# Concentration, Market Power, and Misallocation: The Role of Endogenous Customer Acquisition<sup>\*</sup>

Hassan Afrouzi<sup>†</sup> Columbia University and NBER Andres Drenik<sup>‡</sup> The University of Texas at Austin Ryan Kim<sup>§</sup> Johns Hopkins University

First Draft:August 2020This Draft:May 2023

#### Abstract

This paper explores how different margins of market share are related to markups. Using merged microdata on producers and consumers, we document that a firm's market share is mainly related to its number of customers, while its price-cost markup is associated only with its average sales per customer. We develop a new model that reflects this empirical evidence and the endogenous nature of customer acquisition. When calibrated, this model predicts a higher degree of markup dispersion, which suggests greater efficiency losses due to customer misallocation. An analysis of the efficient allocation in this model reveals that compared with the equilibrium, aggregate TFP and output are 10.8% and 14% higher, respectively.

*JEL Codes:* D61, D24, D43, E22 *Key Words:* Misallocation, Customer acquisition, Markups, Concentration

<sup>\*</sup>We thank David Argente, Ariel Burstein, Doireann Fitzgerald, and seminar participants at various institutions and conferences for valuable comments and suggestions. Luigi Caloi provided superb research assistance. Previous versions of this manurscript were circulated under the title "Growing by the Masses: Revisiting the Link between Firm Size and Market Power." Researcher(s)' own analyses calculated (or derived) based in part on data from Nielsen Consumer LLC and marketing databases provided through the NielsenIQ Datasets at the Kilts Center for Marketing Data Center at The University of Chicago Booth School of Business. The conclusions drawn from the NielsenIQ data are those of the researcher(s) and do not reflect the views of NielsenIQ. NielsenIQ is not responsible for, had no role in, and was not involved in analyzing and preparing the results reported herein.

<sup>&</sup>lt;sup>†</sup>Department of Economics. Email: hassan.afrouzi@columbia.edu.

<sup>&</sup>lt;sup>‡</sup>Department of Economics. Email: andres.drenik@austin.utexas.edu.

<sup>&</sup>lt;sup>§</sup>School of Advanced International Studies (SAIS). Email: rkim59@jhu.edu.

## 1 Introduction

In standard macroeconomic models of endogenous markups—e.g., Atkeson and Burstein (2008), Klenow and Willis (2016)—firms with larger market shares charge higher pricecost markups. At a macro level, such mechanisms generate a close connection between market share distribution and markup dispersion, which is a source of resource misallocation due to firm-level wedges (Restuccia and Rogerson, 2008, Hsieh and Klenow, 2009). However, whereas in these models market share is driven by firms' sales to a representative household (the intensive margin of demand), in practice a firm's market share is affected by more than one margin. Specifically, firms spend a vast amount of resources on expanding their customer bases (the extensive margin of demand).<sup>1</sup> This leads us to ask: Do both the intensive and extensive margins of market share have the same relationship with a firm's market power (markups)? If not, how does this change our understanding of the ties between concentration, market power, and misallocation?

In this paper, we investigate how these two demand margins are related to firms' market shares and markups. Merging individual product-level consumption data and producer-level data, we find that firms' markups correlate *only* with their average sales per customer (the intensive margin), but not with the size of their customer base (the extensive margin). Yet only about 22% of the variation in firms' market share is tied to the margin associated with markups. To explore the macroeconomic implications of these findings, we develop a model with endogenous customer acquisition that is consistent with these facts. The model predicts an equilibrium relationship between market share and market power as in standard models; however, when calibrated to match the same distribution of firm size, it generates noticeably higher variation in markups across firms. This higher markup dispersion hints at higher efficiency losses due to misallocation of demand. To measure these losses rigorously, we characterize and quantify the first best allocation in our model: Relative to the efficient allocation, equilibrium aggregate TFP and output are 10.8% and 14% lower, respectively.

Our first contribution is to empirically investigate the relationship between different margins of market share and markups by merging the Nielsen Homescan Panel and Compustat datasets. First, our most novel finding is that firms' markups are not correlated with the size of their customer bases. Instead, they are positively associated *only* with their average sales per customer. Second, we find that around three-quarters of the variation across firms' market shares is explained by the margin that is not correlated with markups;

<sup>&</sup>lt;sup>1</sup>Arkolakis (2010) reports total spending on marketing as high as 5% of GDP in the US.

i.e., the extensive margin. This is consistent, both qualitatively and quantitatively, with the contemporary findings in Einav, Klenow, Levin, and Murciano-Goroff (2022) and Argente, Fitzgerald, Moreira, and Priolo (2021), who, using different data, show that most of the variation in firms' sales is driven by the size of their customer base. Third, we find that firms' non-production expenses are positively associated with the acquisition of new customers but not with their average sales per customer, which suggests that firms engage in activities in order to expand their customer bases by acquiring new customers.

These empirical findings inform the theoretical model we develop to understand the links between varying sources of market share and market power. In our model, at each period, a set of new firms draw different productivity levels and decide whether to enter the economy subject to fixed operating costs. Conditional on entry, these firms are monopolistically competitive and can spend resources to acquire new customers. These firms face semi-kinked demand curves a la Kimball (1995) from each customer, which indicates that each customer's demand is more elastic when a firm's relative price is larger (also known as Marshall's second law of demand). Therefore, whereas a firm's total customer base merely shifts its demand, as in Phelps and Winter (1970), the elasticity of this demand is determined by each customer's individual demand curve. Thus, the model generates a comovement between markups and average sales per customer, but not between markups and the number of customers (consistent with our first fact). Moreover, it allows firms to grow through both margins of demand (second fact), which implies an endogenous relationship between sales and non-production costs (third fact).

By allowing firms to grow through the extensive margin, our model breaks the *direct* relationship between market share and demand elasticity generated by the exogenous shape of the demand curve in standard models. Since market share in our model is determined through two separate margins—but markups are only correlated with one of those margins—conditional on each size group, there is a whole range of markups that firms within that group charge. However, on average, our model still creates a positive correlation between markups and market shares as an *equilibrium outcome* bearing unique counterfactual implications. Although the extensive margin does not correlate with markups conditional on sales per customer, this channel allows for a relationship between size and market power through the costs and benefits of customer acquisition: More productive firms expect to charge higher markups to their customers, anticipating higher gains from additional customers. Therefore, higher-markup firms also invest more in their customer bases. Hence, the model is consistent with a positive relationship

between market shares and markups, as in conventional models, but has notably different macroeconomic implications.

We next calibrate the model to investigate these implications quantitatively. One of the key challenges in this analysis is to identify model parameters that determine the equilibrium allocation of customers across firms. To do so, we devise a strategy based on the model's predictions that we implement with available data on firms' sales and cost structures. At the core of this strategy is the comovement between a firm's sales and its non-production expenses (conditional on production expenses that control for confounding factors), which is informative of returns to scale in the customer acquisition technology.

With the calibrated model at hand, we ask how does the model change our understanding of the relationship between market share and market power? To answer this question, we compare the equilibrium allocation with the one obtained in a version of the model that corresponds to a specification of conventional models that is *recalibrated* to match the same moments, including the size distribution of firms.<sup>2</sup> Comparing the two, we find that our model associates the same moments with a higher degree of markup dispersion, which anticipates a higher degree of welfare loss due to misallocation.

Motivated by this higher markup dispersion, we characterize and quantify the efficient allocation in our model. Under this allocation, the social planner increases aggregate productivity by allocating more customers to more productive firms while equalizing the relative demand per customer across weakly substitutable varieties. This result contrasts with the efficient allocation in conventional models, in which the planner uses the intensive margin of demand to target two mutually exclusive objectives: concentrate demand among more productive firms to increase aggregate productivity versus equalize demand across varieties to eliminate utility losses from demand dispersion. In our model, this trade-off is nonexistent because the planner uses both margins of demand as instruments to achieve both objectives.

Even though our model features a margin of demand that is not necessarily associated with market power, welfare losses are potentially large. This result follows from the observation that the endogenous allocation of customers pushes the Pareto frontier of our economy beyond what models without this extensive margin would suggest. While

<sup>&</sup>lt;sup>2</sup>Since this model does not have the parameter that governs endogenous customer acquisition, it has one fewer parameter to calibrate. As a result, we drop the correlation between sales and non-production costs that our identification strategy relates to this parameter.

the uniform allocation of customers across firms is still feasible, the planner improves on this allocation by concentrating customers among more productive firms. Hence, welfare losses could be large if the equilibrium allocation of customers is sufficiently distorted. Thus, our analysis unveils a novel source of efficiency losses due to the *misallocation of customers*.

We find that the misallocation of demand has large negative effects on efficiency and welfare: The consumption equivalent welfare gains of the representative household under the efficient allocation is 13.6%. The majority of this gain comes from the efficiency gains in aggregate TFP under the planner's allocation, quantified at 10.8% higher than in the equilibrium. The planner achieves higher aggregate TFP by reallocating customers from low-productivity firms to the most productive ones. Indeed, in the efficient allocation, the top 5% sales share increases by almost 40%, and the number of operating firms declines by 11%. Finally, we verify that these results are mainly driven by customer misallocation, which in the equilibrium is determined by the degree of decreasing returns to advertising. To do so, we show that by moving halfway from the calibrated model to an economy with constant returns to advertising, the differences across allocations become much less pronounced. For example, compared with this alternative equilibrium, the efficient allocation generates only 3.2% higher TFP, 4% higher welfare, and 15% higher concentration.

Literature review Our paper is closely connected to the literature that emphasizes the macroeconomic significance of the customer margin in firm growth and market share (Foster, Haltiwanger, and Syverson, 2015, Hottman, Redding, and Weinstein, 2016). No-tably, recent research by Einav, Klenow, Levin, and Murciano-Goroff (2022) documents that about 80% of firms' sales variation arises from the customer margin, while Fitzgerald, Haller, and Yedid-Levi (2016), Argente, Fitzgerald, Moreira, and Priolo (2021) demonstrate that firms mainly grow through expansionary activities, rather than through reducing markups early in their life cycle. Our main contribution to this literature is to investigate the relationship between markups and different demand margins. Specifically, we document that once the market share is decomposed to the extensive and intensive margins, markups correlate only with firms' *average sales per customer*. Collectively, these facts paint a wholesome picture of firm growth: Firms mainly grow via the customer margin through non-price-related expansionary activities, with their markups tied only to their average sales per customer.

Based on these facts, our theoretical framework contributes to the literature on vari-

able markups by connecting models of firm growth through expansionary activities (e.g., Arkolakis, 2010) to models of endogenous markups at the intensive margin (e.g., Atkeson and Burstein, 2008).<sup>3</sup> Our contribution is to show that these two ingredients interact in a nontrivial way: Variable markups create differential incentives for firms to invest in their customer bases through non-price activities. These incentives lead to an equilibrium relationship between market shares and markups, but one that has different implications for the misallocation of resources relative to models that lack either ingredient. Finally, a notable recent work in this area is by Cavenaile, Celik, Perla, and Roldan-Blanco (2023), who provide microfoundations for the role of advertising in targeting different types of customers and showing how that leads to higher markups when advertising is more targeted.

Given our focus on misallocation, our paper is also related to the literature that analyzes the role of misallocation of production inputs across firms in affecting aggregate TFP (Restuccia and Rogerson, 2008, Hsieh and Klenow, 2009). In particular, our focus on the misallocation of customers *across* firms with variable markups relates our work to that of Bornstein and Peter (2022), who study misallocation of customers *within* firms, as well as Edmond, Midrigan, and Xu (2022), Peters (2020), who study the misallocation consequences of variable markups in settings without customer acquisition. We contribute to this literature by highlighting a new source of distortions in aggregate productivity that stems from the misallocation of customers across firms.

**Layout** Section 2 presents our empirical analysis. Section 3 describes the model. Section 4 discusses the model calibration. Section 5 quantifies the efficiency losses, and Section 6 concludes.

### 2 Motivating Facts

This section documents two new motivating facts using micro-level data that emphasize the importance of customer bases for firm dynamics and price-cost markups. Briefly, we document that price-cost markups are correlated only with average sales per customer and are unrelated to the size of firms' customer bases, even though these firms mainly grow through their customer bases.

<sup>&</sup>lt;sup>3</sup>For growth through expansionary activities, see also Drozd and Nosal (2012), Kaplan and Zoch (2020), Einav, Klenow, Levin, and Murciano-Goroff (2022), Argente, Fitzgerald, Moreira, and Priolo (2021). There is also an extensive literature on growth through pricing activities; see, e.g, Phelps and Winter (1970), Rotemberg and Woodford (1999). For the most recent work in this area, see Bornstein (2021) and its review of that literature.

#### 2.1. Data Description

We construct a detailed customer-firm-matched dataset to decompose firms' sales into the size of their customer bases and average sales per customer. Formally, we consider the following exact decomposition of log sales of firm *i*:

$$\ln S_i = \ln m_i + \ln \left( p_i q_i \right), \tag{2.1}$$

where  $m_i$  denotes the number of firm *i*'s customers,  $p_i$  its price, and  $q_i$  the average quantity purchased per customer.

We use the Nielsen Homescan Panel, which is one of the few sources of data that allows us to measure the number of each firm's customers.<sup>4</sup> The data contain approximately 4.5 million barcode-level product sales recorded from an average of 55,000 households per year in the United States. Nielsen samples households and provides in-home scanners so that households can record their purchases of products with barcodes. A barcode is a unique universal product code (UPC) allocated to each unique product and is used to scan and store its information. Each household is assigned a sample weight—or a projection factor—by Nielsen based on 10 demographic variables to make the sample nationally representative. Nielsen assigns a broad product-group label for each product, such as pet food and school supplies, and records information about the retailer a household visited to purchase products at a given time. According to Nielsen, the Homescan Panel covers approximately 30% of all household expenditures on goods in the consumer price index (CPI) basket. The data we use cover the period 2004-2016.

Next, we incorporate firm-level balance sheet information from Compustat to analyze firms' cost structures. With the caveat that this dataset covers only publicly listed firms, it constitutes the main source of panel data for firm-level analysis in the US and has been used in the recent literature on price-cost markups (see, for example, De Loecker, Eeckhout, and Unger, 2020, Traina, 2019). Throughout the analysis, we focus on two measures of a firm's costs. From an accounting perspective, a firm's costs associated with the running of the firm are captured in the Operating Expense (OPEX), which is divided into the Cost of Goods Sold (COGS, production costs) and Selling, General, and Administrative Expenses (SGA, non-production costs). According to Compustat, COGS includes "expenses that are directly related to the cost of merchandise purchased or the cost of goods manufactured that are withdrawn from finished goods inventory and sold

<sup>&</sup>lt;sup>4</sup>The dataset was made available by the Kilts Marketing Data Center at the University of Chicago Booth School of Business.

to customers". It records costs attributable to the production of the goods sold by a firm, and its typical categories are the cost of labor and intermediate inputs used in production. On the other hand, SGA expenses include "commercial expenses of operation (such as expenses not directly related to product production) incurred in the regular course of business...". This includes the costs incurred to sell and deliver products and services and the costs to manage the company; typical categories are advertising, marketing, shipping, and research and development, among others.

Finally, we combine the Nielsen database with GS1 US Data Hub to group individual products according to their producing firms and to merge in firm-level information from Compustat. GS1 is the business entity that provides barcodes for products and records the firm name for each UPC available in the Nielsen data. The definition of a firm is based on the unit that purchased barcodes from GS1. Therefore, a firm in our data corresponds to either a manufacturer or a retailer. This merge procedure provides a unique link between customer- and producer-level data for each firm.<sup>5</sup> Although we only have 332 firms in the Nielsen-Compustat matched dataset, these cover approximately 25% of total sales in Nielsen. Section SM.1 and Table SM.1.1 in Supplemental Materials provide detailed descriptions of the data-cleaning procedure and the merged dataset.

#### 2.2. The Relationship between Firm Size and Markups

We start by revisiting the predictions of a large class of models that relate firms' relative market power to their relative size. These models predict that larger firms charge higher markups (see Edmond, Midrigan, and Xu, 2022, Burstein, Carvalho, and Grassi, 2020, for supporting evidence). Since in most of these models every firm produces only one product, we can take this prediction to the data at either the firm level or the product level. We do both: We first present results at the firm level and then provide evidence that these results extend to the product level.

**Evidence Based on Firm-level Markups** Our decomposition of firms' sales in Equation (2.1) raises the following question: Which margin of sales captures the relationship between relative size and markups? To answer this question we first need to measure firm-level markups. We follow the methodology of De Loecker, Eeckhout, and Unger

<sup>&</sup>lt;sup>5</sup>We match the Nielsen-GS1 database with the Compustat database, following a procedure similar to Argente, Lee, and Moreira (2018). We use the "reclink" STATA software command based on the company name after standardizing it with the "std\_compname" command (Wasi and Flaaen 2015). Once Stata reports the matching rate for each observation, we keep those having higher than a 0.99 matching rate. We manually check the company name for every observation and drop inconsistent matches.

(2020), in which markups are measured by the inverse variable cost share of sales multiplied by the output elasticity of those variable inputs. Because this methodology does not require information on all variable costs, we follow De Loecker, Eeckhout, and Unger (2020) and use data on COGS from Compustat as a measure of variable costs. In addition, since we are interested in *relative* markups within industries at a given point in time, we absorb the output elasticity term with sector-year fixed effects.<sup>6</sup> Thus, our regression specification is given by:

$$\ln (\text{Sales/COGS})_{it} = \alpha_1 \ln p_{it} q_{it} + \alpha_2 \ln m_{it} + \lambda_{s,t} + \varepsilon_{it}, \qquad (2.2)$$

where  $(Sales/COGS)_{it}$  is the Sales-to-COGS ratio of firm *i* at time *t*. The sector-year fixed effects  $\lambda_{s,t}$  absorb all the variation at the sector-year level, which allows us to interpret the markup measure in *relative* terms and the size variables (average sales per customer and the number of customers) in terms of *market shares*.

	(1)	(2)	(3)	(4)	(5)
$\ln p_{\rm it} q_{\rm it}$	0.092***	0.091***	0.060***	0.059***	0.060**
	(0.033)	(0.033)	(0.022)	(0.022)	(0.024)
ln <i>m</i> <sub>it</sub>	-0.002	-0.002	0.002	0.002	0.003
	(0.006)	(0.006)	(0.007)	(0.007)	(0.007)
Observations	2433	2433	2433	2433	2433
$R^2$	0.046	0.047	0.311	0.313	0.338
Year FE		$\checkmark$		$\checkmark$	
SIC FE			$\checkmark$	$\checkmark$	
SIC-year FE					$\checkmark$

Table 1: Markups, Sales per Customer, and Number of Customers

*Notes:* \* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01; standard errors are clustered at the firm level. Markups are measured as the Sales-to-COGS ratio. The variable  $\ln p_{it}q_{it}$  denotes the log of the average sales per customer and  $\ln m_{it}$  the log number of customers. SIC industries correspond to a two-digit SIC code. All Nielsen variables are projection-factor adjusted.

Our first motivating fact is that *firms' markups are highly correlated with their average sales per customer but are unrelated to their number of customers*, as reported in Table 1.

<sup>&</sup>lt;sup>6</sup>The main challenge in the estimation of markups lies in estimating the output elasticity using only data on firms' revenues (see Bond, Hashemi, Kaplan, and Zoch, 2021). Given our focus on relative markups and our log regression specification, we do not need to estimate output elasticities. Instead, our regression specification incorporates a set of fixed effects that absorb these output elasticities. The underlying assumption we make, which is standard in the literature, is that the output elasticity with respect to COGS remains constant across firms within an industry and/or time.

Results are robust to including different combinations of year and sector fixed effects, which shows that the identified relationships hold within sectors by year.<sup>7</sup> With the inclusion of industry-time fixed effects, our results show that firms that charge higher markups have a higher market share in terms of average sales per customer; this is expected, given existing theories and previous evidence. However, the most novel result is that markups are not associated with firms' market shares in terms of the number of customers.<sup>8</sup> Thus, the relevant notion of relative size for markups is based on average sales per customer.<sup>9</sup>

**Evidence Based on Product-level Markups** Our results in Table 1 are limited by the scope of the Compustat dataset, which focuses on large public firms and firm-level markups. To demonstrate the external validity of these results, we perform two complementary analyses by using alternative measures of markups and sales. In the first approach, we measure the retailer-product-level markup as the difference between the retailer-product-level price available from the Nielsen Homescan Panel data and the wholesale cost obtained from the Nielsen PromoData, which is based on Gopinath, Gourinchas, Hsieh, and Li (2011) and Stroebel and Vavra (2019). In the second approach, we analyze the relationship between markups and each margin of demand by exploiting the rich dimensions of the data that allow us to control for marginal costs through an extensive set of fixed effects. This approach is similar to that of Fitzgerald, Haller, and Yedid-Levi (2016) and is valid under the assumption of common marginal costs across different subsets of observations. Both approaches include product fixed effects, which allow us to analyze the relationships within a given product. Appendix A.3 provides detailed descriptions of

<sup>&</sup>lt;sup>7</sup>The sector fixed effects allow us to show that firms that sell more per customer have higher price markups *within* their sector. The advantage of this firm-level analysis is that it allows us to measure relative markups directly. However, it raises the question of whether this relationship is coming from within-firm products or is due to compositional effects across products. Since this approach measures price markups at the firm level in a limited sample, we cannot include fixed effects for products available in the Nielsen data or for more disaggregated SIC codes. However, as discussed below, we use two alternative approaches to measure markups at the product (UPC) level and include more disaggregated product fixed effects in Appendix A.3. We find that the reported relationships hold within products, which supports the view that they do not arise from compositional effects.

<sup>&</sup>lt;sup>8</sup>One source of concern might be measurement error in a firm's customer base leading to attenuation bias. Two points alleviate this concern. First, the estimated coefficients are precisely estimated; i.e., we are finding a precise zero association. Second, in Appendix A.1, we use the lagged number of customers  $(\ln m_{it-1})$  as an instrument for  $\ln m_{it}$ , which provides consistent estimates under classical measurement error. We find a strong first stage and the point estimates are similar to those in Table 1.

<sup>&</sup>lt;sup>9</sup>Appendix A.2 presents results of an alternative specification that includes both total sales and average sales per customer as regressors. Given the decomposition in Equation (2.1), it is not surprising that markups are only associated with the average sales per customer but not with total sales.

both approaches.

Although we switch the focus to product-level markups and a broader sample of products and firms, both alternative approaches generate results consistent with those in Table 1: Markups are positively associated with average sales per customer, but not with the number of customers.

#### 2.3. The Role of Customer Base for Sales Growth

Armed with the empirical evidence that shows how each margin of demand is associated with markups, we now document that the main source of variation in firm sales is the variation in the number of customers rather than average sales per customer. Table SM.1.2 in Supplemental Materials presents summary statistics of the customer-firm matched data. A quick glance of the Nielsen-GS1 data already reveals that much of the firm-product group-year sales are driven by the number of customers, not by average sales per customer: More than 500,000 customers spend only approximately \$10 for each product group and firm per year on average.

To formalize this point, we follow Equation (2.1) and decompose the variance of log sales into the variances of log average sales per customer and log number of customers, as well as the covariance between these two components. Table 2 documents our second motivating fact: *Firms mainly grow by acquiring new customers instead of increasing their average sales per customer.* The number of customers accounts for approximately 80% of the variation in sales across firms. Average sales per customer, on the other hand, accounts for approximately 11% of the variance of sales, with the covariance accounting for the rest.<sup>10</sup> These results parallel the findings in contemporary work by Einav, Klenow, Levin, and Murciano-Goroff (2022). Using transaction-level data for a broad set of industries, they document that differences in customer bases account for 74% of sales variation across merchants. This shows that our finding extends beyond the consumer packaged goods sector and is representative of similar patterns in a wider set of industries.

<sup>&</sup>lt;sup>10</sup>Our decomposition results are similar when we instead use the first-difference of log sales; approximately 78%, 20%, and 2% of the variation are explained by the  $\Delta$  log number of customers,  $\Delta$  log average sales per customer, and the covariance between the two, respectively.

Table 2: Decomposing the Variance of Sales

Var(ln S <sub>igt</sub> )	$Var(\ln p_{igt}q_{igt})$	Var(ln $m_{igt}$ )	$2$ Cov(ln $p_{igt}q_{igt}$ , ln $m_{igt}$ )
7.5807	0.8672	6.1146	0.5989

*Notes:*  $S_{igt}$  denotes sales,  $p_{igt}q_{igt}$  average sales per customers, and  $m_{igt}$  the number of customers. We use 557,820 firm-group-year-level observations in Nielsen-GS1 data. All variables are projection-factor adjusted.

In addition to decomposing sales in the cross-section of firms, we find that the acquisition of new customers is also the main driver of firms' sales growth. As firms enter the economy, they can grow either by selling more per customer or by selling to more customers. To document this fact, we analyze firm growth patterns after entry. We mark a firm's entry as the time when it appears in our data for the first time. To be conservative, we drop entry events that occurred in the first 4 years in the dataset.<sup>11</sup> To quantify the importance of each margin, we estimate the following equation:

$$\ln S_{it} = \sum_{a=1}^{8} \delta_a \mathbf{1} (age_{it} = a) + \lambda_i + \lambda_t + \varepsilon_{it}, \qquad (2.3)$$

where  $S_{it}$  stands for sales and its components of firm *i* in year *t*,  $age_{it}$  is the number of years firm *i* stayed in the economy after entry in year *t*, and  $\lambda_i$  and  $\lambda_t$  are the firm and year fixed effects, respectively (see Argente, Lee, and Moreira (2019) for a similar analysis of the life-cycle of individual products). The parameters of interest are  $\delta_a$ , which measure the dynamics of average sales and its components over the life-cycle of the firm.

Figure 1 plots the log sales as a function of firm age ( $\hat{\delta}_a$  as a function of *a*) and decomposes it into the log number of customers and log average sales per customer. As a firm's sales grow over time, both margins of demand also increase. However, regardless of the firm's age, sales growth is mostly attributed to the increase in the number of customers. At age 1, differences in the number of customers explain approximately 78% of differences in sales, whereas average sales per customer explain approximately 22% of sales. Although the importance of the number of customers decreases as firms become older, on average, it still accounts for approximately 70% of sales for the maximum firm age observed in the data. Results are robust to including only those firms that survive at least 3 or 5 consecutive years and analyzing average monthly sales as the dependent variable, which accounts for staggered entry throughout the year. Since the degree of

<sup>&</sup>lt;sup>11</sup>There is a large increase in the number of households and firms in the Nielsen Homescan Panel data in the years 2006 and 2007. We drop the years 2004-2007 to render the analysis conservative. Thus, the maximum firm age in our sample of entrants is 8 years.



Figure 1: Decomposition of Firm Sales Growth by Firm Age

*Notes*: This figure plots the average firm sales, sales per customer, and the number of customers for each firm-age based on Equation (2.3), after controlling for firm and year fixed effects. The blue circled line shows the results for log sales, the red diamond line for the log average sales per customer, and the green triangle line for the log number of customers. There are 40,442 observations and 9,990 firms that newly enter the economy starting from the year 2008 in the Nielsen-GS1 data. All estimates are normalized based on age 0. All variables are projection-factor adjusted.

durability of a product might affect the ability of firms to grow through different margins, we repeat the analysis by splitting products according to their durability and find similar patterns within both subsamples. See Appendix A.4 for further details.

The fact that firms mainly grow through the extensive margin, which is not associated with their markups, begs the question: To what extent do firms control their growth through different sales margins? To investigate this, in Appendix A.5 we use data on SGA expenses, which in part capture firms' expenditures on expansionary activities, and find that (1) firms that spend more non-production costs have higher sales, and (2) these costs are associated with the number of new customers firms acquire, but not with the number of customers they retain or average sales per customer. In summary, the evidence shows that firms' SGA expenses contribute to relative firm size only through customer acquisition.

All these facts together suggest that by spending on expansionary activities, firms grow through a margin of demand that *does not* correlate with markups—a distinction not previously studied in the macroeconomics literature. In the next section, we develop a model that can account for these facts and explore its implications.

# 3 Model

In this section, we present a model with variable markups and endogenous customer acquisition that is consistent with our motivational facts in Section 2. We then use this model to quantitatively investigate the implications of endogenous customer acquisition for markup dispersion, misallocation, and welfare.

#### 3.1. Setup

Time is discrete and is indexed by  $t \in \{0, 1, 2, ...\}$ . There is a representative household with a continuum of individual members denoted by  $j \in [0, 1]$ . A continuum of firms, indexed by  $i \in N_t$ , produce weakly substitutable goods in a representative industry. With slight abuse of notation, we use  $N_t$  to denote both the set and the measure of these firms.

