

Towards Privacy-Preserving and Personalized Smart Homes via Tailored Small Language Models

Xinyu Huang, *Student Member, IEEE*, Leming Shen, *Student Member, IEEE*, Zijing Ma, *Student Member, IEEE*, Yuanqing Zheng, *Senior Member, IEEE*,

Abstract—Large Language Models (LLMs) have showcased remarkable generalizability in language comprehension and hold significant potential to revolutionize human-computer interaction in smart homes. Existing LLM-based smart home assistants typically transmit user commands, along with user profiles and home configurations, to remote servers to obtain personalized services. However, users are increasingly concerned about the potential privacy leaks to the remote servers. To address this issue, we develop *HomeLLaMA*, an on-device assistant for privacy-preserving and personalized smart home serving with a tailored small language model (SLM). *HomeLLaMA* learns from cloud LLMs to deliver satisfactory responses and enable user-friendly interactions. Once deployed, *HomeLLaMA* facilitates proactive interactions by continuously updating local SLMs and user profiles. To further enhance user experience while protecting their privacy, we develop *PrivShield* to offer an optional privacy-preserving LLM-based smart home serving for those users, who are unsatisfied with local responses and willing to send less-sensitive queries to remote servers. For evaluation, we build a comprehensive benchmark *DevFinder* to assess the service quality. Extensive experiments and user studies ($M = 100$) demonstrate that *HomeLLaMA* can provide personalized services while significantly enhancing user privacy.

Index Terms—Smart Home, Large Language Model, User privacy, Personalization

I. INTRODUCTION

THE proliferation of smart homes has significantly facilitated the development of intelligent living spaces [1], [2]. Typically, a smart home is a residence equipped with various interconnected devices and systems [3], [1] that can be controlled remotely or autonomously to enhance efficiency and convenience through technologies such as IoT and AI-based chatbots [4]. The long-term goal of smart homes is to achieve seamless user-assistant interaction, allowing systems to deeply comprehend user intents and deliver satisfactory and personalized responses [5].

Existing commercial-off-the-shelf (COTS) smart home assistants, like Amazon Alexa [6] and Apple Siri [7], are task-specific models pretrained on various instruction datasets which may lead to degraded performance on unseen tasks. For instance, when users provide an under-specified command (e.g., “Let guests in”) without mentioning specific devices, the system might struggle to generate a reasonable action plan involving smart devices due to the lack of a predefined

command. Consequently, users and developers must add new command-action pairs manually to customize the assistant. The configuration process [8] mainly involves setting up triggers (e.g., specific times and conditions) and defining corresponding actions (e.g., starting appliances and predefined routines), which are complex and time-consuming for application developers, let alone novice users.

To overcome these limitations, recent works integrate LLMs [9], [10] to revolutionize smart home services, enabling assistants to understand user intents beyond predefined commands. Among them, Sasha leverages ChatGPT [11] to generate action plans in response to under-specified user commands. SAGE further enhances user experience by storing the conversation histories for personalized plan generation. Nevertheless, these cloud LLM-based assistants introduces substantial privacy risks. In practice, users typically need to register for API keys to access cloud services, providing identifiable details such as a user ID and email address. During operation, in-home commands, user profiles, and home device states (e.g., temperature settings, lighting conditions) are transmitted to cloud for processing. Under the *honest-but-curious* threat model [12], [13], [14], [8], this workflow may expose user privacy [15], [16], including daily routines (e.g., cooking, exercising), personal preferences, and detailed home configurations, to third parties.

To protect user privacy from being exposed, an alternative approach is to exploit open-source models to serve smart homes locally. Though promising, local devices can only support small-size language models (SLMs) due to resource constraints. [17], [18], [19] Our preliminary study (§ III-B) reveals that SLMs often fall short in fully comprehending user intents [20]. Therefore, users are facing a *performance-privacy dilemma*: while cloud LLMs excel in delivering high-quality services, they may raise privacy concerns; conversely, local SLMs secure user privacy but fail to generate satisfactory responses due to limited model capabilities.

To address this dilemma, we propose *HomeLLaMA*, a privacy-preserving local home assistant that delivers personalized and satisfactory services through continuous learning. The key insight of *HomeLLaMA* is empowering local SLMs with the capabilities of cloud LLMs to shift most privacy-sensitive query processing tasks from the cloud to the local, thereby achieving a balance between model performance and user privacy. The powerful cloud services are consulted with users’ explicit approval only when necessary (e.g., unsatisfactory responses of the local SLMs). While the basic idea is simple, several technical challenges must be addressed.

X. Huang, L. Shen, Z. Ma, Y. Zheng are with the Department of Computing, The Hong Kong Polytechnic University, Hong Kong SAR, China. E-mail: unixy-xinyu.huang@connect.polyu.hk, leming.shen@connect.polyu.hk, zijing.ma@connect.polyu.hk, and csyqzheng@comp.polyu.edu.hk.

Corresponding author: Yuanqing Zheng.

- *SLMs perform poorly compared to LLMs and lack high-quality datasets for effective enhancement.*

Preliminary experiments (§ III-B) reveal the key bottleneck of SLMs in delivering high-quality smart home plans lies in their limited capabilities to accurately associate relevant devices with user commands compared with cloud LLMs. Yet further experiments show that directly tuning SLMs on existing command-action pairs only yields slight performance gains due to inadequate generalizability across heterogeneous home configurations. To address it, we propose a novel labor-free data augmentation method with a tailored inference paradigm. Specifically, we instruct powerful cloud LLMs to synthesize a generalizable command-device dataset based on available crowdsourced data. We fine-tune local SLMs on such a synthesized dataset and guide them using the consistent inference pipeline with well-crafted prompts. As a result, the fine-tuned local SLMs can effectively generate higher-quality action plans that are applicable across diverse homes with varied device configurations.

- *Even a well-enhanced SLM may not consistently provide cloud LLM-level services for users.*

As shown in the preliminary results (Fig. 2(b)), even after fine-tuning with our well-constructed dataset, there still remains a substantial performance gap between the local SLM and the cloud LLM. This gap in serving smart homes may undermine user experience limited by local SLMs. To resolve this issue, we design a privacy-preserving local-cloud collaboration paradigm, providing users with the option to consult cloud assistance for higher-quality responses. During this collaboration, *HomeLLaMA* retains user preference and home configuration data locally, and further obfuscates the commands sent for assistance to preserve user privacy. The process is entirely user-driven, meaning that the privacy-sensitive commands will only be processed and then transmitted to remote servers for performance enhancement upon explicit user approval.

- *Limited local space for guaranteeing long-term personalized services.*

Prior work [21] embeds entire historical user-assistant conversations into prompts for personalization. However, open-source SLMs have a shorter context length (e.g., 8K tokens for LLaMA3) than commercial models and cannot incorporate long conversation histories into prompts. While the popular retrieval-augmented generation (RAG) [22] is promising in reducing context length by fetching relevant information from a database, merely storing all historical conversations in the database can lead to the continuous accumulation of preference-related data, resulting in redundancy and hindering the efficient retrieval of highly correlated information. To mitigate this, at the end of each conversation, we instruct the local SLM to distill the current chat into a concise user profile containing topics, preferences, the current command, and its final approved plan. Following the digestion of historical data, we design a dynamic profile updating mechanism based on similarity to reduce redundancy.

We implement and deploy *HomeLLaMA* on a local server concerning specific smart home layouts and evaluate its per-

formance across multiple commonly used scenarios (e.g., atmosphere adjustment, and energy management). Both quantitative experiments and sufficient user studies ($M = 100$) reveal *HomeLLaMA* significantly enhances user-centered privacy while maintaining an acceptable level of performance, alleviating the raised performance-privacy dilemma for smart home users. In short, the **contributions** are as follows:

- To the best of our knowledge, *HomeLLaMA*¹ is the first on-device smart home assistant to support privacy-preserving and personalized services via user-in-the-loop.
- *HomeLLaMA* features three key novel technical modules: *Local SLM Enhancement* for effectively enhancing the performance of local assistants with a tailored inference paradigm, *Local-Cloud Collaboration* for maximizing user experience via a user-centered local-cloud collaborative workflow with privacy considerations, and *User Preference Learning* for efficient locally-hosted long-term personalized services.
- We build a comprehensive smart home benchmark *DevFinder*² to quantitatively evaluate the plan quality of smart home assistants. Extensive experiments and user studies demonstrate *HomeLLaMA* can effectively offer satisfactory and privacy-enhanced services while boosting user-oriented privacy confidence.

II. RELATED WORK

A. Smart Homes

In recent years, smart homes have emerged as a significant area of interest within the broader domain of Internet of Things (IoT) [25]. These systems integrate various connected devices to automate and enhance home living, offering functionalities [26] such as energy management [1] and personalized services [27]. A typical scenario contains several smart devices, a user interface component, and a central processing unit that connects the smart home with cloud servers [28]. On the smart device side, recent research has focused on enhancing device capabilities through machine learning algorithms. For instance, [29] explores activity recognition for home automation by developing a deep learning algorithm that identifies user activities based on accelerometer data collected by devices. On the user interface side, voice-based assistants are increasingly preferred due to their ability to facilitate natural language interactions and hands-free control. Commercial products like Google Assistant [30], and Alexa [6] exemplify this trend, offering intuitive interfaces capable of managing various commands, such as shopping and setting reminders, to streamline automated device control. However, these modern home assistants usually struggle with implicit and complex commands [31], as demonstrated in our preliminary study. On the other hand, recent advances in LLMs have shown excellent performance in open-vocabulary question answering, which can better comprehend user intentions with under-specified commands. *HomeLLaMA* enhances user experiences with improved system performance by integrating LLMs with smart home devices to overcome the aforementioned challenges.

¹The trained model is available in Huggingface with an anonymous account: <https://huggingface.co/USER9724/HomeLlama-8B>.

²The dataset is available in Huggingface with an anonymous account: <https://huggingface.co/datasets/USER9724/SmartHome-Device-QA>.

TABLE I
A COMPREHENSIVE COMPARISON WITH OTHER LLM-BASED ASSISTANTS.

