

# Sticky Consumers and Cloud Welfare \*

Chuqing Jin<sup>†</sup>

Sida Peng<sup>‡</sup>

Peichun Wang<sup>§</sup>

September 14, 2023

**ABSTRACT:** We estimate welfare benefits of the public cloud and study the impact of customer inertia on welfare. We develop a novel demand model that allows for both multiple product choices and continuous usage, and estimate the model using proprietary customer-level data. We find the average consumer surplus from cloud usage to be 216% of its cost, and that smaller customers disproportionately benefit from public cloud. We also find significant inertia on the cloud, reducing welfare benefits by 62%. Finally, we show that cloud migration services and introductory discounts can improve both consumer surplus and provider revenue.

*Key words:* cloud computing, inertia, demand estimation, multiple discrete-continuous choices

---

\*We thank Juan Camilo Castillo, Ulrich Doraszelski, Marc Rysman, numerous colleagues at Microsoft, and seminar participants at IIOC, Boston University, Nanyang Technological University and Zhejiang University for helpful comments. All remaining errors are our own.

<sup>†</sup>Toulouse School of Economics, Email: [cjin@bu.edu](mailto:cjin@bu.edu)

<sup>‡</sup>Office of the Chief Economist, Microsoft, Email: [sida.peng@microsoft.com](mailto:sida.peng@microsoft.com)

<sup>§</sup>Unity Technologies, Email: [willw@unity.com](mailto:willw@unity.com)

# 1 Introduction

The public cloud, i.e., third-party cloud computing services, has become an increasingly important part of the economy. In 2021, the Infrastructure-as-a-Service (IaaS) market alone generated \$91 billion in revenue, and grew further to \$115 billion in 2022 (Gartner 2022, 2023). The rest of the economy is also moving to the cloud, making it a critical infrastructure of the digital economy.<sup>1</sup> A key benefit of the cloud is that it allows firms to access computing resources without owning physical hardware, enabling them to easily adjust their usage of computing products to suit their needs. However, in practice, cloud customers may develop inertia on existing products and fail to make these product adjustments optimally, which undermines the cloud’s benefit.<sup>2</sup> Hence, despite its rapid growth, the cloud’s overall welfare impact remains an open question.

In addition, the cloud may have a distributional effect across small and large firms. Small firms, with budget constraints and demand uncertainty, are less likely to own computing hardware that are sufficient for peak demand and have the latest technology. By lowering the upfront setup costs, the cloud may disproportionately benefit small firms and alleviate concerns over industry concentration and dynamism in the economy (Bloom and Pierri 2018).

In this paper, we estimate both the overall and distributional welfare impact of cloud usage, taking into account the impact of inertia. To that end, we develop a structural model of cloud demand and estimate it using customer-product-month level data from a major U.S.-based cloud provider. We then conduct counterfactuals to compare the welfare benefit of cloud usage with and without inertia, as well as potential remedies for inertia.

The challenge in modeling cloud demand is that cloud customers use multiple products simultaneously and incur continuous usage on each product, deviating from classical discrete choice models. Our proposed model keeps this demand pattern for welfare evaluation. We posit that in each period, each customer draws multiple computing tasks exogenously. Then, the customer chooses the best product for each task based on the task size, i.e., the amount of computing resources needed to complete the task, as well as prices and characteristics of the available products. We allow customer

---

<sup>1</sup>92% of organization’s IT environment is at least partially on the cloud. See <https://www.idg.com/tools-for-marketers/2020-cloud-computing-study/>, accessed Dec. 12, 2021.

<sup>2</sup>Inertia is observed in many other low-maintenance service markets such as mobile, broadband, cash savings, insurance, and mortgages. See [https://www.citizensadvice.org.uk/Global/CitizensAdvice/Consumerpublications/Super-complaint-Excessivepricesfordisengagedconsumers\(1\).pdf](https://www.citizensadvice.org.uk/Global/CitizensAdvice/Consumerpublications/Super-complaint-Excessivepricesfordisengagedconsumers(1).pdf), a super-complaint filed to the UK Competition and Markets Authority, accessed May 3, 2023.

tastes to be heterogeneous both at the customer level and at the task level, capturing how different products may be better suited for different customers and tasks.

We show that the model is identified given our customer-product-month level data, i.e., the econometrician only observes the total usage of each customer on each product but not the number of tasks or the task sizes. Three components of the model need to be identified. First, the distribution of the number of tasks is identified from the sparsity of customers’ product choices: if a customer has many tasks, then it is unlikely that we observe zero usage on any product for this customer. Second, given the distribution of the number of tasks, the distribution of the task sizes is identified from each customer’s monthly cloud usage. Finally, customer taste parameters are identified from usage and choice patterns across products and markets.

Our data provides customer-product-month level usage history from November 2015 to June 2018 for 32 months. We focus on virtual machine (VM) products — the main part of IaaS — which are bundles of different combinations of computing resources such as CPU, RAM, storage, and network. The VM products may have one of two operating systems (Windows or Linux) and may be deployed in one of 24 geographic locations (data centers across the world). We define combinations of these geographic regions and operating systems as markets. While our data comes from one cloud provider, products and prices across all major U.S. cloud providers are nearly identical, so our results can be interpreted as representative of the cloud market.

Before introducing our structural estimates, we first document three data patterns that emphasize the importance of inertia in the cloud market. First, cloud technology is rapidly growing and new products are often launched with better performance and lower prices ([Kilcioglu and Rao 2016](#)), but there remains significant demand for the dominated old products. Second, cloud providers sometimes launch promotional versions of existing regular products (with identical characteristics but discounted prices), while the regular products continue to be offered. One such promotional product is launched in our sample. We find that 90% of existing customers from the regular product did not move to the new promotional product in the 12 months after its launch, even though these customers would have saved 22% of their total cloud spending. Third, using customers’ usage histories, we can distinguish between existing customers who are affiliated with the regular product, versus new customers who have to make an active choice between all products. We find that new customers are significantly more likely to adopt the new promotional product than existing customers. These data

patterns combined strongly suggest that existing cloud customers are likely choosing sub-optimal products due to inertia.

Given these data patterns, we model inertia as an adoption cost for new products that customers have not previously used and remain agnostic about its exact mechanism.<sup>3</sup> As [Handel and Schwartzstein \(2018\)](#) point out, identification of the specific mechanism of inertia is often difficult and may be irrelevant for the welfare evaluation of many policies.<sup>4</sup> We study the overall welfare impact of inertia and explore remedies that directly steer customers towards adopting new products.

Our structural estimates reveal several features of cloud demand. First, we find low price elasticities for most of the cloud provider’s products, i.e., substantial price changes are required to change customers’ usage patterns. Second, cloud customers face significant costs in adopting new products offered by the same provider. The average estimated cost of adopting a new product is equivalent to increasing the price of an average product by 12 times for small customers and 7 times for large customers for an average task. Third, small customers are estimated to have both lower number of tasks and smaller task sizes, but their number of tasks grows faster over time. While small customers have higher disutility from prices per unit of usage, they are less price elastic compared to large customers once usage is taken into account.

With our model and estimates, we first simulate customer product choice and usage histories from the beginning of the sample and show that the model and estimates capture the rich heterogeneity in usage and product choices, and fit the data well. We then compute consumer surplus and find large welfare gains from cloud usage. On average, we estimate cloud customers’ return on investment (ROI), defined as the ratio between consumer surplus and their cloud spending (i.e., provider revenue), to be 216%, or roughly 2.2x. The IaaS market’s annual revenue averaged about \$26 billion over the sample period ([Gartner 2017, 2018, 2019](#)). Assuming the same ROI across cloud providers, our estimate suggests average consumer welfare gains from the IaaS market alone to be \$56 billion annually between 2016 and 2018, or \$248 billion in 2022 if further extrapolated.

---

<sup>3</sup>This specification of inertia ignores customers’ switching costs to products they have previously used. We do not model switching costs between products for two reasons. First, switching costs are typically identified under the assumption that each customer chooses one product at a time so switching behaviors are clearly observed, which does not apply in our setting. Second, modeling continuous usage also adds to the difficulty: when a customer’s usage on a product is decreased, we cannot distinguish whether her total demand has decreased or she has switched to another product.

<sup>4</sup>[Handel and Schwartzstein \(2018\)](#) argue that it is sufficient to identify the combined effect of inertia to evaluate policies that strongly steer customers to specific actions (“allocation policies”). Identification of the specific mechanism is needed only for policies that target specific mechanisms (“mechanism policies”).

Furthermore, the average ROI for small customers is 2.7x and for large customers is 2x, consistent with the hypothesis that the public cloud disproportionately benefits small firms.

We then re-compute customers’ cloud usage when there is no inertia, holding the supply side fixed.<sup>5</sup> We find that customers lose 62% of consumer surplus due to inertia. The provider also loses 58% of revenue from customer inertia, due to consumers’ slower adoption of new cloud products, and hence, lower overall cloud usage. To dissect whether the loss in consumer surplus is real welfare loss, we decompose it into a direct adoption cost and an indirect cost from sub-optimal product choices. We find that sub-optimal product choices due to inertia account for 98% of the loss in consumer surplus, whereas the direct adoption cost only accounts for 2%. While the direct adoption cost is mechanism-agnostic and may include real costs needed to adopt new products (thus not welfare relevant), the indirect cost from sub-optimal product choices provides a lower bound for the true welfare loss due to inertia.<sup>6</sup>

In our counterfactual analyses, we explore two remedies to customer inertia. First, we implement a full subsidy for customers to adopt new products: customers choose products as if there were no inertia, and their adoption costs are fully subsidized whenever they choose a new product they have not used before. Consumer surplus under this scenario is equivalent to the one without customer inertia. In practice, the subsidy can be implemented via “white-glove” services that help customers migrate to new cloud products. If such a subsidy is funded by the social planner, we estimate that each dollar spent would generate about 2 dollars in total welfare. If the provider were to fund the subsidy, we find it to be unprofitable for the provider in the short run, which may explain why in practice only limited “white-glove” migration services are offered. In the long run, however, we find it profitable for the provider to subsidize new product adoption costs, due to increased revenue from faster adoption of cloud products.

Another common practice to encourage customers to adopt new products is introductory discounting, e.g., offering discounts for new product launches (“new product preview”) or for customers’ first-month usage on any product (“personalized product trial”). We explore both forms of introductory pricing using the latest product launch in our sample. We further explore the effect of allowing

---

<sup>5</sup>We do not model cloud customers’ choices of cloud providers, which drive providers’ pricing decisions. We also do not observe cloud providers’ hardware purchase and maintenance costs.

<sup>6</sup>This is conditional on customers being myopic, which we assume given the complexity of our demand model. If customers were forward-looking, the adoption cost estimates would be higher and the direct adoption cost would likely account for a higher percentage of the loss in consumer surplus.

the provider to target large and small customers with different discounts. We find that, in the long run, personalized product trial yields 12.8% higher revenue for the provider *and* 16.5% higher consumer surplus than the baseline with no discounts, compared to a 1.3% and 1.9% improvement in revenue and consumer surplus from new product preview, respectively. Moreover, the provider’s optimal discount for personalized product trial is 225%, suggesting that more subsidy than free first-month usage (i.e., 100% discount) can improve provider revenue. Finally, when the provider can offer different discounts, we find that the optimal discount for small customers is higher than that of large customers, contrary to what managers typically do in practice<sup>7</sup> and despite small customers’ lower price elasticity in our setting. This is because smaller customers’ usage grows faster than large customers and it is thus relatively cheaper for the provider to discount small customers upfront and increase their cloud usage in the long run.

This paper is related to a long literature on measuring the benefits of IT adoption. An extensive literature attempts to directly measure the effect of IT adoption on firm performance (see [Brynjolfsson and Hitt \(2000\)](#) for a review). Our paper instead follows an approach in the industrial organization literature to infer firms’ gains from cloud usage by estimating their willingness to pay (i.e., demand curve) with transaction data of cloud products. Using the same approach, [Bresnahan \(1986\)](#) estimates a 3.3x-6.1x ROI from the adoption of mainframe computers in the financial services industry. [Hendel \(1999\)](#) estimates a 0.92x ROI from personal computer (PC) adoptions in the same industry. Our paper is the first to assess welfare benefits of cloud usage. Furthermore, with growing concern over decreasing dynamism in the global economy as well as rising industry concentration ([Decker et al. 2016](#); [De Loecker, Eeckhout and Unger 2020](#)), evidence has pointed to access to information technology (IT) as a potential cause ([Bessen 2020](#); [Tambe et al. 2020](#)). Our analysis covers cloud customers from all sectors, and we study different gains from cloud by small and large firms separately.<sup>8</sup>

Our demand model fits in the literature on discrete choices. The extant literature mostly models each consumer purchasing one unit of a single product ([Berry 1994](#); [Berry, Levinsohn and Pakes 1995](#)). However, in many settings, customers choose more than one product and use multiple units

---

<sup>7</sup>In practice, in the cloud market, large customers typically get higher discounts from cloud providers.

<sup>8</sup>Other studies on the heterogeneous impact of IT adoption on different firms include [Forman \(2005\)](#), [Tambe and Hitt \(2012\)](#), [Jin and McElheran \(2017\)](#), and [Peng et al. \(2021\)](#). In particular, [Jin and McElheran \(2017\)](#), using expenditures on outsourced IT services from the US Census, argue that young plants benefit more from the cloud as it reduces their learning costs.

of each product, e.g., credit card payment, grocery shopping, and firm procurement. To estimate welfare and conduct counterfactual experiments, it is important to retain these demand patterns. A few exceptions in the literature also model multiple discrete choice and usage. [Burda, Harding and Hausman \(2012\)](#) use a Poisson mixture model to model households’ discrete number of supermarket visit counts. Our model, on the other hand, does not require the econometrician to observe either the number of tasks or the size of each task, and handles continuous usage rather than discrete counts. These features allow for a broader range of applications (e.g., credit card payment and firm procurement). An identification challenge is how to disentangle determinants of product choice from those of usage. [Hendel \(1999\)](#) and [Koulayev et al. \(2016\)](#) both rely on observing characteristics that only affect product choice but not usage. In contrast, we notice that if a customer has zero usage on a product, it must be that the product is never chosen, and we use this sparsity pattern in the data to infer the number of discrete choices a customer has to make, hence identifying choice from usage without additional data. A similar sparsity moment is used in [Quan and Williams \(2018\)](#) to identify demand heterogeneity across markets in a single-discrete-choice setting, whereas our paper models multiple discrete choices with continuous usage and thus requires exactly zero usage on some products to identify the number of discrete choices.