**3.1.1.** Households. The representative household supplies labor in a competitive labor market and demands different varieties produced by firms at given prices. We let  $m_{i,t}$  denote both the measure and the set of a variety *i*'s customers and write  $j \in m_{i,t}$  when member *j* is a customer of *i*. Household members jointly maximize their utility when the utility of their consumption baskets is aggregated by a *Kimball aggregator*, subject to their budget constraint:

$$\max_{\{C_t, L_t, (c_{i,j,t})\}_{j \in m_{i,t}}^{i \in N_t}\}} \sum_{t=0}^{\infty} \beta^t \left[ \frac{C_t^{1-\gamma}}{1-\gamma} - \xi \frac{L_t^{1+\psi}}{1+\psi} \right]$$
(3.1)

$$s.t.\int_{0}^{N_{t}}\int_{0}^{1}\mathbf{1}_{\{j\in m_{i,t}\}}\Upsilon\left(\frac{c_{i,j,t}}{C_{t}}\right)djdi = 1$$
(3.2)

$$\int_{0}^{N_{t}} \int_{0}^{1} p_{i,t} c_{i,j,t} dj di \leq W_{t} L_{t} + \int_{0}^{N_{t}} \Pi_{i,t} di - T_{t}.$$
(3.3)

Here  $c_{i,j,t}$  is the consumption of member *j* from variety *i*;  $p_{i,t}$  is the price of variety *i*;  $C_t$  is aggregate consumption;  $L_t$  is the total labor supply;  $W_t$  is the wage;  $\Pi_{i,t}$  is the profit of *i*; and  $T_t$  is a lump-sum tax. The function  $\Upsilon(.)$  is strictly increasing and concave with  $\Upsilon(1) = 1$ .<sup>12</sup> It follows that all customers of *i* choose to purchase the same amount implied by the following demand function:

$$\frac{c_{i,j,t}}{C_t} = q_{i,t} \equiv \Upsilon'^{-1} \left( \frac{p_{i,t}}{P_t D_t} \right) \mathbf{1}_{\{j \in m_{i,t}\}}.$$
(3.4)

Here  $q_{i,t}$  denotes the relative demand per matched customer of variety *i*. Moreover,

$$D_t \equiv \left[\int_{i \in N_t} \int_0^1 \mathbf{1}_{\{j \in m_{i,t}\}} \frac{c_{i,j,t}}{C_t} \Upsilon'\left(\frac{c_{i,j,t}}{C_t}\right) dj di\right]^{-1}$$

<sup>&</sup>lt;sup>12</sup>In the case of the CES aggregator,  $\Upsilon(x) = x^{1-\sigma^{-1}}$ , where  $\sigma$  is the elasticity of substitution across varieties.

is an *aggregate demand index* and  $P_t$  is the price of the aggregate consumption good, which, henceforth, we normalize to one.<sup>13</sup> The homogeneity of  $q_{i,t}$  across all customers of firms follows from the homogeneity of preferences. In Section SM.2 in Supplemental Materials, we introduce an extension of the model with heterogeneity in tastes and provide motivational evidence for why we abstract from them. Therefore, the household's total demand for variety *i* is *proportional* to the number of its customers, and given by the demand function:

$$c_{i,t} \equiv \int_0^1 \mathbf{1}_{\{j \in m_{i,t}\}} c_{i,j,t} dj = m_{i,t} q_{i,t} C_t$$
  
$$\implies \ln(S_{i,t}) \equiv \ln(p_{i,t} c_{i,t}) = \ln(m_{i,t}) + \ln(p_{i,t} q_{i,t}) + \ln C_t$$
(3.5)

where  $c_{i,t}$  is total demand for variety *i*, and  $q_{i,t}$  is the relative demand per customer in Equation (3.4). Also, the expression on the right shows how the model delivers the same decomposition of sales to the number of customers and demand per customer as in our empirical analysis in Equation (2.1).<sup>14</sup>

For the functional form of  $\Upsilon(.)$ , we use the Kimball aggregator of Klenow and Willis (2016):

$$\Upsilon(q) = 1 + (\sigma - 1)e^{\frac{1}{\eta}}\eta^{\frac{\sigma}{\eta} - 1} \left[ \Gamma\left(\frac{\sigma}{\eta}, \frac{1}{\eta}\right) - \Gamma\left(\frac{\sigma}{\eta}, \frac{q^{\frac{\eta}{\sigma}}}{\eta}\right) \right],$$
(3.6)

where  $\sigma > 1$  and  $\eta > 0$  control the demand elasticities and super-elasticities, as we discuss below, and  $\Gamma(.,.)$  is the incomplete Gamma function.<sup>15</sup> This specification for  $\Upsilon(.)$  is a generalization of a CES aggregator with substitution elasticity  $\sigma$ , which is nested when  $\eta = 0$ . Using this functional form in Equation (3.4), we obtain the following relative demand per customer for firm *i* at time *t*:

$$q_{i,t} = \left[1 - \eta \ln\left(\frac{p_{i,t}}{D_t(1 - \sigma^{-1})}\right)\right]^{\frac{\nu}{\eta}}.$$
(3.7)

Intuitively, this demand function is a smoothed version of a kinked demand curve (Dotsey and King, 2005, Basu, 2005), in which a customer's demand is more price sensitive at

$$\int_0^{N_t} m_{i,t} \Upsilon \left( \Upsilon'^{-1} \left( \frac{p_{i,t}}{P_t D_t} \right) \right) di = 1, \qquad \int_0^{N_t} m_{i,t} \frac{p_{i,t}}{P_t} \Upsilon'^{-1} \left( \frac{p_{i,t}}{P_t D_t} \right) di = 1.$$

<sup>14</sup>In the empirical analysis, the  $\ln C_t$  term would be absorbed by the industry fixed effects.

<sup>&</sup>lt;sup>13</sup>In the special case where the aggregator is CES, this demand index takes a value of  $1/(1 - \sigma^{-1})$ ; however, with the generalized Kimball aggregator this quantity is not necessarily a constant. Moreover, we could characterize the equations that pin down  $P_t$  and  $D_t$  in terms of prices. The equations that determine  $P_t$  and  $D_t$  are:

<sup>&</sup>lt;sup>15</sup>The incomplete Gamma function is given by  $\Gamma(s, x) \equiv \int_{x}^{\infty} t^{s-1} e^{-t} dt$ .

higher relative prices—i.e., demand satisfies Marhalls's second law of demand. As we show below, firms with larger demand per customer face lower elasticities and charge higher markups. To observe this, let us consider the demand elasticity  $\varepsilon_{i,t}$ , and super-elasticity,  $\varepsilon_{i,t}^{\varepsilon}$ :

$$\varepsilon_{i,t} \equiv -\frac{\partial \ln(c_{i,t})}{\partial \ln(p_{i,t})} = \sigma q_{i,t}^{-\frac{\eta}{\sigma}}, \quad \varepsilon_{i,t}^{\varepsilon} \equiv \frac{\partial \ln(\varepsilon(q_{i,t}))}{\partial \ln(p_{i,t})} = -\frac{\eta}{\sigma} \varepsilon_{i,t} \le 0, \tag{3.8}$$

where we see that demand elasticity  $\varepsilon_{i,t}$  is a decreasing function of the relative demand per customer, formally shown by the negative sign of the super-elasticity  $\varepsilon_{i,t}^{\varepsilon}$  as long as  $\eta > 0$ .

Finally, to conclude households' optimality conditions, the household's labor supply is characterized by the following standard intratemporal Euler equation:  $\xi L_t^{\psi} = W_t C_t^{-\gamma}$ .

**3.1.2.** Dynamics of Customer Bases. We model the dynamics of customer bases following our empirical findings in Section 2. In particular, as we document in Appendix A.5, firms' SGA expenses correlate with the acquisition of new customers (but not the retention of old customers).

Motivated by this evidence, we assume that firms can engage in expansionary activities, such as advertising campaigns or increasing the availability of their goods, to attract *new customers*. In addition, two processes in the model separate customers from firms: at the end of each period, (1) all customers of exiting firms separate, and (2) customers of incumbent firms separate at an exogenous rate of  $\delta \in [0, 1]$ . We assume that the total mass of matches is fixed over time and, without loss of generality, normalize this mass to one. This implies that while expansionary activities affect the distribution of customers across firms, it does not increase the total number of customers who buy from an industry (as in Einav, Klenow, Levin, and Murciano-Goroff, 2022). We view this as a conservative benchmark, as we discuss in Section 5.

As for the dynamics of new matches, we assume that operating firm *i* at time *t* posts  $a_{i,t} \ge 0$  ads to acquire new customers. Every unmatched member then draws an ad from the pool of all available ads and is matched to the firm they draw. Therefore, the number of new customers firm *i* acquires at time *t* is proportional to the number of ads it posted relative to the total number of ads posted by all firms.<sup>16</sup> Hence, firm *i*'s customer base

<sup>&</sup>lt;sup>16</sup>Note that we have modeled customer acquisition through expansionary activities that operate independent of firms' pricing decisions, and thus abstract away from customer acquisition through lower markups/prices. This is motivated by the evidence in Fitzgerald, Haller, and Yedid-Levi (2016), who conclude that firms do not manipulate prices to shift demand by documenting that after entering a new market, their markups remain the same while their quantities grow. Our findings in Table 1 and Appendix A.5 also

evolves according to

$$m_{i,t} \le (1-\delta)m_{i,t-1} + \frac{a_{i,t}}{P_{m,t}},\tag{3.9}$$

where the inequality captures the notion that there is free disposal of customers, should the firm choose to exercise that option. Moreover,  $P_{m,t}$  is the endogenous conversion rate of ads to customers. This is the number of ads needed for a firm to get one new customer, determined in the equilibrium to clear the matching market:

$$\int_{i \in N_t} m_{i,t} di = 1 \Longrightarrow P_{m,t} = \frac{\int_{i \in N_t} a_{i,t} di}{1 - (1 - \delta) \int_0^{N_t} m_{i,t-1} di}.$$
(3.10)

This expression shows that the cost of a match decreases with the total number of separated customers and increases with the total number of posted ads by all firms.

**3.1.3. Firms.** On the firm side, we assume endogenous entry and exit with an order of events as summarized in Figure 2. We provide a detailed description of these decisions below.

**Entry and Exit Decisions** At each period *t*, a measure  $\lambda$  of potential entrants are born, each with an initial productivity  $z_{i,t}$  drawn from a log-normal distribution:

$$\ln(z_{i,t}) \sim \mathcal{N}(\bar{z}_{ent}, \sigma_z^2). \tag{3.11}$$

We let  $\Lambda_t$  denote the set of these potential entrants at *t*. Incumbents—i.e., firms that entered at least one period ago—also draw new productivities according to the following AR(1) process:

$$\ln(z_{i,t}) = \rho \ln(z_{i,t-1}) + \epsilon_{i,t}, \quad \epsilon_{i,t} \sim \mathcal{N}(0, \sigma_z^2). \tag{3.12}$$

With new productivities drawn, each incumbent or potential entrant then decides whether to stay in the economy or to drop out.<sup>17</sup> We refer to this decision by  $\mathbf{1}_{i,t} \in \{0,1\}$ , with 1 being an indicator for entering or staying. Finally, all the incumbent firms who decided to stay draw i.i.d. Bernoulli survival shocks,  $v_{i,t}$ , that are equal to 1 with probability  $v \in [0,1]$ , and drop out if  $v_{i,t} = 0$ .

support the importance of customer acquisition through expansionary activities by showing that markups are not correlated with the size of firms' customer bases conditional on sales per customer, but expenses related to expansionary activities are.

<sup>&</sup>lt;sup>17</sup>Following Clementi and Palazzo (2016) and Ottonello and Winberry (2020), we allow the mean of incumbents' productivity distribution—normalized to 0—to differ from that of entrants,  $\bar{z}_{ent}$ . This introduces a natural trend in firms' productivity based on their age and allows us to account for differences in size across age groups, as reported in the Business Dynamics Statistics (BDS).

#### Figure 2: Order of Events



Notes: The figure shows the timing of firms' decisions in the model.

**Expansionary Activities, Pricing, and Production** Firms that stay or enter the economy pay an overhead cost of  $\chi > 0$  in units of labor at each period. Also, firms use labor to produce ads using the technology  $a_{i,t} = l_{i,s,t}^{\phi} \ge 0$ , where  $l_{i,s,t}$  denotes the amount of labor allocated to advertising activities. The firm's customer base then evolves according to Equation (3.9), where  $m_{i,t-1} \equiv 0$  for firms that entered at time *t*. Moreover,  $\phi \in [0, 1]$  is the degree of decreasing returns to advertising. Once new customers are acquired, firms' demands are realized as in Equation (3.5). Taking this demand as given, firms then choose the prices. Each firm *i* then produces to meet its realized demand using the production function  $y_{i,t} = z_{i,t} l_{i,p,t}^{\alpha}$ , where  $z_{i,t}$  is the firm's productivity, and  $l_{i,p,t}$  is its labor demand for production. Finally,  $\alpha \in [0, 1]$  is the degree of decreasing returns to productivity.

**Firms' Problem** Given an initial level of productivity and customer base, firm *i*'s problem is given by

$$v_{t}(m_{i,t-1}, z_{i,t}) \equiv \max_{\substack{(q_{i,\tau}, l_{i,s,\tau}, \\ l_{i,p,\tau}, \mathbf{1}_{i,\tau})_{\tau=t}^{\infty}}} \mathbb{E}_{t} \sum_{\tau=t}^{\infty} (\beta v)^{\tau-t} \left( \prod_{h=t}^{\tau} \mathbf{1}_{i,h} \right) \left( \frac{C_{\tau}}{C_{t}} \right)^{-\gamma} \left[ \underbrace{\underbrace{\mathcal{D}_{\tau} \Upsilon'(q_{i,\tau}) y_{i,\tau}}_{\text{total sales}} - \underbrace{W_{\tau} l_{i,p,\tau}}_{\text{COGS}} - \underbrace{W_{\tau}(l_{i,s,\tau} + \chi)}_{\text{SGA expenses}} \right]$$
(3.13)

subject to 
$$y_{i,\tau} = m_{i,\tau} q_{i,\tau} C_{\tau} = z_{i,\tau} l_{i,p,\tau}^{\alpha}$$
 (3.14)

$$m_{i,\tau} \le (1-\delta)m_{i,\tau-1} + \frac{l_{i,s,\tau}^{\phi}}{P_{m,\tau}}, \quad l_{i,s,t} \ge 0.$$
(3.15)

The problem states that firm *i* maximizes the expected discounted stream of its profits subject to its demand curve in Equation (3.14), the law of motion for customers in Equation (3.15), and the nonnegativity of labor allocated to advertising. We have also labeled the terms in the profit function as *total sales*, *COGS*, and *SGA expenses*, which we use later to map the model to the data.

#### 3.2. Characterization of Firms' Decisions

In this section, we characterize the firms' optimal decision rules for pricing, expansionary activities, and entry and exit. All proofs are in Appendix **B**.

**Prices and Markups** Conditional on decisions for entry/exit and customer acquisition, firms' pricing decisions have a static nature. Formally, firm *i*'s optimal price at *t* is:

$$p_{i,t} = \underbrace{\frac{\varepsilon_{i,t}}{\varepsilon_{i,t} - 1}}_{\text{markup}} \times \underbrace{\alpha^{-1} \frac{W_t l_{i,p,t}}{y_{i,t}}}_{\text{marginal cost}}.$$
(3.16)

This expression shows that despite the presence of variable customer acquisition costs, the usual relationship between the labor share and the markup also holds in this model. This verifies our use of De Loecker, Eeckhout, and Unger (2020) methodology to identify markups from the Compustat data. Moreover, it is also important to note that the firm's elasticity of demand,  $\varepsilon_{i,t}$ , is itself a function of demand per customer in Equation (3.7) and varies with the firm's pricing choice. Therefore, as long as  $\eta$  is not zero, the optimal markup of the firm varies with its marginal cost, which leads to the following lemma.

**Lemma 1.** At a given time *t*, firms with higher marginal costs charge higher prices and lower markups. Formally, let  $\mu_{i,t}$  denote a firm's markup and  $mc_{i,t}$  its marginal cost. Then, the elasticities of markups and prices to marginal costs are:

$$\frac{\partial \ln(p_{i,t})}{\partial \ln(mc_{i,t})} = \frac{1}{1 + \eta \sigma^{-1} \varepsilon_{i,t}(\mu_{i,t} - 1)} \ge 0$$
(3.17)

$$\frac{\partial \ln(\mu_{i,t})}{\partial \ln(mc_{i,t})} = -\frac{\eta \sigma^{-1} \varepsilon_{i,t}(\mu_{i,t}-1)}{1 + \eta \sigma^{-1} \varepsilon_{i,t}(\mu_{i,t}-1)} \le 0.$$
(3.18)

Equation (3.17), which is also known as the *incomplete pass-through* property of Kimball demand (see, e.g., Gopinath and Itskhoki, 2010, Amiti, Itskhoki, and Konings, 2019), shows that a 1% increase in the marginal cost of a firm increases their price by less than 1%. The intuition is that firms with higher marginal costs need to charge higher relative prices to keep their positive margin. But at such new prices, demand is more elastic and optimal markups are lower, in contrast to a CES demand in which demand elasticities are constant and pass-through is complete. With these variable markups, we obtain the following proposition.

**Proposition 1.** At any *t*, firms with *higher relative sales per customer* charge *higher markups*:

$$\frac{d\ln(\mu_{i,t})}{d\ln(p_{i,t}q_{i,t})}\Big|_{t} = \eta \sigma^{-1} \mu_{i,t}(\mu_{i,t}-1) \ge 0$$
(3.19)

Proposition 1 shows that firms with higher sales per customer charge higher markups, which links the model to our empirical fact on this relationship in Table 1. Intuitively, firms with higher sales per customer must have lower relative prices and thus can charge higher markups because demand is less elastic at such prices. Note, however, that firms with higher sales per customer do not necessarily have to be larger in terms of total sales, since the size in our model also depends on the number of firms' customers. Due to the dynamic evolution of the customer base, for any given market share there is a distribution of firms with different sales per customer holding that market share. As a result, fixing market share, firms with more customers must also sell less per customer, and as a corollary of Proposition 1, charge lower markups.

**Corollary 1. Conditional on relative total sales**, firms with more customers charge lower markups:

$$\frac{d\ln(\mu_{i,t})}{d\ln(m_{i,t})}\Big|_{\frac{p_{i,t}y_{i,t}}{\int_{i\in N_t}p_{i,t}y_{i,t}di},t} = -\eta\sigma^{-1}\mu_{i,t}(\mu_{i,t}-1) \le 0$$
(3.20)

Corollary 1 highlights the main departure of our paper from the literature on variable markups, in which customer acquisition is not modeled explicitly and customers are *homogeneously* distributed across firms. In such models, the fact that all firms are on the same demand curve implies that relative size (market share) is a sufficient statistic for firms' markups. However, in our model, firms can hold the same relative size either because they sell more per customer or because they have more customers, which implies that relative size is no longer a sufficient statistic for market power.

But what determines the *unconditional* relationship between relative size and markups in this model? To answer this question, we need to characterize firms' optimal expansionary activities.

**Optimal Expansionary Activities** A key feature of our model is that firms internalize the decision to acquire customers. For this decision, while the marginal cost of a new customer is determined by the amount of labor the firm needs to employ to find it, its benefit is closely linked to the firm's market power and the amount of additional demand from that customer. The following proposition formulates the optimality condition for firms' advertising decisions in terms of this cost-benefit analysis. Proposition 2. The optimal customer acquisition strategy of a firm is characterized by

$$\underbrace{\phi^{-1} \frac{W_t l_{i,s,t}}{m_{i,t} - (1 - \delta) m_{i,t-1}}}_{\text{marginal cost of a new customer}} = \mathbb{E}_t \sum_{\tau=t}^{\infty} \underbrace{\left[ (v(1 - \delta))^{\tau - t} \prod_{h=t}^{\tau} \mathbf{1}_{i,\tau} \right]}_{\text{probability of match survival}} \times \underbrace{\beta^{\tau - t} \left( \frac{C_\tau}{C_t} \right)^{-\gamma} (\mu_{i,\tau} - 1) mc_{i,\tau} q_{i,\tau} C_\tau}_{\text{discounted (gross) marginal profit per customer}}$$
(3.21)

Equation (3.21) shows that firms optimally equate the cost of acquiring the marginal customer to the net present value of the gross profits earned from them for the duration of the match. Since the marginal profits generated by a new customer are increasing in the firm's markups, firms that charge higher markups (or expect to charge higher markups on average for the duration of a match) anticipate a higher return on investing in their customer base. Hence, our model predicts a positive *but endogenous* relationship between markups and the size of firms' customer bases in the equilibrium. This is in contrast to a model with an exogenous customer base or a representative household in which any relationship between markups and relative size is dictated by the shape of the exogenous demand curves. The endogenous nature of this relationship in a model with customer acquisition hints at its different counterfactual implications, which we will discuss in later sections.

**Entry and Exit Policies** A potential entrant enters and an incumbent stays if their value, specified in Equation (3.13), is positive:  $v_t(m_{i,t-1}, z_{i,t}) > 0$ . It follows that for any  $m_{-1}$ , there is a threshold  $z^*(m_{-1})$  such that firms with higher productivity than  $z^*(m_{-1})$  stay or enter (see Hopenhayn, 1992).

#### 3.3. Equilibrium

An equilibrium is defined as (a) an allocation for the households  $\{(c_{i,j,t})_{j\in[0,1]}, C_t, L_t\}_{t\geq 0}$ ; (b) a set of exit decisions for potential entrants and incumbents  $\{(\mathbf{1}_{i,t})_{i\in\Lambda_t\cup N_{t-1}}\}_{t\geq 0}$ ; (c) an allocation for operating firms  $\{(p_{i,t}, y_{i,t}, m_{i,t}, l_{i,p,t}, l_{i,s,t})_{i\in N_t}\}_{t\geq 0}$ ; and (d) a sequence of aggregate prices  $\{W_t, P_{m,t}\}_{t\geq 0}$  and a sequence of sets  $\{N_t\}_{t\geq 0}$  such that

- 1. given (c) and (d), household's allocation in (a) solves their problem in Equation (3.1),
- 2. given (a) and (d), firms' allocations in (b) and (c) solve their problems in Equation (3.13),
- 3. labor and matching markets clear:  $L_t = \int_{i \in N_t} (l_{i,p,t} + l_{i,s,t} + \chi) di$ ,  $1 = \int_{i \in N_t} m_{i,t} di$ ,
- 4. the set of operating firms,  $N_t$ , evolves according to

$$N_t = \{i \in \Lambda_t \cup N_{t-1} : \mathbf{1}_{i,t} v_{i,t} = 1\}, N_{-1} \text{ given.}$$
(3.22)

**Solution Method** We solve the model globally by combining collocation methods and nonstochastic simulation to approximate the distribution of firms over  $(z_{i,t}, m_{i,t-1})$ . Section SM.3 in Supplemental Materials describes the recursive formulation of the firm's problem and the computational algorithm that finds the steady state of this economy.

# 4 Quantitative Analysis

To quantify the implications of customer allocation with variable markups, we calibrate the steady state of the model using the Simulated Method of Moments and matching several micro and macro moments related to firm dynamics in the US economy.

#### 4.1. Calibration Strategy

To provide an overview of our calibration strategy, the new and most relevant parameter to calibrate is  $\phi$ , the returns to scale in advertising—which influences the relationship between a firm's relative size and market power by determining the size of its customer base.<sup>18</sup> As we discuss below, this parameter is identified by the relationship between sales and firms' expansionary activities included within the SGA expenses. The fact that the latter is only available from Compustat poses a challenge because it contains only a subset of the firms in the economy. To address this challenge, whenever possible, we calibrate the model to the aggregate US economy in 2012 by matching moments from the Business Dynamics Statistics (BDS) and Statistics of US Businesses (SUSB) provided by the Census Bureau. When matching moments based on data from Compustat, we apply a filter to the model-simulated data to account for the selection into Compustat based on firm size and age.<sup>19</sup>

**Fixed Parameters** We set the length of a period to 1 year. Panel A of Table 3 presents the set of parameters that are externally fixed. We set the subjective discount factor  $\beta$  to 0.96. The elasticity of intertemporal substitution  $\gamma$  is set to 2. We set the inverse of the Frisch elasticity of labor supply to  $\psi = 1$  and the labor coefficient in the production function to  $\alpha = 0.64$ . In the calibration exercise, we normalize the measure of potential entrants  $\lambda$ 

<sup>&</sup>lt;sup>18</sup>Note that our model does not require firms to spend on customer acquisition and grow through the extensive margin of demand, and nests the conventional model with exogenous customer bases as a special case. When  $\phi \rightarrow 0$ , every firm receives the same flow of customers in each period, without having to spend on  $l_{i,s,t}$ . If, in addition,  $\delta = 1$ , then all firms have the same stock of customers in every period.

<sup>&</sup>lt;sup>19</sup>That is, to compute equivalent moments, we restrict the simulated sample of firms to those that are at least 7 years old, as in Ottonello and Winberry (2020), and have sales above 19% of the average sales in the simulated economy (which corresponds to the ratio of the 5th percentile of the sales distribution in Compustat to average sales in SUSB). Average firm sales in the 2012 US economy were USD5.7 million (SUSB) and the 5th percentile of the sales distribution in Compustat was USD1.06 million.

and the disutility of labor supply  $\xi$  to generate a steady-state output of Y = 1 and wage of W = 1.

We set the retention rate of customers to  $1 - \delta = 0.72$ , which corresponds to the repurchase probability of customers in the Nielsen-GS1 matched dataset in 2012.<sup>20</sup> This value of the repurchasing probability is similar in other industries based on evidence from the marketing literature.<sup>21</sup>

**Calibrated Parameters** We jointly calibrate the remaining 8 parameters by the simulated method of moments (SMM).<sup>22</sup> These parameters can be grouped in three sets: those shaping firms' cost structure ( $\phi$  and  $\chi$ ), their demand ( $\sigma$ , and  $\eta$ ), and their life cycle and shock structure ( $\rho_z, \sigma_z, \bar{z}_{ent}$  and v). Although these parameters are jointly identified by all moments, we provide below a discussion of which moment should intuitively be more relevant to identify each parameter. We formalize this discussion in Section SM.4.1 in Supplemental Materials by analyzing the local elasticities of model moments with respect to each parameter and the sensitivity measure developed by Andrews, Gentzkow, and Shapiro (2017).

To calibrate the overhead cost  $\chi$ , we target the cross-sectional average COGS-to-OPEX ratio. The model counterpart of this ratio for firm *i* is  $\frac{W_t l_{i,p,t}}{W_t l_{i,p,t}+W_t} \equiv \frac{COGS_{i,t}}{COGS_{i,t}+SGA_{i,t}}$ . Intuitively, a larger fixed cost  $\chi$ , ceteris paribus, should increase a firm's total costs and drive down this ratio. To measure this ratio, we use data from Compustat.

To identify the elasticity  $\phi$ , we exploit the observed relationship between SGA and sales in Compustat. Proposition 3 illustrates the source of identification in the special case with  $\delta = 1$ .

<sup>&</sup>lt;sup>20</sup>More specifically, define Sales<sub>*i*,*g*,*t*</sub> as the total expenditure of (projection-factor adjusted) households that purchase products made by firm *i* in group *g* at time *t*. Define the probability of repurchasing firm's products as  $s_{i,g,t} = \frac{\text{Sales}_{i,g,t-1,t}}{\text{Sales}_{i,g,t-1}}$ , where  $\text{Sales}_{i,g,t-1,t}$  is the total expenditure of (projection-factor adjusted) households who purchase products made by firm *i* in group *g* in both periods t - 1 and *t*. Then, we take a weighted average of  $s_{i,g,t}$  across firms and groups, in which the weights are the expenditure in firm-group bins across all years.

<sup>&</sup>lt;sup>21</sup>For example, the repurchase probability is 0.7 in the automotive industry based on survey data used by Mittal and Kamakura (2001). According to Bolton, Kannan, and Bramlett (2000), the loyalty program member share is 0.693 and the cancellation probability is 0.187 for the financial services industry. Finally, Bornstein (2021) estimates an annual retention probability of 0.85 for the top two largest firms in each product category from the Nielsen data. If we also restrict the sample to the top two firms, our retention measure increases to 0.84.

<sup>&</sup>lt;sup>22</sup>More specifically, we calibrate the model by choosing a set of parameters  $\mathcal{P} = (\phi, \chi, \sigma, \eta, \rho_z, \sigma_z, \bar{z}_{ent}, v)$  that minimizes the SMM objective function  $\left(\frac{m_m(\mathcal{P})}{m_d} - 1\right)' W\left(\frac{m_m(\mathcal{P})}{m_d} - 1\right)$ , where  $m_m$  and  $m_d$  are a vector of model-simulated moments and data moments, respectively, and W is a diagonal matrix. Section SM.3.2 in Supplemental Materials provides the computational details of the calibration exercise.

**Proposition 3.** Suppose  $\delta = 1$ . Then, the total  $SGA_{i,t}$  expenses of a firm can be decomposed into a fixed ( $SGAF_{i,t}$ ) and a variable ( $SGAV_{i,t}$ ) component:

$$SGA_{i,t} = SGAF_{i,t} + SGAV_{i,t} = W_t \chi + \phi Sales_{i,t} - \frac{\phi}{\alpha} COGS_{i,t}$$
(4.1)

Equation (4.1) is obtained from the firm's optimality condition Equation (3.21) regarding customer acquisition, which shows that firms' expenditures on customer acquisition are increasing in their markups because they directly determine the returns from customer acquisition. Since  $\phi$  is the returns to scale in acquiring new customers, it naturally follows that it should be identified from the sensitivity of SGA expenditures to firms' markups, which in the model are proportional to firms' Sales-to-COGS ratios. Proposition **3** formalizes this intuition in the special case of  $\delta = 1$ , and in Section SM.4 in Supplemental Materials we find a high sensitivity of  $\phi$  to the same relationship; this indicates that our intuition also carries on for the general case with  $\delta < 1$ . Thus, we identify the sensitivity of SGA expenses to markups by regressing firms' SGA expenses on their sales while controlling for COGS and time fixed effects (which capture the fixed components of SGA expenses). We calibrate  $\phi$  so that the regression coefficient on *Sales*<sub>i,t</sub> in the model-simulated data matches the regression coefficient obtained from Compustat.

To calibrate the parameters that shape firms' demand, we set the elasticity of substitution  $\sigma$  to match a COGS-weighted average markup of 1.25 computed from Compustat as in Edmond, Midrigan, and Xu (2022). While the level of markups identifies  $\sigma$ , the sensitivity of markups to relative size identifies  $\eta$ . Since we do not directly observe markups, following Edmond, Midrigan, and Xu (2022), we pin down the super-elasticity of demand  $\eta$  by the relationship between a firm's relative revenue productivity of labor and its relative sales. In a model without customer acquisition, the revenue productivity of labor  $p_{i,t}y_{i,t}/W_t l_{i,p,t}$  is directly proportional to the production markup  $\mu_{i,t}$ , which validates this approach. The following proposition confirms that a similar relationship holds in our model (see Section SM.4 in Supplemental Materials for the numerical mapping between  $\eta$  and this relationship).