System	Base Model	Plan Quality	Personalization	Privacy Protection	User Engagement
HomeGPT [23]	GPT-3.5 (Cloud)	Medium	Limited	Low	Limited
Sasha [5]	GPT-4 (Cloud)	High	Limited	Low	Limited
SAGE [21]	GPT-4 (Cloud)	High	High	Low	Limited
TT-Gemma [24]	Gemma (Local)	Low	Limited	High	Limited
TT-Phi-2 [24]	Phi-2 (Local)	Low	Limited	High	Limited
HomeLLaMA	LLaMA3 (Local)	Medium	High	High ↑	Proactive

B. Integrate LLMs with Smart Homes

Recognizing the strong generalizability and language processing capabilities of LLMs [32], [33], [34], [35], researchers are attempting to integrate them with smart homes for enhanced user experiences. A pioneering work, HomeGPT [23], directly prompts LLMs to generate a series of routines for the smart devices by providing the user command with detailed device states. The routines will further be parsed to adjust the states of the smart devices accordingly. Sasha [5] further optimizes the inference and control procedures by dividing the entire process into five steps: clarifying, filtering, planning, execution, and feedback. Nonetheless, it cannot adapt to user habits to generate personalized action plans, lowering long-term user satisfaction. To address this, SAGE [21] and Jordan *et al.* [36] enable LLMs to incorporate user profiles for generating personalized plans. However, these systems transmit user data and smart home configurations to the cloud LLM for processing, raising privacy concerns for users as the data exits the local environment. To provide satisfactory and personalized plans while enhancing privacy, *HomeLLaMA* tailors a locally deployed SLM via fine-tuning, focusing on providing satisfactory and personalized smart home plans while enhancing user privacy through our designed *PrivShield*.

In summary, Table I qualitatively presents a comprehensive comparison between *HomeLLaMA* and other LLM-based smart home assistants across multiple dimensions. Each of these dimensions corresponds to specific quantitative metrics discussed in the evaluation section (§ VI), for example, *plan quality* is measured using the defined *Device Relevance Score*. Compared with cloud-based assistants, *HomeLLaMA* offers personalized services while significantly enhancing user privacy with minimal performance trade-offs. On the other hand, compared with local-based solutions, *HomeLLaMA* excels in providing superior personalization and higher-quality plans. Additionally, it favors an innovative interaction paradigm that promotes proactive user engagement through a user-driven user-assistant interaction chain, enabling users to actively personalize responses for a more adaptive experience.

III. MOTIVATION AND CHALLENGES

A. Limitations of Existing Smart Home Assistants

Existing solutions for smart assistants can be categorized into task-based and LLM-based assistants. Task-based assistants are trained on predefined command-action pairs, while

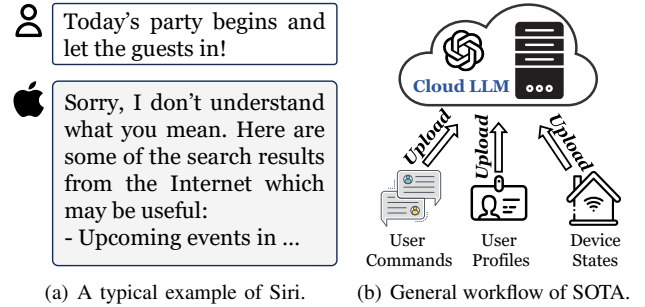


Fig. 1. Illustrations of existing works, including the typical example of task-specific models and the workflow of existing LLM-based models.

LLM-based assistants utilize the robust capabilities of LLMs to understand user intents in various smart home scenarios.

Limitations of existing task-based assistants. As a notable task-based assistant trained on a vast human-annotated dataset, Siri [37] can deliver excellent responses to predefined tasks [7]. However, its performance degrades when encountering unseen and complex commands. Fig. 1(a) illustrates a typical failure scenario in a conversation between Apple Siri and a smart home user. When the user inputs a command such as "Party begins and let all the guests in!" Siri fails to provide an appropriate response and instead directly returns the search results from the Internet, leading to a poor user experience.

Privacy concerns of LLM-based assistants. Recent advancements in LLM-based smart home assistants, such as Sasha and SAGE, allow users to issue commands more freely and receive responses that go beyond predefined tasks. Specifically, Sasha prompts the LLM using a designed pipeline with steps like clarifying, filtering, and planning to generate satisfactory action plans in response to user commands. Sasha fails to provide personalized services. On the other hand, SAGE further enhances personalization by storing conversation histories and summarizing them into user profiles.

Despite the improvements, as illustrated in Fig. 1(b), this workflow may pose significant risks to user privacy. In practice, users are required to transmit commands along with constructed user profiles and detailed device states (*e.g.*, a JSON file indicating the status of an air conditioner) to cloud servers for processing. Assuming an *honest-but-curious* cloud adversary [12], this workflow may lead to the exposure of:

- Sensitive personal information, including personally identifiable information (PII) and user preferences [38];

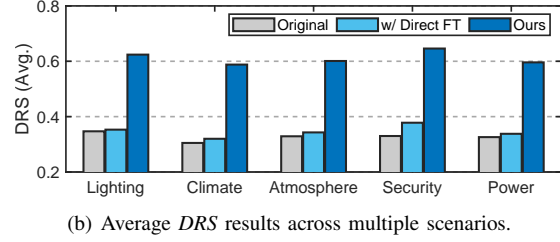
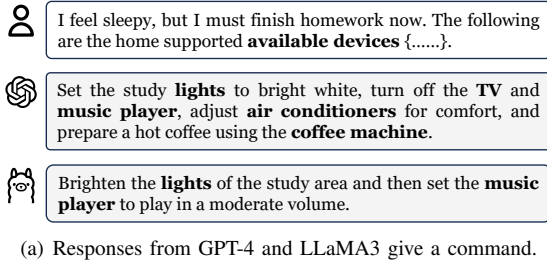


Fig. 2. Preliminary results of (a) responses from GPT-4 and LLaMA3 and (b) *DRS* after setting GPT-4 as references.

- Home configurations, *e.g.*, real-time home device states [39];
 - Users' daily in-home activities/routines, *e.g.*, exercising [40].
- These privacy risks hinder existing LLM-based assistants.

B. Challenges

As a straightforward solution to mitigate privacy concerns of existing LLM-based assistants, we conduct preliminary experiments by deploying the open-source LLaMA3-8B [41] on a local server. From the preliminary study, we report several technical challenges that further inspire *HomeLLaMA*.

Challenge 1: Vanilla SLMs perform poorly in identifying relevant devices and lack high-quality datasets for effective fine-tuning. To first uncover the underlying reasons why SLMs underperform LLMs in smart homes qualitatively, we input an under-specified command along with a set of available devices to both GPT-4 and LLaMA3 to generate responses. As shown in Fig. 2(a), GPT-4 involves a comprehensive list of relevant devices whereas LLaMA3 generates a simpler response, mentioning only lights and the music player. The result suggests that SLMs mainly lack the inherent capability and domain knowledge in identifying the latent semantic correlation between user commands and relevant devices.

Given this observation, a user-configured dataset from smart home platforms [42] is then collected for fine-tuning the SLM. Once tuned, we input the prepared test commands into both the original and the fine-tuned SLM in the same prompt format, generating two sets of relevant devices as responses. For a fair comparison, the same test commands are also processed using GPT-4 to produce reference device outputs. All responses are generated based on a predefined device set, thereby constraining the models from generating outputs in a freestyle manner. To quantify each model's ability to associate relevant devices with input commands, we adopt the *device relevance score (DRS)* defined in [43], with a detailed metric definition provided in § VI. As shown in Fig. 2(b), the comparison reveals that *DRS* values only exhibit a slight improvement (less than 10%) across various scenarios after tuning on the dataset. The minimal performance gain from the existing crowd-sourced dataset drives us to construct a high-quality dataset tailored for fine-tuning SLMs in smart homes. **Challenge 2: Even a well-enhanced SLM may not consistently generate cloud LLM-level responses.** We explore potential strategies to address **Challenge 1** and finally construct a fit-for-purpose dataset for effectively enhancing SLMs in the context of smart homes (§ IV-B). However, as demonstrated in Fig. 2(b), though fine-tuning SLMs with our tailored dataset significantly improves performance, there still remains

a gap between local SLMs and cloud LLMs. In practice, this implies that even with enhancement, SLMs may still fail to consistently deliver high-quality services to users in smart home environments. Such inconsistencies can diminish the user experience, particularly when compared to the robust cloud-based assistants. The gap highlights the need for a user-driven cloud-assisted mechanism—one that incorporates privacy-preserving measures—enabling users to obtain higher-quality responses when local outputs fall short of expectations. **Challenge 3: Simply storing all interaction history for personalization necessitates an excessively long context.** Existing approaches [21] directly incorporate raw conversation history or accumulated user profiles into prompts for cloud LLMs to enhance personalization. However, applying this method to local SLMs faces a unique challenge: local SLMs have a much shorter context length (*e.g.*, only 8K tokens for LLaMA3), making it infeasible to include lengthy user profiles in prompts. While the widely adopted retrieval-augmented generation (RAG) method [22] presents a promising solution for conserving context length by retrieving relevant information from a local database, its long-term use in smart homes may lead to the continuous accumulation of preference-related files. This accumulation can result in increased data redundancy over time, thereby hindering the effective retrieval of highly-correlated information. This practical limitation calls for innovative solutions to optimize the use of historical data.

IV. DESIGN OF HOMELLAMA

A. System Overview

To address the aforementioned challenges, we propose *HomeLLaMA*, with an overview outlined in Fig. 3. Before deployment, it begins with the offline *Local SLM Enhancement* module (§ IV-B) and with the enhanced model, user commands are further processed through the online stage, consisting of the *Multi-party Interaction* module (§ IV-C) and the *User Preference Learning* module (§ IV-D).

- **Local SLM Enhancement** enables the local SLM to generate plans for various user commands. Before deploying the local assistant, it is necessary to enhance the SLM so that it can identify relevant devices based on user commands. We begin by selecting seed commands from an open-source command-action dataset, covering various scenarios such as lighting, environment control, and security [42]. We then propose a tailored data augmentation method by feeding seed commands into a cloud LLM (GPT-4) to generate new commands, incorporating different expression styles (*user diversity*) and scenarios (*application diversity*). This process

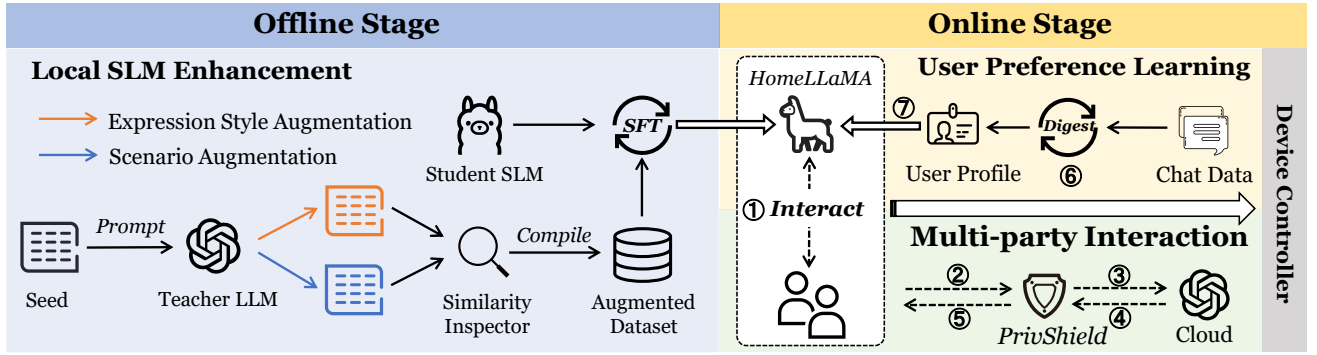


Fig. 3. System overview of HomeLLaMA. The system begins with an offline stage to enhance service quality within smart homes. Once deployed, it enters the online stage, where it continuously learns and updates user profiles in real time, with optional cloud assistance upon user request.

synthesizes a large set of user commands. Then, the teacher LLM labels the commands with comprehensive relevant devices and compiles them into an augmented dataset, which is further used to fine-tune the local SLM.

