

# Does Firm Size Influence the Collection of Sensitive Data?: A Study of Child-Orientated Apps

Cecere, G.<sup>\*</sup>, Lefrere, V.<sup>†</sup>, Tucker, C.<sup>‡</sup>

June 10, 2025

## Abstract

How does firm size affect the privacy protections offered to customers? On the one hand, it could be that larger firms use their size to amass more data. On the other hand, smaller firms may be less careful in their data protection practices, because they have a different perception of risk. Using data from the Google Play Store over a three-year period, we explore this empirical question in the U.S. children's app market. Our findings indicate that larger app developers consistently implement stronger privacy protections, requesting less sensitive data compared to smaller developers. These results hold across empirical approaches, including instrumental variables and the propensity-score matching approach. Additionally, our analysis shows that mergers between developers and sudden increases in size of the user-bases of the product are associated with reduced data collection. We show that newly created and updated apps produced by large developers collect less data compared to existing apps. Our findings indicate a trend toward standardized privacy practices across different national regulatory regimes. This research highlights the potential for growth-driven improvements in data privacy practices among app developers, regardless of their regulatory context.

**Keywords:** Data and Competition, Economics of Privacy, Apps for Young Children

JEL CODE: D82, D83, M31, M37

---

<sup>\*</sup>Institut Mines-Telecom, Business School. Email: grazia.cecere@imt-bs.eu.

<sup>†</sup>Institut Mines-Telecom, Business School. Email: vincent.lefrere@imt-bs.eu.

<sup>‡</sup>Massachusetts Institute of Technology (MIT) - Management Science (MS). Email: cetucker@mit.edu.

# 1 Introduction

Many antitrust cases focus on allegations of larger firms collecting excessive data. In 2024, the U.S. Federal Trade Commission (FTC) filed a major lawsuit against Meta (formerly Facebook), challenging its acquisitions of Instagram and WhatsApp. The FTC claims that these acquisitions reduced competition and limited privacy options for consumers, forcing them to share more personal data.<sup>1</sup> In 2019, the German regulator also challenged Facebook, arguing that its market dominance compelled consumers to provide personal data to access its services.<sup>2</sup> However, from an economic perspective, it is unclear whether larger or smaller firms have stronger incentives to collect privacy-intrusive data. On the one hand, larger firms may collect more data on a given subject because they have the size and scale to use data more effectively, and may offer better products that result in being able to request more data from consumers to service them. On the other hand, smaller firms may collect more data as a result of being less cautious about the negative risks of consumer-data collection and believing that they need more data to compete. Since theoretical arguments could go either way, this paper investigates this empirical question of how firm size relates to data collection in a case where privacy protection undoubtedly matters: the data collection of sensitive information from very young children.

Children may not understand the potential negative outcomes of revealing personal information online, so regulation often requires parental consent to collect children’s data (Bleier *et al.*, 2020; Banerjee *et al.*, 2024; Johnson *et al.*, 2024). There are evident reasons to want to safeguard the data of toddlers and preschoolers, and this is also a useful market to study because of the amount of discretion developers have in choosing what data to collect from their users (Kircher and Foerderer, 2024a). Apps that target very young children tend to be simple and provide content based primarily on images and sound. They do not require large swathes of user data to perform better. Furthermore, the simplicity of these apps not only makes them low-cost to develop (Ghose and Han, 2014), but also attracts a high number of

---

<sup>1</sup>For more on the FTC vs. Meta case, see <https://www.theverge.com/2024/11/13/24295637/meta-must-face-ftc-antitrust-trial-instagram-whatsapp>, Last accessed December 2, 2024.

<sup>2</sup>For details on the German case, see [https://www.bundeskartellamt.de/EN/Home/home\\_node.html](https://www.bundeskartellamt.de/EN/Home/home_node.html), Last accessed December 2, 2024.

competing developers internationally.

We collected weekly data on apps available in the U.S. Google Play Store, over the period July 2017 to January 2021. To identify child-targeted apps, we collect apps in the “Designed for Families” category, which helps parents identify child-appropriate content,<sup>3</sup> and we use a keyword search with terms such as “preschool” and “toddler.” Our dataset includes 27,119 apps leading to 1,498,645 observations. These apps are produced by 11,090 developers located in 127 countries. The Online Privacy Protection Act (COPPA) protects the privacy of American children under 13 years of age and defines what is sensitive data in the case of children. We use the COPPA definition of sensitive data to determine whether an app requires sensitive data.

Our analysis explores the link between developer size and the collection of sensitive data. Descriptive evidence shows 40.7% of apps produced by small developers requested at least one type of sensitive data, compared to only 20.4% of the apps produced by larger developers. Empirical evidence shows that child apps produced by larger developers are less likely to collect sensitive data. We use several empirical strategies to demonstrate the robustness of our results. Although our main models include a large set of controls and fixed effects, we address potential identification challenges; we use an instrumental variable approach that exploits variation across developers’ countries in the costs of starting business. This is because unobserved and confounding changes in app characteristics and data collection over time may be correlated with a firm’s size. For instance, large firms are more likely to be compliant, as they can afford the compliance costs. The instrumental variable approach suggests that the effects of developers on sensitive data collection outcomes may be even more significant than what is implied by the correlations observed in our panel regressions.

Establishing a causal relationship between developer size and data collection is challenging. As we have a large panel of data, we exploit the variation in the number of apps by developers, as we observe that a subset of developers increases in size throughout our sample period. First, we consider external developer growth due to mergers and acquisitions. In this setting, we compare the apps produced by developers that merge with developers that do

---

<sup>3</sup>Developers who opt in to the program self-declare that the app complies with Google Play Store’s internal Designed for Families policy and the USA’s COPPA.

not benefit from this external growth. Second, using a propensity-score matching approach, we match developers that increase in size at a given time to developers that did not increase in size. We exploit non-experimental variation due to differences in developer size. We consider that our treatment group is composed of apps produced by developers that show an increase in size at a given point in time, and our control group comprises developers who never increase throughout our sample period. We compare the intensity of data collection before and after across the two groups using a difference-in-differences analysis. We show that apps produced by developers that increase in size are likely to collect sensitive data.

We then present suggestive evidence as to the mechanism. First, we show that large developers decrease their data collection over time when they create new apps and update existing ones, suggesting that there are spillovers in data collection that are consistent with diminishing returns to data collection in the market of children’s apps (Bajari *et al.*, 2019; Farboodi and Veldkamp, 2023; Goldfarb and Que, 2023). Second, we show that data collection is related to the business models of apps (Markovich and Yehezkel, 2024). The results suggest that large developers that rely on advertising business models are less likely to collect sensitive data compared to small developers, suggesting that they are better able to extract value from fewer pieces of data compared to small advertisers. Third, we investigate whether privacy regulation also influences developers’ behaviour (Marthews and Tucker, 2019; Peukert *et al.*, 2022).

The study builds on four main streams of academic literature: privacy regulation, the economics of mobile apps, the relationship between data and market power, and the relation between data and market power. By integrating these perspectives, we aim to investigate the impact of firm size on privacy practices in the context of apps targeted at young children.

The first stream of literature is on privacy regulation. Most of these articles have documented a trade-off between protecting privacy and innovation, in sectors such as health (Miller and Tucker, 2009, 2011, 2017) and advertising (Goldfarb and Tucker, 2011, 2012; Montes *et al.*, 2019; Jia *et al.*, 2020; Johnson *et al.*, 2020; Miller and Skiera, 2024). Several articles have documented distortions in terms of firm location (Rochelandet and Tai, 2016) and creating incentives for firms to collect more data (Adjrid *et al.*, 2015) which may also alter consumer trust Brough *et al.* (2022). Privacy regulation could also play a role in the

quantity and quality of content produced (Lefrere *et al.*, 2024; Miller *et al.*, 2024; Congiu *et al.*, 2022) or consumed (Shen *et al.*, 2024; Yan *et al.*, 2022). Using app data, the article of Cheyre *et al.* (2023) show that after the introduction of privacy preserving policy App Tracking Transparency policy, app developers did not exit the market but instead adjusted their strategy by implementing a more protective privacy framework. Relevant to our theoretical framework is research on platform governance related to curating video content for children. The study by Kircher and Foerderer (2024b) demonstrates that ad bans on YouTube for children’s content led to a decline in the quality of child-oriented content, which in turn reduced the audience for this content. Similarly, Johnson *et al.* (2024) found that the ad ban introduced by Google resulted in a decrease in both the production and overall quality of content directed at children by child-focused creators. By contrast, in this paper we focus on the question of what drives whether firms collect data from vulnerable individuals, and how this appears to be shaped by firm size.

The second stream of literature is the app market. This literature has focused on app-developer strategies to gain attention through distorting popularity information (Bresnahan *et al.*, 2014a,b), using free apps to build demand for paid apps (Deng *et al.*, 2023), overcoming search costs and navigation costs (Yin *et al.*, 2014; Ershov, 2024), and offering low price points in return for user data (Kummer and Schulte, 2019). This literature has also documented how app store policy affects app developer strategies, for example through its product rating system. Leyden (2025) shows that this policy change led to higher-quality products but less frequent product updates. Comino *et al.* (2019) show how a developer’s ability to post updates influences downloads. Bian *et al.* (2021) show that consumers reduce the demand for apps that disclose data collection practices after the platform’s privacy policy change. Based on app data, Mayya and Viswanathan (2024) show that when app stores allow developers to delay updating an app’s privacy policy, it results in a decline in downloads and user ratings. There is a growing literature that has attempted to characterize the market for child apps. Kesler *et al.* (2017) document that apps targeting age categories of 13+ and 16+ tend to be more intrusive compared to other age groups. In addition, Liu *et al.* (2016) and Reyes *et al.* (2018) demonstrate that most apps do not comply with U.S. child privacy regulation. Our paper builds on this literature by trying to uncover what shapes app developers’ decisions to

collect sensitive data from children.

The third stream of literature is that of the relationship between privacy and competition. The literature suggests that there is a trade-off between privacy regulation and competition—something that has been alluded to in theoretical work (Athey, 2015; Campbell *et al.*, 2015; Fuller, 2017; Tucker, 2019; de Cornière and Taylor, 2021; Krämer and Shekhar, 2024) and empirical work (Marthews and Tucker, 2019; Jia *et al.*, 2021; Peukert *et al.*, 2022; Pinto *et al.*, 2024). The literature also shows that privacy regulation may have differential effects depending on firms’ business models (Markovich and Yehezkel, 2024). Our results are important for competition authorities, because to our knowledge this is one of the first papers that explicitly asks whether larger or smaller firms collect more personal or intrusive data. Our results suggest that privacy protection designed to limit data collection adversely affects smaller firms more than larger firms. The final stream of literature we contribute to is that which tries to understand the relationship between data and market power. Data allows us to know consumers’ preferences, which then permits us to design better business models. Much of this economics literature has been devoted to the question of whether there are economies of scale and scope in data. Most of these papers have found evidence instead of diminishing returns to data (Chiou and Tucker, 2017; Bajari *et al.*, 2019; Peukert *et al.*, 2024; Farboodi *et al.*, 2019). The literature suggests that due to a cold start problem, data is valuable at the beginning, but in the long run, data have diminishing returns to improving predictions (Farboodi and Veldkamp, 2021; Goldfarb and Que, 2023). By contrast, we ask whether firm size appears to influence the amount of sensitive data collected.

Our results are important for regulators because of the importance of protecting children’s privacy and because of some of the intricacies of global competition in the digital space. Children’s privacy issues are particularly pressing, as in the United States 32% of the children between 7 and 9 years old use apps, and 49% of the children between 10 and 12 years old use social media apps.<sup>4</sup>

These results have several implications. First, many theories of competitive harm by large digital platforms are based on the idea that their size allows them to collect more sensitive

---

<sup>4</sup><https://www.statista.com/statistics/1293278/us-children-use-of-apps-by-age-group/> Last accessed December 2, 2024.

data. But we see no evidence of such a pattern in our data. Second, our results support the view that regulatory interventions should not only be imposed on larger companies, but also encourage compliance by small companies. Small developers often face significant constraints in implementing robust privacy protections. These include limited financial resources, lack of technical expertise, and the high cost of compliance. For instance, smaller firms may not afford dedicated privacy officers or advanced data protection technologies. To support small developers, regulatory bodies could offer assistance programs, and industry collaborations could facilitate shared access to privacy compliance resources. The findings suggest that a one-size-fits-all regulatory approach may not be effective in addressing privacy protections across different firm sizes. Policymakers should consider graduated compliance requirements that scale with firm size, providing smaller developers with feasible pathways to enhance privacy practices without imposing undue burdens. Additionally, incentives such as tax breaks or grants for small firms investing in privacy technologies could promote better compliance. It is also crucial to address the trade-offs involved, balancing the need for stringent privacy protections with the operational realities of small developers. On a global scale, international cooperation could help harmonize privacy standards, ensuring consistent protection for users regardless of the platform or country of origin.

The paper is structured as follows. Section 2 describes the data sources, presents the descriptive statistics and our variables of interest. Section 3 presents our empirical strategy and our main estimates. Section 4 shows the econometric results based on several different specifications and provides robustness checks. Section 5 presents the potential underlying mechanisms. The conclusion follows.

