

# Are You Killing Time? Predicting Smartphone Users' Time-killing Moments via Fusion of Smartphone Sensor Data and Screenshots

Yu-Chun, Chen\*

National Yang Ming Chiao Tung  
University  
Hsinchu, Taiwan  
lesley.yc.chen@gmail.com

Yu-Jen, Lee\*

National Yang Ming Chiao Tung  
University  
Hsinchu, Taiwan  
daniel08099@gmail.com

Kuei-Chun, Kao

National Yang Ming Chiao Tung  
University  
Hsinchu, Taiwan  
johnson0213.cs07@nycu.edu.tw

Jie, Tsai

National Yang Ming Chiao Tung  
University  
Hsinchu, Taiwan  
jessic462071@gmail.com

En-Chi, Liang

National Yang Ming Chiao Tung  
University  
Hsinchu, Taiwan  
n7enchi.ee08@nycu.edu.tw

Wei-Chen, Chiu

National Yang Ming Chiao Tung  
University  
Hsinchu, Taiwan  
walon@cs.nctu.edu.tw

Faye, Shih

Bryn Mawr College  
Bryn Mawr, United States  
shihfaye@gmail.com

Yung-Ju, Chang<sup>†</sup>

National Yang Ming Chiao Tung  
University  
Hsinchu, Taiwan  
armuro@nycu.edu.tw

## ABSTRACT

Time-killing on smartphones has become a pervasive activity, and could be opportune for delivering content to their users. This research is believed to be the first attempt at time-killing detection, which leverages the fusion of phone-sensor and screenshot data. We collected nearly one million user-annotated screenshots from 36 Android users. Using this dataset, we built a deep-learning fusion model, which achieved a precision of 0.83 and an AUROC of 0.72. We further employed a two-stage clustering approach to separate users into four groups according to the patterns of their phone-usage behaviors, and then built a fusion model for each group. The performance of the four models, though diverse, yielded better average precision of 0.87 and AUROC of 0.76, and was superior to that of the general/unified model shared among all users. We investigated and discussed the features of the four time-killing behavior clusters that explain why the models' performance differ.

## CCS CONCEPTS

• **Human-centered computing** → **Ubiquitous and mobile computing systems and tools**; *Smartphones*.

\*Both authors contributed equally to this research.

<sup>†</sup>Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

CHI '23, April 23–28, 2023, Hamburg, Germany

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-1-4503-9421-5/23/04...\$15.00  
<https://doi.org/10.1145/3544548.3580689>

## KEYWORDS

Time-killing; Screenshot; Deep Learning; Opportune Moment; Mobile Devices

### ACM Reference Format:

Yu-Chun, Chen, Yu-Jen, Lee, Kuei-Chun, Kao, Jie, Tsai, En-Chi, Liang, Wei-Chen, Chiu, Faye, Shih, and Yung-Ju, Chang. 2023. Are You Killing Time? Predicting Smartphone Users' Time-killing Moments via Fusion of Smartphone Sensor Data and Screenshots. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*, April 23–28, 2023, Hamburg, Germany. ACM, New York, NY, USA, 19 pages. <https://doi.org/10.1145/3544548.3580689>

## 1 INTRODUCTION

Researchers have leveraged smartphones' capabilities to engage individuals in a variety of tasks, including mobile learning exercises [12], just-in-time interventions [19], mobile self-reports [66], and crowdsourcing tasks [18]. In recent years, commercial platforms have also started doing so to obtain crowdsourced data, such as locale information<sup>1</sup> [3, 94] and labeled data<sup>2</sup> [17, 18]. However, given human beings' limited attentional resources, a crucial problem for anyone delivering content to phones is how to make it stand out from the feast of other incoming information. One mainstream approach to achieving this is to predict moments at which users are receptive to such content, e.g., the content related to notifications [63, 70, 73], questionnaires [70], and reading material [22, 70] explored in prior studies.

Moments of "attention surplus" [72] constitute another opportunity for such detection attempts. Pielot et al. [72], for example, attempted to detect one kind of "attention surplus" state – boredom. In such situations, some people tend to seek stimulation on

<sup>1</sup><https://maps.google.com/localguides>

<sup>2</sup><https://play.google.com/store/apps/details?id=com.google.android.apps.village.boond>

their phones to alleviate boredom. Beyond boredom, however, research has shown that mobile phone usage is not always driven by a specific purpose [32], but is often accompanied by, or is primarily, "time-killing" behavior: i.e., filling periods of perceived free time [56, 65], such as waiting for a train or during an uninteresting speech [41]. Time-killing is not necessarily linked to boredom, as the user may not be feeling unfulfilled [30], unengaged [25], lacking meaning [91], or having low arousal [20] – which are often used to define boredom [26]. Rather, this type of usage can simply be a result of habit [56, 65]. Habitual phone usage is widespread and can occur even while a person is performing work-related tasks. However, when the trigger of habitual phone usage is to fill free time, even if it may occasionally include productive tasks such as checking emails and messages [13, 21], it is typically viewed as aimless [32], and if it occurs frequently without regulation, as addictive [51] and compulsive [87]. In an effort to improve users' overall productivity [102] or well-being [36, 53, 93], previous studies have aimed to decrease this type of phone usage or make it more productive [12, 23, 37]. Meanwhile, phone users may not necessarily feel bored during the phone usage, since they may feel interested, entertained, and amused, once they have started using the phone and been exposed to entertaining and interesting content, such as playing games, watching funny videos.

However, despite these conceptual and operational differences between the feeling of boredom and phone usage for time-killing, other than boredom detection [72] which is the closest to this purpose, there has been limited effort in detecting the occurrence of such prevalent phone usage intended for time-killing, making the feasibility of detecting this phone usage unclear. Nevertheless, detecting this phone usage provides not only opportunities for researchers and practitioners to deliver contents that help users fill time by engaging in more productive activity during these periods [70], but also enables researchers to identify when and how often in a day such phone usage periods take place.

In light of these benefits, we aim to develop a model for effectively detecting phone usage intended for time-killing. To achieve this aim, we developed an Android research application that automatically collected smartphone screenshot and phone-sensor data, and an interface that allowed its users to efficiently annotate their screenshot data with time-killing periods and their availability during those periods for viewing notifications. Screenshot data were collected because we expected that it would reveal rich temporal, textual, graphical, and topical information about people's phone usage [9] for detection.

Data collection with 36 participants over 14 days yielded a dataset of 967,466 pairings of annotated phone-sensor data with screenshots, covering 1,343.7 hours of phone usage. Moreover, the participants self-reported being more available for notifications during time-killing phone use (82.2%) than non-time-killing phone use (63.3%). Using this dataset, we built a deep-learning-based fusion model that achieved a precision of 0.83 and an Area Under the Receiver Operating Characteristics (AUROC) of 0.72. To further improve the model's performance by taking account of differences in the participants' time-killing behaviors, we employed two-stage clustering that grouped people with similar phone usage behaviors into four groups, and built a fusion model for each group. The four resulting models' collective average precision and AUROC

went up to 0.87 and 0.76, respectively: i.e., better than those of the general model (i.e., the one shared among all users). However, the four models achieved quite different performance on many metrics, and to obtain insights into these differences, we delved into the characteristics of each user group's phone-usage behavior as well as the important features learned by their respective models that were positively and negatively correlated with time-killing moments. The results of that investigation help explain both how and why the effectiveness of sensor data and phone screenshots for detecting time-killing moments varied across user clusters. To facilitate future research in this area, we also release open source code and model at GitHub<sup>3</sup>.

This paper makes the following three major contributions to HCI.

1. It presents the development of a deep-learning-based fusion model that detects smartphone users' time-killing phone usage, enabling researchers to deliver productive content or intervention accordingly during such phone usage.
2. It demonstrates that building such models for user groups clustered according to their phone-usage behaviors can achieve better overall model performance, and that all group-specific models may achieve significantly better performance than the general model.
3. It shows how and why the effectiveness of sensor data and phone screenshots for detecting time-killing moments vary across different time-killing behavioral patterns.

## 2 RELATED WORK

### 2.1 Interruptibility, Breakpoint, and Opportune Moment Prediction

Many studies have employed machine-learning techniques to predict interruptible moments, breakpoints, and opportune moments. For instance, Pejovic et al. [68] achieved the predictions of mobile interruptibility with a precision of 0.72. Others have focused on predicting opportune moments for receiving calls and notifications. For example, Fisher et al. [28] built personalized models to predict such moments in the case of incoming phone calls, and achieved an average accuracy above 0.96 (see also Smith et al. [80]); and Pielot et al. [71] applied machine-learning techniques to predict whether users would view an incoming message notification within the next few minutes or not. Mehrotra et al. [59], for instance, proposed a system based on machine-learning algorithms that automatically extracted rules for phone users' preferences about receiving notifications. A similar study by Visuri et al. [92] reported that 81.7% of phone-user interactions with alert dialogs could be accurately predicted based on user clusters.

Among researchers seeking to identify opportunities based on breakpoints, Ho et al. [34] detected postural and ambulatory activity transitions in real-time. Iqbal and Bailey [39] showed that scheduling notifications at breakpoints reduced frustration and reaction times. Okoshi et al. [63], who also developed a breakpoint-detection system for mobile devices, showed that notifications delivered during breakpoints required 33% less cognitive load than those delivered randomly. Later, the same authors [64] showed

<sup>3</sup>[https://github.com/johnsonkao0213/kill\\_time\\_detection](https://github.com/johnsonkao0213/kill_time_detection)

that delaying notification delivery until an interruptible moment significantly reduced in user response time. Adamczyk et al. [1] divided breakpoints in tasks into two types, coarse and fine, and showed that delivering notifications at their predicted best points for interruptions consistently produced less annoyance, frustration, and time pressure. Iqbal et al. [38] applied it to statistical models that mapped interaction features to each breakpoint type, based on task-execution data and video footage. And Park et al. [67] used built-in sensors to detect social contexts, enabling them to identify four distinct types of breakpoints suitable for delivering deferred smartphone notifications.

Detecting moments when device users want to engage with content has also been a focus of considerable research effort. Sarker et al. [78], for example, sought to identify moments for delivering notifications that would result in maximum engagement. Similarly, Choi et al. [19] built a mobile intervention system for preventing prolonged sedentary behaviors, and showed that contextual factors and cognitive/physical states were good predictors of decision points. Turner et al. [88] decomposed notification interaction into three stages – reachability, engageability, and receptivity – and developed models for predicting when phone users reached each of them. Pielot et al. [70] built a model that predicted whether their participants would engage with different types of content they were offered, which achieved a success rate 66.6% higher than the baseline. Steil et al. [81] predicted whether people’s primary attentional focus was on their handheld mobile devices, and proposed “attention forecasting”, similar in spirit to user-intention prediction.

Another strand of research on attention prediction involves identifying “attention surplus” moments and timing the delivery of specific content and tasks accordingly. Such content and tasks have thus far included reading material [22, 70], learning material [12, 23, 37], interventions [19, 60, 78], questionnaires [33, 70], and crowdsourcing tasks [18], among others. For example, Pielot et al. [72] deemed moments of boredom to be moments of attention surplus, and detected them using phone logs: an approach that achieved 0.83 AUROC. However, they obtained a high number of false positives, which they felt would lead to user annoyance, and therefore tuned their model to strike an optimal balance between recall and precision. Based on boredom levels detected via phone-sensor data, Dingler et al. [23] delivered micro-learning reminders to language learners, and their results suggested the feasibility of identifying moments of boredom as mobile learning opportunities. Cai et al. [12] developed WaitSuite, which detects various types of moments when its users are waiting for something to happen, and delivers micro-learning tasks during them. Similarly, Inie and Lungu [37] detected that when users were about to become unproductive due to visiting time-wasting websites, blocked such visits, and delivered learning exercises instead. Chen et al. [16] used screenshot data instead of phone sensor data to study the detection of time-filling phone use. However, their work did not leverage both sensor and screenshot data. In contrast, this work uses both types of data and builds models for different user groups.

## 2.2 Phone-usage Research

Several studies have utilized self-report methods such as interviews and diaries. For instance, Palen et al. [66] investigated mobile usage

via a voicemail diary study. However, because self-report methods are subject to recall biases [24, 29], quantitative analysis of phone-usage logs is becoming increasingly popular [27, 97, 99]. For example, Böhmer et al.’s [8] large-scale study based on logged application usage found that news applications were most popular in the morning; and that game-playing mostly occurred at night. Xu et al. [97] also found differential patterns by app type, e.g., that sports apps were more frequently used in the evening. Falaki et al. [27] distinguished between two broad types of intentional use activities – user/phone interaction, and app use – and found that strong diversity in users’ behavior was linked to different purposes for using phones. Canneyt et al. [90] revealed how app-usage behavior was disrupted during major political, social, and sporting events. And Li et al. [54] studied the long-term evolution of mobile-app usage, and found that the diversity of app-category usage declined over time, whereas the diversity of the individual apps used increased.

Lukoff et al. [56] identified situations in which people felt a lack of meaning while using their phones, which prominently included passively browsing social media, consuming entertainment, and habitual use. Hiniker et al. [32] likewise reported “ritualistic” uses of phones, which tended to be habitual. Another habitual phone usage is “phubbing”, i.e., the habit of snubbing someone in favour of a mobile phone. As Al-Saggaf et al. [5] have suggested, individuals engage in phubbing while they are experiencing negative emotions such as boredom, loneliness, and fear of missing out. In a different study, Al-Saggaf and colleagues [4] reported that trait boredom could predict phubbing frequency.

A growing body of work involves attempts to construct models of phone usage. Kostakos et al. [48], for instance, developed a Markov state transition model of smartphone screen use. Jesdabodi et al. [42] identified phone users’ behavioral states, and showed that morning and evening routines were both mostly marked by communication and gaming activities. Some other work has focused on understanding differences in usage features across distinct user clusters. Zhao et al. [101] studied app usage with a two-step clustering approach and revealed clusters of users including “night communicators”, “evening learners”, and “screen checkers”, among others. Jones et al. [43], on the other hand, identified three clusters of users: “checkers”, “waiters”, and “responsives”. And Katevas et al. [44], based on a combination of phone-use log data and experience-sampling method data, identified five types of mobile-phone use: “limited use”, “business use”, “power use”, “personality-induced problematic use”, and “externally induced problematic use”.

