

Android Mobile OS Snooping By Samsung, Xiaomi, Huawei and Realme Handsets

Haoyu Liu¹, Paul Patras¹, Douglas J. Leith²

¹University of Edinburgh, UK

²Trinity College Dublin, Ireland

6th October 2021

Abstract—The privacy of mobile apps has been extensively studied, but much less attention has been paid to the privacy of the mobile OS itself. A mobile OS may communicate with servers to check for updates, send telemetry and so on. We undertake an in-depth analysis of the data sent by six variants of the Android OS, namely those developed by Samsung, Xiaomi, Huawei, Realme, LineageOS and /e/OS. We find that even when minimally configured and the handset is idle these vendor-customized Android variants transmit substantial amounts of information to the OS developer and also to third-parties (Google, Microsoft, LinkedIn, Facebook etc) that have pre-installed system apps. While occasional communication with OS servers is to be expected, the observed data transmission goes well beyond this and raises a number of privacy concerns. There is no opt out from this data collection.

I. INTRODUCTION

The analysis of whether mobile apps disclose sensitive information to their associated back-end servers has been the focus of much research [1], [2], [3], [4], [5], especially with a view to risks such user de-anonymization, location tracking, behaviour profiling, and cross-linking of data by different stakeholders in the device/software supply chain. In contrast, the disclosure of information at operating system level has received almost no attention and is not well understood. Mobile OS behaviour has come to the fore only recently, with analyses of the Google-Apple Exposure Notification (GAEN) system that underpins Covid contract tracing apps [6] and following revelations of mass surveillance of journalists, politicians, and human rights activists though spyware exploiting zero-touch vulnerabilities (see the Pegasus project [7]).

We report on an in depth measurement study of the data shared by a range of popular proprietary variants of the Android OS, namely those developed by Samsung, Xiaomi, Huawei and Realme¹. In addition, we report on the data shared by the LineageOS and /e/OS open-source variants of Android. Samsung currently has by far the largest share of this market, followed by Xiaomi, Huawei and Oppo (the parent company of Realme) [8]. LineageOS is probably the most popular open-source Android variant, currently used on around 30M handsets,² while /e/OS is a new privacy-focused fork of LineageOS.

¹Note that we study the European models of handsets from Samsung, Xiaomi, Huawei and Realme and use the handsets within Europe. The data collection behaviour on models targeted at other regions may, or may not, differ.

²<https://stats.lineageos.org/>, accessed 31st July 2021

It is worth noting that much of the functionality of the Android OS³ is provided by so-called system apps. These are privileged pre-installed apps that the OS developer bundles with the OS. System apps cannot be deleted (they are installed on a protected read-only disk partition) and can be granted enhanced rights/permissions not available to ordinary apps such as those that a user might install. It is common for Android to include pre-installed third-party system apps, i.e. apps not written by the OS developer. One example is the so-called GApps package of Google apps (which includes Google Play Services, Google Play store, Google Maps, Youtube etc). Other examples include pre-installed system apps from Microsoft, LinkedIn, Facebook and so on.

We intercept and analyse the data traffic sent by the Android OS, including by pre-installed system apps, in a range of scenarios. We focus on defining simple scenarios that can be applied uniformly to the handsets studied (so allowing direct comparisons) and that generate reproducible behaviour. We assume a privacy-conscious but busy/non-technical user, who when asked does not select options that share data but otherwise leaves handset settings at their default value. This means that the user has opted out of diagnostics/analytics/user experience improvement data collection and has not logged in to an OS vendor user account. The user also does not make use of optional services such as cloud storage, find my phone etc. Essentially, the handset is just being used to make and receive phone calls and texts. This provides a baseline for privacy analysis, and we expect that the level of data sharing may well be larger for a less privacy-conscious user and/or a user who makes greater use of the services on a handset.

We find that the Samsung, Xiaomi, Huawei and Realme Android variants all transmit a substantial volume of data to the OS developer (i.e. Samsung etc) and to third-party parties that have pre-installed system apps (including Google, Microsoft, Heytap, LinkedIn, Facebook). LineageOS sends similar volumes of data to Google as these proprietary Android variants, but we do not observe the LineageOS developers themselves collecting data nor pre-installed system apps other than those of Google. Notably, /e/OS sends no information to Google or other third parties and sends essentially no information to the /e/OS developers.

While it is perhaps unsurprising that a privacy-focused OS such as /e/OS collects almost no data, it nevertheless provides a useful baseline and establishes that extensive data collection

³By Android OS we mean the distribution as installed on a handset, not just the kernel.

TABLE I
SUMMARY OF DATA COLLECTION BY EACH ANDROID OS VARIANT.

	Samsung	Xiaomi	Realme	Huawei	LineageOS	/e/OS	Google
<i>Long-lived Device Identifiers</i>	IMEIs, hardware serial numbers	IMEIs, Secure DeviceID, MD5 hash of Wifi MAC address	IMEI, deviceID, guid	hardware serial number, device RSA cert	-	-	IMEI, hardware serial number, Wifi MAC address
<i>Resettable Identifiers Relinkable to Device</i>	Samsung Consumer ID, Firebase IDs	VAID, Google Ad ID	VAID, OAID, device_id, registrationId, Google Ad ID, Firebase IDs	-	-	-	AndroidID, Google Ad ID
<i>Third-Party System App Data Collectors</i>	Google, Mobile Operator, Microsoft, LinkedIn, Hiya	Google, Mobile Operator, Facebook	Google, Heytap	Google, Daily Motion, Avast, Qihoo 360, Microsoft	Google	-	
<i>Main Telemetry Collectors (By Data Volume)</i>	Google, Samsung, Microsoft	Google, Xiaomi	Google, Heytap	Google, Microsoft	Google	-	
<i>Loggers of App Usage Over Time</i>	Samsung	Google, Xiaomi	-	Google, Microsoft	-	-	
<i>Loggers of Apps Installed On Handset</i>	Google, Samsung	Google, Xiaomi	Google, Realme, Heytap	Google, Huawei	Google	-	

by a mobile OS is neither necessary nor essential, but rather a choice made by the OS developer. Although occasional data transmission to the OS developer to check for updates, etc. is to be expected, as we will see the observed data transmission by the Samsung, Xiaomi, Huawei, Realme and LineageOS Android variants goes well beyond this.

Table I summarises the data collected by each of the Android OS variants studied.

Re-linkability of advertising identifiers. Samsung, Xiaomi, Realme and Google all collect long-lived device identifiers, e.g. the hardware serial number, as well as user-resettable identifiers, such as advertising IDs. By analysing the identifiers sent together in connections, we find that a long-lived device identifier is sent alongside the resettable identifier on these handsets. This means that when a user resets an identifier the new identifier value can be trivially re-linked back to the same device. This largely undermines the use of user-resettable advertising identifiers. See the second row of Table I for a list of resettable identifiers that can be re-linked to the handset in this way.

Data ecosystem. We also find that typically multiple parties collect data from each handset and that considerable potential exists for cross-linking of data collected by these different parties. On every handset, apart from the /e/OS handset, Google collects a large volume of data. On the Samsung handset the Google Advertising ID is sent to Samsung servers, a number of Samsung system apps use Google Analytics to collect data and the Microsoft OneDrive system app uses Google's push service. On the Huawei handset the Microsoft Swiftkey keyboard sends the Google Advertising ID to Microsoft servers. On the Xiaomi handset the Google Advertising ID is sent to Xiaomi servers, while on the Realme handset the Google Advertising ID is sent to Heytap (who partner with Realme/Oppo to provide handset services, so linkage of data collected by Heytap and Realme is also possible).

Recording of user interactions with handset. System apps on several handsets upload details of user interactions with the apps on the handset (what apps are used and when, what app screens are viewed, when and for how long). The effect is analogous to the use of cookies to track users across web sites. On the Xiaomi handset the system app com.miui.analytics uploads a time history of the app windows viewed by the handset user to Xiaomi servers. This reveals detailed information on user handset usage over time, e.g. timing and duration of phone calls. Similarly, on the Huawei handset the Microsoft Swiftkey keyboard (the default system keyboard) logs when the keyboard is used within an app, uploading to Microsoft servers a history of app usage over time. Again, this is revealing of user handset usage over time e.g. writing of texts, use of the search bar, searching for contacts. Several Samsung system apps use Google Analytics to log user interactions (windows viewed etc). On the Xiaomi and Huawei handsets the Google messaging app (the system app used to send and receive SMS texts) logs user interactions, including when an SMS text is sent. In addition, with the notable exception of the /e/OS handset, Google Play Services and the Google Play store upload large volumes of data from all of the handsets (at least 10× that uploaded by the mobile OS developer). This has also been observed in other recent studies [6], which also note the opaque nature of this data collection.