**Proposition 4.** Suppose  $\delta = 1$ . Then, a firm's average revenue productivity of labor is given by

$$\frac{p_{i,t}y_{i,t}}{W_t(l_{i,p,t}+l_{i,s,t})} = \frac{\mu_{i,t}}{\alpha + \phi(\mu_{i,t}-1)}$$

which is strictly increasing in the production markup  $\mu_{i,t}$  if and only if  $\alpha > \phi$ .

When  $\eta = 0$ , markups and the revenue productivity of labor are constant and indepen-

dent of sales. However, markups and revenue productivity become positively correlated with sales when  $\eta > 0$  because the demand elasticity decreases with size in proportion to  $\eta$ . Therefore, the strength of the relationship between labor productivity and sales is informative about  $\eta$ , holding the other parameters fixed. We summarize this relationship with the regression coefficient of a sales-weighted OLS regression of relative revenue productivity of labor on relative sales of 0.036 for firms with relative sales greater than 1, as computed by Edmond, Midrigan, and Xu (2022) using aggregate data from SUSB.<sup>23</sup>

Finally, the parameters of the AR(1) productivity process for incumbent firms,  $\sigma_z$  and  $\rho_z$ , are set to match a standard deviation of annual employment growth of 0.415 from Elsby and Michaels (2013) and the unweighted distribution of within-industry relative sales from Edmond, Midrigan, and Xu (2022). The mean of the productivity distribution of entrants  $\bar{z}_{ent}$  is set to match the fact that old firms (those older than 11 years) are, on average, six times larger in terms of employment than 1-year old firms (BDS). The exogenous separation probability v is calibrated to match an average exit rate of 7.3% (BDS).

**Calibration Results** The set of calibrated parameters is shown in Panel B of Table 3. The process for the productivity shock is quite persistent and volatile, although in line with estimates from Lee and Mukoyama (2015). The calibrated elasticity and super-elasticity of demand are 6.49 and 4.95, respectively. The value for the elasticity is standard and the value for the super-elasticity is close to estimates found in the literature (see Nakamura and Zerom, 2010). Finally, note that the calibrated value for the elasticity of the matching function  $\phi = 0.533$  is close to a model-generated regression coefficient of 0.474. This similarity lends support to the identification argument provided in Proposition 3. Section SM.4.2 in Supplemental Materials shows that our calibrated model matches the targeted moments reasonably well.

<sup>&</sup>lt;sup>23</sup>In our definition of model revenue productivity of labor, we include the variable component of SGA ( $l_s$ ) but not the fixed component of SGA ( $\chi$ ). The former is due to the fact that the SUSB reports information on the total wage bill across firms in a size group, without distinguishing between types of labor (e.g., production and advertising labor). The decision not to include  $\chi$  is due to the fact that part of overhead costs are, in reality, not associated with labor costs (e.g., rent) and thus not included in the wage bill reported by SUSB. Ideally, we would use data on the subcomponents of SGA expenses in Compustat to compute the share of labor costs within SGA expenses. Unfortunately, a full disaggregation of SGA expenses is not available. To alleviate concerns about this choice, note that we target a moment based on a sample of relatively large firms (those with relative sales greater than 1), for which arguably the fixed overhead cost represents a smaller fraction of total costs.

Parameter	Description	Value	
Panel A: Fixed Parameters			
β	Annual discount factor	0.960	
γ	Elast. of intertemporal substitution	2.000	
$\psi$	Frisch elasticity	1.000	
α	Decreasing returns to scale	0.640	
δ	Prob. of losing customer	0.280	
Panel B: Calibrated Parameters			
$\phi$	Elasticity matching function	0.533	
χ	Overhead cost	0.307	
$\sigma$	Avg. elasticity of substitution	6.490	
η	Superelasticity	4.956	
ν	Exog. survival probability	0.964	
$ ho_z$	Persistence of productivity shock	0.973	
$\sigma_z$	SD of productivity shock	0.218	
$\bar{z}_{ent}$	Mean productivity of entrants	-1.453	
λ	Mass of entrants	0.137	
ξ	Disutility of labor supply	1.981	

Table 3: Model Parameters

*Notes:* This table shows the calibration of the model. Panel A contains parameters externally chosen. Panel B contains parameters internally calibrated to match moments presented in Table SM.4.1 and Figure SM.4.2 in Supplemental Materials.

#### 4.2. Model Validation

We also provide overidentifying tests of the calibrated model regarding its ability to match relevant untargeted moments. We have previously documented that, in the data, the major source of cross-firm differences in sales is the size of their customer bases. This is verified in Table 4, which compares the variance decomposition of log sales in the model with the decomposition from the data. In the data, differences in the log of average sales per customer account for 11.4% of the variance of log sales. The model closely matches this fact, with a fraction of 15.2%. Also, in the model, the largest contributor to the dispersion in log sales is the variance of the log number of customers, as in the data. However, since differences across firms are ultimately driven by only one source of heterogeneity (i.e., the productivity shocks), the model naturally overpredicts the size of

the covariance term.<sup>24</sup>

	Var(ln sales per customer)	Var(ln n. of customers)	Covariance
Data	11.44	80.66	7.90
Model	15.17	47.54	37.29

Table 4: Sources of Sales Dispersion across Firms

*Notes:* This table provides a variance decomposition of firms' log sales. The first column reports the variance of the log sales per customer,  $var(\ln p_{i,t}\Upsilon'(p_{i,t}/D_t))$ , relative to the overall variance of log sales. The second column reports the relative variance of the log number of customers,  $var(\ln m_{i,t})$ . The last column reports the covariance between both terms,  $cov(\ln m_{i,t}, \ln p_{i,t}\Upsilon'(p_{i,t}/D_t))$ . The first row reports the results obtained from the Nielsen Homescan Panel. Sales and the number of customers are adjusted with household sample weights. The second row reports the results obtained from model-simulated data.

Relatedly, we have shown that, despite not being the main driver of sales growth, average sales per customer are strongly associated with market power (see Table 1). Next, we show that our model quantitatively reproduces this untargeted fact by regressing firms' markups on sales per customer and the size of their customer bases using model-simulated data:

$$\ln(\mu_{i,t}) = \theta_0 \ln p_{i,t} q_{i,t} + \theta_1 \ln m_{i,t} + \varepsilon_{i,t}.$$

Table 5 presents the results. The data show a significant relationship between markups and sales per customer and an economically as well as statistically insignificant relationship between markups and the size of the customer base. The model matches these facts fairly well and predicts that 1% higher average sales per customer are associated with 0.11% higher markups. This point estimate is between the baseline estimate of 0.06 reported in Table 1 and the estimate of 0.187 reported in the additional analysis in Appendix A.3. On the other hand, a 1% increase in the size of the customer base increases markups by only 0.02%.<sup>25</sup> Therefore, the model captures the differential roles of the intensive and extensive margins of demand on firms' markups documented in the data.

Finally, the model is able to generate firm dynamics similar to those observed in the data. Figure SM.4.6 in Supplemental Materials shows two model-based moments that were not explicitly targeted in the calibration exercise: a decreasing average exit rate

<sup>&</sup>lt;sup>24</sup>The difference between the model and the data could be explained, for instance, by (unmodeled) orthogonal preference shocks that affect the size of the customer base.

<sup>&</sup>lt;sup>25</sup>In our model, the size of the customer base is a demand shifter and does not directly affect the elasticity of demand. The only reason the regression coefficient on the size of the customer base is not 0 in the model-simulated data is the minor nonlinear nature of the relationship between these variables.

	Data	Model
$\ln p_{\rm it} q_{\rm it}$	0.060***	0.111
	(0.024)	
$\ln m_{\rm it}$	0.003	0.022
	(0.007)	
Observations	2433	
$R^2$	0.338	0.869
Year FE	$\checkmark$	
SIC FE	$\checkmark$	

Table 5: Sources of Dispersion in Sales and Markups

*Notes*: This table reports the results of an OLS regression of a firm's log markup  $(\ln(\mu_{i,t}))$  on log sales per customer  $(\ln p_{i,t}q_{i,t})$  and log size of the customer base  $(\ln m_{i,t})$ . Column (1) reproduces the empirical estimates from Table 1. Column (2) reports estimates based on model-simulated data. The model-simulated panel is restricted to mimic selection into Compustat (see Section 4 for details). In the model, we do not include SIC FE or Year FE because we model a single "representative" industry in steady state.

by age and a decreasing average employment growth by age (as in the data; see, e.g., Haltiwanger, Jarmin, and Miranda, 2013). Intuitively, this is explained by the fact that entrants enter the economy with lower average productivity and no customer base. This makes young firms more likely to exit when faced with negative productivity shocks due to overhead costs and, conditional on staying, to grow more rapidly than older firms due to frontloaded focus on customer acquisition.

#### 4.3. The Role of Endogenous Customer Acquisition

In this section, we investigate how endogenous customer acquisition changes the implications of the relationship between the size distribution of firms and their markups. To do so, we compare our calibrated model to an alternative model without customer heterogeneity that is calibrated to match the same data moments.

This alternative model—labeled "homogeneous customers" hereafter—corresponds to a version of our model in which all firms are exogenously matched with the representative household (i.e.,  $m_{i,t} = 1, \forall i$ ) and it has the same demand structure as in Klenow and Willis (2016), Edmond, Midrigan, and Xu (2022).<sup>26</sup> The results of this calibration, as well as results for its goodness of fit, are reported in Section SM.5 in Supplemental Materials.

<sup>&</sup>lt;sup>26</sup>Since this model does not feature endogenous customer acquisition, it lacks the parameter  $\phi$  that determines the returns to advertising. For this reason, in our recalibration exercise we drop the moment on the relationship between SGA expenses and sales, which was used to pin down  $\phi$ .

Here, we discuss the main implication of this exercise. In Section SM.4.3 in Supplemental Materials, we also provide comparative statics results—i.e., comparisons under the same parameter values—between these two models to illustrate the mechanisms at play.

How does shutting down endogenous customer acquisition distort our interpretation of the data through the lens of the model? As shown in Figure 3, the recalibrated homogeneous customers model generates much lower markup dispersion relative to our baseline model while matching the same level of cost-weighted markup.<sup>27</sup>

The intuition behind this difference between the two models is closely related to how the two models match the distribution of sales across firms under demand curves that allow for heterogeneous markups. In this class of models, markup variation stems from the fact that at the intensive margin, a consumer's demand is less elastic at lower relative prices. Therefore, the maximum variation in markups these models can generate depends on how much variation there is in this elasticity across the demand curve, which is governed by the ratio  $\eta/\sigma$  as illustrated in Equation (3.8). Higher values of  $\eta/\sigma$  allow for higher variation in markups across firms, and both models nest CES demand (i.e., no variation in markups) when  $\eta = 0$ . Note, however, that a higher value of  $\eta/\sigma$  means that the demand elasticity of each customer declines more with relative prices, which implies that their quantity demanded increases by less as prices decline. Visually, a higher  $\eta/\sigma$ bends the right tail of the demand curve downward. In fact, it is well known that this property leads to a "choke quantity" at the intensive margin—i.e., a maximum quantity that is bought as relative prices approach zero (Edmond, Midrigan, and Xu, 2022). It is only when  $\eta \rightarrow 0$  that this quantity goes to infinity, and we obtain the CES demand. By drawing this link, we would like to emphasize that the existence of a choke quantity is closely related to the amount of variation a model implies for markups across firms.

These choke quantities are exactly the root cause of why the homogenous customers model implies a much lower markup dispersion than our baseline model. To see why, note that *all* demand in the Homogenous customers model comes from the intensive margin; i.e., this model locates all the firms in the economy on the same demand curve. So, any variation in the size of firms in this model comes from these firms being at different parts of this single demand curve. Thus, a choke quantity puts an upper bound on how much variation such a model can generate in the size distribution of firms. As a result, the Homogenous customers model faces a fundamental trade-off between generating high

<sup>&</sup>lt;sup>27</sup>Recall that with  $\eta = 0$ , the Kimball aggregator converges to the CES aggregator and markup dispersion is zero.

dispersion in sales and generating high dispersion in markups.

Since the choke quantity decreases with  $\eta/\sigma$ , by matching the same distribution of sales as in the baseline model (Figure SM.5.1 in Supplemental Materials), the homogeneous customers model assigns a much lower calibrated value to this ratio (similar to the one reported by Edmond, Midrigan, and Xu (2022)), which generates much lower dispersion in markups across firms. This is, however, not the case with our baseline model, in which firms' size distribution is also affected by the distribution of customers across firms—which is not directly related to their markups, and thus is not constrained by the choke quantity demanded by each customer.



Figure 3: Distribution of Markups: Baseline vs Homogeneous Customers Model

*Notes:* The figure plots the distribution of production markups in the baseline and *recalibrated* homogeneous customers models. The vertical dashed lines show the average cost-weighted production markup in each model.

Two observations follow. First, the homogeneous customers model overestimates markups for small firms and underestimates markups for large firms relative to our baseline model. Second, and more importantly, by generating a higher markup dispersion than the homogeneous customers model, the baseline model implies much higher welfare losses due to misallocation, as we show in the next section.<sup>28</sup>

<sup>&</sup>lt;sup>28</sup>In the Compustat data, the overall dispersion of markups is larger than the dispersion obtained from

# 5 Efficient Allocation

What are the implications of endogenous customer acquisition for welfare and, in particular, misallocation? In this section we start by developing a welfare decomposition result that highlights the potential sources of welfare losses or gains around the equilibrium allocation. We then characterize the efficient allocation in our economy by solving the problem of a social planner who faces the same constraints as agents do in the equilibrium. Using our decomposition result, we then decompose the equilibrium welfare losses to aggregate TFP losses (misallocation), losses from total underutilization of labor due to *aggregate* markups, and losses from different entry/exit policies for firms. Our core finding is that welfare losses from allocative efficiency in TFP are around 7.8% in consumption-equivalent terms and constitute around 57% of total welfare losses under the equilibrium allocation.

#### 5.1. A Welfare Decomposition Result

Before stating our welfare decomposition result, we need to derive the model-consistent notions of *aggregate TFP* and *aggregate markup*. In particular, for any allocation of production inputs and demand  $(l_{i,p,t}, m_{i,t}, q_{i,t})_{i \in N_t}$  across a set of operating firms  $N_t$ , we derive the aggregate TFP as the productivity implied by an aggregated production function, and the aggregate markup as the wedge between the marginal product of labor and the marginal rate of substitution between leisure and consumption.

**Aggregate Production Function** Let  $L_{p,t}$  denote the total amount of labor allocated to production across all firms. Then, we have:

$$L_{p,t} \equiv \int_{i \in N_t} l_{i,p,t} di = \int_{i \in N_t} \left( \frac{C_t m_{i,t} q_{i,t}}{z_{i,t}} \right)^{\alpha^{-1}} di,$$
(5.1)

where the second equality follows from equating firms' demand to their supply.

Defining aggregate output as the aggregate consumption,  $Y_t \equiv C_t$ , and rearranging Equation (5.1), we arrive at the aggregate production function:

$$Y_t = Z_t L_{p,t}^{\alpha}, \tag{5.2}$$

either our baseline model or the homogeneous customers model. This is expected, since our model focuses only on the dispersion of markups that arises from differences in relative sales (per customer). In reality, markups or wedges in the marginal revenue of products can arise due to other distortions (see, e.g., Hsieh and Klenow, 2009) from which our model abstracts away.

where  $Z_t$ , the aggregate TFP, is given by

$$Z_t \equiv \left[ \int_{i \in N_t} \left( \frac{z_{i,t}}{q_{i,t} m_{i,t}} \right)^{-\alpha^{-1}} di \right]^{-\alpha}.$$
(5.3)

Notice how the allocation of relative demand  $(q_{i,t}, m_{i,t})_{i \in N_t}$  affects aggregate productivity. Higher relative demand means higher production relative to other firms, which decreases aggregate productivity through a lower marginal product of labor (due to decreasing returns to scale) and dispersion in demand per customer (due to imperfect substitutability of goods).

**Aggregate Markup** We define the aggregate markup,  $\mathcal{M}_t$ , as the wedge between the marginal product of aggregated production labor and the real wage  $W_t/P_t$  (which, from the household's labor supply condition, is equal to the marginal rate of substitution between leisure and consumption). Recall also that we have normalized the aggregate price to one ( $P_t = 1$ ) so that the real wage is  $W_t$ . Formally, having derived the aggregate production function in Equation (5.2), aggregate markup is defined as

$$\mathcal{M}_t \equiv \frac{\partial Y_t / \partial L_{p,t}}{W_t} = \alpha \frac{Y_t}{W_t L_{p,t}}.$$
(5.4)

We can also define the firm-level markup as the analog of this wedge for firm *i*:

$$\mu_{i,t} \equiv \alpha \frac{p_{i,t} y_{i,t}}{W_t l_{i,p,t}},\tag{5.5}$$

which corresponds to the equilibrium relationship between the markup and the labor share in Equation (3.16)—with the exception that here we are defining this wedge for an arbitrary allocation of inputs and prices. By combining the last two equations, we obtain the aggregate markup as the production cost-weighted average of firm-level markups (as in Edmond, Midrigan, and Xu, 2022):

$$\mathcal{M}_t = \int_{i \in N_t} \omega_{i,t} \mu_{i,t} di, \qquad (5.6)$$

where the weight  $\omega_{i,t}$  is the *production cost share* of firm *i* or, as referred to by Baqaee and Farhi (2019), the cost-based Domar weight of firm *i*:

$$\omega_{i,t} \equiv \frac{W_t l_{i,p,t}}{\int_{i \in N_t} W_t l_{i,p,t}}.$$
(5.7)

**Decomposition of Welfare** Having defined aggregate TFP and markup, we obtain the following proposition that emphasizes their relevance by showing that, up to the first order, allocations affect welfare only through these and other aggregate objects.

**Proposition 5.** For small perturbations around the equilibrium allocation, changes in household's welfare at a given time *t* are given, up to a first-order approximation, by

$$\underbrace{\frac{\Delta U_t}{U_{c,t}C_t}}_{\text{AWelfare (C.E.)}} \approx \underbrace{\Delta \ln(Z_t)}_{\Delta \text{TFP}} + \underbrace{\alpha(1 - \mathcal{M}_t^{-1})\Delta \ln(L_{p,t})}_{\Delta \text{ from Aggregate Markup}} - \alpha \mathcal{M}_t^{-1} \left[ \underbrace{\chi \frac{N_t}{L_{p,t}}\Delta \ln(N_t)}_{\Delta \text{ from Entry/Exit}} + \underbrace{\frac{L_{s,t}}{L_{p,t}}\Delta \ln(L_{s,t})}_{\Delta \text{ from Advertising}} \right], \quad (5.8)$$

r

1

where  $Z_t$  is the aggregate TFP in Equation (5.3),  $\mathcal{M}_t$  is the equilibrium aggregate (costweighted) markup in Equation (5.6),  $L_{p,t}$  and  $L_{s,t}$  are the aggregate amounts of labor allocated to production and advertising, and  $N_t$  is the equilibrium measure of operating firms.

Equation (5.8) decomposes the consumption-equivalent welfare changes of the household around the equilibrium allocation to four separate terms: (1) allocative and distributional changes that lead to changes in aggregate TFP, (2) losses due to underutilization of labor that arise from aggregate market power—and demonstrates itself as a wedge between the marginal product of labor and the marginal rate of substitution between consumption and leisure, (3) changes in the aggregate labor supply that is allocated to the overhead costs of operating firms, and (4) changes in the aggregate labor supply that is allocated to advertising in the equilibrium.

Proposition 5 lays out the road map for the rest of our analysis. As we move on to quantify the efficient allocation, our main objective is to measure these welfare changes under counterfactual demand allocations.

#### 5.2. Characterization of the Social Planner's Problem

Endogenous customer acquisition creates a new channel for the relationship between relative size and markups. Not only does this new channel affect this relationship in the equilibrium, but it also defines a new Pareto frontier, since the social planner also chooses the distribution of customers across firms. This section characterizes this efficient allocation in our economy.

Given an initial distribution of firms, the social planner of this economy maximizes the household's lifetime utility by choosing (1) which incumbent firms should exit and which potential entrants should enter at each period, (2) how many customers each operating firm should get—which can be achieved either by depreciating their customer base if the firm has too many customers or by launching advertising campaigns if the firm needs to grow—and, finally, (3) how much each operating firm should produce. A formal statement of the planner's problem is included in Appendix B.8. Here, we discuss the key properties of the solution to this problem.

Any Distribution of Customers is Attainable There are two sources of inefficiencies regarding customer acquisition and the allocation of customers in the equilibrium. First, the planner might choose to allocate customers differently across firms than the equilibrium (misallocation of customers). A second source of inefficiency is the business-stealing externality of advertising, which leads to an overuse of labor for advertising in the equilibrium.<sup>29</sup> To focus on the *misallocation* of customers, we shut down this second source of inefficiency by restricting the social planner to spend the same amount of aggregate labor for advertising as in the equilibrium. More precisely, the restriction shuts down Channel 4 ( $\Delta$  from Advertising) in Proposition 5.

A potential concern is that shutting down the second channel in this manner might affect the distribution of customers that are attainable for the social planner and, thus, our conclusions about the misallocation of customers. This is a nontrivial concern, but it turns out that it is nonbinding: In the following lemma, we show that shutting down the second source of inefficiencies is inconsequential for the implications of customer misallocation. That is, by restricting the planner to using a certain amount of aggregate labor for advertising, we are not imposing any restriction on the planner's choices regarding the allocation of customers across firms.

# **Lemma 2.** Any desired distribution of customers across a set of operating firms can be achieved by any strictly positive level of aggregate labor allocated to advertisement.

This result follows from the advertising technology, which requires that returns to advertising are fully relative in labor allocated to posting ads. Given this result, we solve the planner's problem in two steps. First, for any set of operating firms, we characterize the optimal allocation of demand in terms of how many customers each operating firm

<sup>&</sup>lt;sup>29</sup>This assumption is irrelevant for the calibration of the model because the measure of matches can be normalized in that exercise. But it is important for counterfactuals in which it is important how the total number of matches changes with firms' advertising. Assuming that the total number of matches are fixed *across* counterfactuals implies that while advertising changes the distribution of customers across firms, it does not expand the industry's total customer base as in (Drozd and Nosal, 2012, Einav, Klenow, Levin, and Murciano-Goroff, 2022). This assumption is corroborated by marketing and IO literature (Bagwell, 2007). Studies of the publishing (Garthwaite, 2014) and prescription drugs (Sinkinson and Starc, 2019) sectors indicate that advertising boosts own sales without enlarging the total size of the market. Moreover, while we have made the extreme assumption in the model that advertising has no effect on the total size of the customer base in the industry, we view this as a conservative benchmark because the planner also faces this limitation, ensuring that welfare gains come from customer reallocation only and ignoring any potential welfare gains by increasing the size of the customer base for an industry. Finally, by restricting the social planner to using the same amount of labor allocated to advertising, we ensure that our reported welfare gains are not due to the planner's using fewer resources for advertising.

should get and how much they should produce. Second, we characterize the optimal entry and exit rule that determines the sets of operating firms over time.

**Optimal Allocation of Demand Maximizes Aggregate TFP** Here, we characterize the efficient allocation of demand for a given set of operating firms. Formally, by an allocation of demand, we mean the choice of (a) allocating customers across firms in  $N_t$ —i.e.,  $\mathbf{1}_{\{j \in m_{i,t}\}}, \forall j, \forall i \in N_t$ —and (b) choosing the relative demand of every matched customer,  $q_{i,j,t}, \forall j \in m_{i,t}, \forall i \in N_t$ . Naturally, any allocation of demand must:

1. be consistent with the Kimball aggregator:

$$\int_{i \in N_t} \int_0^1 \mathbf{1}_{\{j \in m_{i,t}\}} \Upsilon(q_{i,j,t}) dj di = 1,$$
(5.9)

2. respect the constraint for the total number of matches:

$$\int_{i \in N_t} \int_0^1 \mathbf{1}_{\{j \in m_{i,t}\}} dj di \le 1.$$
 (5.10)

**Proposition 6.** Fix a choice for the set of operating firms,  $N_t$ . Then, the efficient allocation of demand is the one that maximizes aggregate TFP subject to constraints in Equations (5.9) and (5.10). This allocation equalizes demand per *matched* customer across all firms, and matches more customers to more productive firms:

$$q_{i,j,t} = 1, \forall j \in m_{i,t}^*, \quad m_{i,t}^* = \frac{z_{i,t}^{\frac{1}{1-\alpha}}}{\int_{i \in N_t} z_{i,t}^{\frac{1}{1-\alpha}} di}.$$
(5.11)

Proposition 6 shows that the planner equalizes the level of consumption per customer across all firms. Furthermore, to capitalize on the higher efficiency of more productive firms, the planner gives them *more customers*.<sup>30</sup> This is in contrast to the equilibrium, in which more productive firms have both higher sales per customer and more customers than other firms (but potentially fewer customers than what is efficient).

Our result in Proposition 6 is also at odds with the trade-off the social planner faces in the homogenous customers model, in which all firms are assumed to serve the representative consumer. On the one hand, the social planner would like more productive firms to produce more in order to increase *aggregate consumption* (i.e., to equalize the

<sup>&</sup>lt;sup>30</sup>Note that the allocation of customers does not depend on the initial distribution of matches. This follows from Lemma 2. Since the implementation cost of all distributions is the same for the planner, we can assume without loss of generality that the planner exercises the free disposal of matches at the beginning of every period and re-matches all customers based on firms' new productivities. This might lead some firms to lose customers beyond those exogenously depreciated. We note that the same option is also available to firms in the competitive equilibrium, but they choose not to exercise it (see the proof of Proposition 2).

marginal product of inputs across firms). On the other hand, since demand comes from the intensive margin in those models, instructing more productive firms to produce more creates dispersion in *relative consumption* across varieties, which is costly due to the weak substitutability of goods. Therefore, the social planner has to balance these two opposing forces when choosing the optimal allocation of inputs.

However, in our model, the social planner does not face such a trade-off due to the existence of an extensive margin of demand. The optimal allocation equalizes relative consumption across all customers and, instead, equalizes the marginal product of inputs across firms by giving more customers to more productive firms.

More generally, endogenizing the allocation of customers has two important macroeconomic implications. First, it widens the Pareto frontier of the economy because, in our model, the social planner always has the option to replicate the homogeneous allocation of customers across firms. Second, both the magnitude of losses from misallocation and the distance of the equilibrium allocation from this new frontier depend on how effective the equilibrium advertising technology is in replicating the efficient allocation of customers, rather than the efficient allocation of sales per customer, as is the case in conventional models.

**Social Value of a Firm** In characterizing the efficient allocation, we show that to make entry or exit decisions, the social planner considers the lifetime valuation of a firm for the welfare of the representative household, which we denote by the social value of a firm. In order to compare how the social planner values individual firms relative to the equilibrium, the following proposition characterizes the social value of firms under the efficient allocation.

**Proposition 7.** The social value of any firm at any given time depends only on its current productivity and is given by

$$\nu_{t}^{*}(z_{i,t}) = \max_{\substack{\{\mathbf{1}_{i,\tau}, m_{i,\tau} \\ q_{i,\tau}, l_{i,p,\tau}\}_{\tau \ge t}}} \mathbb{E}_{t} \sum_{\tau=t}^{\infty} (\beta \nu)^{\tau-t} \left(\prod_{h=t}^{\tau} \mathbf{1}_{i,h}\right) \left(\frac{C_{\tau}^{*}}{C_{t}^{*}}\right)^{-\gamma} \left(D_{\tau}^{*} \frac{\Upsilon(q_{i,\tau})}{q_{i,\tau}} y_{i,\tau} - W_{t}^{*}(l_{i,p,\tau} + \chi) - C_{\tau}^{*}(D_{\tau}^{*} - 1)m_{i,\tau}\right)$$
(5.12)

s.t. 
$$y_{i,\tau} = m_{i,\tau} q_{i,\tau} C_{\tau}^* = z_{i,\tau} l_{i,p,\tau}^{\alpha}$$
 (5.13)

where  $(C_{\tau}^*, D_{\tau}^*, W_{\tau}^*)_{\tau \ge t}$  are the aggregate consumption, aggregate demand index, and decentralized wage under the planner's allocation.

The first observation is that the social value of a firm does not depend on its initial customer base at a given time *t*. This follows from Lemma 2: Since any distribution
of matches is attainable for the social planner using the matching technology of the economy, the efficient allocation is not restricted by the initial distribution of matches at any given time. Accordingly, the social planner can choose the optimal number of customers for any firm given their current productivity. If a firm has too many customers, the planner simply exercises the free disposal condition of matches (i.e., choose  $\delta_{i,t} \ge \delta$ ), and if a firm has too few customers, the planner allocates more ads for that firm in relative terms to attain the optimal number of matches.

Moreover, the expression in Equation (5.12) clarifies how the planner's allocation for a given firm differs from its equilibrium counterpart in Equation (3.13). There are three notable differences. First, the social planner maximizes a firm's value based on the *average consumer surplus*  $(\Upsilon(q)/q)$ . In contrast, in the equilibrium, the revenue firms receive is based on the *marginal valuation*  $(\Upsilon'(q))$ . Second, the advertising labor cost of allocating customers does not appear in this value because aggregate advertising labor is a sunk cost for the planner and, at the margin of keeping a firm, it does not affect her decision. Third, there is a new term  $(C_{\tau}^*(D_{\tau}^* - 1)m_{i,\tau})$  that captures the externalities of allocating customers across firms, which is endogenous to the planner but neglected by individual firms. This externality arises from the fact that, fixing the total number of customers in the economy, if one firm is being matched to a particular customer, then another firm must lose one.<sup>31</sup>

## 5.3. Quantifying the Efficient Allocation

This section presents the differences between the equilibrium and efficient allocations and, in particular, quantifies the welfare gains under the efficient allocation of demand. In doing so, we also revisit our exante decomposition of welfare in Proposition 5 and quantify the contribution of each of the three channels (TFP, aggregate markups, and overhead costs). To isolate the role of endogenous customer acquisition, we conclude this section by repeating the same exercise for counterfactual values of  $\phi$ , the parameter that governs returns to scale in customer acquisition for firms.