Finally, our estimation of inertia is related to the literature on frictions in product choice. Most of this literature focus on consumers ([Hortaçsu and Syverson 2004](#); [Handel and Kolstad 2015](#); [Erdem and Keane 1996](#); [Bartoš et al. 2016](#); [Bhargava, Loewenstein and Sydnor 2017](#); [Hanna, Mullainathan and Schwartzstein 2014](#)), whereas we study firms as customers, who may have additional frictions: when changing between products, they may face engineering costs ([Greenstein 1993](#)) or other types of switching costs ([Burnham, Frels and Mahajan 2003](#)); there may also be organizational slacks ([Cyert and March 1963](#)) such as how budgets are planned and spent ([Liebman and Mahoney 2017](#)). Inertia in firms’ IT procurement may become more prevalent as IT products are offered as services with low maintenance needs, and may also have implications for their downstream competition. In terms of estimation, [Farrell and Klemperer \(2007\)](#) discuss the challenge of identifying inertia from unobserved consumer heterogeneity. With market-level data, the literature typically needs to assume that a consumer is affiliated with a single product at any time ([MacKay and Remer 2022](#)). In contrast, we allow customers to use multiple products simultaneously, and leverage customer-level usage history and product launches to identify inertia.

The rest of the paper is organized as follows. Section 2 gives an overview of the public cloud market and describes our data and evidence of inertia. Section 3 presents our demand model and discusses identification. Section 4 describes our estimation strategy and Section 5 presents the results. Section 6 discusses welfare and conducts counterfactual analyses. Section 7 concludes.

## 2 Data and Descriptive Evidence

In this section, we begin by providing an overview of the public cloud market and our data. We then turn to describing patterns in the data that motivate our model and identification.

### 2.1 The Public Cloud Market

The public cloud broadly includes three sets of markets: Infrastructure-as-a-Service (IaaS, e.g., VMs), Platform-as-a-Service (PaaS, e.g., App Services), and Software-as-a-Service (SaaS, e.g., Office 365). This paper focuses on the IaaS market, specifically VMs, which is the main infrastructure for the rest of cloud services and the digital economy. The IaaS market has four major providers globally, Amazon Web Services, Microsoft Azure, Google Cloud Platform, and Alibaba Cloud, accounting for almost 80% of market share ([Gartner, 2022](#)). These providers build large data centers globally with clusters of computing hardware purchased from upstream component makers such as Intel, AMD, and Nvidia. These hardware, with CPU, RAM, storage, and network bundled together, are then rented out via virtualization technology as VMs on an hourly basis. There are three purchase models of VMs: on-demand, by reservation, and preemptible. We focus on the most popular model: on-demand VMs. As the name suggests, customers can start or stop an on-demand VM at any time with a high level of performance guarantee and will be charged based on the amount of computing resources used.<sup>9</sup>

Cloud customers, mostly firms, purchase and configure VMs to substitute or complement traditional on-premise computing resources (e.g., PCs and private data centers) managed by their IT departments. The cloud allows these firms to minimize upfront setup costs and reduce maintenance costs, as they no longer need to purchase and manage physical machines. Auto-scaling technology

---

<sup>9</sup>In contrast, reservations typically require one year or three years of commitment in exchange for a per-hour discount. Preemptible VMs do not have any performance guarantees, e.g., customers' workloads can be preempted by the provider at any time.



then helps these firms efficiently scale the amount of computing resources they need based on their fluctuating demand. As a result, cloud customers can achieve higher utilization of the computing resources they purchase and afford scaling to higher demand with the latest technology.

## 2.2 Data

Our data comes from a major cloud provider headquartered in the U.S. and includes prices and characteristics of each VM stock-keeping unit (SKU), as well as histories of customer-SKU-month level usage for 32 months from November 2015 to June 2018. For confidentiality, we retain a random sub-sample of customers and re-scale the usage and price data separately by an unknown constant. To initialize usage history for existing customers at the beginning of our sample, e.g., to define whether a product is new to a customer, we hold out the first six months of data in our sample from estimation, leaving us with 26 months in our estimation sample. Finally, we drop customers with less than 6 months of usage, leaving us a total of 3,233 customers who account for 90% of total usage in our sample.

A SKU is defined by the provider hierarchically as follows: (i) region: location of the data center where customers' workloads would be hosted; (ii) operating system (OS): Linux or Windows; (iii) VM series: product category defined by function, e.g., Azure's D series includes SKUs optimized for CPU-intensive applications; (iv) VM family: group of SKUs within a VM series defined by the underlying hardware, e.g., in Azure's D series, the 3rd generation Dv3 family is based on Intel's Broadwell chips and the 2nd generation Dv2 family is based on the older Haswell chips; (v) VM size: SKUs within a VM family with specified number of virtual CPU cores, gigabytes of RAM and storage, and other resources, e.g., the D2 v3 SKU is in the Dv3 family and comes with 2 virtual CPU cores, 8 gigabytes of RAM, and 50 gigabytes of temporary storage.

We define a region-OS combination as a market, because most customers consistently choose SKUs within the same region running the same OS.<sup>10</sup> There are 47 markets in our sample. We define a product at the VM family level, grouping different VM sizes within the same family as one product because they share the same technology and same price per unit of computing resource (e.g., D4 v3 has double the amount of CPU cores and gigabytes of RAM and storage than D2 v3, and is

---

<sup>10</sup>Reasons for little substitution across regions may include data sovereignty regulations, latency, or simply where the bulk of customers' data are stored.

charged twice as much). In the rest of this paper, we use “product” and “VM” interchangeably. We further focus on the most popular VM series (henceforth X for confidentiality), which is designed for compute-intensive workloads based on Intel CPUs, and accounts for 42% of total usage.<sup>11</sup> There are four products in the X series, X1-X4, and one “other” product which we define as the group of VM series other than X. Finally, X1 and X2 exist throughout our sample, whereas X3 and X4 are launched during the sample, with variations in launch dates across markets.

Cloud usage is measured in compute units, which are based on the amount of computing resources used, specifically in terms of core-hours. The number of core-hours used by a VM is calculated by multiplying the VM’s number of CPU cores with the number of hours used. As the performance of each CPU core may differ for different VMs, cloud providers use compute units to normalize core-hours so that usage across different VMs are comparable.<sup>12</sup> We then measure the price of each product in dollars per compute unit. Table 1 presents summary statistics of our sample.

Table 1: Summary statistics

Variable	N	Mean	SD	1st quartile	Median	3rd quartile	Max
Price, per product-market-month	3755	0.565	0.213	0.371	0.573	0.732	1.079
Market share, per product-market-month	3755	0.282	0.273	0.050	0.190	0.481	1
Number of products, per market-month	1209	3.106	1.446	3	3	4	5
Usage, per customer-market-month	96878	0.463	1.447	0	0.120	0.465	75.966

*Notes:* Table shows summary statistics of our estimation sample. Market is defined as a region-OS combination. Products include X1-X4, and the “other” product. Price is measured in re-scaled dollars per compute unit. Market share represents the usage share of each product, excluding the outside option, in each market-month. Usage is measured in re-scaled compute units.

Prices mainly vary across markets and products. Over time, there is only one major price drop that lowered prices for X2 in about one third of the markets at the same time. Across regions, prices vary significantly. The most expensive region has a 56% average price uplift compared to the cheapest region. There is also a price premium for Windows VMs compared to Linux VMs. Across products, the most expensive X-series product is priced 35% higher than the cheapest on average. We leverage the variation across products and regions to identify cloud customers’ price elasticities. To that end, we supplement our data with three cost shifters obtained from the cloud provider: (i) electricity prices across regions; (ii) power consumption ratings across product-regions; and (iii) hardware costs across product-regions (more details in Appendix C).

<sup>11</sup>Focusing only on the X series significantly reduces the computational burden. Moreover, compute-intensive workloads are different from memory-intensive, high performance computing, or GPU-based workloads.

<sup>12</sup>For example, see <https://docs.microsoft.com/en-us/azure/virtual-machines/acu> for Azure compute unit, accessed Dec. 26, 2021.

Finally, we supplement data from the provider with industry reports on the total market size to calculate customers’ computing needs outside of the provider. According to [Cisco \(2018\)](#) and [Gartner \(2017\)](#), public cloud usage accounts for 58% of total cloud usage, which accounts for 83% of total computing workloads worldwide. Multiplying these with the provider’s share in the public cloud market gives us the proportion of computing done by the provider. For each customer in our data, we divide their observed usage by this proportion to calculate their total computing demand.<sup>13</sup>

## 2.3 Customer Usage Patterns

We now turn to patterns in the data that motivate our model and identification. First, we compare trends of cloud usage between large and small customers. We define customer size based on whether a customer’s first six-month’s usage is above or below the median. On average, large customers’ cloud usage is five times that of small customers. Both grow significantly on average during our sample, with large customers growing by 237%, from 0.365 to 1.230, and small customers by 452%, from 0.057 to 0.315. Small customers’ cloud usage grows faster, and by the end of our sample, the average small customer’s cloud usage exceeds 24% of large customers.

Second, there is rich heterogeneity in customers’ cloud product choice and usage. Table 2 shows the frequency of a customer choosing any number of products and incurring different amount of usage in a market-month. There is significant heterogeneity in both the number of distinct products a customer uses as well as their usage. In particular, 16.6% of customers use more than one product in a market-month and these customers account for 45.6% of total cloud usage in sample. Customers with more usage typically use more distinct products, possibly due to a larger number of diverse computing tasks. To estimate welfare from cloud usage and conduct counterfactuals, our model needs to retain these rich patterns of demand.

## 2.4 Evidence of Inertia

In this section, we provide descriptive evidence of inertia on cloud. We proceed in three steps, looking first at the aggregate transition from  $X_1$  to  $X_2$ , then  $X_2$  to  $X_3$ , and finally the impact of usage history on product choices using a regression discontinuity design.

---

<sup>13</sup>If a customer does not have any usage in a month with the provider, we attribute all her usage to her outside option (on premise or other cloud providers), where her total usage is interpolated from her usage before and after that month.

Table 2: Frequency distribution of number of products chosen and usage

Usage / Products	1	2	3	4	5	Total
0-0.1	10316	479	73	3	1	10872
0.1-1	38469	4469	523	17	1	43479
1-10	5002	3613	1130	159	18	9922
10-100	63	97	93	17	3	273
Total	53850	8658	1819	196	23	64546

*Notes:* Table shows the number of customer-market-months tabulated by the number of distinct products chosen and the total usage within the same market-month, excluding customer-market-months with zero products chosen. Usage is measured in re-scaled compute units. Products include X1-X4, and the “other” product.

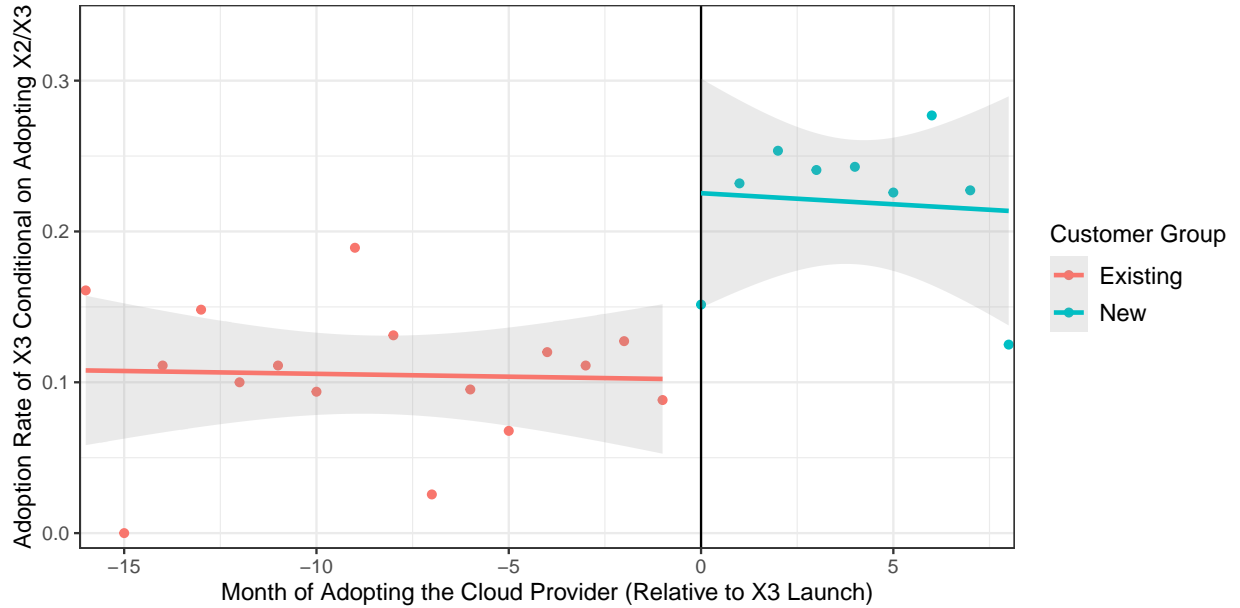
In the cloud market, new products are typically offered with higher performance and lower prices than the previous generation, while the older products continue to attract demand. For our provider, X2 is launched one year after X1 (both before our sample), with better CPU technology and at least 10% lower prices when compared in the same market and measured during our sample. However, we find that, for customers who exist from the beginning of our sample, X1 usage accounts for 22% of total usage at the beginning of our sample, and 11% even at the end of the sample, while X2’s usage share increases from 7% to 18%.