Details of installed apps. Samsung, Xiaomi, Realme, Huawei, Heytap and Google collect details of the apps installed on a handset. Although less worrisome than tracking of user interactions with apps, the list of installed apps is potentially sensitive information since it can reveal user interests and traits, e.g. a muslim prayer app, an app for a gay magazine, a mental health app, a political news app. It also may well be unique to one handset, or a small number of handsets, and so act as a device fingerprint (especially

when combined with device hardware/system configuration data, which is also widely collected). See, for example, [9], [10] for recent analyses of such privacy risks and we note that in light of such concerns, Google recently introduced restrictions on Play Store apps collection of this type of data⁴, but such restrictions do not apply to system apps since these are not installed via the Google Play store.

No opt-out. As already noted, this data collection occurs even though privacy settings are enabled. Handset users therefore have no easy opt out from this data collection.

Where Data Is Sent. On most handsets data appears to be sent to servers located within Europe. A notable exception is the Xiaomi handset which sends data from Europe to servers estimated to be located in Singapore⁵. The Samsung handset also sends data to server capi.samsungcloud.com which appears to be located in the US.

In summary, we find that *e*/OS collects essentially no data and in that sense is by far the most private of the Android OS variants studied. On all of the other handsets the Google Play Services and Google Play store system apps send a considerable volume of data to Google, the content of which is unclear, not publicly documented and Google confirm there is no opt out from this data collection. LineageOS collects no data beyond this data collected by Google and so is perhaps the next most private choice after *e*/OS. We observe the Realme handset collecting device data, including details of installed apps, but nothing more. The Samsung, Xiaomi and Huawei handsets collect details of user interactions with the handset, in addition to device/app data. Of these, Xiaomi collects the most extensive data on user interactions, including the timing and duration of every app window viewed by a user. On the Huawei handset it is the Microsoft Swiftkey keyboard that collects details of user handset interactions with apps, Huawei themselves are only observed to collect device/app data. We observe Samsung collecting data on user interaction with their own system apps, but not more generally.

A. Ethical Disclosure

The mobile OS's studied here are in active use by many millions of people. We informed Samsung, Xiaomi, Huawei, Realme, Microsoft/SwiftKey and Google of our findings and delayed publication to allow them to respond. Huawei and Google responded with some clarifications, which we have included.

II. THREAT MODEL: WHAT DO WE MEAN BY PRIVACY?

The transmission of user data from mobile handsets to back-end servers is not intrinsically a breach of privacy. For instance, it can be useful to share details of the device model/version and the locale/country of the device when checking for software updates. This poses few privacy risks if the data is common to many handsets and therefore cannot be easily linked back to a specific handset/person [11], [12].

⁴<https://thehackernews.com/2021/04/google-limits-which-apps-can-access.html>

⁵Including tracking.intl.miui.com, api.ad.intl.xiaomi.com, data.mistat.intl.xiaomi.com. Server location estimated from IP address using the <https://ipinfo.io/> service, and verified using ping times/trace route.

Two major issues in handset privacy are (i) release of sensitive data, and (ii) handset deanonymisation i.e. linking of the handset to a person's real world identity.

Release of sensitive data. What counts as sensitive data is a moving target, but it is becoming increasingly clear that data can be used in surprising ways and that so-called metadata can be sensitive data. One example of potentially sensitive metadata is the name, timing and duration of the app windows viewed by a user. This can be used to discover the time and duration of phone calls, when texts/messages are sent and received, when a prayer or dating app is used, and so on. More generally, such data reveals what apps a user spends most time viewing and which windows within the app they look at most. Another example is the list of apps installed on a handset. This can reveal user interests and traits [9], [10]. The list of apps can also act as a handset fingerprint, unique to only a small number of handsets, and so be used for tracking.

Data which is not sensitive in isolation can become sensitive when combined with other data, see for example [13], [14], [15]. This is not a hypothetical concern since large vendors including Google, Samsung, Huawei, and Xiaomi operate mobile payment services and supply custom web browsers with the handsets they commercialize.

It is important to be note, however, that the transmission of user data from mobile handsets to back-end servers is not intrinsically a breach of privacy. For instance, it can be useful to share details of the device model/version and the locale/country of the device when checking for software updates. This poses few privacy risks if the data is common to many handsets and therefore cannot be easily linked back to a specific handset/person [11], [12].

The key requirement for privacy is that the data is common to many handsets. Risk factors therefore include whether data is tagged with identifiers that can be used to link different data together and to link it to a specific handset or person. Tagging data with the handset hardware serial number immediately links it to a single handset. Other long-lived device identifiers include the IMEI (the unique serial number of a SIM slot in a handset) and the SIM IMSI (which uniquely identifies a SIM on the mobile network). To mitigate such risks, Google provides a Google Advertising ID that a user can reset to a new value. The idea is that data tagged with the new value cannot be linked to data tagged with the old value, and so resetting the identifier creates a break with the past. However, this is undermined if the new and old values can both be tied back to the same device and so linked together. It is worth noting that there already exist commercial services that given a Google Advertising ID offer to supply the name, address, email etc of the person using the handset⁶.

Deanonymisation. Android handsets can be directly tied to a person's identity in at least two ways, even when a user takes active steps to try to preserve their privacy. Firstly, via the SIM. When a person has a contract with a mobile operator then the SIM is tied to that contract and so to the person. In addition, several countries require presentation of photo ID to buy a SIM. Secondly, via the app store used. On Android handsets

⁶<https://www.vice.com/en/article/epnmvz/industry-unmasks-at-scale-maid-to-pii>, accessed 18th Aug 2021.

the Google Play store is the main way that people install apps. Use of the Google Play store requires login using a Google account, which links the handset to that account since Google collect device identifiers such as the hardware serial number and IMEI along with the account details [6], [16].

A handset can also become linked to a person's identity when data is collected that allows their identity to be inferred or guessed with high probability. One way that this might happen is via a handset's location time history. Many studies have shown that location data linked over time can be used to de-anonymize users, see e.g. [17], [18] and later studies. This is unsurprising since, for example, knowledge of the work and home locations of a user can be inferred from such location data (based on where the user mostly spends time during the day and evening), and when combined with other data this information can quickly become quite revealing [18]. It is worth noting that every time a handset connects with a back-end server, it necessarily reveals its IP address, which acts as a rough proxy for user location via existing geoIP services. With this in mind, the frequency with which connections are made becomes relevant, e.g. observing an IP address/proxy location once a day has much less potential to be revealing than observing one every few minutes.

III. THE CHALLENGES OF SEEING WHAT DATA IS SENT

It is generally straightforward to observe packets sent from a mobile handset. Specifically, we configure the handsets studied to use a WiFi connection to a controlled access point, on which we use `tcpdump` to capture outgoing traffic. However, this is of little use for privacy analysis because (i) packet payloads are almost always encrypted – not just due to the widespread use of HTTPS to transfer data but, as we will see, also because the message data is often further encrypted by the sender using a cipher that may not be explicitly specified through meta-data, particularly when the data may be sensitive (end-to-end encryption); (ii) prior to message encryption, data is often encoded in a binary format for which there is little or no public documentation; and (iii) for proper attribution, we need to be able link a message to the sending process/app on the handset.

A. Reverse Engineering

A fairly substantial amount of non-trivial reverse engineering is generally required in order to decrypt messages and to at least partially decode the binary plaintext.