**Welfare Gains** We start by quantifying the three channels of welfare gains from Proposition 5 in our calibrated model per Proposition 5 (recall that we shut down welfare gains from advertising labor—Channel 4—by restricting the social planner to using the same amount of labor).

<sup>&</sup>lt;sup>31</sup>Note that this effect would still exist even if advertising increases the total number of customers because, conditional on having spent on advertising labor and creating a match, the planner still has to decide which firm to allocate the match to.

$$\underbrace{\Delta U_{t}}{U_{c,t}C_{t}} \approx \underbrace{\Delta \ln(Z_{t})}_{\Delta \text{Welfare (C.E.)} = 13.6\%} \text{ TFP gains = 10.8\%} \underbrace{-\alpha \mathcal{M}_{t}^{-1} \chi \frac{N_{t}}{L_{p,t}} \Delta \ln(N_{t})}_{\text{Gains from Entry/Exit = 1.6\%}} \underbrace{+\alpha(1 - \mathcal{M}_{t}^{-1}) \Delta \ln(L_{p,t})}_{\text{Losses from Underutilization of Labor = 0.78\%}} (5.14)$$

There are two main takeaways from this decomposition: (1) the consumption-equivalent welfare gains of the household under the efficient allocation are substantial and quantified at 13.6%, and (2) the majority of this gain is coming from the efficiency gains in aggregate TFP under the planner's allocation, quantified at 10.8% higher than the equilibrium TFP. In addition to this substantial gain in TFP, the planner also achieves 1.6% higher welfare by reducing the amount of labor allocated to the overhead cost of operating firms and 0.78% higher welfare by correcting for the underutilization of labor due to aggregate market power.

Moreover, the "Baseline" column in Table 6 presents the implied changes in other quantities that arise from these gains. As a result of higher TFP and higher production labor, the output is 14.6% higher under the efficient allocation, even though the number of firms is 11.3% lower. This higher production with fewer firms is made possible by the fact that the concentration of sales among the top 5% largest firms is 39.2% larger than in the equilibrium. In addition to the baseline model, Table 6 presents similar results for two counterfactual values of  $\phi$ . In the remainder of this section, we dissect these changes and study the underlying forces that shape these gains.

**Implications for Misallocation** By explicitly modeling the extensive margin of demand in a way consistent with the data, we find large TFP gains of 10.8%. This shows that efficiency losses from the misallocation of customers are large and go well beyond the social costs of markups (for instance, Edmond, Midrigan, and Xu (2022) estimate the efficiency losses from markups to be around 0.8% to 1.8%).

To further analyze the increase in aggregate productivity, we consider the decomposition of TFP derived by Baqaee and Farhi (2019) and separate *allocative efficiency gains* from *technological change*. For us, allocative efficiency refers to how differently the planner allocates resources across firms, while technological change is a manifestation of the different entry and exit policies the planner adopts. Formally, let  $Z(N_t, \mathcal{A}_t)$  denote the aggregate productivity implied by the set of operating firms  $N_t$  with an allocation rule  $\mathcal{A}_t \equiv (l_{i,p,t})_{i \in N_t}$  among them. Then, we can decompose the difference in TFPs across two

	Endogenous $m_{i,t}$					
	$\phi = 0.25$	Baseline	$\phi = 0.75$			
TFP	24.1	10.8	3.2			
Output	27.5	14.6	7.7			
Number of firms	-41.9	-11.3	-2.6			
Employment	-5.0	2.1	4.4			
Production	5.3	6.0	7.0			
Welfare	37.9	13.6	4.0			
Agg. markup	-27.8	-22.8	-19.1			
Top 5% sales share	88.8	39.2	15.5			

Table 6: Comparison with Efficient Allocation

*Notes:* The table compares aggregate variables between the social planner's allocation and the equilibrium allocation. Differences are reported as percent deviations from equilibrium allocations. Three comparisons are presented by varying the value of  $\phi$  while keeping the remaining parameters fixed at the values in the baseline calibration.

allocations as

$$\underbrace{\ln\left(\frac{Z(N_t^*, \mathscr{A}_t^*)}{Z(N_t, \mathscr{A}_t)}\right)}_{\Delta \text{ TFP} = 10.8\%} = \underbrace{\ln\left(\frac{Z(N_t, \mathscr{A}_t^*)}{Z(N_t, \mathscr{A}_t)}\right)}_{\Delta \text{ Allocative Efficiency = 7.8\%}} + \underbrace{\ln\left(\frac{Z(N_t^*, \mathscr{A}_t^*)}{Z(N_t, \mathscr{A}_t^*)}\right)}_{\Delta \text{ Entry/Exit Efficiency = 3.0\%}}.$$
(5.15)

The first term on the right-hand side of Equation (5.15) shows that, keeping the set of operating firms fixed, almost 75% of the efficiency gains under the planner's allocations are due to allocative efficiency gains. This is the most important consequence of endogenous customer acquisition: Having the ability to reallocate customers across firms, the planner shifts the distribution of customers toward the top of the productivity distribution, and hence is able to allocate higher amounts of production labor to them.

The extensive margin of demand is the key to the high TFP gains in our model: Without the extensive margin, the planner can only achieve a higher aggregate productivity by shifting demand toward more productive firms on the intensive margin. However, since varieties are weak substitutes, distorting the distribution of relative demand is costly. These costs are even higher when demand is more elastic at higher quantities (as with Kimball preferences or any semi-kinked demand system).

However, in this model, the efficient allocation is not restricted by the weak substitutability margin: Completely equalizing relative consumption across individuals  $(q_{i,t}^* = 1)$ , the efficient allocation achieves much higher aggregate productivity by simply allocating *more customers* to more productive firms  $(m_{i,t}^* \propto z_{i,t}^{\frac{1}{1-\alpha}})$ —as seen in Figure SM.4.12 in Supplemental Materials, which shows a comparison of the allocation of customers between the equilibrium and the efficient allocation. As a consequence of this reallocation of customers, the concentration of sales among the top 5% of firms increases by 39.2%. This higher concentration of customers among more productive firms is efficient to the point that marginal costs of production are equalized across all firms.

It is important to note that the efficient allocation of customers across firms is not restricted by the decreasing returns to scale in advertising, even though the planner is subject to the *same* advertising technology as in the equilibrium. This follows from the fact that the planner internalizes the business-stealing externalities of advertising, and by Lemma 2 can implement any desired distribution of customers given the same equilibrium technology.

Finally, while the optimal allocation of resources accounts for around 75% of the change in aggregate TFP, the remaining 25% is explained by sheer compositional changes in the distribution of productivity, i.e., technological change. Under the efficient allocation, the planner is more selective in allowing firms to enter and chooses a higher productivity cutoff for the entry and exit of firms. A more selective policy increases TFP because it increases the average productivity of firms that operate in the economy and implies fewer operating firms.

The Optimal Number of Firms We start by reviewing the usual costs and benefits of having more firms in the economy and then discuss the new mechanism that comes into play in our model. In conventional models, the optimal number of firms is affected by the interaction of three forces: decreasing returns to scale, love of variety, and aggregate overhead costs. On one side, with decreasing returns to scale at the firm level, having more firms increases the aggregate efficiency by dividing resources across a larger number of firms. Moreover, with love-of-variety, even fixing the average output produced by a larger set of firms, the household enjoys the resultant *aggregated* output more, and hence the economy experiences higher productivity.<sup>32</sup> While these forces form the benefits of a

<sup>&</sup>lt;sup>32</sup>Both of these forces can be summarized by the following simple example inspired by Edmond, Midrigan, and Xu (2022). Consider an economy with *N* firms indexed by *i*, where every firm produces with  $y_i = l_i^{\alpha}$ and aggregate output is given by a CES aggregator,  $Y = [\int_0^N y_i^{\theta^{-1}} di]^{\theta}$ . For a given amount of aggregate labor, *L*, every firm gets to produce  $y_i = (L/N)^{\alpha}$  and the aggregate output is given by  $Y = N^{\theta - \alpha} L^{\alpha}$ . Now if we shut down love of variety ( $\theta = 1$ ), productivity is  $N^{1-\alpha}$ , which increases with *N*. If we shut down decreasing returns to scale ( $\alpha = 1$ ), productivity is  $N^{\theta-1}$ , which indicates higher productivity due to love of variety with

higher number of firms, the cost is usually modeled either as a fixed entry cost for every firm or, as we model here, a stream of overhead costs over time, both of which lead to an optimal finite measure of firms in the equilibrium.

Moreover, our model has an additional force that arises from the allocation of customers across firms. Since the number of customers is fixed, a higher concentration of customers at the top comes at the cost of fewer customers at the bottom of the productivity distribution, which in turn reduces the social value of such firms (as shown in Proposition 7). Hence, with this additional instrument, our planner achieves a higher TFP without having to pay for the overhead costs of more firms, which leads to a lower number of firms in the efficient allocation and increases the welfare of the household by 1.6%, as shown in Equation (5.3).

**Aggregate Labor Supply** Two forces work in opposite directions in affecting the differences in aggregate labor supply between the efficient and equilibrium allocations. On the one hand, the more selective policy of the planner for entry and exit reduces the amount of labor required to finance the overhead costs of operating firms. On the other hand, production labor is underutilized in the equilibrium due to the aggregate market power of firms. The "Baseline" column of Table 6 shows that while labor allocated to production goes up by 6% under the efficient allocation—which together with the higher aggregate TFP contributes to the 14.6% increase in output—the aggregate labor goes up by only 2.1%, since it is mitigated by the lower use of labor in financing the entry cost of firms.

The Role of Returns to Scale in Marketing While for the planner the only relevant margin in allocating customers is the returns to scale in production, the efficiency gains from the reallocation of customers depend on the returns to scale for customer acquisition,  $\phi$ . Larger returns to scale in customer acquisition would imply that more productive firms would invest more in customer acquisition, which is desirable from the perspective of the efficient allocation. Figure 4 shows the scatter plot of firms' productivity and output for both the equilibrium and social planner's allocation and for three different values of  $\phi$ (a low value, the calibrated value, and a high value). The figure shows that with a larger  $\phi$ the equilibrium allocation of customers is closer to that of the planner.<sup>33</sup>

Finally, the  $\phi = 0.25$  and  $\phi = 0.75$  columns of Table 6 show how the allocation of customers is solely responsible for the large efficiency gains under the planner's allocation.

larger N.

<sup>&</sup>lt;sup>33</sup>In fact, for the special case of  $\delta = 1$ , the allocation of customers in the equilibrium coincides with the efficient allocation when  $\phi \rightarrow 1$ .



Figure 4: Allocation of Output: Equilibrium vs. Efficient Allocation

*Notes:* This figure shows a scatter plot between relative productivity  $z_{i,t}/\bar{z}$  and relative output  $y_{i,t}/\bar{y}$  for both the equilibrium and the social planner's allocation. We present three plots by varying the value of  $\phi$ , and keeping the remaining parameters fixed at the values in the baseline calibration. Low  $\phi$  corresponds to 0.25, baseline to 0.53, and high to 0.75.

By simply allowing  $\phi$  to be larger, the equilibrium welfare losses drop from 38% in the case of  $\phi = 0.25$  to only 4% with  $\phi = 0.75$ . When  $\phi$  is larger, in the equilibrium, more productive firms grow mainly through acquiring more customers (higher *m*) rather than selling more per customer (higher *q*). As a result, they produce for more customers but sell less per customer, and as a result charge lower markups. Hence, aggregate TFP, output, and concentration increase while the aggregate markup decreases and the economy gets closer to the efficient allocation.

## 6 Conclusion

In this paper, we revisit the role of the extensive and intensive margins of demand in firms' market share and market power. Using a dataset that merges information from the consumer and the producer sides, we document that while firms' sales grow mainly through acquiring more customers, their market power is only correlated with their average sales per customer.

Guided by these empirical findings, we develop and quantify a model that microfounds the relationship between market power and concentration in the extensive and intensive margins. In our model, while firms hold market power over each customer, the total number of customers acts as a demand shifter. The model provides a new perspective on the relationship between relative firm size and market power. Firms that are big due to a larger customer base have lower market power relative to equally big firms with higher sales per customer. We also find substantive welfare gains under the efficient allocation that stems from the new Pareto frontier of the economy under endogenous customer acquisition.

Our analysis sheds light on the effectiveness of policies that target market power through concentration and profits. In particular, our model highlights a new unintended consequence of policies that target only firms' market power. In our model, although market power is distortionary, it compensates more productive firms for their investment in customer acquisition and improves the allocation of customers. Thus, policies that target larger firms disproportionately may have adverse effects through the misallocation of customers. In particular, if more productive firms are taxed for their larger sales due to larger customer bases, on the margin, they will sell to fewer customers at lower prices but higher markups—both of which are inefficient from a social perspective

h

## References

- AMITI, M., O. ITSKHOKI, AND J. KONINGS (2019): "International Shocks, Variable Markups, and Domestic Prices," *Review of Economic Studies*, 86(6), 2356–2402.
- ANDREWS, I., M. GENTZKOW, AND J. M. SHAPIRO (2017): "Measuring the Sensitivity of Parameter Estimates to Estimation Moments," *Quarterly Journal of Economics*, 132(4), 1553–1592.
- ARGENTE, D., D. FITZGERALD, S. MOREIRA, AND A. PRIOLO (2021): "How Do Firms Build Market Share?," Manuscript.
- ARGENTE, D., M. LEE, AND S. MOREIRA (2018): "Innovation and Product Reallocation in the Great Recession," *Journal of Monetary Economics*, 93, 1–20.
- ——— (2019): "The Life Cycle of Products: Evidence and Implications," Manuscript.
- ARKOLAKIS, C. (2010): "Market Penetration Costs and the New Consumers Margin in International Trade," *Journal of Political Economy*, 118(6), 1151–1199.
- ARNOUD, A., F. GUVENEN, AND T. KLEINEBERG (2019): "Benchmarking Global Optimizers," Manuscript.
- ATKESON, A., AND A. BURSTEIN (2008): "Pricing-to-Market, Trade Costs, and International Relative Prices," *American Economic Review*, 98(5), 1998–2031.
- BAGWELL, K. (2007): "The economic analysis of advertising," *Handbook of Industrial Organization*, 3, 1701–1844.

- BAQAEE, D. R., AND E. FARHI (2019): "Productivity and Misallocation in General Equilibrium\*," *Quarterly Journal of Economics*, 135(1), 105–163.
- BASU, S. (2005): "Comment On: "Implications of State-Dependent Pricing for Dynamic MacRoeconomic Modeling"," *Journal of Monetary Economics*, 52(1), 243–247.
- BOLTON, R. N., P. K. KANNAN, AND M. D. BRAMLETT (2000): "Implications of Loyalty Program Membership and Service Experiences for Customer Retention and value," *Journal of the Academy of Marketing Science*, 28(1), 95–108.
- BOND, S., A. HASHEMI, G. KAPLAN, AND P. ZOCH (2021): "Some unpleasant markup arithmetic: Production function elasticities and their estimation from production data," *Journal of Monetary Economics*.
- BORNSTEIN, G. (2021): "Entry and Profits in an Aging Economy: The Role of Consumer Inertia," Mimeo.
- BORNSTEIN, G., AND A. PETER (2022): "Nonlinear Pricing and Misallocation," Mimeo.
- BURSTEIN, A., V. M. CARVALHO, AND B. GRASSI (2020): "Bottom-up Markup Fluctuations," Manuscript.
- CAVENAILE, L., M. A. CELIK, J. PERLA, AND P. ROLDAN-BLANCO (2023): "A model of product awareness and industry life cycles," Manuscript.
- CLEMENTI, G. L., AND B. PALAZZO (2016): "Entry, Exit, Firm Dynamics, and Aggregate Fluctuations," *American Economic Journal: Macroeconomics*, 8(3), 1–41.
- DAVIS, S. J., AND J. HALTIWANGER (1992): "Gross Job Creation, Gross Job Destruction, and Employment Reallocation," *Quarterly Journal of Economics*, 107(3), 819–863.
- DAVIS, S. J., J. C. HALTIWANGER, S. SCHUH, ET AL. (1998): "Job Creation and Destruction," *MIT Press Books*, 1.
- DE LOECKER, J., J. EECKHOUT, AND G. UNGER (2020): "The Rise of Market Power and the MacRoeconomic Implications," *Quarterly Journal of Economics*, 135(2), 561–644.
- DOTSEY, M., AND R. G. KING (2005): "Implications of State-Dependent Pricing for Dynamic MacRoeconomic Models," *Journal of Monetary Economics*, 52(1), 213–242.
- DROZD, L. A., AND J. B. NOSAL (2012): "Understanding International Prices: Customers as Capital," *American Economic Review*, 102(1), 364–395.
- EDMOND, C., V. MIDRIGAN, AND D. Y. XU (2022): "How Costly Are Markups?," Manuscript.
- EINAV, L., P. J. KLENOW, J. D. LEVIN, AND R. MURCIANO-GOROFF (2022): "Customers and Retail Growth," Manuscript.
- ELSBY, M. W., AND R. MICHAELS (2013): "Marginal Jobs, Heterogeneous Firms, and Unemployment Flows," *American Economic Journal: Macroeconomics*, 5(1), 1–48.
- FITZGERALD, D., S. HALLER, AND Y. YEDID-LEVI (2016): "How Exporters Grow," *National Bureau of Economic Research Working Paper Series.*
- FOSTER, L., J. HALTIWANGER, AND C. SYVERSON (2015): "The Slow Growth of New Plants: Learning About Demand?," *Economica*, 83(329), 91–129.
- GARTHWAITE, C. L. (2014): "Demand spillovers, combative advertising, and celebrity endorsements," *American Economic Journal: Applied Economics*, 6(2), 76–104.

- GOPINATH, G., P.-O. GOURINCHAS, C.-T. HSIEH, AND N. LI (2011): "International Prices, Costs and Mark-up differences," *American Economic Review*, 101(6), 2450–86.
- GOPINATH, G., AND O. ITSKHOKI (2010): "Frequency of Price Adjustment and Pass-Through," *Quarterly Journal of Economics*, 125(2), 675–727.
- HALTIWANGER, J., R. S. JARMIN, AND J. MIRANDA (2013): "Who Creates Jobs? Small Versus Large Versus Young," *Review of Economics and Statistics*, 95(2), 347–361.
- HOPENHAYN, H. A. (1992): "Entry, Exit, and Firm Dynamics in Long Run Equilibrium," *Econometrica*, pp. 1127–1150.
- HOTTMAN, C. J., S. J. REDDING, AND D. E. WEINSTEIN (2016): "Quantifying the Sources of Firm Heterogeneity," *Quarterly Journal of Economics*, 131(3), 1291–1364.
- HSIEH, C.-T., AND P. J. KLENOW (2009): "Misallocation and Manufacturing TFP in China and India," *Quarterly Journal of Economics*, 124(4), 1403–1448.
- KAPLAN, G., AND P. ZOCH (2020): "Markups, Labor Market Inequality and the Nature of Work," *National Bureau of Economic Research Working Paper Series*.
- KIMBALL, M. (1995): "The Quantitative Analytics of the Basic Neomonetarist Model," *Journal of Money, Credit and Banking*, 27(4), 1241–77.
- KLENOW, P. J., AND J. L. WILLIS (2016): "Real Rigidities and Nominal Price Changes," *Economica*, 83(331), 443–472.
- LEE, Y., AND T. MUKOYAMA (2015): "Productivity and Employment Dynamics of US Manufacturing Plants," *Economics Letters*, 136, 190–193.
- MITTAL, V., AND W. A. KAMAKURA (2001): "Satisfaction, Repurchase Intent, and Repurchase Behavior: Investigating the Moderating Effect of customer characteristics," *Journal of Marketing Research*, 38(1), 131–142.
- NAKAMURA, E., AND D. ZEROM (2010): "Accounting for Incomplete Pass-Through," *Review* of *Economic Studies*, 77(3), 1192–1230.
- NEIMAN, B., AND J. S. VAVRA (2019): "The Rise of Niche Consumption," Manuscript.
- OTTONELLO, P., AND T. WINBERRY (2020): "Financial heterogeneity and the investment channel of monetary policy," *Econometrica*, 88(6), 2473–2502.
- PETERS, M. (2020): "Heterogeneous markups, growth, and endogenous misallocation," *Econometrica*, 88(5), 2037–2073.
- PHELPS, E. S., AND S. G. WINTER (1970): "Optimal Price Policy Under Atomistic Competition," *Microeconomic foundations of employment and inflation theory*, pp. 309–337.
- RESTUCCIA, D., AND R. ROGERSON (2008): "Policy Distortions and Aggregate Productivity With Heterogeneous Establishments," *Review of Economic Dynamics*, 11(4), 707–720.
- ROTEMBERG, J. J., AND M. WOODFORD (1999): "The Cyclical Behavior of Prices and Costs," *Handbook of Macroeconomics*, 1, 1051–1135.
- SINKINSON, M., AND A. STARC (2019): "Ask your doctor? Direct-to-consumer advertising of pharmaceuticals," *Review of Economic Studies*, 86(2), 836–881.
- STROEBEL, J., AND J. VAVRA (2019): "House Prices, Local Demand, and Retail Prices," *Journal of Political Economy*, 127(3), 1391–1436.

- TRAINA, J. (2019): "Is Aggregate Market Power Increasing? Production Trends Using Financial Statements," Manuscript.
- WASI, N., AND A. FLAAEN (2015): "Record Linkage Using Stata: Preprocessing, Linking, and Reviewing Utilities," *The Stata Journal*, 15(3), 672–697.
- YOUNG, E. R. (2010): "Solving the Incomplete Markets Model With Aggregate Uncertainty Using the Krusell–Smith algorithm and non-stochastic simulations," *Journal of Economic Dynamics and Control*, 34(1), 36–41.

# **Online Appendix**

Concentration, Market Power, and Misallocation: The Role of Endogenous Customer Acquisition

by Hassan Afrouzi, Andres Drenik, and Ryan Kim

## Table of Contents for the Online Appendix

A	Add	itional Empirical Results	<b>48</b>
	A.1	Measurement error in $\ln m_{it}$	48
	A.2	Markups, Sales per Customers, and Sales	49
	A.3	Evidence Based on Product-level Markups	49
	A.4	Firm Sales Growth Decomposition	52
	A.5	Customer Acquisition and Firms' Non-production Costs	54
B	Der	ivations	60
	<b>B.</b> 1	Lemma 1	60
	B.2	Proposition 1	60
	B.3	Corollary 1	61
	B.4	Proposition 2	61
	B.5	Proposition 3	62
	B.6	Proposition 4	62
	<b>B.7</b>	Proposition 5	62
	B.8	Planner's Problem	63
	B.9	Lemma 2	64
	B.10	Proposition 6	64
	B.11	Proposition 7	68

## A Additional Empirical Results

This section provides additional results to confirm the robustness of our main empirical findings.

## A.1. Measurement error in $\ln m_{it}$

	(1)	(2)	(3)	(4)	(5)
ln <i>m</i> <sub>it</sub>	-0.002	-0.002	0.001	0.001	0.002
	(0.007)	(0.007)	(0.007)	(0.007)	(0.008)
$\ln p_{\rm it} q_{\rm it}$	$0.107^{***}$	$0.107^{***}$	0.072***	$0.071^{***}$	0.073***
	(0.037)	(0.037)	(0.023)	(0.023)	(0.026)
Observations	2012	2012	2012	2012	2012
$R^2$	0.062	0.061	0.033	0.031	0.032
Year FE		$\checkmark$		$\checkmark$	
SIC FE			$\checkmark$	$\checkmark$	
SIC-year FE					$\checkmark$
First-stage F statistics	1.6e+04	1.6e+04	1.1e+04	1.1e+04	1.1e+04

Table A.1: Markups, Sales per Customer, and Number of Customers: IV Approach

*Notes:* This table replicates Table 1 by instrumenting  $\ln m_{it}$  with its lagged value  $(\ln m_{it-1})$ . \* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01; standard errors are clustered at the firm-level. Markups are measured as the Sales-to-COGS ratio. The variable  $\ln p_{it}q_{it}$  denotes the log of the average sales per customer and  $\ln m_{it}$  denotes the log number of customers. SIC industries correspond to a two-digit SIC code. All Nielsen variables are projection-factor adjusted.

	(1)	(2)	(3)	(4)	(5)
$\ln m_{\rm i,t-1}$	0.970***	0.971***	0.958***	0.958***	0.963***
	(0.008)	(0.008)	(0.009)	(0.009)	(0.009)
$\ln p_{ m it} q_{ m it}$	-0.010	-0.010	0.005	0.008	0.016
	(0.026)	(0.027)	(0.032)	(0.032)	(0.034)
Observations	2012	2012	2012	2012	2012
$R^2$	0.948	0.948	0.950	0.950	0.957
Year FE		$\checkmark$		$\checkmark$	
SIC FE			$\checkmark$	$\checkmark$	
SIC-year FE					$\checkmark$

Table A.2: Markups, Sales per Customer, and Number of Customers: First Stage

*Notes:* This table presents the first-stage results behind the IV estimation in Table A.1.

#### A.2. Markups, Sales per Customers, and Sales

Table A.3 replicates Table 1 by replacing the number of customers with a firm's total sales. We replace the number of customers with sales, so that our regressors have the same unit. Regardless of whether we control for total sales or the number of customers, our results show a strong correlation between markups and sales per customer, which suggests the importance of sales per customer in understanding firm-level markups.

	(1)	(2)	(3)	(4)	(5)
$\ln p_{\rm it} q_{\rm it}$	0.094***	0.093***	0.058**	0.056**	0.057**
	(0.031)	(0.032)	(0.023)	(0.023)	(0.025)
ln C	0.000	0.000	0.000	0.000	0.002
$III S_{it}$	-0.002	-0.002	0.002	0.002	0.005
	(0.006)	(0.006)	(0.007)	(0.007)	(0.007)
Observations	2433	2433	2433	2433	2433
$R^2$	0.046	0.047	0.311	0.313	0.338
Year FE		$\checkmark$		$\checkmark$	
SIC FE			$\checkmark$	$\checkmark$	
SIC-year FE					$\checkmark$

Table A.3: Markups, Sales, and Sales per Customer

*Notes:* \* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01; standard errors are clustered at the firm-level. Markups are measured as Sales-to-COGS ratio. SIC is a two-digit SIC code, and all Nielsen variables are projection-factor adjusted.

## A.3. Evidence Based on Product-level Markups

This section revisits our analysis of the relationship between markups and the relative size of different margins of a firm's demand by combining two alternative approaches to measure pricecost markup and account for marginal costs. Our baseline analysis focuses on firm-level markups and public firms. Here, we show that we find similar results when expanding the scope to product-level markups and analyzing a much broader set of firms in the Nielsen data.

As a first alternative approach, we use the difference between the retail-level price and the wholesale cost at the product level to measure markups. For this, we use microdata on the UPC-market-year-level wholesale cost from the Nielsen PromoData available through the Chicago Booth Kilts Center. The PromoData record wholesale costs by UPC, market, and year and are collected from 12 national wholesalers in the period 2006-2012 (we restrict our sample to the years 2006-2011 since there is a substantial number of missing observations in 2012). In total, we use data from 45 markets (examples of markets are Chicago, Los Angeles, and Atlanta). Given that the data lack detailed sales information, we take a simple geometric average of wholesale costs by UPC, market, and year after adjusting for the package size. We then (1) combine this information with the retail-level price, sales, and sales per household from the Nielsen Homescan Panel data and (2) measure the retail-level markup at the product (UPC) level as the difference between the retail price and the

wholesale cost (we drop negative values of the measured markups). A similar approach has been followed by Gopinath, Gourinchas, Hsieh, and Li (2011) and Stroebel and Vavra (2019).

This measure of markups assumes that other costs, such as wages and capital expenditures, do not confound the observed relationship among markups, average sales per customer, and the number of customers. However, there could exist sources of variation in marginal costs across retailers as well, in addition to the variation arising from differences in product-level wholesale costs. To alleviate these concerns, we exploit the rich variation in the data and control for such differences in marginal costs by incorporating various sets of fixed effects in the regression, which is the approach followed by Fitzgerald, Haller, and Yedid-Levi (2016). Our markup measure could vary at the UPC, year, market, and retailer level, so we progressively include different combinations of fixed effects at those levels to absorb differences in marginal costs that could potentially confound the relationship of interest. For example, the fact that retailers sell the same UPC in multiple locations allows us to control for differences in marginal costs at the retailer-year-UPC level. The underlying assumption is that the retailer's marginal cost of selling a given UPC is the same across markets. Similarly, we can account for any time-invariant costs by including a set of UPC-market-retailer fixed effects and any differential distribution costs by incorporating a set of UPC-market-year fixed effects. The final sample contains data from 7,802 UPCs, 45 markets, 6 years, and 167 retailers.

Table A.4 confirms our previous results regarding the relationship among markups, average sales per customer, and the number of customers reported in Table 1. Column (1) shows that markups are strongly positively correlated with average sales per customer. The relationship between markups and the number of customers is negative but an order of magnitude smaller in size. Column (2) includes UPC-, market-, and retailer-fixed effects interacted with year-fixed effects. We find similar results, although the coefficients become smaller (in absolute value). Columns (3), (4), (5), and (6) allow for additional sets of fixed effects that absorb any differences in marginal costs at those levels (including the variation in wholesale costs, which is at the UPC-market-year-level). Once we exploit the geographic variation in the data and include retailer-year-UPC fixed effects in column (3), the coefficient on the number of customers declines to -0.016, while the coefficient on the average sales per customer remains stable at 0.166. Controlling for time-invariant costs in column (4) shows similar changes in the coefficients: There is still a positive and significant relationship between markups and average sales per customer, but no economically significant relationship with the number of customers (the coefficient is estimated to be close to zero). Adding the additional sets of fixed effects in columns (5) and (6) shows similar relationships among variables. Table A.5 replaces the number of customers in Table A.4 with total sales so that our two main regressors are expressed in the same unit. These results confirm that the relevant margin for the relationship between relative size and markups is the average sales per customer.