Finally, because log data are limited to system events like screen events and app states, some researchers have used screenshots and video recordings to study phone usage. For example, Brown et al. [10] combined screen-captures of iPhone use with recordings from wearable video cameras, and showed that video data illuminated various aspects of people’s interactions with their phones. Subsequently, Brown et al. [11] collected screen recordings of phone use and audio recordings of ambient talk, and identified various situations in which people engaged in phone usage with their “free” attention, such as engaged in quick games or social-media checking while waiting for a friend to arrive or for an event to start. Reeves et al. [76] showed how screenshots could be used to unobtrusively collect valuable data on individuals’ digital life experience. Later, Reeves et al. [9] explored how textual and graphical features

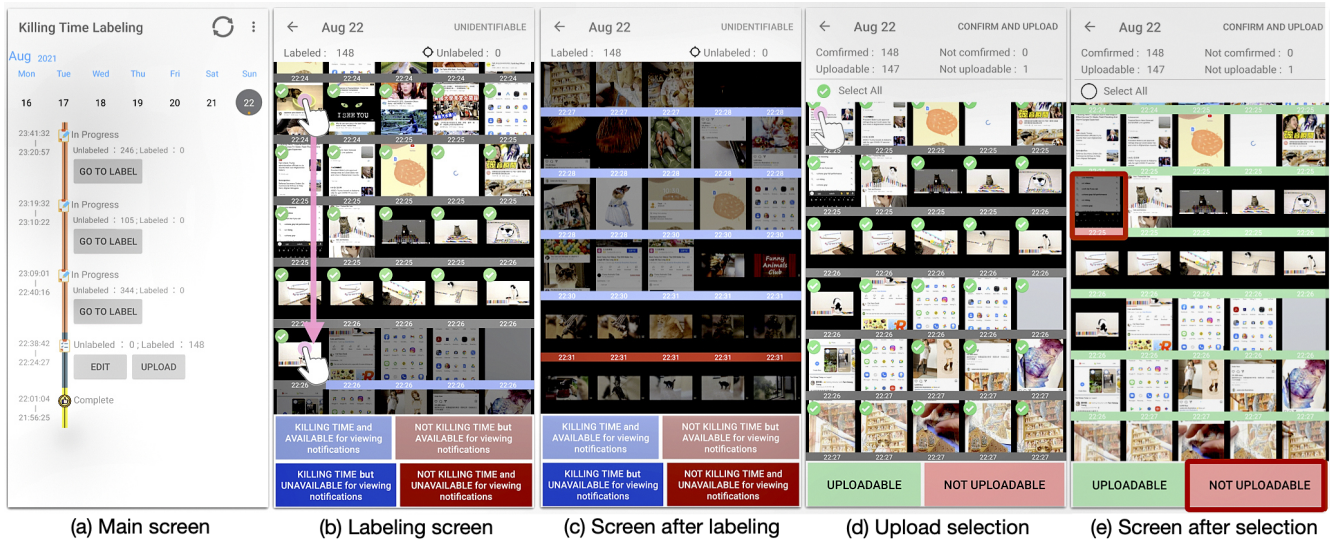


Figure 1: User interfaces for the main functions of the Killing Time Labeling application

changed during sessions and found that word and image velocity both decreased late at night.

Some other researchers have used deep-learning models trained on large amounts of Graphical User Interface (GUI) data to detect screenshots. For instance, Beltramelli’s [7] Pix2Code applies an end-to-end neural image captioning model to generate code from a single input image, with better than 0.77 accuracy across various platforms. Similarly, Chen et al. [15] utilized a CNN-RNN model to generate GUI skeletons from screenshots. Other work focused on locating UI elements on screens, such as by White et al. [96], has used YOLOv2 [75] to automatically identify GUI widgets in screenshots. Chen et al. [14] built a gallery of large scale of GUI designs by applying a Faster RCNN model [77]; and Zhang et al. [100] proposed an on-device model capable of detecting UI elements.

### 3 DATA COLLECTION

#### 3.1 Input-data Selection

Screenshot collection has become a popular method in HCI research, because it allows researchers to collect quantitative and qualitative data simultaneously [45, 49, 50, 83] in high granularity and rich detail [9]. Along with information about people’s interactions with their phones, it can help researchers reconstruct both moment-to-moment phone use and wider usage patterns [58, 74, 76, 98]. Due to these advantages, we aimed to leverage screenshot data, along with phone-sensor interaction information (including user/phone interaction and phone status), to extract features that characterized our participants’ app usage and switching patterns. We then attempted to associate such usage information and patterns with time-killing vs. non-time-killing moments.

#### 3.2 Research Instrument

We developed an Android research application, called Killing Time Labeling (KTL), to collect annotated screenshots and phone-sensor

data (i.e., Android accessibility events, screen status, network connections, phone volume, application usage, and type of transportation). KTL also captures the notifications its users receive, the times at which they receive them, and how they are dealt with. The background service that automatically collects data is activated within a 12-hour timeframe every day, the default being from 10:00 a.m. to 10:00 p.m., but the start time and end time are both user-adjustable, meaning that the data might be collected for more than 12 hours per day in some cases. During whatever 12+-hour window the user has chosen, his/her phone-sensor data is collected every five seconds. Screenshots are also captured every five seconds, but only when the phone screen is on.

We designed a user interface for KTL that allowed our participants to easily select groups of screenshots via drag-and-drop for data labeling (see Fig. 1). A detailed demonstration of this data-labeling procedure is provided in our supplemental video. The participants were instructed to review and annotate screenshots in accordance with the situations in which they were taken. Specifically, they were instructed to annotate screenshots with whether their specific phone use periods were intended to kill time, as well as whether they were available for notifications during those periods. Notably, whereas the former information was used as the ground truth for developing the model for time-killing detection, the latter was to examine whether the participants would be more available for notifications during time-killing phone use than during non-time-killing phone use. Therefore, for each screenshot, participants had five annotation options: 1) killing time and available for viewing notifications; 2) not killing time but available for viewing notifications; 3) killing time but unavailable for viewing notifications; 4) not killing time and unavailable for viewing notifications; and 5) unidentifiable, i.e., the participant could not be certain of his/her time-killing state or had forgotten it. Each time s/he manually selected and annotated a series of screenshots, the

participant was to report his/her actual activities<sup>4</sup> at the time those screenshots were taken. We instructed the participants to annotate them as “killing time” as long as they felt or subjectively deemed that their mobile-phone usage at the time was to pass the time, and otherwise to annotate it as “not killing time.” Regarding the availability label for viewing notifications, we instructed them to annotate screenshots as “unavailable for viewing notifications” if they positively did not want to be interrupted or to see any notifications when using the app, and otherwise to annotate them as “available.” To reduce labeling bias due to recall errors, we instructed participants to mark screenshots as “unidentifiable” when feeling unsure, to annotate screenshots whenever possible at their convenience, and to complete annotation before going to bed every day. In addition, KTL also sent a reminder at night and invalidated screenshots that had not received any annotations after two days, which participants could no longer annotate.

All screenshots were reduced in size and temporarily stored in the local storage of the participants’ respective phones before they were reviewed, labeled, and manually uploaded to our server. The participants had the right not to upload any given screenshot, e.g., because it contained sensitive information. Phone-sensor data, on the other hand, was automatically uploaded by KTL whenever a participant’s phone was connected to the Internet, to avoid such data taking up too much storage space. Also, to avoid impacting the participants’ data plans, KTL only did so via WiFi networks, unless a user overrode this feature and chose to upload using the cellular network. The participants were informed of all these rules in a pre-study meeting (the other purposes are detailed in section 3.3 below).

KTL also delivered notifications linked to experience sampling method (ESM) questionnaires and various other content types. That other content consisted of 1) crowdsourcing tasks<sup>5</sup> [17, 18], 2) non-ESM questionnaires<sup>6</sup> [70], 3) advertisements [70], and 4) news items [69, 70, 72]. KTL only sent such notifications within the user’s chosen 12+-hour timeframe and only when his/her screen was on. Each notification was randomly selected from among the four types listed above, and delivered at random intervals of not less than one or more than three hours. Five minutes after each notification arrived, an ESM questionnaire was also sent, asking the participant to report his/her awareness of and receptivity to that notification, as well as what context s/he was in when it arrived.

### 3.3 Study Procedure

Prior to data collection, due to the COVID-19 pandemic, we allowed our participants to choose between remotely and physically attending a pre-study meeting, during which the researchers helped them install KTL on their phones, explained how to use it, and walked them through the study procedure. We told them that we expected them to annotate all screenshots automatically captured by KTL every day, and that 14 days of active participation were needed for their data to be useful to us. Therefore, for every day they failed

to submit annotated screenshots, their participation deadline was postponed by one additional day. On their respective final days of participation, to aid future analysis, they completed four additional questionnaires that measured their boredom proneness [82], smartphone addiction [55], inattention [46], and perceived acceptability of time-killing detection being deployed on their phone. In addition, we invited all participants to two optional semi-structured interviews, the first of which was held after they had contributed data for seven full days, and the second, after their participation was complete. In those interviews, we asked them about their labeling processes, time-killing behaviors and preferences, and how they killed time (both typically and during the study). Those who completed 14 days of data collection were paid NT\$1,350 (approximately US\$44). Those who participated in the mid-study interview were paid an additional NT\$150 (US\$5), and those who were interviewed after the study, another NT\$250 (US\$8). The study was approved by our university’s Institutional Review Board (IRB).

### 3.4 Recruitment and Participants

We selected participants with various occupations, expecting they would have different time-killing patterns. Also, to ensure that sufficient data were collected, we selected participants who used their mobile phones more than one hour a day, according to their self-reporting in a screening questionnaire. We recruited participants primarily via several Facebook groups aimed at matching researchers with study participants in our country, but also posted a recruiting message on Facebook pages for the local community in the hope of further diversifying our subjects’ backgrounds. A total of 55 participants were recruited for this study, including 12 participants who participated in a pilot study that helped us improve the KTL design, annotation mechanism, and study procedure. Of the remaining 43 participants, one withdrew before data collection commenced, two did not complete the experiment, and four others were excluded as being outliers (i.e., they had annotated more than 95% of their data as “killing time”). As a result, data from 36 people were used for training our time-killing detection model. Of those 36, 32 took part in both optional interviews, two only in the mid-study interview, and two others, only in the post-study interview. All 36 participants were aged between 20 and 54 ( $M = 27.4$ ,  $SD = 6.8$ ), with 16 identifying as male and 20 as female. Half were students, and the other half in employment.

### 3.5 Data Collection

Most participants provided data on 12 hours of phone usage per day, but six voluntarily extended this to 13-15 hours; one, to 17.5 hours; and another, to the whole day. In total, 1,186,345 screenshots were annotated (per-participant  $M = 32,954.0$ ,  $SD = 15,557.9$ ), which represented approximately 1,633.8 hours of phone use. Among these 1,186,345 annotated data points, 1,062,780 (89.6%) screenshots were uploaded; a per-participant average of 29,521.7 screenshots ( $SD = 13,544.9$ ). Thus, the initial dataset that we collected for analysis consisted of 1,062,780 annotated screenshots and the phone-sensor data associated with the moments at which they were captured. Two-thirds ( $n = 773,401$ ) of uploaded and non-uploaded screenshots were annotated as “killing time”, and somewhat over a quarter ( $n = 346,792$ ) as “not killing time”, with the remaining 5.6% ( $n = 66,152$ )

<sup>4</sup>This question was adopted from previous research [52].

<sup>5</sup>The crowdsourcing questions were inspired by Google Crowdsourcing and Local Guide, two platforms that aim to improve Google Maps and various other Google services through user-oriented training of multiple algorithms.

<sup>6</sup>The questionnaire was inspired by Google Opinion Rewards, which offers rewards to its users who answer surveys and opinion polls on various topics.

**Table 1: Summary of data collection**

Labels	Uploaded	Not uploaded	Total
Killing time and available for viewing notifications	606,760 (51.1%)	29,160 (2.5%)	635,920 (53.6%)
Killing time but unavailable for viewing notifications	135,380 (11.4%)	2,101 (0.2%)	137,481 (11.6%)
Not killing time but available for viewing notifications	202,327 (17.1%)	17,081 (1.4%)	219,408 (18.5%)
Not killing time and unavailable for viewing notifications	118,313 (10.0%)	9,071 (0.8%)	127,384 (10.7%)
Unidentifiable	0 (0.0%)	66,152 (5.6%)	66,152 (5.6%)
Total	1,062,780 (89.6%)	123,565 (10.4%)	1,186,345 (100.0%)

**Table 2: The sensor features used in the study**

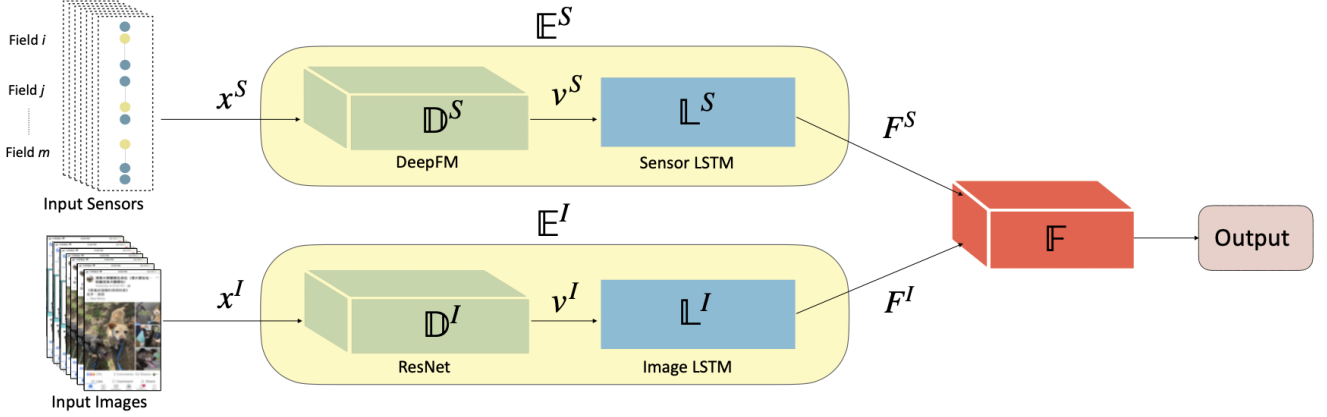
Phone Context	Current Characteristics	Current session characteristics (accumulated up to the current screenshot record)
Transportation Mode	Physical activity (i.e., not moving, on foot, in vehicle, or on bicycle)	Cumulative time of {not moving, on foot, in vehicle, on bicycle}
	Was moving (i.e., on foot, in vehicle, or on bicycle)	Majority of physical activity
Type of Day	Day of the week (0-6)	
	Was weekend (i.e., Saturday, Sunday)	
Time of a Day	Hour of the day in 24-hour notation (0-23)	
	Was meal time (11:00 a.m.-12:59 p.m., 5:00 p.m.-6:59 p.m.)	
Battery Status	Phone battery level	{AVG, STD, MIN, MAX, MED} Phone-battery level
	Phone was charging / not charging	Charging count
	If charging over AC or USB	Cumulative charging time
Screen Time		{AVG, STD, MIN, MAX, MED, SUM} Screen time
Screen Orientation	Portrait / landscape mode	
Foreground App	Name of the app in the foreground	Count and frequency of app switches
	Package name of the app in the foreground	Count of used apps
	Category of the app in the foreground	Cumulative usage time of the 15 most frequently used app categories and all remaining app categories combined into one category group.
Network Info	{WiFi, Mobile} network was available / unavailable	Cumulative time the phone was connected to the {WiFi, Mobile} network
	{Type, operator} of the network the phone connected to	Cumulative time the phone was not connected to any network
	Was connected to the network	
Ringer Mode	Silent / vibrate / normal	Cumulative time of {silent, vibrate, normal} Was adjusted
Audio Mode	Ringling / in call / in communication / normal	Cumulative time of {ringing, in call, in communication, normal} Call count
Stream Volume	Volume of streams, e.g., music playback, notification, phone calls, phone ring, system sounds	{AVG, STD, MIN, MAX, MED} Volume of stream {music playback, notification, phone calls, phone ring, system sounds} Volume of stream {music playback, notification, phone calls, phone ring, system sounds} was adjusted
Call Status	Device call state: idle / off-hook / ringing	
Usage	Current Characteristics	Current session characteristics (accumulated up to the current screenshot record)
Screen-on Events	Count of Screen-on events during the past 180/300/600/900/1,800/3,600 seconds	{count, frequency} of screen-on events
Accessibility Events	Count of {clicking, long-clicking, scrolling, hover enter/exit, setting-input focus, changing-the-text, selecting} events during the past 30/60/180/300/600/900/1,800/3,600 seconds	{count, frequency} of {clicking, long-clicking, scrolling, hover enter/exit, setting-input focus, changing-the-text, selecting} events

Note. \* All time-related calculations were in seconds

being “unidentifiable” (see Table 1). The above distribution cannot perfectly represent the participants’ actual phone usage, insofar as some screenshots were not annotated and/or not uploaded. Nevertheless, we are confident in its general outlines, e.g., that there were more time-killing moments than non-time-killing ones, and

that the participants reported being available to view notifications more frequently during time-killing moments (82.2%) than during non-time-killing ones (63.3%).