- **Multi-party Interaction** further enhances the user experience in the loop of user-assistant-cloud interactions. Users can interact with HomeLLaMA by freely expressing their requirements (①). If the response generated by HomeLLaMA falls short of expectations, users can give feedback or allow the local assistant to seek advice from cloud LLMs (②). To enhance user privacy, PrivShield obfuscates user commands by blending them with adversarial commands generated by the SLM before sending the mixture to a cloud-based LLM for processing (③). The cloud LLM, upon receiving these mixed queries, generates a set of responses and returns them to the local PrivShield (④). The real response corresponding to the original user command is then identified and recovered as advice, which the SLM integrates to provide users with a refined action plan (⑤).
- **User Preference Learning** ensures the assistant continuously learns and adapts to user preferences. Specifically, HomeLLaMA records each user-assistant interaction, which are then digested into structured user profiles with a pre-defined format (⑥). And the maintaining user profiles will dynamically update based on profile similarity. The module allows the assistant to retrieve user profiles to generate personalized plans for similar commands in the future (⑦). Over time, with the accumulated user profiles, HomeLLaMA becomes increasingly attuned to user preferences.

Remarks. HomeLLaMA focuses on generating satisfactory personalized action plans. In practice, the action plans can be automatically translated to executable files (e.g., JSON) to control smart devices via APIs (e.g., Apple API [44]) and we omit the automatic translation part in this paper.

B. Local SLM Enhancement

To improve SLMs in smart homes, a viable approach can be applying supervised fine-tuning (SFT) [41] on a tailored "user command → relevant devices" dataset. Inspired by recent advances in data augmentation methods (e.g., WizardLM [45]), we investigate the potential of leveraging powerful cloud LLMs (e.g., GPT-4) to automatically synthesize a customized dataset with higher quality. This approach effectively transfers

the knowledge embedded within the cloud LLM (teacher) to the local SLM (student) through the fine-tuning process.

1) *Understanding the dataset:* Serving different smart homes with diverse user groups is not a straightforward one-input-to-one-output mapping problem and two types of diversity need to be considered: **1) Command diversity.** It arises from two main aspects: user diversity and scenario diversity. User diversity refers to the fact that different users may express their requests in various ways, while scenario diversity refers to different types of home scenarios. **2) Device diversity.** It refers to the fact that different smart homes may have varying sets of available devices, leading to multiple possible responses for the same command.

2) *Command augmentation:* Concerning the issue of **command diversity**, it is essential to construct a dataset that includes a wide range of high-quality commands across different user groups and various scenarios. We begin this process by manually selecting a set of under-specified commands as the seed from crowd-sourcing platforms (e.g., IFTTT) [42] based on their popularity, i.e., overall adoption frequency among users. The selected commands encompass several commonly used smart home scenarios, such as climate control and lighting control. Each scenario contains 10 commands and we obtain 90 seed commands in total.

Synthesis of new commands. Harnessing the strong generative capabilities of cloud LLMs allows us to expand the dataset without the need for manual data collection. During each iteration of synthesis, we randomly sample five commands from the command pool as a starting point. Motivated by the two aspects of command diversity, we proceed to augment the original commands along the following two directions:

- *Vertical synthesis* generates new commands for different smart home scenarios. With the sampled seed commands, we first instruct GPT-4 to generate a new yet relevant command considering a different scenario via our carefully designed prompts (Fig. 4(a)). With the newly obtained command, we feed it back into GPT-4 to verify whether the command is indeed relevant to the smart home context. If not, we discard the command and proceed to the next iteration.
- *Horizontal synthesis* aims to generate new commands with varied expression styles. Similar to the vertical synthesis process, we instruct (Fig. 4(b)) GPT-4 to change the expression style of the original command while maintaining

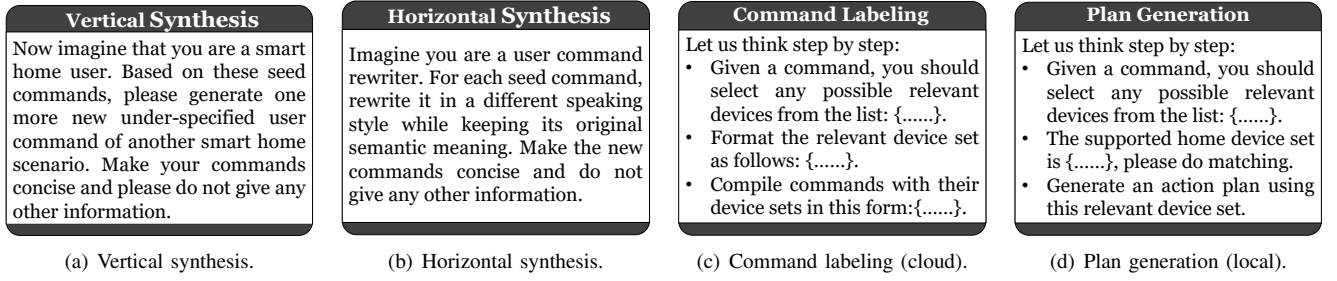


Fig. 4. The prompt template for (a) vertical and (b) horizontal synthesis, (c) command labeling, and (d) plan generation.

its original meaning. Once generated, we add the new command to the candidate command pool for further verification.

Similarity inspector. To ensure the quality of augmented commands, it is necessary to remove redundant commands with similar semantic meanings from the candidate pool. Specifically, given any newly generated command s_{new} , let the set of existing commands in the command pool be $S = \{s_1, s_2, \dots, s_n\}$. The ROUGE-L score function [46], denoted as $R(s, s')$, measures the similarity between commands. Then the retention condition for the new command is

$$\text{Retain } s_{\text{new}} \iff \max_{i \in \{1, 2, \dots, n\}} R(s_{\text{new}}, s_i) < \alpha \quad (1)$$

where α is a predefined threshold that controls the portion of overlap in semantic similarity between the new and existing commands. This means that a new command will be preserved only if the similarity between the new command and any existing command is less than the predefined threshold.

3) *Command labeling:* Given the augmented command pool, the next critical step is accurately labeling these commands to construct a comprehensive command-device dataset. We leverage the cloud LLM to label the commands with comprehensive device sets, encompassing all potential relevant devices. Specifically, we first simulate a virtual and large-scale smart home deployed with a comprehensive set of COTS devices (39 devices in total) collected from a smart home platform [42]. With the prompt designed in Fig. 4(c), we instruct the cloud LLM to identify a subset of relevant devices from the comprehensive set for each user command in the augmented dataset. The labeling process can be expressed as:

$$\mathcal{D}_a = \{s_i \rightarrow G(s_i, \mathbf{D})\}, \forall s_i \in \mathcal{S}_a \quad (2)$$

where \mathcal{D}_a is the augmented dataset, s_i is a user command from the augmented command pool \mathcal{S}_a , \mathbf{D} is the comprehensive device set we build, and $G(\cdot)$ represents the black-box LLM. **Remarks.** For the uncommon situation where a smart home contains a device not included in the comprehensive set, the user can explicitly suggest the missed device and specify her preference on how to adjust the device. Such explicit user feedback will be recorded and retrieved for future reference § IV-C. Note that the labeling process is agnostic to distinct device configurations and does not require transmitting specific user data to the cloud LLM.

4) *Training the adapter as the device identifier:* After obtaining the tailored command-device dataset, we proceed to fine-tune the local SLM to enhance its capabilities. Specifically, we utilize the QLoRA technique [47], a widely adopted

parameter-efficient fine-tuning (PEFT) [48] method. Instead of fine-tuning all model parameters, which is both resource-intensive and time-consuming, QLoRA trains a lightweight adapter integrated into the target model. In the smart home context, this process involves training a LoRA adapter to act as a device identifier for the local SLM. By combining this adapter with the original SLM, which retains extensive world knowledge, *HomeLLaMA* becomes more adept at accurately identifying and interacting with various smart devices.

5) *Inference paradigm:* With the enhanced SLM, we then propose a tailored inference paradigm for serving each individual home concerning the **device diversity** based on Chain-of-Thoughts (CoTs) [49]. Fig. 4(d) illustrates our designed prompt that instructs the SLM to generate the corresponding plans in a step-by-step manner. The inference paradigm can be divided into two steps outlined in Fig. 5:

- Initially, we consider a large home equipped with almost all COTS devices as mentioned before. Then, we prompt SLM to generate a comprehensive list of relevant devices given a command. We denote the comprehensive relevant device set as \mathbf{D}_l , as shown in the red box of Fig. 5.
- Then, the generated results are adapted to a specific home by performing a matching process. Specifically, with the available device set in home i denoted as \mathbf{D}_i , we prompt the SLM with the instructions in Fig. 4(d) to execute the task, matching the common devices of \mathbf{D}_l with \mathbf{D}_i to obtain the matched set \mathbf{D}_f^i for home i (as shown in the blue box of Fig. 5). The matching operation via the SLM is:

$$\mathbf{D}_f^i = \mathbf{D}_l \cap \mathbf{D}_i. \quad (3)$$

Remarks: In the initial step, while it is feasible to directly input the device set of a target home to generate the action plan, this approach may result in performance degradation. The main reason is that the local SLM is fine-tuned on our tailored dataset with a predefined input format. Therefore, the enhanced ability in relevant device identification may only be activated when the prompt aligns with the expected format.

C. Multi-party Interaction

As mentioned in § III-B, the enhanced SLM may still fail to consistently offer high-quality services in practice. To further improve user experience, we propose a *multi-party interaction* module that facilitates user feedback as well as consultation of cloud LLMs when necessary. This module supports two types of interactions: 1) the user-assistant interaction, which allows users to explicitly specify their requests and preferences, and

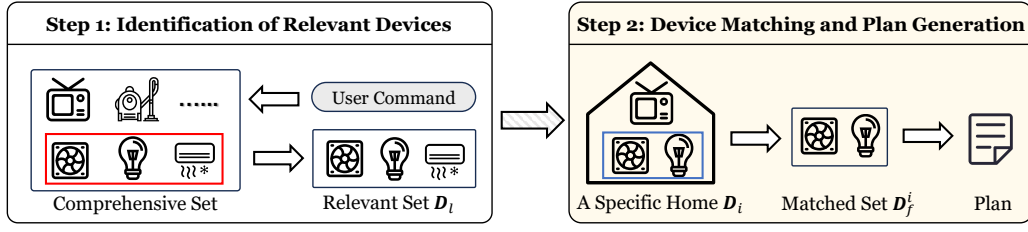


Fig. 5. The designed inference paradigm of the local SLM.