A more direct illustration lies in the transition from X2 to X3. X3 is launched two years after X2 as its promotional version. That is, X3 shares the same technical characteristics as X2, except for its lower prices. Although marketed as a promotional product, X3 is launched with an indefinite end date and ends up being available well beyond the launch of the next product X4. Thus, X3 should be a perfect substitute for X2 and dominates X2 due to its lower prices. However, in the 12 months after the launch of X3, 90% of X2 customers have not adopted X3, and 91% of total usage between X2 and X3 are from X2. Moreover, these choices are not financially trivial: these customers would have saved 22% of their spending on X2 (or 7% of their total cloud spending) had they converted their usage from X2 to X3. These facts suggest that customers may face significant inertia in adopting new products.

One may argue that a promotional product can be perceived as lower quality by customers, and thus the low adoption rate cannot be attributed to inertia. To alleviate this concern, we further compare adoption rates of X3 between existing X2 customers—those who have used X2 before the launch of X3—and new customers to the provider after X3’s launch, using a regression

discontinuity (RD) design. Figure 1 presents the comparison in a standard RD plot. The running variable is each customer-cohort's month of first adopting the provider and centered around X3's launch month. Customer cohorts to the right of the vertical line are new customers after X3's launch, whereas customers to the left are existing X2 customers. We calculate the outcome variable as the probability of adopting X3, conditional on adopting X2 or X3, in the 12 months after X3's launch. We find that new customers, who are otherwise similar but do not have a natural affiliation with any existing products and thus face similar adoption costs for all products, have a significantly higher probability of adopting X3 compared to existing customers who have used X2.<sup>14</sup> Moreover, the overall adoption rate of X3 is low for all customers, suggesting that the perceived quality of the promotional product may indeed be low.

Figure 1: X3 adoption rates: New vs. existing customers



*Notes:* Regression discontinuity plot comparing adoption rates of X3 between existing X2 customers and new customers to the provider after X3's launch. The running variable is each customer-cohort's month of first adopting the provider relative to X3's launch month. The outcome variable is each customer-cohort's probability of adopting X3, conditional on adopting X2 or X3, in the 12 months after X3's launch. Lines are linear regressions of the adoption rates on the running variable. Gray shaded area represents 95% confidence interval.

<sup>14</sup>We do not observe any significant change in the total number of new customers or those who use X2 or X3 after the launch of X3, so it is unlikely that the new customers are different from existing customers due to X3's launch.

### 3 Model and Identification

In this section, we first present a hierarchical demand model that allows for multiple product choices and continuous usage on each product. Then, we discuss the identification of this model. For tractability, the model is static. We refer to market-months simply as markets and denote them by  $m$ . Customers are indexed by  $i$ , and model parameters with subscript  $i$  vary by customer size to capture their usage differences and welfare implications.

Each month, cloud customer  $i$  faces a random number of computing tasks  $n_{im}$ , drawn i.i.d. from a Poisson distribution, i.e.,  $n_{im} \sim \text{Poisson}(\lambda_{im})$ .  $\lambda_{im}$  is the expected number of tasks in market  $m$  and has a time trend to capture growing computation needs, i.e.,  $\lambda_{im} = \lambda_{i0} + \lambda_{i1}t_m$ . Each task is endowed with task size  $q_{im}$  drawn from an i.i.d. exponential distribution, i.e.,  $q_{im} \sim \text{Exp}(\gamma_i)$ , where  $\gamma_i$  is the scale parameter and is equal to the expected task size. The task size represents the amount of computing resources needed to complete the task, which is in compute units and thus the same regardless of the product it runs on. Different from [Burda, Harding and Hausman \(2012\)](#), we allow usage to be continuous, so our model applies to settings with any usage-based services.

In the rest of this section, we first describe customer utility and product choice at the task level, and then discuss the aggregation to the product level.

For each task, the customer chooses between  $J$  available VMs and the outside option of either computing on premise or using other cloud providers. The customer's utility from using VM  $j$  in market  $m$  for task  $k$  is:

$$u_{ijmk} = \alpha_i p_{jm} q_{im} + X_j \beta_i + \delta_i \text{New}_{ijm} + \xi_{ijm} + \epsilon_{ijmk}, \quad (1)$$

where  $p_{jm}$  is the per compute unit price for VM  $j$  in market  $m$  and  $\alpha_i$  is the price coefficient;  $\alpha_i p_{jm} q_{im}$  measures the customer's dis-utility from the cost of the task.  $X_j$  is a vector of VM  $j$ 's observable characteristics and  $X_j \beta_i$  measures the customer's utility from these characteristics. In our setting,  $X_j$  is comprised of dummies for each product in the X series and the "other" product. Because X3 shares the same technical characteristics as X2, we specify the utility for X3 as the sum of the X2 dummy plus an indicator  $\text{Promo}_j$ , which equals 1 only for X3.

To model inertia, we define  $\text{New}_{ijm}$  as a new product indicator that equals 1 if customer  $i$  has

never used product  $j$  in market  $m$ , i.e., the product is new to the customer.<sup>15</sup> This specification agnostically captures inertia that prevents customers from adopting new products. Depending on the specific source of inertia,  $\delta_i$  may be interpreted as information frictions, organizational slacks, or real engineering and operational costs associated with adopting a new IT product. In our model,  $\delta_i$  captures the joint effect of these different types of inertia, and we call it the adoption cost for new products.

$\xi_{ijm}$  is an unobservable demand shock for customer  $i$  using VM  $j$  in market  $m$ . While there is little unobservable characteristics at the VM  $j$  level,  $\xi_{ijm}$  captures any demand shocks that vary at the customer-VM-market level that are unobservable to the econometrician. The unobservable demand shocks may thus create endogeneity in prices, which we account for using the control function approach following [Petrin and Train \(2010\)](#) in the next section.

Finally,  $\epsilon_{ijmk}$  is an idiosyncratic task-level preference shock, distributed i.i.d. type-I extreme value, for all  $j = 0, 1, \dots, J$ . If the outside option is chosen, the customer receives utility  $u_{i0mk} = \epsilon_{i0mk}$ . We normalize the mean utility of the outside option to be zero, so  $u_{ijmk}$  measures incremental utility from this provider's VM  $j$  relative to the outside option.

Putting it together, the task-level choice probability is given as follows. Let  $d_{ijmk}(q_{im})$  be a dummy variable, which equals one if VM  $j$  is chosen for task  $k$ , i.e.,  $u_{ijmk} > u_{ilmk}$  for all  $l \neq j$ . Then, the probability of choosing VM  $j$  for task  $k$  of size  $q_{im}$  is given by

$$\mathbb{P}(d_{ijmk}(q_{im}) = 1) \equiv P_{ijm}(q_{im}) = \begin{cases} \frac{\exp(\alpha_i p_{jm} q_{im} + X_j \beta_i + \delta_i New_{ijm} + \xi_{ijm})}{1 + \sum_{l=1}^J \exp(\alpha_i p_{lm} q_{im} + X_l \beta_i + \delta_i New_{ilm} + \xi_{ilm})} & \text{if } j \neq 0 \\ \frac{1}{1 + \sum_{l=1}^J \exp(\alpha_i p_{lm} q_{im} + X_l \beta_i + \delta_i New_{ilm} + \xi_{ilm})} & \text{otherwise.} \end{cases} \quad (2)$$

### 3.1 Identification

To build intuition, we first discuss the identification of a simplified model without inertia, customer heterogeneity, or time trend in the number of tasks. Then, we conclude this section by adding them back. In the simplified model, there are at most  $J + 3$  parameters to identify: at most  $J$  parameters ( $\beta$ 's) for product characteristics, one price coefficient ( $\alpha$ ), and two parameters from the task generation process ( $\lambda, \gamma$ ).

---

<sup>15</sup>For tractability, we model inertia at the task level instead of customer level. If inertia is modeled at the customer level, a customer's product choice for one task will depend on her choices for other tasks, greatly complicating the choice problem.

Denote  $y_{ijm}$  as the total usage for customer  $i$  on product  $j$  in market  $m$  across all her tasks, i.e.,  $y_{ijm} = \sum_{k=1}^{n_{im}} q_{im} d_{ijk}(q_{im})$ . To identify the model, we require that for any customer  $i$ , her total usage on every product  $(y_{i0m}, y_{i1m}, \dots, y_{iJm})$  is observed. This requirement is natural in most markets with continuous usage and more flexible than [Burda, Harding and Hausman \(2012\)](#), as we do not require the number of tasks or the size of each task to be observable. Unlike [Burda, Harding and Hausman \(2012\)](#), the likelihood function for our hierarchical model cannot be easily constructed. This is because our observable  $y_{ijm}$  is a continuous variable and, without observing task sizes, its density depends on an infinite sum of convoluted integration (see Appendix A for details). Therefore, we instead use generalized method of moments (GMM) and consider the following two sets of moments for identification.

**Proposition 1** (Zero Usage Probability and Expected Usage Moments). *The probability of customer  $i$  having zero usage on product  $j$  in market  $m$  is given by*

$$\mathbb{P}(y_{ijm} = 0) = \sum_{n_{im}=0}^{\infty} \underbrace{\frac{\exp(-\lambda)\lambda^{n_{im}}}{n_{im}!}}_{Pr(\text{Number of Tasks})} \cdot \int_{q_{im}} \underbrace{(1 - P_{ijm}(q_{im}))^{n_{im}}}_{Pr(\text{Product Choice})} \underbrace{\frac{1}{\gamma} \exp(-\frac{1}{\gamma} q_{im})}_{Pr(\text{Task Size})} dq_{im}; \quad (3)$$

and the expected usage is given by

$$\mathbb{E}(y_{ijm}) = \underbrace{\lambda}_{\mathbb{E}(\text{Number of Tasks})} \underbrace{\int_{q_{im}} q_{im} P_{ijm}(q_{im}) \frac{1}{\gamma} \exp(-\frac{1}{\gamma} q_{im}) dq_{im}}_{\mathbb{E}(\text{Task Size})}. \quad (4)$$

*Proof.* See Appendix B. □

In Proposition 1, the first moment  $\mathbb{P}(y_{ijm} = 0)$ , the probability of zero usage on each product, captures the sparsity of customers' product choices and helps identify the number of tasks. To see how it is calculated, for any given number of tasks  $n_{im}$  and task size  $q_{im}$ ,  $\mathbb{P}(y_{ijm} = 0)$  has three components, as shown in equation (3). First,  $Pr(\text{Number of Tasks})$  is the probability of receiving  $n_{im}$  tasks. Second,  $Pr(\text{Product Choice})$  is the probability of not choosing product  $j$  in any of these  $n_{im}$  tasks of size  $q_{im}$ . Third,  $Pr(\text{Task Size})$  is the probability of receiving a task of size  $q_{im}$ . Multiplying the three probabilities gives the probability of zero usage on a product for a given number of tasks and task size. We then integrate over all task numbers and sizes to obtain the zero usage moment.



To better understand why the sparsity of customers' product choices help identify the number of tasks, consider an extreme case of zero usage on *all* products. The probability that customer  $i$  in market  $m$  has zero usage on all products is given by

$$\mathbb{P}(\sum_{j=0}^J y_{ijm} = 0) \equiv \mathbb{P}(n_{im} = 0) = \exp(-\lambda).$$

If the sample analogue of this probability is non-zero, we can invert it to back out the average number of tasks  $\lambda$  directly. However, because we do not observe customers' total computation needs (including outside option) in our data, we do not use the probability of zero usage across all products  $\mathbb{P}(\sum_{j=0}^J y_{ijm} = 0)$  as a moment. Instead, we use the probability of having zero usage on each product  $\mathbb{P}(y_{ijm} = 0)$  as moments.

The second moment in Proposition 1 is the expected usage on each product  $\mathbb{E}(y_{ijm})$  and helps identify task size. As shown in equation (4), product  $j$ 's expected usage is calculated as the product of the expected number of tasks and the expected task size if product  $j$  is chosen. To see how each product's expected usage helps identify task size, we again take the sum of this moment across all products for customer  $i$  in market  $m$ :

$$\sum_{j=0}^J \mathbb{E}(y_{ijm}) \equiv \mathbb{E}(y_{im}) = \lambda\gamma,$$

which depends only on the task generation process. So if the average number of tasks  $\lambda$  is identified (by the first moment), this sum pins down the average task size  $\gamma$ .

Besides the task generation process, the zero usage probability and the expected usage on each product also help identify taste parameters  $\alpha$ 's and  $\beta$ 's in the utility function via the market share function. Together, these two sets of moments jointly identify all parameters in the simplified model. To see that, in a market with  $J$  products and an outside option, Proposition 1 gives  $J + 1$  moment conditions each for the zero usage probability and expected usage moments. These  $2(J + 1)$  moment conditions are sufficient to parametrically identify the at most  $J + 3$  parameters of the simplified model for any  $J \geq 1$ .

Finally, going back to our full model, in order to identify inertia, customer heterogeneity, and time trend in the number of tasks, we construct the zero usage probability and expected usage

moments conditional on customer and market characteristics. We discuss these conditional moments in detail in the next section.

## 4 Estimation

Following the identification argument, we estimate the model with GMM. We begin by describing all the moments used in estimation. Then, we discuss how we account for potential price endogeneity. Finally, we discuss our estimation procedure.

### 4.1 Moments

We start with the two moments, zero usage probability  $\mathbb{P}(y_{ijm} = 0)$  and expected usage  $\mathbb{E}(y_{ijm})$ , from Proposition 1. We condition both moments on customer and market characteristics to further identify heterogeneity in the model. Table 3 shows all the moment conditions we use in estimation.