**Figure 2: Illustration for the architecture of our proposed model, which takes the input composed of the phone-sensor data and the screenshots (collected within a certain time window, e.g., 30 seconds) and predicts the users’ time-killing moments.**

Because the focus of this paper is on how to predict time-killing moments, it will not systematically discuss the interview data, collected notification data, ESM results, or the results of the three questionnaires that were not related to our approach’s user acceptance. Those other datasets will instead be used in future research.

### 3.6 Feature Selection and Extraction

To predict time-killing moments, we extracted two kinds of feature sets from the phone-sensor data: phone context and user interactions. For each of these feature sets, we created two temporal ranges, one describing the phone at the moment when a screenshot was taken, and the other, describing the characteristics of the phone-use session during which it was taken. We defined a phone-use session as a continuous use of the phone during which any brief screen-off interval was not longer than 45 seconds, based on the findings of van Berkel et al. [89], that using the 45-second threshold separating two sessions was more accurate than the others. Thus, if more than 45 seconds had passed since the last screen-off event, the current usage was considered a new session. In addition, inspired by our interview data and prior research findings [72] suggesting that some phone events or user actions occur intensively during time-killing, we created features that measured the frequency of various types of phone and interaction events during nine past-time windows, ranging from a minimum of 30 seconds to a maximum of 3,600 seconds (e.g., frequency of scrolling within the previous 30 minutes). We excluded data from the first hour of each person’s participation day, because a large portion of such data could not allow us to compute these features. As a result, the final dataset for developing the model consisted of 967,466 annotated screenshots, from which 183 features were derived, as shown in Table 2. The 1,181 apps used during the study by our participants were placed in 56 categories based on their Google Play Store categorizations and prior literature [101].

## 4 MODEL DESIGN

We adopt deep-learning, which learns the pattern in an end-to-end manner. Specifically, our proposed model (shown in Fig. 2) is composed of three main subnetworks: 1) an encoder  $\mathbb{E}^S$  built

upon DeepFM [31] and an LSTM [35] that extracts *sensor features* from phone-sensor data, 2) an encoder  $\mathbb{E}^I$  based on the ResNet and an LSTM that encodes the sequences of screenshots into *visual features*, and 3) a fusion subnetwork  $\mathbb{F}$  that adopts an attention mechanism followed by fully-connected layers to fuse the sensor features and the visual features into the final prediction outcome, i.e., time-killing vs. non-time-killing. Please note that, for the feature encoder  $\mathbb{E}^S$  of phone-sensor data, we chose to use a DeepFM network for our model because its architecture allows it to learn the various interactions among features without requiring extensive feature engineering. Additionally, DeepFM has demonstrated superior performance when modeling sensory data composed of multiple feature fields, making it well-suited for our purposes [31]). For the feature encoder  $\mathbb{E}^I$  for screenshots, on the other hand, we chose to use a ResNet for our model because it is able to effectively address the gradient vanishing problem, allowing the model to be more stable and robust during training. Moreover, for both  $\mathbb{E}^S$  and  $\mathbb{E}^I$  encoders, we leveraged LSTM rather than RNN to handle the sequential input because LSTM is a popular choice for modeling sequential data because of its forget-gate that allows it to focus on the important parts of the input data. More details of these subnetworks are provided in the following sections.

### 4.1 Encoder $\mathbb{E}^S$ of Phone-sensor Data

Given a sequence of phone-sensor data collected at several time steps within a certain time window (ideally these time steps are evenly distributed within a given time window), denoted as  $\mathcal{X}^S = \{x_k^S\}_{k=1}^K$ , where  $K$  is the number of time steps, the encoder  $\mathbb{E}^S$  which is built upon a DeepFM module  $\mathbb{D}^S$  and a 3-layer LSTM module  $\mathbb{L}^S$  turns  $\mathcal{X}^S$  into the sensor feature  $\mathcal{F}^S$ . As our phone-sensor data  $x_k^S$  contain both continuous and categorical values (e.g., a phone battery level is a continuous value, whereas a ringer mode is a categorical value), our DeepFM module  $\mathbb{D}^S$  adopts the DeepFM [31] framework that extracts a feature representation  $v_k^S = \mathbb{D}^S(x_k^S)$  for each  $x_k^S$ . Note that the architecture of our DeepFM module  $\mathbb{D}^S$  is almost identical to the one proposed in [31], except that it uses a 128-dimensional vector in the last fully-connected layer in order

to fit into the size of  $v_k^S$ . Specifically, the feature vectors  $\{v_k^S\}_{k=1}^K$  extracted from the sensor data  $\{x_k^S\}_{k=1}^K$  are sequentially fed into the LSTM module  $\mathbb{L}^S$  to model the temporal variations in  $\{x_k^S\}_{k=1}^K$ , which then generates a 256-dimensional sensor-feature vector  $\mathcal{F}^S$ .

## 4.2 Encoder $\mathbb{E}^I$ of Screenshots

The visual encoder  $\mathbb{E}^I$  which extracts the visual feature  $\mathcal{F}^I$  from a stack of  $K$  screenshots  $\mathcal{X}^I = \{x_k^I\}_{k=1}^K$  is composed of a ResNet module  $\mathbb{D}^I$  and a 3-layer LSTM module  $\mathbb{L}^I$ . All the screenshots are resized to  $224 \times 224$  pixels, regardless of whether they were taken horizontally or vertically; then they are fed into the ResNet module  $\mathbb{D}^I$  to extract the feature representation  $v_k^I = \mathbb{D}^I(x_k^I)$ , where  $\mathbb{D}^I$  adopts the ImageNet-pretrained Resnet-101 backbone and the size of  $v_k^I$  is  $7 \times 7 \times 2048$ . Similar to the procedure of encoding phone-sensor data, these extracted features  $\{v_k^I\}_{k=1}^K$  are taken as a sequential input for the LSTM module  $\mathbb{L}^I$  to derive their visual feature  $\mathcal{F}^I$  (which is 256-dimensional) of  $\mathcal{X}^I$ . For both LSTM modules  $\mathbb{L}^S$  and  $\mathbb{L}^I$ , the dimensions of all the hidden state, cell state, and the hidden layer are set to 512 respectively. Note that although  $\mathbb{L}^S$  and  $\mathbb{L}^I$  have a similar architecture, they are trained independently and do not share any weight.

## 4.3 Fusion Subnetwork $\mathbb{F}$ over Sensor and Visual Features

After obtaining the sensor feature  $\mathcal{F}^S$  and visual feature  $\mathcal{F}^I$  from phone-sensor data  $\mathcal{X}^S$  and screenshots  $\mathcal{X}^I$ , respectively, we used a fusion subnetwork  $\mathbb{F}$  that jointly considers the high-level information from these two features in order to detect participants' time-killing behaviors. To achieve this, instead of concatenating two features and utilizing a simple classifier to perform a multi-modal fusion, we introduced an additional multi-fusion layer that takes both features as inputs to predict the reweighting coefficients  $\alpha^S$  and  $\alpha^I$  (i.e., analogous to the importance) for both feature dimension  $\mathcal{F}^S$  and  $\mathcal{F}^I$ ; The reweighted features, denoted as  $\tilde{\mathcal{F}}^S = \alpha^S \otimes \mathcal{F}^S$  and  $\tilde{\mathcal{F}}^I = \alpha^I \otimes \mathcal{F}^I$ , are then concatenated with the original  $\mathcal{F}^S$  and  $\mathcal{F}^I$ , which are further intertwined by several fully-connected layers to generate the final classification outcome of time-killing or not.

**Training Details.** We adopted a stage-wise training procedure, in which we first trained the encoders,  $\mathbb{E}^S$  and  $\mathbb{E}^I$ , independently, followed by training the fusion subnetwork. Specifically, we first attached a fully connected layer to the end of the encoder  $\mathbb{E}^S$  and  $\mathbb{E}^I$  individually. Then, the layer maps the sensor feature  $\mathcal{F}^S$  and the visual feature  $\mathcal{F}^I$  to the output of time-killing detection respectively, i.e., the whole encoder together with the attached fully connected layer becomes a classification model and can be pre-trained via using our collected dataset and a classification objective of cross-entropy. After pre-training both encoders till they converged, we removed the attached fully connected layers and fixed the weights of encoders. Then we trained the fusion subnetwork  $\mathbb{F}$  via the cross-entropy loss. We chose to follow a stage-wise training procedure because it performs better than training from scratch. We adopted the Adam optimizer [47] for training the model. In pretraining the encoder  $\mathbb{E}^S$ , we set the batch size 512 and the learning rate  $10^{-3}$ ,

while for pretraining the encoder  $\mathbb{E}^I$ , we set a batch size 196 and the learning rate  $10^{-5}$ . Lastly, for training the fusion subnetwork  $\mathbb{F}$ , we set a batch size 196 and the learning rate  $10^{-5}$ . Our model is implemented with PyTorch and trained using 8 Tesla V100 GPU cores.

# 5 THE FUSION MODEL FOR PREDICTING TIME-KILLING MOMENTS

## 5.1 Experiment

**5.1.1 Dataset.** We paired each labeled screenshot with phone-sensor data according to the time at which that screenshot was taken. To predict whether a screenshot was labeled as time-killing or non-time-killing, we used features derived from the screenshots and their paired sensor data 30 seconds (i.e., seven screenshots) prior to the predicted one (and without using any data other than such 30-seconds period). In other words, a sequence of data contained seven data pairs, including the predicted screenshot and the data for predicting it. We made sure that such sequences did not overlap with one another; and that, if a sequence contained fewer than seven data pairs, we padded it to that length seven by using zero padding, i.e., a whole black image.

Each participant contributed a different amount of data. Therefore, to prevent our model from being overly biased towards particular participants who contributed much more data than others, we sampled 20,000 screenshots from each participant to create our training dataset. Such sampling was random, except insofar as we ensured that it contained 1) data collected on both weekends and weekdays, and 2) exactly equal numbers of time-killing and non-time-killing instances. For the testing dataset, on the other hand, we did not seek to strike this balance, but instead followed the original distribution, such that the evaluation of the model would more accurately reflect the time-killing distribution that one would observe in the real world.

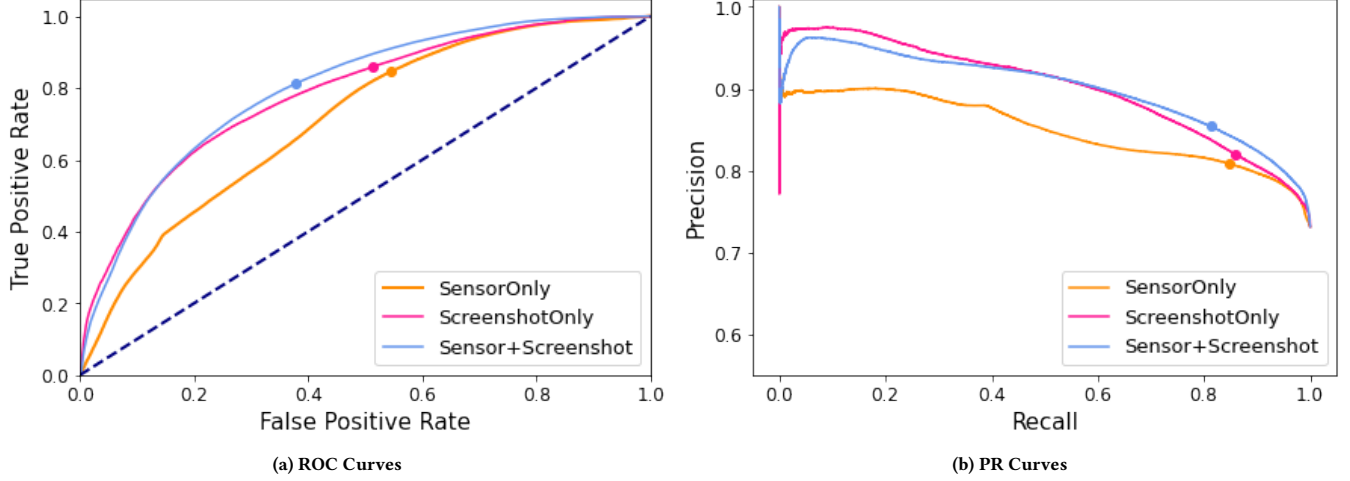
**5.1.2 Evaluation Metrics.** Our testing dataset had more time-killing instances than non-time-killing ones, in the ratio 7:3. We made many computations to compare model performance, but here, we will focus on ROC-curve (Receiver Operating Characteristics) and PR-curve (Precision Recall). The ROC curve plots the true positive rate against the false positive rate at various classification thresholds for time-killing classification, and AUROC, i.e., the area under the ROC curve, indicates better performance where its values are higher. The PR curve allowed us to observe the precision score against the recall score at various classification thresholds. We prioritized the precision of the prediction over recall, because the higher the former is, the fewer non-time-killing moments will be falsely predicted as time-killing moments, and thus, fewer notifications will be mistakenly sent to the user at these moments. For the same reason, we also assessed specificity, which measures the prediction's true negative rate.