	(1)	(2)	(3)	(4)	(5)	(6)
$\ln p_{\rm urmt} q_{\rm urmt}$	0.550***	0.186***	0.166***	0.185***	0.185***	0.187***
	(0.035)	(0.027)	(0.023)	(0.023)	(0.021)	(0.019)
ln <i>m</i> <sub>urmt</sub>	-0.063***	-0.040***	-0.016***	-0.004*	-0.003	-0.004
	(0.018)	(0.005)	(0.005)	(0.002)	(0.002)	(0.002)
Observations	426032	426032	426032	426032	426032	426032
$R^2$	0.126	0.550	0.742	0.870	0.875	0.928
UPC-year FE		$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	
market-year FE		$\checkmark$	$\checkmark$	$\checkmark$		
retail-year FE		$\checkmark$				
retail-year-UPC FE			$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
UPC-market-retail FE				$\checkmark$	$\checkmark$	$\checkmark$
year-market-retail FE					$\checkmark$	$\checkmark$
UPC-market-year FE						$\checkmark$

Table A.4: UPC-market-retailer-year-level analysis

*Notes*: \* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01; standard errors are three-way clustered by product group, retail, and market. Markup<sub>*urmt*</sub> is measured as the difference between retail-level price and wholesale cost,  $p_{urmt}q_{urmt}$  denotes the average sales per customer, and  $m_{urmt}$  the number of customers, where the subscript *u* refers to a particular UPC, *r* the retailer, *m* the market, and *t* the year. All variables are projection-factor adjusted. We balance the sample across columns based on the tightest specification in column (6); our final sample includes 7,802 UPCs, 45 markets, 6 years, and 167 retailers.

			-	-		
	(1)	(2)	(3)	(4)	(5)	(6)
$\ln p_{\rm urmt} q_{\rm urmt}$	0.612***	0.226***	0.182***	0.189***	0.188***	0.191***
	(0.046)	(0.028)	(0.025)	(0.023)	(0.022)	(0.020)
ln S <sub>urmt</sub>	-0.063***	-0.040***	-0.016***	$-0.004^{*}$	-0.003	-0.004
	(0.018)	(0.005)	(0.005)	(0.002)	(0.002)	(0.002)
Observations	426032	426032	426032	426032	426032	426032
$R^2$	0.126	0.550	0.742	0.870	0.875	0.928
UPC-year FE		$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	
market-year FE		$\checkmark$	$\checkmark$	$\checkmark$		
retail-year FE		$\checkmark$				
UPC-market-year FE						$\checkmark$
retail-year-UPC FE			$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
year-market-retail FE					$\checkmark$	$\checkmark$
UPC-market-retail FE				$\checkmark$	$\checkmark$	$\checkmark$

Table A.5: UPC-market-retailer-year-level analysis

*Notes:* The regression specification is the same as the one estimated in Table A.4, with the exception that we replace the number of customers with total sales.

## A.4. Firm Sales Growth Decomposition

One concern in Figure 1 is that some firms might only appear temporarily in our data—not because of their actual behavior but due to sampling error. For example, it could be that households in our sample do not appear to purchase a firm's product even though the product was actually purchased and not recorded. In this case, the average value of sales, number of customers, and sales per customer of young firms in our analyses might be confounded with those of old firms.

To address concerns regarding the sampling error, Figure A.1 uses only those firms that appear for at least 3 or 5 consecutive years. The results still show that the number of customers is a primary factor that generates an increase in firms' sales over time. There is a steeper increase in sales in the firm's early stage than our baseline results in Figure 1. The results are intuitive, since firms that survive for several years are likely to generate more sales at the beginning relative to firms that could not survive. Overall, the robustness analysis suggests that the sampling errors are not first-order concerns in our analyses.



Figure A.1: Decomposition of Firm Sales Growth by Firm Age: Survivors

(a) Survive for 3+ Consecutive Years

(b) Survive for 5+ Consecutive Years

*Notes:* Figures A.1a and Figure A.1b replicate Figure 1 by using firms that appear in the sample for at least 3 and 5 consecutive years, respectively. There are 32,242 number of observations and 6,400 firms used in Figures A.1a and 19,603 number of observations and 2,997 firm used in Figure A.1b.

Another concern is that firms might sell their products in a different number of months over the following years. For example, some firms might enter the market in late November or December but sell their products over many months in subsequent years. To adjust for these differences, we calculate the average monthly sales over a year per firm and redo the decomposition exercise in Figure A.2. There is a smaller increase in firms' sales at age 1, which suggests that some firms enter during the late months of the initial year. However, the relative importance of the number of customers in explaining sales remains the same, and accounts for approximately 70% of sales on average.



Figure A.2: Decomposition of Firm Sales Growth by Firm Age: Monthly Sales

Notes: Figure A.2 replicates Figure 1 by using average monthly sales per firm and year instead of yearly sales.

Also, one might be worried that the empirical pattern we observe might not apply to products outside of our sample, which is restricted to products with a barcode. For example, it could be that for more durable products that customers purchase occasionally, firms might not be able to grow as much through the extensive margin of demand because they may not face the same customers every year.

Given that there is no other consumer-producer matched dataset (to the best of our knowledge), we use our data—which cover a substantial fraction of consumer goods with a wide variety of products—to understand the underlying differences between durable and non-durable products. We closely follow Argente, Lee, and Moreira (2019) and define product group-level durability by using information on the number of shopping trips. We count the average yearly number of trips customers made to purchase products in each product group and divide product groups into durable and non-durable products based on the median value of average trips. The set of durable products include, for example, "LIGHT BULBS, ELECTRIC GOODS," "HARDWARE, TOOLS," and "AUTOMOTIVE," and the non-durable products include "MILK," "SNACKS," and "BEER."

Figure A.3 presents the results. Regardless of whether we are analyzing durable or non-durable products, firms mainly grow by expanding their customer bases. The relevance of the number of customers for firm growth remains when redefining durable goods based on the 75th percentile of the trips distribution or analyzing the variance decomposition of total sales of durable or non-durable products.



Figure A.3: Decomposition of Firm Sales Growth by Firm Age: By Durability

*Notes:* Figures A.3a and Figure A.3b replicate Figure 1 by dividing firms based on the durability of the products they sell. There are 19,988 observations and 5,050 firms used in Figure A.3a and 29,816 observations and 7,323 firms used in Figure A.3b.

## A.5. Customer Acquisition and Firms' Non-production Costs

In this section, we study the extent to which firms are able to control their growth through different sales margins. To do so, we estimate the following specification:

$$\ln S_{igt} = \gamma \ln \text{SGA}_{it} + X'_{it} v + \lambda_{ig} + \lambda_{st} + \lambda_{gt} + \varepsilon_{igt}, \qquad (A.1)$$

where  $S_{igt}$  stands for sales and its components of firm *i* in product group *g* and year *t*,  $X'_{it}$  is a vector of firm-time-level control variables,  $\lambda_{ig}$  are firm-product-group fixed effects,  $\lambda_{st}$  are 2-digit SIC-year fixed effects, and  $\lambda_{gt}$  are product-group-year fixed effects.<sup>34</sup> The vector of controls  $X'_{it}$  includes lagged total sales and lagged total number of customers, which allow us to compare firms with similar relative sizes and customer bases. The coefficient of interest is  $\gamma$ , which captures the correlation between total sales (and its components) and SGA expenses.

As shown in the first column of Table A.6, firms that spend more on SGA expenses have larger sales. Moreover, the second and third columns show that approximately 95% (0.090/0.095) of the correlation between sales and SGA expenses is due to the correlation between non-production costs and the number of customers, not to the correlation with the average sales per customer. Finally, the last two columns further decompose the correlation of SGA expenses with the size of firms' customer bases into the acquisition of new customers and the retention of old customers. To measure these outcomes, we only include households that appear in the Nielsen data in at least two consecutive periods.<sup>35</sup> We find that while there is a strong correlation between SGA expenses and

<sup>&</sup>lt;sup>34</sup>Since sales in the Nielsen data vary across both a detailed product category in the Nielsen data ("product group") and the major industry code available in Compustat data ("SIC"), we allow for both product-group fixed effects and firm-SIC-code fixed effects to compare products within fine product categories.

<sup>&</sup>lt;sup>35</sup>There is a change in the number of surveyed households every year in the Nielsen data, especially in 2006 and 2007.

	Decom	position o	of ln S <sub>igt</sub>	ln $m_{igt}$ : New vs. Old		
	(1) ln <i>S</i>	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$		(4) ln <i>m</i> <sup>New</sup>	(5) In <i>m</i> <sup>Old</sup>	
ln SGA <sub>it</sub>	0.095***	0.005	0.090***	0.095***	0.016	
	(0.036)	(0.014)	(0.028)	(0.032)	(0.027)	
Observations	13131	13131	13131	13131	13131	
$R^2$	0.962	0.909	0.965	0.943	0.961	
Firm-year Controls	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	
Group-year FE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	
SIC-year FE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	
Group-firm FE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	

Table A.6: Sales and SGA: A Decomposition

*Notes*: \* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01; standard errors are two-way clustered at the firm and product-group level. *S* denotes total sales, pq the average sales per customers, *m* the number of customers,  $m^{\text{New}}$  the number of new customers, and  $m^{\text{Old}}$  the number of old customers. New customers are defined as customers who do not purchase firm *i*'s products in group *g* at time t - 1 but start to purchase those products at time *t*, whereas old customers are the customers who consecutively purchase firm *i*'s products in group *g* in both t - 1 and t ( $m = m^{\text{New}} + m^{\text{Old}}$ ). We only use households that consecutively appear in t - 1 and t in the Nielsen data when we measure  $m^{\text{New}}$  and  $m^{\text{Old}}$ . SIC industries correspond to a two-digit SIC code. All Nielsen variables are projection-factor adjusted.

the number of new customers, the regression coefficient for old customers is neither economically nor statistically significant. These results show that the non-production costs of firms are associated with the acquisition of new customers rather than retaining the existing customer base.

**Robustness** Below, we redo the analysis for subcomponents of SGA expenses. The results are consistent for subcomponents that are related to expansionary activities (advertising and rent expenses). We find a positive and statistically significant relationship between each of these SGA subcomponents and sales, which arises from the customer acquisition margin.<sup>36</sup>

Given our results in the theory, the extent to which firms can grow through customer acquisition depends on a *variable* component of SGA expenses, which has been a topic of discussion in the recent literature (see Traina, 2019, De Loecker, Eeckhout, and Unger, 2020). The correlation between sales and SGA expenses in Table A.6 is indicative of such a variable nature of these costs. However, as we show below, this contemporaneous correlation is weaker than the one between sales and COGS, which is commonly considered to be a measure of variable production costs. Therefore, total SGA expenses seem to be composed of both variable and fixed components. Although our

One concern in using these data is that the entry and exit of customers could arise from entry and exit in the sample. For example, the entry of new customers may not reflect the actual customer acquisition of a firm but may instead reflect an increase in the number of surveyed customers who already purchased this firm's products. To address this concern, new customers are defined as households that are present in the sample in periods t - 1 and t but only buy the product at t. The exit of customers is similarly defined.

<sup>&</sup>lt;sup>36</sup>On the other hand, other SGA subcomponents do not show a statistically significant relationship with sales at conventional levels, which supports the view that the correlation between sales and SGA expenses is due to firms' expansionary activities.

reduced-form empirical analyses do not allow us to separate these components, in our model we incorporate both components and provide a strategy to measure how much firms can grow through investing in their customer bases.

**A.5.1. Sales and SGA Expenses: Decomposition by Durability of Products.** Using the same durability measure we constructed and used in Figure A.3, we analyze the potential heterogeneity in the relationship between SGA expenses and sales based on a product's durability. Table A.7 reports the results. We find no statistically significant differences in this relationship by the durability of products.

	Decon	nposition	ln $m_{igt}$ : New vs. Old		
	(1) ln <i>S</i>	(2) In <i>pa</i>	(3) ln <i>m</i>	(4) In <i>m</i> <sup>New</sup>	(5) In <i>m</i> <sup>Old</sup>
ln SGA <sub>it</sub>	0.079**	-0.006	0.085***	0.078***	0.024
	(0.036)	(0.017)	(0.030)	(0.030)	(0.038)
ln SGA <sub>it</sub> x Durability	0.040	0.026	0.013	0.042	-0.018
	(0.046)	(0.028)	(0.032)	(0.039)	(0.046)
Observations	13131	13131	13131	13131	13131
$R^2$	0.962	0.909	0.965	0.943	0.961
Firm-year Controls	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Group-year FE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
SIC-year FE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Group-firm FE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$

Table A.7: Sales and SGA Expenses: Decomposition by Durability of Products

*Notes:* The regression specification is the same as the one estimated in Table A.6, with the exception that we include as additional regressors the durability measure and its interaction with the log of SGA expenses.

**A.5.2.** Sales and the Components of SGA Expenses. This section revisits the correlation of sales and SGA expenses by analyzing the subcomponents of the latter. In the main body of the paper, we rely on total SGA expenses as the main cost variable since it is reported for almost all the firms in the dataset. To be more explicit on the nature of these expenses, we further investigate the subcomponents of SGA expenses that are available for a subset of our sample. Among all the components reported in the data, the ones that are strongly correlated with firms' sales are related to expansionary activities: advertising (e.g., awareness of the existence of products/firms) and rent expenses (geographic proximity). We analyze these expansionary components to provide further evidence on the correlation between SGA expenses and sales.

Tables A.8 and A.9 show the correlation between sales and its components with advertising and rent expenses, respectively. Both tables show the same empirical relationship reported in Table A.6: The correlation between sales and the expansionary components of SGA costs arises from the correlation with the number of customers and, in particular, with the acquisition of new customers.

Overall, these results corroborate the empirical evidence of endogenous customer acquisition.<sup>37</sup>

Our analysis also reveals that other components of SGA (such as Foreign Exchange Income, Staff Expense, Receivables, and State Income Taxes) are not correlated with sales (results not reported), consistent with the view that total SGA expenses correlate with a firm's sales through their expansionary components.

	Decom	position	$\ln m_{igt}$ : New vs. Old		
	(1)	(2)	(3)	(4)	(5)
	ln S	ln <i>pq</i>	$\ln m$	$\ln m^{ m New}$	ln m <sup>Old</sup>
ln ADV <sub>it</sub>	0.073**	-0.005	$0.078^{***}$	0.077**	0.019
	(0.032)	(0.012)	(0.028)	(0.031)	(0.028)
Observations	11239	11239	11239	11239	11239
$R^2$	0.964	0.920	0.966	0.947	0.966
Firm-year Controls	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Group-year FE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
SIC-year FE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Group-firm FE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$

Table A.8: Sales and Advertising Expenses

*Notes:* \* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01. The regression specification is the same as the one estimated in Table A.6, with the exception that we replace total SGA expenses with advertising expenses.

	Decom	position o	$\ln m_{igt}$ : New vs. Old		
	(1)	(2)	(3)	(4)	(5)
	ln S	ln <i>pq</i>	ln <i>m</i>	$\ln m^{ m New}$	ln m <sup>Old</sup>
ln Rent <sub>it</sub>	0.086***	-0.004	0.090***	0.099***	0.027
	(0.030)	(0.015)	(0.025)	(0.029)	(0.030)
Observations	12552	12552	12552	12552	12552
$R^2$	0.962	0.907	0.965	0.943	0.960
Firm-year Controls	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Group-year FE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
SIC-year FE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Group-firm FE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$

### Table A.9: Sales and Rent Expenses

*Notes:* \* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01. The regression specification is the same as the one estimated in Table A.6, with the exception that we replace total SGA expenses with rent expenses.

<sup>&</sup>lt;sup>37</sup>For the particular case of R&D expenses, which could also be interpreted as an expansionary activity, the point estimates obtained are similar in magnitude to the ones reported in Tables A.8 and A.9. However, the small fraction of firms reporting these expenses decreases the precision of the estimates.

**A.5.3. The Semi-variable Nature of SGA.** This section establishes that the SGA expenses have a semi-variable nature and correlate with firms' short-run sales. Previous studies that examined the non-production cost of firms (SGA) made polar opposite assumptions on the variable nature of this cost. For instance, in measuring price-cost markups, Traina (2019) includes non-production costs as variable costs, whereas De Loecker, Eeckhout, and Unger (2020) interpret non-production costs as fixed costs in their baseline approach. We empirically assess the validity of such assumptions by comparing the comovement of sales and SGA expenses with the comovement of sales and other costs that are commonly assumed to be variable and fixed in the short run in the literature: COGS expenses and investment.





#### (a) Restricted Sample

#### (b) Full Sample

*Notes:* The figure shows the binned scatter plots of the correlation between the quarterly change in log sales and the quarterly change in (i) log SGA expenses, (ii) log COGS expenses, and (iii) log stock of capital for firms in the quarterly Compustat dataset. We also plot the best linear fit for each variable. The correlations control for quarter and firm fixed effects. In Panel (a), we restrict the sample to observations with a change in the log of sales between -0.1 and +0.1 (the 25th and 75th percentiles of the quarterly change of log sales are -8% and 11%, respectively). Panel (b) plots the binned scatter plot using the full sample. We adopt the perpetual inventory method following Traina (2019) and use Gross and Net Capital (PPEGT and PPENT) and deflate investment with NIPA's non-residential fixed investment good deflator to measure the capital stock. There are 17,168 firms in the 1964-2016 period used in this analysis.

Our results suggest that SGA expenses have both variable and fixed components; They are more variable than the capital expenditure but less than COGS expenses. Figure A.4a reports the binned scatter plot of changes in sales against changes in firm's costs for a range of  $\Delta \ln S_{i,t}$  between -10% and 10%, which are approximately the 25th and 75th percentiles of the  $\Delta \ln S_{it}$  distribution. Consistent with the view in the literature (e.g., De Loecker, Eeckhout, and Unger 2020), sales exhibit the largest comovement with production costs ( $\beta = 0.894$ ; SE 0.008) and the lowest comovement with investment ( $\beta = 0.081$ ; SE 0.008). Figure A.4b shows that similar relationships hold in the full sample.

We consider other empirical specifications to confirm the semi-variable nature of SGA expenses.

Table A.10 presents the regression results that correspond to Figure A.4. The semi-variable nature of SGA expenses is clear in this table, both with and without fixed effects. We also show that R&D expenses (a subcomponent of SGA costs) are more variable than the stock of capital but less variable than total SGA expenses. Finally, Figure A.5 reports the coefficients of a regression of SGA expenses on leads and lags of total sales, which further supports the short-run variability of SGA expenses: Although there is a strong correlation between SGA expenses and contemporaneous sales, we find no large correlation with future or past sales.

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	$\Delta ln(\text{COGS}_{it})$	$\Delta ln(COGS_{it})$	$\Delta ln(SGA_{it})$	$\Delta ln(SGA_{it})$	$\Delta ln$ (Capital <sub>it</sub> )	$\Delta ln$ (Capital <sub>it</sub> )	$\Delta ln(\text{R\&D}_{it})$	$\Delta ln(\text{R\&D}_{it})$
$\Delta lnS_{it}$	0.920***	$0.894^{***}$	0.473***	$0.405^{***}$	0.133***	0.081***	0.278***	0.200***
	(0.008)	(0.008)	(0.009)	(0.010)	(0.008)	(0.008)	(0.028)	(0.031)
$R^2$	0.055	0.162	0.011	0.103	0.002	0.155	0.002	0.083
Quarter FE	No	Yes	No	Yes	No	Yes	No	Yes
Firm FE	No	Yes	No	Yes	No	Yes	No	Yes
N	293962	292739	292785	291547	134743	133745	69845	69363

Table A.10: The Semi-variable Nature of SGA

*Notes:* Dependent variables are the quarterly change in log COGS expenses, SGA expenses, and the stock of capital. The estimation method used in all columns is OLS. Standard errors (in parentheses) are clustered at the firm level.



Figure A.5: Correlation between SGA expenses and leads and lags of sales

(a) Trimmed Sample

(b) Full Sample

*Notes:* Figure A.5a uses the trimmed sample presented in Figure A.4b and Figure A.5b uses the full sample. 95% confidence intervals are presented for every estimate. Figure A.5a replicates column (4) in Table A.10 at quarter = 0.

## **B** Derivations

## B.1. Lemma 1

*Proof.* We start from the expression for the optimal price of the firm:

$$\ln\left(\frac{p_{i,t}}{D_t}\right) = \ln(\varepsilon_{i,t}) - \ln(\varepsilon_{i,t} - 1) + \ln\left(\frac{mc_{i,t}}{D_t}\right),\tag{B.1}$$

where

$$\varepsilon_{i,t} = -\frac{\partial \ln(q_{i,t})}{\partial \ln(p_{i,t})} = \frac{\sigma}{1 - \eta \ln(p_{i,t}) + \eta \ln(D_t(1 - \sigma^{-1}))}.$$
(B.2)

The last equality in Equation (B.2) follows from the expression of demand per match in Equation (3.7). Differentiating Equation (B.1) we have

$$d\ln\left(\frac{p_{i,t}}{D_t}\right) = (1-\mu_{i,t})d\ln(\varepsilon_{i,t}) + d\ln\left(\frac{mc_{i,t}}{D_t}\right) = \frac{1}{1+\eta\sigma^{-1}\varepsilon_{i,t}(\mu_{i,t}-1)}d\ln\left(\frac{mc_{i,t}}{D_t}\right),$$

where  $\mu_{i,t} \equiv \frac{\varepsilon_{i,t}}{\varepsilon_{i,t}-1} \ge 1$  is the firm's markup. Then, it follows that

$$d\ln(\mu_{i,t}) = d\ln(p_{i,t}) - d\ln(mc_{i,t}) = -\frac{\eta\sigma^{-1}\varepsilon_{i,t}(\mu_{i,t}-1)}{1 + \eta\sigma^{-1}\varepsilon_{i,t}(\mu_{i,t}-1)}d\ln\left(\frac{mc_{i,t}}{D_t}\right).$$
(B.3)

Restricting  $d \ln(D_t) = 0$  and considering partial changes in  $\ln(mc_{i,t})$ —e.g., moving in the crosssection of firms toward firms with higher marginal costs within a particular time when  $D_t$  is fixed—give us the expressions of interest. The sign restrictions follow from  $\mu_{i,t} \ge 1$  and  $\varepsilon_{i,t} \ge 0$ .

## **B.2.** Proposition 1

*Proof.* Consider the sales per match of firm *i* normalized by the demand index  $D_t$ ,  $p_{i,t}q_{i,t}/D_t$ . Differentiating the log of this quantity, we have

$$d\ln\left(\frac{p_{i,t}q_{i,t}}{D_t}\right) = (1 - \varepsilon_{i,t})d\ln\left(\frac{p_{i,t}}{D_t}\right) = -\frac{\varepsilon_{i,t} - 1}{1 + \eta\sigma^{-1}\varepsilon_{i,t}(\mu_{i,t} - 1)}d\ln\left(\frac{mc_{i,t}}{D_t}\right).$$
(B.4)

Therefore, combining this expression with Equation (B.3), we have

$$d\ln(\mu_{i,t}) = \eta \sigma^{-1} \mu_{i,t} (\mu_{i,t} - 1) d\ln\left(\frac{p_{i,t} q_{i,t}}{D_t}\right).$$
(B.5)

Again, restricting  $d \ln(D_t) = 0$  and considering partial changes in sales per customer—e.g., moving in the cross-section of firms toward firms with higher sales per customer within a particular period when  $D_t$  is fixed—we get

$$\frac{d\ln(\mu_{i,t})}{d\ln(p_{i,t}q_{i,t})}\Big|_{d\ln(D_t)=0} = \eta \sigma^{-1} \mu_{i,t}(\mu_{i,t}-1) \ge 0, \tag{B.6}$$

where the sign restriction follows from  $\mu_{i,t} \ge 1$ .

#### **B.3. Corollary 1**

*Proof.* Recall that a firm's relative total sales within a period is given by

$$\frac{p_{i,t}y_{i,t}}{\int_{i\in N_t} p_{i,t}y_{i,t}di} = \frac{p_{i,t}y_{i,t}}{C_t} = p_{i,t}q_{i,t}m_{i,t}.$$
(B.7)

Now, restricting  $p_{i,t}y_{i,t}/C_t = \bar{s}$ , we have

$$0 = d\ln(p_{i,t}y_{i,t})\Big|_{p_{i,t}y_{i,t}/C_t = \bar{s}} = d\ln(p_{i,t}q_{i,t})\Big|_{p_{i,t}y_{i,t}/C_t = \bar{s}} + d\ln(m_{i,t})\Big|_{p_{i,t}y_{i,t}/C_t = \bar{s}}$$
(B.8)

$$\Rightarrow d\ln(p_{i,t}q_{i,t})\Big|_{p_{i,t}y_{i,t}/C_t=\bar{s}} = -d\ln(m_{i,t})\Big|_{p_{i,t}y_{i,t}/C_t=\bar{s}}.$$
(B.9)

Now, using Equation (B.6) we have

$$\frac{\partial \ln(\mu_{i,t})}{\partial \ln(m_{i,t})}\Big|_{p_{i,t}y_{i,t}/C_t=\bar{s}} = -\frac{\partial \ln(\mu_{i,t})}{\partial \ln(p_{i,t}q_{i,t})} = -\eta\sigma^{-1}\mu_{i,t}(\mu_{i,t}-1) \le 0, \tag{B.10}$$

where the sign restriction follows from  $\mu_{i,t} \ge 1$ .

## B.4. Proposition 2

*Proof.* This relationship is obtained directly from the first-order condition of the firm's problem with respect to  $m_{i,t}$ . For the rest of the proof, we derive this first-order condition.

We start by showing that the firm's customer acquisition constraint always binds (meaning that the firm never disposes of its existing customers). To show this, note that it cannot be the case that  $l_{i,s,t} > 0$  but  $m_{i,t} < (1-\delta)m_{i,t-1} + \frac{l_{i,s,t}^{\phi}}{P_{m,t}}$  since the firm can keep the same  $m_{i,t}$  with a lower  $l_{i,s,t}$ . Thus, if  $m_{i,t} < (1-\delta)m_{i,t-1} + \frac{l_{i,s,t}^{\phi}}{P_{m,t}}$ , then optimality requires that  $l_{i,s,t} = 0$ . Now, suppose  $l_{i,s,t} = 0$  but  $m_{i,t} < (1-\delta)m_{i,t-1}$ . Note that in this case, the slope of the firm's "production" profit function with respect to  $m_{i,t}$  is given by

$$\frac{\partial}{\partial m_{i,t}}(p_{i,t}y_{i,t} - W_t l_{i,p,t}) = (p_{i,t} - mc_{i,t})\frac{y_{i,t}}{m_{i,t}} > 0,$$
(B.11)

where the last equality follows from the fact that for any choice of  $q_{i,t} > 0$ , the firm's markup is always strictly larger than 1 and hence  $p_{i,t} > mc_{i,t}$ . Thus, the firm's profit is strictly increasing in  $m_{i,t}$  and since  $m_{i,t} < (1 - \delta)m_{i,t-1}$ , then the firm can increase its  $m_{i,t}$  at no cost and gain more profits at time *t*. Moreover, this will not affect firms' profits in the future, since the firms can always dispose of the increase in  $m_{i,t}$  in the next period at no cost. Hence, optimality requires that  $m_{i,t} = (1 - \delta)m_{i,t-1} + \frac{l_{i,s,t}^{\phi}}{P_{m,t}}$ .

Now, in writing firm *i*'s problem at time *t*, replace  $l_{i,p,\tau} = (y_{i,\tau}/z_{i,\tau})^{\alpha^{-1}}$ ,  $y_{i,\tau} = m_{i,\tau}q_{i,\tau}C_{\tau}$ , and  $l_{i,s,\tau} = P_{m,\tau}^{\phi^{-1}}(m_{i,\tau} - (1-\delta)m_{i,\tau-1})^{\phi^{-1}}$  to obtain the problem as

$$\max_{\{p_{i,\tau},m_{i,\tau},q_{i,\tau}\}_{\tau\geq t}} \mathbb{E}_t \sum_{\tau\geq t} (\beta \nu)^{\tau-t} C_{\tau}^{-\gamma} \left(\prod_{h=t}^{\tau} \mathbf{1}_{i,h}\right) \times$$

$$\begin{bmatrix} p_{i,\tau} m_{i,\tau} q_{i,\tau} C_{\tau} - W_{\tau} \left( \frac{m_{i,\tau} q_{i,\tau} C_{\tau}}{z_{i,\tau}} \right)^{\alpha^{-1}} - W_{\tau} P_{m,\tau}^{\phi^{-1}} (m_{i,\tau} - (1-\delta) m_{i,\tau-1})^{\phi^{-1}} - W_{\tau} \chi$$
  
s.t.  $q_{i,\tau} = \left[ 1 - \eta \ln \left( \frac{p_{i,\tau}}{D_{\tau} (1 - \sigma^{-1})} \right) \right]^{\frac{\sigma}{\eta}}.$ 

Next, if  $\mathbf{1}_{i,t} = 0$ , then  $l_{i,s,t} = 0$ . However, conditional on  $\mathbf{1}_{i,t} = 1$ , the FOC with respect to  $m_{i,t}$  is

$$0 = \mathbf{1}_{i,t} (p_{i,t} - \alpha^{-1} \frac{W_t l_{i,p,t}}{y_{i,t}}) q_{i,t} C_t - \mathbf{1}_{i,t} \phi^{-1} \frac{W_t l_{i,s,t}}{m_{i,t} - (1 - \delta) m_{i,t-1}} + \beta \nu (1 - \delta) \mathbb{E}_t \left[ \left( \frac{C_{t+1}}{C_t} \right)^{-\gamma} \mathbf{1}_{i,t+1} \phi^{-1} \frac{W_{t+1} l_{i,s,t+1}}{m_{i,t+1} - (1 - \delta) m_{i,t}} \right].$$

Replacing  $mc_{i,t} = \alpha^{-1} \frac{W_t l_{i,p,t}}{y_{i,t}}$  and iterating the FOC forward gives us the expression of interest.