Table 3: Moment conditions

Moments	Characteristics conditioned on	# Moment conditions
Zero usage probability	- region, customer size, choice set	178
	- $New_{ijm}$ (for $j \neq 0$ )	20
Expected usage	- region, customer size, choice set	222
	- $New_{ijm}$ (for $j \neq 0$ )	20
	- month ( $t_m$ ), customer size	52

We first condition both moments on each combination of region, customer size, and choice set. Conditioning on regions helps identify the taste parameters, in particular the price coefficient, as prices vary mostly across markets. The substitution patterns across regions also identify the distribution of the unobservable demand shocks  $\xi_{ijm}$ . Conditioning on customer size identifies heterogeneity between small and large customers. Choice sets reflect variation in the products available to customers across markets. For example, the same market’s choice set is different after a product launch, so we compute moments before and after separately. Given different product availability across markets, there are eight unique choice sets in total. We drop the zero usage probability moment for the outside option.<sup>16</sup> We also drop moments with few observations or

<sup>16</sup>We only observe aggregate shares of outside option usage from the industry reports and apply them to each customer in each market. As a result, the outside option always has positive usage in our data.

outliers.<sup>17</sup>

Second, we condition both moments on whether a cloud product (excluding the outside option) is new to a customer in a market. These moments reflect differences in choice probabilities and usage with and without inertia, thus identifying the adoption cost parameters.

Finally, we condition the sum of the expected usage moment across all products on customer size and each month in our sample, identifying the time trend in the number of tasks for small and large customers separately.

## 4.2 Price Endogeneity

A common identification challenge in demand estimation is that prices may be correlated with unobserved demand shocks, i.e., in our setting,  $p_{jm} \not\perp \xi_{ijm}$ . To account for potential price endogeneity, we use the control function approach with cost shifters as instruments.<sup>18</sup>

Let  $z_{jm}$  denote cost shifters at the product-market level, which affect prices of the corresponding product-markets and are independent of demand shocks  $\epsilon_{ijmk}$  and  $\xi_{ijm}$ . Following [Petrin and Train \(2010\)](#), we assume price is additive in its observed and unobserved covariates, i.e.,

$$p_{jm} = W(z_{jm}, \zeta) + \theta_{jm}, \quad (5)$$

where  $\theta_{jm}$  is the unobserved covariate that is independent of  $z_{jm}$  but correlated with  $\xi_{ijm}$ . We write  $\xi_{ijm} = \kappa_i \theta_{jm} + \tilde{\xi}_{ijm}$  and substitute  $\xi_{ijm}$  in equation (1). Then, rewriting  $\theta_{jm}$  as  $p_{jm} - W(z_{jm}, \zeta)$  following equation (5), we can derive the utility function with control function as

$$u_{ijmk} = \alpha_i p_{jm} q_{im} + \beta_i X_j + \delta_i New_{ijm} + \kappa_i (p_{jm} - W(z_{jm}, \zeta)) + \tilde{\xi}_{ijm} + \epsilon_{ijmk}, \quad (6)$$

where  $p_{jm} - W(z_{jm}, \zeta)$  is the control function, and  $\tilde{\xi}_{ijm}$  is an independent demand shock.

---

<sup>17</sup>We drop moments with fewer than 200 observations. We also drop one expected usage moment for one region-product where all but 3 out of 311 observations are zeroes.

<sup>18</sup>The characteristics being conditioned on in Table 3 are sometimes called GMM “instruments”. To avoid confusion, we only refer to the control function instruments as instruments.

### 4.3 Estimation Procedure

As  $W(z_{jm}, \zeta)$  is unknown to the econometrician, we estimate the model in two stages. In the first stage, we regress prices  $p_{jm}$  on the instruments  $z_{jm}$ , and take the residuals,  $\hat{\theta}_{jm} = p_{jm} - W(z_{jm}, \hat{\zeta})$ , as estimates for the unobserved covariate. Our first two instruments include electricity prices, varying across regions, and power consumption ratings, varying across products and regions because VMs are supported on different hardware in different regions. Besides, we also include the provider's hardware procurement costs as the third instrument, which also vary across products and regions. Detailed construction of the instruments and results of the first stage are presented in Appendix C.

We assume that the part of the unobserved demand shock that is independent of prices,  $\tilde{\xi}_{ijm}$ , is normally distributed with mean zero and standard deviation  $\sigma_j$ . We allow the standard deviation to vary across products to capture different degrees of demand heterogeneity. This is similar to adding a random coefficient to each product dummy  $X_j$ .

In the second stage, we treat  $\hat{\theta}_{jm}$  as a product characteristic according to equation (6). We also exclude every customer's first-month data from estimation given our focus on cloud customers' product choices and usage rather than their decisions on adopting the provider.

The rest of the second stage follows the standard GMM procedure using moment conditions listed in Table 3. We first compute these moments in the data. Then, for any set of parameter values, we use numerical integration to compute the model moments. Finally, we run a quasi-Newton algorithm to minimize the weighted sum of the squared differences between data moments and model moments. Details of the numerical integration and the optimization algorithm are presented in Appendix D. The weighting matrix is a diagonal matrix containing the inverse of the moments' sample variance. We compute standard errors using the GMM asymptotic variance and taking the square root of its diagonal. In Appendix E, we show convergence of this estimator in finite samples with Monte Carlo simulations.

## 5 Results

Table 4 presents our estimated demand parameters. Parameters governing product choice (mean coefficients and standard deviation of the random coefficients) are presented separately from parameters governing task generation (number of tasks and task size). The two columns show estimates

separately for small and large customers.

Table 4: Demand estimates

Variables		Coefficient estimates	
		Small	Large
<b>Product choice</b>	X1 ( $\beta_i^1$ )	-4.324 (0.063)	-3.460 (0.041)
	X2/X3 ( $\beta_i^2$ )	-3.884 (0.059)	-3.174 (0.040)
	X4 ( $\beta_i^4$ )	-4.020 (0.089)	-3.134 (0.059)
	Other ( $\beta_i^O$ )	-2.609 (1.472)	-2.343 (1.466)
	Promo ( $\beta_i^P$ )	-0.850 (0.072)	-0.303 (0.039)
	Price $\times$ task size ( $\alpha_i$ )	-7.857 (0.947)	-4.226 (0.264)
	Adoption cost ( $\delta_i$ )	-5.535 (0.077)	-6.053 (0.118)
	Control function ( $\kappa_i$ )	3.359 (0.346)	0.911 (0.162)
	X1 ( $\sigma_1$ )		1.149 (0.061)
	X2 ( $\sigma_2$ )		1.018 (0.034)
<b>Random coefficient</b>	X3 ( $\sigma_3$ )		1.311 (0.095)
	X4 ( $\sigma_4$ )		1.381 (0.075)
	Other ( $\sigma_O$ )		0.002 (8.620)
<b>Number of tasks</b>	Intercept ( $\lambda_{i0}$ )	43.135 (1.388)	80.848 (2.906)
	Time trend ( $\lambda_{i1}$ )	1.846 (0.141)	0.460 (0.127)
<b>Task size</b>	Exponential scale ( $\gamma_i$ )	0.106 (0.003)	0.371 (0.008)
<b>Observations</b>		172,223	167,244

*Notes:* Table presents estimates of demand parameters. Standard errors are in parentheses and computed using estimated GMM asymptotic variance.

The estimates of the product dummies are in line with the evolution of the cloud products' technical specifications. Quality increases from X1 to X2 and X3 significantly due to significant

improvements of the underlying CPU technology. On the other hand, the CPU improvement from X2 and X3 to X4 is on a minor Intel release, resulting in similar performance between these products.<sup>19</sup> Moreover, X4 is launched with hyperthreading, which reduces its cost but slightly hurts its performance. Also, consistent with the earlier descriptive evidence, the perceived quality of a promotional product is significantly lower. Note that our estimated product dummies are all negative, because during our sample period, cloud usage still only accounts for a small fraction of total computing needs. Finally, the estimated standard deviations of the random coefficients are statistically significant and sizeable for X1-X4, suggesting significant heterogeneity in preferences for the X-series products across customers and markets.

The price coefficients are estimated to be negative and statistically significant, i.e., cloud customers dislike higher prices. The control function coefficient is also statistically significant, suggesting that prices are indeed endogenous and positively correlated with the unobserved demand shocks (first stage results reported in Appendix Table C1). Conditional on the same amount of computing needs, we find small customers to be more price sensitive than large customers. Taking into account differences in computing needs, however, yields higher price elasticities for large customers (Table 5). The own-price elasticity for VM  $j$  is computed as the change in expected usage on VM  $j$  in response to a unit change in its price, i.e.,  $\frac{\partial \mathbb{E}(y_{ijm})}{\partial p_{jm}} \frac{p_{jm}}{\mathbb{E}(y_{ijm})}$  (detailed formula in Appendix F). The estimated own-price elasticities are generally inelastic, which are lower than most consumer products, though slightly higher than estimates of short-run elasticity for electricity (Deryugina, MacKay and Reif 2020). Our estimates of low price elasticities are consistent with other studies in emerging digital markets, where providers are often concerned with expanding the overall market size and long-run growth objectives and do not increase prices despite inelastic demand (Castillo 2020).

We find that customers face significant inertia when adopting a new product. The estimated adoption cost is equivalent to increasing the average product’s price for an average task size by 12 and 7 times for small and large customers, respectively, or 13 and 21 times the significant quality improvement from X1 to X2 and X3.<sup>20</sup> These large adoption costs prevent cloud customers from

<sup>19</sup>Intel has a tick-tock model of CPU innovation, alternating between major CPU improvements (“tock”) and minor cost-reducing die shrinks (“tick”) in its product launches. See <https://www.intel.com/content/www/us/en/silicon-innovations/intel-tick-tock-model-general.html>, accessed on Feb. 28, 2021.

<sup>20</sup>The equivalent price increases are calculated as  $\frac{\delta_i}{0.56 \cdot \alpha_i \gamma_i}$ , where 0.56 is the average price across product-markets in the data and  $i$  varies between small and large. The equivalent quality increases are calculated as  $\frac{\delta_i}{\beta_i^2 - \beta_i^1}$ .

Table 5: Own-price elasticities

Product	Avg.	Small customers	Large customers
X1	-0.824	-0.671	-0.988
X2	-0.732	-0.585	-0.887
X3	-0.629	-0.487	-0.766
X4	-0.744	-0.588	-0.887
Other	-0.941	-0.788	-1.102

*Notes:* Table presents estimated own-price elasticities for different products and customer sizes.

adopting new products, which are typically more powerful and cheaper, thus limiting the potential welfare benefits that could be realized from cloud usage. We illustrate more on this point in the counterfactual section. Finally, small customers face lower adoption costs than large customers when adopting a new product.

Turning to the task generation process, we find that small customers have lower number of tasks as well as smaller task sizes than large customers. On average, at the beginning of our estimation sample in May 2016, we estimate small customers have 43 tasks per month while large customers have 81 tasks. These tasks represent the number of distinct computing projects for a customer in a region-OS per month. The number of tasks for small customers, however, grows about 4 times as fast as that of the large customers over time, capturing faster growth rates of small cloud customers. By June 2018, at the end of our sample, the average small customer has 89 tasks per month while the average large customer has 92 tasks, suggesting that customer size differences are shrinking on the cloud. Small customers also have smaller task sizes than large customers, on average 71% smaller in compute units for each task.

## 5.1 Model Fit

To show how our model fits the data, we initialize customers' product adoption history using each customer's first month usage in our data, and simulate the model forward for each customer with the rest of the data. Importantly, whether a product is new to any customer ( $New_{ijt}$ ) after the first month is simulated within the model. We then compare the simulated usage with usage in the data.

The simulation proceeds as follows: (i) take each customer's first month's data as given and initialize which products the customer has used; (ii) starting in the second month, for every customer

in every market, draw their number of tasks and task sizes from the estimated distributions; (iii) for each task, draw i.i.d. type-I extreme value shocks and random coefficients for each product, and simulate product choice based on estimated utilities; (iv) aggregate usage for each customer-product, and every time a customer uses a new product, update the new product indicator in her utility function; (v) iterate until the end of the sample.

Figure 2 shows that our estimated model fits the data well. For each simulation, we calculate the distribution of average market-month usage for each product, and then average across 200 simulations. The simulated usage captures well the variation in usage across products and markets. For example, the simulated usage tracks the low usage of X4 in the data despite its high estimated utility, suggesting that our model captures the effect of customer inertia on newer products. The model also predicts well zero usage for each product in the sample. Moreover, Appendix Figure H3 shows that we fit the aggregate time series of total computing demand well, and Appendix Figure H4 shows the fit of cloud usage over time by product and customer size.

## 6 Counterfactuals

With our model and estimates, we conduct two sets of counterfactual analyses. First, we compute the welfare benefits of cloud usage under current market conditions and compare it to a model where customer inertia is eliminated. Second, we explore two potential remedies that incentivize customers to overcome inertia and adopt new products.