**5.1.3 Model Evaluation.** To evaluate the model's performance, we performed three-fold cross-validation on the dataset. As noted earlier, two-thirds of the data from each participant were used for re-sampling, and formed a training dataset, with the rest forming the test dataset. We ensured that when we divided the dataset, the



**Table 3: The three models' time-killing prediction task performance**

Model	Accuracy	Precision	Recall	AUROC	Specificity
Fusion (Sensor+Screenshot)	0.76	0.83	0.81	0.72	0.62
SensorOnly	0.74	0.8	0.85	0.65	0.45
ScreenshotOnly	0.76	0.81	0.86	0.67	0.49

**Figure 3: Two performance measurements of our proposed fusion model (i.e., *Sensor+Screenshot*), its variants (i.e., *SensorOnly* and *ScreenshotOnly*). Note. Point on the curves represents a classification threshold equal to 0.5.**

order among the screenshot and phone-sensor pairs was maintained. In evaluating the performance of the fusion model for predicting time-killing moments, we also compared it against two other models, which respectively used only phone-sensor data and only screenshot data. We describe all three models in more detail below.

- **Fusion (*Sensor+Screenshot*)** - Used both phone-sensor data and screenshot data; model design as described earlier.
- ***SensorOnly*** - Used the phone-sensor data encoded by  $\mathbb{E}^S$  to perform time-killing prediction, with an additional fully connected layer attached to  $\mathbb{E}^S$  acting as the linear classifier.
- ***ScreenshotOnly*** - Used phone-screenshot data encoded by  $\mathbb{E}^I$  to perform time-killing prediction, with an additional fully connected layer attached to  $\mathbb{E}^I$  as a linear classifier.

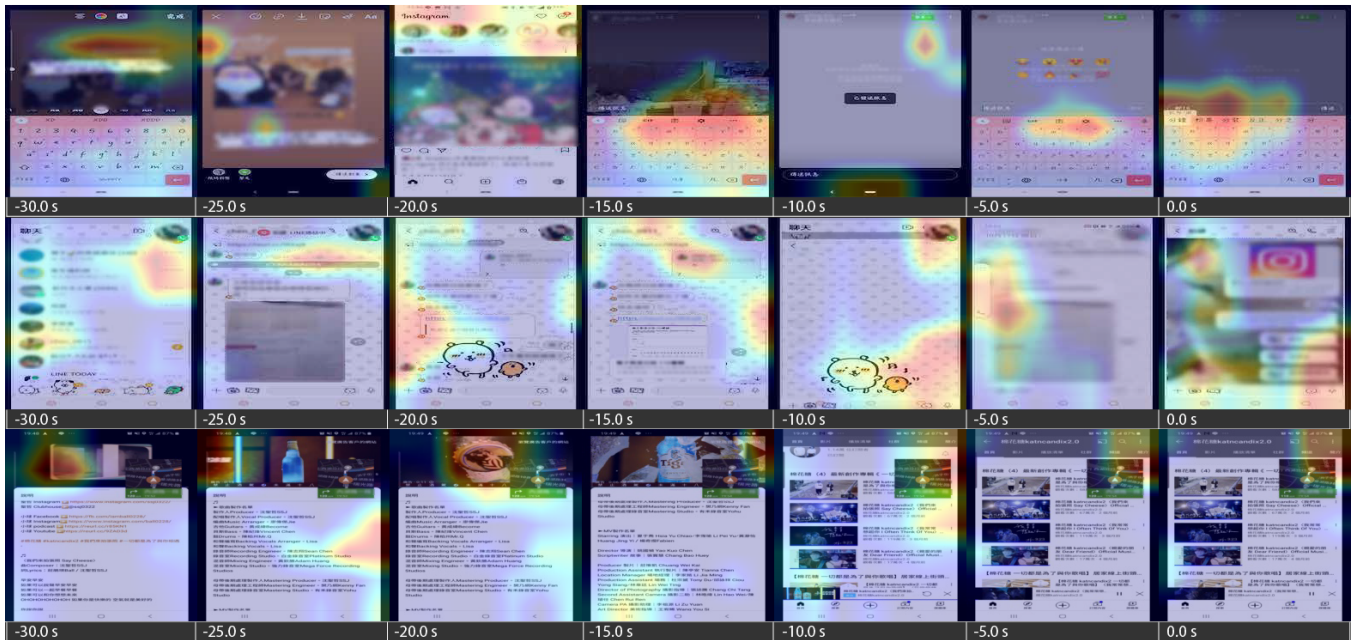
## 5.2 Result

The models' overall performance metrics are presented in Table 3, which uses a classification threshold of 0.5. Fig. 3a and 3b show their ROC curves and PR curves. Overall, the fusion model achieved the best AUROC among the three models, as shown in both Table 3 and Fig. 3a. The fusion model's prediction of a given moment as being a time-killing one was the most accurate among the three models. Moreover, as shown by the PR curves, the fusion model achieved higher precision with high recall than the other two models, and its specificity score was also significantly higher than theirs. These results imply that taking account of both sensor data and screenshot data makes it less likely to falsely predict a non-time-killing moment

as a time-killing one than when only one source or the other is considered. The *SensorOnly* model achieved the lowest performance across all metrics except recall. As shown in both Fig. 3a and Fig. 3b, it had notably lower precision across classification thresholds than the other two models, suggesting that many of the moments it predicted as time-killing were incorrect. This was because some phone states or interactions that occurred mainly during time-killing by one group of users often occurred during the non-time-killing-moments of another group, making it difficult to differentiate these two kinds of moments across users with different behavior patterns: a phenomenon that will be explored in Section 6. The *ScreenshotOnly* model, on the other hand, had a better ability to distinguish between them, suggesting that phone-screenshot data were more informative about time-killing moments than sensor data were. That being said, the inclusion of phone-sensor data improved the performance of the fusion model.

### 5.2.1 Investigation for Potential Concern on Model Memorization.

Here, we provided additional investigation into the potential concern that our model may have simply memorized the training data labels rather than learning to predict time-killing moments from multi-modal inputs. To examine the issue, we trained our proposed fusion model using randomly assigned labels from our dataset. However, the model did not converge, which implies that it could not discern meaningful patterns in the data and thus was less susceptible to overfitting (for further details on the model's performance, please refer to Appendix A.) Furthermore, our fusion model trained



**Figure 4: Example attention maps, produced by Grad-CAM [79] and the *ScreenshotOnly* model, comprising a sequence of time-killing screenshots in the top row, and two sequence of non-time-killing ones in the middle and bottom row. The sensitive content on the images have been blurred for privacy reasons.**

on our dataset with participants’ annotations showed small differences between its training and test performance (only 18% in terms of AUROC), indicating a certain degree of generalizability. Based on these observations, we conclude that our proposed fusion model did not suffer from the issue of model memorization.

### 5.3 Examples of How Fusing Phone-sensor Data and Screenshots Helped us Recognize Time-killing vs. Non-time-killing Behaviors

In our view, the fact that fusing phone-sensor data and screenshots yielded the best performance in detecting time-killing moments implies that these two data sources, to some extent, complemented each other. To explore this possible phenomenon, we inspected cases in our test dataset in which a time-killing moment was correctly detected by the fusion model, but incorrectly detected by either or both of the *SensorOnly* and *ScreenshotOnly* models.

To facilitate this exploration and our sense-making of these cases, we created attention maps from the final convolution layer of the *ScreenshotOnly* model, using a popular technique called Grad-CAM [79]. These attention maps helped us to identify regions in the screenshots that the fusion/*ScreenshotOnly* model considered influential in its time-killing behavior detection. For instance, the top row of Fig. 4 provides examples in which both the *ScreenshotOnly* and fusion models correctly recognized a time-killing moment that was mistaken as a non-time-killing one by the *SensorOnly* model. We suspect that the *SensorOnly* model incorrectly recognized such data sequences because a series of text-changed events were detected, despite in an Instagram application, which was more likely to occur when not killing time. On the other hand, we suspect that

the *ScreenshotOnly* model detected it correctly because it recognized the layout of the user interface of Instagram’s Story feature, which tended to be associated with time-killing moments.

The middle row in Fig. 4, meanwhile, shows a distinctive case in which both the *SensorOnly* and fusion models correctly predicted a non-time-killing moment that was incorrectly predicted by the *ScreenshotOnly* model as a time-killing one. We suspect that the *ScreenshotOnly* model misinterpreted this screenshot sequence as a time-killing moment because it recognized the layout of LINE. In this case, the participant was discussing an assignment with others via text conversation; however, the participant was talking to her friend (prompted by the communication icon in the upper-right corner) while, which was often associated with time-killing moments. The *ScreenshotOnly* model did not attend to the communication icon in all sequences of the screenshots, but instead relied mostly on the layout of the chat room. Nevertheless, we observed that the relevant information was captured in the user’s phone-sensor data: specifically, by the call status and the change of the call volume (as the sixth screenshot shows). Knowing these pieces of information enabled the fusion model to recognize this non-time-killing moment correctly.

Finally, an example where both the *ScreenshotOnly* model and the *SensorOnly* model detected incorrectly, but the fusion model still detected correctly is shown in the bottom row of Fig. 4. In this non-time-killing case, the participant was sitting in a moving vehicle and was using a navigation app while simultaneously watching a video on Youtube. We suspect that the *ScreenshotOnly* model misinterpreted this series of screenshots as a time-killing moment because it recognized Youtube’s layout but ignored the

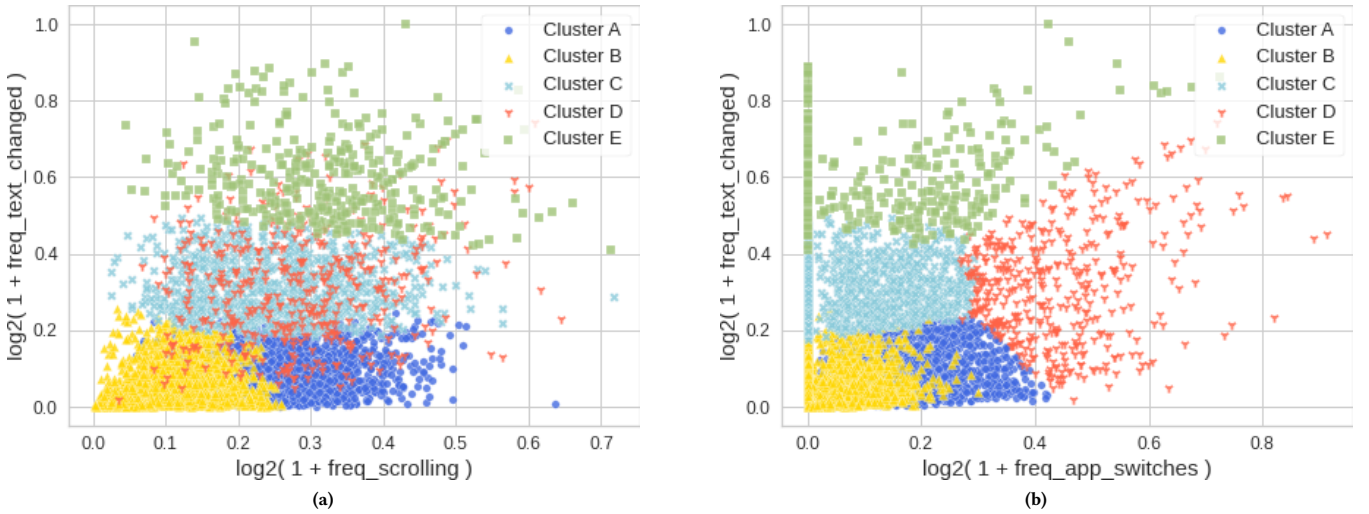


Figure 5: Scatter plot of session clusters, grouped based on in-session behavioral characteristics

navigation map in picture-in-picture mode, possibly due to the similar color of the video screen to that of the map. We suspect that the *SensorOnly* model mistook this situation as time-killing based on the application and physical-activity information. However, the fusion model correctly identified this moment as a non-time-killing moment, possibly because it considered both the physical-activity information and the use of the navigation functionality. There were many similar instances; however, these three vivid examples should suffice to explain why the fusion model performed best at detecting time-killing moments across nearly all metrics.

## 6 TAILORING FUSION MODELS TO USERS CLUSTERED BY PHONE-USAGE BEHAVIOR

We learned from the interviews that various distinct time-killing patterns existed among our participants, who could be grouped based on similarities in their phone interactions, task choices, task switching, audio modes, and so on. Because we could not group participants based on their time-killing behaviors, assuming that such a label might not be obtainable during system runtime, we instead grouped them based on their phone-usage behavior. Below, we present the group-based model we arrived at using clustering, followed by model evaluation and our observations about the features of these individual models.

### 6.1 Clustering Participants Based on their Phone-usage Behavior

We employed two stages of the k-means method [57] to group users hierarchically. First, inspired by Isaacs et al. [40], we employed clustering to identify distinct phone-usage behavioral patterns. Then, we clustered participants according to how often their phone use belonged to each of the identified phone-usage patterns, based on the assumption that a user was likely to display more than one such pattern.

**6.1.1 Clustering Phone-usage Behavior.** Inspired by previous work [40] that used the concept of *sessions* to cluster phone usage, we generated participants' sessions based on the rule suggested by van Berkel et al. [89], that is, we divided pairs of sessions using a separation threshold of 45 seconds. This approach resulted in a total of 5,266 phone-usage sessions. For each of them, inspired by our interview, we computed nine features: 1) session duration, 2) screen-switching frequency, 3) application-switching frequency, 4) scroll-event frequency, 5) text-change event frequency, 6) maximum and 7) minimum gap durations for scroll events, and 8) maximum and 9) minimum gap durations for text-change events. We then applied k-means, and used the Elbow method [86] to determine the number of clusters. This revealed the optimal number of clusters as five. The 5,266 phone-usage sessions were grouped into these five clusters, named A, B, C, D, and E, in descending order by cluster size, whose sizes were 1,882, 1,664, 942, 417, and 361, respectively.

