

Productivity and Quality of Multi-product Firms*

Mauro Caselli[†] University of Trento
Arpita Chatterjee[‡] Federal Reserve Board & UNSW
Shengyu Li[§] University of New South Wales

August 25, 2024

Abstract

This paper proposes a novel method to estimate productivity and quality of multi-product firms at the firm-product level. The method utilizes firm optimization conditions to establish a one-to-one mapping between observed data and unobserved productivity and quality. It eliminates the need for imputing firm-product input shares or imposing productivity evolution processes. Our method is scalable to accommodate numerous products and can address the bias caused by unobserved heterogeneous intermediate input prices. We apply this method to a set of Mexican manufacturing industries. Multi-product firms' top-performing products exhibit both higher productivity and quality, but they face a trade-off between the two, referred to as the cost of quality. In a counterfactual exercise, we show that a reduction in the cost of quality can lead to substantial firm-level productivity gains and that, on average, about 30.1 percent of these gains are due to the within-firm reallocation of production.

Keywords: *multi-product firms, production function, productivity, output quality, intra-firm reallocation*

JEL classification: *D24, L11, L15, O47.*

*The authors thank Eleni Aristodemou, Zhiyuan Chen, Jan De Loecker, Erwin Diewert, Kevin Fox, Andrea Fracasso, Moyu Liao, Matthias Mertens, Scott Orr, Ariell Reshef, Mark Roberts, Stefano Schiavo, Petr Sedlacek, Nikos Theodoropoulos, Andreas Tryphonides, Nelli Valmari, Eric Verhoogen, Frederic Warzynski, Daniel Xu, Hongsong Zhang, and many seminar and conference participants for very helpful comments. All errors are the authors' responsibility. All views/opinions are authors' own and do not reflect views of the Federal Reserve Board.

[†]School of International Studies & Department of Economics and Management, University of Trento. Email: mauro.caselli@unitn.it.

[‡]Federal Reserve Board & School of Economics, UNSW. Email: chatterjee.econ@gmail.com.

[§]Corresponding author: School of Economics & Centre for Applied Economic Research, Business School, the University of New South Wales, Australia. Email: shengyu.li@unsw.edu.au.

1 Introduction

The production landscape of many manufacturing industries is dominated by multi-product firms, which operate across a diverse range of product lines. However, existing empirical studies that explore the determinants of firm performance have primarily focused on analyzing variations across different firms, such as heterogeneity in productivity levels and demand characteristics (e.g., [Foster et al., 2008](#); [Pozzi and Schivardi, 2016](#); [Kumar and Zhang, 2019](#)). Consequently, there remains a considerable gap in the understanding of the factors that drive within-firm heterogeneity and resource reallocation, as well as their subsequent impact on firm performance and growth. This knowledge gap is mainly due to methodological limitations and data constraints, which hinder the accurate estimation of different aspects of heterogeneity at the firm-product level.

This paper introduces an innovative method to estimate productivity and quality (product appeal) at the firm-product level, along with the transformation function and demand parameters. This method constructs a unique one-to-one mapping from observed data to unobservable variables by leveraging firm optimization conditions. This provides distinct advantages compared to the recent methods (e.g., [Dhyne et al., 2022](#); [Orr, 2022](#); [Valmari, 2022](#)). First, it eliminates the need for imputing intra-firm input allocations. Second, it does not impose restrictions on productivity evolution, allowing for flexibility in exploring complex productivity dynamics after estimation. Third, it is scalable to handle a large number of products. Fourth, it addresses the estimation bias caused by heterogeneous firm-level intermediate input prices, which are usually unobservable in commonly available data sets.

Drawing on comprehensive firm-product-level data from three major industries in the Mexican manufacturing sector, we employ this method to study the trade-off between productivity and quality within firms and the role of product scope in shaping firm growth through intra-firm resource reallocation. While we implement the approach for a setup with a specific functional form, our methodology is readily extendable to more general settings and broader applications.