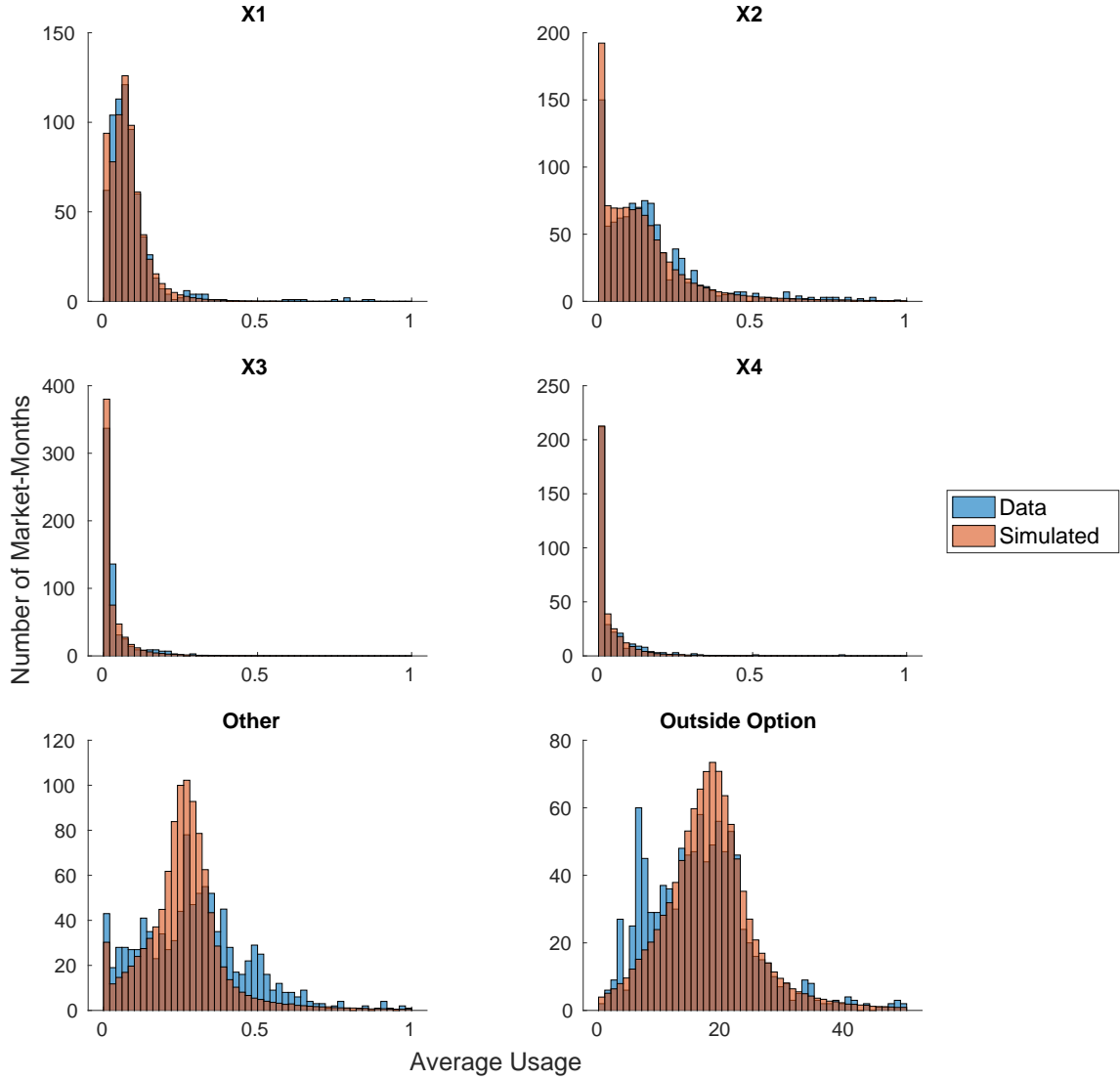
### 6.1 Welfare

Using our model and estimates, we first compute the welfare benefits of cloud usage for customers during our sample with the estimated inertia. To do that, we follow the same procedure described in Section 5.1 when evaluating model fit by simulating each customer’s product choices and usage forward over time. We divide the utility received for each task by the price coefficient to convert it into dollars, and then aggregate across all tasks to obtain consumer surplus. Provider revenue is calculated as the total usage on each product multiplied by the corresponding prices. We repeat these simulations 200 times for each customer and compute the averages.

The first row in Table 6 shows our estimates of welfare benefits from cloud usage. The dollar



Figure 2: Model fit: Distribution of average market-month usage for each product



*Notes:* Figure shows histograms of average usage in a market-month for each product based on data vs. simulations. The blue histograms are based on the data, and the orange histograms are based on simulations. Each bar represents the number of market-months that fall within the corresponding average usage for that product. Results are averaged across 200 simulations.

amounts are on a normalized scale due to data anonymization, so we focus on interpreting cloud customers’ return on investment (i.e.,  $ROI \equiv \text{consumer surplus} / \text{provider revenue}$ ). We find large welfare gains from cloud usage. The average ROI for all cloud customers is 216% ( $= 51.6/23.9$ ) or roughly 2.2x. Given the average annual revenue of the IaaS market of \$26 billion over the sample period from 2016 to 2018, assuming the same ROI across all cloud providers, this estimate suggests an average annual \$56 billion of consumer welfare gains from the IaaS market alone, or \$248 billion if extrapolated to 2022. Small customers have higher ROIs (2.7x) than large customers (2.0x), even though their total consumer surplus is lower given their lower total usage. This is consistent with the hypothesis that cloud usage disproportionately benefits small firms by allowing them to easily scale their computing needs and access the latest technology.

Table 6: Welfare benefits of cloud usage and cost of inertia

		<b>Consumer Surplus (10<sup>3</sup>)</b>		<b>Provider Revenue (10<sup>3</sup>)</b>		<b>Customer ROI (CS/Revenue)</b>		<b>Cloud Usage (10<sup>3</sup>)</b>	
		(Small)	(Large)	(Small)	(Large)	(Small)	(Large)	(Small)	(Large)
Welfare: With inertia	By size	13.2 (0.1)	38.4 (0.3)	5.0 (0.04)	18.9 (0.1)	2.7 (0.02)	2.0 (0.01)	6.7 (0.1)	29.2 (0.2)
	Total	51.6 (0.3)		23.9 (0.1)		2.2 (0.01)		35.8 (0.2)	
Counterfactual: Without inertia	By size	33.4 (0.2)	103.5 (0.5)	11.0 (0.1)	45.9 (0.2)	3.0 (0.02)	2.3 (0.01)	17.5 (0.1)	82.2 (0.4)
	Total	136.9 (0.5)		56.9 (0.2)		2.4 (0.01)		99.7 (0.4)	
$\Delta$ Welfare (without - with)	By size	20.2 (0.2)	65.1 (0.6)	6.0 (0.1)	27.0 (0.2)	0.3 (0.03)	0.3 (0.02)	10.8 (0.1)	53.1 (0.5)
	Total	85.3 (0.6)		33.0 (0.3)		0.2 (0.02)		63.9 (0.5)	

*Notes:* Table shows consumer surplus, provider revenue, cloud customer ROI, and total cloud usage in the estimated model and in a model where inertia is eliminated. Consumer surplus and provider revenue are shown in re-scaled dollars. Customer ROI is calculated as consumer surplus divided by provider revenue. Cloud usage is shown in re-scaled compute units. Results are based on 200 simulations, with standard deviations in parentheses.

Our estimated ROI from cloud usage is comparable to ROI estimates of other computing technologies: it is higher than that of PCs (0.92x, [Hendel 1999](#)) and lower than that of mainframes (3.3-6.1x, [Bresnahan 1986](#)). Like previous studies, our estimate should be taken as an upper bound because it only reflects the welfare of using cloud for existing cloud customers — there may be other costs (e.g., labor training) associated with transitioning to cloud that are not accounted for in our model. Moreover, cloud providers, similar to providers in other digital markets, often optimize for consumer surplus for long-run growth rather than short-run profits ([Castillo 2020](#)), and thus early cloud adopters may extract more surplus.

We then re-do the simulations with our estimated model but with the adoption cost eliminated

to evaluate the welfare benefits of cloud usage in a world without inertia, while holding the supply side fixed. The second row in Table 6 shows the results. The third row then shows the difference between the first and the second row, illustrating the welfare cost of inertia. Consumer surplus is 62% ( $= 85.3/136.9$ ) lower than what it could be if there were no inertia. Correspondingly, if there were no inertia, cloud customers' ROI would be 2.4x compared to 2.2x with inertia. Interestingly, provider revenue is also 58% ( $= 33.0/56.9$ ) lower in the presence of inertia, despite customers using older and more expensive products.<sup>21</sup> This is mainly due to differences in overall cloud usage: without inertia, customers would adopt new products from the cloud provider faster and move more usage to the provider from on-premise computing and other cloud providers. As the last column shows, total cloud usage would more than double if inertia is fully eliminated. Appendix Figure H5 further decomposes time series of cloud usage by product and customer size with and without inertia: while all cloud products' usage increase without inertia, new product launches—X3 and X4—gain popularity much more quickly.

How much of the 62% consumer surplus lost due to inertia are true welfare loss? Because we agnostically estimate inertia as adoption costs for new products, we cannot speak to specific mechanisms driving customer inertia. It is possible that some of the adoption costs we estimate are real costs firms need to incur when adopting a new VM (e.g., engineering cost), and thus they do not represent true welfare loss. However, we can bound the true welfare loss from inertia by decomposing the consumer surplus difference with and without inertia into a direct adoption cost and an indirect cost from sub-optimal product choices. While the former is mechanism agnostic, the latter provides a lower bound for the true welfare loss due to inertia.

To first calculate the direct adoption cost, in the simulations without inertia, every time a new product is chosen, we sum up the adoption cost it would have incurred, and average across simulations. We then attribute the rest of the consumer surplus lost from inertia to sub-optimal product choices. Table 7 presents the results separately for small and large customers, which are similar. On average, only 2% of the loss in consumer surplus come from direct adoption costs, and the rest 98% are lost due to customers not choosing their most preferred products had there been no inertia. These results show that the vast majority of consumer surplus lost from inertia are true welfare loss, and thus we may be able to improve total welfare by incentivizing customers to adopt

---

<sup>21</sup>The average price per unit of cloud usage is  $23.9/35.8 = 0.67$  with inertia and  $56.9/99.7 = 0.57$  without inertia.

new products.

Table 7: Decomposition of consumer surplus lost due to inertia

	Small	Large
$\Delta$ Consumer surplus ( $10^3$ )	20.2	65.1
- Direct adoption cost	2.4%	1.6%
- Indirect cost from sub-optimal product choices	97.6%	98.4%

*Notes:* Table shows the difference between consumer surplus without inertia and with inertia, and decomposes this difference into direct adoption costs and indirect costs from sub-optimal product choices.

Finally, one may argue that the result of this decomposition depends on the length of the sample and results from Table 7 only apply to our particular sample. As a robustness check, in Appendix Table H3, we re-do the decomposition exercise for different lengths of samples by simulating our sample forward with the same set of customers, choice sets, and prices. We find that the direct adoption cost consistently accounts for a similar percentage of total consumer surplus lost from inertia.

## 6.2 Remedies

Given the large welfare cost of inertia, in this section, we explore two remedies that incentivize customers to adopt new products.

### 6.2.1 Subsidy for Cloud Migration

Cloud providers and third-party service providers often provide cloud migration services to help customers move their computing workloads. These services range from basic information provision to automated tools, or in some cases “white-glove” migration services for typically the largest customers. We first explore a counterfactual scenario where such services fully subsidize and eliminate customer inertia during our sample period. Welfare under this scenario is thus the same as if there were no inertia as shown in the second row of Table 6.

One way to implement this subsidy of adoption cost is via public financing. Table 8 shows how much it would cost (in re-scaled dollars) to fully subsidize adoption costs for each product and customer size. The subsidy cost is calculated in the same way as the direct adoption cost by adding up the adoption cost every time a customer chooses a new product in the simulations

without inertia. This assumes that all of the adoption costs are real costs, and thus the subsidy costs in Table 8 should be viewed as upper bounds. The total cost to fully subsidize adoption costs for all products and customers amounts to at most 59.2k. The total welfare gain (consumer surplus plus provider revenue) from eliminating adoption costs is 118.3k ( $=85.3k + 33.0k$ , from Table 6). In other words, the social planner would gain at least 2 dollars in welfare for each dollar spent on the subsidy. Besides welfare, this subsidy would also more than double the total amount of cloud usage (Table 6).

Table 8: Cost of subsidizing adoption costs ( $10^3$ )

Product	All Customers	Small	Large
X1	9.9 (0.2)	2.5 (0.1)	7.4 (0.2)
X2	10.8 (0.3)	3.1 (0.1)	7.7 (0.2)
X3	15.2 (0.4)	3.1 (0.1)	12.1 (0.4)
X4	18.7 (0.5)	4.5 (0.1)	14.2 (0.5)
Other	4.5 (0.1)	0.8 (0.03)	3.7 (0.1)
Total	59.2 (0.7)	14.0 (0.2)	45.1 (0.7)

*Notes:* Table presents the costs of fully subsidizing adoption costs in re-scaled dollars during our sample, broken down by products and customer sizes. Results are based on 200 simulations, with standard deviations in parentheses.

Another way to implement this subsidy is via cloud providers themselves since revenue also goes up in the absence of inertia. However, at least during our sample, the cost of this subsidy (59.2k) exceeds the revenue benefits (33.0k) cloud providers would gain from eliminating inertia. This is also true for any specific product and customer size. This suggests that, left on their own, cloud providers would not subsidize migration services, which may explain why in practice only limited “white-glove” migration services are offered. However, it is possible that cloud providers stand to gain more revenue benefits from eliminating inertia in the long run because they pay for customers’ adoption costs early on and receive revenue benefits from more cloud usage over time. In the next section, we compare long-run benefits of the full subsidy together with other remedies.

### 6.2.2 Introductory Discount

One reason subsidizing adoption costs may not be profitable for providers is that an adoption cost is paid for regardless of the size of the benefits to providers (e.g., usage or revenue). Another common approach online service providers use to help customers overcome inertia and adopt new products is in the form of introductory discounts for new products or new users of a product. Introductory discounts also reduce providers’ short-run revenue and potentially increase their long-run revenue by incentivizing customers to adopt new products, but do so by only subsidizing part of adoption costs proportional to revenue.

We first study a product-level introductory discount, often referred to as “public new product preview” or “open beta”, which discounts a new product after its initial launch. A second form of introductory discount we explore is a “personalized product trial” that discounts a product for any customers who have not used it before. We assume that both types of discounts only last for one month. Finally, we explore how much finer level of targeting by customer size may improve the outcome of these introductory discounts.