The five groups mainly differed in terms of how actively their members used their phones. For example, Fig. 5a shows the distribution of the frequency of the participants' scrolling by the frequency of text-changes in a session, colored according to the cluster they belonged to; and Fig. 5b, the distribution of the same frequency by the frequency of app switching. For example, Cluster B contained inactive phone-usage sessions characterized by low frequency of text changes, scrolling, and app switching. In contrast, the sessions in Cluster A were marked by low-frequency text changes and relatively low-frequency app switching, but with high-frequency scrolling. Furthermore, the sessions in Cluster D exhibited the highest frequency of app switching among all the clusters.

**6.1.2 Clustering Users by the Proportions of Five Behavioral Outcomes.** Having clustered similar phone-usage behaviors as described above, we observed that most users performed all five behaviors, but in varying proportions. Therefore, to group users with similar overall mobile-phone usage, we calculated the proportions of each user's five outcome behaviors, and used those proportions to

**Table 4: Experimental Results: Clustering Participants by Behavioral and Temporal Characteristics**

Group	Accuracy			Precision			Recall			AUCROC			Specificity		
Group 1	0.70	0.73	0.73	0.81	0.81	0.88	0.85	0.81	0.80	0.68	0.76	0.77	0.38	0.54	0.59
Group 2	0.75	0.77	0.77	0.85	0.89	0.91	0.84	0.87	0.84	0.70	0.75	0.78	0.46	0.44	0.55
Group 3	0.80	0.77	0.78	0.91	0.95	0.93	0.87	0.82	0.82	0.68	0.75	0.72	0.39	0.48	0.50
Group 4	0.72	0.74	0.77	0.71	0.74	0.77	0.78	0.78	0.79	0.70	0.73	0.77	0.63	0.69	0.74
Average	0.74	0.75	0.76	0.82	0.84	0.87	0.83	0.82	0.81	0.69	0.75	0.76	0.47	0.54	0.60
General model	0.74	0.76	0.76	0.80	0.81	0.83	0.85	0.86	0.81	0.65	0.67	0.72	0.45	0.49	0.62

Note. The white, light gray, and dark gray backgrounds indicate the results for *SensorOnly*, *ScreenshotOnly*, and Fusion (*SensorOnly+ScreenshotOnly*) models, respectively.

cluster users. The same k-means and Elbow methods as described above were performed, and the resulting k-value for user clustering was 4. Thus, we divided our participants into four groups, with the number of participants being 11, 11, 9, and 5 respectively. The positive (time-killing) and negative (non-time-killing) instance ratios of those four groups were 13:6, 3:1, 81:19, and 3:2, respectively.

## 6.2 Overall Performance of the Cluster-based Models

We built the same fusion model for each of the four user groups, and examined each one’s average performance separately via the same three-fold cross-validation approach mentioned in Section 5.1. Table 4, which presents the respective performance of those four models along with their average performance, shows that both their average AUROC (0.76) and precision (0.87) were higher than those of the general model (AUROC: 0.72, precision: 0.83). In terms of individual model performance, all four models’ AUROC values were at least as good as that of the general model, with three significantly higher than it; and three models’ precision values were also higher than the general model’s. These results suggest that it is beneficial to divide users into groups according to their phone-usage behavior and build a time-killing prediction model for each such user group.

We also looked at the correlations between time-killing moments and phone-sensor features for each user groups. Table 5 shows the 15 non-category features most highly correlated (either positively or negatively) with time-killing moments, by user group. In each such group, some features were more correlated with time-killing moments than their counterparts in the general model, suggesting that clustering users into behavioral groups was also beneficial to time-killing prediction: i.e., doing so revealed features correlated with time-killing moments specifically for certain participants, which would not have been revealed had they not been divided into groups. That being said, the results in Table 4 also show that the performances of the four models varied, suggesting that some user groups’ time-killing moments might be more difficult than others’ to predict. We discuss each user group’s model performance and time-killing behaviors in the next section.

## 6.3 Model Performance and Behavior by User Group

First, Group 2’s fusion model achieved the best AUROC among the four user groups. It is also worth noting that Group 2’s *ScreenshotOnly* model achieved better performance than its *SensorOnly*

model for all metrics except specificity, suggesting that it was accurate in predicting time-killing moments but less so in predicting non-time-killing moments. When observing features correlated with time-killing moments in Group 2, we found that screen-on events, the number of calls, the volume of communication, and the volume of ringtones negatively correlated with the members’ time-killing moments. In other words, when participants in this group were not killing time, they tended to increase the audio volume of their phones and frequently turned their screens on and off. Their switching to normal ringer mode was also positively correlated with time-killing moments; this reflected their higher usage of the two relatively quiet modes, vibrate and silent, when they were not killing time. All of this implies that these participants’ non-time-killing moments were more often associated with making calls. As prior research has reported a high association between quiet ringer modes and proactive phone-checking behaviors [13], the Group 2 behaviors we observed could have indicated participants checking their phones frequently to avoid missing calls and/or notifications. The fact that these behaviors might have been captured better by sensor data than by screenshot data could explain why – in this group alone – the *SensorOnly* model performed better at identifying non-time-killing moments (i.e., higher specificity; true negative rate) than the *ScreenshotOnly* model did.

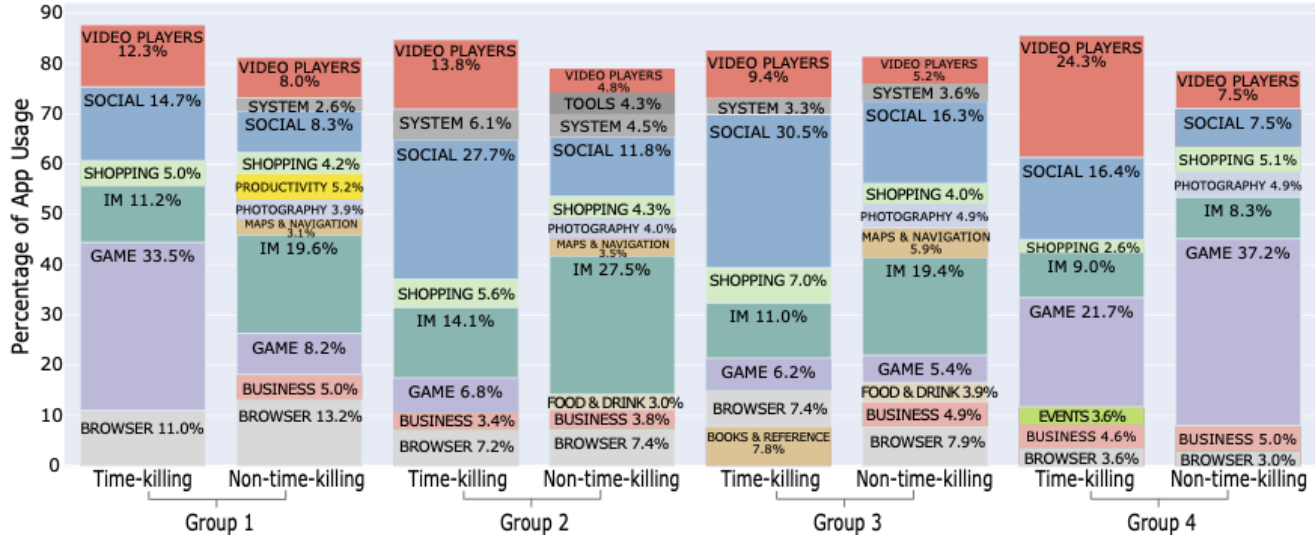
Secondly, Group 1’s and Group 4’s fusion models both achieved AUROCs of 0.77. However, the reasons for these two models achieving this same value differed dramatically, as shown by the significant differences in their other performance metrics. Specifically, whereas Group 1’s fusion model achieved significantly higher precision (0.88) than Group 4’s fusion model (0.77), Group 4’s fusion model performed particularly well in specificity (0.74): significantly higher than any of the other models. In other words, Group 1’s fusion model was better at predicting its members’ time-killing moments, whereas Group 4’s fusion model was better at predicting its members’ non-time-killing moments. As shown in Table 5, Group 4’s key features for prediction were predominantly battery-related ones, which were negatively correlated with time-killing moments. Also, while the feature number of charging events is not displayed in Table 5, its correlation was -0.27 – higher than many other features in other user groups – suggesting that this group’s members’ non-time-killing moments were associated with high values of battery-related features, very likely linked to battery-charging at non-time-killing moments. We further observed the app-usage distribution of Group 4’s members, as shown in Fig. 6, and found that they played games much more often during non-time-killing



**Table 5: The 15 non-category features most highly correlated (either positively or negatively) with time-killing moments, by user group**

Group 1	corr.	Group 2	corr.	Group 3	corr.	Group 4	corr.	General Model	corr.
call_count	-0.25	screen-on_past_900s	-0.22	T_photography_apps	-0.18	battery_level	-0.40	T_vibration	-0.17
is_adjusted_vol_noti	-0.25	screen-on_past_600s	-0.22	scrolling_past_3600s	0.15	AVG_battery	-0.40	scrolling_past_3600	0.15
is_adjusted_vol_ring	-0.25	screen-on_past_300s	-0.21	screen-on_past_600s	-0.15	MED_battery	-0.40	call_count	-0.15
T_Silent	0.24	screen-on_past_1800s	-0.21	screen-on_past_1800s	-0.15	MIN_battery	-0.39	scrolling_past_1800s	0.14
is_adjusted_vol_voicecall	-0.24	call_count	-0.21	screen-on_past_900s	-0.14	MAX_battery	-0.37	T_InComm.	-0.14
is_adjusted_vol_sys	-0.24	screen-on_past_3600s	-0.21	scrolling_past_1800s	0.14	MAX_vol_music	0.36	MIN_battery	-0.14
T_game_apps	0.24	screen-on_past_180s	-0.21	T_normal_ringer	0.14	AVG_vol_music	0.35	T_ringer_silent	0.13
MAX_vol_ring	-0.21	T_InComm.	-0.19	screen-on_past_300s	-0.14	MED_vol_music	0.33	MED_battery	-0.13
MAX_vol_noti	-0.21	T_normal_audio	0.19	T_map_apps	-0.13	MIN_vol_ring	0.32	AVG_battery	-0.13
MAX_vol_sys	-0.20	T_ringtones	-0.16	scrolling_count	0.13	strm_vol_music	0.32	scrolling_past_900s	0.13
STD_vol_sys	-0.19	MAX_vol_sys	-0.16	long-clicking_count	0.13	AVG_vol_ring	0.31	T_photography_apps	-0.12
STD_vol_noti	-0.19	MAX_vol_noti	-0.16	T_social_apps	0.13	MED_vol_ring	0.31	scrolling_past_600s	0.12
STD_vol_ring	-0.19	T_mobile_network	0.15	scrolling_past_900s	0.13	strm_vol_ring	0.31	battery_level	-0.12
MIN_vol_voicecall	0.18	freq_text_changed	-0.15	scrolling_past_600s	0.13	AVG_vol_sys	0.31	focus_event_past_3600s	0.12
T_InComm.	-0.16	MAX_vol_ring	-0.15	screen-on_past_180s	-0.12	MAX_vol_sys	0.30	MAX_vol_music	0.12

Note. The T prefix indicates the cumulative time; the green and blue backgrounds indicate positive and negative correlations, respectively, with darker colors indicating higher correlations.

**Figure 6: Percentage of application categories used by each user group when killing time and not killing time Note. Categories 1) related to the launcher and 2) with percentages <2.5% are not displayed.**

moments than during time-killing ones (37.2% vs. 21.7%); this percentage was also the greatest among the four groups. When we took a closer look at the games they played, we found that 88.6% of their game time during non-time-killing moments was taken up by Pokémon Go, and 95% of the time, they were correctly predicted by the model to be non-time-killing moments. Possibly because of the large quantity of this distinctive behavior during non-time-killing moments, the Group 4 fusion model's true negative rate was particularly high. Interestingly, Group 1 was another group whose members spent considerable time playing games, but in contrast to the Group 4 members, they were much more likely to do so during time-killing moments, and rarely did so in non-time-killing ones. The Group 1 participants also often used social-media applications,

watched videos, and engaged in IM during their time-killing moments, but seldom did so during their non-time-killing moments. It is noteworthy that Group 1's *SensorOnly* model achieved much poorer specificity than its *ScreenshotOnly* model, suggesting that the fusion model relied heavily on screenshot data to recognize non-time-killing moments.

Finally, Group 3's fusion model achieved the lowest AUROC (0.72) among the four groups' fusion models, an outcome even worse than that of its *ScreenshotOnly* model (0.75). This was because, despite having the highest precision among the four groups, it had a particularly low true-negative rate. In part, this distinctive characteristic of the model might be attributed to it having the most unbalanced dataset: 80% of the instances were time-killing moments, which might have made it tend to predict Group 3 members'

moments as time-killing ones. The chief reason this user group's dataset was unbalanced was that its members used their phones mainly to kill time. Notably, correlations between features and time-killing moments were also lowest for Group 3, suggesting that its members' time-killing behaviors tended to be diverse and not associated with solid patterns. Also, when we looked into the Group 3 members' app-usage distribution in their time-killing vs. non-time-killing moments, we found it to be likewise highly diverse and evenly distributed. In short, a lack of clear patterns in phone usage during time-killing moments might explain the relatively low performance of this user group's *SensorOnly* model, which in turn seemed to lead the fusion model astray.

## 7 DISCUSSION

To take advantage of the benefits of distinguishing phone use for killing time from use for specific purposes, we developed an Android app that collects smartphone users' phone use data, including screenshots and sensor data, as well as time-killing annotations from users. We then used them to build a model to detect time-killing phone use and evaluated its performance. We found that a deep-learning model fusing screenshot and phone-sensor data could achieve a precision of 0.83 and an AUROC of 0.72. However, there are two even more important takeaways from our results.

First, leveraging both phone-sensor and screenshot data in time-killing detection can achieve better performance than using either of these data sources by itself, including Chen et al. [16]'s screenshot-based model, particularly good at distinguishing non-time-killing moments from time-killing-ones. This vital capability could help prevent a future system from sending users digital content or messages at falsely detected time-killing moments. This is particularly important for avoiding sending intervention messages that remind users to reduce or pause their current phone use during productive phone use.

We suspect that the fusion model has this capability because, to some extent, sensor features and the visual information extracted from screenshots complement each other. For example, while screenshots do not inform us about various aspects of phone status such as battery, voice, and network, and are thus unhelpful in recognizing certain time-killing moments characterized by these features, they contain rich and unambiguous contextual information about the activity a user is undertaking during time-killing and non-time-killing-moments alike. It may be possible that the complementary nature of the two data sources might be also helpful in the detection of other behavior/moments on phones and other devices, such as interruptible moments [2, 62, 64, 95], moments of boredom [72], moments of micro-waiting [12, 41], moments of normative dissociation or absentminded use of the phone [6, 84], and/or breakpoint [1, 34, 63]. In addition to identifying these moments, we believe that our approach can usefully be employed in future research more generally in fields that have already leveraged screenshot data to analyze broader behavior patterns, such as smartphone users' media consumption [27].