In modelling the production side, our method is designed to address the challenges commonly faced in estimating multi-product production functions. The recent strand of production function estimation methodologies implicitly assumes that each firm produces a single product (e.g., [Olley and Pakes, 1996](#); [Levinsohn and Petrin, 2003](#); [Akerberg et al., 2015](#); [Gandhi et al., 2020](#)). In this context, the input allocation is observable to researchers and each firm only has a single dimension of unobservable productivity, which can be controlled for by an observable proxy. Multi-product firms, on the contrary, may produce different products and thus have different levels of productivity in these products, even within the same firm.

Extending the proxy-based methods to the context of multi-product firms requires at least the same number of proxies as the number of products (cf., [Dhyne et al., 2022](#)). Moreover, researchers do not observe the within-firm division of inputs used to produce different products because firms usually only report total inputs. The potentially heterogeneous ability of input sharing (e.g., machinery and workers) across product lines within firms (e.g., [Cairncross et al., 2023](#)), which is observable to the firms but is unobservable to researchers, further complicates the problem.¹ Finally, intermediate input prices, which significantly vary across firms and over time due to various reasons such as quality differentiation, bargaining power in the input market, transport costs, and suppliers’ marginal cost as documented by [Atalay \(2014\)](#) using US Census Bureau data, should be controlled for in the estimation to avoid biased estimates of production function and input elasticities (i.e., input price bias as emphasized in [Ornaghi, 2006](#); [De Loecker et al., 2016](#); [Grieco et al., 2016](#)). However, only input expenditure (rather than input price and input quantity) is observable to researchers at the firm level in commonly available data sets.

To address these issues, we model the production technology using a constant elasticity of substitution (CES) transformation function, which transforms inputs into different products. The inputs can be shared in production across products within the same firm. Each product is associated with a potentially different level of physical productivity (i.e., quantity-based productivity, as in [Foster et al., 2008](#)).² The firm observes these productivity levels before making input and output decisions to maximize profits. In the spirit of [Grieco et al. \(2016\)](#), we show that the optimization conditions implied from our model can always be inverted to form an explicit one-to-one mapping from observed input and output decisions to unobserved productivity at the firm-product level (regardless of the number of products), while controlling for unobserved intermediate input prices. We use the inverted relationship to substitute unobserved productivity to estimate the parameters of the transformation function. Once the parameters are estimated, we compute productivity at the firm-product level from the one-to-one mapping.

In modelling the demand side, we adopt a commonly used CES demand function.³ The firm’s products are chosen from a set of horizontally differentiated categories. Within each product category, each firm’s product variety is vertically differentiated according to its quality level. Because the optimal product prices are chosen after the firm’s decisions on

¹For example, a printing firm may use the same design software to create multiple products, such as logos and product labels; workers with specialized skills, such as pattern makers and shoe designers, may be used across different product lines within the same footwear firm; in pharmaceutical industries, a firm may use the same reactors and mixing tanks to produce different products, by adjusting the process parameters.

²We refer to physical productivity as simply “productivity” in this paper unless explicitly stated otherwise.

³Our method is extendable to a general demand function which explicitly allows for the complementarity or substitutability of products, as described in Online Appendix [A](#).

the product quality levels, we face a classic endogeneity problem in estimating the price elasticity of demand. The traditional solution is to use cost shifters (such as capital stock) as instrumental variables (IVs) for the price to estimate the demand function directly. However, if the cost of producing high-quality products is higher as suggested by the recent literature (e.g., [Grieco and McDevitt, 2017](#); [Forlani et al., 2023](#); [Li et al., 2023](#)), then cost shifters may still be correlated with quality. Thus, we depart from the existing literature to examine the advantage offered by intra-firm decisions in multi-product firms – profit maximization of the firm implies a relationship between the revenues of products within the same firm. This relationship depends on the demand elasticities and intra-firm relative quality *differences* (as opposed to quality *levels*), which can be instrumented by the commonly used firm-level cost shifters. Therefore, we use this relationship to help identify the demand elasticities. We use Monte Carlo exercises to demonstrate that this approach is able to recover the true parameters well. After the estimation, we compute quality as the residual of demand after controlling for price in the spirit of [Khandelwal \(2010\)](#).⁴

We apply our method to establishment-level panel production data from three major Mexican manufacturing industries (i.e., footwear, printing, and pharmaceuticals) that record prices and quantities at the firm-product level along with rich input data at the firm level. Multi-product production is an essential feature of the firms in our sample. Multi-product firms account for around 65% of the total number of firms and 86% of total revenues, and their average number of products is 6.7 per year, albeit with differences between industries. Within each industry, the markets for different products (e.g., women’s shoes vs. men’s shoes in the footwear industry) are largely segmented. Nevertheless, for each product, firms’ output is vertically differentiated, as evidenced by the large dispersion in prices. These features are consistent with the model’s assumption of a monopolistically competitive market structure with vertically differentiated products.