1) *Handset Rooting*: The first step is to gain a shell on the handset with elevated privileges, i.e. in the case of Android to root the handset. This allows us then to (i) obtain copies of the system apps and their data, (ii) use a debugger to instrument and modify running apps (e.g. to extract encryption keys from memory and bypass security checks), and (iii) install a trusted SSL root certificate to allow HTTPS decryption, as we explain below. Rooting typically requires unlocking the bootloader to facilitate access to the so-called fastboot mode, disabling boot image verification and patching the system image. Unlocking the bootloader is often the hardest of these steps, since many handset manufacturers discourage bootloader unlocking. Some, such as Oppo, go so far as

to entirely remove fastboot mode (the relevant code is not compiled into the bootloader). The importance of this is that it effectively places a constraint on the handset manufacturers/mobile OSes that we can analyse. Xiaomi and Realme provide special tools to unlock the bootloader, with Xiaomi requiring registering user details and waiting a week before unlocking. Huawei require a handset-specific unlock code, but no longer supply such codes. To unlock the bootloader on the Huawei handset studied here, we needed to open the case and short the test point pads on the circuit board, in order to boot the device into the Huawei equivalent of Qualcomm's Emergency Download (EDL) mode. In EDL mode, the bootloader itself can be patched to reset the unlock code to a known value (we used a commercial service for this), and thereby enable unlocking of the bootloader.

2) *Decompiling and Instrumentation*: On a rooted handset, the Android application packages (APKs) of the apps on the `/system` disk partition can be extracted, unzipped and decompiled. While the bytecode of Android Java apps can be readily decompiled, the code is almost always deliberately obfuscated in order to deter reverse engineering. As a result, reverse engineering the encryption and binary encoding in an app can feel a little like exploring a darkened maze. Perhaps unsurprisingly, this is frequently a time-consuming process, even for experienced researchers/practitioners. It is often very helpful to connect to a running system app using a debugger, so as to view variable values, extract encryption keys from memory, etc. On most of the handsets studied we used `Frida`⁷ to provide a convenient debug interface, allowing dynamic hooking of running code to extract variable values, overwrite function return values and indeed replace the implementation of whole functions. However, on the Huawei handset studied, this approach is not possible since a protected memory model appears to be used, which causes an app to crash when a debugger attaches to it. The protected memory model is likely a write-rarely one – essentially the memory can be modified during the initial startup of an app, but not thereafter [19]. To work around this, we used the fact that on Android all Java apps are cloned/forked from a single Zygote process that is started early after the system boots. We used `Riru`⁸ to modify the Zygote process to allow code injection, and `edXposed`⁹ to provide an interface to Riru that loads user specified code. Riru works by replacing a dynamic library loaded by Zygote, and since this occurs at Zygote startup, it is compatible with the Huawei protected memory model. Once Zygote is modified, the changes propagate to all apps, since they run in clones of the Zygote process, and so all apps can be instrumented/modified. This is less convenient than Frida since changes require a reboot plus Java Native Interface (JNI) C code cannot be instrumented.

3) *Decrypting Data*: A number of system apps on the Xiaomi, Realme and Huawei handsets first encrypt data, generally using either AES/ECB or AES/CBC, before transmitting it over an SSL connection. In more detail:

⁷<https://frida.re/>

⁸<https://github.com/RikkaApps/Riru>

⁹<https://github.com/ElderDrivers/EdXposed>

i) Xiaomi. The app `com.miui.analytics` sends extensive telemetry to the server `tracking.intl.miui.com`. The data sent is AES/ECB encrypted. The key exchange protocol between handset and server involves the handset generating a random 128-bit AES key, encrypting this using an RSA public key and transmitting it base64 encoded to the server specified in `/track/key_get` endpoint. The server responds by sending a second AES key encrypted using the first, together with a SID value that is sent along with later encrypted messages to identify the key used for encryption. The handset decrypts the received key, generates an RSA private/public key pair in the handset Secure Element, and uses the public key to encrypt the AES key before storing it on disk as a SharedPreference data entry. Since the RSA private key is held within the secure element, it is only accessible to the app. This approach means that the AES key is never unencrypted at rest and so it is necessary to extract the key from the memory of the running app. We do this using Frida to intercept the entry points to the various functions used to carry out AES encryption and record the key as it is passed in. A similar key exchange protocol is used by other Xiaomi system apps. In particular, the app `com.miui.msa.global` sends encrypted data to the server `api.ad.intl.xiaomi.com` which appears to be associated with ad management. A number of user-facing system apps, e.g. the file manager `com.mi.android.globalFileexplorer`, the Settings app `com.xiaomi.misettings` and the Security Center app `com.miui.securitycenter`, use a similar approach to encrypt data sent to `data.mistat.intl.xiaomi.com`. Since the user agent header value is the same for all of these apps, to determine the app associated with a connection to `data.mistat.intl.xiaomi.com` (so that we can extract the AES key from its memory) we monitor the handset TCP sockets in `/proc`.

ii) Realme. The app `com.heytag.mcs`, which appears to implement the main Heytag services on the Realme handset, encrypts data with AES/CBC before sending it to `deuex.push.heytagmobile.com`. The 128-bit AES key and IV are hard-coded in the app and so can be readily extracted and used to decrypt the data sent. The plaintext is encoded as a protobuf. Messages sent to `ifrus-eu.coloros.com` by app `com.nearme.romupdate` are AES/CTR encrypted base-64 encoded JSON. A token that helps reconstruct the AES key using a custom encoding scheme is appended to the end of the base-64 message. Using this, the message can be decrypted.

iii) Huawei. Data sent to `query.hicloud.com` by app `com.huawei.android.hwouc` has an `extra_info` field with encrypted information. The `extra_info` field consists of three sections, the first is AES encrypted by a custom obfuscated JNI C library, the second section is AES encrypted in Java, and the third section is the AES key encrypted using an RSA public key. Since we do not have access to the RSA private key, we cannot decrypt this third section to obtain the AES key. Instead, we use Riru/edXposed to extract the key from the memory of the running app and then use it to decrypt the data in the second section. The C code that encrypts the first section uses AES encryption, but the key is generated by heavily obfuscated code (symbol names in the code appear to refer to so-called white-box cryptography, i.e. where the crypto algorithm remains secure even when

the software implementation can be inspected). Due to the protected memory implementation on the Huawei handset, we cannot instrument this C code (Riru/edXposed can only be used with Java code). Instead, we use Riru/edXposed to extract the plaintext data sent into the JNI library by the Java app. The `com.huawei.systemmanager` contains embedded SDKs: `com.avast.android.sdk` from Avast plus `com.qihoo.cleandroid.sdk` and other SDKs from Qihoo 360. These encrypt the data sent, respectively, to `avast.com` and `360safe.com`. The Avast SDK uses 128-bit AES/CBC encryption and a key exchange protocol with rotating keys. To decrypt the data, we used Riru/edXposed to extract the AES key and IV from the app memory – since the keys frequently rotate, we do this on an ongoing basis and dump the keys to the handset log where they can be viewed using `logcat`. The plaintext is a binary encoded protobuf. The Qihoo 360 SDK periodically (every 1-2 days) sends data to `mvconf.cloud.360safe.com/safeupdate` and `mclean.cloud.360safe.com/CleanQuery`. The data is sent in a custom binary data format with the payload encrypted using a JNI C library. To decrypt the data we therefore extracted the plaintext from the app memory using Riru/edXposed.

It goes without saying that the reverse engineering involved was time consuming and required quite some persistence.

4) Decoding Data: Sometimes the plaintext data (i.e. after decryption, if needed) is human-readable, e.g. json. However, frequently it is encoded, often with multiple nested encodings. Common encodings that are straightforward to detect and decode include: JWT tokens¹⁰, base64, hexstring and URL encoding of binary data, gzipping. More complex data is often binary encoded in the Google Protobuf serialization format¹¹. Protobuf's can be decoded without knowledge of the scheme, although this means that field names are missing and there is sometimes with ambiguity as to interpretation of field types. We used the Google Protobuf compiler for this, with the `--decode_raw` option when a protobuf schema was unavailable. Google apps often encode data in a Protobuf array format, namely as a sequence of `length/varint;protobuf;` entries, from which the individual Protobufs need to be extracted and decoded. For Firebase Analytics we manually reconstructed the protobuf schema from the decompiled Firebase code. Other encoding formats that we less commonly observed include Snappy¹², Avro¹³, Bond¹⁴ and also some proprietary formats. In particular, the Microsoft Swiftkey system app sends telemetry data encoded in gzipped Avro serialisation format. Unlike protobufs, Avro cannot be decoded without knowledge of the schema used for encoding. We therefore extracted the schema from the app by executing a `getSchema()` call on app startup (by dynamically patching the app using edXposed) and then dumping the large (about 200KB) json response to disk. The Microsoft OneDrive system app sends telemetry data encoded in Microsoft's Bond Compact Binary format. Again the schema is needed to decode Bond data. Bond works by compiling the schema to Java code, and so we