2) the user-driven local-cloud collaboration, where the local SLM is triggered by users to consult cloud LLMs for better services with privacy protection.

1) *User-assistant interaction*: As the core component of *HomeLLaMA*, the user-assistant interaction acts as an interface for users to express their intents. In the flowchart illustrated in Fig. 6, upon receiving a user command, the assistant will generate action plans and respond to the users for confirmation. For every generated response to users, they may accept, reject, or follow up with a piece of advice.

- **Accept**: If the user is satisfied with the generated action plan, the action plan will be translated into the command to smart devices to control relevant devices.
- **Advice**: The user can provide natural language feedback to further refine the user intent. For example, suppose the user inputs a command like "Brighten the bedroom," and the assistant responds with "Turn on all the lights in the bedroom." If the user only wants to turn on the bedside lamp, she can follow up with a detailed instruction (e.g., "Bedside lamp only, please."). The assistant will then re-generate the action plan by incorporating the user's advice into the prompt.
- **Reject**: If the user is not satisfied with the local response, she may reject the action plan. For instance, if the user wants to hold a home party and inputs "Let the party begin," but the assistant responds with a simple action like "Turn on the lights and adjust the room temperature," the user might reject the response. In that case, the assistant leverages the cloud LLM to generate an improved action plan with the user-driven local-cloud collaboration module (§ IV-C2).

Remarks. Note that after each generated response, including revised plans resulting from "Advice" or "Reject," user can further interact with the assistant. Only when the user explicitly "Accept" a proposed action plan, the plan will be translated to control smart devices accordingly (Fig. 6). Subsequently, the command and the final approved plan will be saved as an interaction record, which will be further utilized by the preference learning module (§ IV-D).

2) *Local-cloud collaboration*: When user rejects a plan, *HomeLLaMA* will ask user for permission to consult a cloud model (e.g., GPT-4). If approved, the assistant will proceed the process to generate an improved response [50].

Potential privacy risks. However, directly querying the cloud LLM via registered API calling with the raw user command and home details may raise privacy concerns [51] since user activities and personal information may be inferred and monitored [52] indirectly (e.g., through differential attacks [53]) by the curious cloud servers. For example, suppose a user first

requests, "At 9 pm, make my living room chilly and turn on the TV," followed by, "At 10 pm, check if the doors and windows are locked and make my bedroom comfortable." From these commands, it can be easily inferred that the user might be watching TV from 9 pm to 10 pm and then go to bed.

Role of *HomeLLaMA* during collaboration. While the goal of local-cloud collaboration is to enhance user experience, it must not come at the cost of unacceptable privacy compromises. To this end, *HomeLLaMA* functions as both a processing center and a privacy guardian: it retains all smart home configurations and user profiles locally, transmitting only the current user command to the cloud for assistance. To further mitigate potential privacy risks embedded in raw commands, we incorporate *PrivShield*, a lightweight obfuscation module designed to anonymize user queries before cloud interaction, as illustrated in Fig. 7.

***PrivShield*.** Essentially, the *PrivShield* operates within a **SLM-in-the-middle** framework. In practice, *PrivShield* safeguards user privacy through procedures including user command rewriting, adversarial command generation, and plan recovery.

- *User command rewriting*. An original user command may contain personal information (e.g., names, locations) and many colloquial expressions (e.g., modal particles). These components not only introduce information redundancy but also provide opportunities for third-parties to infer the user's actual command through continuous pattern recognition in subsequent processes. To address this, we direct the local SLM to first filter sensitive personal information [54], and then paraphrase the original user command using the customized prompt illustrated in Fig. 8(a).
- *Adversarial command generation*. Given a paraphrased command, the *PrivShield* prompts the local SLM with the designed instructions in Fig. 8(b) to generate other N adversarial commands across various unrelated scenarios to obscure the original command. Each of the commands is assigned a unique command ID and shuffled, with only the original command's ID t being locally recorded. These commands are subsequently combined into a single query along with their respective command IDs, as shown in Fig. 8(c). The combined query will be transmitted to a cloud LLM to generate action plans for all the commands.
- *Action plan recovery*. Upon receiving the response from the cloud LLM, the *PrivShield* extracts the comprehensive action plan associated with the right order. This extracted action plan is then fed into the local assistant as advice for generating a tailored plan for the user. The tailored plan is subsequently delivered to the user as the updated plan.

Remarks. *PrivShield* enables users to access cloud services

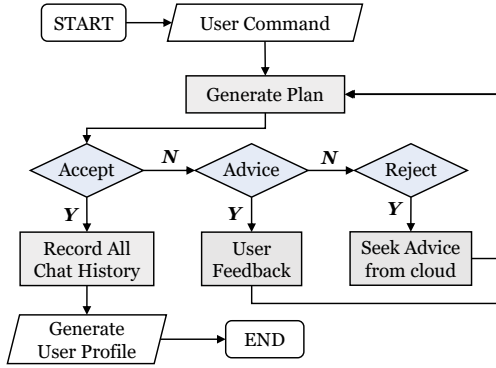


Fig. 6. The user-assistant interaction flow.

with privacy protection in an easily understandable manner. However, the assistant primarily operates locally whenever possible. The reasons are twofold: 1) User profiles and home configurations are stored locally and will not be transmitted to the cloud for processing due to privacy concerns. 2) The cost of constantly querying the cloud may be prohibitive. Before deployment, users are allowed to customize the number of adversarial commands, *i.e.*, N , to achieve user-oriented privacy-cost balance.

D. User Preference Learning

Due to the restricted context length and information redundancy, local SLMs cannot simply store all the chat history for generating personalized responses. To address this challenge, we develop a lightweight user profiling method, enabling the assistant to efficiently retrieve dynamically updating user profile database for reference. In practice, the user preference learning module operates in three key stages: user profile generation, profile updating, and personalized plan generation.

1) *User Profile Generation*: The interaction records between the user and the assistant are locally recorded. A structured prompt, as illustrated in Fig. 9(a), guides the local SLM to digest and generate a well-organized user profile for each conversation. These profiles include details on ① **topics** (*i.e.*, the keywords of conversations summarized by the SLM), ② **preferences**, ③ **commands**, and ④ **final action plans** in a concise way. These profiles are then transformed into vector representations and stored in a text embedding database \mathcal{E} .

2) *Profile Updating*: The user profile database follows a carefully designed updating mechanism to maintain its effectiveness over time. When a new user profile is generated, it is compared with all existing profiles via cosine similarity. If the similarities between the new profile and all the existing profiles are below a pre-defined threshold, the new profile will be saved as a distinct entry in the database. Otherwise, it is constructively merged with the most similar existing profile. Specifically, given the embedding of a newly generated user profile denoted as p_n , the condition for inserting this profile into the embedding database is determined by:

$$\text{Insert } p_n \text{ into } \mathcal{E} \iff \max_{p_i \in \mathcal{E}} C(p_n, p_i) < \beta \quad (4)$$

where $C(\cdot)$ is the cosine similarity function and β is a pre-defined similarity threshold. If the maximum cosine similarity

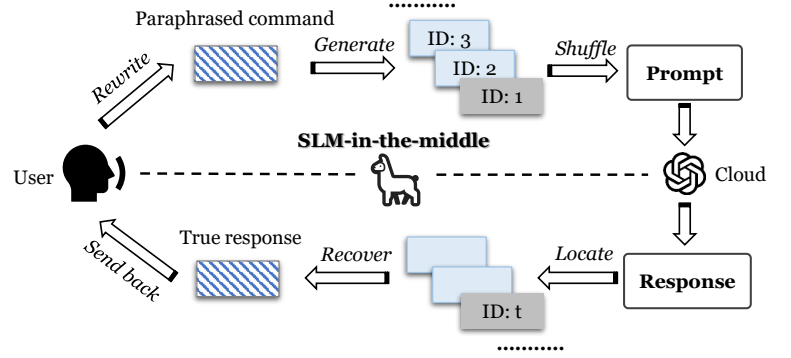


Fig. 7. Workflow of the PrivShield.

between p_n and any existing profile p_i in the database is less than β , the new profile will be inserted as a distinct entry. Otherwise, the two similar profiles are merged into a new, consolidated profile, which replaces the original profile by prompting the SLM with the designed prompts shown in Fig. 9(b). This process enhances the storage efficiency of the user profile database.

3) *Personalized Plan Generation*: During the inference stage, given a new query q_i , the assistant retrieves the top-3 user profiles in the form of text embedding (denoted as p_m , p_n , and p_p) that have the highest cosine similarity to the query. We then convert the selected embedding into text (denoted as u_m , u_n , and u_p) through decoding:

$$(p_m, p_n, p_p) \xrightarrow{\text{Decode}} (u_m, u_n, u_p) \quad (5)$$

By concatenating the decoding result with the user query q_i , the personalized output plan is generated based on these profiles and the current home configuration H_i :

$$\text{Plan} \leftarrow \mathcal{L}(u_m, u_n, u_p, H_i, q_i) \quad (6)$$

where \mathcal{L} represents the action plan generation function of the local SLM. The generated action plan is subsequently used to control the corresponding smart devices.

Remarks: The user preference learning module is designed to maintain a dynamically updating user profile database for personalization. Additionally, our designed user profile learning paradigm can be extended to multi-user smart homes by constructing separate databases for each user. During service, the assistant will first perform voice recognition [39] to determine the user identity before processing.