The second key takeaway of our results is the benefits of clustering users according to their phone-usage behaviors and then

tailoring fusion models to the resulting clusters. That is, the user-cluster-based model outperformed both Chen et al. [16]'s non-cluster-based model and the general model that was based on all users' data. However, Chen et al. [16]'s sample only comprised six smartphone users, implying that clustering may not be necessary for such a small dataset. Compared to the general model, on the other hand, we attribute the superior performance achieved via this group-based approach to the diverse time-killing patterns of our participants, which sometimes were even opposite to each other, confusing the general model. A vivid example of this phenomenon was that participants in Group 1 tended to play games during time-killing moments, whereas those in Group 4 tended to do so at non-time-killing ones. Unsurprisingly, after these participants were separated, both their groups' respective models achieved significantly higher AUROC than the general model did.

The profound benefits of building user-cluster-based models were also manifested in the complementarity between sensor data and screenshot data. This was because some participants' behavior changes were associated more with changes in sensor data than phone-screen data, while others were the opposite. For example, Groups 1, 2, and 4 exhibited phone-usage behavior that was clearly associated with time-killing moments (see Table 5). Thus, the extra information from sensors complemented that from screenshots, because each captured some aspect(s) of time-killing moments that the other missed. In contrast, Group 3's fusion model achieved lower AUROC than its *ScreenshotOnly* model. This may provide an example of conflicting instead of complementary information provided by the two data sources: i.e., the sensor information collected from this group of participants did not assist the fusion model in distinguishing time-killing moments from non-time-killing ones. This can also be seen from the low correlations between sensor features and this group's time-killing behaviors.

These results suggest that the effectiveness of phone sensor data for predicting time-killing moments depends heavily on phone users' behavior patterns. They also imply that decisions about whether it is worthwhile to engage in the privacy-intrusive and phone-resource-demanding process of capturing users' screenshots should take into account the objective of such detection. For example, the *SensorOnly* models of both Group 1 and Group 3 achieved higher recall than their fusion models; so, if one's objective were to capture as many time-killing moments as possible, capturing only sensor information on the phones of users of the Group 1 and Group 3 types would be adequate to purpose. On the other hand, if one's main aim was to reduce falsely detected time-killing moments, leveraging screenshot data would generally be more helpful.

In sum, we believe the approach we have presented in this paper will help researchers and practitioners interested in leveraging screenshot data for predicting or detecting specific smartphone-user behavior and moments. In particular, we expect it to be useful for those interested in detecting time-killing moments for delivering content to which people may not be receptive at other moments, such as crowdsourcing tasks [61] or questionnaires<sup>7</sup>. On the other hand, if the researchers intend to make these moments more productive, they can send productive content such as reading and/or material [22]; if the purpose is to reduce habitual, addictive, or

<sup>7</sup><https://play.google.com/store/apps/details?id=com.google.android.apps.paidtasks>



compulsive phone use [87], to or reduce phubbing in social contexts [85], which are likely to be associated with time-killing moments, researchers can also send intervention prompts or reminders during these detected periods. Moreover, researchers can also analyze when and how frequently time-killing phone use occurs in users' daily lives.

## 8 LIMITATIONS

This research has several limitations. First, it is possible that there is an overlap between the moments when people are killing time and the moments when they feel bored. However, our study did not use the experience sampling approach to collect data on participants' boredom, so we cannot provide quantitative evidence on the differences between these two moments. Second, despite the same approach having been leveraged by other research for studying phone usage [9, 74], given the sensitive data this current research collected, our participants might feel conscious of the data being collected and adjust their phone use behavior during the study period. Moreover, while our participants were allowed to choose which screenshots not to upload, they were not provided with the same option for sensor data. Although all of our participants were fully informed and aware of this limitation, it was still likely that this limitation might have also impacted the participants' behavior. Future research should consider these ethical considerations, for example, by blurring screenshots on the phone before uploading them to the server or only uploading abstract features.

Third, our data collection inherently relied on participants' in-the-wild annotations, which might not always be reliable. Even though we have tested the employed annotation interface and mechanism with twelve pilot participants, we could not deny the possibility that, given participants' annotations being added post hoc rather than in situ, participants' annotated data might be subject to recall errors. Also, although we strove to ease our participants' screenshot-annotation burdens – on the grounds that otherwise, their compliance would have been much lower – it is possible that the user-friendly drag-and-drop interface we developed to address this problem facilitated mislabeling. That is, some subjects might have considered it more efficient, at least in some cases, to label a whole block of data at once. Indeed, our observations of the dataset indicated that some screenshots were mistakenly labeled, which could account for some of our models' apparent inaccuracies.

Fourth, we did not include other features that can be used to infer the users' activity on the phone screen such as activity name, picture-in-picture mode, or names of UI elements. It is possible that these features might have helped the *SensorOnly* model recognize the user's activity, which can be explored in future work.

Fifth, our dataset was established based on a small ( $n=36$ ) sample of smartphone users in Taiwan; all our participants were under 55 years old, and half of them were students. As a result, it is unclear whether our models' detection performance can be generalized to populations that display even more diverse time-killing behaviors or different phone-usage patterns. For example, we believe that such behaviors may be clustered into more types than the four that our small sample suggested. Thus, longer-term and larger-scale data collection could lead to more reliable results.

Sixth, we utilized the raw screenshots and sensor data collected directly from participants' phones without applying additional privacy protection techniques. This may have influenced some participants' willingness to share the captured screenshots. In future investigations, encryption techniques can be employed to decrease the identifiable information in the gathered screenshot and sensor data for ethical and data-completeness considerations. For screenshots in particular, blurring or pixelizing can be applied to screenshots before they are uploaded. It should be noted, however, that the performance of the proposed models may drop on blurred or pixelated screenshots, and the extent to which performance is affected would require further analysis. As a result, instead of modifying the fidelity of the images, an alternative is to extract high-level information from the screenshots on participants' phones or even run the models on their smartphones. However, if these operations take place on the phone, researchers would need to take the phone's computation power and battery into account.

Finally, although we collected other aspects of the participants' tendencies and characteristics that might have affected their time-killing behaviors, such as their demographic characteristics and occupations, we did not include them in this paper. We also did not analyze their notification-attendance behavior during time-killing moments. These aspects should be given greater attention in future studies.

## 9 CONCLUSION

In this paper, we leveraged both phone-sensor and screenshot data to predict time-killing moments using deep-learning techniques. We collected 967,466 pairs of annotated phone-sensor data and screenshots from 36 participants over 14 days for training our time-killing models. We have shown that phone-sensor and screenshot data each have their advantages in such detection tasks; and that integrating these two data sources can yield better model performance than using either of them by itself can. We also have shown that separating users into groups according to their phone-usage patterns and building individual time-killing models for each group can achieve strong overall performance, with most group-specific models also achieving better performance than a general model. Additionally, we have provided insights into how and why the effectiveness of sensor data and phone screenshots as a basis for detecting time-killing moments vary across different user groups. We believe this paper offers a good starting point for researchers and practitioners who are interested in leveraging both screenshot and sensor data in their prediction tasks, and that it will be useful for practitioners who want to incorporate this detection into their applications.

## ACKNOWLEDGMENTS

We greatly thank our study participants for their contributions and National Center for High-performance Computing (NCHC), Taiwan, for providing computational and storage resources. This project is supported by National Science Technology Council, Taiwan (110-2222-E-A49-008-MY3, 107-2218-E-009 -030 -MY3, 111-2628-E-A49 -018 -MY4), as well as the Higher Education Sprout Project of National Yang Ming Chiao Tung University and Ministry of Education (MOE), Taiwan.

## REFERENCES

- [1] Piotr D. Adamczyk and Brian P. Bailey. 2004. *If Not Now, When? The Effects of Interruption at Different Moments within Task Execution*. Association for Computing Machinery, New York, NY, USA, 271–278. <https://doi.org/10.1145/985692.985727>
- [2] Piotr D. Adamczyk, Shamsi T. Iqbal, and Brian P. Bailey. 2005. A method, system, and tools for intelligent interruption management. In *Proceedings of the 4th international workshop on Task models and diagrams (TAMODIA '05)*. Association for Computing Machinery, New York, NY, USA, 123–126. <https://doi.org/10.1145/1122935.1122959>
- [3] Elena Agapie, Jaime Teevan, and Andrés Monroy-Hernández. 2015. Crowdsourcing in the field: A case study using local crowds for event reporting. In *Third AAAI Conference on Human Computation and Crowdsourcing (HCOMP '15)*.
- [4] Yeslam Al-Saggaf, Rachel MacCulloch, and Karl Wiener. 2019. Trait Boredom Is a Predictor of Phubbing Frequency. *Journal of Technology in Behavioral Science* 4 (09 2019). <https://doi.org/10.1007/s41347-018-0080-4>
- [5] Yeslam Al-Saggaf and Sarah B O'Donnell. 2019. Phubbing: Perceptions, reasons behind, predictors, and impacts. *Human Behavior and Emerging Technologies* 1, 2 (2019), 132–140.
- [6] Amanda Baughan, Mingrui Ray Zhang, Raveena Rao, Kai Lukoff, Anastasia Schaadhardt, Lisa D Butler, and Alexis Hiniker. 2022. “I Don’t Even Remember What I Read”: How Design Influences Dissociation on Social Media. In *CHI Conference on Human Factors in Computing Systems*. 1–13.
- [7] Tony Beltramelli. 2018. Pix2code: Generating Code from a Graphical User Interface Screenshot. In *Proceedings of the ACM SIGCHI Symposium on Engineering Interactive Computing Systems (Paris, France) (EICS '18)*. Association for Computing Machinery, New York, NY, USA, Article 3, 6 pages. <https://doi.org/10.1145/3220134.3220135>
- [8] Matthias Böhmer, Brent Hecht, Johannes Schöning, Antonio Krüger, and Gernot Bauer. 2011. Falling Asleep with Angry Birds, Facebook and Kindle: A Large Scale Study on Mobile Application Usage. In *Proceedings of the 13th International Conference on Human Computer Interaction with Mobile Devices and Services (Stockholm, Sweden) (MobileHCI '11)*. Association for Computing Machinery, New York, NY, USA, 47–56. <https://doi.org/10.1145/2037373.2037383>
- [9] Miriam Brinberg, Nilam Ram, Xiao Yang, Mu-Jung Cho, S Shyam Sundar, Thomas N Robinson, and Byron Reeves. 2021. The idiosyncrasies of everyday digital lives: Using the Human Screenome Project to study user behavior on smartphones. *Computers in Human Behavior* 114 (2021), 106570.
- [10] Barry Brown, Moira McGregor, and Eric Laurier. 2013. *iPhone in Vivo: Video Analysis of Mobile Device Use*. Association for Computing Machinery, New York, NY, USA, 1031–1040. <https://doi.org/10.1145/2470654.2466132>
- [11] Barry Brown, Moira McGregor, and Donald McMillan. 2014. 100 Days of iPhone Use: Understanding the Details of Mobile Device Use. In *Proceedings of the 16th International Conference on Human-Computer Interaction with Mobile Devices & Services (Toronto, ON, Canada) (MobileHCI '14)*. Association for Computing Machinery, New York, NY, USA, 223–232. <https://doi.org/10.1145/2628363.2628377>
- [12] Carrie J. Cai, Anji Ren, and Robert C. Miller. 2017. WaitSuite: Productive Use of Diverse Waiting Moments. *ACM Trans. Comput.-Hum. Interact.* 24, 1, Article 7 (March 2017), 41 pages. <https://doi.org/10.1145/3044534>
- [13] Yung-Ju Chang and John C. Tang. 2015. Investigating Mobile Users’ Ringer Mode Usage and Attentiveness and Responsiveness to Communication. In *Proceedings of the 17th International Conference on Human-Computer Interaction with Mobile Devices and Services (Copenhagen, Denmark) (MobileHCI '15)*. Association for Computing Machinery, New York, NY, USA, 6–15. <https://doi.org/10.1145/2785830.2785852>
- [14] Chunyang Chen, Sidong Feng, Zhenchang Xing, Linda Liu, Shengdong Zhao, and Jinshui Wang. 2019. Gallery D.C.: Design Search and Knowledge Discovery through Auto-Created GUI Component Gallery. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 180 (Nov. 2019), 22 pages. <https://doi.org/10.1145/3359282>
- [15] Chunyang Chen, Ting Su, Guozhu Meng, Zhenchang Xing, and Yang Liu. 2018. From UI Design Image to GUI Skeleton: A Neural Machine Translator to Bootstrap Mobile GUI Implementation. In *Proceedings of the 40th International Conference on Software Engineering (Gothenburg, Sweden) (ICSE '18)*. Association for Computing Machinery, New York, NY, USA, 665–676. <https://doi.org/10.1145/3180155.3180240>
- [16] Yu-Chun Chen, Keui-Chun Kao, Yu-Jen Lee, Faye Shih, Wei-Chen Chiu, and Yung-Ju Chang. 2021. Killing-Time Detection from Smartphone Screenshots. In *Adjunct Proceedings of the 2021 ACM International Joint Conference on Pervasive and Ubiquitous Computing and Proceedings of the 2021 ACM International Symposium on Wearable Computers (Virtual, USA) (UbiComp '21)*. Association for Computing Machinery, New York, NY, USA, 15–16. <https://doi.org/10.1145/3460418.3479295>
- [17] Pei-Yu Peggy Chi, Matthew Long, Akshay Gaur, Abhimanyu Deora, Anurag Batra, and Daphne Luong. 2019. Crowdsourcing Images for Global Diversity. In *Proceedings of the 21st International Conference on Human-Computer Interaction with Mobile Devices and Services (Taipei, Taiwan) (MobileHCI '19)*. Association for Computing Machinery, New York, NY, USA, Article 79, 10 pages. <https://doi.org/10.1145/3338286.3347546>
- [18] Chia-En Chiang, Yu-Chun Chen, Fang-Yu Lin, Felicia Feng, Hao-An Wu, Hao-Ping Lee, Chang-Hsuan Yang, and Yung-Ju Chang. 2021. “I Got Some Free Time”: Investigating Task-Execution and Task-Effort Metrics in Mobile Crowdsourcing Tasks. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, Article 648, 14 pages. <https://doi.org/10.1145/3411764.3445477>
- [19] Woohyeok Choi, Sangkeun Park, Duyeon Kim, Youn-kyung Lim, and Uichin Lee. 2019. Multi-Stage Receptivity Model for Mobile Just-In-Time Health Intervention. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 3, 2, Article 39 (June 2019), 26 pages. <https://doi.org/10.1145/3328910>
- [20] Mihaly Csikszentmihalyi. 2000. *Beyond boredom and anxiety*. Jossey-bass.
- [21] Tilman Dingler and Martin Pietlot. 2015. I’ll Be There for You: Quantifying Attentiveness towards Mobile Messaging. In *Proceedings of the 17th International Conference on Human-Computer Interaction with Mobile Devices and Services (Copenhagen, Denmark) (MobileHCI '15)*. Association for Computing Machinery, New York, NY, USA, 1–5. <https://doi.org/10.1145/2785830.2785840>
- [22] Tilman Dingler, Benjamin Tag, Sabrina Lehrer, and Albrecht Schmidt. 2018. Reading Scheduler: Proactive Recommendations to Help Users Cope with Their Daily Reading Volume. In *Proceedings of the 17th International Conference on Mobile and Ubiquitous Multimedia (Cairo, Egypt) (MUM 2018)*. Association for Computing Machinery, New York, NY, USA, 239–244. <https://doi.org/10.1145/3282894.3282917>
- [23] Tilman Dingler, Dominik Weber, Martin Pietlot, Jennifer Cooper, Chung-Cheng Chang, and Niels Henze. 2017. Language Learning On-the-Go: Opportune Moments and Design of Mobile Microlearning Sessions. In *Proceedings of the 19th International Conference on Human-Computer Interaction with Mobile Devices and Services (Vienna, Austria) (MobileHCI '17)*. Association for Computing Machinery, New York, NY, USA, Article 28, 12 pages. <https://doi.org/10.1145/3098279.3098565>
- [24] Trinh Minh Tri Do, Jan Blom, and Daniel Gatica-Perez. 2011. Smartphone Usage in the Wild: A Large-Scale Analysis of Applications and Context. In *Proceedings of the 13th International Conference on Multimodal Interfaces (Alicante, Spain) (ICMI '11)*. Association for Computing Machinery, New York, NY, USA, 353–360. <https://doi.org/10.1145/2070481.2070550>
- [25] John D Eastwood, Alexandra Frischen, Mark J Fenske, and Daniel Smilek. 2012. The unengaged mind: Defining boredom in terms of attention. *Perspectives on Psychological Science* 7, 5 (2012), 482–495.
- [26] Andreas Elpidorou. 2018. The bored mind is a guiding mind: Toward a regulatory theory of boredom. *Phenomenology and the Cognitive Sciences* 17 (2018), 455–484.
- [27] Hossein Falaki, Ratul Mahajan, Srikanth Kandula, Dimitrios Lymberopoulos, Ramesh Govindan, and Deborah Estrin. 2010. Diversity in Smartphone Usage. In *Proceedings of the 8th International Conference on Mobile Systems, Applications, and Services (San Francisco, California, USA) (MobiSys '10)*. Association for Computing Machinery, New York, NY, USA, 179–194. <https://doi.org/10.1145/1814433.1814453>
- [28] Robert Fisher and Reid Simmons. 2011. Smartphone Interruptibility Using Density-Weighted Uncertainty Sampling with Reinforcement Learning. In *2011 10th International Conference on Machine Learning and Applications and Workshops, Vol. 1*. 436–441. <https://doi.org/10.1109/ICMLA.2011.128>
- [29] Jon Froehlich, Mike Y. Chen, Sunny Consolvo, Beverly Harrison, and James A. Landay. 2007. MyExperience: A System for in Situ Tracing and Capturing of User Feedback on Mobile Phones. In *Proceedings of the 5th International Conference on Mobile Systems, Applications and Services (San Juan, Puerto Rico) (MobiSys '07)*. Association for Computing Machinery, New York, NY, USA, 57–70. <https://doi.org/10.1145/1247660.1247670>
- [30] Ralph R Greenson. 1953. On boredom. *Journal of the American Psychoanalytic Association* 1, 1 (1953), 7–21.
- [31] Hui Feng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: a factorization-machine based neural network for CTR prediction. In *International Joint Conference on Artificial Intelligence (IJCAI)*.
- [32] Alexis Hiniker, Shwetak N. Patel, Tadayoshi Kohno, and Julie A. Kientz. 2016. Why Would You Do That? Predicting the Uses and Gratifications behind Smartphone-Usage Behaviors. In *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing (Heidelberg, Germany) (UbiComp '16)*. Association for Computing Machinery, New York, NY, USA, 634–645. <https://doi.org/10.1145/2971648.2971762>
- [33] Bo-Jhang Ho, Bharathan Balaji, Mehmet Koseoglu, Sandeep Sandha, Siyou Pei, and Mani Srivastava. 2020. Quick Question: Interrupting Users for Microtasks with Reinforcement Learning. *arXiv:2007.09515 [cs]* (July 2020). <http://arxiv.org/abs/2007.09515> arXiv: 2007.09515
- [34] Joyce Ho and Stephen S. Intille. 2005. *Using Context-Aware Computing to Reduce the Perceived Burden of Interruptions from Mobile Devices*. Association for Computing Machinery, New York, NY, USA, 909–918. <https://doi.org/10.1145/>