## **B.5.** Proposition 3

*Proof.* Recall that from Equation (3.16) the relationship between a firm's production labor share and markup is given by

$$\frac{W_t l_{i,p,t}}{p_{i,t} y_{i,t}} = \frac{\alpha}{\mu_{i,t}}.$$
(B.12)

Moreover, assuming  $\delta = 1$ , we can use the characterization of the firm's optimal marketing strategy in Equation (3.21) to write their marketing labor share of an operating firm as

$$\phi^{-1} \frac{W_t l_{i,s,t}}{m_{i,t}} = (p_{i,t} - mc_{i,t}) q_{i,t} C_t \Leftrightarrow \frac{W_t l_{i,s,t}}{p_{i,t} y_{i,t}} = \phi(1 - \mu_{i,t}^{-1}).$$
(B.13)

Combining these two equations, we get that

$$W_t l_{i,s,t} = \phi p_{i,t} y_{i,t} - \phi \alpha^{-1} W_t l_{i,p,t}.$$
 (B.14)

Finally, notice that

$$SGA_{i,t} \equiv W_t \chi + W_t l_{i,s,t} = \underbrace{SGAF_t}_{=W_t \chi} + \underbrace{\phi Sales_{i,t}}_{=p_{i,t} y_{i,t}} - \frac{\phi}{\alpha} \underbrace{COGS_{i,t}}_{=W_t l_{i,p,t}}.$$
(B.15)

## **B.6.** Proposition 4

*Proof.* This result can be derived from combining Equations (B.12) and (B.13):

$$\frac{W_t(l_{i,p,t}+l_{i,s,t})}{p_{i,t}y_{i,t}} = \alpha \mu_{i,t}^{-1} + \phi(1-\mu_{i,t}^{-1}).$$
(B.16)

Notice that this is strictly decreasing in  $\mu_{i,t}$  if and only if  $\alpha > \phi$ . Hence, the firm's revenue productivity of labor, the inverse of the equation above, is increasing in  $\mu_{i,t}$  if and only if  $\alpha > \phi$ .

## **B.7.** Proposition 5

*Proof.* Let  $(C_t, L_{p,t}, L_{s,t}, L_t, N_t)_{t \ge 0}$  denote the equilibrium allocation. A log-linearization of  $U(C_t, L_t) =$ 

 $\frac{C_t^{1-\gamma}}{1-\gamma} - \xi \frac{L_t^{1+\psi}}{1+\psi}$  around this allocation gives

$$\Delta U(C_t, L_t) = C_t^{1-\gamma} \Delta \ln(C_t) - \xi L_t^{\psi} L_{p,t} (\Delta \ln(L_{p,t}) + \frac{L_{s,t}}{L_{p,t}} \Delta \ln(L_{s,t}) + \chi \frac{N_t}{L_{p,t}} \Delta \ln(N_t)) + \mathcal{O}(\|.\|^2).$$
(B.17)

Next, divide by  $U_{c,t}C_t = C_t^{1-\gamma}$  and use the household's optimal labor supply condition  $\xi \frac{L_t^{\psi}}{C_t^{-\gamma}} = W_t$  to get

$$\frac{\Delta U(C_t, L_t)}{U_{c,t}C_t} = \Delta \ln(C_t) - \frac{W_t L_{p,t}}{C_t} (\Delta \ln(L_{p,t}) + \frac{L_{s,t}}{L_{p,t}} \Delta \ln(L_{s,t}) + \chi \frac{N_t}{L_{p,t}} \Delta \ln(N_t)) + \mathcal{O}(\|.\|^2).$$
(B.18)

Finally, using the aggregate production function in Equation (5.2), replace  $\Delta \ln(C_t) = \Delta \ln(Z_t) + \alpha \Delta \ln(L_{p,t})$ , and using the definition of the aggregate markup in Equation (5.6), replace the labor share in terms of the cost-weighted markup  $\left(\frac{W_t L_{p,t}}{C_t} = \frac{\alpha}{\mathcal{M}_t}\right)$  to get

$$\frac{\Delta U(C_t, L_t)}{U_{c,t}C_t} \approx \Delta \ln(Z_t) + \alpha (1 - \mathcal{M}_t^{-1}) \Delta \ln(L_{p,t}) - \alpha \mathcal{M}_t^{-1} (\frac{L_{s,t}}{L_{p,t}} \Delta \ln(L_{s,t}) + \chi \frac{N_t}{L_{p,t}} \Delta \ln(N_t)).$$
(B.19)

## **B.8.** Planner's Problem

The planner's problem for this economy is given by

$$\max_{\substack{\{(c_{i,j,t})_{j\in[0,1]},(\mathbf{1}_{i,t})_{i\in N_{t}-1}\cup\Lambda_{t},\\(\delta_{i,t},m_{i,t},l_{i,p,t},l_{i,s,t})_{i\in N_{t}},C_{t}\}_{t\geq0}}\sum_{t=0}^{\infty}\beta^{t}\left[\frac{C_{t}^{1-\gamma}}{1-\gamma}-\xi\frac{L_{t}^{1+\psi}}{1+\psi}\right]$$
(B.20)

subject to

$$\int_{0}^{1} \mathbf{1}_{\{j \in m_{i,t}\}} c_{i,j,t} dj = z_{i,t} l_{i,p,t}^{\alpha}, \quad \forall i \in N_t,$$
(B.21)

$$\int_{i \in N_t} \int_0^1 \mathbf{1}_{\{j \in m_{i,t}\}} \Upsilon\left(\frac{c_{i,j,t}}{C_t}\right) dj di = 1$$
(B.22)

$$\int_{i \in N_t} (l_{i,p,t} + l_{i,s,t} + \chi) di = L_t$$
(B.23)

$$\int_{i \in N_t} m_{i,t} di = 1 \tag{B.24}$$

$$N_{t} = \{i \in N_{t-1} \cup \Lambda_{t} : \mathbf{1}_{i,t} v_{i,t} = 1\}, \quad N_{-1} \text{ given.}$$

$$\mathbf{1} = \{i \in N_{t-1} \cup \Lambda_{t} : \mathbf{1}_{i,t} v_{i,t} = 1\}, \quad N_{-1} \text{ given.}$$
(B.25)

$$m_{i,t} = (1 - \delta_{i,t})m_{i,t-1} + \frac{1 - \int_{i \in N_t} (1 - \delta_{i,t})m_{i,t-1}}{\int_{i \in N_t} l_{i,s,t}^{\phi} di} l_{i,s,t}^{\phi}, \quad \forall i \in N_t,$$
(B.26)

$$\delta_{i,t} \in [\delta, 1], l_{i,s,t} \ge 0, \quad \forall i \in N_t, \tag{B.27}$$

$$\int_{i \in N_t} l_{i,s,t} di = \bar{L}_{s,t} > 0.$$
(B.28)

Here, Equation (B.21) requires that for every firm, their supply meets their allocated demand; Equation (B.22) is the Kimball aggregator that implicitly defines  $C_t$  given the planner's allocation of demand; Equation (B.23) requires that labor supply meets demand for labor from production, advertising and overhead costs; Equation (B.24) requires that the matching market clears; Equation (B.25) is the law of motion for the set of operating firms given an entry/exit policy by the planner; Equation (B.26) determines firm *i*'s evolution of customers given an allocation for advertising; Equation (B.27) requires the nonnegativity of labor for advertising and the constraint that while the planner can separate customers from firms, the separation rate should be at least  $\delta$ ; and finally Equation (B.28) requires that the planner at least spends  $\bar{L}_{s,t}$  on advertising. This last constraint is for an arbitrary but strictly positive  $\bar{L}_{s,t}$ —in Lemma 2 we show that the level of this quantity does not matter for the optimal distribution of customers, which is also well defined in the limit when  $\bar{L}_{s,t} \rightarrow 0$ .

## B.9. Lemma 2

*Proof.* Suppose that at any given time t, a choice for  $N_t$  is fixed. Suppose next that the planner desires to allocate matches according to a rule

$$\mathscr{A}: (i \to m_{i,t}^*)_{i \in N_t}. \tag{B.29}$$

Note that this can be any arbitrary allocation of matches as long as it is feasible:

$$m_{i,t}^* \ge 0, \quad \forall i \in N_t \tag{B.30}$$

$$\int_{i \in N_t} m_{i,t}^* di = 1.$$
(B.31)

To show that the allocation  $\mathscr{A}$  is implementable on  $N_t$  for any given level of  $\bar{L}_{s,t}$ , we need to show that (1) it is generated by a choice of  $(\delta \leq \delta_{i,t} \leq 1, l_{i,s,t} \geq 0)_{i \in N_t}$ , and (2) it is feasible  $\int_{i \in N_t} l_{i,s,t} di = \bar{L}_{s,t}$ . We show this by construction. In particular, consider the choice

$$\left(\delta_{i,t}^{*} = 1, l_{i,s,t}^{*} = \bar{L}_{s,t} \frac{m_{i,t}^{*\phi^{-1}}}{\int_{i \in N_{t}} m_{i,t}^{*\phi^{-1}} di}\right)_{i \in N_{t}}$$
(B.32)

That is, first, let the planner separate all the matches from their corresponding firms ( $\delta_{i,t}^* = 1$ ) and then reallocate them based on  $\mathscr{A}$ . It follows that Equations (B.30) and (B.31) from above hold by construction. Next, to verify that these values implement  $\mathscr{A}$ , observe that

$$m_{i,t} \equiv (1 - \delta_{i,t}) m_{i,t-1} + \left(1 - \int_{i \in N_t} (1 - \delta_{i,t}) m_{i,t-1} di\right) \frac{l_{i,s,t}^{*\phi}}{\int_{i \in N_t} l_{i,s,t}^{*\phi} di} = m_{i,t}^*.$$
 (B.33)

Finally, to confirm feasibility, note that  $\int_{i \in N_t} l_{i,s,t}^* di = \overline{L}_{s,t}$ .

## **B.10.** Proposition 6

*Proof.* To prove this Proposition, we proceed in two steps. First, we characterize the optimal demand per customer and the allocation of customers for a given set of operating firms,  $N_t$ , over time. Second, we show that the *optimal* allocation for these two objects maximizes the aggregate TFP, given  $N_t$ , subject to feasibility constraints for  $(q_{i,t}, m_{i,t})_{i \in N_t}$ .

Moreover, throughout this proof we rely on the result from Lemma 2 and directly characterize the distribution of matches  $m_{i,t}$ , ignoring the constraints in Equations (B.26) and (B.27) as well as Equation (B.28). Having characterized this distribution, we can then use the result from Lemma 2 to find the allocation of advertising labor that implements it.

**Step 1: Optimal Allocation of Demand.** The results in this Proposition follow from the first-order conditions of the planner's problem in Equation (B.20), fixing the planner's other choices at an arbitrary allocation. In the remainder of this proof, we characterize these first-order conditions.

Formally, for firm *i* and period *t*, let  $\beta^t \varphi_{c,i,t} di$  be the shadow cost on Equation (B.21); for period *t*, let  $\beta^t \varphi_{Y,t}$ ,  $\beta^t \varphi_{L,t}$  and  $\beta^t \varphi_{m,t}$  be the shadow costs on Equation (B.22), Equation (B.23) and Equation (B.24), respectively. Moreover, similar to the equilibrium allocation, let us define  $q_{i,j,t} \equiv \frac{c_{i,j,t}}{C_t}$ . It is straight forward to show that for  $j \notin m_{i,t}$ ,  $q_{i,j,t} = 0$ . So hereafter, we only refer to  $q_{i,j,t}$  when  $j \in m_{i,t}$ .

The first-order conditions with respect to  $q_{i,j,t}$  are

$$\varphi_{c,i,t}C_t = \Upsilon'(q_{i,j,t})\varphi_{\Upsilon,t}, \quad \forall j \in m_{i,t}.$$
(B.34)

It immediately follows that all households that are matched to a variety consume the same amount:

$$q_{i,j,t} = q_{i,t}, \quad \forall j \in m_{i,t}. \tag{B.35}$$

Replacing this result in Equation (B.21) and Equation (B.22) and taking the first-order condition with respect to  $m_{i,t}$ , we have

$$\varphi_{c,i,t}q_{i,t}C_t + \varphi_{m,t} = \Upsilon(q_{i,t})\varphi_{\Upsilon,t}.$$
(B.36)

Notice that in deriving this first-order condition, we have ignored the constraint in Equation (B.26). The reason we can do this goes back to Lemma 2, which states that any choice of  $(m_{i,t})_{i \in N_t}$  can be implemented without any loss of generality. Therefore, we can ignore the constraint in Equation (B.26) and use Lemma 2 to show that it is satisfied.

Next, replacing Equation (B.35) in Equation (B.34), multiplying it by  $q_{i,t}$  and subtracting it from Equation (B.36), we have

$$\varphi_{m,t} = \left[\Upsilon(q_{i,t}) - q_{i,t}\Upsilon'(q_{i,t})\right]\varphi_{\Upsilon,t}.$$
(B.37)

Since  $\varphi_{\Upsilon,t} \neq 0$ ,<sup>38</sup> it follows that

$$\Upsilon(q_{i,t}) - q_{i,t}\Upsilon'(q_{i,t}) = \frac{\varphi_{m,t}}{\varphi_{\Upsilon,t}}, \quad \forall i \in N_t.$$
(B.38)

Notice that the left-hand-side of this equation is only a function of  $q_{i,t}$  and is strictly monotonic in  $q_{i,t} > 0$ .<sup>39</sup> Moreover, the right-hand-side of the equation is only a function of time-*t* shadow costs

<sup>&</sup>lt;sup>38</sup>To see why, suppose not. Then, by Equation (B.34), either  $C_t = 0$ —which is clearly not optimal since marginal utility approaches infinity as  $C_t \rightarrow 0$ —or  $\varphi_{c,i,t} = 0$ —which means that the household can freely supply infinite labor to firm *i* at *t* and is also a contradiction since it violates the positive disutility of the labor supply.

<sup>&</sup>lt;sup>39</sup>Observe that  $D_x[\Upsilon(x) - \Upsilon'(x)x] = -\Upsilon''(x)x > 0.$ 

and is independent of *i*. Hence, there exists a unique  $q_t^*$  such that

$$q_{i,t} = q_t^* \quad \forall i \in N_t. \tag{B.39}$$

Replacing this last equation into Equation (B.22) we have

$$\int_{i \in N_t} m_{i,t} \Upsilon(q_t^*) di = 1 \Rightarrow \Upsilon(q_t^*) = 1 \Rightarrow q_t^* = 1,$$
(B.40)

where the second statement uses the market-clearing condition for matches in Equation (B.24) and the last statement uses the strict monotonicity of  $\Upsilon(x)$  and the fact that  $\Upsilon(1) = 1$ .

Given that the social planner sets  $q_{i,t} = 1$  for all firms, this implies that firms' production will differ under the efficient allocation only through different numbers of customers. To determine the optimal level of production, we only need to consider the first-order condition with respect to  $l_{i,p,t}$ :

$$\alpha \varphi_{c,i,t} z_{i,t} l_{i,p,t}^{\alpha-1} = \varphi_{L,t}. \tag{B.41}$$

Dividing this equation by the first-order condition for  $m_{i,t}$  in Equation (B.34), we get

$$z_{i,t}l_{i,p,t}^{\alpha-1} = \frac{C_t}{\alpha \Upsilon'(1)} \frac{\varphi_{L,t}}{\varphi_{\Upsilon,t}}.$$
(B.42)

Solving for  $l_{i,p,t}$  from this equation and replacing it Equation (B.21) we have

$$\int_{j\in m_{i,t}} c_{i,j,t} dj = m_{i,t} C_t = z_{i,t} \left( \frac{C_t}{z_{i,t} \alpha \Upsilon'(1)} \frac{\varphi_{L,t}}{\varphi_{\Upsilon,t}} \right)^{\frac{\alpha}{\alpha-1}} \Rightarrow m_{i,t} = \left( \frac{z_{i,t}}{C_t} \right)^{\frac{1}{1-\alpha}} \left( \frac{\varphi_{L,t}}{\alpha \Upsilon'(1)\varphi_{\Upsilon,t}} \right)^{\frac{\alpha}{\alpha-1}}.$$
 (B.43)

Finally, imposing the market-clearing condition for matches in Equation (B.24) we get

$$m_{i,t} = \frac{z_{i,t}^{\frac{1}{1-\alpha}}}{\int_{i \in N_t} z_{i,t}^{\frac{1}{1-\alpha}} di}, \quad \forall i \in N_t.$$
(B.44)

**Step 2: Optimal Demand Allocation Maximizes Aggregate TFP.** First, let us derive aggregate TFP. Consider an allocation of  $l_{i,p,t}$  across firms in  $N_t$ . Then, aggregate production labor is given by

$$L_{p,t} = \int_{i \in N_t} l_{i,p,t} di = \int_{i \in N_t} \left( \frac{C_t \int_{j \in m_{i,t}} q_{i,j,t} dj}{z_{i,t}} \right)^{\alpha^{-1}} di$$
(B.45)

where the second equality follows from the fact that for every firm, demand must meet supply. Rearranging this gives us an expression for TFP as a function of demand allocation within  $N_t$ :

$$C_{t} = \underbrace{\left[\int_{i \in N_{t}} \left(\frac{z_{i,t}}{\int_{j \in m_{i,t}} q_{i,j,t} dj}\right)^{-\alpha^{-1}} di\right]^{-\alpha}}_{Z_{t} = \text{Aggregate TEP}} \times L_{p,t}^{\alpha}$$
(B.46)

 $Z_t \equiv \text{Aggregate TFP}$ 

Now, maximizing  $Z_t$  is equivalent to minimizing  $Z_t^{-\alpha^{-1}}$ . Thus, the problem of maximizing  $Z_t$  subject to the Kimball aggregator becomes

$$\min_{\{q_{i,j,t},\mathbf{1}_{\{j\in m_{i,t}\}}\}} \int_{i\in N_t} \left(\frac{\int_0^1 \mathbf{1}_{\{j\in m_{i,t}\}} q_{i,j,t} dj}{z_{i,t}}\right)^{\alpha^{-1}} di$$
(B.47)

$$s.t. \int_{i \in N_t} \int_0^1 \mathbf{1}_{\{j \in m_{i,t}\}} \Upsilon(q_{i,j,t}) dj di = 1$$
(B.48)

$$\int_{i \in N_t} \int_0^1 \mathbf{1}_{\{j \in m_{i,t}\}} dj di = 1$$
(B.49)

First, fix an allocation for  $m_{i,t}$ . Notice that since  $\alpha^{-1} > 1$ , the objective function is convex in each  $q_{i,j,t}$  for any  $j \in m_{i,t}$  so the first-order condition for  $q_{i,j,t}$  is sufficient for optimality. Allowing  $\eta_{Y,t}$  and  $\eta_{m,t}$  to be the multipliers on Equations (B.48) and (B.49), the first-order condition for  $q_{i,j,t}$ ,  $j \in m_{i,t}$  reads

$$\frac{1}{\alpha z_{i,t} \eta_{\Upsilon,t}} \left( \frac{\int_0^1 \mathbf{1}_{\{j \in m_{i,t}\}} q_{i,j,t} dj}{z_{i,t}} \right)^{\alpha^{-1}-1} = \Upsilon'(q_{i,j,t}).$$
(B.50)

Since the left-hand-side of this equation depends only on firm/time level variables, it follows that all the consumers in the  $m_{i,t}$  consume the same  $q_{i,j,t}$ —i.e.,  $q_{i,j,t} = q_{i,t}^*$  where  $q_{i,t}^*$  solves the above equation. Substituting this into the objective, we can rewrite the choice of  $\mathbf{1}_{\{j \in m_{i,t}\}}$  as simply choosing the size of the customer base for firms, which with slight abuse of notation we also denote as  $m_{i,t}$ . Formally,

$$\int_0^1 \mathbf{1}_{\{j \in m_{i,t}\}} q_{i,j,t} dj = m_{i,t} q_{i,t}^*$$
(B.51)

Again, it is straightforward to see that the objective is convex in  $m_{i,t}$  and the first-order condition for it reads

$$\frac{q_{i,t}^*}{\alpha z_{i,t}} \left(\frac{m_{i,t} q_{i,t}^*}{z_{i,t}}\right)^{\alpha^{-1}-1} + \eta_{m,t} = \Upsilon(q_{i,t}^*)\eta_{\Upsilon,t}.$$
(B.52)

Subtracting the first-order condition for  $m_{i,t}$  from the first-order condition for  $q_{i,j,t}$  we arrive at

$$\frac{\eta_{m,t}}{\eta_{\Upsilon,t}} = \Upsilon(q_{i,t}^*) - q_{i,t}^* \Upsilon'(q_{i,t}^*).$$
(B.53)

Note that the right-hand-side is a monotonic function of  $q_{i,t}^*$  (since  $\Upsilon$  is increasing and concave) and the left-hand-side does not depend on *i*. So  $q_{i,t}^* = q_t^*$  for some  $q_t^*$ . Substituting  $q_{i,j,t} = q_t^*$  in Equation (B.48) then yields that  $\Upsilon(q_t^*) = 1$ , which implies  $q_t^* = 1$ . Substituting this into the FOCs and using the constraints, we then retrieve the optimal  $m_{i,t}^*$  as

$$m_{i,t}^{*} = \frac{z_{i,t}^{\frac{1}{1-\alpha}}}{\int_{i \in N_{t}} z_{i,t}^{\frac{1}{1-\alpha}} di}$$
(B.54)

Hence, choosing the allocation of demand to maximize the aggregate TFP gives the optimal allocation of demand under the social planner's problem.

## B.11. Proposition 7

*Proof.* Our goal here is to derive the social value of a firm that enters the economy at a certain point in time. To do so, we start by writing a (partial) Lagrangian for the planner's problem and rearrange it to isolate all the terms in the Lagrangian that are related to a particular firm. In writing this Lagrangian, we employ two of the results we have established so far. First, given our result in Lemma 2, we ignore the constraints in Equations (B.26) and (B.27) as well as Equation (B.28). Instead, we can allow the social planner to directly choose the distribution of matches at any given time and use the results in the proof of Lemma 2 to substitute for the optimal advertising labor. Second, to render the comparison to the equilibrium value of firms natural, we take advantage of Equation (B.35), which establishes that the social planner chooses the same  $q_{i,j,t}$  for all households that are matched to firm *i* at time *t*. Aside from these two results, however, we refrain from imposing any other optimality conditions at the firm level to maximize the comparability of the firms' social value with their equilibrium value.

Thus, we can form the following partial Lagrangian for the social planner (with the Lagrange multipliers defined as in Appendix **B.10**):

$$\max_{\substack{(1_{i,t})_{i\in N_{t-1}\cup\Lambda_{t},}\\(q_{i,t},m_{i,t},l_{i,p,t})_{i\in N_{t}},C_{t},L_{t}}}\sum_{t=0}^{\infty}\beta^{t}\left[\frac{C_{t}^{1-\gamma}}{1-\gamma}-\xi\frac{L_{t}^{1+\psi}}{1+\psi}+\varphi_{\Upsilon,t}\left(\int_{i\in N_{t}}m_{i,t}\Upsilon(q_{i,t})-1\right)\right.\\\left.+\varphi_{L,t}\left(L_{t}-\int_{i\in N_{t}}(l_{i,p,t}+\bar{L}_{s,t}+\chi)di\right)+\varphi_{m,t}\left(1-\int_{i\in N_{t}}m_{i,t}di\right)\right] \tag{B.55}$$

subject to the constraints in Equation (B.21) (i.e., demand meets supply) and Equation (B.25), which captures the law of motion for the set of firms active in the economy. To rearrange this partial Lagrangian so that it resembles the problem of the firms in the equilibrium, consider the following two operations: First, rearrange the terms within the brackets to put all integrals of the form  $\int_{i \in N_t} di$  in one bracket. Second, factor out  $\varphi_{Y,t}$  from these integrals. These two operations lead to the following expression for the partial Lagrangian:

$$\max_{\substack{(\mathbf{l}_{i,t})_{i\in N_{t-1}\cup\Lambda_{t},}\\(q_{i,t},m_{i,t},l_{i,p,t})_{i\in N_{t}},C_{t},L_{t}}}_{(q_{i,t},m_{i,t},l_{i,p,t})_{i\in N_{t}},C_{t},L_{t}}} \mathbb{E}_{0} \sum_{t=0}^{\infty} \beta^{t} \left[ \frac{C_{t}^{1-\gamma}}{1-\gamma} - \xi \frac{L_{t}^{1+\psi}}{1+\psi} - \varphi_{\Upsilon,t} + \varphi_{L,t} \left(L_{t} - \bar{L}_{s,t}\right) + \varphi_{m,t} \right] \\
+ \mathbb{E}_{0} \sum_{t=0}^{\infty} \beta^{t} \int_{i\in N_{t}} \varphi_{\Upsilon,t} \left( m_{i,t}\Upsilon(q_{i,t}) - \frac{\varphi_{L,t}}{\varphi_{\Upsilon,t}} (l_{i,p,t} + \chi) - \frac{\varphi_{m,t}}{\varphi_{\Upsilon,t}} m_{i,t} \right) di$$
(B.56)

The second line is what we are after, once we switch the order of the sum and the integration: Since this expression first integrates over the set of all operating firms at time *t* and then sums over time, we can isolate the contribution of every firm to this Lagrangian separately by switching the order of integration and summation. In doing so, we will take advantage of the law of motion for the

set  $N_t$ , Equation (B.25), as a function of the planner's choice for the entry/exit of firms as well as endogenous exit shocks. Formally, pick a firm *i* that enters the economy for the first time at time *t* (i.e.,  $t = \min_{\tau \ge 0} \{i \in N_\tau\}$ , with  $t \equiv -1$  indicating firms that were in the initial distribution of firms in the economy). Then, *i* will be in any  $N_\tau$ , for  $\tau \ge t$  as long as  $\prod_{h=t}^{\tau} v_{i,h} \mathbf{1}_{i,h} = 1$ . Therefore, the second line in the equation above can be written as

$$\mathbb{E}_{0}\sum_{t=0}^{\infty}\beta^{t}\int_{i\in N_{t}}\varphi_{\Upsilon,t}\left(m_{i,t}\Upsilon(q_{i,t})-\frac{\varphi_{L,t}}{\varphi_{\Upsilon,t}}(l_{i,p,t}+\chi)-\frac{\varphi_{m,t}}{\varphi_{m,t}}m_{i,t}\right)di$$
(B.57)

$$=\mathbb{E}_{0}\sum_{t=0}^{\infty}\beta^{t}\int_{i\in\tilde{\Lambda}_{t}}\mathbb{E}_{t}\sum_{\tau=t}^{\infty}(\beta\nu)^{\tau-t}\left(\prod_{h=t}^{\tau}\mathbf{1}_{i,h}\right)\varphi_{\Upsilon,\tau}\left(m_{i,\tau}\Upsilon(q_{i,\tau})-\frac{\varphi_{L,\tau}}{\varphi_{\Upsilon,\tau}}(l_{i,p,\tau}+\chi)-\frac{\varphi_{m,\tau}}{\varphi_{\Upsilon,\tau}}m_{i,\tau}\right)di,\quad(B.58)$$

where  $\tilde{\Lambda}_0 = \Lambda_0 \cup N_{-1}$  and  $\tilde{\Lambda}_t = \Lambda_t$  for all  $t \ge 1$ .

Finally, substituting Equation (B.57) into Equation (B.56), we arrive at the following reformulation of the social planner's problem:

$$\max_{\{C_{t},L_{t}\}_{t\geq0}} \mathbb{E}_{0} \sum_{t=0}^{\infty} \beta^{t} \left\{ \frac{C_{t}^{1-\gamma}}{1-\gamma} - \xi \frac{L_{t}^{1+\psi}}{1+\psi} - \varphi_{Y,t} + \varphi_{L,t} \left(L_{t} - \bar{L}_{s,t}\right) + \varphi_{m,t} + \varphi_{Y,t} \int_{i\in\tilde{\Lambda}_{t}} \max_{\left\{ \substack{\mathbf{1}_{i,\tau}, m_{i,\tau} \\ q_{i,\tau}, l_{i,p,\tau} \right\}_{\tau\geq t}} \mathbb{E}_{t} \sum_{\tau=t}^{\infty} (\beta v)^{\tau-t} \left( \prod_{h=t}^{\tau} \mathbf{1}_{i,h} \right) \frac{\varphi_{Y,\tau}}{\varphi_{Y,t}} \left( m_{i,\tau} Y(q_{i,\tau}) - \frac{\varphi_{L,\tau}}{\varphi_{Y,\tau}} (l_{i,p,\tau} + \chi) - \frac{\varphi_{m,\tau}}{\varphi_{Y,\tau}} m_{i,\tau} \right) di \right\}$$
  
Social Value of Firm  $i \equiv G_{i,t}$ 

s.t. 
$$m_{i,\tau}q_{i,\tau}C_{\tau} = z_{i,\tau}l^{\alpha}_{i,p,\tau}, \quad \forall i \in \tilde{\Lambda}_t, \forall \tau \ge t,$$
 (B.60)

where  $G_{i,t}$  is the culmination of all of the terms in the social planner's problem in which a firm *i* shows up and thus captures the value of the firm for the planner.