¹⁰<https://jwt.io>

¹¹<https://developers.google.com/protocol-buffers/>

¹²<https://google.github.io/snappy/>

¹³<https://avro.apache.org/>

¹⁴<https://github.com/microsoft/bond>

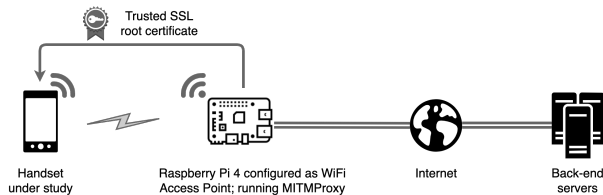


Fig. 1. Measurement setup. Mobile handset configured to access the Internet using a WiFi access point hosted on a Raspberry Pi. A system certificate is installed on the phone to be able to decrypt outgoing traffic. The laptop pretends to any process running on the handset to be the destination server, creates a connection to the actual target, and relays requests and their replies between handset and server while logging the traffic.

decompiled the app, manually reconstructed the schema from the decompiled code and then compiled a C++ programme based on the reconstructed schema using Microsoft’s Bond compiler to yield a decoder that can deserialise the observed POST payload data, then re-serialise to json so that its human readable. The Qihoo 360 SDK uses a proprietary binary format that we reconstructed by decompiling the SDK and inspecting the code.

Once decoded, known values such as the handset IMEI, hardware serial number, Google Advertising Id can often be readily identified. Otherwise, we manually examined the decompiled app to find the code that writes each value and so establish how the value is generated. This is necessary, for example, to identify values that are hashes of device identifiers.

B. Decrypting HTTPS Connections

Almost all of the data we observe is sent over HTTPS connections and so encrypted using TLS/SSL (in addition to any other encryption used by the app). However, decrypting SSL connections is relatively straightforward. We route handset traffic via a WiFi access point (AP) that we control, configure this AP to use `mitmdump` as a proxy [20] and adjust the firewall settings to redirect all WiFi HTTP/HTTPS traffic to `mitmdump` so that the proxying is transparent to the handset. When a process running on the handset starts a new network connection, the `mitmdump` proxy pretends to be the destination server and presents a fake certificate for the target server. This allows `mitmdump` to decrypt the traffic. It then creates an onward connection to the actual target server and acts as an intermediary, relaying requests and their replies between the app and the target server while logging the traffic. The setup is illustrated schematically in Figure 1.

System processes typically carry out checks on the authenticity of server certificates received when starting a new connection and abort the connection when these checks fail. Installing the `mitmproxy` CA cert as a trusted certificate causes these checks to pass, except on the Huawei handset. Installing a trusted cert is slightly complicated in Android 10, since the system disk partition, on which trusted certs are stored, is read-only and security measures prevent it being mounted as read-write. Fortunately, folders within the system disk partition can be overridden by creating a new mount point corresponding to the folder, and in this way the `mitmdump` CA cert can be added to the `/system/etc/security/cacerts` folder. On the Huawei handset each system app contains embedded server

certificate SHA256 hashed and when starting an HTTPS connection checks that the certificate offered by the server matches one of these hashes. It is thus necessary to bypass these checks on each app individually (installing a system-wide trusted cert is not enough). We used `Riru/edxposed` for this.

IV. EXPERIMENTAL SETUP

A. Hardware and Software Used

Mobile handsets: (i) Samsung Galaxy S9 (model SM-G960F)/Android 10 (build QP1A.190711.020, One UI v2.0), (ii) Xiaomi Redmi Note 9 (model M2003J15SG)/Android 10 (build QP1A.190711.020, MIUI Global 12.0.7 QJOMIXM), (iii) Realme 6 Pro (model RMX2063)/Android 10 (build RMX2063_11_A.38, realme UI v1.0), (iv) Huawei P10 Lite (model MAR-LX1B)/Android 9¹⁵ (build 9.1.0.372, EMUI 9.1.0), (v) Google Pixel 2/Android 10 (LineageOS build 17.1-20210316, opengapps 10.0-nano-20210314), (vi) Google Pixel 2/Android 10 (eos build e-0.11-q-20200917). Rooted using `Magisk v20.4` and `Magisk Manager v7.5.1`.

WiFi access point: Raspberry Pi 4 Model B Rev 1.2/Raspbian GNU Linux 11/Mitmproxy 6.0.2 with iptables firewall configured to redirect HTTP/S traffic to port 8080 (on which `mitmproxy` listens) and also to block UDP traffic on HTTPS port 443 (so as to force any Google QUIC traffic to fall back to using TCP since we have no tools for decrypting QUIC).

B. Device Settings

At the start of each test we removed any SIM card and carried out a hard factory reset of the handset, i.e. we used `TWRP` to manually wipe the data partition, thereby forcibly removing all user data and settings, all user installed apps and resetting any disk encryption. Note that we observed that simply clicking on the “factory reset” option in the UI sometimes did not fully remove user data and settings.

Following this factory reset, the handset reboots to a welcome screen and the user is then typically asked to agree to terms and conditions, and presented with a number of option screens. We note that all of the option toggle switches default to the opt-in choice, and so it is necessary for the user to actively select to opt-out. To mimic a privacy conscious user, we unchecked any of the options that asked to share data and only agreed to mandatory terms and conditions. *Samsung:* we unchecked the Sending of Diagnostic Data, Information Linking, Receipt of Marketing Information components of the terms and conditions, skipped the Protect Your Phone screen, did not sign into a Samsung account (which raises a warning that it disables Samsung Cloud, Bixby, Galaxy Themes, Find My Mobile, Samsung Pass, Galaxy Store, Secure Folder). *Xiaomi:* we unchecked the Location, Send Diagnostic Data Automatically, Automatic System Updates, Personalised Ads, User Experience Programme options. *Realme:* we unchecked the User Experience Programme and Uploading Device Error Log Data components of the terms of service, unchecked the WiFi Assistant and Auto-update Overnight options. *Huawei:* we selected No Thanks on the Enhanced Services screen, Later on the User Experience Improvement Programme screen,

¹⁵Following US trade sanctions against Huawei, Android 9 is the latest version of Android available on a Huawei handset that we could root.

Update Manually on the Keep Your Software Up To Date screen. *LineageOS*: we unchecked the Help Improve LineageOS, Location Services options. */e/OS*: we unchecked the Location Services option, skipped Fingerprint Setup, Protect You Phone and */e/* account setup. All of the mobile OSes, apart from */e/OS*, also displayed a Google services screen on first startup. On this we unchecked the Use Location, Allow Scanning, Send Usage and Diagnostic Data options, and we did not log in to a Google user account.

During this startup process, we left WiFi disabled and since no SIM was inserted, there was also no cellular data connection. This allowed us to install the mitmproxy CA cert, and on the Huawei handset Riru/edXposed modules to disable HTTPS cert checks by individual system apps, before the handset made any network connections. WiFi access was then enabled after these steps were completed.

C. Test Design

We seek to define simple experiments that can be applied uniformly to the handsets studied (so allowing direct comparisons) and that generate reproducible behaviour. Mobile OS developers commonly provide add-on services that can be used in conjunction with their handsets, e.g. Samsung offer Cloud storage, Bixby, the Samsung Store; Huawei offer Cloud storage, the AppGallery store; Xiaomi offer Xiaomi Cloud, Mi Coin and Credit. Here we try to keep these two aspects separate and to focus on the handset as a device in itself, separate from optional services such as these. We also assume a privacy-conscious but busy/non-technical user, who when asked, does not select options that share data but otherwise leaves handset settings at their default values.¹⁶

On Android the Settings app must be used to e.g. enable location and WiFi. Since use of the Settings app is not optional for handset users, we include them in our tests. In addition, while on Android apps may be sideloaded over adb, all of the handsets provided include the Google Play store and for most users this is the primary way to install apps. Other than on */e/OS*, use of the Google Play store requires the user to sign in to a Google account and so disclose their email address and perhaps other personal details. We therefore also include opening of the handset Google Play store app and login to a Google account in our tests.