## 2 Data

We collect data on apps published in the Google Play Store. This is the largest worldwide platform that distributes apps for the Android ecosystem. We study children’s apps published in the U.S. Google Play Store. App descriptions are automatically released worldwide with automated translation unless the developer specifies otherwise.<sup>5</sup> We collect app data on

---

<sup>5</sup>Certain countries may impose additional requirements on developers to comply with local regulations.

children’s apps from mid-July 2017 to January 2021, tracking each app starting from its first appearance to the end of the sample period. We collect data on average every two weeks. Our final sample includes 106 weeks because we keep only weeks that contain the full sample of data. The final sample includes 1,498,645 observations<sup>6</sup> with 27,119 apps and 11,090 developers. This large number of apps reflects the fact that it is easy to produce and commercialize apps worldwide for children and especially those under five, since these apps are based mainly on images, sounds, and colors. This is something that has been estimated by Ghose and Han (2014) as part of a broader demand estimation exercise.

We aim to collect the broad market of children’s apps. We collect apps that belong to the Designed for Families program.<sup>7</sup> We complete this data with keyword searches with the aim of including broad apps that appeal to children. We identify the list of keywords most frequently associated with children’s apps such *kids* and *toddlers* using the Google Adwords keyword planning tool, a similar data collection is used in Cecere *et al.* (2025).

The search was repeated every two weeks on average to identify new apps. According to the FTC, general-audience content should comply with COPPA rules if it appeals to children, as underlined in the case against TikTok.<sup>8</sup>

Table 1 shows descriptive statistics. We collect all publicly available data over time, such as type of permissions required by apps, developer location, and app characteristics such as *Designed for Families*, *Average Stars Rating*, *Freemium Pricing*, *Paid App*. The Google Play Store provides 21 ranges of installs for each app, from 0 to 5 installs to more than 5 billion installs. We include a set of dummies representing each range presented in Table A3 in the Appendix C). We use an unbalanced panel, which accounts for entry and exit.

We study children’s apps published in the U.S. Google Play Store, but which have been developed worldwide. In our dataset, developers originate from 127 countries. We exploit geographical information disclosed by each developer to identify each developer’s country. Overall, a plurality of the apps in the U.S. market are produced by U.S. developers—24.78%

---

<sup>6</sup>An observation is at app and week level.

<sup>7</sup>The Designed for Families program includes three broad age categories aimed at children ages 0-5, 6-8 and 9+, with an additional six categories: Action & Adventure, Brain Games, Creativity, Education, Music & Video, and Pretend Play. While the choice of thematic category is optional, developers must choose appropriate age categories.

<sup>8</sup><https://www.ftc.gov/news-events/news/press-releases/2024/08/ftc-investigation-leads-lawsuit-against>  
Last accessed December 2, 2024.



of the sample. After the United States, the largest producers of children’s apps are India with 7.72%, and the United Kingdom with 6.32%.

Table 1: **Panel Data Summary Statistics**

	Mean	SD	Min	Max
Sensitive Data	0.589	1.122	0	11
Sharing	0.081	0.317	0	3
Location data	0.189	0.552	0	4
Identity Information	0.275	0.511	0	2
User Surveillance	0.043	0.283	0	5
Binary Sensitive Data	0.329	0.470	0	1
# Apps by Developer	24.162	46.305	1	289
1-2 Apps	0.351	0.477	0	1
3-6 Apps	0.164	0.371	0	1
7-24 Apps	0.241	0.428	0	1
25-67 Apps	0.145	0.352	0	1
68+ Apps	0.099	0.298	0	1
Increase: + Three Apps	0.047	0.213	0	1
Merger & Acquisition	0.096	0.294	0	1
Increase Only for New Apps	0.048	0.213	0	1
Designed for Families	0.705	0.456	0	1
Average Star Rating	4.177	0.592	1	5
Contains Ad	0.535	0.499	0	1
Freemium Pricing	0.377	0.485	0	1
Paid App	0.260	0.439	0	1
Log # User Reviews	5.194	3.640	0	18.7
Category Education	0.196	0.397	0	1
Category Educational	0.254	0.435	0	1
# Distinct Apps	27,119			
# Distinct Developers	11,090			
Observations	1,498,645			

*Notes:* This table presents the summary statistics for the panel data used in the analysis. The variables captured include measures of data collection practices, app characteristics, and developer attributes. The download and category dummies are presented in Table A3 and Table A4, respectively, in Appendix C.

## 2.1 Dependent Variable: Sensitive Data

COPPA regulation defines the list of child-sensitive data collection covered by the law. It includes geolocation details (sufficiently precise to identify street name and city), photos, videos, and audio files that contain children’s images or voices, usernames, and persistent

identifiers to recognize an app user over time and across different apps.<sup>9</sup> User data can be requested and collected using the permissions system implemented by the Google Play Store. To measure whether children’s apps violate COPPA, we identify the Google Play Store permissions and interactive elements that allow apps to collect these sensitive data on children.

We identify eleven permissions and three interactive elements that require personal data covered by the COPPA regulation. We created the variable *Sensitive Data*, which counts the types of sensitive data covered. We identify four broad categories of sensitive data: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. Table A1 in Appendix A presents the permissions and interactive elements required to construct the main dependent variable *Sensitive Data*.

Table 1 presents the descriptive statistics of the main dependent variable. The average number of pieces of sensitive data required by an app is 0.589. We also construct a dummy variable *Binary Sensitive Data* measuring whether the app requests at least one piece of sensitive data; 32.9 % of apps belong to this category.

## 2.2 Developer Size

Conceptually, developer size could affect the likelihood of sensitive data collection through channels, through compliance costs, and through shaping the underlying demand for sensitive data. In terms of compliance costs, on the one hand, larger developers may find it easier to internalize compliance costs and therefore may have lower marginal costs of collecting more sensitive data. The fixed cost of compliance may be substantial. In 2013 (when COPPA was last revised), the estimated average cost of compliance according to TechFreedom (working on behalf of the FTC) was around \$6,200 per year but up to \$18,670 a year for newly

---

<sup>9</sup>The law requires verifiable parental consent for the collection, use, and disclosure of personal information on children under age 13. This information is not available to researchers: only developers and users who actually use the app have access to this information. Thus, we are only able to measure the type of permissions required by each app. The complete list of children’s personal data is available in the FTC rule-making regulatory reform proceedings (<https://www.ftc.gov/enforcement/rules/rulemaking-regulatory-reform-proceedings/childrens-online-privacy-protection-rule>). Last accessed January 8, 2018.

created companies.<sup>10</sup> On the other hand, smaller developers may be more likely to take a less risk-averse approach and a non-robust approach to compliance, and consequently to have lower compliance costs for collecting more sensitive data. There is substantial legal risk from collecting sensitive data. Recent FTC and state cases show that the FTC imposes high settlements on firms that do not comply with COPPA, as shown in Table A2 Appendix B.

In terms of underlying demand, larger developers may find it desirable to collect more data because their scale of operations and data-sophistication means they can extract the most value from it. Smaller developers may find it desirable to collect more data because ultimately the incremental value of data is larger for smaller firms, given that data is often duplicative.

To effectively measure developer size, we count the number of apps each developer has available each week, referred to as *# Apps by Developer*. The average number of apps per developer is 24.16.

### 2.2.1 Alternative Measures of Developer Size

We use alternative measures of developer size. The marginal effect of producing one more app may impact smaller and larger developers differently. To account for this effect, we split the continuous variable *# Apps by Developer* into five categories, ranging from 1 app to over 68 apps using the 25th, 50th, 75th, and 90th percentile distribution. The categories are defined as follows: *1-2 Apps* indicates that at time  $t$  the developer has only 1 or 2 apps, *3-6 Apps* indicates that the developer has between 3 and 6 apps, *7-24 Apps* indicates that the developer has between 7 and 24 apps, *25-67 Apps* indicates that the developer has between 25 and 67 apps, and *68+ Apps* indicates that the developer has more than 68 apps (top decile). The smallest group of app developers represents the largest share, at 35.1%. We also use the variable *Log # Apps by Developer* which measures the log of number of apps by developer.

The literature on big data suggests that data performance does not depend linearly on the amount of data collected (Tucker, 2019). We want to investigate whether developers with

---

<sup>10</sup>These figures do not include additional costs and reduced revenue from ads. <https://www.lexology.com/library/detail.aspx?g=0b6d68a9-5d17-4d52-9b30-54d356ddb08a>. Last accessed May 31, 2020.

a large number of installs experience increasing returns to data collection. In this section, we use several metrics to capture alternative measures of developer size, based on the number of consumers (downloads) and new products (new apps). The number of downloads is another important measure of developer size, which is also considered by competition authorities in recent cases.<sup>11</sup> We construct the binary variable *Large # Installs* which takes value 1 if the developer has at least one app with more than 5 million downloads. This market is characterized by a high degree of skew in the size distribution of app demand (Bresnahan *et al.*, 2014a).

## 2.3 Initial Model-Free Evidence

Before turning to the regression analyses, we explore the raw data. Table 2 presents the average number of types of sensitive data collected by developer size. We find that 40.7% of apps produced by small developers request at least one type of sensitive data, but that percentage drops to 20.4% for larger developers. In all rows, the amount of sensitive data collected declines as the developer size increases. This descriptive evidence suggests that large developers are less likely to collect sensitive data.

Table 2: Sensitive Data Collected by Developer Size

	1-2 App (1)	3-6 Apps (2)	7-24 Apps (3)	25-67 Apps (4)	68+ Apps (5)
<b>Sensitive Data</b>	0.856	0.598	0.456	0.349	0.298
Sharing	0.135	0.088	0.055	0.027	0.027
Location data	0.308	0.188	0.122	0.099	0.063
Identity Information	0.325	0.276	0.269	0.213	0.206
User Surveillance	0.088	0.047	0.011	0.010	0.002
<b>Binary Sensitive Data</b>	0.407	0.332	0.309	0.254	0.204
# Distinct Apps	10,740	4,834	6,433	3,600	2,385
# Distinct Developers	9,356	1,280	561	99	21
Observations	525,433	246,490	361,492	217,091	148,139

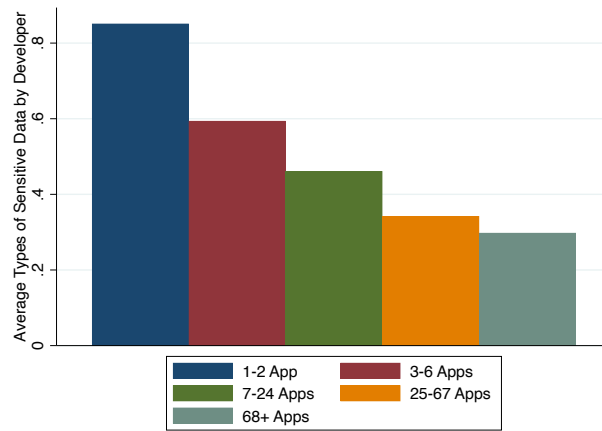
*Notes:* This table provides descriptive statistics segmented by developer size. Developer size is categorized based on the number of apps. The statistics include the mean and standard deviation for key variables, such as the number of permissions requested and types of sensitive data collected.

Figure 1 shows the average types of sensitive data requested by developer size. The

<sup>11</sup>[https://www.ftc.gov/system/files/documents/cases/musical.ly\\_complaint\\_ecf\\_2-27-19.pdf](https://www.ftc.gov/system/files/documents/cases/musical.ly_complaint_ecf_2-27-19.pdf)  
Last accessed June 5, 2019.

histogram reveals a clear and consistent trend across developer size categories: as developer size increases, the amount of sensitive data collected decreases. This pattern is evident with each step up in size category, from the smallest developers in the 25th percentile to the largest in the 90th percentile. The histogram effectively illustrates a negative correlation between developer size and the type of sensitive data collected.

**Fig. 1** *Average Types of Sensitive Data by Developer Size*



*Notes:* The y-axis indicates the average number of pieces of sensitive data collected by developer.

## 3 Empirical Analysis

### 3.1 Model Specification

We investigate the trade-offs between promoting competition and protecting children’s privacy. Strong privacy protections can protect children, but may adversely affect smaller developers. The decision to request sensitive data given the app quality can be correlated with developer size; this is a proxy for a developer’s ability to extract value from data and the ability to internalize compliance costs. Our empirical work aims to measure the effect of developer size on collecting sensitive data.

Building on our conceptual framework, we model how developer size is likely to influence the types of sensitive data requested. Our dependent variable, *Sensitive Data*, measures the pieces of sensitive data requested by each app  $i$  ( $i = 1$  to  $N = 27,119$ ) in week  $t$  ( $t = 1$  to  $T=106$ ). We use our panel data to estimate an OLS model with individual app fixed effects and time fixed effects and standard errors clustered on the app level.

We model the intensity of data collection using the following specification:

$$\text{Sensitive Data}_{it} = \alpha_0 + S_{it}\beta + \theta_{it} + \zeta_i + \rho_t + \epsilon_{it} \quad (1)$$

Our primary variable of interest is  $S$ , which indicates the developer size of app  $i$  at time  $t$ .  $\theta$  is a vector of other time-varying app characteristics, including the following variables *Contains Ad*, *Freemium Pricing*, *Paid App*, *Designed for Families*, *Average Star Rating*, and a vector of dummy variables indicating the intensity of download,<sup>12</sup> as well as a vector of dummy for apps categories.  $\zeta$  is the vector of app  $i$  fixed effects. Adding the app fixed effects ensures that identification of the coefficient is based on within-app variation over time rather than cross-app variation. The equation also includes time (week) effects  $\rho_t$ , which capture market trends related to privacy over time in our sample.  $\epsilon_{it}$  is the error term. Standard errors are clustered at the app level.

