *Proceedings of the SIGCHI conference on Human Factors in Computing Systems*. ACM, 1, 345–352.

[17] Eugene Yujun Fu, Tiffany CK Kwok, Erin You Wu, Hong Va Leong, Grace Ngai, and Stephen CF Chan. 2017. Your mouse reveals your next activity: towards predicting user intention from mouse interaction. In *2017 IEEE 41st Annual Computer Software and Applications Conference (COMPSAC)*, Vol. 1. IEEE, IEEE, 1, 869–874.

[18] Yu Gu, Xiang Deng, and Yu Su. 2022. Don't Generate, Discriminate: A Proposal for Grounding Language Models to Real-World Environments. *arXiv preprint arXiv:2212.09736* 1, 1 (2022), 1–1.

[19] Sakib Haque, Zachary Eberhart, Aakash Bansal, and Collin McMillan. 2022. Semantic similarity metrics for evaluating source code summarization. In *Proceedings of the 30th IEEE/ACM International Conference on Program Comprehension*. IEEE/ACM, 1, 36–47.

[20] Kasper Hornbæk. 2006. Current practice in measuring usability: Challenges to usability studies and research. *International journal of human-computer studies* 64, 2 (2006), 79–102.

[21] Zhiming Hu, Andreas Bulling, Sheng Li, and Guoping Wang. 2022. EHTask: recognizing user tasks from eye and head movements in immersive virtual reality. *IEEE Transactions on Visualization and Computer Graphics* 29, 4 (2022), 1992–2004.

[22] Forrest Huang, Gang Li, Tao Li, and Yang Li. 2024. Automatic Macro Mining from Interaction Traces at Scale. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. ACM, Honolulu, 1–16.

[23] Peter C Humphreys, David Raposo, Tobias Pohlen, Gregory Thornton, Rachita Chhaparia, Alistair Muldal, Josh Abramson, Petko Georgiev, Adam Santoro, and Timothy Lillicrap. 2022. A data-driven approach for learning to control computers. In *International Conference on Machine Learning*. PMLR, 1, 1, 9466–9482.

[24] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7B. *arXiv preprint arXiv:2310.06825* 1, 1 (2023), 1–1.

[25] Kazuki Kawamura and Jun Rekimoto. 2024. FastPerson: Enhancing Video-Based Learning through Video Summarization that Preserves Linguistic and Visual Contexts. In *Proceedings of the Augmented Humans International Conference 2024*. 1, 1, 205–216.

[26] Saskia Koldijk, Mark Van Staalduinen, Mark Neerincx, and Wessel Kraaij. 2012. Real-time task recognition based on knowledge workers' computer activities. In *Proceedings of the 30th European Conference on Cognitive Ergonomics*. ACM, Edinburgh, 152–159.

[27] Tiffany CK Kwok, Eugene Yujun Fu, Erin You Wu, Michael Xuelin Huang, Grace Ngai, and Hong-Va Leong. 2018. Every little movement has a meaning of its own: Using past mouse movements to predict the next interaction. In *Proceedings of the 23rd International Conference on Intelligent User Interfaces*. ACM, Berlin, 397–401.

[28] Haoran Li and Wei Lu. 2021. Mixed cross entropy loss for neural machine translation. In *International Conference on Machine Learning*. PMLR, 1, 1, 6425–6436.

[29] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*. PMLR, 1, 1, 19730–19742.

[30] Jiahao Nick Li, Yan Xu, Tovi Grossman, Stephanie Santosa, and Michelle Li. 2024. OmniActions: Predicting Digital Actions in Response to Real-World Multimodal Sensory Inputs with LLMs. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. ACM, Honolulu, 1–22.

[31] Rui Li, Guoyin Wang, and Jiwei Li. 2024. Are Human-generated Demonstrations Necessary for In-context Learning?. In *The Twelfth International Conference on Learning Representations*, Vol. 1. 1, 1, 1–1.

[32] Yang Li, Jiacong He, Xin Zhou, Yuan Zhang, and Jason Baldridge. 2020. Mapping natural language instructions to mobile UI action sequences. *arXiv preprint arXiv:2005.03776* 1, 1 (2020), 8198–8210.

[33] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*. 1, 1, 74–81.