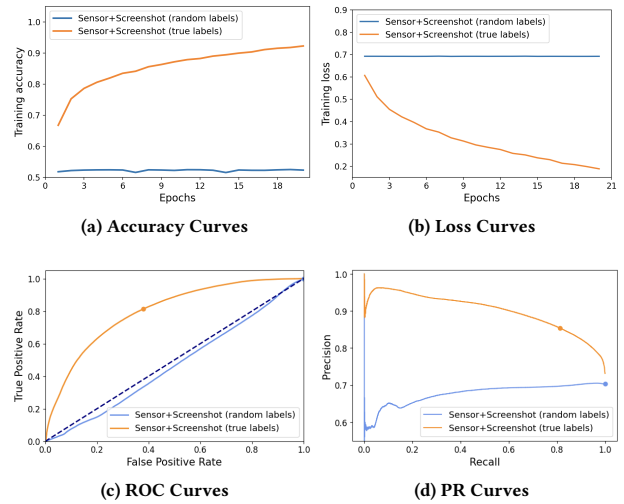
- 1054972.1055100
- [35] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [36] Cynthia A. Hoffner and Sangmi Lee. 2015. Mobile Phone Use, Emotion Regulation, and Well-Being. *Cyberpsychology, Behavior, and Social Networking* 18, 7 (2015), 411–416. <https://doi.org/10.1089/cyber.2014.0487> arXiv:<https://doi.org/10.1089/cyber.2014.0487> PMID: 26167841.
- [37] Nanna Inie and Mircea F Lungu. 2021. Aiki - Turning Online Procrastination into Microlearning. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (CHI '21). Association for Computing Machinery, New York, NY, USA, Article 369, 13 pages. <https://doi.org/10.1145/3411764.3445202>
- [38] Shamsi T. Iqbal and Brian P. Bailey. 2007. Understanding and developing models for detecting and differentiating breakpoints during interactive tasks. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 697–706. <https://doi.org/10.1145/1240624.1240732>
- [39] Shamsi T. Iqbal and Brian P. Bailey. 2011. Oasis: A framework for linking notification delivery to the perceptual structure of goal-directed tasks. *ACM Transactions on Computer-Human Interaction* 17, 4 (Dec. 2011), 15:1–15:28. <https://doi.org/10.1145/1879831.1879833>
- [40] Ellen Isaacs, Alan Walendowski, Steve Whittaker, Diane J. Schiano, and Candace Kamm. 2002. The Character, Functions, and Styles of Instant Messaging in the Workplace. In *Proceedings of the 2002 ACM Conference on Computer Supported Cooperative Work* (New Orleans, Louisiana, USA) (CSCW '02). Association for Computing Machinery, New York, NY, USA, 11–20. <https://doi.org/10.1145/587078.587081>
- [41] Ellen Isaacs, Nicholas Yee, Diane J Schiano, Nathan Good, Nicolas Ducheneaut, and Victoria Bellotti. 2009. Mobile microwaiting moments: The role of context in receptivity to content while on the go. *PARC white paper* (2009) 10 (2009).
- [42] Chakajkla Jesdabodi and Walid Maalej. 2015. Understanding Usage States on Mobile Devices. In *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing* (Osaka, Japan) (UbiComp '15). Association for Computing Machinery, New York, NY, USA, 1221–1225. <https://doi.org/10.1145/2750858.2805837>
- [43] Simon L. Jones, Denzil Ferreira, Simo Hosio, Jorge Goncalves, and Vassilis Kostakos. 2015. Revisitation Analysis of Smartphone App Use. In *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing* (Osaka, Japan) (UbiComp '15). Association for Computing Machinery, New York, NY, USA, 1197–1208. <https://doi.org/10.1145/2750858.2807542>
- [44] Kleomenis Katevas, Ioannis Arapakis, and Martin Pielot. 2018. Typical Phone Use Habits: Intense Use Does Not Predict Negative Well-Being. In *Proceedings of the 20th International Conference on Human-Computer Interaction with Mobile Devices and Services* (Barcelona, Spain) (MobileHCI '18). Association for Computing Machinery, New York, NY, USA, Article 11, 13 pages. <https://doi.org/10.1145/3229434.3229441>
- [45] Jürgen Kawalek, Annegret Stark, and Marcel Riebeck. 2008. A New Approach to Analyze Human-Mobile Computer Interaction. *J. Usability Studies* 3, 2 (Feb. 2008), 90–98.
- [46] Ronald C Kessler, Lenard Adler, Minnie Ames, Olga Demler, Steve Faraone, EVA Hiripi, Mary J Howes, Robert Jin, Kristina Secnik, Thomas Spencer, et al. 2005. The World Health Organization Adult ADHD Self-Report Scale (ASRS): a short screening scale for use in the general population. *Psychological medicine* 35, 2 (2005), 245–256.
- [47] Diederik P. Kingma and Jimmy Ba. 2017. Adam: A Method for Stochastic Optimization. arXiv:1412.6980 [cs.LG]
- [48] Vassilis Kostakos, Denzil Ferreira, Jorge Goncalves, and Simo Hosio. 2016. Modelling Smartphone Usage: A Markov State Transition Model. In *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing* (Heidelberg, Germany) (UbiComp '16). Association for Computing Machinery, New York, NY, USA, 486–497. <https://doi.org/10.1145/2971648.2971669>
- [49] Philipp Krieter. 2019. Can I Record Your Screen? Mobile Screen Recordings as a Long-Term Data Source for User Studies. In *Proceedings of the 18th International Conference on Mobile and Ubiquitous Multimedia* (Pisa, Italy) (MUM '19). Association for Computing Machinery, New York, NY, USA, Article 23, 10 pages. <https://doi.org/10.1145/3365610.3365618>
- [50] Philipp Krieter and Andreas Breiter. 2018. Analyzing Mobile Application Usage: Generating Log Files from Mobile Screen Recordings. In *Proceedings of the 20th International Conference on Human-Computer Interaction with Mobile Devices and Services* (Barcelona, Spain) (MobileHCI '18). Association for Computing Machinery, New York, NY, USA, Article 9, 10 pages. <https://doi.org/10.1145/3229434.3229450>
- [51] Nuning Kurniasih. 2017. Internet Addiction, Lifestyle or Mental Disorder? A Phenomenological Study on Social Media Addiction in Indonesia. *KNE Social Sciences* 2, 4 (Jun. 2017), 135–144. <https://doi.org/10.18502/kss.v2i4.879>
- [52] Hao-Ping Lee, Kuan-Yin Chen, Chih-Heng Lin, Chia-Yu Chen, Yu-Lin Chung, Yung-Ju Chang, and Chien-Ru Sun. 2019. Does Who Matter? Studying the Impact of Relationship Characteristics on Receptivity to Mobile IM Messages. Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3290605.3300756>
- [53] Jian Li, Andrew Lepp, and Jacob E. Barkley. 2015. Locus of control and cell phone use: Implications for sleep quality, academic performance, and subjective well-being. *Computers in Human Behavior* 52 (2015), 450–457. <https://doi.org/10.1016/j.chb.2015.06.021>
- [54] Tong Li, Mingyang Zhang, Hancheng Cao, Yong Li, Sasu Tarkoma, and Pan Hui. 2020. "What Apps Did You Use?": Understanding the Long-Term Evolution of Mobile App Usage. In *Proceedings of The Web Conference 2020* (Taipei, Taiwan) (WWW '20). Association for Computing Machinery, New York, NY, USA, 66–76. <https://doi.org/10.1145/3366423.3380095>
- [55] Yu-Hsuan Lin, Li-Ren Chang, Yang-Han Lee, Hsien-Wei Tseng, Terry BJ Kuo, and Sue-Huei Chen. 2014. Development and Validation of the Smartphone Addiction Inventory (SPAI). *PLoS one* 9, 6 (2014), e98312.
- [56] Kai Lukoff, Cissy Yu, Julie Kientz, and Alexis Hiniker. 2018. What Makes Smartphone Use Meaningful or Meaningless? *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 2, 1, Article 22 (March 2018), 26 pages. <https://doi.org/10.1145/3191754>
- [57] James MacQueen et al. 1967. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, Vol. 1. Oakland, CA, USA, 281–297.
- [58] Donald McMillan, Moira McGregor, and Barry Brown. 2015. From in the Wild to in Vivo: Video Analysis of Mobile Device Use. In *Proceedings of the 17th International Conference on Human-Computer Interaction with Mobile Devices and Services* (Copenhagen, Denmark) (MobileHCI '15). Association for Computing Machinery, New York, NY, USA, 494–503. <https://doi.org/10.1145/2785830.2785883>
- [59] Abhinav Mehrotra, Robert Hendley, and Mirco Musolesi. 2016. PrefMiner: Mining User's Preferences for Intelligent Mobile Notification Management. In *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing* (Heidelberg, Germany) (UbiComp '16). Association for Computing Machinery, New York, NY, USA, 1223–1234. <https://doi.org/10.1145/2971648.2971747>
- [60] Varun Mishra, Florian Künzler, Jan-Niklas Kramer, Elgar Fleisch, Tobias Kowatsch, and David Kotz. 2021. Detecting receptivity for mhealth interventions in the natural environment. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 5, 2 (2021), 1–24.
- [61] Aditi Misra, Aaron Gooze, Kari E. Watkins, M. Asad, and Christopher A. Le Dantec. 2014. Crowdsourcing and Its Application to Transportation Data Collection and Management. *Transportation Research Record* 2414 (2014), 1 – 8.
- [62] Christopher Monk, Deborah Boehm-Davis, and J. Trafton. 2002. The Attentional Costs of Interrupting Task Performance at Various Stages. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 46 (Sept. 2002). <https://doi.org/10.1177/154193120204602210>
- [63] Tadashi Okoshi, Julian Ramos, Hiroki Nozaki, Jin Nakazawa, Anind K Dey, and Hideyuki Tokuda. 2015. Attelia: Reducing user's cognitive load due to interruptive notifications on smart phones. In *2015 IEEE International Conference on Pervasive Computing and Communications (PerCom)*. IEEE, 96–104. <https://doi.org/10.1109/PERCOM.2015.7146515>
- [64] Tadashi Okoshi, Kota Tsubouchi, Masaya Taji, Takanori Ichikawa, and Hideyuki Tokuda. 2017. Attention and engagement-awareness in the wild: A large-scale study with adaptive notifications. *2017 IEEE International Conference on Pervasive Computing and Communications (PerCom)* (2017), 100–110.
- [65] Antti Oulasvirta, Tye Rattenbury, Lingyi Ma, and Eva Raita. 2012. Habits Make Smartphone Use More Pervasive. *Personal Ubiquitous Comput.* 16, 1 (Jan. 2012), 105–114. <https://doi.org/10.1007/s00779-011-0412-2>
- [66] Leysia Palen and Marilyn Salzman. 2002. Voice-Mail Diary Studies for Naturalistic Data Capture under Mobile Conditions. In *Proceedings of the 2002 ACM Conference on Computer Supported Cooperative Work* (New Orleans, Louisiana, USA) (CSCW '02). Association for Computing Machinery, New York, NY, USA, 87–95. <https://doi.org/10.1145/587078.587092>
- [67] Chunjong Park, Junsung Lim, Juho Kim, Sung-Ju Lee, and Dongman Lee. 2017. Don't Bother Me. I'm Socializing! A Breakpoint-Based Smartphone Notification System. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing* (Portland, Oregon, USA) (CSCW '17). Association for Computing Machinery, New York, NY, USA, 541–554. <https://doi.org/10.1145/2998181.2998189>
- [68] Veljko Pejovic and Mirco Musolesi. 2014. InterruptMe: Designing Intelligent Prompting Mechanisms for Pervasive Applications. In *Proceedings of the 2014 ACM International Joint Conference on Pervasive and Ubiquitous Computing* (Seattle, Washington) (UbiComp '14). Association for Computing Machinery, New York, NY, USA, 897–908. <https://doi.org/10.1145/2632048.2632062>
- [69] Martin Pielot, Linas Baltrunas, and Nuria Oliver. 2015. Boredom-Triggered Proactive Recommendations. In *Proceedings of the 17th International Conference on Human-Computer Interaction with Mobile Devices and Services Adjunct* (Copenhagen, Denmark) (MobileHCI '15). Association for Computing Machinery, New York, NY, USA, 1106–1110. <https://doi.org/10.1145/2786567.2794340>