---

<sup>12</sup>The omitted category in the set of dummies that measures the number of downloads is apps with fewer than 50 downloads.

### 3.2 Results from Panel Data: Sensitive Data Collection from Children

We present our initial results when we examine how data collection is affected by developer size. Table 3 presents the main estimates and incrementally builds up to the final specification, Equation (1), in column (7). Column (1) presents the baseline pooled OLS. Columns (2)-(7) add the set of time-varying app characteristics, including the variables *Designed for Families*, *Average Star Rating*, *Contains Ad*, *Freemium Pricing* and *Paid App*, a vector of dummy variables measuring download intensity and a vector of dummy variables measuring app categories. We introduce country fixed effects in Columns (2)-(4). We observe a consistent negative correlation between the number of apps and the intensity of sensitive data requested. The estimates in Columns (4) and (5) include developer fixed effects respectively with and without time trends, illustrating the distinct behaviors of developers over time regarding sensitive data collection.

Columns (6) and (7) refine the model by adding app-specific fixed effects. Although we add multiple app characteristics, it is possible that unobserved factors vary between app and over time. We add a full set of app fixed effects<sup>13</sup> to absorb cross-sectional differences and week fixed effects. Column (6) includes a time trend at the app level. There is a negative association between increase in developer size and intensity of data collection. Column (7) includes time trend isolating the effects of app characteristics on data request practices over time.

There are many potential explanations for this finding. One is the theoretical findings in Campbell *et al.* (2015) that privacy regulation imposes costs on all firms, but larger firms are more likely to internalize these costs. For example, larger firms can benefit from economies of scale on the fixed compliance costs. In this case, regulation might distort competition against small companies. Another possibility is that companies may benefit from having large quantities of data but with diminishing returns to scale (Bajari *et al.*, 2019).

In Appendix E, we explore alternative measures of the dependent variable. We check whether a given set of sensitive data is driving our results. Appendix F provides further

---

<sup>13</sup>This corresponds to the week the data was scraped.

robustness checks by employing alternative measures of developer size such as the categorical variables *1-2 App*, *3-6 Apps*, *7-24 Apps*, *25-67 Apps*, *68+ Apps*, *Log # Apps by Developer* and the variable *Big # Installs* with our main estimate. Overall, our results are robust to these alternative measures of developer size. The pattern that larger developers are less likely to request sensitive data is replicated across these estimates. The coefficient of the variable *Big# Installs* is negative and statistically significant, suggesting that developers with at least one app with a large number of users are less likely to collect sensitive data. This pattern is consistent with a literature that has documented diminishing returns to data (Chiou and Tucker, 2017; Bajari *et al.*, 2019; Peukert *et al.*, 2024; Farboodi *et al.*, 2019). This result challenges the recent approach of different competition authorities of targeting larger firms. In Appendix G, we shift our focus from app-level to developer-level analysis by averaging data collection metrics across all apps managed by a single developer.

Table 3: **OLS Estimates: Drivers of Requests for Sensitive Data**

<i>Sensitive Data</i> as Dependent Variable	Pooled OLS	Country FE	Country × Time	Developer FE		Apps FE	
	(1)	(2)	(3)	Time Trend (4)	Time FE (5)	Time Trend (6)	Time FE (7)
# Apps by Developer	-0.002*** (0.000)	-0.002*** (0.000)	-0.002*** (0.000)	-0.003 (0.003)	-0.022*** (0.003)	-0.004** (0.002)	-0.009*** (0.003)
Apps Characteristics	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Country FE	No	Yes	Yes	No	No	No	No
Country FE × Times Trend	No	No	Yes	No	No	No	No
Category FE	No	No	No	Yes	Yes	Yes	Yes
Developer FE	No	No	No	Yes	Yes	No	No
Developer × Times Trend	No	No	No	Yes	No	No	No
Apps FE	No	No	No	No	No	Yes	Yes
Apps × Times Trend	No	No	No	No	No	Yes	No
Time FE	Yes	Yes	Yes	No	Yes	No	Yes
Mean Dependent	0.589	0.589	0.589	0.589	0.589	0.589	0.589
Adjusted $R^2$	0.145	0.173	0.182	0.834	0.808	0.971	0.942
Observations	1,498,645	1,498,645	1,498,645	1,498,633	1,498,633	1,498,645	1,498,645

*Notes:* OLS estimates. The dependent variable *Sensitive Data* includes four broad categories of data: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. Column (1) reports the results of a pooled OLS regression with app characteristics and time fixed effects. Column (2) adds country fixed effects to measure country-specific effects. Column (3) introduces country-by-time interaction terms. Column (4) includes developer fixed effects with a time trend to account for developer-specific effects and temporal trends. Column (5) includes developer fixed effects with time fixed effects. Column (6) introduces app fixed effects with a time trend, controlling for app-specific characteristics and trends over time. Column (7) includes app fixed effects with time fixed effects, providing the most granular control over both app-specific characteristics and fixed temporal effects. Robust standard errors are clustered at app level and reported in parentheses. Significance levels: \* $p < .10$ , \*\* $p < .05$ , \*\*\* $p < .01$

### 3.3 Identification Through Instrumental Variable Approach

Our initial specification assumes that the developer size is independent of other unobserved factors influencing decisions on data collection. This may be reasonable if data collection



is largely determined independently of developer size. Despite including numerous control variables and app fixed effects, potential unobserved biases may still exist. For example, if particular kinds of apps become more or less popular at one point in time, leading to an unobserved change in strategy on the part of developer. In this case, inadequate controls for app self-selection could lead to an underestimate of the positive impact of the developer’s size in data collection.

To address this concern, we employ an instrumental variable approach. We exploit the variation in business costs in the country of each developer. We select instruments that can influence developer size but are unlikely to directly impact data collection intensity. We collect data from the *World Bank’s Doing Business*<sup>14</sup> which measures the quality of business regulations in more than 150 countries using a set of indicators that influence firms’ activities (Regulations, 2019). Our chosen instruments are associated with the costs and regulatory burdens in the developers’ operating environments:

1. *RegisterScore-Cost % of Property Value*: It measures the cost of registering a business in a given country for a specific year.<sup>15</sup> This variable acts as a proxy for the regulatory and financial environment impacting developers, where higher registration costs could disproportionately affect smaller developers due to the fixed nature of such costs. This variable ranges from 0 to 100, 0 represents the worst regulatory performance and 100 the best regulatory performance.
2. *Price of Electricity (U.S. cents per kWh)*: This index reflects price of electricity in the largest business city of the country.<sup>16</sup> We hypothesize that price electricity increase with firm size due to tax codes that tend to favor smaller firms, affecting larger developers more significantly.

---

<sup>14</sup><https://subnational.doingbusiness.org/en/data/exploretopics/registering-property/score>, Last accessed January, 2025.

<sup>15</sup>This variable is defined by the World Bank as ”The score for cost benchmarks economies with respect to the regulatory best practice on the indicator.” <https://archive.doingbusiness.org/en/data/exploretopics/registering-property>, last accessed January 2025.

<sup>16</sup>This variable is defined by the World Bank as ”The price of electricity is measured in U.S. cents per kWh. A monthly electricity consumption is assumed, for which a bill is then computed for a warehouse based in the largest business city of the economy for the month of March. The bill is then expressed back as a unit of kWh.” <https://archive.doingbusiness.org/en/data/exploretopics/registering-property>, last accessed January 2025.

These instruments are expected to be correlated with firm size due to varying administrative burdens, but not with the firm’s choice of data collection practices, except through their effect on the size of the developer and their ability to grow.

We estimate a 2SLS IV estimate. App fixed effects are included in all regressions to account for app intrinsic characteristics. We also include the full set of time-varying app characteristics and the time fixed effects. Table 4 shows the main IV estimates. We report the estimates of the first stage specification. Consistent with our OLS findings, the IV results confirm our previous results: larger developers are likely to collect sensitive data. There is an increase in the magnitudes of the size variable as well as its associated standard error in the generalized method of moment, GMM-IV estimate, compared to OLS estimates.

The P-value from the F-test, which tests the joint significance of the instruments, confirms their validity. This satisfies the first necessary condition for instrumental validity. We also test the over-identification restrictions implied by using multiple instruments for a single endogenous regressor and report the Hansen J statistic and its associated P-value in the tables. These tests consistently fail to reject the null hypothesis, assuming that at least one instrument is exogenous, which provides further reassurance of the instruments’ validity. To ensure that the effect is not being driven by the span of the instruments, further checks are conducted. We conduct robustness checks by omitting one instrument at a time to ensure the reliability and validity of the IV estimates. The separate contributions of each instrument are explored in Table A10 in Appendix H.

Table 4: **Identification Strategy: IV Approach**

	Sensitive Data (1)	# Apps by Developer (2)
# Apps by Developer	-0.068*** (0.016)	
RegisterScore-Cost (% of property value)		0.085*** (0.026)
Price of Electricity (U.S. cents per kWh)		-0.008* (0.005)
Underidentification (LM)	13.351	
P-value (LM-Stat)	0.001	
Weak identification	5.455	
Hansen's J statistic	0.813	
P-value (J-Stat)	0.367	
Apps Characteristics	Yes	Yes
Apps FE	Yes	Yes
Time FE	Yes	Yes
Adjusted $R^2$	-0.173	
Number of groups	26,829	26,829
Observations	1,485,631	1,485,631

*Notes:* *Sensitive Data* is the dependent variable in Column (1), which presents the 2SLS estimates. Column (2) presents the first-stage estimates, where *RegisterScore-Cost (% of property value)* and *Price of Electricity (U.S. cents per kWh)* are used as instrumental variables for the number of apps by developer. Robust standard errors clustered at the app level are reported in parentheses. Significance levels: \* $p < .10$ , \*\* $p < .05$ , \*\*\* $p < .01$

## 4 Robustness Checks

Our empirical approach has a key challenge in identifying the causal effect of developer size. Thus, as our dataset covers three-year periods, we can conduct several robustness tests leveraging variations in developer size. First, we leverage sudden increases in developers’ size and exploit external growth driven by mergers and acquisitions. Second, we exploit non-experimental variation to identify the effect of developer size. Using the propensity-score matching method, we match apps of developers who experience growth at a specific point in time to apps with developers who do not grow over time.

### 4.1 How Increase in Developer Size Affects Data Collection

We exploit two sources of variation to estimate the causal effect of increase in developer size on developer data collection strategies. This robustness check is independently interesting because it can contribute to a debate related to recent anti-trust cases. First, we identify sudden increases in the number of apps by a given developer from one period to another. We create the variable *Increase + Three Apps* which takes value one if we observe an increase of more than three apps for a given developer from one period to another. This corresponds with up to the 25th percentile of app size distribution. This sudden increase in the number of apps might indicate a merger or acquisition or, alternatively, a substantial investment by a developer, both of which can affect the quantity and type of sensitive data collected. Second, we leverage changes in the developer’s name to formally identify mergers and acquisitions. For this purpose, we create the binary variable *Merger & Acquisition* which measures this change. This allows us to evaluate whether institutional changes can impact developer strategy in term of data collection. Third, we identify the newly created apps that are produced by developers from one week to another.

Column (1) of Table 5 shows a negative correlation between the sudden increase in number of apps and data collection, suggesting a possible consolidation of data practices or enhancements in privacy measures after a merger or significant investment in apps.

Column (2) relies on changes in the developer’s institutional details such as name, link,

or unique ID, which could indicate a merger or acquisition. Alongside this, the interaction variable  $\text{Merger} \times \# \text{ Apps by Developers}$  explores how changes in the ownership structure, combined with the size of the developer, impact data collection practices. The negative coefficient suggests a negative correlation between an increase in app count post-merger and data collection.

Column (3) narrows the focus to new apps added from one week to another, using the variable *Increase Only for New Apps*. Similar to the findings in Column (1), the increase in number of new apps is associated with stricter data privacy practices, potentially reflecting a strategic alignment with enhanced privacy standards.

Table 5: **Large Increase & Mergers**

<i>Sensitive Data</i> as Dependent Variable	Increase: + Three Apps (1)	Merger & Acquisition (2)	Increase Only for New Apps (3)
Increase: + Three Apps	-0.06570*** (0.019)		
Merger & Acquisition		0.00792 (0.010)	
# Apps by Developer		-0.00915*** (0.003)	
Merger & Acquisition $\times$ # Apps by Developer		-0.00029** (0.000)	
Increase Only for New Apps			-0.07388*** (0.021)
Apps Characteristics	Yes	Yes	Yes
App FE	Yes	Yes	Yes
Time FE	Yes	Yes	Yes
Mean Dependent	0.589	0.589	0.589
Adjusted $R^2$	0.942	0.942	0.942
Observations	1,498,645	1,498,645	1,498,645

*Notes:* The table presents OLS estimates with app and week fixed effects. The dependent variable *Sensitive Data* includes four broad categories of data: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. The top of the column indicates the main variable of interest estimated in each regression. Robust standard errors are clustered at the app level and reported in parentheses. Significance levels: \* $p < .10$ , \*\* $p < .05$ , \*\*\* $p < .01$

## 4.2 Using Non-experimental Variation to Tease Apart the Effect of Size

We are able to observe developers that increase in size over time during our data collection and those that do not increase in size. We use non-experimental variation to tease apart the difference in size, but this is obviously less clean than being able to directly randomize the two. Thus, we use the propensity-score matching method to match developers that

have increased in size to developers (that we observe at a given point in time) who do not increase in size during our data collection. This allows us to create a balanced group between developers who grew significantly and those who did not grow by matching them based on similar observable characteristics at the beginning of our data collection. We use the nearest neighbor matching method with one neighbor and no replacement, ensuring that each treated app produced by a developer that increases in size over time is matched to only one app produced by a developer that does not increase over time. We use different app characteristics to perform the propensity-score matching method *Designed for Families*, *Average Star Rating*, *Contains Ad*, *Freemium Pricing*, *Paid Apps*, *Log Nbr Reviews*, *Category Education*, and *Category Educational*. *Contains Ad*, *Freemium Pricing*, *Paid Apps* account for the app business models. *Log Nbr Reviews* measures the demand for apps. The set of variables *Designed for Families*, *Category Education*, and *Category Educational* identifies the type of content offered by the apps. We perform the matching based on the average values of the variables over the five periods following the entry of a given app into our panel.