To compare  $G_{i,t}$  with the equilibrium value of the firm, the final step is to relate the Lagrange multipliers to aggregate objects that are similar to the equilibrium wage  $W_t$  and demand index  $D_t$ . Let us derive the first-order conditions for the aggregate consumption  $C_t$  and aggregate labor supply  $L_t$ . For  $C_t$ , we get

$$C_t^{*-\gamma} = \varphi_{c,i,t} \int_{i \in N_t} m_{i,t}^* q_{i,t}^* di = \varphi_{\Upsilon,t} \int_{i \in N_t} m_{i,t}^* q_{i,t}^* \Upsilon'(q_{i,t})^* di C_t^{*-1},$$
(B.61)

(B.59)

where stars denotes that these equations hold under the optimal choice of these variables and the right-hand-side follows from substituting the FOC for  $c_{i,j,t}$  from above. Moreover, note that the integral on the right-hand-side is exactly how we defined the demand index in the equilibrium. So let us define it similarly for the social planner as

$$D_t^* = \left[\int_{i \in N_t} m_{i,t}^* q_{i,t}^* \Upsilon'(q_{i,t}^*)\right]^{-1} \Rightarrow \varphi_{\Upsilon,t} = C_t^{*1-\gamma} D_t^*.$$
(B.62)

Furthermore, the FOC for  $L_t$  gives:  $\xi L_t^{*\psi} = \varphi_{L,t}$ . If we define the notion of the "wage" for the planner

as the marginal rate of substitution between leisure and consumption,  $W_t^* \equiv \xi L_t^{*\psi} C_t^{*\gamma}$ , we can write

$$\frac{\varphi_{L,t}}{\varphi_{\Upsilon,t}} = \frac{W_t^*}{C_t^* D_t^*}.$$
(B.63)

Now, to derive an expression for  $\varphi_{m,t}/\varphi_{Y,t}$ , we use Equation (B.38). Multiplying both sides of that equation by  $m_{i,t}$  and integrating over the set  $N_t$ , we arrive at

$$\frac{\varphi_{m,t}}{\varphi_{\Upsilon,t}} = \frac{\overbrace{\int_{i \in N_t} m_{i,t} \Upsilon(q_{i,t}) di}^{=1} - \overbrace{\int_{i \in N_t} m_{i,t} q_{i,t} \Upsilon'(q_{i,t}) di}^{=D_t^{*-1}}}{\underbrace{\int_{i \in N_t} m_{i,t} di}_{=1}} = 1 - D_t^{*-1}.$$
 (B.64)

The last step in comparing  $G_{i,t}$  with the equilibrium value of firms defined in Equation (3.13) is to render their units consistent. The equilibrium value of a firm entering at time t in Equation (3.13) is defined in the units of aggregate consumption at time t, where  $G_{i,t}$  is in time t utils; to convert time t utils to time t aggregate consumption units we need to multiply  $G_{i,t}$  by  $C_t^* D_t^*$ . Let  $v_{i,t}^*$  denote this converted value; then it is given by

$$\nu_{i,t}^{*} \equiv C_{t}^{*} D_{t}^{*} G_{i,t}$$

$$= \max_{\substack{\left\{1_{i,\tau}, m_{i,\tau}\right\}\\q_{i,\tau}, l_{i,p,\tau}\right\}_{\tau \ge t}}} \mathbb{E}_{t} \sum_{\tau=t}^{\infty} (\beta \nu)^{\tau-t} \left(\prod_{h=t}^{\tau} \mathbf{1}_{i,h}\right) \left(\frac{C_{\tau}^{*}}{C_{t}^{*}}\right)^{-\gamma} \left(D_{\tau}^{*} \frac{\Upsilon(q_{i,\tau})}{q_{i,\tau}} y_{i,\tau} - W_{t}^{*}(l_{i,p,\tau} + \chi) - C_{\tau}^{*}(D_{\tau}^{*} - 1) m_{i,\tau}\right)$$
(B.65)

s.t. 
$$y_{i,\tau} = m_{i,\tau} q_{i,\tau} C_{\tau}^* = z_{i,\tau} l_{i,p,\tau}^{\alpha}$$
 (B.66)

## **Supplemental Materials**

Concentration, Market Power, and Misallocation: The Role of Endogenous Customer Acquisition

by Hassan Afrouzi, Andres Drenik, and Ryan Kim
# SM.1 Data Description

This section further describes the process to clean the Nielsen and Compustat data used in this paper. For this, we follow the previous studies that use these two datasets (Hottman, Redding, and Weinstein 2016 for Nielsen and De Loecker, Eeckhout, and Unger 2020, Traina 2019 for Compustat).

## SM.1.1. Nielsen and GS1: Variables and Data Cleaning

We construct the following variables from the Nielsen Homescan Panel, Promodata, and GS1 company data:

- UPC: scanner-level product identifier available in Nielsen Homescan Panel and Promo data.
- *GS1 Company Identifier*: GS1 company ID. GS1 provides both UPC and firm identifiers along with company name and headquarters location. This information allows us to identify firm boundaries in the Nielsen data and further merge the data with Compustat.
- *Sales*: We define sales as the sum of the total expenditures of households at different levels of aggregation: UPC-year, firm-year, and group-firm-year. We use sample weights (projection factor) to render the Nielsen household sample representative at the national level.
- *Number of Customers*: For each firm, product (UPC), or group-firm, we aggregate the number of households with the sample weight adjustment.
- *Sales per Customer*: At each level of aggregation, we define average sales per customer as total sales divided by the number of customers.

## Sample Selection of Homescan Panel data The Homescan Panel data we use covers 2004-2016.

- 1. Following Neiman and Vavra (2019), we balance the product modules across years to exclude the effect of entry and exit of modules, which mainly arises from name changes and potentially adds measurement error.
- 2. To render the sample representative, we drop "magnet products" in Nielsen data, which are fresh produce and other items without barcodes.
- 3. We drop products that do not have a product group identifier to render the analysis consistent across different specifications.
- 4. There is a small number of observations for which the sampling year of the household is different from the year their purchases took place. These reflect the fact that households were sampled in late December. While the corresponding purchases are recorded as the current year, their household panel years are recorded as the following year. We drop these observations to use coherent sample weights across households and years.

We restrict the sample in the Promodata to the years 2006-2011 since the data are incomplete for the other years. Following the manual for these data, we use both active and inactive files and drop duplicate observations in inactive files. We adjust for multi-package and unit size when using the UPC-level information; we do the same when we merge with price information from Homescan Panel data. We exclude the small number of markets that are not common across the Homescan Panel and Promodata.

## SM.1.2. Compustat: Variables and Data Cleaning

We download and construct the following variables from Compustat:

- *Global company key* (mnemonic gvkey): Compustat's firm ID.
- *Year* (mnemonic fyear): the fiscal year.
- *Selling, general and administrative expense* (mnemonic XSGA): the SG&A sums "all commercial expenses of operation (such as expenses not directly related to product production) incurred in the regular course of business pertaining to the securing of operating income." They include expenses such as marketing and advertising expenses, research and development, accounting expenses, delivery expenses, etc.
- *Costs of goods sold* (mnemonic COGS): the COGS sums all "expenses that are directly related to the cost of merchandise purchased or the cost of goods manufactured that are withdrawn from finished goods inventory and sold to customers." They include expenses such as labor and related expenses (including salary, pension, retirement, profit sharing, provision for bonus and stock options, and other employee benefits), operating expense, lease, rent, and loyalty expense, write-downs of oil and gas properties, and distributional and editorial expenses.
- *Operating expenses, total* (mnemonic XOPR): OPEX represents the sum of COGS, SG&A, and other operating expenses.
- *Sales (net)* (mnemonic SALE): this variable represents gross sales, for which "cash discounts, trade discounts, and returned sales and allowances for which credit is given to customer" are discounted from the final value.
- *Capital*: we calculate capital in two ways. First, we simply set capital to be equal to the gross property, plant, and equipment value (mnemonic PPEGT) deflated by the investment goods deflator from NIPA's nonresidential fixed investment good deflator (line 9). For our second measurement of capital, we use the perpetual inventory method—we set the first observation of each firm to be equal to the gross property, plant, and equipment value and for susbsequent years we add the difference from  $netPPE_t$  (mnemonic PPENT) and  $netPPE_{t-1}$ . We also deflate the difference in PPENT by the investment goods deflator from NIPA's nonresidential fixed investment good deflator.
- Company's initial public offering date (mnemonic IPODATE).
- *Age*: given that the initial public offering date was missing for a large portion of our dataset, we calculated age as the fiscal year of a given observation minus the first year we observe a firm in the dataset. According to Compustat, a firm enters the dataset after it starts providing consistent accessible annual reports trading on a U.S. exchange market, i.e., after its IPO. Following Haltiwanger, Jarmin, and Miranda (2013), we exclude the first 15 years of the dataset for all analyses using age and we group together firms older than 16 years, because we do not

know for certain a firm's IPO date for firms that were in the Compustat data since the first year.

We used NIPA Table 1.1.9. GDP deflator (line 1) to generate the real value for the variables sale, COGS, XOPR, and XSGA.

**Sample Selection** We downloaded the dataset "Compustat Annual Updates: Fundamentals Annual," from Wharton Research Data Services, from Jan 1950 to Dec 2016. The following options were chosen:

- Consolidated level: C (consolidated)
- Industry format: INDL (industrial)
- Data format: STD (standardized)
- Population source: D (domestic)
- Currency: USD
- Company status: active and inactive

We took the following steps in the cleaning process:

- 1. To select American companies, we filtered the dataset for companies with Foreign Incorporation Code (FIC) equal to "USA."
- 2. We replace industry variables (sic and naics) with their historical values whenever the historical value is not missing.
- 3. We drop utilities (sic value in the range [4900, 4999]) because their prices are very regulated and financials (sic value in the range [6000,6999]) because their balance sheets are notably different than the other firms in the analysis.
- 4. To ensure the quality of the data, we drop missing or nonpositive observations for sales, COGS, OPEX, sic 2-digit code, gross PPE, net PPE, and assets. We also exclude observations in which acquisitions are more than 5% of the total assets of a firm.
- 5. A portion of the data missing for sales, COGS, OPEX, and capital between years for firms. We input these values using a linear interpolation, but we do not interpolate for gaps longer than 2 years. This exercise inputs data for 4.6% of our sample.

## SM.1.3. Coverage of the Nielsen-Compustat Sample

Approximately 300 firms identified in Compustat can be matched with the Nielsen data for 2004-2016. Although the number of firms we matched is small, they account for a significant fraction of total sales, number of UPCs, and observations in the Nielsen Homescan Panel data, as shown in Table SM.1.1.

	Sales (b)	# of UPCs (k)	# of Obs. (m)
Nielsen-Compustat Sample	94.5	114.9	12.1
Nielsen Sample	421.2	698.9	51.6
Share (%)	22.4	16.4	23.5

Table SM.1.1: Coverage of the Nielsen-Compustat Sample

Note: Sales is the projection-factor-weighted sales in Nielsen data and is denoted in billions US dollars. # of UPCs is in thousand UPCs, and # of Obs. is in millions of observations. All variables are annual averages.

### SM.1.4. Summary Statistics

Variable	Ν	Mean	SD	p10	p50	p90
	Panel A:	Nielsen-G	S1, Firm-Pro	oduct G	roup-Year `	Variables
$S_{igt}$ (in thousands \$)	557,820	6,708.16	64,961.12	3.97	126.08	5,170.55
$p_{igt}q_{igt}$ (in \$)	557,820	10.03	20.14	1.96	5.88	19.50
$m_{igt}$ (in thousands)	557,820	500.79	2,789.82	0.82	19.83	639.87
$m_{igt}^{\text{New}}$ (in thousands)	557,820	250.34	988.95	0.48	16.03	424.89
$m_{igt}^{\text{Old}}$ (in thousands)	557,820	250.45	1,963.09	0.00	1.60	194.18
	Pan	el B: Nielse	en-Compust	at, Firm	-Year Varia	ables
$SGA_{it}$ (in millions \$)	2,101	2,009.17	4,993.82	7.63	299.09	4,882.24
COGS <sub>it</sub> (in millions \$)	2,299	7,147.11	18,558.75	17.17	1,123.66	17,251.26
OPEX <sub><i>it</i></sub> (in millions \$)	2,299	9,147.10	21,337.48	25.72	1,620.66	24,212.49
SGA-to-OPEX <sub>it</sub>	2,101	0.30	0.20	0.08	0.27	0.58
Sales-to-COGS <sub>it</sub>	2,299	1.79	1.06	1.14	1.49	2.61

#### Table SM.1.2: Summary Statistics

*Notes:* The Nielsen-GS1 data in Panel A has 40,418 firms and 109 product groups in the period of 2004-2016.  $S_{igt}$  denotes sales of firm *i* in product group *g* and time *t*,  $p_{igt}q_{igt}$  average sales per customer,  $m_{igt}$  number of customers,  $m_{igt}^{\text{New}}$  new customers in year *t* who did not purchase products in year *t* – 1, and  $m_{igt}^{\text{Old}}$  customers who purchase the products consecutively in year *t* – 1 and *t* ( $m_{igt} = m_{igt}^{\text{New}} + m_{igt}^{\text{Old}}$ ).  $S_{igt}$  is measured in thousands of US dollars, and  $m_{igt}$ ,  $m_{igt}^{\text{New}}$ , and  $m_{igt}^{\text{Old}}$  are in thousands of customers. All Nielsen variables are projection-factor adjusted. The Nielsen-Compustat matched data in Panel B have 332 firms in the period 2004-2016. All cost-side variables are in millions of US dollars and are deflated by the GDP deflator.

## SM.2 Model with Heterogeneous Tastes and Sorting

In this section, we provide two model extensions to our households' preferences. In the first part, we allow for idiosyncratic preference shocks that affect the optimal quantities consumed by each consumer. In the second, we allow for shifters that affect the perceived "quality" of each variety. We discuss how each extension affects the total demand of a firm as well as the relationship between the number of customers and average sales per customer. We also provide empirical evidence for why we abstract away from these extensions in the main analysis.

#### SM.2.1. Taste for Quantity

Consider the following extension of the Kimball aggregator in Equation (3.1):

$$\int_{0}^{N_{t}} \int_{0}^{1} \mathbf{1}_{\{j \in m_{i,t}\}} \xi_{i,j,t}^{-1} \Upsilon\left(\frac{\xi_{i,j,t} c_{i,j,t}}{C_{t}}\right) dj di = 1,$$
(SM.2.1)

where  $\xi_{i,j,t}$  is now household *j*'s "quantity" taste for variety *i* at time *t*. The implied relative demand for firm *i* at *t* is then given by

$$\frac{c_{i,t}}{C_t} = \underbrace{m_{i,t}}_{\text{number of customers}} \times \underbrace{\mathbb{E}[\xi_{i,j,t} | j \in m_{i,t}] \Upsilon'^{-1}\left(\frac{p_{i,t}}{D_t}\right)}_{q_{i,t} \equiv \text{ demand per customer}},$$
(SM.2.2)

where  $\mathbb{E}[\xi_{i,j,t} | j \in m_{i,t}]$  is the average quantity taste of firm *i*'s customers. The introduction of such shocks will affect firm *i*'s demand depending on the nature of sorting between customers and firms. If there is no sorting or selection regarding who is matched to a firm, then  $\mathbb{E}[\xi_{i,j,t} | j \in m_{i,t}] = \mathbb{E}[\xi_{i,j,t}]$  and this extension replicates the demand function in the main text. However, sorting generates a correlation between the number of customers and average sales per customer. For example, with positive sorting, firms with a larger customer base should sell less per customer on average (because the marginal customer always buys less than the average customer). Instead, with negative sorting (e.g., when marginal consumers are the ones who experiment with the product and buy more than the average consumer), firms with a larger customer base should sell more per customer on average. More generally, this model implies the following relationship between average sales per customer and a firm's number of customers:

$$\ln(p_{i,t}q_{i,t}) = \beta_0 \ln(m_{i,t}) + \ln\left(p_{i,t} \Upsilon'^{-1}\left(\frac{p_{i,t}}{D_t}\right)\right),$$
 (SM.2.3)

where the sign and magnitude of  $\beta_0$  determine the type and strength of sorting, respectively. To test this, we aggregate the Nielsen Homescan Panel data at the UPC-year-level. For this analysis, we use the whole sample available in the Homescan Panel data. We compute the price of each product as sales divided by the unit-adjusted quantity. Table SM.2.1 shows the results of OLS regressions of the log average sales per customer of UPC u at time t on the log number of customers and polynomials of log price that approximate the nonlinear function  $\Upsilon'^{-1}(\cdot)$ . Across all specifications, we consistently find a small and positive relationship between the size of the customer base and average sales per customer. Firms with 1% more customers sell on average 0.03% more per customer. Given the economic insignificance of this estimate, we do not consider this kind of preference heterogeneity in our model, which, as discussed below, is a conservative choice.

	(1)	(2)	(3)	(4)	(5)
ln m <sub>it</sub>	0.026***	0.025***	0.023**	0.029***	0.029***
	(0.009)	(0.009)	(0.009)	(0.002)	(0.002)
ln p <sub>it</sub>	0.162***	0.180***	0.132***	0.680***	0.688***
	(0.028)	(0.037)	(0.027)	(0.017)	(0.017)
ln p <sup>2</sup> <sub>it</sub>		0.021*	0.025	-0.003	-0.003
I II		(0.012)	(0.015)	(0.005)	(0.005)
ln n <sup>3</sup> :			0.005	-0.001	-0.001
p n			(0.003)	(0.001)	(0.001)
Observations	9097076	9097076	9097076	8452707	8452707
$R^2$	0.086	0.096	0.102	0.851	0.852
UPC FE				$\checkmark$	$\checkmark$
Year FE				$\checkmark$	
Product Group-Year FE					$\checkmark$

Table SM.2.1: Average Sales per Customer and the Size of the Customer Base

*Notes:* \* p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01; Standard errors are clustered by product group. The variable  $p_{ut}q_{ut}$  denotes sales per customer of UPC u at time t,  $m_{ut}$  is the number of customers, and  $p_{ut}$  is the price.

**Implications of Sorting for Misallocation Results** A potential concern in our analysis of efficiency losses from the misallocation of demand is that the sorting of customers might exacerbate or reduce welfare losses from misallocation. For example, since with positive sorting the marginal customer values a variety less than the average customer, the marginal value of allocating customers to more productive firms might not be as high as in our baseline model. However, regression results in Table SM.2.1 show that, if anything, there is negative—albeit small—sorting (i.e., the marginal customer buys more than the average customer). Viewed through the lens of this model extension, these results indicate that the quantified welfare losses from the misallocation of demand in Section 5 provide a *lower* bound on the actual welfare losses once this small negative sorting is taken into account.

#### SM.2.2. Taste for Quality

Another source of preference heterogeneity can be the different perception of the "quality" of a variety across customers. For this, consider the following extension of the Kimball aggregator in Equation (3.1):

$$\int_{0}^{N_{t}} \int_{0}^{1} \mathbf{1}_{\{j \in m_{i,t}\}} \zeta_{i,j,t} \Upsilon\left(\frac{c_{i,j,t}}{C_{t}}\right) dj di = 1$$
(SM.2.4)

where  $\zeta_{i,j,t}$  is now household *j*'s "quality" taste for variety *i* at time *t*. We assume for any *i* and *t* that  $\zeta_{i,j,t}$  is i.i.d. and its distribution is scaled so that its unconditional mean is 1 ( $\mathbb{E}[\zeta_{i,j,t}] = 1$ ). The

implied relative demand of firm *i* is then given by

$$q_{i,t} \equiv \frac{c_{i,t}}{C_t} = \int_0^1 \mathbf{1}_{\{j \in m_{i,t}\}} \Upsilon'^{-1} \left(\frac{p_{i,t}}{\zeta_{i,j,t} D_t}\right) dj.$$
(SM.2.5)

The curvature of  $\Upsilon'^{-1}(\cdot)$ , and thus the elasticity of demand, now interacts with the distribution of  $\zeta_{i,j,t}$ . However, we can show that up to a first-order approximation, sorting would imply a correlation between average demand per customer and the number of customers, similar to the one derived for quantity shocks. To see this, note that up to a first-order approximation around the point  $\frac{p_{i,t}}{\mathbb{E}[\zeta_{i,j,t}]D_t}$ :

$$\Upsilon'^{-1}\left(\frac{p_{i,t}}{\zeta_{i,j,t}D_t}\right) \approx \Upsilon'^{-1}\left(\frac{p_{i,t}}{D_t}\right) \left(1 + \sigma_{i,t}(\zeta_{i,j,t} - 1)\right),\tag{SM.2.6}$$

where  $\sigma_{i,t}$  is the elasticity of demand evaluated at the average quality taste,  $p_{i,t}/D_t$ , and hence depends only on the price. Now, using this approximation, we can write demand as

$$\frac{c_{i,t}}{C_t} \approx \underbrace{m_{i,t}}_{\text{number of customers}} \times \underbrace{\Upsilon'^{-1}\left(\frac{p_{i,t}}{D_t}\right)(1 + \sigma_{i,t}(\mathbb{E}[\zeta_{i,j,t}|j \in m_{i,t}] - 1)),}_{q_{i,t} \equiv \text{demand per customer}}$$
(SM.2.7)

where  $\mathbb{E}[\zeta_{i,j,t}|j \in m_{i,t}]$  is the average quality taste of firm *i*'s customers. Similar to quantity shocks, if there is positive (negative) sorting, then  $\mathbb{E}[\zeta_{i,j,t}|j \in m_{i,t}]$  is a decreasing (increasing) function of  $m_{i,t}$  and we should see a negative (positive) relationship between  $q_{i,t}$  and  $m_{i,t}$ , which brings us to the same argument in the previous section for quantity shocks.

## SM.3 Computational Appendix

In this section, we present the details of the computation algorithm that solves and calibrates the model. We first describe the recursive representation of the firm's problem. Then, we describe the law of motion of firms and characterize the stationary distribution. Next, we describe the algorithm that solves the model and the algorithm used in the calibration.

### SM.3.1. Solution Method

**Firm's Recursive Problem** In period *t*, a firm that decided to operate with a customer base of  $m_{-1}$  and productivity *z* solves the following dynamic programming problem (which corresponds to the sequential representation in Equation (3.13)):

$$v_{t}(m_{-1}, z) = \max_{l_{s}, l_{p}, p} \left\{ py - W_{t}l_{p} - W_{t}(l_{s} + \chi) + \beta v \frac{U_{c, t+1}}{U_{c, t}} \mathbb{E} \left[ V_{t+1}(m, z') | z \right] \right\}$$
  
s.t.  $q = \left[ 1 - \eta \ln \left( \frac{p}{D_{t}(1 - \sigma^{-1})} \right) \right]^{\frac{\sigma}{\eta}}$   
 $y = mqC_{t} = zl_{p}^{\alpha}$   
 $m = (1 - \delta)m_{-1} + \frac{l_{s}^{\phi}}{P_{m, t}},$ 

where  $V_t(m_{-1}, z) \equiv \max\{0, v_t(m_{-1}, z)\}$  denotes the endogenous exit choice.

**Stationary Distribution** Let  $\mathcal{N}_t : M \times Z \to [0, 1]$  denote the cdf of incumbent firms measured after the realization of idiosyncratic productivity shocks, but before exit decisions are made. The law of motion of the distribution of firms is given by

$$\mathcal{N}_{t}(m, z') = \int_{M \times Z} F(z'|z) \mathbf{1}_{\{m_{t}^{*}(m_{-1}, z) \leq m\}} \mathbf{v} \mathbf{1}_{\{v_{t}(m_{-1}, z) \geq 0\}} d\mathcal{N}_{t-1}(m_{-1}, z)$$
  
+  $\lambda \int_{M \times Z} F(z'|z) \mathbf{1}_{\{m_{t}^{*}(0, z) \leq m\}} \mathbf{1}_{\{v_{t}(0, z) \geq 0\}} dF^{e}(z),$ 

where F(z'|z) is the Markov chain given by the AR(1) productivity process,  $F^e(z)$  is the productivity distribution of potential entrants, and  $m^*(m_{-1}, z)$  denotes the optimal policy for customer acquisition.

**Solution Algorithm** In steady state consumption is constant, so that  $U_{c,t+1}/U_{c,t} = 1$ . The algorithm for the numerical solution of the steady state of the model is as follows:

**Step 0:** Set up a grid for firm's state  $S = M \times Z$ . We choose 15 collocation points in each dimension. For *Z*, we use the 0.0001 and 0.9999 percentiles of the ergodic distribution of the AR(1) productivity process as the grid bounds. For *M*, the lower bound of the grid is 0, and the upper bound is chosen so that the largest customer base in the solution of any version of the model is smaller than the bound. Given these bounds, we construct power grids to concentrate grid points at lower values for  $m_{-1}$  and *z*.

**Step 1:** Guess values for *C*,  $\tilde{W} \equiv W/(CD)$  and  $P_m$ .

**Step 2:** Solve firm's problem given  $(C, \tilde{W}, P_m)$  by scaling the value function by 1/(CD) (this reduces the number of aggregate variables we need to solve for by one). We solve this problem by using projection methods to approximate both the value function  $v_t(m_{-1}, z)$  and its expected value  $\mathbb{E}[V_{t+1}(m, z')|z]$ . We approximate these functions with the tensor product of a linear spline in the *z* dimension and a cubic spline in the  $m_{-1}$  dimension. We follow a two-step procedure to compute optimal policies. First, for a given candidate *m*, we compute *q*,  $l_s$ , and  $l_p$  by solving the nonlinear FOC for *q* and using the production function and the law of motion of matches. Second, to optimize the value function with respect to *m*, we use the golden search method. Having approximated these values and guessed a vector of the spline's coefficients, we combine an iteration procedure and a Newton solver to find the coefficient of the basis function. To compute the expectation in  $\mathbb{E}[V_{t+1}(m, z')|z]$ , we rely on the following approximation:

$$\mathbb{E}\left[V_t(m, z')|z\right] = \sum_{i=1}^{50} \omega_i V_t\left(m, \exp(\rho \ln z + \varepsilon_i)\right).$$
(SM.3.1)

To construct the nodes  $\varepsilon_i$ , we generate an equidistant grid of 50 points from 0.0001 to 0.9999 and invert the CDF of the  $\mathcal{N}(0, \sigma_z^2)$  distribution. To construct the weights  $\omega_i$ , we discretize the normal distribution with a histogram centered around the nodes.

**Step 3:** To approximate the ergodic distribution of firms, we construct a finer grid with 100 and 500 points in the  $m_{-1}$  and z direction, respectively. Then, we solve the firm's problem once on the new grid using the approximation to the value functions from the previous step. To find the ergodic distribution, we rely on the nonstochastic simulation approach of Young (2010). This method approximates the distribution of firms on a histogram based on the finer grid. Since both optimal policies and productivity shocks are allowed to vary continuously, we assign values of m and z that do not fall on points in the grid in the following way. Let  $s \equiv (m_{-1}, z)$  denote a firm's state. Then, the transition matrix for a firm's customer base can be constructed as

$$Q_M(s, m'(s)) = \left[\mathbf{1}_{m'(s)\in[m_{j-1}, m_j]} \frac{m'(s) - m_j}{m_j - m_{j-1}} + \mathbf{1}_{m'(s)\in[m_j, m_{j+1}]} \frac{m_{j+1} - y'(s)}{m_{j+1} - m_j}\right]$$
(SM.3.2)

for all states *s* in the grid. That is, the transition matrix allocates firms in the histogram based on the proximity of the optimal policy to each point in the finer grid. The transition matrix for productivity shocks is approximated as  $Q_Z = \sum_{i=1}^{200} \omega_i Q_{z,i}$ , where  $Q_{z,i}$  is similarly constructed as in Equation (SM.3.2) for  $z'(s) = \exp(\rho \ln z + \varepsilon_i)$ . The overall transition matrix is then given by  $Q = Q_Z \otimes Q_M$ . Finally, the distribution of firms is obtained by iterating until convergence the approximation to the law of motion

$$\mathcal{N} = Q' \left( v \mathbf{1}_{v(\mathbf{s}) \ge 0} \mathcal{N} + \lambda \mathbf{1}_{v(\mathbf{s}) \ge 0} F^e \right),$$

where  $F^e$  is an approximation of the distribution of entrants on the finer grid.

**Step 4:** Compute aggregate variable *X* from firms' vectorized policies x(s) as  $X = (v \mathbf{1}_{v(s) \ge 0} \mathcal{N} + \lambda \mathbf{1}_{v(s) \ge 0} F^e)' x(s)$ . Compute the residual vector

$$1 = \int_0^N m_i di, \quad 1 = \int_0^N m_i \Upsilon(q_i) di, \text{ and } \quad 1 = \frac{W}{\xi C^{\gamma} L^{\psi}}.$$

If the distance is small, stop. Otherwise, update  $(C, \tilde{W}, P_m)$  with a Newton method and go to **Step 2**.

## SM.3.2. Estimation routine

We estimate the parameters of the model via the Simulated Method of Moments (SMM). More specifically, we choose a set of parameters  $\mathcal{P}$  that minimizes the SMM objective function

$$\left(\frac{\boldsymbol{m}_m(\mathscr{P})}{\boldsymbol{m}_d}-1\right)' W\left(\frac{\boldsymbol{m}_m(\mathscr{P})}{\boldsymbol{m}_d}-1\right),$$

where  $m_m$  and  $m_d$  are a vector of model-simulated moments and data moments, respectively, and W is a diagonal matrix. To compute the model-simulated moments, we follow these steps:

- **Step 1:** Given a vector of parameters  $\mathscr{P}$ , we find the steady state of the model. For this, we slightly modify the previous algorithm. Since in the estimation we normalize aggregate output C = Y = 1 and the normalized wage  $\tilde{W} = 1$  (see Step 1 of the solution algorithm) with the free parameters  $(\lambda, \xi)$ , we need to solve for only one aggregate variable,  $P_m$ .
- **Step 2:** We simulate 100,000 firms for 150 periods and compute model moments using data from the last 25 periods. When matching moments based on the entire US economy, we use data from all simulated firms. When matching moments based on Compustat data, we impose a filter that mimics selection into Compustat based on firm age and size. On the age dimension, we restrict the simulated sample to those firms that are at least 7 years old, as in Ottonello and Winberry (2020). On the size dimension, we restrict the sample to firms with sales above 19% of the average sales in the simulated economy. This cutoff corresponds to the ratio of the 5th percentile of the sales distribution in Compustat (USD1.06 million) to the average firm sales in SUSB (USD5.7 million) in 2012.