For each type of discount, our goal is to first assess whether it benefits the provider in the long run, and if so, what the optimal discount is for the provider. Taking those optimal discounts, we then evaluate their impact on consumer surplus and total welfare. To do so, we consider the launches of X4 across markets and solve for the provider’s optimal discounts that maximize long-run revenue.<sup>22</sup> Specifically, we conduct a grid search over a wide range of discounts for X4 for the first month after its launch, holding everything else constant. For each discount, we fix characteristics of each market (customers, choice sets, and prices) and simulate demand forward, allowing the number of tasks to grow following estimated time trends. Finally, we average across 1000 simulations to find the provider’s optimal discounts and evaluate their impact on welfare and cloud usage. As the provider trades off short-run vs. long-run benefits, the optimal discounts depend on how far out we simulate. In Appendix G, we show that the optimal discounts increase as the provider’s optimization horizon increases, and stabilize once we simulate 96 months after the initial X4 launch, which we present as our main results in this section.

Table 9 presents the results, relative to a baseline where there are no discounts at all. In the

---

<sup>22</sup>We assume a monthly discount factor of 0.99 for up to 96 months after X4 launch.

baseline, we simply simulate our model with inertia for 96 months, holding market characteristics fixed from the launch of X4, and calculate consumer surplus, provider revenue, and total cloud usage.

In the first row of Table 9, we find the optimal discount for a uniform product-level introductory discount (“new product preview”) at the launch of X4 to be 250%, i.e., the provider pays back to customers (e.g., possibly in store credits) 1.5 times their cloud spending on X4 in the first month after its launch.<sup>23</sup> The optimal discount is high because of the low price elasticity and high adoption cost we estimate in the model. With this product-level introductory discount, we find that the provider can improve its revenue by 1.3%, and at the same time, increase consumer surplus by 1.9%, improving total welfare by 1.7%. So the introductory discount results in a Pareto improvement of total welfare. Total cloud usage also goes up by 1.9%.

Table 9: Effect of different remedies for X4 launch in comparison to the baseline

	Optimal Discount	% Difference from Baseline			
		Provider Revenue	Consumer Surplus	Total Welfare	Cloud Usage
1. New product preview	250%	1.3%	1.9%	1.7%	1.9%
2. New product preview (by customer size)	Small: 325% Large: 250%	1.3%	2.1%	1.9%	2.0%
3. Personalized product trial	225%	12.8%	16.5%	15.4%	16.7%
4. Personalized product trial (by customer size)	Small: 275% Large: 200%	13.2%	16.0%	15.2%	15.8%
5. Subsidizing X4 adoption costs		11.8%	24.5%	20.8%	21.4%

*Notes:* Table compares welfare and cloud usage in the baseline (i.e., no discount) vs. product-level introductory discounts (i.e., new product preview) vs. customer-product level introductory discounts (i.e., personalized product trial) vs. fully subsidizing adoption costs whenever a customer chooses X4 for the first time. Optimal discounts are calculated based on grid searches over a wide range of discounts and 96 months of simulation after X4’s initial launch. Provider revenue, consumer surplus, total welfare, and cloud usage are shown in percentage improvements relative to the baseline. All results are averages across 1000 simulations.

In the second row, we allow the provider to give targeted differential discounts for small vs. large customers. We find that the optimal discount is 325% for small customers vs. 250% for large customers. There are two reasons for why the provider should give a higher discount to smaller

<sup>23</sup>Another interpretation of the magnitude of the discount is that it represents the combined discount in one month from a multi-month discount campaign.

customers in our setting, which is contrary to what managers typically do in practice in this market. First, small customers are less price elastic after taking into account their usage and thus need a higher discount to overcome similar adoption costs. Second, because small customers have lower usage in the beginning but grow faster, it is cheaper for the provider to discount small customers upfront to encourage new product adoptions and reap the benefits of their higher usage and revenue later. Finally, the targeted discounts only yield a small increase in provider revenue, consumer surplus, total welfare, and cloud usage compared to the uniform discount.

We now turn to “personalized product trials”—introductory discounts at the customer-product level. The product-level introductory discounts only yield small improvements in revenue for the provider due to its restrictive timing during the month of initial product launch. Personalized product trials relax this restriction by giving discounts to customers whenever they try a new product for the first time.

The third row in Table 9 shows the results for a uniform customer-product level introductory discount. We find that the optimal discount for this more flexible introductory discount is 225%, slightly lower than the optimal product-level introductory discount. In other words, under this optimal discount, the provider would pay its customers 1.25 times their first month’s cloud spending on X4 if they have not previously used it. While the optimal discount is lower, the flexibility of the customer-product level discount leads to a significant 12.8% revenue increase compared to the baseline. Consumer surplus also increases by 16.5%, and total welfare by 15.4%. Thus, the uniform personalized product trial dominates new product previews in terms of both consumer surplus and provider revenue. Total cloud usage also increases by 16.7%.

In the fourth row, we again show the results for a personalized product trial with targeted differential discounts for small vs. large customers. The optimal discounts are, again, higher for small customers than large customers. While the provider revenue (necessarily) increases under the targeted discounting compared to the uniform personalized product trial, consumer surplus decreases, and so do total welfare and cloud usage. This is because with differential discounts, the provider optimally chooses a smaller discount for large customers than under uniform discounting, which slows down cloud adoption and decreases consumer surplus and total welfare.

Finally, to facilitate comparison, we repeat the exercise from the previous section and simulate demand and welfare under a full subsidy of adoption costs for X4 for 96 months after its initial launch



(rather than during our sample period in the previous section). The last row in Table 9 shows the results. Consumer surplus is the highest compared to other remedies, representing a 24.5% increase compared to the baseline. Provider revenue net of the costs of subsidizing X4’s adoption costs is lower than those under the personalized product trials but higher than new product previews, representing a 11.8% increase compared to the baseline. So while it is not profitable for the provider to fully subsidize adoption costs during our sample period, it is profitable in the long run as the provider continues to receive revenue from customers who overcome their inertia to adopt new cloud products. This may explain why cloud migration and management services, although limited to only selected customers, are offered in the market. Total cloud usage, similar to consumer surplus, is also the highest under the full subsidy. Taken together, we find that fully subsidizing adoption costs yields the highest total welfare; though it is also profitable for the provider, it generates lower revenue compared to personalized product trials.

## 7 Conclusion

The public cloud is one of the most important technological innovations in the 21st century. It has grown to be a significant industry and enabled much of the rest of the economy to digitize. Both digitization and the cloud itself are still fast growing and becoming increasingly more important. Cloud alleviates customers’ needs to buy and maintain physical computing hardware, and thus lowers the upfront cost of owning computing resources and allows customers of all sizes access to the latest technology. But, at the same time, the cloud’s lack of maintenance needs may induce customers to develop inertia that prevents them from adopting newer and better products.

In this paper, we estimate that the public cloud provides significant returns to its customers, with an average ROI of 2.2x (consumer surplus divided by cloud spending), based on data from the IaaS market. This is comparable to estimates for other significant IT inventions from the 20th century—mainframes and PCs—and translates to an average consumer welfare gains of \$56 billion annually from 2016 to 2018, or \$248 billion in 2022 if extrapolated. Cloud computing is also claimed to act as a democratizing force, as it allows firms access to computing resources without hardware ownership, empowering smaller firms ([Bloom and Pierri 2018](#)). We empirically show that, different from those earlier technologies, the cloud indeed disproportionately benefits its small customers.

Despite the significant welfare benefits, we estimate that more than half of the potential welfare benefits from cloud are lost due to inertia, which lowers both consumer surplus and provider revenue, as customers are slower to adopt cloud products. To help customers overcome inertia and adopt new products, we find that both subsidizing adoption costs and offering introductory discounts are effective. The subsidy, if financed publicly, can generate substantial returns, and is even profitable for the provider to pay for in the long run (though not in the short run). Personalized product trials, which offer first-month discounts for customers to adopt a new product, yield the highest revenue gain for the provider, while increasing consumer surplus at the same time, resulting in a Pareto improvement compared to the baseline.

One major contribution of the paper is the multiple-choice continuous-usage demand model itself. The model only requires the researcher to observe customer-product level usage, and thus can be applied broadly to any other market with usage-based services and where customers may choose more than one products. One limitation of the model in our setting is that demand is static, whereas customers facing inertia may exhibit forward-looking behavior, depending on their level of sophistication. We do not capture such dynamics given the already complex static model, in line with the literature on consumer inertia.

Another limitation of our paper is that we focus on existing customers' demand for products from one cloud provider due to data limitation. While this limitation does not have significant impact on our welfare estimates for existing customers, it does not allow us to study new customers' adoption choices between cloud providers as well as its implications for providers' pricing decisions. Policymakers have expressed concerns over migration frictions between cloud providers as a potential barrier to competition.<sup>24</sup> Our paper serves as a first step — looking at frictions within one provider and subsidies and introductory discounts across products — to understanding inertia in the cloud computing market, and further highlights its importance for future research.

Finally, an alternative way to measure welfare benefits of cloud is to analyze the impact of cloud adoption and usage on direct measures of firm performance, which is out of the scope of this paper due to data limitations but would be complementary to our paper and an interesting avenue for future research.

---

<sup>24</sup>See [https://www.ftc.gov/system/files/ftc\\_gov/pdf/Cloud-RFI-June-21-2023.pdf](https://www.ftc.gov/system/files/ftc_gov/pdf/Cloud-RFI-June-21-2023.pdf) and [https://www.ofcom.org.uk/\\_\\_data/assets/pdf\\_file/0029/256457/cloud-services-market-study-interim-report.pdf/](https://www.ofcom.org.uk/__data/assets/pdf_file/0029/256457/cloud-services-market-study-interim-report.pdf/), accessed Sept. 4, 2023.

## References

- Bartoš, Vojtěch, Michal Bauer, Julie Chytilová and Filip Matějka. 2016. “Attention discrimination: Theory and field experiments with monitoring information acquisition.” *American Economic Review* 106(6):1437–75.
- Berry, Steven, James Levinsohn and Ariel Pakes. 1995. “Automobile prices in market equilibrium.” *Econometrica* 63(4):841–890.
- Berry, Steven T. 1994. “Estimating discrete-choice models of product differentiation.” *The RAND Journal of Economics* 25(2):242–262.
- Bessen, James. 2020. “Industry concentration and information technology.” *The Journal of Law and Economics* 63(3):531–555.
- Bhargava, Saurabh, George Loewenstein and Justin Sydnor. 2017. “Choose to lose: Health plan choices from a menu with dominated option.” *The Quarterly Journal of Economics* 132(3):1319–1372.
- Bloom, Nicholas and Nicola Pierri. 2018. “Cloud computing is helping smaller, newer firms compete.” *Harvard Business Review* 94(4).
- Bresnahan, Timothy F. 1986. “Measuring the spillovers from technical advance: mainframe computers in financial services.” *American Economic Review* 76(4):742–755.
- Brynjolfsson, Erik and Lorin M Hitt. 2000. “Beyond computation: Information technology, organizational transformation and business performance.” *Journal of Economic Perspectives* 14(4):23–48.
- Burda, Martin, Matthew Harding and Jerry Hausman. 2012. “A Poisson mixture model of discrete choice.” *Journal of Econometrics* 166(2):184–203.
- Burnham, Thomas A, Judy K Frels and Vijay Mahajan. 2003. “Consumer switching costs: a typology, antecedents, and consequences.” *Journal of the Academy of Marketing Science* 31(2):109–126.
- Castillo, Juan Camilo. 2020. “Who benefits from surge pricing?” *Available at SSRN 3245533*.
- Cisco. 2018. Cisco Global Cloud Index: Forecast and Methodology, 2016-2021. White paper Cisco.

- Cyert, Richard M and James G March. 1963. *A behavioral theory of the firm*. Englewood Cliffs, N.J.: Prentice-Hall.
- De Loecker, Jan, Jan Eeckhout and Gabriel Unger. 2020. “The rise of market power and the macroeconomic implications.” *The Quarterly Journal of Economics* 135(2):561–644.
- Decker, Ryan A, John Haltiwanger, Ron S Jarmin and Javier Miranda. 2016. “Declining business dynamism: What we know and the way forward.” *American Economic Review* 106(5):203–07.
- Deryugina, Tatyana, Alexander MacKay and Julian Reif. 2020. “The long-run dynamics of electricity demand: Evidence from municipal aggregation.” *American Economic Journal: Applied Economics* 12(1):86–114.
- Erdem, Tülin and Michael P Keane. 1996. “Decision-making under uncertainty: Capturing dynamic brand choice processes in turbulent consumer goods markets.” *Marketing Science* 15(1):1–20.
- Farrell, Joseph and Paul Klemperer. 2007. Coordination and lock-in: Competition with switching costs and network effects. Vol. 3 of *Handbook of Industrial Organization* Elsevier.
- Forman, Chris. 2005. “The corporate digital divide: Determinants of Internet adoption.” *Management Science* 51(4):641–654.
- Gartner. 2017. “Gartner says worldwide IaaS public cloud services market grew 31 percent in 2016.”. Accessed Nov. 27, 2021, <https://www.gartner.com/en/newsroom/press-releases/2017-09-27-gartner-says-worldwide-iaas-public-cloud-services-market-grew-31-percent-in-2016>.
- Gartner. 2018. “Gartner says worldwide IaaS public cloud services market grew 29.5 percent in 2017.”. Accessed April 7, 2023, <https://www.gartner.com/en/newsroom/press-releases/2018-08-01-gartner-says-worldwide-iaas-public-cloud-services-market-grew-30-percent-in-2017>.
- Gartner. 2019. “Gartner says worldwide IaaS public cloud services market grew 31.3 % in 2018.”. Accessed April 27, 2023, <https://www.gartner.com/en/newsroom/press-releases/2019-07-29-gartner-says-worldwide-iaas-public-cloud-services-market-grew-31point3-percent-in-2018>.
- Gartner. 2022. “Gartner Says Worldwide IaaS Public Cloud Services Market Grew 41.4% in

- 2021.”. Accessed July 8, 2022, <https://www.gartner.com/en/newsroom/press-releases/2022-06-02-gartner-says-worldwide-iaas-public-cloud-services-market-grew-41-percent-in-2021>.
- Gartner. 2023. “Gartner Forecasts Worldwide Public Cloud End-User Spending to Reach Nearly \$600 Billion in 2023.”. Accessed April 19, 2023, <https://www.gartner.com/en/newsroom/press-releases/2023-04-19-gartner-forecasts-worldwide-public-cloud-end-user-spending-to-reach-nearly-600-billion-in-2023>.
- Greenstein, Shane M. 1993. “Did installed base give an incumbent any (measureable) advantages in federal computer procurement?” *The RAND Journal of Economics* 24(1):19–39.
- Handel, Benjamin and Joshua Schwartzstein. 2018. “Frictions or mental gaps: what’s behind the information we (don’t) use and when do we care?” *Journal of Economic Perspectives* 32(1):155–78.
- Handel, Benjamin R and Jonathan T Kolstad. 2015. “Health insurance for" humans": Information frictions, plan choice, and consumer welfare.” *American Economic Review* 105(8):2449–2500.
- Hanna, Rema, Sendhil Mullainathan and Joshua Schwartzstein. 2014. “Learning through noticing: Theory and evidence from a field experiment.” *The Quarterly Journal of Economics* 129(3):1311–1353.
- Hendel, Igal. 1999. “Estimating multiple-discrete choice models: An application to computerization returns.” *The Review of Economic Studies* 66(2):423–446.
- Hortaçsu, Ali and Chad Syverson. 2004. “Product differentiation, search costs, and competition in the mutual fund industry: A case study of S&P 500 index funds.” *The Quarterly Journal of Economics* 119(2):403–456.
- Jin, Wang and Kristina McElheran. 2017. “Economies before scale: survival and performance of young plants in the age of cloud computing.” *Rotman School of Management working paper* (3112901).
- Kilcioglu, Cinar and Justin M Rao. 2016. Competition on price and quality in cloud computing. In *Proceedings of the 25th International Conference on World Wide Web*. International World Wide Web Conferences Steering Committee pp. 1123–1132.