V. IMPLEMENTATION

We implement the designed *HomeLLaMA* on a local server using the PyTorch framework [55], with the entire workflow supported by LangChain [56]. The implementation details of each key component are as follows:

Data Augmentation: We prepare seed commands from IFTTT and access OpenAI's services via the OpenAI API. During the augmentation process, we select GPT-4-Turbo as the default LLM and appropriately prompt it to generate the required data. By default, we set α in Equation 1 to 0.7. The augmented dataset contains 14K action-command pairs in total, covering

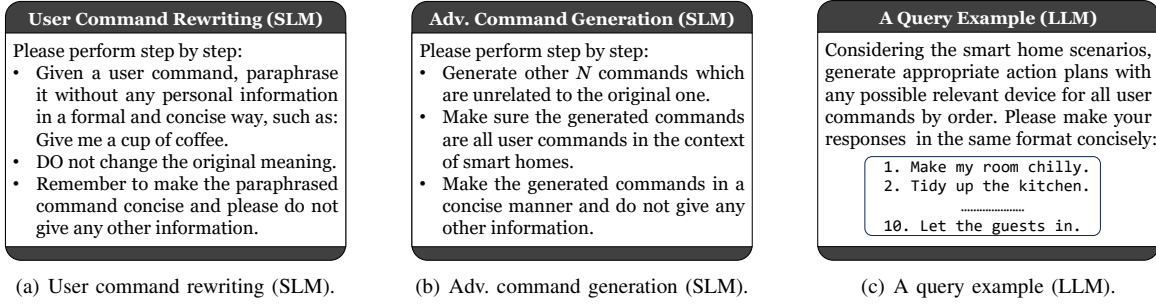
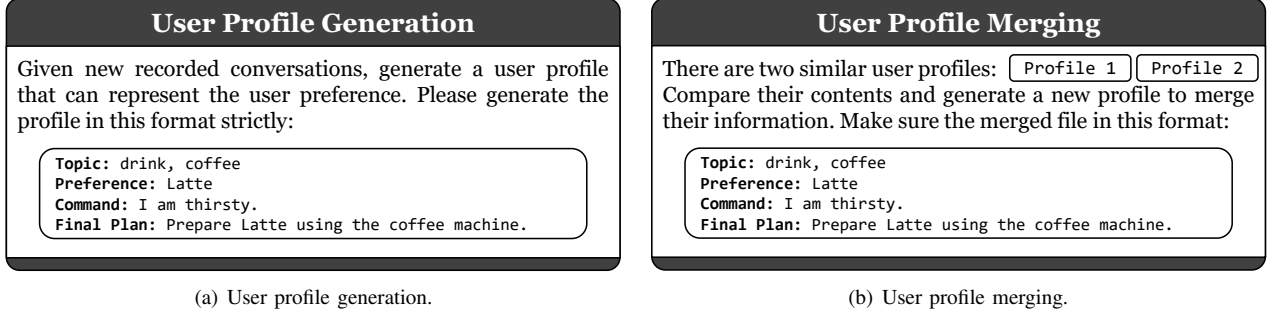
Fig. 8. The prompt templates of the designed *PrivShield*.

Fig. 9. The prompt templates for user profile generation and merging.

9 common smart home scenarios (*e.g.*, atmosphere adjustment, power management, *etc.*).

Fine-Tuning: We choose Meta-LLaMA3-8B [57] as our base model downloaded from Hugging Face³. To fine-tune the base model with our augmented dataset, we utilize the QLoRA [47] technique with 8-bit quantization. The rank r is set to 64 and lora_alpha to 128. The learning rate is initialized at 3×10^{-5} , and a dropout rate of 0.1 is applied to alleviate the over-fitting problem. The number of fine-tuning epochs is 3, with a train-test split ratio of 80-20. The model is fine-tuned on a server installed with Ubuntu 22.04 LTS with a single NVIDIA RTX 4090 GPU, taking around 8 hours.

User Profile Database: We deploy and maintain the user profile database via FAISS [58], which is a library for efficient similarity search and clustering of dense vectors. The conversation histories collected from the interactions are saved in the text format and summarized into user profiles. Those well-summarized user profiles are stored in the text embedding database, and ready to be retrieved for the generation of personalized plans during the inference stage. By default, the β in Equation 4 is set to 0.6.

VI. PERFORMANCE EVALUATION

In this section, we conduct comprehensive quantified experiments to evaluate the effectiveness of *HomeLLaMA* in addressing the performance-privacy dilemma. Specifically, we aim to answer the following questions:

- **Q1 - Performance:** Can *HomeLLaMA* provide high-quality services locally?
- **Q2 - Privacy:** Does *HomeLLaMA* quantitatively enhance user privacy?

- **Q3 - System Overhead:** Is *HomeLLaMA* affordable to be deployed locally?
- **Q4 - Sensitivity:** How do system configurations (*i.e.*, base models) impact performance?

A. Model Capacity (Q1)

1) **DevFinder Benchmark:** Considering the comprehensive home setup described in § IV-B, which is equipped with commonly used smart devices, we select 100 test commands with human-annotated device labels from the IFTTT dataset [42] with more details outlined in the open-source dataset. These commands with labels encompass a wide range of scenarios, including environmental control, atmosphere adjustment, power management, *etc.* The commands are input into the smart home assistants to generate responses, which are then compared with the annotated labels to quantify the quality of the generated action plans.

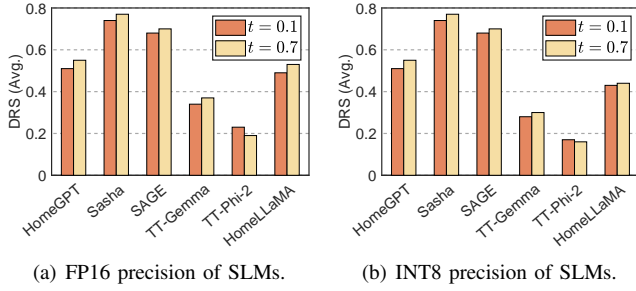
To quantify *HomeLLaMA*'s capability in identifying relevant devices, we adopt the Device Relevance Score (*DRS*) as the evaluation metric [5]. Suppose the ground truth device set is G_l and the device set generated by the local SLM is G_r , we compute the relevance score as:

$$DRS = \frac{|G_l \cap G_r| - |G_r - G_l|}{|G_r|} \quad (7)$$

where $|G_l \cap G_r|$ represents the number of overlapping devices, and $|G_r - G_l|$ refers to the number of devices in the response that are not included in the ground truth. Then, the relevance score is normalized to $[-1, 1]$. The higher the relevance score, the better the performance in identifying relevant devices.

2) **Baselines:** To contextualize the model performance of *HomeLLaMA*, we compare our system with several other LLM-powered smart home baselines:

³<https://huggingface.co>

Fig. 10. Avg. *DRS* after setting (a) FP16 and (b) INT8 precision.

- **HomeGPT** [23] directly prompts the LLM to generate smart home plans in response to user commands.
- **Sasha** [5] modifies HomeGPT by introducing a revised pipeline consisting of five procedures, including filtering, planning, *etc.*, to enhance the quality of plans.
- **SAGE** [21] generates personalized plans by inserting all conversation history into prompts.
- **Thoughtful Things (TT)** [24] utilizes the on-device SLMs (Google Gemma-7B [59] and Microsoft Phi-2-3B [60]) to generate routines and control smart devices.

Note that for a fair comparison, the cloud-assisted module is disabled during the evaluation of device relevance, and all results are computed purely based on local operations.

3) *Results*: We test *HomeLLaMA* and baselines on the proposed *DevFinder* and report the average score of each system. Additionally, we set the temperature t of these models to 0.1 and 0.7 for evaluation, respectively. As illustrated in Fig. 10, the designed *HomeLLaMA* significantly outperforms the other two on-device assistants but still lags behind cloud-based LLM assistants. To investigate the reasons behind this, we examine and analyze the generated action plans and uncover the following insights:

- Despite the superior performance of cloud-based assistants enabled by larger models [20], *HomeLLaMA* achieves comparable *DRS* to GPT-3.5 while ensuring user privacy through *PrivShield* and local processing, demonstrating high performance without compromising user privacy.
- Among on-device assistants, *HomeLLaMA* delivers the best local service, outperforming TT-Gemma and TT-Phi-2 in device relevance. This is due to its fine-tuning on smart home-specific data and its hybrid design, which enables selective cloud assistance to further boost *DRS* when needed.
- Raising the model temperature increases creativity and can improve relevance (e.g., *HomeLLaMA* from 0.51 to 0.55 at $t=0.7$ under FP16), but may also cause hallucinations [66]. While mitigations exist, smaller models like TT-Phi-2-3B still suffer performance drops due to weaker reasoning when temperature rises.

B. Privacy Protection (Q2)

1) *Qualitative risk analysis*: Before quantifying *HomeLLaMA*'s privacy protection, it is important to first examine potential privacy risks qualitatively. In typical cloud-based smart home systems, users must register for API keys and

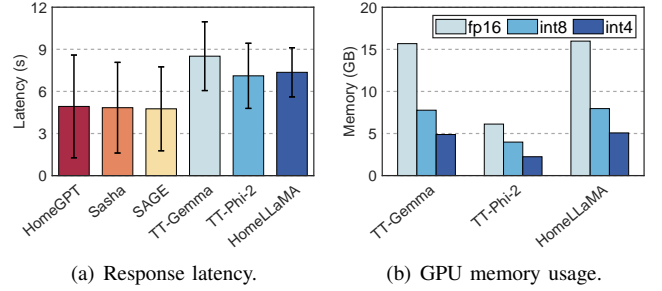


Fig. 11. Results for (a) latency and (b) memory usage.

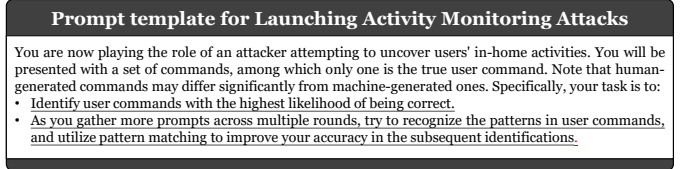


Fig. 12. The prompt template for launching activity monitoring attacks.

transmit commands and device states to remote servers, exposing them to risks during data storage, network transmission, and inference, as detailed in Table II. These include unauthorized access, interception, and inference attacks such as PII extraction, attribute inference, and activity monitoring. Unlike such systems, *HomeLLaMA* operates locally without sharing user profiles or home configurations, thereby mitigating storage and transmission risks. Additionally, its *PrivShield* module filters sensitive content before sending prompts to the cloud, preventing PII exposure. However, obfuscated prompts remain susceptible to activity monitoring via API traceability. The next section presents quantitative experiments to evaluate resilience against this residual threat.

2) *Quantitative Analysis*: We focus on examining the in-home activity monitoring threat in *HomeLLaMA*. During the use of *PrivShield*, the real commands are obfuscated with N other SLM-generated adversarial commands before being sent to the cloud LLM for processing.

Threat Model. We assume the cloud LLM operates on the *honest-but-curious* remote server [31], where adversaries deliver the correct inference results but investigate all transmitted user queries. The goal is to identify the real query from the mixture using pretrained classifiers, thereby enabling real-time monitoring of users' in-home activities. Following the procedures in [67], we launch the activity monitoring attack by instructing cloud GPT-4 using well-designed prompts shown in Fig. 12 to infer user in-home activities based on the continuously received commands. We then use *attack success rate* [40] to assess the system's privacy level, where a higher rate indicates a greater threat to user privacy and reduced system protection.