With these considerations in mind, for each handset we carry out the following experiments:

- 1) Start the handset following a factory reset (mimicking a user receiving a new phone), recording the network activity.
- 2) Insert a SIM, recording the network activity.
- 3) Following startup, leave the handset untouched for several days (with power cable connected) and record the network activity. This allows us to measure the connections made

¹⁶There is also an important practical dimension to this assumption. Namely, each handset has a wide variety of settings that can be adjusted by a user and the settings on each handset are generally not directly comparable. Exploring all combinations of settings between a pair of handsets is therefore impractical. A further reason is that the subset of settings that a user is explicitly asked to select between (typically during first startup of the handset) reflects the design choices of the handset developer, presumably arrived at after careful consideration and weighing of alternatives. Note that use of non-standard option settings may also expose the handset to fingerprinting.

when the handset is sitting idle. This test is repeated with the user being logged in and logged out, and with location enabled/disabled.

4) Open the pre-installed Google Play app and log in to a user account, recording the network activity. Then log out and close the app store app.

5) Open the settings app and view every option but leave the settings unchanged, recording the network activity. Then close the app.

6) Open the settings app and enable location, then disable. Record the network activity.

7) Make and receive a phone call, send and receive a text. Record the network activity.

D. Additional Material: Connection Data

The content of connections is summarised and annotated in the additional material available anonymously at https://www.dropbox.com/s/b137n94i9rpp177/additional_material_neversleepingears.pdf.

V. RESULTS

As already noted, Table I gives an overview of the data collection observed on the handsets studied. It is helpful to consider this in light of four basic questions: (i) who is collecting data, (ii) what sort of data is being collected, (iii) can resettable identifiers be relinked to the device, (iv) what is the potential for cross-linking of data collected by different parties.

A. Who Is Collecting Data?

1) *Mobile OS Developers*: We observe that Samsung, Xiaomi, Realme and Huawei all collect data from user handsets, despite the user having opted out of data collection/telemetry/analytics and making no use of services offered by these companies. This data is tagged with long-lived identifiers that tie it to the physical device, including across factory resets.

In contrast, *LineageOS* and */e/OS* were not observed to collect handset data. The latter is notable because a case might be made for the necessity of mobile OS operators collecting handset data in order to monitor software operation and catch problems early (i.e. devops). However, it is hard to justify the necessity of such data collection, i.e. that users should have no opt-out, when two mobile OSes adopt an opt-in approach. It is also worth noting that it can be hard to distinguish between diagnostics for existing software and beta testing (or A/B testing) for new or updated software/features. Traditionally, beta testing has always been opt-in. Finally, it is worth noting that it is hard to see why data collection for diagnostics cannot be carried out in a fully anonymous manner, without any use of long-lived identifiers.

2) *Pre-installed Third-Party System Apps*: System apps are pre-installed on the `/system` partition of the handset disk. Since this partition is read-only, these apps cannot be removed. They are also privileged in the sense that they can be assigned permissions without needing user consent, be silently started, etc. The Settings app is, for example, a system app. All of the mobile OSes studied, apart from */e/OS*, have pre-installed Google system apps. We discuss these further below, but first we consider pre-installed system apps from other companies.

The Samsung handset studied also contains pre-installed system apps from Microsoft that send handset telemetry data to `mobile.pipe.aria.microsoft.com`, `app.adjust.com` (a third-party analytics company¹⁷) and use Firebase push messaging. A LinkedIn (now owned by Microsoft) system app also sends telemetry to `www.linkedin.com/li/track`. This third-party data collection occurs despite no Microsoft/LinkedIn apps were ever opened on the device, and no popup or request to send data was observed.

The Samsung and Xiaomi handsets studied also contain pre-installed system apps from mobile operators (SFR/Altice in France, Deutsch Telekom in Germany), which were observed to send telemetry. Note that our handsets were bought second-hand on the Internet and a more controlled study of operator installed system apps may well be warranted. As well as sending telemetry directly, the SFR/Altice app on the Samsung handset also uses Google Analytics to log events.

The Realme handset studied contains pre-installed system apps from Heytap, a Singapore-based private company. It appears that Realme partners with Heytap, who provide account management, cloud data, an app store, etc.

Huawei also appear to partner with a number of third parties to provide handset system services. The Huawei handset studied contains a pre-installed `com.huawei.systemmanager` app which has embedded within it components from third-party scanning/anti-virus services Avast (based in the Czech Republic) and Qihoo 360 (based in China). App data is sent to `avast.com` when an app is installed on the handset. Periodic connections are also observed to `360safe.com` (associated with Qihoo 360) that send device data. The `com.huawei.himovie` overseas system app sends handset data to servers associated with Dailymotion, even though no video app was ever opened on the handset (perhaps these connections prefetch news/topical videos). The Microsoft Swiftkey keyboard app `com.touchtype.swiftkey` is pre-installed on the Huawei handset and sends crash data to `in.appcenter.ms/logs` and telemetry data to `telemetry.api.swiftkey.com`.

In addition to mobile operator system app sharing data on the Xiaomi handset, a pre-installed Facebook app collects data.

Apart from Google’s GApps, no third-party system apps on the LineageOS handset were observed to perform data collection. On /e/OS, we observed no data collection by system apps.

3) *Google System Apps (GApps)*: The Samsung, Xiaomi, Realme and Huawei handsets studied all have pre-installed Google system apps, the so-called GApps package. These include Google Play Services,¹⁸ Google Play Store, YouTube, Gmail, Maps, Drive, Wallet, Chrome. On LineageOS it is necessary to install GApps to use the Google Play store, but this is not necessary with /e/OS (which uses the open-source MicroG re-implementation of Google Play Services and the Google Play app). It is known that Google Play Services and

¹⁷Their website says “Adjust offers a number of analytics tools designed to give you the deepest insight into your user interaction, your marketing channels, and your campaign performance”.

¹⁸Google Play Services provides the API for Google Firebase services such as Google Analytics and Crashlytics to apps on the handset, but also performs device logging/telemetry on behalf of Google.

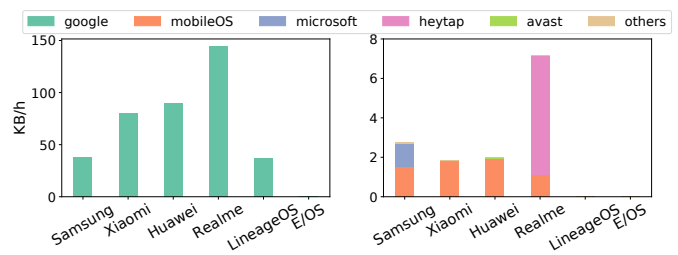


Fig. 2. The average volume of the network traffic generated on each handset by each data collector.

the Google Play store send large volumes of handset data to Google and collect long-lived device identifiers, although until recently there has been a notable lack of measurement studies (see [6], [16]). Other Google apps such as YouTube and Gmail also send handset data and telemetry to Google.

It is worth noting that the volume of data uploaded by Google is considerably larger than the volume of data uploaded to other parties. For example, Figure 2 shows the average rate at which data is uploaded from each handset when lying idle, broken down by data source. The volume of data sent to Google is broken out into a separate plot to make it easier to see the volumes sent to other companies.

It can be seen that no data is uploaded to the LineageOS or /e/OS developers. On the Realme handset Heytap uploads around 3-4× more data than Samsung, Xiaomi and Huawei. Realme themselves collect far less data than Heytap, about half of that collected by Samsung, Xiaomi and Huawei. On the Samsung handset the Microsoft system app uploads a similar volume of data as Samsung.

The volume of data uploaded by Google varies across the handsets. It is zero for /e/OS, since it uses the MicroG open source re-implementation of Google GApps. LineageOS and Samsung send similar volumes of data, Xiaomi and Huawei about twice as much and Realme about three times as much. These differences are likely related to different configurations of Google GApps e.g. on LineageOS the so-called nano version of GApps was installed (other options includes micro, mini, full, stock¹⁹). In all cases the volume of data uploaded to Google is at least 10× that uploaded by the mobile OS developer. For Xiaomi, Huawei and Realme the volume rises to around 30×. Recall that this is despite the “usage & diagnostics” option being disabled for Google services on all handsets (and also the diagnostics/analytics options also being disabled for the mobile OS developers, see Section IV-B).

Note however that from a privacy viewpoint it is not the volume of data that is primarily of concern, but rather the contents of that data and the frequency with which it is sent.

B. What Sort Of Data Is Being Collected?

The data that we observe being sent from handsets can be roughly categorised as: (i) device/user identifiers, (ii) device configuration data and (iii) event logging data/telemetry.