- Koulayev, Sergei, Marc Rysman, Scott Schuh and Joanna Stavins. 2016. “Explaining adoption and use of payment instruments by US consumers.” *The RAND Journal of Economics* 47(2):293–325.
- Liebman, Jeffrey B and Neale Mahoney. 2017. “Do expiring budgets lead to wasteful year-end spending? Evidence from federal procurement.” *American Economic Review* 107(11):3510–49.
- MacKay, Alexander and Marc Remer. 2022. “Consumer inertia and market power.” *Available at SSRN 3380390*.
- Peng, Sida, Peichun Wang, Eric Auerbach, Hongwei Liang and Andy Wu. 2021. “Digitization and employment in the pandemic: Evidence from seventy billion emails.” Working Paper.
- Petrin, Amil and Kenneth Train. 2010. “A control function approach to endogeneity in consumer choice models.” *Journal of Marketing Research* 47(1):3–13.
- Quan, Thomas W. and Kevin R. Williams. 2018. “Product variety, across-market demand heterogeneity, and the value of online retail.” *The RAND Journal of Economics* 49(4):877–913.
- Tambe, Prasanna, Lorin Hitt, Daniel Rock and Erik Brynjolfsson. 2020. Digital capital and superstar firms. Working Paper 28285 National Bureau of Economic Research.
- Tambe, Prasanna and Lorin M Hitt. 2012. “The productivity of information technology investments: New evidence from IT labor data.” *Information Systems Research* 23(3-part-1):599–617.

# Appendices

## A Maximum Likelihood Estimation

In this section, we discuss challenges of maximum likelihood estimation for our model. Suppose for now that task sizes are known and equal to one, i.e.,  $q_{im} = 1$ . Then our model collapses to that of

Burda, Harding and Hausman (2012), and the likelihood of  $(y_{i0m}, y_{i1m}, \dots, y_{iJm})$  can be written as

$$L_i(\lambda_{im}, \alpha_i, \beta_i, \delta_i | y_{i0m}, y_{i1m}, \dots, y_{iJm}) \\ = \frac{\exp(-\lambda_{im}) \lambda_{im}^{\sum_{j=0}^J y_{ijm}}}{(\sum_{j=0}^J y_{ijm})!} \cdot \prod_{j=0}^J P_{ijm}(1)^{y_{ijm}},$$

which is a product of two components: 1) the probability of receiving  $\sum_{j=0}^J y_{ijm}$  tasks; and 2) the probability that VM  $j$  is chosen for exactly  $y_{ijm}$  tasks, for all  $j$ . This likelihood is easy to compute because the task size is known. However, if the task size is unknown, the likelihood becomes the sum of the probabilities of all scenarios that generate usages  $(y_{i0m}, y_{i1m}, \dots, y_{iJm})$ . For example, if we observe usage of  $y_{ijm} = 10$ , it could be generated by one task of size 10, two tasks of sizes (1, 9), (2, 8), and so on, amounting to infinite possibilities with continuous usage. In other words, the likelihood becomes an infinite sum of convoluted integration, which is computationally infeasible. Therefore, we take the GMM approach instead.

## B Proof of Proposition 1

*Proof.* For the simplified model without inertia, customer heterogeneity, or time trend in the number of tasks, let the probability of customer  $i$  choosing product  $j$  in market  $m$  for a task of size  $q_{im}$  be denoted by  $P_{ijm}(q_{im})$ , i.e.,

$$P_{ijm}(q_{im}) = \begin{cases} \frac{\exp(\alpha p_{jm} q_{im} + X_j \beta)}{1 + \sum_{l=1}^J \exp(\alpha p_{lm} q_{im} + X_l \beta)} & \text{if } j \neq 0, \\ \frac{1}{1 + \sum_{l=1}^J \exp(\alpha p_{lm} q_{im} + X_l \beta)} & \text{otherwise.} \end{cases}$$

**Zero Usage Probability Moment** For any given number of tasks  $n_{im}$  and task size  $q_{im}$ , the probability of customer  $i$  not using product  $j$  in market  $m$  is equal to the probability of customer  $i$  not choosing product  $j$  in market  $m$  for any of the  $n_{im}$  tasks. Because tasks are independent, this probability is equal to  $(1 - P_{ijm}(q_{im}))^{n_{im}}$ , which we call  $Pr(\text{Product Choice})$ . Now, to obtain the probability of zero usage at the customer-product level, we first integrate the probability of product choice over the exponential distribution of task sizes, and then the Poisson distribution of

the number of tasks:

$$\mathbb{P}(y_{ijm} = 0) = \sum_{n_{im}=0}^{\infty} \underbrace{\frac{\exp(-\lambda)\lambda^{n_{im}}}{n_{im}!}}_{Pr(\text{Number of Tasks})} \cdot \int_{q_{im}} \underbrace{(1 - P_{ijm}(q_{im}))^{n_{im}}}_{Pr(\text{Product Choice})} \underbrace{\frac{1}{\gamma} \exp(-\frac{1}{\gamma}q_{im})}_{Pr(\text{Task Size})} dq_{im}.$$

**Expected Usage Moment** Again, first fix the number of tasks to be  $n_{im}$ . By integrating over the exponential distribution of task sizes and the probability of product  $j$  being chosen, we obtain the average task size on product  $j$ :

$$\int_{q_{im}} q_{im} P_{ijm}(q_{im}) \frac{1}{\gamma} \exp(-\frac{1}{\gamma}q_{im}) dq_{im}.$$

Then, we multiply it by  $n_{im}$  to get the expected usage for each product given  $n_{im}$  tasks, because task sizes are the same across tasks for the same customer in the same market. Finally, integrating over the Poisson distribution of the number of tasks, we obtain the expected usage at the customer-product level:

$$\begin{aligned} \mathbb{E}(y_{ijm}) &= \sum_{n_{im}=0}^{\infty} \frac{\exp(-\lambda)\lambda^{n_{im}}}{n_{im}!} \cdot \left( n_{im} \int_{q_{im}} q_{im} P_{ijm}(q_{im}) \frac{1}{\gamma} \exp(-\frac{1}{\gamma}q_{im}) dq_{im} \right) \\ &= \left( \sum_{n_{im}=0}^{\infty} \frac{\exp(-\lambda)\lambda^{n_{im}}}{n_{im}!} \cdot n_{im} \right) \int_{q_{im}} q_{im} P_{ijm}(q_{im}) \frac{1}{\gamma} \exp(-\frac{1}{\gamma}q_{im}) dq_{im} \\ &= \underbrace{\lambda}_{\mathbb{E}(\text{Number of Tasks})} \underbrace{\int_{q_{im}} q_{im} P_{ijm}(q_{im}) \frac{1}{\gamma} \exp(-\frac{1}{\gamma}q_{im}) dq_{im}}_{\mathbb{E}(\text{Task Size})}, \end{aligned}$$

where  $\lambda$  is exactly the expected number of tasks of the Poisson distribution. □

## C Prices, Instruments, and First Stage Results

We first construct a price index for the prices of the “other” VM across markets, which include all of the provider’s VM products other than the X series. Specifically, in each market, we calculate the price index as the average price weighted by the demand of each product from 2017, so that the variation of the price index over time is driven by changes in product availability and price changes,



rather than demand.

We use three cost shifters as instruments for prices. Cloud providers’ variable costs vary across regions and products due to regional differences in electricity costs and because the same VMs are supported on different hardware across regions. Our first instrument uses 2019 electricity prices for each region to capture the provider’s regional variation in costs.

The second instrument uses power ratings of different types of hardware and the third instrument uses costs of procuring an additional cluster of different types of hardware. Because different VMs can run on different hardware and the allocation of VMs onto different hardware clusters differ across regions, we transform the hardware costs and power ratings at the hardware type level to be at the VM level, weighting them by demand from 2017. As a result, both the power rating and the hardware cost instruments vary at the region-product level.

The first stage results of the control function, introduced in Section 4, are shown in Table C1. We regress prices of all cloud products (X1-X4 and Other) at the product-market level on the three cost shifters, as well as the product characteristics  $X_j$ . The estimated coefficients of the cost shifters are statistically significant and have the expected signs: VM prices increase with electricity prices and VM hardware costs. VM prices also increase with VM hardware power rating, i.e., VMs with higher energy consumption (higher rating) are more expensive.

## D Computation Details

To calculate the GMM objective function, we numerically integrate over the distributions of task sizes, random coefficients, as well as the number of tasks. For task sizes, we use the Gauss-Laguerre quadrature with 10 nodes. For the random coefficients, we use the Halton sequence with 20 nodes. For the number of tasks, we compute moments for any  $k$  number of tasks, for  $k = 0, 1, 2, \dots, 120$ . For large customers at the end of our sample (i.e., largest number of tasks), 120 tasks cover 99.8% of the distribution of the number of tasks.

To speed up the computations, we parallelize  $k = 0, 1, 2, \dots, 120$  across 28 cores. In total, this integration procedure yields  $10 \times 20 \times 121/28 = 864$  points to compute per core. On a 2.6 GHz Intel Gold 6132 CPU with 28 cores, one evaluation of the objective function with gradient takes about 162 seconds, using 508G of memory. We are thus limited by the curse of dimensionality to

Table C1: Control function first stage results

	VM price
Electricity price	0.092*** (0.015)
VM (hardware) power rating	0.766*** (0.078)
VM (hardware) cost	1.151*** (0.065)
X1	0.195*** (0.025)
X2/X3	0.136*** (0.027)
X4	−0.050*** (0.018)
Promo	−0.112*** (0.007)
Constant	−0.933*** (0.106)
Observations	3,790
R <sup>2</sup>	0.600

*Notes:* Table presents results from the first stage in the control function approach. We regress VM prices at the product-market level on the three cost shifters as well as product characteristics. \*p<0.1; \*\*p<0.05; \*\*\*p<0.01

further increase the number of nodes for the numerical integrations. Finally, we minimize the GMM objective function using the BFGS quasi-Newton algorithm supplied with the analytical gradient.

## E Monte Carlo Simulations

To illustrate identification of our model and our estimation procedure in the finite sample, we conduct Monte Carlo simulations for the simplified model without inertia, customer heterogeneity, or time trend in the number of tasks.

The data generating process is as follows. To generate tasks, we let both the Poisson distribution for the number of tasks and the exponential distribution for task sizes to be parameterized by one parameter each, i.e., their means  $\lambda$  and  $\gamma$ . For product choices, we consider markets each with four products and an outside option, where the utility of customer  $i$  choosing product  $j$  in market  $m$  for task  $k$  is given by

$$u_{ijmk} = \alpha p_{jm} q_{im} + X_j \beta + \epsilon_{ijmk},$$

where  $X_j$ 's are product dummies, and the idiosyncratic shocks  $\epsilon_{ijmk}$ 's are distributed type-I extreme value. Mean utility for the outside option is normalized to zero.

To generate estimating samples, we first fix parameters at some true values (see Table E2). We generate two samples, both with 100,000 customers per market, but one with 100 markets and the other with 500 markets. We draw prices independently from a normal distribution at the product-market level. Then, we draw each customer's number of tasks and the size of each task, and simulate their product choices and usage based on the prices and product dummies in their utilities.

Finally, we take the resulting estimating samples through our GMM estimation procedure and present the results in Table E2. We accurately recover true values of the parameters when initial values are different from the truth. Moreover, the accuracy slightly improves with more markets (and thus more price variation).

Table E2: Monte Carlo simulation results

		True Value	Initial Value	Spec1	Spec2
	No. Customers			100k	100k
	No. Markets			100	500
<b>Product choice</b>	Product1FE	1.000	1.000	1.002	1.001
	Product2FE	1.500	1.000	1.503	1.500
	Product3FE	2.000	1.000	2.001	2.000
	Product4FE	2.500	1.000	2.501	2.501
	Price	-10.000	-15.000	-9.995	-9.998
<b>Task size</b>	Mean	0.300	0.100	0.300	0.300
<b>Number of tasks</b>	Mean	2.000	1.000	1.999	2.000
<b>Runtime (s)</b>				3285	18567

*Notes:* Table presents results from Monte Carlo simulations. We generate two samples, both with 100,000 customers per market, but one with 100 markets and the other with 500 markets. Both samples are generated based on true values of the parameters. Prices are drawn from an i.i.d. normal distribution at the product-market level. Estimation is based on our GMM procedure and started at the initial values of the parameters.