Results. We use the constructed *DevFinder* as test user queries, setting the number of adversarial commands N to 2, 4, 9, and 19. We also vary the base SLMs in *HomeLLaMA* to examine their impacts on the quality of adversarial command generation. The average attack accuracy results across query rounds obtained through extensive experiments are presented in Fig. 13, and we report the following key findings:

TABLE II
QUALITATIVE RISK ANALYSIS OF LLM-BASED ASSISTANTS INCLUDING CLOUD-BASED SYSTEMS, LOCAL TT AND *HomeLLaMA*. HERE "●" INDICATES COMPROMISING PRIVACY REGARDING THIS THREAT, WHILE "○" SIGNIFIES IT REQUIRES QUANTITATIVE EVALUATION (§ VI-B2).

Threat Type	Cloud-Based Systems	TT (Local)	HomeLLaMA (Hybrid)
Data Storage [61]	●	○	○
Network Transmission [62]	●	○	○
PII Extraction [64]	●	○	○
Inference [63] Attribute Inference [65]	●	○	○
Activity Monitoring [40]	●	○	●

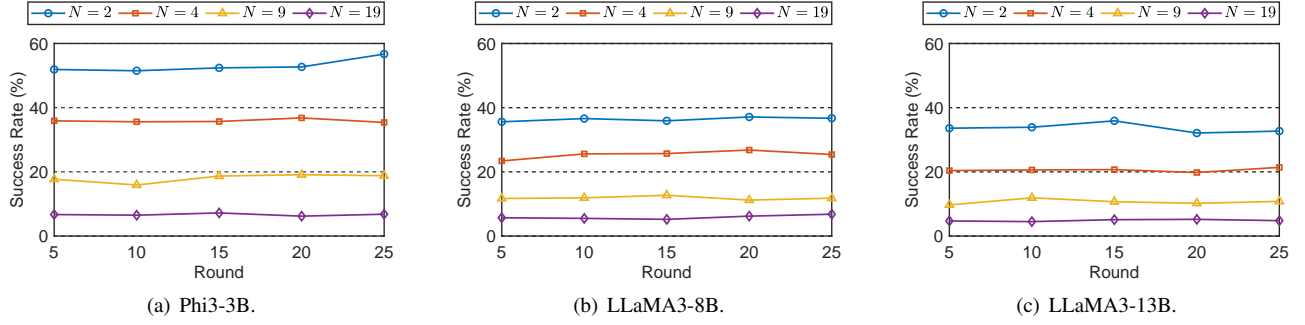


Fig. 13. Average attack success rate across query rounds for different values of N and base models.

- *PrivShield* effectively safeguards user privacy by maintaining significantly lower attack success rates compared to direct queries without its protection. As shown in Fig. 13, leveraging *PrivShield* with different N values and various base models results in a substantial reduction in attack accuracy, far below the 100% success rate of direct queries. Furthermore, while an increase in query rounds allows adversaries to accumulate more information, as illustrated in all sub-figures, this does not enhance their ability to accurately infer user prompts. In fact, denote the average attack success rate of the *PrivShield* as SR_p , and the overall attack success rate of *HomeLLaMA* SR_h should be

$$SR_h = \epsilon \cdot SR_p \quad (8)$$

where ϵ represents *PrivShield*'s frequency of use. During practical daily usage, the frequency ϵ will gradually become lower with user profiles being progressively constructed, as discussed in § VII-B. Consequently, the overall privacy protection of *HomeLLaMA* will strengthen over time, as users will increasingly rely on local operations without cloud LLM assistance.

- Privacy protection strengthens as the number of adversarial commands increases. These adversarial commands can **introduce noises in the text space**, making it harder for malicious attackers to identify the real queries. Moreover, generating more adversarial commands incurs higher latency and cost, users should be allowed to decide their preferred trade-off between privacy and performance.
- Privacy protection also benefits from the use of stronger SLMs. As demonstrated in Table II, replacing base SLMs with larger and more powerful models significantly reduces attack accuracy. This is because **identifying real user queries becomes equivalent to distinguishing AI-generated text from the mixture**. Stronger SLMs are more adept at generating high-quality adversarial commands, further obscuring the real query from identification. However,

deploying larger models may be impractical due to resource constraints, which will be discussed further in § VI-D.

Remarks. Although adversaries may deploy advanced pre-trained classifiers to distinguish user commands from obfuscated mixtures. In such cases, *PrivShield* could be enhanced by strengthening query obfuscation (e.g., selecting a larger N) for stronger privacy protection according to users' requirements.

C. System Cost (Q3)

1) *Metrics*: We evaluate the system cost in terms of the following aspects: **1) Response latency**: We measure the time cost from the moment the user inputs a command to the generation of the final action plan as the response latency of each system. **2) Memory usage**: We measure the system overhead by tracking the GPU memory usage (in GB) via a Python package named memory-profiler [68] during usage.

2) *Baselines*: We keep the same baselines as selected in § VI-A2. Note that the first three systems (i.e., HomeGPT, Sasha, and SAGE) are based on cloud LLMs which cannot be directly accessed, while only TT explores the integration of SLMs into smart home assistants. Consequently, we only measure the memory usage of TT-Gemma, TT-Phi-2, and our proposed *HomeLLaMA*.

3) *Results*: To measure the system overhead of *HomeLLaMA* and the baselines, we input each test command in the *DevFinder* benchmark into them and report the average response latency with its variance. As shown in Fig. 11(a), the cloud-based assistants exhibit relatively faster average response time (around 4.97 seconds) compared to the local-based assistants, primarily due to the performance optimizations of OpenAI services [11]. However, cloud-LLM-based systems are highly susceptible to network conditions and server stability, leading to significant variance in response latencies, which can negatively impact user experience [69]. In contrast, *HomeLLaMA* exhibits less variance in response

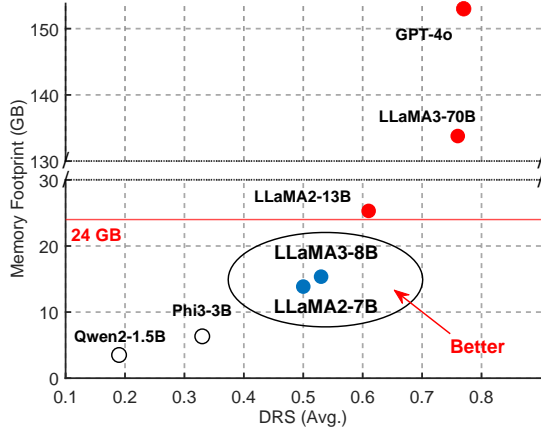


Fig. 14. Impacts of different base models.

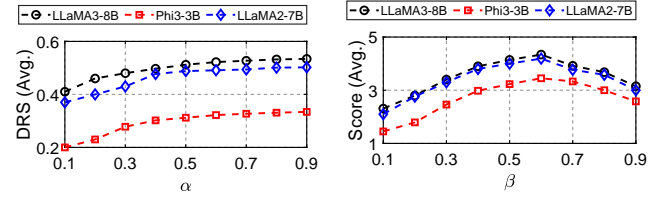
latency, albeit with a slightly longer response time. We also track the maximum GPU memory usage of the two local SLM-based systems (*i.e.*, TT and *HomeLLaMA*) when adopting different quantization precisions (*i.e.*, fp16, int8, and int4) during inference. As shown in Fig. 11(b), the maximum GPU memory requirement for these systems remains under 16 GB, which is affordable and manageable for a typical household⁴.

D. Sensitivity Analysis (Q4)

1) *Different base models*: We evaluate the impact of base models by deploying Qwen2-1.5B [70], Phi3-3B [71], LLaMA2-7B/13B [72], LLaMA3-8B [57], and LLaMA3-70B [72] in *HomeLLaMA*. Given the test inputs from *DevFinder*, we record their average *DRS* and GPU memory usage. As shown in Fig. 14, models fall into three groups: (1) lightweight models (white) such as Qwen2-1.5B and Phi3-3B offer minimal memory usage (~ 8 GB) but poor performance; (2) high-end models (red) like LLaMA2-13B, LLaMA3-70B, and GPT-4o yield the best *DRS* but require ≥ 24 GB GPU memory; (3) mid-range models (blue) such as LLaMA2-7B and LLaMA3-8B balance performance and cost. Thus, we choose LLaMA3-8B as the default base model.

2) *Different α during augmentation*: To study the effect of the augmentation threshold α (Eq. 1), we vary it from 0.1 to 0.9 and fine-tune LLaMA3-8B, Phi3-3B, and LLaMA2-7B on corresponding augmented datasets. As shown in Fig. 15(a), increasing α improves average *DRS* by promoting data diversity. However, the gain saturates at higher values, while training cost rises due to dataset expansion. We therefore set $\alpha = 0.7$ by default.

3) *Different β during profile updating*: The profile update threshold β (Eq. 4) determines when to save a new user profile. Using 10 participants, we evaluate personalization ratings (1–5) under varying β from 0.1 to 0.9. As shown in Fig. 15(b), higher β initially improves satisfaction by preventing premature profile merging, but overly high values lead to redundant entries, impairing retrieval. Satisfaction peaks around $\beta = 0.6$, which we adopt as the default.

Fig. 15. Impacts of selecting different thresholds α and β .

VII. USER STUDY

We conduct user studies to gather user feedback. The study can be divided into two parts: an *online survey* and an *onsite interview*, aiming to answer those questions:

- **Q5 - User Privacy Confidence**: Does *HomeLLaMA* lift user-perceived privacy confidence?
- **Q6 - User Satisfaction**: How do users feel about the overall service quality of *HomeLLaMA*?
- **Q7 - Long-term Personalization**: Is *HomeLLaMA* capable of adapting to users' preferences continuously?

A. Online Survey: Cold-start Evaluation (Q5 & Q6)

1) *Preparation works*: We select two representative base-lines: a cloud-based assistant SAGE and a local-based assistant TT-Gemma. For each scenario, we select a test command from *DevFinder* and input them into systems to generate initial action plans. Subsequently, we manually select the "Advice" option (§ IV-C1) and provide the feedback to all systems. These systems then perform preference learning, if applicable, and regenerate action plans which are recorded for further analysis. To ensure fairness, the order of system presentation is randomized before being rated by participants.

2) *Participants*: We created an online survey using Microsoft Forms⁵ and distributed it via emails and social media to recruit participants. Before participating, all respondents were presented with an informed consent form outlining the purpose of the study, data handling procedures, and their rights as participants. The study protocol was reviewed and approved by our institution's ethics review board (IRB). Participants were informed that their responses would be anonymized and used solely for academic research purposes. Participation was entirely voluntary, and no monetary compensation was provided. After 12 days, we received **100 responses** from volunteers in total, and the data of participants shows a relatively balanced distribution of participants in terms of gender, age, educational background, English proficiency, and familiarity with smart home assistants.