1) *Device/User Identifiers*: We observe that most of the connections from a handset are tagged with an identifier of some sort. Single-use identifiers can be used to avoid

¹⁹See <https://github.com/opengapps/opengapps/wiki/Package-Comparison>

duplicate messages being received and session identifiers can be used to link together groups of connections, e.g. when accessing authenticated resources. These types of identifier are ephemeral, i.e. they are short-lived, hard to link to a particular device and so carry little privacy risk. Longer-lived identifiers can also be used, e.g. to maintain state, and so long as the same identifier is shared by many devices, this carries little privacy risk. Google's Safe Browsing service is a good example of such an approach [21].

Unfortunately we observe little use of such privacy-friendly identifiers in our handset measurements. Instead we find that sending persistent identifiers in connections is ubiquitous. Table I lists the main identifiers sent in connections on each handset. Some of these identifiers are long-lived, e.g. the IMEI (which is typically engraved on the SIM slot), hardware serial number and, on Huawei handsets, the device RSA cert [22]. These identifiers persist across factory resets of the device and are effectively permanent and indelible. Others, such as the Google Advertising Id and VAID, are user-resettable either manually or by a factory reset of the phone. But in practice that means they rarely change and act as strong device identifiers. Further, as we discuss in more detail below, most of these resettable identifiers can be relinked back to the device since long-lived identifiers are sent alongside them.

This means that connections from the same handset can generally be easily linked together over time, which has several consequences. One is that data on device and user behaviour is linked over time, with obvious privacy implications. Another is that every time a handset connects with a back-end server it necessarily reveals the handset IP address, which acts as a rough proxy for user location via existing geoIP services. Many studies have shown that location data linked over time can be used to de-anonymise, e.g. see [17], [18] and later studies. This is unsurprising since, for example, knowledge of the work and home locations of a user can be inferred from such location data (based on where the user mostly spends time during the day and evening), and when combined with other data this information can quickly become quite revealing [18].

2) *Device Configuration Data*: Sharing device hardware/system configuration data such as the device model, screen size, operating system version, radio version generally carries little privacy risk since these are common to many devices (e.g. all devices of the same model). Such data is needed when checking for software updates and selecting the right version of an app to install. Samsung, Xiaomi, Realme and Huawei all collect this type device configuration data, as do Google and many third-party system apps.

Additionally, Samsung, Xiaomi, Realme, Huawei, Heytap and Google also collect details of all apps installed on a handset. This is potentially more sensitive information since the set of apps installed is more likely to be unique to one handset, or a small number of handsets, and so act as a device fingerprint (especially when combined with device hardware/system configuration data). It is not clear why this data collection is needed (if just to check for app updates or to scan for malware then that could be carried out anonymously and without revealing the full set of apps installed on a handset).

```
1 POST https://tracking.intl.miui.com/track/v4
2 Headers
3   OT_SID: 1904b90...536c63d4
4   OT_ts: 1627029461128
5   OT_net: WIFI
6   OT_sender: com.miui.analytics
7     "seq": [
8       {
9         "event": 1,
10        "pkg": "com.google.android.dialer",
11        "class": "com.android.incallui.InCallActivity",
12        "ts": 1627028918422,
13        "vn": "67.0.383690429",
14        "stat": "app_start"
15      },
16      {
17        "event": 2,
18        "pkg": "com.google.android.dialer",
19        "class": "com.android.incallui.InCallActivity",
20        "ts": 1627028934973,
21        "vn": "67.0.383690429",
22        "duration": 16551,
23        "stat": "app_end",
24        "app_duration": 16551
25      }
26    ]
```

Fig. 3. Xiaomi telemetry logs the user interaction with the dialer app when receiving a phone call, including the start and end times of the call.

3) *Event Logging Data/Telemetry*: Samsung and Xiaomi both log data that can reveal user interactions occurring on a handset. Third-party system apps by Google and Microsoft also carry our event logging that can reveal user interactions. Heytap, Daily Motion and the mobile operator log events related to operation of their specific app.

Some logging of events is probably reasonable, e.g. to allow early detection of app performance issues (excessive battery drain, slow operation, etc.). But ongoing detailed logging of the activity on a handset, particularly user activity, can quickly become intrusive and a serious privacy concern. The last row of Table I lists the companies carrying out ongoing and frequent telemetry/event logging on each handset.

Notr that this occurs despite the user opting out of diagnostics/analytics collection on the handsets during onboarding following factory reset.

Xiaomi collects extensive event logging data/telemetry. This is mainly sent to tracking.intl.miui.com. The data sent is doubly-encrypted i.e. the data is first AES encrypted and then sent over an encrypted HTTPS connection. After quite some work reverse engineering the AES key management scheme used, we managed to decrypt the data. The data consists of both timestamped individual events and timestamped sequences of events grouped together. The events logged include, for example, every opening and closing of an app window (“activities” in Android parlance) plus the duration a window is open. Since all window events appear to be logged, this can easily reveal detailed information on user handset usage. For example, Figure 3 shows decrypted logging data sent to Xiaomi when a phone call is received. The dialer app opens its InCallActivity window when the call arrives and closes it when the call ends. Timestamps of the open and close events, plus the duration, are sent to tracking.intl.miui.com. Xiaomi system apps com.miui.msa.global, com.xiaomi.discover, com.android.thememanager also log events using Google Analytics.

Microsoft's Swiftkey keyboard (used on the Huawei hand-

```

1 {'event': {'metadata':
2 {'installId': 'b'\xe7\x19\xec\xa8KD\xff\xal&E\xa3\x066G\
   xf6['', 'appVersion': '7.8.3.5', 'timestamp':
3   {
4     'utcTimestamp': 1628165014657, 'utcOffsetMins': 0}, '
   vectorClock': {'major': 103, 'minor': 482, 'order': 100}
5   },
6   'application': 'com.google.android.apps.messaging', '
   durationMillis': 6891,
7   'typingStats':
8     {'totalTokensEntered': 0, 'tokensFlowed': 0, '
   tokensPredicted':
9     0, 'tokensCorrected': 0, 'tokensVerbatim': 0, '
   tokensPartial': 0, 'netCharsEntered': 3, 'deletions': 1,
   'characterKeystrokes': 0, 'predictionKeystrokes': 0, '
   remainderKeystrokes': 0, 'predictionSumLength': 0, '
   typingDurationMillis': 837, 'emojisEntered': 0, '
10  totalTokensEnteredEdited': 0, 'tokensFlowedEdited': 0, '
   tokensPredictedEdited': 0, 'tokensCorrectedEdited': 0, '
   tokensVerbatimEdited': 0, 'tokensPartialEdited': 0},
11  'languagesUsed': 0, 'termsPerLanguage': {}, '
   tokensPerSource': {}, 'tokensShownPerSource': {'': 6, '
   en_GB/en_GB.lm': 16, 'user/dynamic.lm': 6},
12  'userHandle': 0
  }}

```

Fig. 4. The Microsoft Swiftkey keyboard logs user interaction with the messaging app when sending a text.

set) also carries out extensive event logging, sending this data to telemetry.api.swiftkey.com. In particular, when the keyboard is used within an app then the app name, number of characters entered and an event timestamp are sent. In this way use, for example, of the searchbar, contacts and messaging apps is logged and so can easily reveal detailed information on user handset usage. See, for example, Figure 4. Interactions with the keyboard, e.g. opening the clipboard, viewing/modifying the settings, are also logged. Information on Swiftkey app crashes, including stack traces, is sent to in.appcenter.ms.

Several Samsung system apps use Google Analytics to log user interaction events, including windows/activities viewed plus duration and timestamp. System apps instrumented in this way include com.wssyncmldm, com.samsung.android.samsungpass, com.samsung.android.authfw, com.samsung.android.bixby.agent, com.samsung.android.game.gamehome, com.sec.android.app.samsungapps. The app api.omc.samsungdm.com logs when a SIM is inserted and samsung-directory.edge.hiyaapi.com logs making/receiving of a phone call.

We did not observe any substantial event logging by Huawei, Realme (including Heytap), LineageOS or /e/OS.

On the Xiaomi and Huawei handsets the Google messaging app com.google.android.apps.messaging uses Google Analytics to log user interaction, including screens/activities viewed plus duration and timestamp, and logs the event that text is sent. In addition, with the notable exception of the /e/OS handset, Google Play Services and the Google Play store collect large volumes of data from all of the handsets (see Figure 2). This has also been observed in other recent studies [6], which also note the opaque nature of this data collection (no documentation, binary encoded payloads, obfuscated code). From our discussions with Google we understand that they plan to publish documentation on this data collection/telemetry, but to date that has not happened.