[34] Jing Lin, Ailing Zeng, Shunlin Lu, Yuanhao Cai, Ruimao Zhang, Haoqian Wang, and Lei Zhang. 2024. Motion-x: A large-scale 3d expressive whole-body human motion dataset. *Advances in Neural Information Processing Systems* 36, 1 (2024), 1–1.

[35] Yang Liu and Mirella Lapata. 2019. Hierarchical transformers for multi-document summarization. In *The 57th Annual Meeting of the Association for Computational Linguistics (ACL)*, Vol. 1. 1, 1, 1–1.

[36] Zhe Liu, Chunyang Chen, Junjie Wang, Mengzhuo Chen, Boyu Wu, Yuekai Huang, Jun Hu, and Qing Wang. 2024. Unblind Text Inputs: Predicting Hint-text of Text Input in Mobile Apps via LLM. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. ACM, Honolulu, 1–20.

[37] Ilya Loshchilov and Frank Hutter. 2016. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983* 1, 1 (2016), 1–1.

[38] Ariana Martino, Michael Iannelli, and Coleen Truong. 2023. Knowledge injection to counter large language model (LLM) hallucination. In *European Semantic Web Conference*. Springer, 1, 1, 182–185.

[39] Runliang Niu, Jindong Li, Shiqi Wang, Yali Fu, Xiyu Hu, Xueyuan Leng, He Kong, Yi Chang, and Qi Wang. 2024. Screenagent: A vision language model-driven computer control agent. In *The 33rd International Joint Conference on Artificial Intelligence*, Vol. 1. 1, 1, 1–1.

[40] Kishore Papineni. 2001. BLEU: a method for automatic evaluation of MT. *Research Report, Computer Science RC22176 (W0109-022)* 1, 1 (2001), 1–1.

[41] Shiva Kumar Pentyala, Zhichao Wang, Bin Bi, Kiran Ramnath, Xiang-Bo Mao, Regunathan Radhakrishnan, Sitaram Asur, et al. 2024. PAFT: A Parallel Training Paradigm for Effective LLM Fine-Tuning. *arXiv preprint arXiv:2406.17923* 1, 1 (2024), 1–1.

[42] Derek J. Phillips, Tim A. Wheeler, and Mykel J. Kochenderfer. 2017. Generalizable intention prediction of human drivers at intersections. In *2017 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 1, 1665–1670. https://doi.org/10.1109/IVS.2017.7995948

[43] Cheng Qian, Bingxiang He, Zhong Zhuang, Jia Deng, Yujia Qin, Xin Cong, Yankai Lin, Zhong Zhang, Zhiyuan Liu, and Maosong Sun. 2024. Tell Me More! Towards Implicit User Intention Understanding of Language Model Driven Agents. *arXiv preprint arXiv:2402.09205* 1, 1 (2024), 1–1.

[44] Dana Rezazadegan, Shlomo Berkovsky, Juan C Quiroz, A Baki Kocaballi, Ying Wang, Liliana Laranjo, and Enrico Coiera. 2020. Automatic speech summarisation: A scoping review. *arXiv preprint arXiv:2008.11897* 1, 1 (2020), 1–1.

[45] Noam Rotstein, David Bensaïd, Shaked Brody, Roy Ganz, and Ron Kimmel. 2024. Fusecap: Leveraging large language models for enriched fused image captions. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. IEEE, 1, 5689–5700.

[46] Bryan Wang, Gang Li, and Yang Li. 2023. Enabling conversational interaction with mobile ui using large language models. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. ACM, Hamburg, 1–17.

[47] Bryan Wang, Gang Li, Xin Zhou, Zhourong Chen, Tovi Grossman, and Yang Li. 2021. Screen2words: Automatic mobile UI summarization with multimodal learning. In *The 34th Annual ACM Symposium on User Interface Software and Technology*. ACM, 1, 498–510.

[48] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. 2024. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *Advances in Neural Information Processing Systems* 36, 1 (2024), 1–1.

[49] Jason Wei and Kai Zou. 2019. Eda: Easy data augmentation techniques for boosting performance on text classification tasks. *arXiv preprint arXiv:1901.11196* 1, 1 (2019), 1–1.