3) *Survey design*: The survey evaluated smart home systems from the following four metrics: ① **General Plan Satisfaction**, assessing satisfaction with the quality of initial action plans irrespective of personalization; ② **User Profile Correctness**, measuring how accurately the generated user profile reflects personal preferences and behaviors during the current interaction; ③ **Personalization Fit Score**, evaluating how well the responses and recommendations align with

⁴An NVIDIA RTX 4070 Ti SUPER GPU (16 GB) costs around \$840.

⁵<https://forms.office.com/>

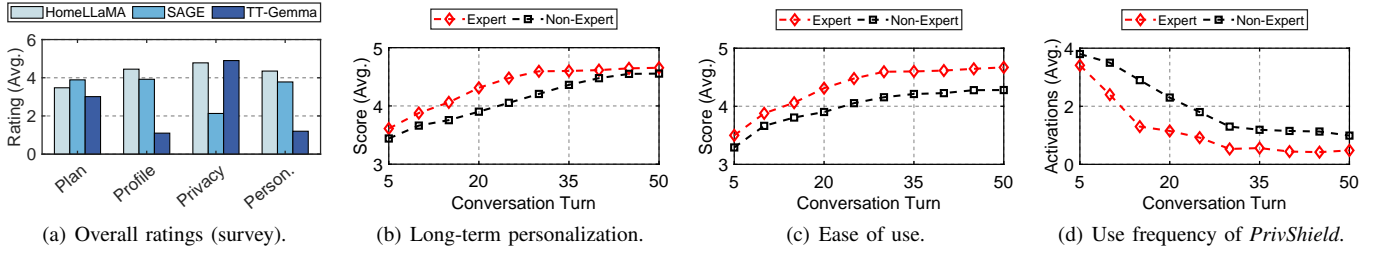


Fig. 16. User study results, including (a) online survey ratings, and (b)–(d) show average interview results.

historical feedback within the interaction round; and ④ **Privacy Assurance Score**, gauging participants’ confidence in the system’s ability to safeguard personal data and maintain privacy. Participants were asked to rate these metrics on a Likert scale [73] from 1 (*e.g.*, not at all) to 5 (*e.g.*, completely).

4) *Overall results*: Fig. 16(a) visualizes the overall rating results, and key observations of each aspect include: 1) Participants expressed high satisfaction (close to cloud-based assistant SAGE) with the quality of services provided by *HomeLLaMA*, reflecting its ability to deliver effective and reliable action plans tailored to user needs. This may be attributed to the tailored enhancement method of SLMs proposed in this paper. 2) *HomeLLaMA* excelled in delivering personalized responses and generating precise user profiles compared with other baselines, achieving the highest average score of around 4.35 points. This may be attributed to the *User Preference Learning* module, which accurately characterizes well-structured user profiles to adapt plans effectively. 3) *HomeLLaMA* received high ratings for its privacy-preserving features, surpassing cloud-based solutions and approaching the level of the fully local system, TT-Gemma. This demonstrates its ability to enhance user-perceived privacy by operating locally and minimizing data transmission to the cloud. Additionally, the local-cloud collaboration paradigm, supported by the designed *PrivShield*, boosts users’ confidence in privacy, thereby improving the overall usability of the system.

B. Onsite Interview: Long-term Evaluation (Q7)

To evaluate the long-term performance of *HomeLLaMA*, we deploy the system locally for conducting an onsite interview. The interview lasts for 25 days with 50 conversation turns, and the actual usage time of volunteers is approximately 45.6 minutes on average. Each turn represents a complete user-assistant-cloud interaction (Fig. 6), starting with the user command and ending with the final action plan.

1) *Procedures*: We deployed *HomeLLaMA* on a laboratory PC using the configuration in § V. From 100 survey respondents described in § VII-A1, 10 participants were invited and divided into two groups: 5 experts and 5 non-experts, based on their familiarity with smart homes. Prior to interviews, all participants provided informed consent and received compensation in the form of supermarket coupons valued at approximately \$10. Participants were first introduced to the interaction workflow—including the accept, advise, and reject options—to ensure familiarity. Additionally, two evaluation metrics were explained to them before beginning the interaction: ① **Long-term Personalization** assesses how

effectively the system infers and retains user preferences over time. ② **Ease of Use** evaluates how efficiently the system minimizes user efforts for delivering satisfactory responses as usage continues. Participants were then invited to freely interact with *HomeLLaMA*, issuing smart home commands in natural language through UI without constraints. After every five conversation turns, they rated the system using the two predefined metrics on a 5-point scale. Throughout the session, we recorded all evaluation scores and, every five turns, also tracked the number of times *PrivShield* was activated for cloud assistance. Finally, the average results for each group were computed and analyzed separately.

2) *Results*: As illustrated in Fig. 16, both expert and non-expert participants show a steady increase in ratings for *HomeLLaMA* across personalization and ease of use as the number of conversation turns grows, reflecting the system’s ability to adapt through dynamically maintained user profiles. Experts tend to assign higher scores earlier, likely due to their clearer articulation of preferences, which accelerates profile quality and system adaptation. A consistent gap remains in ease-of-use ratings, with experts maintaining an advantage of about 0.35 points by the 50th turn, though both groups converge above 4.2, indicating strong usability. Additionally, the activation frequency of *PrivShield* declines for both groups, approaching zero by the 50th turn, highlighting *HomeLLaMA*’s reduced reliance on cloud-based support as it better internalizes user preferences, thereby enhancing privacy protection.

VIII. CONCLUSION

This paper presents an on-device smart home assistant that balances privacy and performance for users. The designed *HomeLLaMA* comprises three technical modules: *Local SLM Enhancement*, *Multi-Parity Interaction*, and *User Preference Learning*, enabling seamless and privacy-enhanced interactions involving user-in-the-loop. We construct the *DevFinder* benchmark to assess the quality of the generated responses. Comprehensive user studies demonstrate that *HomeLLaMA* delivers satisfactory plans with enhanced personalization, while alleviating privacy concerns. Additionally, extensive quantitative experiments verify the effectiveness of *HomeLLaMA* in enhancing user privacy in data storage, network transmission, and inference-related (*e.g.*, PII extraction) attacks.

REFERENCES

- [1] B. L. R. Stojkoska and K. V. Trivodaliev, “A review of internet of things for smart home: Challenges and solutions,” *Journal of cleaner production*, vol. 140, pp. 1454–1464, 2017.