To conduct the propensity-score matching method, we consider different threshold levels to measure the increase of developer size. First, we identify larger developers as those who introduced seven or more apps during the first period of analysis, representing developers up to the 50th percentile of the size distribution. We label them as “Apps Increase Seven”. This transition highlights developers who exhibit significant growth in size. Second, we identify a separate group of larger developers as those who introduced three or more apps during the same period, corresponding to developers up to the 25th percentile of the size distribution. This group is labeled as “Apps Increase Three”.

Table 6 presents the estimates of the main equation considering a different subsample of apps issued, using the propensity-score matching using the nearest-neighbor matching without replacement.

Column (1) shows the estimate of the subsample of apps considering a matched sample of developers who introduced more than seven apps compared to those who do not increase in size. The findings indicate a significant decrease in the intensity of data collection. Column (2) considers only developers that increase by more than three apps compared to those that do not increase in size over time. It shows a consistent reduction in data collection. As a

robustness check, we employ two thresholds jointly. Column (3) compares apps produced by developers that increase by more than seven apps or developers that increase by more than three apps to apps that do not increase over time. These results suggest that as developers grow, whether by increasing the number of apps substantially or crossing a higher threshold, they tend to collect less data per app.

Table 6: **Impact of Developer Growth on Data Collection: Propensity-Score Matching Approach**

<i>Sensitive Data</i> as Dependent Variable	Matching Based Apps Increase Seven	Matching Based Apps Increase Three	Matching Based On Either Condition
	(1)	(2)	(3)
# Apps by Developer	-0.004*** (0.001)	-0.009*** (0.001)	-0.008*** (0.001)
Apps Characteristics	Yes	Yes	Yes
Week FE	Yes	Yes	Yes
Apps FE	Yes	Yes	Yes
Mean Dependent	0.414	0.547	0.441
Adjusted $R^2$	0.911	0.946	0.934
Observations	117,717	90,207	168,657

*Notes:* This table presents OLS estimates with app and week fixed effects. The samples are matched based on nearest-neighbor matching without replacement. The dependent variable *Sensitive Data* includes four broad categories of data: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. Column (1) presents the estimates based on matched samples, where we compare apps produced by developers who increased their number of apps to over seven with those produced by developers who did not experience a similar increase in size. Column (2) presents the estimates based on matched samples, where we compare apps produced by developers who increased their number of apps to over three with those produced by developers who did not experience a similar increase in size. Column (3) combines apps produced by developers who either increase their number of apps to over seven or by more than four. Robust standard errors clustered at the app level are reported in parentheses. Significance levels:  $*p < .10$ ,  $**p < .05$ ,  $***p < .01$

## 5 Underlying Mechanisms

In this section, we explore various potential mechanisms to explain the observed differences in data collection between larger and smaller developers. Larger developers might adopt more sophisticated and responsible data-handling practices. The risks of not complying with data regulations are larger for large firms due to their increased exposure. Furthermore, larger developers are often better positioned to absorb the costs associated with stringent data regulations, enabling them to implement international standards more effectively. This capability allows them to innovate methods of app development continually, ensuring that each new product improves upon previous privacy and data management practices.

We present three types of evidence. First, we see whether firms that are larger behave differently than smaller developers when it comes to producing new apps. Second, we explore whether the data collection practice is likely to be associated with more ad-based business models. Third, we explore whether big developers located in laxer privacy regulation countries are less likely to collect sensitive data.

### 5.1 Do larger developers have different policies towards their new apps?

We investigate data collection by new apps introduced into the market and whether larger developers introducing new apps are more likely to collect sensitive data. We split our sample and run the regressions separately for the newly created apps and apps that have been updated since our data collection. Table 7 examines how developer size impacts data collection practices when developers produce new apps or update existing apps.

Column (1) estimates the main equation on the subsample of new apps released within the last six months; it shows a negative correlation between developer size and data collection. Column (2) analyses the subsample of new apps produced up to the six months. We find a negative correlation between developer size and data collection with a large magnitude compared to column (1). This indicates that newly created apps produced by larger developers significantly decrease their data collection.



We investigate the mechanisms that can drive the different patterns, and we examine whether apps that have been updated during the data collection period are similar to new apps. Column (3) shows the estimates on a subsample of app updates within the six months. The results show a similar trend to newer apps, with a negative correlation between developer size and data collection. Column (4) extends the analysis to updates on apps older than 6 months, observing a consistent negative effect.

Overall, the results suggest that larger developers tend to reduce data collection over time, particularly when they introduce updates or new apps.

Table 7: **New Apps and Updates**

<i>Sensitive Data</i> as Dependent Variable	New Apps		Updates	
	Six Months or less	More than Six Month	Six Months or less	More than Six Month
	(1)	(2)	(3)	(4)
# Apps by Developer	-0.001** (0.001)	-0.012*** (0.003)	-0.010*** (0.003)	-0.009** (0.004)
Apps Characteristics	Yes	Yes	Yes	Yes
Time FE	Yes	Yes	Yes	Yes
Apps FE	Yes	Yes	Yes	Yes
Mean Dependent	0.535	0.596	0.703	0.496
Adjusted $R^2$	0.978	0.945	0.936	0.964
Observations	193,531	1,304,895	668,977	829,189

*Notes:* This table presents OLS estimates with app and week fixed effects. The dependent variable *Sensitive Data* includes four broad categories of data: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. Each column estimates the impact of developer size on data collection strategies for different subsamples of apps based on the entry in the market and the most recent updates. At the top of each column, we indicate the type of subsample used in the estimates. Robust standard errors are clustered at the app level and reported in parentheses. Significance levels:  $*p < .10$ ,  $**p < .05$ ,  $***p < .01$ .

## 5.2 How Developers Evolve their Data Collection Practices in the Children’s App Market

We use non-experimental variation due to the market dynamics of new app creation to explore whether existing data protection practices of larger developers explain our findings. Table 8 analyzes the impact of the size of the developer portfolio on data collection practices from their initial apps to newly created apps, employing various analytical methods in different subsamples.

Columns (1) and (2) estimate the correlation between size and data collection on the subsample of all first apps by a given developer and first apps by developers who are not already considered big (below seven apps). We exclusively use the initial entries of apps

at time  $t$  without considering any subsequent time-varying factors. This method ensures that the observations reflect the data collection practices at the precise moment these apps were first introduced, providing information on the initial data management strategies of the developers. These show a consistent negative correlation between the number of apps a developer has and the sensitivity of data collected, with a more pronounced effect observed among developers who were not already big.

Columns (3) and (4) use a two-way fixed effect (TWFE) on the subsample of the first apps and the first seven apps by developers over time. The negative coefficients indicate that larger developers, even in the early stages of app releases, tend to collect less sensitive data, reflecting a strategic or compliance-driven approach to data privacy from the onset.

Columns (5) and (6) present results from a matching analysis, where developers are matched based on a propensity-score matching approach using the same of variables presented in Section 4.2. This method further confirms the trends observed in the OLS and TWFE models, with slightly weaker but still significant effects, emphasizing a more cautious approach to data collection by larger developers right from their first few apps.

Overall, the results show a broader trend that larger developers are likely to collect less sensitive data even as they launch their initial apps, possibly due to better resources, established data handling protocols, or a stronger inclination towards compliance with privacy standards, and they might become big because consumers value their apps.

Table 8: **First Apps of Developers**

<i>Sensitive Data as Dependent Variable</i>	OLS		TWFE		Matching	
	First Apps	First Apps Not already Big developers	First Apps	First Seven Apps	First Apps	First Seven apps
	(1)	(2)	(3)	(4)	(5)	(6)
# Apps by Developer	-0.006*** (0.001)	-0.017*** (0.004)	-0.004 (0.003)	-0.009** (0.003)	-0.003 (0.003)	-0.003* (0.002)
Apps Characteristics	Yes	Yes	Yes	Yes	Yes	Yes
Apps FE	No	No	Yes	Yes	Yes	Yes
Week FE	Yes	Yes	Yes	Yes	Yes	Yes
Mean Dependent	0.776	0.788	0.796	0.708	0.556	0.488
Adjusted $R^2$	0.344	0.348	0.954	0.948	0.946	0.937
Observations	11,090	10,683	571,867	831,196	19,579	59,822

*Notes:* This table presents results from OLS and TWFE regressions, and the matching estimates with app and week fixed effects. The dependent variable *Sensitive Data* includes four broad categories of data: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. Columns (1) and (2) use OLS focusing on the first apps and first apps by developers who are not already big. Columns (3) and (4) use TWFE estimate focusing on the first apps and the first five apps by developers. Columns (5) and (6) present results from a matching analysis based on nearest-neighbor matching method without replacement, focusing on the first apps and the first five apps by developers. Robust standard errors are clustered at the app level and reported in parentheses. Significance levels: \* $p < .10$ , \*\* $p < .05$ , \*\*\* $p < .01$ .

### 5.2.1 Sensitive Data collection and Large Developers Based on Business Models

Another concern is that the results are driven by differences between the type of business models. We aim to investigate whether larger developers exhibit a strategic preference for more complex monetization models, which might be associated with more data collection.

We rely on the framework proposed by Markovich and Yehezkel (2024) which highlights the strategic choices that platforms make between monetizing user data or protecting it. This theoretical framework suggests that the choice of the business models depends on competitive dynamics and the perceived value of the data. In our context, we test whether larger developers implement business models that allow them to leverage their scale and data management capabilities more efficiently. This could mean that, despite having the capacity to collect more data, these developers choose to collect less data, possibly due to efficiency gains in data usage (Bajari *et al.*, 2019). This strategic shift towards less intrusive data collection practices as a competitive advantage can attract privacy-conscious users. To investigate this further, we estimated the main equation on different subsamples to shed light on the underlying demand for data collection.

Table 9 presents the estimate of data collection strategies across different app monetization models, and we split the sample of apps according to the apps' business models. Column (1) focuses on apps that do not display any ads, showing that the impact of developer size on data collection is not significant, suggesting that ad-free apps might not scale down data collection as developer size increases. Column (2) includes apps that display ads to users, where the size coefficient is negatively correlated with data collection, suggesting that larger developers in this category manage data more efficiently, possibly optimizing their use of data for targeted advertising without needing to collect as extensively, echoing the trade-offs described by Markovich and Yehezkel (2024).

Column (3) excludes apps that use the freemium pricing business model and reveals a similar pattern as in the ad-supported apps, with larger developers collecting less data. Column (4) examines apps that implemented the freemium pricing model, again showing that larger developers are likely to collect less data, consistent with the efficiency hypothesis. Column (5) considers the subsample of apps using both contain ads and freemium models.

The estimate shows a negative correlation between developer size and data collection.

Finally, Column (6) estimates the main equation on the subsample of paid apps. The estimates show a similar trend, indicating that regardless of the app’s monetization strategy, larger developers tend to collect less data, likely due to their ability to leverage economies of scale and advanced analytics to extract value from smaller data sets. This pattern suggests that larger developers, who typically have more resources and sophisticated data processing technologies, may not need to collect data as aggressively because they can extract more value from what they already possess. This efficiency in data use could lead to less intrusive data collection practices, which might be a competitive advantage in attracting users concerned about privacy. On the other hand, smaller developers, especially those in highly competitive niches like freemium and ad-supported apps, seem to collect more data, potentially to maximize the effectiveness of their limited resources. Overall, our results align with the finding of Markovich and Yehezkel (2024), as they model how platforms choose between commercializing user data or protecting it based on competitive conditions and the nature of data benefits.

Table 9: **Sensitive Data Collection and Large Developers Based on Business Model**

<i>Sensitive Data</i> as Dependent Variable	Without Ads (1)	Contain Ads (2)	Without Freemium Pricing (3)	Freemium Pricing (4)	Contains Ads+ Freemium (5)	Free Pricing (6)	Paid Apps (7)
# Apps by Developer	-0.004 (0.003)	-0.012*** (0.002)	-0.008*** (0.002)	-0.014*** (0.002)	-0.015*** (0.002)	-0.009*** (0.003)	-0.014*** (0.003)
Paid App	0.006 (0.023)	0.155 (0.098)	0.013 (0.034)	0.081 (0.055)			
Freemium Pricing	0.054* (0.032)	0.024 (0.018)				0.035** (0.017)	-0.034 (0.047)
Contains Ad			0.015 (0.012)	0.004 (0.017)		0.008 (0.011)	0.047** (0.023)
Apps Characteristics	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Week FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Apps FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Mean Dependent	0.590	0.587	0.528	0.688	0.647	0.680	0.330
Adjusted $R^2$	0.966	0.922	0.954	0.934	0.924	0.940	0.954
Observations	696,640	801,701	933,967	564,544	393,847	1,108,135	390,362

*Notes:* The table presents OLS estimates with app and week fixed effects. The dependent variable *Sensitive Data* includes four broad categories of data: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. Each column estimates the impact of developer size on data collection strategies for different subsamples of apps based on their monetization models. App characteristics include *App Star Rating*, the vector of app categories and the vector of the number of downloads. Column (1) includes apps that do not display ads. Column (2) includes apps that display ads. Column (3) includes apps that do not use the freemium pricing model. Column (4) includes apps that use the freemium pricing model. Column (5) includes hybrid apps that combine features of both contain ads and freemium pricing strategies. Column (6) includes free apps, and Column (7) includes paid price apps. Robust standard errors are clustered at the app level and reported in parentheses. Significance levels:  $*p < .10$ ,  $**p < .05$ ,  $***p < .01$ .