Other event logging/telemetry that we observed is confined

to operation of specific apps. On the Samsung handset the Microsoft OneDrive app sends data with device details and installed Microsoft apps to mobile.pipe.aria.microsoft.com and app.adjust.com, and uses Firebase push messaging. Events and data related to the mobile operator app com.altice.android.myapps are logged to sun-apps.sfr.com and via Google Analytics (e.g. duration app has been active, errors, stack traces). On the Realme handset events related to app com.heytag.mcs (launch etc) are logged to dceux.push.heytagmobile.com. On the Huawei handset events related to app com.huawei.himovie.overseas are logged to pebed.dmevent.net, and when a new app is installed the app details are sent to a scanning service at apkrep.ff.avast.com²⁰.

C. Can Resettable Identifiers Be Relinked to Device?

In response to privacy concerns, identifiers used to track user behaviour are now often resettable [23]. For example, the Google Advertising Identifier (GAID) can be reset via the Settings app on an Android handset. The idea is that by resetting such an identifier a person effectively unlinks themselves from the data collected about them in the past and starts afresh. However, this aim is largely subverted as the data collected allows relinking of the new identifier to the same physical user/handset. We find that data collection allowing the potential for relinking is commonplace.

Note that we are not in a position to know whether such relinking actually takes place. However, by observing identifiers sent together in the same data connection, we can see whether such relinking could be easily carried out, if desired.

It can be seen from Table I that Samsung, Xiaomi, Realme, Huawei and Google all collect long-lived identifiers from the handset, e.g. the IMEI (which is typically engraved on the SIM slot) or hardware serial number. These identifiers persist across factory resets of the device and are effectively permanent and indelible. If a long-lived identifier is sent in the same connection as a resettable identifier, then relinking of the resettable identifier to the handset is trivial. If one such resettable identifier can be relinked, and is then sent in a connection with other resettable identifiers, then these too can be relinked to the device. Using such an analysis we find that many of the resettable identifiers used by Samsung, Xiaomi, Realme, Huawei and Google can be relinked to the device.

The relevant identifiers are detailed in Table I. Google can potentially relink both the Google AndroidID and Google Advertising Identifier to the device²¹. Xiaomi and Realme can relink the Google Advertising Identifier to the device, as well as all of the other identifiers commonly sent in connections. The same applies to Heytap on Realme handsets. Samsung can relink their Consumer ID, which is sent in many connections to Samsung servers, to the device. Samsung also collect Google Firebase identifiers/authentication tokens (used in conjunction with Google Analytics, etc.) and they can potentially relink

²⁰According to Huawei this can be disabled by opening the Optimiser app, entering the settings screen and unchecking the “Auto-clean junk files” and “online virus scan” option, although we have not verified this.

²¹We note that the Google Play policy <https://support.google.com/googleplay/android-developer/answer/9857753#> prohibits re-linking of advertising identifiers by apps on the Google Play store, and Google have stated to us that internally they also adhere to this policy.

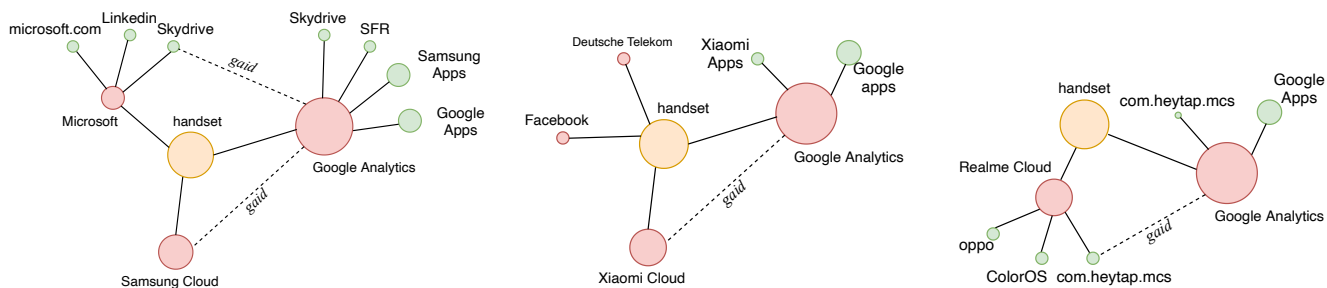


Fig. 5. Potential for cross-linking data collection with different handsets: Samsung (left), Xiaomi (center), Realme (right). Red circles represent data collectors and green circles represent for what specific service instance the data is collected. Observe the potential of cross-linking through the Google Advertising I.

these to the device via Google since the Google AndroidID is sent in Firebase connections and Google can relink the AndroidID to the device. Hence the Google Analytics data collected by Samsung system apps can potentially be relinked to the device. Huawei sends the handset hardware serial number (a long-lived device identifier) in connections, but we observed little use of other identifiers and no potential for relinking of resettable identifiers by Huawei. We also did not observe any potential for relinking on LineageOS and /e/OS.

D. Potential For Cross-linking Data Collection?

We find that typically multiple parties collect data from a handset. For example, on a Samsung handset Samsung, Google and Microsoft/LinkedIn all collect data. That raises the question of whether the data collected separately by these parties can be linked together (and of course combined with data from other sources). While we are not in a position to know whether such linking actually takes place, by inspection of the identifiers jointly collected by the parties we can see whether the potential exists for data linking.

Figure 5 illustrates these potential linkages as a graph.

Samsung record the Google Advertising Id, as do Google and there is therefore immediately potential for Samsung and Google to link their separate data. It is also worth noting that a number of Samsung system apps use Google Analytics to log data. Google already make some of their own data visible to third parties via the Google Analytics dashboard interface, e.g. user demographics, and so limited data sharing from Google to Samsung is likely taking place via that channel.

On the Samsung handset the Microsoft system app sends data to Microsoft servers and to app-adjust.com, and presumably Microsoft have access to the data that their app sends to app-adjust.com. The Google Advertising ID is sent to app-adjust.com, potentially allowing linkage to Google handset data. A LinkedIn system app also collects data. Since Microsoft own LinkedIn they may have access to that data, as well as other data held by LinkedIn.

Xiaomi records the Google Advertising Id, as do Google, and so linking of their data is possible. Xiaomi can display adverts within handset system apps and the UI and so some limited data sharing from Google to Xiaomi may be occurring via the that channel. We also note that a Facebook system app is installed in the Xiaomi handset and the Facebook Ad SDK is embedded in a number of Xiaomi system apps, and so

there appear to be connections between Xiaomi and Facebook although we saw no evidence of sharing of identifiers in our measurements.

On the Realme handset Heytap records the Google Advertising Id as do Google, and so linking of Google and Heytap data is again possible. In its connections Realme sends an identifier supplied by a Heytap server (the registrationId is sent by shortheux.push.heytapmobile.com) and so linkage of data collection by Realme and Heytap is possible, and via Heytap with Google.

On the Huawei handset a hash of the handset android_id is sent to avast.com and a uuid is sent to 360safe.com²² but neither seem easily linked to the hardware serial number sent to Huawei servers. The Swiftkey keyboard sends the Google advertising is to telemetry.api.swiftkey.com, but we did not observe this id being sent to Huawei servers.

VI. RELATED WORK

While the Android ecosystem continues to evolve, most smartphone users remain largely unaware of the personal identifiable information (PII) disclosed by their devices and the apps they run [24]. This has motivated extensive privacy and security over recent years, e.g. see [3], [4] and references therein, and triggered data protection legislation with nearly 100 articles laying out privacy requirements [25].

As nearly a quarter of mobile apps with over 1 billion downloads are known to monetize private data [26], Android privacy analyses have been largely focused on the app ecosystem. Data collection purposes by mobile apps have been classified in [1]. Ren et al. document systematic collection of (PII) over time by different apps and the ability of third-parties to link user activity and locations across apps [2]. Further work examines over 500 apps on the Google Play Store and shows that 76% of them collect and transmit PII insecurely, while 34% of these send PII to third parties [27]. Gamba et al. reveal that the Android open-source model facilitates harmful behaviours and backdoors to sensitive data without user consent, while uncovering potential relationships between manufacturers, network operators and third-parties [28]. Privacy leaks due to misuse of Inter-component Communications (ICC) in Android apps are documented in [29]. With most Android users being based in China, Wang et al. take a look at the degree of domestic

²²According to Huawei this uuid value is changed daily.

mobile app tracking, showing a distinctive mobile tracking market where 10% of users send PII [5].