## F Own-Price Elasticity

In this section, we derive the formula for the own-price elasticity,  $\frac{d\mathbb{E}(y_{ijm})}{dp_{jm}} \frac{p_{jm}}{\mathbb{E}(y_{ijm})}$ . Following Proposition 1, the expected usage of customer  $i$  on VM  $j$  in market  $m$  is

$$\mathbb{E}(y_{ijm}) = \lambda_{im} \int_{q_{im}} q_{im} P_{ijm}(q_{im}) \cdot \frac{1}{\gamma_i} \exp\left(-\frac{1}{\gamma_i} q_{im}\right) dq_{im}.$$

Differentiating the expected usage with respect to its price:

$$\frac{d\mathbb{E}(y_{ijm})}{dp_{jm}} = \lambda_{im} \int_{q_{im}} \alpha_i q_{im}^2 \cdot P_{ijm}(q_{im}) (1 - P_{ijm}(q_{im})) \frac{1}{\gamma_i} \exp\left(-\frac{1}{\gamma_i} q_{im}\right) dq_{im}.$$

Finally, the own-price elasticity is given by

$$\frac{d\mathbb{E}(y_{ijm})}{dp_{jm}} \frac{p_{jm}}{\mathbb{E}(y_{ijm})} = \frac{\alpha_i p_{jm} \int_{q_{im}} q_{im}^2 \cdot P_{ijm}(q_{im}) (1 - P_{ijm}(q_{im})) \exp\left(-\frac{1}{\gamma_i} q_{im}\right) dq_{im}}{\int_{q_{im}} q_{im} P_{ijm}(q_{im}) \exp\left(-\frac{1}{\gamma_i} q_{im}\right) dq_{im}}.$$

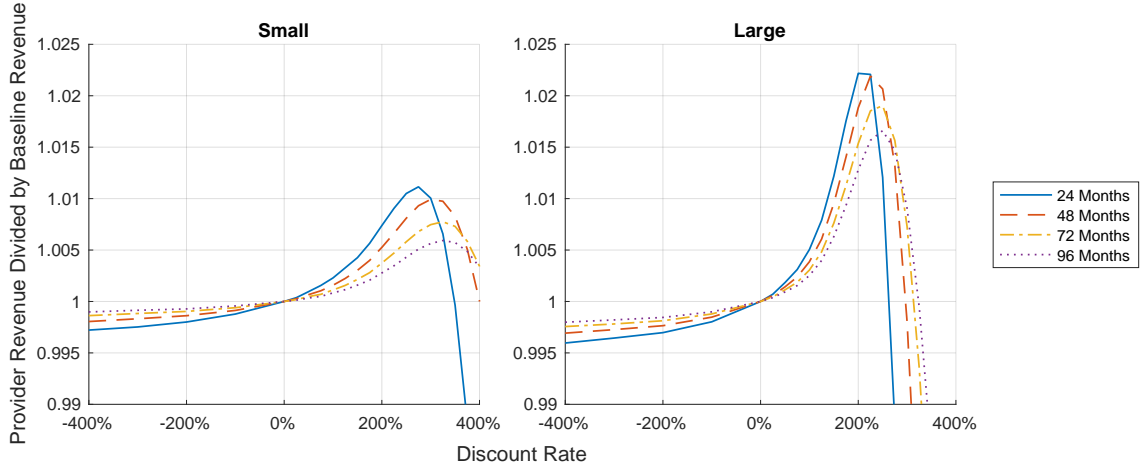
## G Optimal Introductory Discounts

In this section, we introduce how we find the provider's optimal introductory discounts and compare the short-run vs. long-run optimal discount rates. For each type of discount (new product preview

or personalized product trial), we conduct a grid search over discount rates from -1000% to 1000%, where a discount rate of  $x$  on price  $p_{jm}$  implies a discounted price of  $(1 - x)p_{jm}$ . We first conduct a coarse search over this range with 100% intervals to narrow down to a 200% range for the optimal discount. Then, within this 200% range, we refine the grid search to 25% intervals. For each discount, we fix characteristics of each market (customers, choice sets, and prices) and simulate demand forward following the same procedure as in Section 5.1 to compute consumer surplus, provider revenue, and cloud usage. We average across 1000 simulations to find the provider's optimal discount rates.

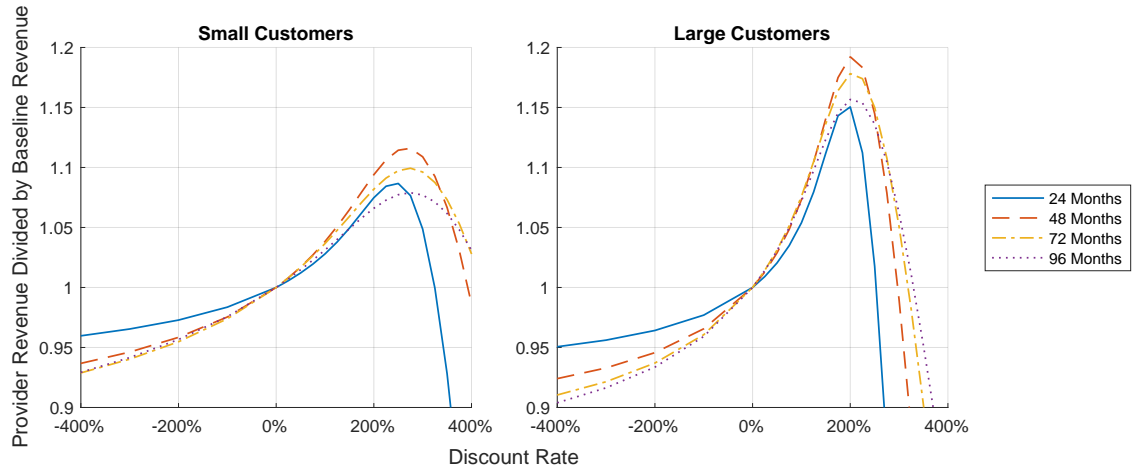
To highlight the difference between the short-run vs. long-run trade-off for the provider, in Figures G1 and G2, we present one line each representing the relationship between the introductory discount rates and the provider's revenue, for when the provider's total revenue objective is defined for 24, 48, 72, and 96 months after X4's launch. The peak of each curve represents the corresponding optimal discount. As the provider's horizon becomes longer, the optimal discount becomes larger as the benefit of overcoming customer inertia increases. The optimal discount stabilizes between 72 and 96 months as later months do not contribute as much to the provider's objective.

Figure G1: New product preview: Different horizons



*Notes:* Figure shows the ratio of provider revenue with different introductory discounts relative to the baseline, separately for small and large customers. Different lines represent different horizons in the provider's revenue calculations. Results are averages across 1000 simulations. Discount rates, shown on the x-axis, are in absolute levels: A discount rate of  $x$  means that, in the first month after X4's launch, customers would pay  $(1-x)$  times their original spending on X4.

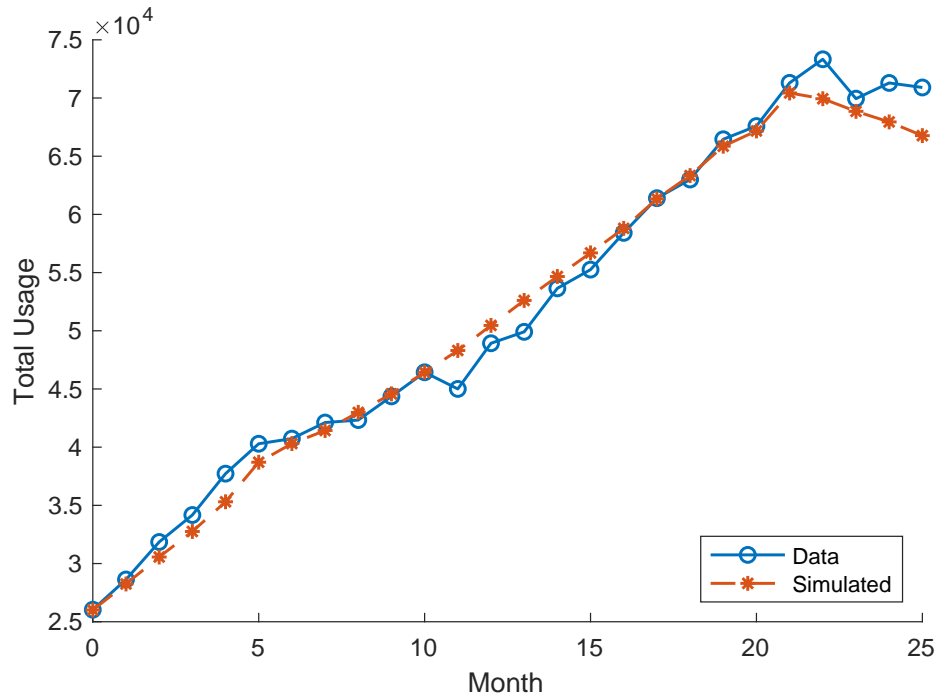
Figure G2: Personalized product trial: Different horizons



*Notes:* Figure shows the ratio of provider revenue with different introductory discounts relative to the baseline, separately for small and large customers. Different lines represent different horizons in the provider's revenue calculations. Results are averages across 1000 simulations. Discount rates, shown on the x-axis, are in absolute levels: A discount rate of  $x$  means that, in the first month that a customer uses X4, the customer would pay  $(1-x)$  times her original spending on X4.

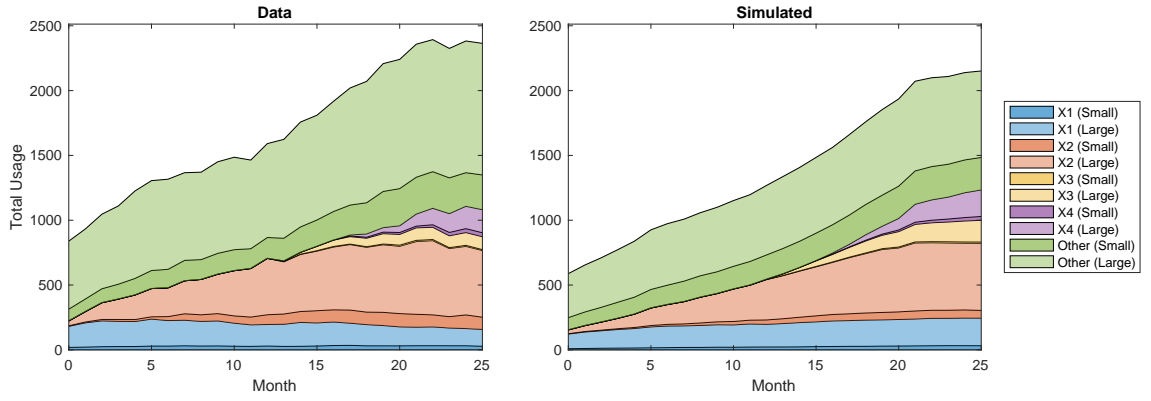
## H Additional Figures and Tables

Figure H3: Time series fit of total computing demand



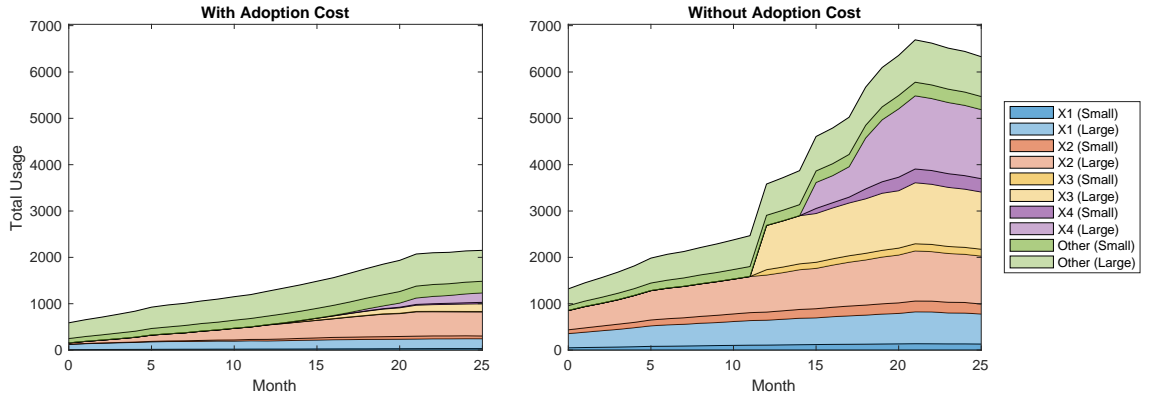
*Notes:* Figure shows total computing demand (including outside option) over time (months from beginning of sample) from the data (blue solid line with circles) and from simulations at our estimates (red dashed line with stars).

Figure H4: Time series fit of cloud demand by product and customer size



*Notes:* Figure shows cloud usage over time (months since beginning of sample) from the data (left) and from simulations at our estimates (right), by product and customer size.

Figure H5: Impact of inertia on usage by product and customer size



*Notes:* Figure shows cloud usage over time (months since beginning of sample) simulated with adoption cost (left) and without adoption cost (right), by product and customer size.



Table H3: Share of direct adoption cost by different sample lengths

# of months after X4 launch	Small	Large
12	2.40%	1.65%
24	2.40%	1.65%
36	2.40%	1.65%
48	2.39%	1.64%
60	2.38%	1.64%
72	2.38%	1.64%
84	2.37%	1.63%
96	2.37%	1.63%

*Notes:* Table shows the share of direct adoption cost as a percentage of total consumer surplus lost from inertia when the sample is simulated forward for different lengths, separately for small and large customers. The sample is simulated with the same market characteristics (customers, choice sets, and prices) from the time of the X4 launch. Results are averages across 200 simulations.