- [2] S. Solaimani, W. Keijzer-Broers, and H. Bouwman, "What we do—and don't—know about the smart home: an analysis of the smart home literature," *Indoor and Built Environment*, vol. 24, no. 3, pp. 370–383, 2015.
- [3] X. Yu, Z. Zhou, L. Zhang, and X.-Y. Li, "Thumbup: Secure smartwatch controller for smart homes using simple hand gestures," *IEEE TMC*, vol. 23, no. 1, pp. 865–878, 2022.
- [4] J. Cui, J. Xiao, H. Zhong, J. Zhang, L. Wei, I. Bolodurina, and D. He, "Lh-ids: Lightweight hybrid intrusion detection system based on differential privacy in vanets," *IEEE TMC*, vol. 23, no. 12, pp. 12 195–12 210, 2024.
- [5] E. King, H. Yu, S. Lee, and C. Julien, "Sasha: creative goal-oriented reasoning in smart homes with large language models," *ACM IMWUT*, vol. 8, no. 1, pp. 1–38, 2024.
- [6] Z. Ramadan, M. F. Farah, and L. El Essrawi, "From amazon. com to amazon. love: How alexa is redefining companionship and interdependence for people with special needs," *Psychology & Marketing*, vol. 38, no. 4, pp. 596–609, 2021.
- [7] J. R. Bellegarda, "Spoken language understanding for natural interaction: The siri experience," *Natural Interaction with Robots, Knowbots and Smartphones: Putting Spoken Dialog Systems into Practice*, pp. 3–14, 2013.
- [8] D. Xu, W. Yin, H. Zhang, X. Jin, Y. Zhang, S. Wei, M. Xu, and X. Liu, "Edgellm: Fast on-device llm inference with speculative decoding," *IEEE Transactions on Mobile Computing*, vol. 24, no. 4, pp. 3256–3273, 2025.
- [9] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, "A survey of large language models," *arXiv preprint arXiv:2303.18223*, 2023.
- [10] C. Liu and J. Zhao, "Enhancing stability and resource efficiency in llm training for edge-assisted mobile systems," *IEEE Transactions on Mobile Computing*, pp. 1–18, 2025.
- [11] Openai, "Gpt-4." [Online]. Available: <https://chat.openai.com/chat/>
- [12] Q. Chai and G. Gong, "Verifiable symmetric searchable encryption for semi-honest-but-curious cloud servers," in *IEEE ICC*, 2012, pp. 917–922.
- [13] P. Sun, Z. Wang, L. Wu, Y. Feng, X. Pang, H. Qi, and Z. Wang, "Towards personalized privacy-preserving incentive for truth discovery in mobile crowdsensing systems," *IEEE TMC*, vol. 21, no. 1, pp. 352–365, 2020.
- [14] T. Li, H. Wang, Q. Li, Y. Jiang, and Z. Yuan, "Cl-shield: A continuous learning system for protecting user privacy," *IEEE Transactions on Mobile Computing*, vol. 24, no. 4, pp. 3148–3162, 2025.
- [15] W. Wang, Y. Wang, P. Duan, T. Liu, X. Tong, and Z. Cai, "A triple real-time trajectory privacy protection mechanism based on edge computing and blockchain in mobile crowdsourcing," *IEEE TMC*, vol. 22, no. 10, pp. 5625–5642, 2022.
- [16] H. Chi, Q. Zeng, X. Du, and J. Yu, "Cross-app interference threats in smart homes: Categorization, detection and handling," in *2020 50th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*. IEEE, 2020, pp. 411–423.
- [17] L. Shen, Q. Yang, K. Cui, Y. Zheng, X.-Y. Wei, J. Liu, and J. Han, "Fedconv: A learning-on-model paradigm for heterogeneous federated clients," in *ACM MobiSys*, 2024, pp. 398–411.
- [18] L. Shen and Y. Zheng, "Feddm: Data and model heterogeneity-aware federated learning via dynamic weight sharing," in *2023 IEEE ICDCS*. IEEE, 2023, pp. 975–976.
- [19] L. Shen, Q. Yang, K. Cui, Y. Zheng, X.-Y. Wei, J. Liu, and J. Han, "Hierarchical and heterogeneous federated learning via a learning-on-model paradigm," *IEEE TMC*, pp. 1–16, 2025.
- [20] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei, "Scaling laws for neural language models," *arXiv preprint arXiv:2001.08361*, 2020.
- [21] D. Rivkin, F. Hogan, A. Feriani, A. Konar, A. Sigal, S. Liu, and G. Dudek, "Sage: Smart home agent with grounded execution," *arXiv preprint arXiv:2311.00772*, 2023.
- [22] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel *et al.*, "Retrieval-augmented generation for knowledge-intensive nlp tasks," *Advances in neural information processing systems*, vol. 33, pp. 9459–9474, 2020.
- [23] E. King, H. Yu, S. Lee, and C. Julien, "get ready for a party": Exploring smarter smart spaces with help from large language models," *arXiv preprint arXiv:2303.14143*, 2023.
- [24] E. King, H. Yu, S. Vartak, J. Jacob, S. Lee, and C. Julien, "Thoughtful things: Building human-centric smart devices with small language models," *arXiv preprint arXiv:2405.03821*, 2024.
- [25] S. Aheleroff, X. Xu, Y. Lu, M. Aristizabal, J. P. Velásquez, B. Joa, and Y. Valencia, "Tot-enabled smart appliances under industry 4.0: A case study," *Advanced engineering informatics*, vol. 43, p. 101043, 2020.
- [26] R. H. Jensen, Y. Strengers, J. Kjeldskov, L. Nicholls, and M. B. Skov, "Designing the desirable smart home: A study of household experiences and energy consumption impacts," in *ACM CHI*, 2018, pp. 1–14.
- [27] M. Li, W. Gu, W. Chen, Y. He, Y. Wu, and Y. Zhang, "Smart home: architecture, technologies and systems," *Procedia computer science*, vol. 131, pp. 393–400, 2018.
- [28] D. Marikyan, S. Papagiannidis, and E. Alamanos, "A systematic review of the smart home literature: A user perspective," *Technological Forecasting and Social Change*, vol. 138, pp. 139–154, 2019.
- [29] R. D. Manu, S. Kumar, S. Snehashish, and K. Rekha, "Smart home automation using iot and deep learning," *International Research Journal of Engineering and Technology*, vol. 6, no. 4, pp. 1–4, 2019.
- [30] A. S. Tulshan and S. N. Dhage, "Survey on virtual assistant: Google assistant, siri, cortana, alexa," in *Advances in Signal Processing and Intelligent Recognition Systems: 4th International Symposium SIRS 2018, Bangalore, India, September 19–22, 2018, Revised Selected Papers 4*. Springer, 2019, pp. 190–201.
- [31] J. Li, Z. Li, G. Tyson, and G. Xie, "Characterising usage patterns and privacy risks of a home security camera service," *IEEE TMC*, vol. 21, no. 7, pp. 2344–2357, 2020.
- [32] K. Han, A. Xiao, E. Wu, J. Guo, C. Xu, and Y. Wang, "Transformer in transformer," *NIPS*, vol. 34, pp. 15908–15919, 2021.
- [33] L. Shen, Q. Yang, Y. Zheng, and M. Li, "Autoiot: Llm-driven automated natural language programming for aiots applications," in *ACM MobiCom*, 2025.
- [34] L. Shen, Q. Yang, X. Huang, Z. Ma, and Y. Zheng, "Gpiot: Tailoring small language models for iot program synthesis and development," in *ACM SenSys*, 2025.
- [35] L. Shen and Y. Zheng, "Iotcoder: A copilot for iot application development," in *ACM MobiCom*, 2024, pp. 1647–1649.
- [36] J. Rey-Jouanchicot, A. Bottaro, E. Campo, J.-L. Bouraoui, N. Vigouroux, and F. Vella, "Leveraging large language models for enhanced personalised user experience in smart homes," 2024.
- [37] A. K. Newendorp, M. Sanaei, A. J. Perron, H. Sabouni, N. Javadpour, M. Sells, K. Nelson, M. Dorneich, and S. B. Gilbert, "Apple's knowledge navigator: Why doesn't that conversational agent exist yet?" in *ACM CHI*, 2024, pp. 1–14.
- [38] P. M. Schwartz and D. J. Solove, "The pii problem: Privacy and a new concept of personally identifiable information," *NYUL rev.*, vol. 86, p. 1814, 2011.
- [39] Z. Wang, Y. Sun, D. Liu, J. Hu, X. Pang, Y. Hu, and K. Ren, "Location privacy-aware task offloading in mobile edge computing," *IEEE TMC*, vol. 23, no. 3, pp. 2269–2283, 2023.
- [40] V. Srinivasan, J. Stankovic, and K. Whitehouse, "Protecting your daily in-home activity information from a wireless snooping attack," in *ACM UbiComp*.
- [41] B. Gunel, J. Du, A. Conneau, and V. Stoyanov, "Supervised contrastive learning for pre-trained language model fine-tuning," *arXiv preprint arXiv:2011.01403*, 2020.
- [42] H. Yu, J. Hua, and C. Julien, "Analysis of ifttt recipes to study how humans use internet-of-things (iot) devices," in *ACM SenSys*, 2021, pp. 537–541.
- [43] C.-Y. Lin and F. Och, "Looking for a few good metrics: Rouge and its evaluation," in *Ntcsr workshop*, 2004.
- [44] D. Blazakis, "The apple sandbox," *Arlington, VA, January*, 2011.
- [45] C. Xu, Q. Sun, K. Zheng, X. Geng, P. Zhao, J. Feng, C. Tao, and D. Jiang, "Wizardlm: Empowering large language models to follow complex instructions," *arXiv preprint arXiv:2304.12244*, 2023.
- [46] C.-Y. Lin, "Rouge: A package for automatic evaluation of summaries," in *Text summarization branches out*, 2004, pp. 74–81.
- [47] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, "Qlora: Efficient finetuning of quantized llms," *NIPS*, vol. 36, 2024.
- [48] Z. Han, C. Gao, J. Liu, S. Q. Zhang *et al.*, "Parameter-efficient finetuning for large models: A comprehensive survey," *arXiv preprint arXiv:2403.14608*, 2024.
- [49] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *NIPS*, vol. 35, pp. 24 824–24 837, 2022.
- [50] Y. Zhao, Z. Yang, X. He, X. Cai, X. Miao, and Q. Ma, "Trine: Cloud-edge-device cooperated real-time video analysis for household applications," *IEEE TMC*, vol. 22, no. 8, pp. 4973–4985, 2022.
- [51] T. W. Li, A. Arya, and H. Jin, "Redesigning privacy with user feedback: The case of zoom attendee attention tracking," in *ACM CHIs*, 2024, pp. 1–14.

- [52] E. Al Qahtani, P. Story, and M. Shehab, "The impact of risk appeal approaches on users' sharing confidential information," in *ACM CHI*, 2024, pp. 1–21.
- [53] U. Scherhag, C. Rathgeb, J. Merkle, and C. Busch, "Deep face representations for differential morphing attack detection," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 3625–3639, 2020.
- [54] Q. Li, J. Wen, and H. Jin, "Governing open vocabulary data leaks using an edge llm through programming by example," *ACM IMWUT*, vol. 8, no. 4, pp. 1–31, 2024.
- [55] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems*, vol. 32, 2019.
- [56] H. Chase, "Langchain," 10 2022. [Online]. Available: <https://github.com/langchain-ai/langchain>
- [57] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.
- [58] M. Douze, A. Guzhva, C. Deng, J. Johnson, G. Szilvasy, P.-E. Mazaré, M. Lomeli, L. Hosseini, and H. Jégou, "The faiss library," *arXiv preprint arXiv:2401.08281*, 2024.
- [59] G. Team, T. Mesnard, C. Hardin, R. Dadashi, S. Bhupatiraju, S. Pathak, L. Sifre, M. Rivière, M. S. Kale, J. Love *et al.*, "Gemma: Open models based on gemini research and technology," *arXiv preprint arXiv:2403.08295*, 2024.
- [60] M. Abdin, S. A. Jacobs, A. A. Awan, J. Aneja, A. Awadallah, H. Awadalla, N. Bach, A. Bahree, A. Bakhtiari, H. Behl *et al.*, "Phi-3 technical report: A highly capable language model locally on your phone," *arXiv preprint arXiv:2404.14219*, 2024.
- [61] P. Sun, "Security and privacy protection in cloud computing: Discussions and challenges," *Journal of Network and Computer Applications*, vol. 160, p. 102642, 2020.
- [62] Z. Yang, M. Yang, Y. Zhang, G. Gu, P. Ning, and X. S. Wang, "Appintnet: Analyzing sensitive data transmission in android for privacy leakage detection," in *ACM CCS*, 2013, pp. 1043–1054.
- [63] C. Marcolla, V. Sucasas, M. Manzano, R. Bassoli, F. H. Fitzek, and N. Aaraj, "Survey on fully homomorphic encryption, theory, and applications," *Proceedings of the IEEE*, vol. 110, no. 10, pp. 1572–1609, 2022.
- [64] S. Kim, S. Yun, H. Lee, M. Gubri, S. Yoon, and S. J. Oh, "Propile: Probing privacy leakage in large language models," *NIPS*, vol. 36, 2024.
- [65] A. Frikha, N. Walha, K. K. Nakka, R. Mendes, X. Jiang, and X. Zhou, "Incognitext: Privacy-enhancing conditional text anonymization via llm-based private attribute randomization," *arXiv preprint arXiv:2407.02956*, 2024.
- [66] Z. Xu, S. Jain, and M. Kankanhalli, "Hallucination is inevitable: An innate limitation of large language models," *arXiv preprint arXiv:2401.11817*, 2024.
- [67] Q. Li, J. Hong, C. Xie, J. Tan, R. Xin, J. Hou, X. Yin, Z. Wang, D. Hendrycks, Z. Wang *et al.*, "Llm-pbe: Assessing data privacy in large language models," *arXiv preprint arXiv:2408.12787*, 2024.
- [68] o. Fabian Pedregosa, "Memory-Profiler: Monitor Memory usage of Python code." [Online]. Available: https://github.com/pythonprofilers/memory_profiler
- [69] F. Bang, "Gptcache: An open-source semantic cache for llm applications enabling faster answers and cost savings," in *NLP-OSS*, 2023, pp. 212–218.
- [70] A. Yang, B. Yang, B. Hui, B. Zheng, B. Yu, C. Zhou, C. Li, C. Li, D. Liu, F. Huang *et al.*, "Qwen2 technical report," *arXiv preprint arXiv:2407.10671*, 2024.
- [71] M. Abdin, S. A. Jacobs, A. A. Awan, J. Aneja, A. Awadallah, H. Awadalla, N. Bach, A. Bahree, A. Bakhtiari, H. Behl *et al.*, "Phi-3 technical report: A highly capable language model locally on your phone," *arXiv preprint arXiv:2404.14219*, 2024.
- [72] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan *et al.*, "The llama 3 herd of models," *arXiv preprint arXiv:2407.21783*, 2024.
- [73] A. Joshi, S. Kale, S. Chandel, and D. K. Pal, "Likert scale: Explored and explained," *British journal of applied science & technology*, vol. 7, no. 4, pp. 396–403, 2015.