What information handset operating systems send to their associated back-end servers is not well understood. Probably closest to the present work are recent analyses of the data that web browsers share with their back-end servers [21] and of the data shared by Google Play Services [6], [16]. The latter is motivated in part by the emergence of Covid contact tracing apps based on the Google-Apple Exposure Notification (GAEN) system, which on Android requires that Google Play Services to be enabled. The present study is broader in scope, given that users appear to have no option to disable data collection by the operating system and by the pre-installed system apps. To the best of our knowledge there has been no previous systematic work reporting measurements of the content of messages sent between Android OSes and the associated back-end servers.

VII. CONCLUSIONS

We present an in-depth analysis of the data sent by the Samsung, Xiaomi, Huawei, Realme, LineageOS and /e/OS variants of Android. We find that, with the notable exception of e/OS, even when minimally configured and the handset is idle these vendor-customized Android variants transmit substantial amounts of information to the OS developer and also to third-parties (Google, Microsoft, LinkedIn, Facebook etc) that have pre-installed system apps. While occasional communication with OS servers is to be expected, the observed data transmission goes well beyond this and raises a number of privacy concerns.

REFERENCES

- [1] H. Jin, M. Liu, K. Dodhia, Y. Li, G. Srivastava, M. Fredrikson, Y. Agarwal, and J. I. Hong, "Why are they collecting my data? inferring the purposes of network traffic in mobile apps," *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, vol. 2, no. 4, Dec. 2018.
- [2] J. Ren, M. Lindorfer, D. J. Dubois, A. Rao, D. Choffnes, and N. Vallina-Rodriguez, "Bug fixes, improvements,... and privacy leaks," in *Network and Distributed System Security Symposium (NDSS)*, 2018.
- [3] A. Razaghpanah, R. Nithyanand, N. Vallina-Rodriguez, and S. Sundaresan, "Apps, Trackers, Privacy, and Regulators: A Global Study of the Mobile Tracking Ecosystem," in *Network and Distributed System Security Symposium (NDSS)*, 2018.
- [4] J. Reardon, A. Feal, P. Wijesekera, A. E. B. On, N. Vallina-Rodriguez, and S. Egelman, "50 ways to leak your data: An exploration of apps' circumvention of the android permissions system," in *28th USENIX Security Symposium (USENIX Security 19)*. Santa Clara, CA: USENIX Association, Aug. 2019, pp. 603–620. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity19/presentation/reardon>
- [5] Z. Wang, Z. Li, M. Xue, and G. Tyson, "Exploring the eastern frontier: A first look at mobile app tracking in china," in *Passive and Active Measurement*, A. Sperotto, A. Dainotti, and B. Stiller, Eds., 2020.
- [6] D. J. Leith and S. Farrell, "Contact Tracing App Privacy: What Data Is Shared By Europe's GAEN Contact Tracing Apps," in *Proc IEEE INFOCOM*, 2021.
- [7] Forbidden Stories, "The pegasus project," <https://forbiddenstories.org/case/the-pegasus-project/>, July 2021.
- [8] IDC, "Worldwide quarterly mobile phone tracker," July 2021.
- [9] A. Pham, I. Dacosta, E. Losiouk, J. Stephan, K. Huguenin, and J.-P. Hubaux, "Hidemyapp: Hiding the presence of sensitive apps on android," in *28th USENIX Security Symposium (USENIX Security 19)*. Santa Clara, CA: USENIX Association, Aug. 2019, pp. 711–728. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity19/presentation/pham>
- [10] G. L. Scoccia, I. Kanj, I. Malavolta, and K. Razavi, "Leave my apps alone! a study on how android developers access installed apps on user's device," in *Proceedings of the IEEE/ACM 7th International Conference on Mobile Software Engineering and Systems*, ser. MOBILESoft '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 38–49. [Online]. Available: <https://doi.org/10.1145/3387905.3388594>
- [11] L. Sweeney, "k-anonymity: A model for protecting privacy," *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 10, no. 05, pp. 557–570, 2002.
- [12] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkatasubramanian, "l-diversity: Privacy beyond k-anonymity," *ACM Transactions on Knowledge Discovery from Data (TKDD)*, vol. 1, no. 1, pp. 3–es, 2007.
- [13] M. Cominelli, F. Gringoli, P. Patras, M. Lind, and G. Noubir, "Even black cats cannot stay hidden in the dark: Full-band de-anonymization of bluetooth classic devices," in *2020 IEEE Symposium on Security and Privacy (S&P)*. IEEE, 2020, pp. 534–548.
- [14] A. Di Luzio, A. Mei, and J. Stefa, "Consensus robustness and transaction de-anonymization in the ripple currency exchange system," in *IEEE International Conference on Distributed Computing Systems (ICDCS)*, 2017, pp. 140–150.
- [15] T.-F. Yen, Y. Xie, F. Yu, R. P. Yu, and M. Abadi, "Host fingerprinting and tracking on the web: privacy and security implications," in *Network and Distributed System Security Symposium (NDSS)*, February 2012.
- [16] D. J. Leith, "Mobile Handset Privacy: Measuring The Data iOS and Android Send to Apple And Google," in *Proc SecureComm*, 2021.
- [17] P. Golle and K. Partridge, "On the Anonymity of Home/Work Location Pairs," in *Pervasive Computing*, 2009.
- [18] M. Srivatsa and M. Hicks, "Deanonymizing mobility traces: Using social network as a side-channel," in *ACM Conference on Computer and Communications Security (CCS)*, 2012, pp. 628–637.
- [19] I. Stoppa, "Kernel Hardening: Protecting the Protection Mechanisms," 2018, accessed 31 July 2021. [Online]. Available: <https://events19.linuxfoundation.org/wp-content/uploads/2017/12/Kernel-Hardening-Protecting-the-Protection-Mechanisms-Igor-Stoppa-Huawei.pdf>
- [20] A. Cortesi, M. Hils, T. Kriebchaumer, and contributors, "mitmproxy: A free and open source interactive HTTPS proxy (v5.01)," 2020. [Online]. Available: <https://mitmproxy.org/>
- [21] D. J. Leith, "Web Browser Privacy: What Do Browsers Say When They Phone Home?" *IEEE Access*, 2021.
- [22] "EMUI 11.0 Security Technical White Paper," 2020, accessed 31 July 2021. [Online]. Available: https://consumer-img.huawei.com/content/dam/huawei-cbg-site/common/campaign/privacy/whitepaper/emui_11.0_security_technical_white_paper_v1.0.pdf
- [23] Wired, "A simple way to make it harder for mobile ads to track you," Aug 2019.
- [24] M. Van Kleek, I. Liccardi, R. Binns, J. Zhao, D. J. Weitzner, and N. Shadbolt, "Better the devil you know: Exposing the data sharing practices of smartphone apps," in *CHI Conference on Human Factors in Computing Systems*, 2017, pp. 5208–5220.
- [25] European Parliament and Council of the European Union, "Regulation on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing directive 95/46/ec (data protection directive)," 2016.
- [26] G. Cecere, F. L. Guel, and V. Lefrere, "Economics of free mobile applications: personal data as a monetization strategy," HAL, Post-Print, Sep. 2018. [Online]. Available: <https://ideas.repec.org/p/hal/journal/hal-01988603.html>
- [27] Q. Jia, L. Zhou, H. Li, R. Yang, S. Du, and H. Zhu, "Who leaks my privacy: Towards automatic and association detection with gdpr compliance," in *Wireless Algorithms, Systems, and Applications*, E. S. Biagioni, Y. Zheng, and S. Cheng, Eds., 2019, pp. 137–148.
- [28] J. Gamba, M. Rashed, A. Razaghpanah, J. Tapiador, and N. Vallina-Rodriguez, "An analysis of pre-installed android software," in *IEEE Symposium on Security and Privacy (S&P)*, 2020, pp. 1039–1055.
- [29] D. Zhang, Y. Guo, D. Guo, R. Wang, and G. Yu, "Contextual approach for identifying malicious inter-component privacy leaks in android apps," in *IEEE Symposium on Computers and Communications (ISCC)*, 2017, pp. 228–235.