[50] Hao Wen, Yuanchun Li, Guohong Liu, Shanhui Zhao, Tao Yu, Toby Jia-Jun Li, Shiqi Jiang, Yunhao Liu, Yaqin Zhang, and Yunxin Liu. 2023. Empowering llm to use smartphone for intelligent task automation. *arXiv preprint arXiv:2308.15272* 1, 1 (2023), 1–1.

[51] Xuhai Xu, Chun Yu, Yuntao Wang, and Yuanchun Shi. 2020. Recognizing unintentional touch on interactive tabletop. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 4, 1 (2020), 1–24.

[52] Lin-Ping Yuan, Boyu Li, Jindong Wang, Huamin Qu, and Wei Zeng. 2024. Generating Virtual Reality Stroke Gesture Data from Out-of-Distribution Desktop Stroke Gesture Data. In *2024 IEEE Conference Virtual Reality and 3D User Interfaces (VR)*. IEEE, IEEE, 1, 732–742.

[53] Guanhua Zhang, Matteo Bortoletto, Zhiming Hu, Lei Shi, Mihai Bâce, and Andreas Bulling. 2023. Exploring Natural Language Processing Methods for Interactive Behaviour Modelling. In *The Proceeding of 2023 IFIP TC13 Conference on Human-Computer Interaction (INTERACT)*. IFIP, York, 1–18.

[54] Guanhua Zhang, Susanne Hindennach, Jan Leusmann, Felix Bühler, Benedict Steuerlein, Sven Mayer, Mihai Bâce, and Andreas Bulling. 2022. Predicting Next Actions and Latent Intents during Text Formatting. In *Proceedings of the CHI Workshop Computational Approaches for Understanding, Generating, and Adapting User Interfaces (2022-01-01)*. ACM, New Orleans, 1–6.

[55] Guanhua Zhang, Zhiming Hu, Mihai Bâce, and Andreas Bulling. 2024. Mouse2Vec: Learning Reusable Semantic Representations of Mouse Behaviour. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. ACM, Honolulu, 1–17.

[56] Guanhua Zhang, Zhiming Hu, and Andreas Bulling. 2024. DisMouse: Disentangling Information from Mouse Movement Data. In *Proc. ACM Symposium on User Interface Software and Technology (UIST)*. ACM, Pittsburgh, 1–13.

[57] Haopeng Zhang, Philip S Yu, and Jiawei Zhang. 2024. A Systematic Survey of Text Summarization: From Statistical Methods to Large Language Models. *arXiv preprint arXiv:2406.11289* 1, 1 (2024), 1–1.

[58] Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2024. Gpt-4v (ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614* 1, 1 (2024), 1–1.

[59] Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. 2024. Synapse: Trajectory-as-Exemplar Prompting with Memory for Computer Control. In *The Twelfth International Conference on Learning Representations*, Vol. 1. 1, 1, 1–1.

# A DESCRIBING MIND2WEB ACTIONS IN NATURAL LANGUAGE

As mentioned in Section 4.1, we preprocessed the original actions provided by Mind2Web into natural language descriptions to better leverage the understanding and reasoning capabilities inherent in LLMs. Table 3 presents three examples of the original actions, each from an user operation category (click, select or type).

| Element Category | Element Content | | User Operation | Element Content (Additional) |
|---|---|---|---|---|
| [*button*] | Add to Cart | → | CLICK | – |
| [*combobox*] | Sort By | → | SELECT | Price Low to High |
| [*searchbox*] | Search | → | TYPE | Johannesburg |

Table 3. Three example input action strings from Mind2Web dataset. Every string, representing an input action, contains the information of the interacted UI element (category and content) and user's operation on it.

Following [39], we first split the original action strings to the UI element category, content (the inherent meaning of this element, e.g., its name or the text on it), additional content for type and select operations (specific content the user is interested in, e.g., selected value from a combo box) and the user operation category; and then inserted these components into the natural language template, structured as:

If $Operation = CLICK$ : [$Operation$] the [$Category$] element with text "[$Content$]" on it

If $Operation = SELECT$ : [$Operation$] "[$Content(Additional)$]" from [$Category$] with text "[$Content$]" on it

If $Operation = TYPE$ : [$Operation$] text "[$Content(Additional)$]]" into [$Category$] with text "[$Content$]" on it

Through this templating approach, the example action strings are represented as:

- Click the button element with text "Add to Cart" on it
- Select "Price Low to High" from combobox with text "Sort By" on it
- Type text "Johannesburg" into searchbox with text "Search" on it

    You are a helpful computer task assistant. You can break any high-level tasks that can be performed in graphical user interfaces into sub-goals. The
user will provide the interface to be used, the domain and subdomain of the interface. Along with these general information, the user will also
provide the task they intend to achieve and the actions they perform towards this task.
    You have to analyse the given information and create a list of sub-goals. Each sub-goal should be a summary for several actions, which means the
number of sub-goals should be more than 1 and less than the number of actions. The sub-goals should together lead to the final task. The following
are examples of generating such lists of sub-goals.                                                                    **Overall context**

---

## Example 1 ##
# INPUT #
Website: exploretock
Domain: Travel
Sub-domain: Restaurant                                                                                                **Environment metadata**
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
Task: Check for pickup restaurant available in Boston, NY on March 18, 5pm with just one guest.
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
Actions (Each line is one action):
- Select "Pickup" from combobox with text "Reservation type" on it
- Type text "Boston" into searchbox with text "Find a location" on it
- Click the span element with text "Boston" on it
- Click the button element with text "18" on it
- Select "5 00 PM" from combobox with text "Time" on it
- Click the span element with text "2 guests" on it
- Select "1 guest" from combobox with text "Size" on it
- Click the button element with text "Update search" on it                                                            **Interactive Behaviour**
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
# OUTPUT #
Sub-goals (Each line is one sub-goal):
- Set the reservation type to "Pickup".
- Enter "Boston" as the location.
- Select the reservation date and time to March 18, 5pm.                                                              **Sub-goals**
- Adjust the number of guests to 1.
......                                                                                                                **Example demonstrations**

---

# INPUT #
Website: exploretock
Domain: Travel
Sub-domain: Car rental
Task: Sign Allan Smith for email offers with the email allan.smith@gmail.com and zip code 10001
Actions (Each line is one action):
- Click the a element with text "Return to a different location" on it',
- Type text "Allan" into textbox with text "First Name (required)" on it',
- Type text "Smith" into textbox with text "Last Name (required)" on it',
- Type text "allan.smith@gmail.com" into textbox with text "Email Address (required)" on it',
- Type text "allan.smith@gmail.com" into textbox with text "Confirm Email Address (required)" on it',
- Type text "10001" into textbox with text "ZIP Code (required)" on it',
- Click the button element with text "Submit" on it'

# OUTPUT #                                                                                   **New sample for which the model needs to generate sub-goals**

Fig. 5. Prompt used to generate sub-goals using in-context learning.

## B   PROMPTS

### B.1   In-Context Learning Prompts for Sub-goal Generation

As described in Section 3.1, we apply in-context learning to generate sub-goals from low-level interactive actions.
Figure 5 shows an example of the prompt we used for the pretrained LLM. The prompt has three main components: 1)
an overall context describing the task the LLM needs to solve, the input format and expected output; 2) examples used
in the in-context learning annotated by experts, each including environment metadata, task (only in training samples),
interactive behaviour and sub-goals (we only show one example due to its excessive length); 3) and a new sample for
which the model needs to generate sub-goals, where the text in green should be replaced and updated for each sample.
The texts in blue are related to the ground-truth overall intentions, and thus should be removed from testing samples to

Your task is to understand and summarize a user's intention behind their actions on a user interface. You have a list of information, including the website, domain and sub-domain of the user interface, history actions the user performed, and sub-goals (several low-level summarizations of subsets of history actions). Combine all the information and summarize the intention.                                                    Overall context

## Website:
dmv.virginia.gov

## Domain:
Service

## Sub-domain:
Government                                                                                                                    Environment metadata

## Actions (Each line is one action):
- Click the link element with text "Locations" on it
- Click the link element with text "DMV'S MOBILE OFFICES" on it
- Click the link element with text "View Calendar by Location" on it
- Click the button element with text "Location" on it
- Click the link element with text "HIGHLAND" on it                                                                          Interactive Behaviour

## Sub-goals summarized from these actions:
- Go to DMV Locations page.
- Open Mobile Offices calendar view.
- Filter the locations based on Highland area.                                                                              Sub-goals