After estimation, we first follow the literature (e.g., [Melitz, 2000](#)) to construct a (firm-product) quality-adjusted productivity (ATFP) measure that accounts for heterogeneity in both productivity and quality. We find significant dispersion of ATFP across firms, even conditional on the product. More importantly, both components of ATFP (i.e., productivity and quality), are important for the within-firm performance of multi-product firms. Products closer to the core competence of the firm (defined by the highest revenue within firms) have both higher productivity and higher quality.

⁴The residual of demand is essentially demand heterogeneity which embodies a set of demand shifters. We leverage the rich fixed effects offered by the firm-product level data to refine the demand residual as a measure of quality in the empirical exercises. Nonetheless, we acknowledge that the refined measure of quality may still have different components, such as product appeal perceived by consumers, if they vary at the firm-product-time level.

These different dimensions of within-firm heterogeneity are not, however, unrelated to each other. Within a firm, improving quality at the product level comes at the cost of reducing productivity. This result is broadly consistent with the emerging literature highlighting the trade-off across firms between these two unique dimensions of firm heterogeneity (e.g., [Jaumandreu and Yin, 2014](#); [Grieco and McDevitt, 2017](#); [Roberts et al., 2018](#); [Atkin et al., 2019](#); [Orr, 2022](#); [Eslava et al., 2023](#); [Forlani et al., 2023](#); [Li et al., 2023](#)).⁵ In the industries we consider, a 1 percent increase in quality reduces productivity by 0.234 percent on average, holding all other variables constant. Moreover, this trade-off is heterogeneous – while it is more costly to produce a high-quality product, long experience in producing a particular product allows the firm to improve quality with less sacrifice in efficiency.

Quantitatively, the cost of quality bears significant implications for firm productivity growth and intra-firm resource allocation. A reduction in the cost of quality not only directly increases the firm’s ATFP but also indirectly influences it through the firm’s endogenous reallocation of resources towards the production of higher-quality products. This is due to the positive relationship observed between ATFP and product quality within the firm. In a counterfactual analysis, we find that a 1 percent reduction in the cost of quality corresponds to an average 2.836 percent improvement in firm-level ATFP. Notably, a substantial 30.1 percent of this improvement can be attributed to the within-firm reallocation of production towards high-quality, high-ATFP products.

We show that the impact of the quality cost reduction on firm performance is particularly pronounced for multi-product firms with larger product scope. Their ability to leverage a larger range of products through reallocation allows them to capitalize on the opportunities arising from reduced quality cost, thus boosting their overall productivity. This result uncovers a novel mechanism for productivity growth for multi-product firms, which dominate manufacturing production.

In terms of methodology, our paper builds on recent advances in the estimation of heterogeneous productivity of multi-product firms. In addressing the common data challenge of input data being observable only at the firm level, while outputs and revenues are reported separately by product, the literature has evolved into two main approaches. The first approach, pioneered by [De Loecker et al. \(2016\)](#), characterizes multi-product production as a collection of single-product production functions, coupled with a rule for allocating firm inputs to each of these functions. Subsequent studies have extended this approach. In particular, [Orr \(2022\)](#) models product lines sharing the same technology (i.e., production

⁵Intuitively, producing one unit of a high-quality product may require more (or longer) production processes, better (or more specialized, exclusive) machinery, higher quality (or more) intermediate materials, and higher standards of quality control, all of which lead to a lower quantity of output holding inputs fixed and consequently an increase in *marginal cost* of production (or lower productivity, equivalently).

parameters) but with individual efficiency shocks, and shows how demand data can be used to assist estimation under profit maximization conditions. Valmari (2022) develops a similar framework, incorporating flexible production parameters across product-specific production functions. Chen and Liao (2022) generalize the previous papers by allowing single-product firms and multi-product firms to have different production functions and by estimating both non-parametric and parametric production functions for multi-product firms. In contrast, the second approach, led by Dhyne et al. (2022), departs from the assumption that multi-product production is a collection of single-product firms. They introduce a transformation production function and show how it can be used to recover the production frontier and estimate firm-product-specific marginal costs, taking into account complementarities and spillovers in multi-product production.