### 5.2.2 Advertising Third Parties

In our exploration of underlying mechanisms that affect data collection, we consider the role of third-party advertising services. It is hypothesized that experienced developers might leverage these services to outsource advertising, potentially reducing their direct collection of sensitive data. To test this hypothesis, we conducted estimates for apps that advertise with the top three third-party advertising services—AdMob, OpenIAB, and Unity—separately and compared them with apps not using these services.

The analysis is structured to observe how the use of these third-party services influences the extent of data collection by developer size. For each advertising service, we compared apps that utilize the service (indicated as ‘=1’) against those that do not (‘=0’). Results are presented in Table 10.

Column (1) examines apps that rely on advertising without using AdMob. Column (2) focuses on apps that do use AdMob, where the coefficient for the number of apps per developer is slightly larger compared to column (1). This further reduction in data collection among larger developers using AdMob could indicate that these developers leverage AdMob’s resources and compliance frameworks to enhance their own data privacy measures as they scale up. Columns (3) and (4) shift the focus to developers using or not using OpenIAB. Similar to the previous observations, this suggests that larger developers continue to maintain or enhance data privacy practices as they expand. Columns (5) and (6) analyze the impact of using Unity. We find similar trends. These findings collectively suggest that as developers’ portfolios grow, regardless of their specific use of different third-party advertising platforms, there is a consistent trend towards reducing the collection of sensitive data. This trend might be influenced by enhanced privacy policies or the adoption of third-party services that help manage compliance and user privacy more effectively.

Table 10: **Largest Advertising Thirds Parties & Sensitive Data Collection**

	AdMob		OpenIAB		Unity	
<i>Sensitive Data</i> as Dependent Variable	= 0	= 1	= 0	= 1	= 0	= 1
	(1)	(2)	(3)	(4)	(5)	(6)
# Apps by Developer	-0.010*** (0.003)	-0.013*** (0.002)	-0.012*** (0.002)	-0.039*** (0.007)	-0.010*** (0.002)	-0.014*** (0.002)
Apps Characteristics	Yes	Yes	Yes	Yes	Yes	Yes
Apps FE	Yes	Yes	Yes	Yes	Yes	Yes
Time FE	Yes	Yes	Yes	Yes	Yes	Yes
Mean Dependent	0.546	0.613	0.573	1.444	0.626	0.483
Adjusted $R^2$	0.952	0.917	0.920	0.964	0.933	0.890
Observations	309,709	491,341	788,468	13,211	586,746	214,576

*Notes:* This table presents results from OLS regressions with app and week fixed effects. The dependent variable *Sensitive Data* includes four broad categories of data: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. Each column estimates the impact of developer size on data collection strategies for different subsamples of apps based on the third-party advertising service used. Column (1) includes apps that do not use AdMob. Column (2) includes apps that do use AdMob. Column (3) includes apps that do not use OpenIAB. Column (4) includes apps that do use OpenIAB. Column (5) includes apps that do not use Unity, and Column (6) includes apps that do use Unity. Robust standard errors are clustered at the app level and reported in parentheses. Significance levels:  $*p < .10$ ,  $**p < .05$ ,  $***p < .01$ .

### 5.3 Privacy Regulation Regime

National privacy regime variation across countries is extensive and leads to a wide range of country heterogeneity. We use variation in privacy regulation worldwide to estimate the effect of different kinds of privacy laws on the pieces of sensitive data collected. To explore this effect, we split the sample into groups of countries according to the stringency of the privacy regulation regime.

To assess differences in national regulatory frameworks, we augment our data with a vector of the institutional framework measures associated with the developer’s address. To account for the heterogeneity of countries in terms of privacy regulation, we use the international measure of the national privacy regime constructed by the French Privacy Regulation Authority (CNIL).<sup>17</sup> They categorize countries according to their level of compliance with EU privacy legislation (comparable to the U.S. COPPA legislation). Table A5 in the Appendix D presents countries categorized according to their level of compliance with EU privacy legis-

<sup>17</sup>CNIL, “La protection des données dans le monde”. <https://www.cnil.fr/fr/la-protection-des-donnees-dans-le-monde>. Last accessed January 8, 2018.

lation. The dummy variable *EU* identifies the developer country as part of the European Economic Area (EEA). The dummy variable *Recognized by the EU* indicates that the country's privacy laws are compatible with EU legislation and thus equally stringent as COPPA. The binary variable *With Legislation* indicates that the country has some level of privacy legislation. The binary variable *Independent Authority* indicates the existence of an independent authority regulating privacy. The dummy variable *No Privacy Law* indicates an absence of privacy laws in the developer's country.

The main findings are presented in Table 11, which reports the baseline specifications for different sub-samples and uses continuous measures of developer size to ensure consistent estimates across groups.

Column (1) examines apps from OECD countries, finding that size does not correlate with better privacy practices. This suggests that in economically developed regions with established privacy norms, as shown by the higher baseline, smaller developers are just as compliant with stringent data protection standards as larger developers. Column (2) looks at non-OECD countries, where a negative relationship between developer size and sensitive data collection is also significant, indicating that larger developers in less regulated regions might adopt more stringent data practices possibly to align with international norms. Column (3) focuses on U.S. developers, where findings suggest that developer size does not significantly influence the collection of sensitive data. This pattern might indicate that developers in the United States are well-acquainted with the stringent national privacy regulations, leading to uniformly cautious data practices across developers of all sizes, potentially driven by closer oversight by home regulatory authorities. Column (4) assesses apps from the EU, where developer size does significantly impact data collection, in contrast to the U.S.. Here, the smaller baseline suggests that high regulatory standards lead to better privacy practices across the market. Interestingly, the impact is even more pronounced among larger developers, indicating that while smaller developers adhere well to strict regulations, larger developers may leverage their resources to exceed these compliance standards, thereby enhancing their data protection practices further.

Column (5) and Column (6) report findings from countries whose privacy laws are recognized by the EU and those with an independent privacy authority, respectively. In both

cases, larger developers tend to collect less sensitive data, especially in regions with more formalized privacy regulations.

Columns (7) and (8) analyze the impact in countries with and without privacy legislation. The results show that larger developers in countries with some form of privacy legislation are less likely to collect sensitive data. The effect is more pronounced in countries without any privacy laws, where larger developers significantly reduce data collection, potentially as a way to self-regulate in the absence of formal legal requirements and to join.

Appendix I extends this analysis to investigate the effects of major regulatory changes, including FTC cases in appendix B and the implementation of General Data Protection Regulation (GDPR) in appendix I.1, finding consistent results. This extended analysis provides further evidence that our observed patterns are robust across major regulatory interventions, emphasizing the proactive adaptation of larger developers to evolving global privacy standards.

Table 11: **Privacy Regimes**

<i>Sensitive Data</i> as Dependent Variable	OECD	Non-OECD	U.S.	EU	Rec. EU	Ind. Aut	With leg	No Privacy
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
# Apps by Developer	-0.001 (0.003)	-0.010*** (0.003)	-0.004 (0.006)	-0.010*** (0.003)	-0.004 (0.005)	-0.022*** (0.006)	-0.004*** (0.001)	-0.024*** (0.005)
Apps Characteristics	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Apps FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Week FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Mean Dependent	0.580	0.599	0.691	0.509	0.651	0.639	0.559	0.708
Adjusted $R^2$	0.954	0.926	0.962	0.943	0.959	0.921	0.920	0.944
Observations	836,724	661,921	371,369	454,028	474,468	141,511	350,580	78,058

*Notes:* OLS with app and week fixed effects. The dependent variable *Sensitive Data* includes four broad categories of data: *Sharing, Location Data, Identity Information, and User Surveillance*. Column (1) shows the estimates within the sub-sample of apps produced in OECD member countries. Column (2) shows the estimates within the sub-sample of apps produced in non-OECD member countries. Column (3) reports the estimates of the sub-sample of apps produced in the U.S. Column (4) reports the estimates of the sub-sample of apps produced in the EU. Column(5) reports the estimates of the sub-sample of apps produced in countries with a privacy regulation regime recognized by the EU. Column (6) shows the estimates within the sub-sample of apps produced in countries with an independent privacy authority. Column (7) shows the estimates of apps produced in countries with a privacy legislation. Column (8) shows the estimates of apps produced in countries with no privacy legislation. Robust standard errors are clustered at app level and reported in parentheses. Significance levels:  $*p < .10$ ,  $**p < .05$ ,  $***p < .01$

## 6 Conclusion

Our study finds that as app developers grow, they often reduce their data collection practices, particularly regarding sensitive data. This finding challenges the assumption that



larger developers necessarily collect more sensitive data, suggesting instead a more complex relationship between developer size and privacy practices.

We observe a negative association between developer size and data collection. This suggests that an increase in size correlates with reduced data collection, indicating that larger developers may adopt more privacy-conscious practices. We further validate these results through an instrumental variable approach, addressing potential endogeneity and reinforcing our findings by isolating growth effects in developer size that are less likely to be driven by unobserved factors.

For further robustness checks, we employ a propensity-score matching approach, which confirms that developers who expand their app portfolios tend to collect less data over time. This trend appears consistently across various growth measures, highlighting a shift toward more cautious data management practices as developers scale. Additionally, our merger analysis shows that developers undergoing corporate restructuring, including mergers and acquisitions, exhibit similar reductions in data collection. This effect may stem from integrating privacy best practices across newly combined entities, aligning larger developers with more standardized approaches to data privacy. Additionally, our extended analyses in the appendices using alternative measures of developer size, such as categorical distinctions based on app counts and logarithmic transformations, support the main conclusions. Additionally, our extended analyses in the appendices, using alternative measures of developer size—such as categorical distinctions based on app counts and logarithmic transformations—reinforce our main conclusions. We also conducted checks at the developer level and examined data diversification practices, consistently finding the same effect.

Our investigation into the underlying mechanisms revealed that larger developers consistently collect less sensitive data, irrespective of their business model, suggesting an overarching improvement in data handling practices. This trend persists across developers from various countries, regardless of the stringency of local data protection regulations. Notably, these findings indicate that larger developers are converging towards a unified standard of data privacy, aligning their practices more closely with global norms. Additionally, our analysis found no evidence of data diversification among larger developers. Instead, these developers appear to streamline their data collection processes, enhancing efficiency and re-

ducing the variety of sensitive data gathered. This behavior reflects a strategic approach to data management that prioritizes privacy, potentially as a response to consumer expectations and international regulatory pressures, thereby fostering trust and compliance across diverse markets.

While this study focuses on the Google Play Store, it is essential to consider that privacy practices might differ across other platforms such as the Apple App Store. These differences can be attributed to platform-specific regulations and guidelines that influence developer behaviors. Future research should aim to compare privacy practices across multiple platforms to validate the generalizability of our findings and provide a comprehensive understanding of the digital app ecosystem. Further research is needed to investigate the extent to which privacy protection is also associated with better content for children. A potential limitation of our findings is that we have no information on the objectives of data collection beyond content improvement and expected user behavior. However, this study provides a first attempt to understand the complexity of the child app market and how national privacy regulation affects firms' decisions worldwide.

# References

- Adjerid, I., Acquisti, A., Telang, R., Padman, R. and Adler-Milstein, J. (2015). The impact of Privacy regulation and technology incentives: The case of health information exchanges. *Management Science*. 62(4), 1042–1063.
- Athey, S. (2015). Information, Privacy, and the Internet. *CPB Lecture, CPB Netherlands Bureau for Economic Policy Analysis*.
- Bajari, P., Chernozhukov, V., Hortaçsu, A. and Suzuki, J. (2019). The impact of big data on firm performance: An empirical investigation. *AEA Papers and Proceedings*. 109, 33–37.
- Banerjee, S., Chakraborty, I., Choi, H., Datta, H., Daviet, R., Farronato, C., Kim, M., Lambrecht, A., Manchanda, P., Oery, A. *et al.* (2024). Digital Platforms 2.0: Learnings, Opportunities, and Challenges. *Opportunities, and Challenges (May 31, 2024)*.
- Bian, B., Ma, X. and Tang, H. (2021). The Supply and Demand for Data Privacy: Evidence from Mobile Apps. *Available at SSRN*.
- Bleier, A., Goldfarb, A. and Tucker, C. (2020). Consumer Privacy and the future of data-based innovation and marketing. *International Journal of Research in Marketing*. 37(3), 466–480.
- Bresnahan, T., Davis, J. P. and Yin, P.-L. (2014a). Economic value creation in mobile applications. In *The Changing Frontier: Rethinking Science and Innovation Policy*. (pp. 233–286). University of Chicago Press.
- Bresnahan, T., Orsini, J. and Yin, P.-L. (2014b). *Platform choice by mobile app developers*. Technical report. National Bureau of Economic Research.
- Brough, A. R., Norton, D. A., Sciarappa, S. L. and John, L. K. (2022). The Bulletproof Glass Effect: Unintended Consequences of Privacy Notices. *Journal of Marketing Research*. 59(4), 739–754.
- Campbell, J., Goldfarb, A. and Tucker, C. (2015). Privacy Regulation and Market Structure. *Journal of Economics & Management Strategy*. 24(1), 47–73.
- Cecere, G., Le Guel, F., Lefrere, V., Yin, P. L. and Tucker, C. (2025). Privacy and Platform Governance: The Case of Apps for Young Children. *Academy of Management Perspectives*.
- Cheyre, C., Leyden, B. T., Baviskar, S. and Acquisti, A. (2023). The Impact of Apple’s App Tracking Transparency Framework on the App Ecosystem. *Working Paper*.
- Chiou, L. and Tucker, C. (2017). *Search engines and data retention: Implications for Privacy and antitrust*. Technical report. National Bureau of Economic Research.
- Comino, S., Manenti, F. M. and Mariuzzo, F. (2019). Updates management in mobile applications: iTunes versus Google Play. *Journal of Economics & Management Strategy*. 28(3), 392–419.
- Congiu, R., Sabatino, L. and Sapi, G. (2022). The Impact of Privacy Regulation on Web Traffic: Evidence From the GDPR. *Information Economics and Policy*. 61, 101003.