# Instructions:                                                                                                            Expected output format
## Summarize:
In clear and concise language, summarize the comprehensive goal the user wants to achieve via the history actions. Present the summary after the heading [SUMMARY].
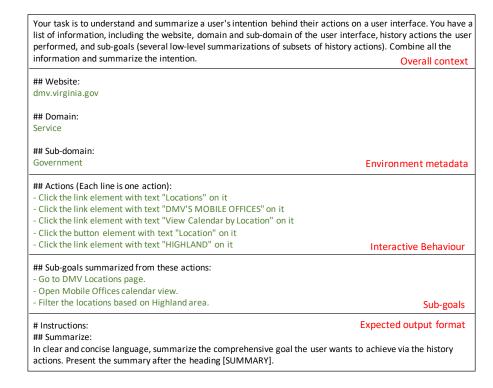
Fig. 6. Prompt used to summarise the overall intention from low-level interactive behaviour and mid-level sub-goals.

avoid data leakage. The example in Figure 5 is from Mind2Web dataset, which provides the metadata including website, domain and sub-domain. When using the prompt on MoTIF dataset, simply replace them with the provided mobile application as the new meta data.

### B.2  Prompts for Detail Enhanced Fine-Tuning

As described in Section 3.2, we fine-tune LLMs to summarise the final, overall intention using both low-level interactive behaviour and mid-level sub-goals. Figure 6 illustrates an example of the prompt we used for fine-tuning. The prompt comprises five parts: 1) an overall context describing the task the LLM needs to solve, the input format and expected output; 2) environment metadata; 3) user input actions; 4) sub-goals; and 5) expected output format. The text in green should be replaced and updated for each sample. The example in Figure 6 is from Mind2Web dataset, which provides the metadata including website, domain and sub-domain. When using the prompt on MoTIF dataset, simply replace them with the provided mobile application as the metadata.

### B.3  Prompts for Next Action Prediction

As described in Section 5.1, we fine-tune an LLM to predict the next action a user may perform based on the history actions and the intentions summarised from them via our SummAct. Figure 7 shows an example of this prompt, consisting of five parts: 1) an overall context describing the task the LLM needs to solve and the allowed action space; 2) input actions the user has performed; 3) DOM of the current webpage related to candidate UI elements; 4) intention

You are an agent in a web-based environment. Your task is to predict the user's next action within the given action space, based on the user's previous actions, the DOM of the current webpage, and intention so far.
# Action space:
1. `CLICK`: Click on an element.
2. `TYPE [value]`: Type a value into an element.
3. `SELECT [value]`: Select a value for an element.                                                                    <span style="color:red">Overall context</span>

# Previous actions (Each line is one action):
- Click the button element with text Hotel on it
- Click the button element with text Restaurant on it
- Select Pickup from combobox with text Reservation type on it
- Type text New into searchbox with text Find a location on it
- Click the span element with text New York on it                                                                       <span style="color:red">Previous actions</span>

# DOM:
(html (div (div (div (p DELICIOUS ) (p id=0 STARTS ) (p HERE. ) ) (div (label id=1 Reservation type ) (div (label id=2 Date ) (div (input text date thu, mar 16
) (div id=3 (button button date, selected value is thu, (svg) ) ) ) ) ) ) ) (a id=4 (div img image of hilton alumni association ) (section (h3 Hilton Alumni
Association ) (p Hilton, NY Non-culinary ) ) ) ) ) ...                                                                   <span style="color:red">DOM of current webpage</span>

# User intention so far: Check for pickup restaurant available in New York.                                             <span style="color:red">Summary of previous actions</span>

What should be the next action? Please select from the following choices by selecting the letter that identifies the correct choice. (If the correct action
is not in the page above, please select A. 'None of the above'):
A. None of the above
B. (p id=0 STARTS )
C. (label id=1 Reservation type )
D. (label id=2 Date )
E. (div id=3 (button button date, selected value is thu, (svg)
F. (a id=4 (div img image of hilton alumni association ) ...

Write your answer as: [ANSWER] your final answer [/ANSWER]                                                              <span style="color:red">Multi-choice question</span>

Fig. 7. Prompt used to forecast interactive behaviour based on previous input actions and the current web page.

summarised from the performed actions; and 5) candidate UI elements as the next target. The text in green should be replaced and updated for each sample.