- [70] Martin Pielot, Bruno Cardoso, Kleomenis Katevas, Joan Serra, Aleksandar Matic, and Nuria Oliver. 2017. Beyond Interruptibility: Predicting Opportune Moments to Engage Mobile Phone Users. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 1, 3 (Sept. 2017), 91:1–91:25. <https://doi.org/10.1145/3130956>
- [71] Martin Pielot, Rodrigo de Oliveira, Haewoon Kwak, and Nuria Oliver. 2014. Didn't You See My Message? Predicting Attentiveness to Mobile Instant Messages. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Toronto, Ontario, Canada) (CHI '14). Association for Computing Machinery, New York, NY, USA, 3319–3328. <https://doi.org/10.1145/2556288.2556973>
- [72] Martin Pielot, Tilman Dingler, Jose San Pedro, and Nuria Oliver. 2015. When attention is not scarce - detecting boredom from mobile phone usage. In *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing (UbiComp '15)*. Association for Computing Machinery, New York, NY, USA, 825–836. <https://doi.org/10.1145/2750858.2804252>
- [73] Benjamin Poppinga, Wilko Heuten, and Susanne Boll. 2014. Sensor-based identification of opportune moments for triggering notifications. *IEEE Pervasive Computing* 13, 1 (2014), 22–29.
- [74] Nilam Ram, Xiao Yang, Mu-Jung Cho, Miriam Brinberg, Fiona Muirhead, Byron Reeves, and Thomas N Robinson. 2020. Screenomics: A new approach for observing and studying individuals' digital lives. *Journal of adolescent research* 35, 1 (2020), 16–50.
- [75] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. 2016. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 779–788.
- [76] Byron Reeves, Nilam Ram, Thomas N Robinson, James J Cummings, C Lee Giles, Jennifer Pan, Agnese Chiatti, Mj Cho, Katie Roehrick, Xiao Yang, et al. 2021. Screenomics: A framework to capture and analyze personal life experiences and the ways that technology shapes them. *Human-Computer Interaction* 36, 2 (2021), 150–201.
- [77] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems* 28 (2015), 91–99.
- [78] Hillol Sarker, Moushumi Sharmin, Amin Ahsan Ali, Md. Mahbubur Rahman, Rummana Bari, Syed Monwar Hossain, and Santosh Kumar. 2014. Assessing the Availability of Users to Engage in Just-in-Time Intervention in the Natural Environment. In *Proceedings of the 2014 ACM International Joint Conference on Pervasive and Ubiquitous Computing* (Seattle, Washington) (UbiComp '14). Association for Computing Machinery, New York, NY, USA, 909–920. <https://doi.org/10.1145/2632048.2636082>
- [79] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization. In *IEEE International Conference on Computer Vision (ICCV)*.
- [80] Jeremiah Smith, Anna Lavygina, Jiefei Ma, Alessandra Russo, and Naranker Dulay. 2014. Learning to Recognise Disruptive Smartphone Notifications. In *Proceedings of the 16th International Conference on Human-Computer Interaction with Mobile Devices and Services* (Toronto, ON, Canada) (MobileHCI '14). Association for Computing Machinery, New York, NY, USA, 121–124. <https://doi.org/10.1145/2628363.2628404>
- [81] Julian Steil, Philipp Müller, Yusuke Sugano, and Andreas Bulling. 2018. Forecasting user attention during everyday mobile interactions using device-integrated and wearable sensors. *Proceedings of the 20th International Conference on Human-Computer Interaction with Mobile Devices and Services* (Sep 2018). <https://doi.org/10.1145/3229434.3229439>
- [82] Andriy A Struk, Jonathan SA Carriere, J Allan Cheyne, and James Danckert. 2017. A short boredom proneness scale: Development and psychometric properties. *Assessment* 24, 3 (2017), 346–359.
- [83] John C. Tang, Sophia B. Liu, Michael Muller, James Lin, and Clemens Drews. 2006. Unobtrusive but Invasive: Using Screen Recording to Collect Field Data on Computer-Mediated Interaction. In *Proceedings of the 2006 20th Anniversary Conference on Computer Supported Cooperative Work* (Banff, Alberta, Canada) (CSCW '06). Association for Computing Machinery, New York, NY, USA, 479–482. <https://doi.org/10.1145/1180875.1180948>
- [84] Nada Terzimehić, Luke Haliburton, Philipp Greiner, Albrecht Schmidt, Heinrich Hussmann, and Ville Mäkelä. 2022. MindPhone: Mindful Reflection at Unlock Can Reduce Absentminded Smartphone Use. In *Designing Interactive Systems Conference*. 1818–1830.
- [85] Naundefneda Terzimehić, Luke Haliburton, Philipp Greiner, Albrecht Schmidt, Heinrich Hussmann, and Ville Mäkelä. 2022. MindPhone: Mindful Reflection at Unlock Can Reduce Absentminded Smartphone Use. In *Designing Interactive Systems Conference* (Virtual Event, Australia) (DIS '22). Association for Computing Machinery, New York, NY, USA, 1818–1830. <https://doi.org/10.1145/3532106.3533575>
- [86] Robert L Thorndike. 1953. Who belongs in the family. In *Psychometrika*. Citeseer.
- [87] Jonathan A. Tran, Katie S. Yang, Katie Davis, and Alexis Hiniker. 2019. Modeling the Engagement-Disengagement Cycle of Compulsive Phone Use. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk) (CHI '19). Association for Computing Machinery, New York, NY, USA, 1–14. <https://doi.org/10.1145/3290605.3300542>
- [88] Liam D. Turner, Stuart M. Allen, and Roger M. Whitaker. 2017. Reachable but not receptive: Enhancing smartphone interruptibility prediction by modelling the extent of user engagement with notifications. *Pervasive and Mobile Computing* 40 (2017), 480–494. <https://doi.org/10.1016/j.pmcj.2017.01.011>
- [89] Niels van Berkel, Chu Luo, Theodoros Anastopoulos, Denzil Ferreira, Jorge Goncalves, Simo Hosio, and Vassilis Kostakos. 2016. A Systematic Assessment of Smartphone Usage Gaps. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (San Jose, California, USA) (CHI '16). Association for Computing Machinery, New York, NY, USA, 4711–4721. <https://doi.org/10.1145/2858036.2858348>
- [90] Steven Van Canneyt, Marc Bron, Andy Haines, and Mounia Lalmas. 2017. Describing Patterns and Disruptions in Large Scale Mobile App Usage Data. In *Proceedings of the 26th International Conference on World Wide Web Companion* (Perth, Australia) (WWW '17 Companion). International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 1579–1584. <https://doi.org/10.1145/3041021.3051113>
- [91] Wijnand AP Van Tilburg and Eric R Igou. 2012. On boredom: Lack of challenge and meaning as distinct boredom experiences. *Motivation and Emotion* 36 (2012), 181–194.
- [92] Aku Visuri, Niels van Berkel, Chu Luo, Jorge Goncalves, Denzil Ferreira, and Vassilis Kostakos. 2017. Predicting Interruptibility for Manual Data Collection: A Cluster-Based User Model. In *Proceedings of the 19th International Conference on Human-Computer Interaction with Mobile Devices and Services* (Vienna, Austria) (MobileHCI '17). Association for Computing Machinery, New York, NY, USA, Article 12, 14 pages. <https://doi.org/10.1145/3098279.3098532>
- [93] Sara Alida Volkmer and Eva Lerner. 2019. Unhappy and addicted to your phone? – Higher mobile phone use is associated with lower well-being. *Computers in Human Behavior* 93 (2019), 210–218. <https://doi.org/10.1016/j.chb.2018.12.015>
- [94] Heli Vääätäjä and Paul Egglesstone. 2012. Briefing news reporting with mobile assignments: perceptions, needs and challenges. In *Proceedings of the ACM 2012 conference on Computer Supported Cooperative Work (CSCW '12)*. Association for Computing Machinery, New York, NY, USA, 485–494. <https://doi.org/10.1145/2145204.2145280>
- [95] Dominik Weber, Alexandra Voit, Gisela Kollotzek, and Niels Henze. 2019. Annotif: A System for Annotating Mobile Notifications in User Studies. In *Proceedings of the 18th International Conference on Mobile and Ubiquitous Multimedia* (Pisa, Italy) (MUM '19). Association for Computing Machinery, New York, NY, USA, Article 24, 12 pages. <https://doi.org/10.1145/3365610.3365611>
- [96] Thomas D. White, Gordon Fraser, and Guy J. Brown. 2019. Improving Random GUI Testing with Image-Based Widget Detection. In *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis* (Beijing, China) (ISSTA 2019). Association for Computing Machinery, New York, NY, USA, 307–317. <https://doi.org/10.1145/3293882.3330551>
- [97] Qiang Xu, Jeffrey Erman, Alexandre Gerber, Zhuoqing Mao, Jeffrey Pang, and Shobha Venkataraman. 2011. Identifying Diverse Usage Behaviors of Smartphone Apps. In *Proceedings of the 2011 ACM SIGCOMM Conference on Internet Measurement Conference* (Berlin, Germany) (IMC '11). Association for Computing Machinery, New York, NY, USA, 329–344. <https://doi.org/10.1145/2068816.2068847>
- [98] Xiao Yang, Nilam Ram, Thomas Robinson, and Byron Reeves. 2019. Using Screenshots to Predict Task Switching on Smartphones. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk) (CHI EA '19). Association for Computing Machinery, New York, NY, USA, 1–6. <https://doi.org/10.1145/3290607.3313089>
- [99] Nalingna Yuan, Heidi M Weeks, Rosa Ball, Mark W Newman, Yung-Ju Chang, and Jenny S Radesky. 2019. How much do parents actually use their smartphones? Pilot study comparing self-report to passive sensing. *Pediatric research* 86, 4 (2019), 416–418.
- [100] Xiaoyi Zhang, Lilian de Greef, Amanda Swearngin, Samuel White, Kyle Murray, Lisa Yu, Qi Shan, Jeffrey Nichols, Jason Wu, Chris Fleizach, Aaron Everitt, and Jeffrey P Bigham. 2021. *Screen Recognition: Creating Accessibility Metadata for Mobile Applications from Pixels*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3411764.3445186>
- [101] Sha Zhao, Julian Ramos, Jianrong Tao, Ziwen Jiang, Shijian Li, Zhaohui Wu, Gang Pan, and Anind K. Dey. 2016. Discovering Different Kinds of Smartphone Users through Their Application Usage Behaviors. In *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing* (Heidelberg, Germany) (UbiComp '16). Association for Computing Machinery, New York, NY, USA, 498–509. <https://doi.org/10.1145/2971648.2971696>
- [102] Eilish Duke and Christian Montag. 2017. Smartphone addiction, daily interruptions and self-reported productivity. *Addictive Behaviors Reports* 6 (2017), 90–95. <https://doi.org/10.1016/j.abrep.2017.07.002>

## 10 APPENDIX

### A MODEL MEMORIZATION EVALUATION

Fig. 7a and 7b depict the learning curves of the fusion model presented in Section 5.2.1, which was trained on our dataset with randomly assigned labels, for the purpose of investigating potential model memorization. In addition, Fig. 7c and 7d show the model performance in terms of ROC and PR curves, respectively. The model trained with true labels is denoted as “true labels” in the figures, while the one trained with random labels is denoted as “random labels.” During the investigation, we found that the model trained on randomly assigned labels failed to converge, despite our efforts to optimize hyperparameters such as the learning rate and weight decay to prevent overfitting. This suggests that the model was unable to identify meaningful patterns in the data and was, therefore, unlikely to overfit. Hence, we conclude that our proposed model learned informative patterns of time-killing from the data rather than memorizing them, as it would have been able to memorize the random labels if that were the case. Therefore, it results in significantly better performance than the model trained on randomized labels, as shown in Figures 7c and 7d.



**Figure 7: (a) The training accuracy and (b) the training loss of our fusion model, *Sensor+Screenshot*, when trained on true or random labels are plotted on the respective curves. The performance of the fusion model on the test set is evaluated using the (c) ROC and (d) PR curves. Note. Points on the ROC and PR curves represent a classification threshold equal to 0.5.**