- de Cornière, A. and Taylor, G. (2021). *Data and competition: A general framework with applications to mergers, market structure, and Privacy policy*. Working Paper n. 20-1076, Toulouse School of Economics, France.
- Deng, Y., Lambrecht, A. and Liu, Y. (2023). Spillover Effects and Freemium Strategy in the Mobile App market. *Management Science*. 29(9), 4973–5693.
- Ershov, D. (2024). Variety-based congestion in online markets: evidence from mobile apps. *American Economic Journal: Microeconomics*. 16(2), 180–203.
- Farboodi, M., Mihet, R., Philippon, T. and Veldkamp, L. (2019). Big data and firm dynamics. *AEA Papers and Proceedings*. 109, 38–42.
- Farboodi, M. and Veldkamp, L. (2021). *A model of the data economy*. Technical report. National Bureau of Economic Research Cambridge, MA, USA.
- Farboodi, M. and Veldkamp, L. (2023). Data and markets. *Annual Review of Economics*. 15(1), 23–40.
- Fuller, C. S. (2017). The perils of Privacy regulation. *The Review of Austrian Economics*. 30(2), 193–214.
- Ghose, P. and Han, S. P. (2014). Estimating demand for mobile applications in the new economy. *Management Science*. 60(6), 1470–1488.
- Goldfarb, A. and Que, V. F. (2023). The economics of digital Privacy. *Annual Review of Economics*. 15, 267–286.
- Goldfarb, A. and Tucker, C. (2012). Privacy and innovation. *Innovation policy and the economy*. 12(1), 65–90.
- Goldfarb, A. and Tucker, C. E. (2011). Privacy regulation and online advertising. *Management Science*. 57(1), 57–71.
- Jia, J., Jin, G. Z. and Wagman, L. (2020). GDPR and the localness of venture investment. *Available at SSRN 3436535*.
- Jia, J., Jin, G. Z. and Wagman, L. (2021). The Short-Run Effects of GDPR on Technology Venture Investment. *Marketing Science*. 40(4), 661–684.
- Johnson, G., Lin, T., Cooper, J. C. and Zhong, L. (2024). COPPAcalypse? The Youtube Settlement’s Impact on Kids Content. *The Youtube Settlement’s Impact on Kids Content*.
- Johnson, G. A., Shriver, S. K. and Du, S. (2020). Consumer Privacy choice in online advertising: Who opts out and at what cost to industry? *Marketing Science*. 39(1), 33–51.
- Kesler, R., Kummer, M. E. and Schulte, P. (2017). *Mobile applications and access to private data: The supply side of the Android ecosystem*. ZEW - Centre for European Economic Research, Discussion Paper # 17-075.
- Kircher, T. and Foerderer, J. (2024a). Ban Targeted Advertising? An Empirical Investigation of the Consequences for App Development. *Management Science*. 70(2), 1070—1092.
- Kircher, T. and Foerderer, J. (2024b). Does Privacy Undermine Content Provision and Consumption? Evidence from Educational YouTube Channels. *Available at SSRN 4473538*.

- Krämer, J. and Shekhar, S. (2024). Regulating Algorithmic Learning in Digital Platform Ecosystems through Data Sharing and Data Siloing: Consequences for Innovation and Welfare. *MISQ*, *forthcoming*.
- Kummer, M. and Schulte, P. (2019). When private information settles the bill: Money and Privacy in Google’s market for smartphone applications. *Management Science*. 65(8), 3470–3494.
- Lefrere, V., Warberg, L., Cheyre, C., Marotta, V. and Acquisti, A. (2024). Does Privacy regulation harm content providers? A longitudinal analysis of the impact of the GDPR. *Management Science*.
- Leyden, B. T. (2025). Platform Design and Innovation Incentives: Evidence from the Product Rating System on Apple’s App Store. *International Journal of Industrial Organization*, 103133.
- Liu, M., Wang, H., Guo, Y. and Hong, J. (2016). Identifying and analyzing the Privacy of apps for kids. In *Proceedings of the 17th International Workshop on Mobile Computing Systems and Applications*. February. ACM, 105–110.
- Markovich, S. and Yehezkel, Y. (2024). Competing for Cookies: Platforms’ Business Models in Data Markets With Network Effects. *Available at SSRN 4770577*.
- Marthews, A. and Tucker, C. (2019). Privacy policy and competition. *Brookings Institution*.
- Mayya, R. and Viswanathan, S. (2024). Delaying Informed Consent: An Empirical Investigation of Mobile Apps’ Upgrade Decisions. *Management Science (Forthcoming)*.
- Miller, A. R. and Tucker, C. (2009). Privacy Protection and Technology Diffusion: The case of Electronic Medical Records. *Management Science*. 55(7), 1077–1093.
- Miller, A. R. and Tucker, C. (2011). Can Health Care Information Technology Save Babies? *Journal of Political Economy*. 119(2), 289–324.
- Miller, A. R. and Tucker, C. (2017). Privacy Protection, Personalized Medicine, and Genetic Testing. *Management Science*. 64(10), 4648–4668.
- Miller, K. M., Schmitt, J. and Skiera, B. (2024). The Impact of the General Data Protection Regulation (GDPR) on Online Usage Behavior. *arXiv preprint arXiv:2411.11589*.
- Miller, K. M. and Skiera, B. (2024). Economic consequences of online tracking restrictions: Evidence from cookies. *International journal of research in marketing*. 41(2), 241–264.
- Montes, R., Sand-Zantman, W. and Valletti, T. (2019). The Value of Personal Information in Online Markets with Endogenous Privacy. *Management Science*. 65(3), 1342–1362.
- Peukert, C., Bechtold, S., Batikas, M. and Kretschmer, T. (2022). Regulatory Spillovers and Data Governance: Evidence from the GDPR. *Marketing Science*. 41(4), 746–768.
- Peukert, C., Sen, A. and Claussen, J. (2024). The editor and the algorithm: Recommendation technology in online news. *Management Science*. 70(9), 5816–5831.
- Pinto, B., Sokol, D. D. and Zhu, F. (2024). The Antitrust and Privacy Interface: Lessons for Regulators from the Data. *Geo. Mason L. Rev.* 31, 1019.

- Regulations, D. B. M. B. (2019). Measuring Business Regulations. *The World Bank Group. Web*.
- Reyes, I., Wijesekera, P., Reardon, J., On, A. E. B., Razaghpanah, A., Vallina-Rodriguez, N. and Egelman, S. (2018). Won't somebody think of the children? Examining COPPA compliance at scale. *Proceedings on Privacy Enhancing Technologies*, 63–83.
- Rochelandet, F. and Tai, S. H. T. (2016). Do Privacy laws affect the location decisions of internet firms? Evidence for Privacy havens. *European Journal of Law and Economics*. 42(2).
- Shen, Y., Miller, K. M. and Li, X. (2024). *How Does Disabling Cookie Tracking Impact Online News Consumption?* Presented at the 2024 Annual Conference (Paper #3258).
- Tucker, C. (2019). Digital data, platforms and the usual [antitrust] suspects: Network effects, switching costs, essential facility. *Review of Industrial Organization*. 54(4).
- Yan, S., Miller, K. M. and Skiera, B. (2022). How does the adoption of ad blockers affect news consumption? *Journal of Marketing Research*. 59(5), 1002–1018.
- Yin, P. L., Davis, J. P. and Muzyrya, Y. (2014). Entrepreneurial innovation: Killer apps in the iPhone ecosystem. *American Economic Review*. 104(5), 255–59.

## Supplementary Appendix A:

### Descriptive Statistics of Permissions and Interactive Elements Used to Construct Sensitive Data

*Sensitive Data* is the major dependent variable because it aggregates all types of COPPA-designated categories of sensitive data. It includes four subsets of sensitive data measures: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. Table A1 presents the detailed descriptive statistics of each piece of sensitive data used to construct the dependent variable. It also provides detailed statistics by developer location.

The variable *Sharing* takes value 1 if the app requests at least one of the interactive elements allowing apps to share users' personal data with other apps and third parties; this includes *Share Location*, *Share Info*, and *Users Interact*. In 2015, the Google Play Store announced the presence of interactive elements to inform consumers about what information the app can access. The binary variable *Users Interact* measures whether the app exchanges sensitive data between users. This feature allows the app to be exposed to unfiltered/uncensored user-generated content including user-to-user communications and media sharing via social media and networks. *Share Info* measures whether the app shares users' personal information with third-parties such as Instagram, Viber, and other social networks. *Share Location* equals 1 if the app shares users' locations to other users of social network likes Facebook and Snapchat.<sup>18</sup>

We identify four permissions that request users' location data to construct the binary variable *Location Data*. *ALEC* (Access Location Extra Commands) indicates whether an app collects user's locations based on various device capabilities, and *ANBL* (Approximate Network Based Location) is used to access approximate location derived from network location sources such as cell towers and Wi-Fi. *MLST* (Mock Location Sources for Testing) is used to facilitate developer access to users' locations, and *Precise GPS Location* provides accurate location data.

The binary variable *Identity Information* includes two permissions to identify unique individual identity. The permission *Read Phone Status and Identity* allows developers to

---

<sup>18</sup>See [esrb.org](http://esrb.org). Last accessed July 21, 2020.

identify a smartphone’s unique IMEI (International Mobile Equipment Identity), which is considered a persistent unique identifier by COPPA and GDPR (Reyes *et al.*, 2018). The IMEI can be used to recognize a user over time and across different online services,<sup>19</sup>, and it could be used to log all kinds of personal data and target the consumer. The IMEI number also permits developers to know which advertising is already seen by a user. A child’s voice can be captured via the permissions *Record Audio*.

*User surveillance* is a binary variable that measures whether at least one permission allows access to user activity and contact information. *Read Your Own Contact Card* allows developers to access users’ contact cards and associate users’ phone numbers with their names. *RCEPCI* (Read Calendar Events Plus Confidential Information) is used to read information stored on users’ phones, including that of friends. *Read Your Contacts* indicates whether the app reads users’ stored contacts, including the frequency with which the user communicates with a given individual. The permission *Read Call Log* allows the app to access data about incoming and outgoing calls. *Read Your Browser History and Bookmarks* gives access to web browser information including internet account information.

---

<sup>19</sup>Complying with COPPA: Frequently Asked Questions. Last accessed September 3, 2020.



Table A1: **List of Permissions and Interactive Elements Used to Construct the Dependent Variable *Sensitive Data***

	Overall (1)	1-2 App (2)	3-6 Apps (3)	7-24 Apps (4)	25-67 Apps (5)	68+ Apps (6)
<b>Sharing</b>	0.081	0.135	0.088	0.055	0.027	0.027
Share Info	0.013	0.016	0.015	0.010	0.011	0.011
Share Location	0.014	0.025	0.015	0.005	0.006	0.004
User Interact	0.055	0.093	0.058	0.040	0.010	0.012
<b>Location data</b>	0.189	0.308	0.188	0.122	0.099	0.063
ALEC <sup>a</sup>	0.003	0.007	0.004	0.001	0.000	0.000
ANBL <sup>b</sup>	0.097	0.150	0.099	0.066	0.059	0.032
MLST <sup>c</sup>	0.001	0.002	0.000	0.000	0.000	0.000
Precise GPS Location	0.088	0.150	0.085	0.054	0.040	0.031
<b>Identity Information</b>	0.275	0.325	0.276	0.269	0.213	0.206
Read Phone Status And Identity	0.199	0.217	0.199	0.208	0.176	0.152
Record Audio	0.076	0.109	0.077	0.061	0.037	0.054
<b>User Surveillance</b>	0.043	0.088	0.047	0.011	0.010	0.002
Read Your Own Contact Card	0.005	0.009	0.006	0.001	0.001	0.000
RCEPCI <sup>d</sup>	0.007	0.010	0.008	0.004	0.004	0.001
Read Your Contacts	0.022	0.048	0.021	0.004	0.004	0.001
Read Call Log	0.005	0.011	0.006	0.001	0.001	0.000
Read Your Browser History and Bookmarks	0.004	0.009	0.006	0.000	0.000	0.000
Observations	1,498,645	525,433	246,490	361,492	217,091	148,139

*Notes:* This table depicts the summary statistics of the permissions and interactive elements used to construct the main dependent variable *Sensitive Data*. Column (1) presents the overall mean. Column (2) presents the mean for apps developed by developers with 1-2 apps. Column (3) presents the mean for apps developed by developers with 3-6 apps. Column (4) presents the mean for apps developed by developers with 7-24 apps. Column (5) presents the mean for apps developed by developers with 25-67 apps. Column (6) presents the mean for apps developed by developers with 68+ apps.