To minimize the SMM objective function and be confident of reaching the global minimum, we follow a two-step procedure in the spirit of Arnoud, Guvenen, and Kleineberg (2019). In the first step, we construct 500 quasi-random vectors of parameters  $\mathscr{P}$  from a Halton sequence, which is a deterministic sequence designed to cover the parameter space evenly. After computing the SMM objective in those points, we choose the 30 parameters vectors with the lowest objective values. In the second step, we initiate a local Nelder-Mead optimizer from each of the 30 starting points and select the local minimum with the lowest objective value.

## SM.4 Additional Model Analysis

## SM.4.1. Identification of Model Parameters

In this section, we formally guide the discussion of the identification of model parameters. Panel A of Figure SM.4.1 shows the *local* elasticity of simulated moments (rows) with respect to parameters (columns), evaluated at the calibrated parameters. In general, the intuition behind the choice of targets is borne out in the model. For example, a higher exogenous exit rate 1 - v mechanically increases the average exit rate. Similarly, a higher super-elasticity of demand  $\eta$  increases the comovement between revenue labor productivity and sales. A higher elasticity of the matching function  $\phi$  makes the positive relationship between a firm's SGA expenses and sales stronger and reduces the comovement between labor productivity on sales, as predicted in the simple version of the model in Propositions 3 and 4. The persistence of productivity shocks  $\rho_z$  affects multiple moments, but it affects most strongly the dispersion of the sales distribution. Finally, a smaller average productivity of entrants  $\bar{z}_{ent}$  increases the relative size of old firms.

We complement this discussion by analyzing the sensitivity measure developed by Andrews, Gentzkow, and Shapiro (2017), which shows the sensitivity of model parameters with respect to targeted moments. To render the numbers more comparable, we convert this measure into elasticities and plot  $(J'(\mathcal{P})WJ(\mathcal{P}))^{-1}J'(\mathcal{P})Wm_m(\mathcal{P})/\mathcal{P}$ , where  $J(\mathcal{P})$  is the Jacobian evaluated at calibrated parameters, W is the weighting matrix, and  $m_m(\mathcal{P})$  are the model moments evaluated at calibrated parameters. Panel B of Figure SM.4.1 shows that overhead cost  $\chi$  is quite sensitive to the average COGS-to-OPEX ratio in the data. Similarly, the elasticity of substitution  $\sigma$  is sensitive to the average production markup, and the standard deviation of productivity shock  $\sigma_z$  is most strongly influenced by the standard deviation of employment growth.

										- 2 00
Dist. relative sales (0.1)	- 13.79	-0.18	1.52	-0.14	1.51	80.07	6.01	0.18 -		2.00
Dist. relative sales (0.5)	- 1.91	-0.02	0.34	-0.05	0.37	13.77	0.95	0.01 -		
Dist. relative sales (1)	- 0.72	-0.01	0.16	-0.03	0.19	5.21	0.35	-0.01 -		1.50
Dist. relative sales (2)	- 0.08	0.00	0.06	-0.02	0.08	1.33	0.08	-0.00 -		
Dist. relative sales (5)	0.14	0.00	0.00	-0.00	0.01	-0.35	-0.03	0.00 -	-	1.00
Dist. relative sales (10)	0.09	0.00	-0.01	0.00	-0.01	-0.50	-0.04	-0.00 -		
Dist. relative sales (50)	0.00	0.00	-0.00	0.00	-0.00	-0.07	-0.00	-0.00 -	-   -	0.50
Dist. relative sales (100)	0.00	0.00	-0.00	0.00	-0.00	-0.02	-0.00	0.00 -		
Slope of OPEX-COGS on sales	2.80	0.34	-0.06	-0.11	1.41	-5.79	-0.09	-0.05 -	-	0.00
Avg. COGS-to-OPEX ratio	- 2.13	-0.17	0.32	-0.01	0.13	12.04	0.66	0.13 -		
Avg. production markup	- 0.29	-0.01	-0.28	0.06	-0.11	2.13	0.16	0.00 -	-	-0.50
Slope of labor productivity on sales	- 4.52	-0.62	-2.05	1.13	-4.22	4.16	0.37	0.08 -		
Avg. exit rate	19.01	0.33	-0.15	0.08	-0.34	-20.82	-1.10	-0.28 -	-	-1.00
Avg. relative employment (2-5)	- 0.59	0.17	-0.02	0.06	-0.43	-8.30	-0.05	-0.31 -		
Avg. relative employment (6-10)	- 1.08	0.25	0.11	0.09	-0.43	-17.88	-0.04	-0.64 -	1	-1 50
Avg. relative employment (11+)	- 10.62	0.17	0.82	-0.06	0.38	-0.32	1.24	-1.02 -		-1.50
Std. employment growth	- 0.60	0.01	0.28	-0.01	0.25	5.41	1.27	0.01 -		
·	ν	χ	σ	η	φ	$\rho_z$	σz	z <sup>ent</sup>		-2.00

## Figure SM.4.1: Parameter Identification

#### 2.00 0.26 0.02 0.62 -0.26 0.02 -0.11 0.49 1.86Dist. relative sales (0.1) -0.25 -2.23 -0.22 0.03-0.03 -0.05 -0.73 Dist. relative sales (0.5) 0.141.50 -0.05 -0.94 -0.48 -2.98 -0.09 0.01 0.07 Dist. relative sales (1) -0.89 -0.04 -1.11 -0.45 -2.57 -0.02 0.00 0.09 -0.72 Dist. relative sales (2) -0.03 -1.28 -0.45 -2.34 0.03 -0.00 0.13 -0.57 1.00 Dist. relative sales (5) -0.00 -0.33 -0.09 -0.41 0.02-0.00 0.04-0.09 Dist. relative sales (10) 0.04 0.21 -0.01 0.50 0.00 0.10 0.00 -0.00 0.05Dist. relative sales (50) 0.00 0.02 0.01 0.03 -0.00 -0.00 -0.00 0.01 Dist. relative sales (100) Slope of OPEX-COGS on sales 0.03 0.86 0.494.850.93 -0.00 -0.29 0.540.00 -0.11 -9.82 0.20 2.40 -0.02 -0.23 -2.28 Avg. COGS-to-OPEX ratio -10.50 -5.42 -0.20 -11.16 1.19 -0.00 1.17-5.76 Avg. production markup -0.50 0.00 0.01 0.730.331.790.01 -0.09 0.32Slope of labor productivity on sales -0.06 -0.40 0.100.510.050.01 -0.04 -0.63 Avg. exit rate -1.00 0.00 0.00 0.00 Avg. relative employment (2-5) 0.00 0.000.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 Avg. relative employment (6-10) -1.50 -0.81 0.01 0.04 0.55-0.00 -0.02 -1.00 Avg. relative employment (11+) 0.20-0.03 -1.77 -0.32 -2.03 -0.05 -0.09 1.250.80Std. employment growth -2.00 z<sup>ent</sup> ν $\chi$ $\sigma$ η $\phi$ $\rho_{\rm z}$ $\sigma_{\rm z}$

#### (a) Sensitivity of Moments

(b) Sensitivity of Parameters

*Notes:* Panel A shows the sensitivity of simulated moments to parameters by computing the local elasticity of moments with respect to parameters. Panel B shows the sensitivity of calibrated parameters to moments by constructing the sensitivity measure of Andrews, Gentzkow, and Shapiro (2017) and converting it into an elasticity. Both measures are evaluated at the calibrated parameters.

### SM.4.2. Calibration Results: Goodness of Fit

Table SM.4.1 and Figure SM.4.2 show the targeted moments and their model counterparts. Overall, the model closely matches the targets. The model reproduces the average cost structure very well, but it slightly under-predicts the relationship between SGA and sales. The model matches the cost-weighted average markup well and generates a similar relationship between revenue productivity of labor and sales. Figure SM.4.2 shows that the model approximates the sales distribution of firms accurately. For example, in the data 33% of firms have sales that are lower than 10% of the average sales in the economy, and 1% of firms have sales that are larger than 10 times the average sales. In the model, these shares are 25% and 1.5%. Figure SM.4.2 also shows that the model is able to replicate the relative size of old firms: 6.07 in the data and 6.4 in the model. In Appendix SM.4, we show the untargeted data and model relationships between average labor productivity and SGA with sales, and the average COGS-to-OPEX ratio by firm age and size. Although we targeted specific moments that summarize these relationships in the calibration exercise, the model matches the data patterns more broadly.

Table SM.4.1:	Targeted	Moments
---------------	----------	---------

Moment	Data	Model
Slope SGA on sales	0.492	0.474
Avg. COGS-to-OPEX ratio	0.660	0.669
Avg. cost-weighted production markup	1.250	1.275
Slope labor prod. on sales	0.036	0.033
Avg. exit rate	0.073	0.071
SD. employment growth	0.416	0.447

*Notes:* This table shows the set of moments targeted in the calibration of the model. Slope SGA on Sales refers to the OLS coefficient of the regression  $SGA_{i,t} = c + \beta Sales_{i,t} + \psi COGS_{i,t} + \varepsilon_{i,t}$ . Avg. COGS-to-OPEX ratio refers to the average of the ratio across firms. Avg. cost-weighted production markup corresponds to the COGS-weighted average markup from Edmond, Midrigan, and Xu (2022). These moments were computed using data from Compustat in 2012. Slope labor prod. on sales corresponds to the OLS coefficient of the sales-weighted regression of relative revenue labor productivity on relative sales from Edmond, Midrigan, and Xu (2022), restricting the sample of firms to those with relative sales above one. This moment was computed using data from the SUSB in 2012. The average exit rate was obtained from the BDS in 2012. The standard deviation of annual employment growth for continuing establishments is obtained from Elsby and Michaels (2013). The growth rate of variable *x* is computed as in Davis and Haltiwanger (1992):  $(x_{i,t} - x_{i,t-1})/(0.5(x_{i,t} + x_{i,t-1}))$ . The last column shows the model counterparts of each moment, which were obtained by simulating a panel of firms and computing each moment with the simulated data. In the model, we account for selection into Compustat by restricting the simulated sample of firms to those that are at least 7 years old and have sales above 19% of the average sales in the simulated economy (which corresponds to the ratio of the 5th percentile of the sales distribution in Compustat to the average sales in SUSB).

Figure SM.4.3 also plots the relationship between relative revenue productivity of labor and relative sales in the data (SUSB) and the model (both the raw data and a linear fit). The model is able to match the positive association between these variables well. Figure SM.4.4 shows the relationship



Figure SM.4.2: Model Fit: Firm Size

*Notes:* This figure shows moments targeted in the calibration of the model. Panel (a) shows the model fit of the distribution of relative sales. The distribution of relative sales is obtained from the SUSB in 2012. Panel (b) shows the model fit of average employment, relative to 1-year-old firms, by firm age. Average firm employment by age group was obtained from BDS in 2012. In the calibration exercise, we only target the relative size of firms older than 10 years.

between relative SGA expenses and relative sales in the data (Compustat) and the model. Although the "Compustat-equivalent" sample from the model does not generate all of the dispersion in relative sales as in the data, the relationship with relative SGA expenses is well matched in the overlapping range of relative sales. Finally, Figure SM.4.5 plots the average COGS-to-OPEX ratio as a function of a firm's age and size in the data (Compustat) and the model. Here, age is normalized as time since entry into Compustat (which in the model occurs after the 7th year). In the model, the composition of firms' costs exhibits a strong size profile and a weak age profile, as in the data.

**Untargeted Moments** Figure SM.4.6 shows two model-based moments that were not explicitly targeted in the calibration exercise: the average exit rate by age and the average employment growth by age. As Panel A shows, in the model, the average exit rate is decreasing in a firm's age (as in the data; see, e.g., Haltiwanger, Jarmin, and Miranda, 2013). The fact that entrants enter the economy with lower average productivity and no customer base makes young firms more likely to exit when faced with negative productivity shocks due to overhead costs. As firms grow larger, their larger customer base and higher productivity allow them to absorb negative productivity shocks without forcing them to exit. It is important to note that in the firm dynamics literature, the decreasing profile of the average exit rate by age is typically used as a target to indirectly calibrate the size of the overhead cost  $\chi$ . Here, we took a more direct approach by matching the observed cost structure of firms. The fact that the model replicates the profile of exit rates provides additional support to the model of the composition of firms' cost structures. Panel B plots decreasing profiles of average employment and sales growth as a function of firms' age. Both patterns are consistent with the empirical evidence of Haltiwanger, Jarmin, and Miranda (2013). For example, in the data, the



## Figure SM.4.3: Model Fit: Labor Productivity and Sales

*Notes:* This figure plots the relationship between the relative revenue productivity of labor and relative sales. Panels A and B show the relationship obtained from the SUSB data and the model-simulated data, respectively. Each figure includes the best linear fit of the data. The x-axis is in log scale.



Figure SM.4.4: Model Fit: SGA and Sales

*Notes*: This figure plots the relationship between relative spending in SGA and relative sales. Panels A and B show the relationship obtained from the Compustat data and the model-simulated data, respectively. The model data are obtained by simulating the model and restricting the sample to firms that are at least 7 years old and have sales above 19% of the average sales in the simulated economy (which corresponds to the ratio of the 5th percentile of the sales distribution in Compustat to the average sales in SUSB). Each figure includes the local linear kernel best fit of the data.

average net employment growth rate of 1-2 year-old and 7-8 year-old continuing firms is close to 12% and 2.5%, respectively. In the model, the average employment growth rates (computed as in Davis, Haltiwanger, Schuh, et al. (1998)) for 2- and 7-year-old firms are 12.5% and 2.1%, respectively.

## Figure SM.4.5: Steady-state COGS/OPEX by Size and Age



*Notes:* This figure plots the average COGS-to-OPEX ratio as a function of a firm's age and size. Panels A and B show the relationship obtained from the Compustat data and the model-simulated data, respectively. The model data are obtained by simulating the model and restricting the sample to firms that are at least 7 years old and have sales above 19% of the average sales in the simulated economy (which corresponds to the ratio of the 5th percentile of the sales distribution in Compustat to the average sales in SUSB). Age is normalized as years since entry into Compustat, which in the model corresponds to year 7.



Figure SM.4.6: Exit and Growth by Age

*Notes:* Panels (a) and (b) plot model predictions regarding two untargeted moments—the average exit rate by age and the average employment growth by age—respectively.

## SM.4.3. Model Mechanisms

In Section 4.3, we compared our calibrated model with a recalibrated homogeneous customer model. Here, we further highlight the mechanisms at play by providing comparative statics with the homogeneous customer model, which shuts down endogenous customer acquisition (i.e., we set

 $\phi \rightarrow 0$  and  $\delta = 1$ , while keeping the remaining parameter values fixed across models).

**Implications for Concentration** Figure SM.4.7 plots the firm dynamics of two entrants that start with productivities equal to the 50th and 99th percentiles of the productivity distribution of entrants. After these initial draws, productivity follows the AR(1) process without further shocks. Panel (a) shows that in the baseline model (depicted with solid lines), firms frontload their efforts to acquire customers when young, with more productive firms spending more on advertising labor,  $l_{i,s,t}$ . Panel (b) shows that the sales of the more productive firm are larger than those of the less productive firm, and also that sales increase initially due to a growing stock of customers but eventually decline due to the mean-reversion of its productivity.

How do firms grow their sales? Panels (c) and (d) of Figure SM.4.7 plot the dynamics of the number of customers and sales per customer, respectively. Both firms build their customer base gradually over time due to decreasing returns to advertising. However, sales per customer are higher when firms are younger due to mean-reverting productivity and the decreasing marginal product of labor in production. Thus, over time, firms shift their sales strategy from selling more to a few customers to selling less to more customers.

How does *endogenous* customer acquisition affect these dynamics? Relative to the homogeneous customers model (depicted with dashed lines), the more productive firm is able to achieve higher sales by selling less per customer but accumulating more than twice as many customers in the first year. Since more productive firms charge lower prices but higher markups, profits per marginal customer are increasing in productivity, which induces the more productive firms to accumulate customers more rapidly. Given a fixed stock of customers, this can only be possible if firms with lower productivity accumulate fewer customers relative to the homogeneous customers model. Thus, endogenous customer acquisition increases the dispersion across firms of the number of customers and decreases the dispersion of sales per customer. Table SM.4.2 shows that the former effect dominates and the overall concentration of sales increases: Whereas in the homogeneous customers 50% of sales.

**Implications for Market Power** Table SM.4.2 shows that despite higher concentration, the baseline model features a lower aggregate markup: 1.26 as opposed to 1.38 in the homogeneous customers model. To understand the sources of this difference, Figure SM.4.8 shows the histogram of markups and the scatter plot between markups and relative employment (the weights used in the construction of the aggregate markup) across the two models. These figures illustrate two forces. On the one hand, in the model with endogenous customer acquisition, the distribution of markups is more concentrated. High-productivity firms charge lower markups than firms with similar productivity in the homogeneous customers model. On the other hand, in our model, those high-productivity firms account for a larger fraction of total employment.

The following decomposition of the difference in aggregate markups, which follows from Equa-



## Figure SM.4.7: Average Firm Dynamics

*Notes*: This figure plots the firm dynamics of two new firms that start with zero customer base and productivities equal to the 50th and 99th percentiles of the distribution of productivities among entrants. After this initial draw, firms follow the AR(1) productivity process without any further shock. Solid lines correspond to the baseline model. Dashed lines correspond to the homogeneous customers model with an exogenous and identical customer base  $(m_{i,t} = 1/N_t)$ . For the latter, we compute the general equilibrium using the calibrated parameters of the baseline model and imposing the parameter restrictions  $\phi \rightarrow 0$  and  $\delta = 1$ . Panels (a)-(d) plot the evolution of labor devoted to customer acquisition  $(l_{i,s,t})$ , sales  $(p_{i,t}y_{i,t})$ , customer base  $(m_{i,t})$ , and output per customer  $(q_{i,t})$ , respectively.

tion (5.6), quantifies each force:

$$\underbrace{\ln(\mathcal{M}_{t}) - \ln(\mathcal{M}_{t}^{HC})}_{-9.71\%} \approx \underbrace{\int_{i \in N_{t}} (\omega_{i,t} - \omega_{i,t}^{HC}) \ln(\mu_{i,t}) di}_{\Delta \text{ Distribution: 9.00\%}} + \underbrace{\int_{i \in N_{t}^{HC}} \omega_{i,t}^{HC} \left(\ln(\mu_{i,t}) - \ln(\mu_{i,t}^{HC})\right) di}_{\Delta \text{ Market power: -18.71\%}}$$

where the superscript *HC* denotes distributions and allocations in the homogeneous customers model.<sup>40</sup> The first term (denoted " $\Delta$  Distribution") captures the contribution of differences in the distribution of relative employment (cost-based Domar weights) across firms while keeping firms' markups fixed at their level in the baseline model. The second term (denoted " $\Delta$  Market power") captures the contribution of differences in markups across models while keeping the distribution of relative employment fixed at its distribution in the homogeneous customers model. Since more

<sup>&</sup>lt;sup>40</sup>While the equation presents the decomposition based on its approximation for expositional purposes, the numbers we present are computed based on the exact decomposition.





(a) Distribution of Markups

(b) Markups vs. Relative Employment

*Notes:* Panel (a) plots the distribution of production markups in the baseline and homogeneous customers models. Vertical dashed lines show the average cost-weighted production markup in each model. Panel (b) shows the scatter plot of relative employment  $(l_{i,p,t}/(L_{p,t}/N_t) \propto \omega_{i,t})$  and production markups  $\mu_{i,t}$ . The homogeneous customers model refers to the model with an exogenous customer base  $(m_{i,t} = 1/N_t)$ .

productive firms charge higher markups, switching the distribution of relative employment from the homogeneous customers model to the baseline model *increases* the average markup by 9 pp. However, the contribution of lower markups in the baseline model *reduces* the aggregate markup by 18.7 pp, so the aggregate markup is on net smaller by 9.71 pp. To summarize, why does the baseline model feature a much higher concentration but a lower aggregate markup? Because in the baseline model, firms grow through larger customer bases  $m_{i,t}$  rather than higher average sales per customer  $p_{i,t}q_{i,t}$ , which reduces their market power (but increases their lifetime profits).

**Aggregate Implications** Beyond market power, endogenous customer acquisition also affects other aggregate outcomes by concentrating consumers among higher-productivity firms. First, relative to the homogeneous customers model, it allocates more of the economy's resources toward more productive firms, which reflects itself in the higher aggregate TFP derived in Equation (5.3). Table SM.4.2 shows that shutting down endogenous customer acquisition in our model reduces aggregate TFP by 28%. This is because in the baseline model, more productive firms employ a higher share of production resources to meet the demand from their larger customer bases.

Second, the concentration of customers among more productive firms leaves fewer customers for firms at the bottom of the productivity distribution and brings them closer to the exit threshold. Therefore, when we restrict customer bases to be equal across firms, the equilibrium number of firms increases by 65%. This additional inflow of firms comes from less productive firms that can now generate positive discounted profits due to a larger (exogenous) customer base.

The higher number of firms in the homogeneous customers model also leads to higher demand

for employment. Table SM.4.2 shows that in the homogeneous customers model, total employment is 8 percent larger than in the baseline model (despite the fact that there is no spending on customer acquisition in the former). Part of this difference stems from the larger number of firms, which requires higher overhead costs. Table SM.4.2 also shows that aggregate *production* labor is 6.3% larger in the homogeneous customers model. This is the result of income effects that increase labor supply due to lower aggregate consumption. Finally, the changes in aggregate TFP and production labor lead to an overall 24% decline in aggregate output.

	Baseline	Homog. Customers
	Model	Model
TFP		-27.9
Output		-23.9
Employment		7.9
Production		6.3
Number of firms		65.1
Agg. markup	1.26	1.38
Top 5% sales share	0.50	0.17

Table SM.4.2: Aggregate Effects of Customer Acquisition

*Notes:* The table reports equilibrium aggregates in the baseline and homogeneous customers versions of the model. The homogeneous customers model refers to the model with an exogenous customer base ( $m_{i,t} = 1/N_t$ ). The second column reports percentage differences with respect to aggregates in the baseline model, with the exception of the aggregate markup and the top 5% sales share, which are reported in levels.

## SM.4.4. Additional Analysis of Model in Steady State

This section further describes how the model works in the steady state. First, we describe firms' optimal policies. Then, we show the average firm dynamics, taking selection into account.

**Firms' Optimal Policies** Figure SM.4.9 shows firms' steady-state optimal policy functions for three productivity levels (the 25th, 50th, and 75th percentiles of the marginal productivity distribution in steady state). The y-axis on the right plots the marginal distribution of the relative customer base.

While optimal spending in  $l_{i,s,t}$  decreases with the size of the customer base, production labor is increasing in a firm's customer base. Thus, when firms have a small customer base, they spend more resources to increase it. However, due to decreasing returns to customer accumulation, firms increase their customer base gradually over time. As firms grow, they spend less on customer acquisition and more on producing goods to satisfy the growing demand. This is reflected in a firm's cost structure: The average COGS-to-OPEX ratio is also increasing in  $m_{-1}$ . As total output increases due to a larger customer base, the marginal cost of production also increases since production is subject to decreasing returns. This raises the price charged by the firm, which in turn reduces the consumption per capita  $q_{i,t}$  and optimal markups. The figure also shows that for a given level of  $m_{-1}$ , spending in  $l_{i,s,t}$  is increasing in a firm's productivity. A higher productivity allows firms to charge lower prices and higher markups. Thus, profits per marginal customer are increasing in productivity, which incentivizes firms to accumulate customers more quickly by spending more on  $l_{i,s,t}$ .



Figure SM.4.9: Firms' Optimal Policies

*Notes:* These figures plot firms' policy functions in steady state. Each figure shows policies as a function of relative customer base for three levels of productivity: the 25th, 50th and 75th percentiles of the stationary productivity distribution. The y-axis on the right plots the stationary marginal distribution of the relative customer base.

Figure SM.4.10 plots firms' optimal exit policies and the stationary joint distribution of  $(m_{-1}, z)$ . Panel A shows the threshold productivity  $z^*(m_{-1})$  such that if  $z < z^*(m_{-1})$ , the firm optimally chooses to exit. The figure shows that  $z'^*(m_{-1}) < 0$ —that is, firms with larger customer bases are able to survive large productivity shocks without the need to exit the market. Although a lower productivity reduces markups and profits per customer, aggregate profits are increasing in a firm's customer base. Panel B shows that there is a positive correlation between firms' productivities and customer bases in the steady state.



Figure SM.4.10: Optimal Exit and Stationary Distribution

*Notes:* Panel A plots the exit threshold  $z^*(m_{-1})$  such that if  $z < z^*(m_{-1})$ , the firm optimally chooses to exit. Panel B shows the contour plot of the stationary joint distribution of  $(m_1, z)$ , censored at the 99th percentile of each variable.

**Average Firm Dynamics with Shocks** Figure SM.4.11 plots the average firm dynamics, taking selection into account. To construct this figure, we simulate a cohort of firms that starts with a zero customer base and draws productivities from the distribution of entrants. Since firms are subject to productivity shocks, some decide to exit over their lifetime. The figure plots the average of each variable across firms that survived up to a given age.

Firms start with lower productivity, which grows over time due to the calibrated lower productivity of entrants and endogenous exit. Conditional on a productivity level, firms frontload spending on customer acquisition, and the average customer base and marginal costs rise rapidly for young firms. This, in turn, increases prices and reduces average output per customer and markups. Over time, only the most productive firms survive, so the average marginal cost and price decline, and output per customer and markups increase.



Figure SM.4.11: Average Firm Dynamics

*Notes:* The figure plots the average firm dynamics, which are obtained by simulating a cohort of firms that start with  $m_{-1} = 0$ , draw *z* from the distribution of entrants, and experience productivity shocks over their lifetime. Each figure plots the average of a variable as a function of firms' age across firms that survived up to that age.



Figure SM.4.12: Allocation of Customers: Equilibrium vs. Efficient Allocation

*Notes:* This figure shows a scatter plot between relative productivity  $z_{i,t}/\bar{z}$  and relative customer bases  $m_{i,t}/\bar{m}$ , for both the equilibrium and the social planner's allocation. We present three plots by varying the value of  $\phi$ , while keeping the remaining parameters fixed at the values in the baseline calibration. Low  $\phi$  corresponds to 0.25, baseline to 0.53, and high to 0.75.

# SM.5 Calibration of the Homogeneous Customers Model

Parameter	Description	Value				
Panel A: Fixed Parameters						
β	Annual discount factor	0.960				
γ	Elast. of intertemporal substitution	2.000				
$\psi$	Frisch elasticity	1.000				
α	Decreasing returns to scale	0.640				
δ	Prob. of losing customer	1.000				
	Panel B: Calibrated Parameters					
X	Overhead cost	0.704				
σ	Avg. elasticity of substitution	5.998				
η	Superelasticity	0.559				
ν	Exog. survival probability	0.960				
$ ho_z$	Persistence of productivity shock	0.981				
$\sigma_z$	SD of productivity shock	0.259				
$\bar{z}_{ent}$	Mean productivity of entrants	-2.379				
λ	Mass of entrants	0.123				
ξ	Disutility of labor supply	2.246				

Table SM.5.1: Model Parameters: Homogeneous Customers Model

*Notes:* This table shows the calibration of the homogeneous customers model with an exogenous customer base. Panel A contains parameters externally chosen. Panel B contains parameters internally calibrated to match moments presented in Table SM.5.2 and Figure SM.5.1.

Moment	Data	Model
Avg. COGS-to-OPEX ratio	0.660	0.667
Avg. cost-weighted production markup	1.250	1.266
Slope labor prod. on sales	0.036	0.035
Avg. exit rate	0.073	0.073
SD. employment growth	0.416	0.436

Table SM.5.2: Targeted Moments: Homogeneous Customers Model

*Notes:* This table shows the set of moments targeted in the calibration of the homogeneous customers model with an exogenous and homogeneous customer base. Avg. COGS-to-OPEX ratio refers to the average of the ratio across firms. Avg. cost-weighted production markup corresponds to the COGS-weighted average markup from Edmond, Midrigan, and Xu (2022). These moments were computed using data from Compustat in 2012. Slope of labor prod. on sales corresponds to the OLS coefficient of the sales-weighted regression of relative revenue labor productivity on relative sales from Edmond, Midrigan, and Xu (2022), restricting the sample to firms with relative sales above one. This moment was computed using data from the SUSB in 2012. The average exit rate was obtained from the BDS in 2012. The standard deviation of annual employment growth for continuing establishments is obtained from Elsby and Michaels (2013). The growth rate of variable *x* is computed as in Davis and Haltiwanger (1992):  $(x_{i,t} - x_{i,t-1})/0.5(x_{i,t} + x_{i,t-1})$ . The last column shows the model counterparts of each moment, which were obtained by simulating a panel of firms and computing each moment with the simulated data. In the model, we account for selection into Compustat by restricting the simulated sample of firms to those that are at least 7 years old and have sales above 19% of the average sales in the simulated economy (which corresponds to the ratio of the 5th percentile of the sales distribution in Compustat to the average sales in SUSB).



Figure SM.5.1: Homogeneous Customers Model Fit: Firm Size

(a) Distribution of Sales

(b) Avg. Size by Age

*Notes:* This figure shows moments targeted in the calibration of the homogeneous customers model with an exogenous and homogeneous customer base. Panel (a) shows the model fit of the distribution of relative sales. The distribution of relative sales is obtained from the SUSB in 2012. Panel (b) shows the model fit of average employment, relative to 1-year-old firms, by firm age. Average firm employment by age group was obtained from BDS in 2012. In the calibration exercise, we only target the relative size of firms older than 10 years.