<sup>a</sup> ALEC: Access Location Extra Commands.

<sup>b</sup> ANBL: Approximate Network Based Location.

<sup>c</sup> MLST: Mock Location Sources for Testing.

<sup>d</sup> RCEPCI: Read Calendar Events Plus Confidential Information.

## Supplementary Appendix B:

### COPPA Enforcement

The FTC ensures compliance with COPPA legislation in the U.S. and in other countries. Since COPPA was implemented, the FTC has investigated more than 30 cases. Table A2 presents the important cases. Some of these cases involve the app developer directly. The FTC imposes strong requirements regarding the type of data that companies can collect and how they should protect children's personal information.<sup>20</sup>

Table A2: **COPPA Enforcement Actions**

Firms	Date	Settlement	Country	Mobile Apps
WW International, Inc.	2022	\$1,500,000	U.S.	Yes
OpenX Technologies, Inc.	2021	\$2,000,000	U.S.	No
Recolor	2021	\$3,000,000	U.S./ Finland	Yes
Musically (TikTok)	2019	\$5,700,000	China	Yes
HyperBeard	2019	\$150,000	U.S.	Yes
YouTube <sup>a</sup>	2019	\$170,000,000	U.S.	-
Inmobi	2016	\$950,000	Singapore	Yes
LAI Systems	2015	\$60,000	U.S.	Yes
Retro Dreamer	2015	\$300,000	U.S.	Yes
TinyCo, Inc.	2014	\$300,000	U.S.	Yes
Path, Inc	2013	\$800,000	U.S.	Yes
Artist Arena LLC	2012	\$1,000,000	U.S.	No
RockYou, Inc.	2012	\$250,000	U.S.	No
Broken Thumbs	2011	\$50,000	U.S.	Yes
Playdom, Inc.	2011	\$3,000,000	U.S.	No
Skidekids.com	2011	\$100,000	U.S.	No
Iconix Brand Group	2009	\$250,000	U.S.	No
Imbee.com	2008	\$130,000	U.S.	No
Sony Music Song BMG	2008	\$1,000,000	U.S.	No
Xanga.com	2006	\$1,000,000	U.S.	No
Ms. Fields Famous Brands	2003	\$100,000	U.S.	No

*Notes:* The table illustrates the amount of settlements imposed by FTC under COPPA rules. All cases can be found on the FTC website.

<sup>a</sup> [https://www.ftc.gov/system/files/documents/cases/youtube\\_complaint.pdf](https://www.ftc.gov/system/files/documents/cases/youtube_complaint.pdf).  
Last accessed May 31, 2020.

<sup>20</sup><https://www.ftc.gov/news-events/blogs/business-blog/2018/10/happy-20th-birthday-coppa>.  
Last accessed July 21, 2020.

## Supplementary Appendix C:

### Downloads and Categories

To measure the market size of a given app, we use the download category provided by the Google Play Store that includes 21 distinct groups. The number of downloads are presented in Table A3 and range from 0 to over five billion downloads. The table shows the mean of apps across download intervals.

Table A3: **Summary Statistics: Distribution of Downloads**

	Mean	Min	Max
Downloads 0	0.001	0	1
Downloads 1	0.014	0	1
Downloads 5	0.013	0	1
Downloads 10	0.059	0	1
Downloads 50	0.035	0	1
Downloads 100	0.099	0	1
Downloads 500	0.047	0	1
Downloads 1k	0.114	0	1
Downloads 5k	0.050	0	1
Downloads 10k	0.112	0	1
Downloads 50k	0.052	0	1
Downloads 100k	0.136	0	1
Downloads 500k	0.063	0	1
Downloads 1000k	0.124	0	1
Downloads 5000k	0.033	0	1
Downloads 10000k	0.034	0	1
Downloads 50000k	0.005	0	1
Downloads 100000k	0.005	0	1
Downloads 500000k	0.0008	0	1
Downloads 1000000k	0.0008	0	1
Downloads 5000000k	0.0001	0	1

*Notes:* The table illustrates the distribution of apps per download range and it indicates the lower range.

To examine the distribution of app categories in the sample, we use the categories provided by the Google Play Store, which include 52 distinct classifications. Table A4 presents the mean, minimum, and maximum values for each category, indicating the frequency and diversity of app types in our dataset.

Table A4: **Summary Statistics: Distribution of Categories**

	Mean	Min	Max
Action	0.013	0	1
Adventure	0.022	0	1
Arcade	0.030	0	1
Art and Design	0.005	0	1
Auto and Vehicles	0.002	0	1
Beauty	0.000	0	1
Board	0.008	0	1
Books and Reference	0.027	0	1
Business	0.002	0	1
Card	0.004	0	1
Casino	0.000	0	1
Casual	0.110	0	1
Comics	0.002	0	1
Communication	0.005	0	1
Dating	0.001	0	1
Education	0.196	0	1
Educational	0.254	0	1
Entertainment	0.040	0	1
Events	0.000	0	1
Finance	0.002	0	1
Food and Drink	0.001	0	1
Health and Fitness	0.016	0	1
House and Home	0.001	0	1
Libraries and Demo	0.000	0	1
Lifestyle	0.011	0	1
Maps and Navigation	0.002	0	1
Medical	0.005	0	1
Music	0.009	0	1
Music and Audio	0.008	0	1
News and Magazines	0.002	0	1
Parenting	0.013	0	1
Personalization	0.004	0	1
Photography	0.004	0	1
Productivity	0.007	0	1
Puzzle	0.067	0	1
Racing	0.012	0	1
Role Playing	0.017	0	1
Shopping	0.001	0	1
Simulation	0.027	0	1
Social	0.004	0	1
Sports	0.009	0	1
Strategy	0.005	0	1
Tools	0.022	0	1
Travel and Local	0.003	0	1
Trivia	0.007	0	1
Video Players and Editors	0.003	0	1
Weather	0.002	0	1
Word	0.012	0	1

*Notes:* The table illustrates the distribution of apps categories in our sample.

## Supplementary Appendix D:

### Developer Location

To explore U.S. regulation spillovers to other countries, we retrieve geographical information disclosed by the developers of apps available in the Google Play Store. Although the FTC requires that firms collecting or maintaining sensitive data from children should indicate in their online notices or information practices their name, address, telephone, and email address, several developers fail to provide a geographical address.<sup>21</sup>

To retrieve developers' countries, we use different strategies. First, we use Maps' APIs to collect the latitudes and longitudes of the given address to identify the country. Second, we used a Python library (Libpostal)<sup>22</sup> to search for a country name in the developer's address. Third, we check the match between the location identified using Google Maps' APIs and the country name identified via Libpostal. Fourth, among the subset of apps without any developer's address, we identify their location using the email extension. Finally, we manually check for certain addresses. We delete apps produced by developers that did not indicate their geographical location since this did not allow us to identify country of origin.

Table A5: **Privacy Regime Based on EU Privacy Regulation: List of Countries Presented in Our Sample**

Category	Countries
EU	EU members, Iceland, Norway, United Kingdom
Recognized by EU	Andorra, Argentina, Canada, Israel, New Zealand, Switzerland, U.S. <sup>a</sup> , Uruguay
Independent Authority	Albania, Australia, Bosnia and Herzegovina, Colombia, Costa Rica, Gabon, Ghana, Hong Kong, Korea, Rep., Kosovo, Macedonia, FYR, Mexico, Moldova, Morocco, Senegal, Serbia, Tunisia, Ukraine
With Legislation	Angola, Armenia, Azerbaijan, Brazil, Chile, China, India, Indonesia, Japan, Kazakhstan, Kyrgyz Republic, Malaysia, Montenegro, Nepal, Nicaragua, Philippines, Qatar, Russian Federation, Seychelles, Singapore, South Africa, Taiwan, China, Thailand, Turkey, Vietnam, Yemen, Rep., Zimbabwe
No Privacy Law	Afghanistan, Algeria, Bahrain, Bangladesh, Barbados, Belarus, Bolivia, Cambodia, Congo, Rep., Cuba, Dominican Republic, Ecuador, Egypt, Arab Rep., El Salvador, Ethiopia, Guatemala, Honduras, Iran, Islamic Rep., Iraq, Jamaica, Jordan, Kenya, Kuwait, Lao PDR, Lebanon, Mongolia, Mozambique, Myanmar, Nigeria, Oman, Pakistan, Palau, Palestine, Panama, Peru, Puerto Rico, Samoa, Saudi Arabia, Sri Lanka, Tanzania, Uganda, United Arab Emirates, Uzbekistan, Venezuela, RB

Notes: This table presents countries categorized according to their level of compliance with EU Privacy legislation.

<sup>a</sup> In July 2020, the EU Court of Justice invalidated the EU-U.S. Privacy Shield Framework. We consider the U.S. as "Recognized by the EU" prior to July 2020.

<sup>21</sup><https://www.ecfr.gov/current/title-16/chapter-I/subchapter-C/part-312>. Last accessed March 2, 2022.

<sup>22</sup><https://github.com/openvenues/pypostal>. Last accessed February 13, 2020.

## Supplementary Appendix E:

### Alternative Measures of Sensitive Data

#### E.1 Data Diversification: Alternative Dependent Variables

Another concern is that larger developers might collect less sensitive data on single user accounts than a one-app developer, as larger developers can collect different pieces of sensitive data for each single app. To address this, we investigate whether large developers diversify their data collection across their apps compared to small developers. Table A6 explores how the number of apps managed by developers influences their practices in collecting various types of sensitive data. Table A6 considers alternative dependent variables. Column (1) uses as dependent variable *Identity Info*. The negative coefficient suggests that developers with more apps tend to collect less identity-related information. Column (2) estimates the main equation using *Location* as the dependent variable; we observe a similar trend where an increase in the number of apps correlates with reduced collection of location information. Column (3) uses as dependent variable *User Surveillance*; the correlation between developer size and data collection is not significant in this estimate. Column (4) uses *Sharing* as dependent variable, which is positive and significant, suggesting a marginal increase in data sharing practices as the number of apps increases. This could suggest that larger developers might have more complex data sharing arrangements, perhaps due to more integrated services or partnerships. We also show robustness to a specification with a binary dependent variable that captures whether an app collects any sensitive data. Column (5) uses *Binary Sensitive Data* as the dependent variable. The coefficient is negative and significant. Finally, we introduce the variable *IMEI Plus One*, which tracks whether apps collect *IMEI* permission along with another type of sensitive data. Column (6) shows that the correlation between size and *IMEI Plus One* is statistically significant and negative. The collection of IMEI numbers can facilitate user tracking across different apps and services, significantly impacting user privacy. The negative coefficient associated with developer size for this metric suggests that larger developers may indeed limit this form of data collection, which could mitigate some concerns about cumulative data diversification.

Table A6: **Analyzing the Probability of Different Types of Data Collection by App Developer Size**

	(1) Binary Identity Information	(2) Binary Location Data	(3) Binary User Surveillance	(4) Binary Sharing	(5) Binary Sensitive Data	(6) IMEI Plus One
# Apps by Developer	-0.00359*** (-3.89)	-0.00354*** (-3.61)	0.0000520 (0.81)	0.000168* (1.91)	-0.00376*** (-3.97)	-0.00276*** (-3.08)
Apps Characteristics	Yes	Yes	Yes	Yes	Yes	Yes
Apps FE	Yes	Yes	Yes	Yes	Yes	Yes
Time FE	Yes	Yes	Yes	Yes	Yes	Yes
Mean Dependent	0.245	0.117	0.028	0.070	0.329	0.086
Adjusted $R^2$	0.886	0.909	0.904	0.951	0.896	0.891
Observations	1,498,645	1,498,645	1,498,645	1,498,645	1,498,645	1,498,645

Notes: This table presents results from LPM regressions with app and week fixed effects. The dependent variable in each column represents a different type of data collection practice. Robust standard errors are clustered at the app level and reported in parentheses. Significance levels: \* $p < .10$ , \*\* $p < .05$ , \*\*\* $p < .01$ .

## E.2 Alternative Definitions of Dependent Variables

We check whether our results hold for different measures of sensitive data. Table A7 checks the robustness of the results to alternative dependent variables. Column (1) shows the robustness to using the variable *Prop Sensitive Data*. The coefficient is negative and significant. One potential critique is that our main dependent variable includes a broad definition of sensitive data. We check whether a given set of sensitive data is driving our results. We use a set of binary variables. Thus, each time we exclude one category of sensitive data. Column (3) excludes the set of data *Prob Sharing* and column (4) excludes *Prob Location Data*. Column (5) reports the estimates when the dependent variable is the main dependent variable excluding *Prob Identity Information*. The estimates show that medium and large size developers are less likely to collect more sensitive data. Column (6) estimates the main dependent variable excluding *Prob User Surveillance*. Overall, we find that larger developers collect less data. Larger developers might be more careful to share and collect information from this vulnerable audience.

Table A7: **Alternative Definitions of Dependent Variables**

	(1) Binary Sensitive Data	(2) Sensitive Data without Sharing	(3) Sensitive Data without Location Data	(4) Sensitive Data without Identity Information	(5) Sensitive Data without User Surveillance
# Apps by Developer	-0.004*** (0.001)	-0.010*** (0.003)	-0.003*** (0.001)	-0.006*** (0.002)	-0.009*** (0.003)
Apps Characteristics	Yes	Yes	Yes	Yes	Yes
Apps FE	Yes	Yes	Yes	Yes	Yes
Time FE	Yes	Yes	Yes	Yes	Yes
Mean Dependent	0.329	0.507	0.400	0.313	0.546
Number of groups	27,119	27,119	27,119	27,119	27,119
Cluster	Apps	Apps	Apps	Apps	Apps
Adjusted R2	0.896	0.930	0.937	0.947	0.934
Observations	1,498,645	1,498,645	1,498,645	1,498,645	1,498,645

*Notes:* The table presents estimates with app and week fixed effects. Each column indicates the dependent variable. Column (1) uses a Linear Probability Model (LPM) to estimate the probability of sensitive data collection. Columns (3) - (6) use OLS. Robust standard errors are clustered at the app level and reported in parentheses, except for Column (2), which uses robust standard errors. Significance levels: \* $p < .10$ , \*\* $p < .05$ , \*\*\* $p < .01$ .



## Supplementary Appendix F:

### Alternative Measures of Developer Size

We also assess the robustness of our results to three alternative measures of developer size.

Table A8 reports the main estimates. This robustness check is independently interesting because it shows that apps produced by larger developers are less likely to collect sensitive data. Overall, the results are consistent with the main estimates in Table 3.

Column (1) explores the effects of categorical developer size, where categories are determined based on the 25th, 50th, 75th, and 90th percentiles of the number of apps developed by developers. The omitted category is the variable *1-2 App*. It illustrates a consistent negative relationship between developer size and sensitive data collection. In Column (2), we check the robustness of our findings to the log-log functional form transformation. The results remain robust to this specification, suggesting that extreme values do not drive our results. Column (3) includes the binary variable *Large # Installs* to investigate whether developers who have access to a large number of users are less likely to collect sensitive data. Overall, our results are robust to these alternative measures of developer size. The pattern that larger developers are less likely to request sensitive data is replicated across these estimates.

Table A8: **Alternative Measure of Size of Developers**

	Categorical Size Measure (1)	Log-Log (2)	Large Installs (3)
3-6 Apps	0.070 (0.046)		
7-24 Apps	0.036 (0.053)		
25-67 Apps	-0.204*** (0.065)		
68+ Apps	-0.545*** (0.155)		
Log # Apps by Developers		-0.065*** (0.017)	
Large # installs			-0.061*** (0.020)
App Installs	Yes	Yes	No
App Characteristics	Yes	Yes	Yes
Apps FE	Yes	Yes	Yes
Time FE	Yes	Yes	Yes
Mean Dependent	0.589	0.313	0.589
Adjusted R2	0.942	0.926	0.942
Observations	1,498,645	1,498,645	1,498,645

*Notes:* The table presents results from OLS regressions with app and week fixed effects. In Columns (1), the dependent variable *Sensitive Data* includes four broad categories of data: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. Columns (2) uses dependent variable the transformation of *Sensitive Data*. In Columns (3), the dependent variable *Sensitive Data* includes four broad categories of data: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. Robust standard errors are clustered at the app level and reported in parentheses. Significance levels:  $*p < .10$ ,  $**p < .05$ ,  $***p < .01$ .

## **Supplementary Appendix G:**

### **Estimate at Developer Level**

The table presented in Section A9 offers a detailed analysis of the relationship between the number of apps by a developer and the average amount of various types of sensitive data collected. This analysis is conducted at the developer level, averaging the collected data types across all apps managed by each developer to provide a broader perspective on data collection practices.

Column (1) focuses on the average amount of sensitive data collected. The coefficient for the number of apps managed by a developer is -0.020, which is statistically significant, indicating that developers managing more apps tend to collect less sensitive data overall. Column (2) examines the average amount of IMEI data collected, showing a similar trend with a coefficient of -0.007. This reduction among larger developers may reflect specific strategies to minimize the collection of highly sensitive data that can uniquely identify devices. Columns (3) and (4) show that as developers with more apps reduce the collection of location and identity-related data. The coefficient is negative, supporting the notion that developers with more apps collect less location data. Columns (5) and (6) shift the focus to read data and shared data, respectively. Unlike the other types of sensitive data, the coefficients in these columns are not significant, suggesting that the number of apps a developer manages does not influence the amount of read or shared data.

Overall, the table indicates a general trend where larger developers tend to reduce the collection of various sensitive data types as they manage more apps.

Table A9: **Estimates at the Developer Level**

	Mean Sensitive data (1)	Mean IMEI (2)	Mean Location (3)	Mean Identity (4)	Mean Read Data (5)	Mean Share (6)
# Apps by Developer	-0.019** (0.008)	-0.006* (0.003)	-0.010* (0.005)	-0.007** (0.003)	0.000 (0.001)	-0.000 (0.002)
Observations	Yes	Yes	Yes	Yes	Yes	Yes
Developer FE	Yes	Yes	Yes	Yes	Yes	Yes
Time Fe	Yes	Yes	Yes	Yes	Yes	Yes
Mean Dependent	0.798	0.215	0.280	0.271	0.051	0.106
Number of groups	11,076	11,076	11,076	11,076	11,076	11,076
Adjusted R2	0.945	0.876	0.927	0.893	0.912	0.941
Observations	609,410	609,410	609,410	609,410	609,410	609,410

*Notes:* The table presents results from OLS regressions with developer and week fixed effects. Each column estimates the average amount of a specific type of sensitive data collected at the developer level. Column (1) focuses on the mean amount of overall sensitive data collected. Column (2) examines the mean amount of IMEI data collected. Column (3) focuses on the mean amount of location data collected. Column (4) looks at the mean amount of identity-related data collected. Column (5) examines the mean amount of read data. Column (6) focuses on the mean amount of shared data. Robust standard errors are clustered at the developer level and reported in parentheses. Significance levels: \* $p < .10$ , \*\* $p < .05$ , \*\*\* $p < .01$ .

## Supplementary Appendix H:

### Estimates Excluding a Single Instrument

We conduct robustness checks by omitting one instrument at a time to ensure the reliability and validity of the IV estimates. Table A10 shows the estimates with a single instrument.

Columns (1) and (2) use *RegisterScore-Cost*. In Column (1), the results indicate that larger developers collect less sensitive data, validating the negative relationship between developer size and data collection practices when considering the costs associated with registering a business. Column (2) shows the first-stage regression where *RegisterScore-Cost* positively predicts developer size. Columns (3) and (4) include *Price of Electricity* as instrument. Column (3) presents a stronger negative effect of developer size on sensitive data collection compared to the first instrument, supporting the consistency of the negative relationship. Column (4) details the first-stage regression with the price of electricity negatively predicting developer size.

Both instrument sets control for variations by including time varying app characteristics, week fixed effects, and app fixed effects. The underidentification tests (LM statistic) in Columns (1) and (3) confirm the relevance of both instruments, indicating that they address the potential endogeneity of developer size.

Table A10: **Estimates Excluding a Single Instrument**

	Sensitive Data (1)	# Apps by Developer First Stage (2)	Sensitive Data (3)	# Apps by Developer First Stage (4)
# Apps by Developer	-0.063*** (0.016)		-0.106** (0.045)	
RegisterScore-Cost (% of property value)		0.085*** (0.026)		
Price of Electricity (U.S. cents per kWh)				-0.010* (0.006)
Apps Characteristics	Yes	Yes	Yes	Yes
Time FE	Yes	Yes	Yes	Yes
Apps FE	Yes	Yes	Yes	Yes
Underidentification (LM)	10.347		12.201	
P-value (LM-Stat)	0.001		0.000	
Weak identification	10.536		2.899	
Observations	1,497,935	1,497,935	1,479,595	1,479,595

*Notes:* The table presents results from separate instrumental variable (IV) estimates with app and week fixed effects. Columns (1) and (3) estimate the impact of developer size on sensitive data collection using two different instruments. Column (1) uses *RegisterScore-Cost* as the instrument, and Column (3) uses the *Price of Electricity*. Columns (2) and (4) show the first-stage regressions for these instruments. Robust standard errors are clustered at the app level and reported in parentheses. Significance levels: \* $p < .10$ , \*\* $p < .05$ , \*\*\* $p < .01$ .

## Supplementary Appendix I: Enforcement and Regulation Change

We explore the impacts of two major FTC cases on data collection practices. Table A11 presents the estimate where we consider how large developers change their data collection strategy after the FTC case against Musically (TikTok) and YouTube. It includes interaction terms between the coefficient associated to large developer and the event associated to the FTC cases. These interactions provide insight into how specific regulatory actions influence data collection practices among different developer sizes and app contexts. Column (1) shows that after the FTC case against Musically large developers did not significantly change their data collection practices. Column (2) includes the interaction terms between the variable measuring developers and the dummy variable indicating the FTC decision against YouTube. The interaction term is significant and negative, suggesting that this decision has significantly influenced the intensity of data collection for big developers. It might be possible that this decision taken against large U.S. developer is more likely to affect large developers.

Table A11: **Effect of FTC Case**

<i>Sensitive Data</i> as Dependent Variable	(1)	(2)
# Apps by Developer	-0.000378* (-1.67)	-0.00188 (-1.54)
Musically=1 × # Apps by Developer	0.00000766 (0.50)	
Youtube=1 × # Apps by Developer		-0.000149*** (-3.13)
App Characteristics	Yes	Yes
Week FE	Yes	Yes
Apps FE	Yes	Yes
Mean Dependent	0.577	0.531
Adjusted $R^2$	0.992	0.971
Observations	154,876	152,153

*Notes:* The table presents results from OLS regressions with app and week fixed effects. The dependent variable *Sensitive Data* includes four broad categories of data: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. Robust standard errors are clustered at the app level and reported in parentheses. Significance levels: \* $p < .10$ , \*\* $p < .05$ , \*\*\* $p < .01$ .

## I.1 Data Collection of the GDPR

Our analysis extends to examine the impact of the GDPR on developers' data collection behaviors. The GDPR was enacted in May 2018. Before and after implementation of GDPR, the Google Play Store introduced several policies aimed at removing suspicious and malicious apps from the Google Play Store. In 2017, more than 700,000 apps were removed and numerous other malicious apps were removed in October 2018.<sup>23</sup> Thus, their simultaneous implementation makes it difficult to disentangle the effects of these platform-specific policies from the potential effect of the GDPR. Additionally, the provisions of the GDPR are very similar to those included in the COPPA legislation, and thus, it is difficult to have separately measurable effects on children's apps. Despite these challenges, we investigated the potential impact of GDPR on data collection practices.

Despite this, the results, detailed in the table A12, underscore a consistent trend of reduced data collection by larger developers, that might be slightly accentuated by the introduction of GDPR. Column (1) explores the intensity of data collection before the implementation of GDPR, indicating that larger developers request less sensitive data than smaller developers. Column (2) shifts the focus to the period after GDPR, showing a continued, albeit smaller, reduction in data collection by larger developers, suggesting that if there is an impact of GDPR, it appears marginal given the pre-existing trends. Column (3) includes all sample in the analysis and it shows a slight intensification in the reduction of data collection post-GDPR among larger developers. Columns (4) and (5) differentiate the effects based on developers' locations, with EU developers showing a lesser adjustment to GDPR compared to non-EU developers, who exhibit a more pronounced compliance response. This analysis demonstrates that while GDPR may have reinforced the trend towards stringent data management among larger developers, particularly impacting non-EU developers by compelling them to enhance their privacy measures to align with EU standards, suggesting that large developer converging towards the same privacy standards.

---

<sup>23</sup><https://mobappdaily.medium.com/google-kicked-out-over-700-000-android-apps-from-play-store-in-2017>  
last accessed January 24, 2025.

Table A12: **The GDPR Effect**

<i>Sensitive Data</i> as Dependent Variable	Before GDPR	After GDPR	Whole Sample	EU	Outside EU
	(1)	(2)	(3)	(4)	(5)
# Apps by Developer	-0.00800*** (-3.18)	-0.00246* (-1.75)	-0.00844*** (-3.31)	-0.00932*** (-3.01)	-0.00895*** (-5.82)
After GDPR=1 $\times$ Nbr App. Dev.			-0.000677*** (-9.21)	-0.000184** (-2.01)	-0.00114*** (-10.61)
App Characteristics	Yes	Yes	Yes	Yes	Yes
Apps FE	Yes	Yes	Yes	Yes	Yes
Time FE	Yes	Yes	Yes	Yes	Yes
Mean Dependent	0.561	0.658	0.589	0.509	0.623
Adjusted $R^2$	0.943	0.978	0.942	0.943	0.942
Observations	1,068,913	429,405	1,498,645	454,028	1,044,617

*Notes:* The table presents results from OLS regressions with app and week fixed effects. The dependent variable *Sensitive Data* includes four broad categories of data: *Sharing*, *Location Data*, *Identity Information*, and *User Surveillance*. Each column indicates the type of subsample used in the estimate. Column (1) shows the effect before GDPR implementation. Column (2) shows the effect after GDPR implementation. Column (3) includes an interaction term for Post GDPR and the number of apps by a developer. Columns (4) and (5) separate the effects for EU developers and Non-EU developers, respectively. Robust standard errors are clustered at the app level and reported in parentheses. Significance levels:  $*p < .10$ ,  $**p < .05$ ,  $***p < .01$ .