Our methodology integrates the strengths of both approaches to overcome their respective limitations. First, we model multi-product production using a transformation production function, similar to Dhyne et al. (2022). This avoids the need to allocate firm-level inputs, as in Orr (2022) and Valmari (2022), and allows for intra-firm input sharing across product lines, which may contribute to economies of scope in multi-product production. Second, in addressing unobserved firm-product productivity, we adopt the profit maximization assumption, similar to Orr (2022) and Valmari (2022). However, instead of imputing input allocation shares, we use the profit-maximizing conditions to establish a one-to-one mapping from observed firm decisions to unobserved productivity, extending the insights of Grieco et al. (2016, 2022), Harrigan et al. (2021) and Li and Zhang (2022) to the context of multi-product firms. Importantly, the number of profit-maximizing conditions, which naturally increase with the number of products, ensures the scalability of our method. This differs from Dhyne et al. (2022), whose method requires a separate proxy for each additional firm-product-level productivity. Rather, it is more similar to recent approaches to identify markdowns (Morlacco, 2020; Caselli et al., 2021; Kirov and Traina, 2023) or factor-augmenting productivity (Demirer, 2022; Raval, 2023) using necessary conditions for optimality with respect to multiple flexible inputs. Third, our method addresses the bias due to unobserved firm-level heterogeneity in input prices without requiring the availability of input price data. This is in contrast to the existing methods (e.g., Orr, 2022; Valmari, 2022), which typically require access to such data. Finally, our method does not rely on modelling the evolution of productivity, which offers a distinct advantage in exploring the evolution of productivity after estimation. Such an advantage is particularly beneficial in studying complex (e.g., interdependent) productivity dynamics, factors that endogenously shape the productivity trajectory (e.g., Chen et al., 2021), and frequent product turnover decisions, such as for exported products, where the observation of products is truncated by latent variables.

Our empirical application contributes to the emerging literature analyzing the trade-off between productivity and quality (i.e., the cost of quality). Focusing on the healthcare industry, [Grieco and McDevitt \(2017\)](#) show that reducing the quality standards of a healthcare center can increase its patient load. [Atkin et al. \(2019\)](#) reveal a reverse correlation between quantity and quality among rug-makers in Egypt, drawing insights from data that include direct quality assessments. [Forlani et al. \(2023\)](#) document a strong negative correlation between demand and quantity-based productivity in various Belgian manufacturing industries. Using an objective measure of output quality, [Li et al. \(2023\)](#) find that about half of the benefits of quality are offset by the cost of producing the quality in the Chinese steel industry. These papers document such a trade-off across firms. Our paper finds a similar trade-off at the firm-product level and shows that the trade-off diminishes as firms gain experience in manufacturing. To this end, our analysis is consistent with the negative relationship between productivity and “product appeal” documented at the same level of disaggregation by [Orr \(2022\)](#). Nonetheless, after taking both the cost and the benefit of quality into account, ATFP is documented to be positively correlated with quality. This result is consistent with endogenous quality choice models ([Kugler and Verhoogen, 2009, 2012](#)) and empirical analysis using an objective quality measure by [Li et al. \(2023\)](#).

Finally, our paper is related to a large literature on resource reallocation, which focuses on across-firm analyses and shows that much of the aggregate productivity growth is attributable to the resource reallocation towards more productive firms (e.g., [Baily et al., 1992](#); [Bartelsman and Doms, 2000](#); [Aw et al., 2001](#); [Foster et al., 2008](#); [Syverson, 2011](#); [Collard-Wexler and De Loecker, 2015](#)). Our counterfactual analysis shows that there can be a substantial contribution to productivity growth due to *within-firm* resource reallocation – a mechanism that is emphasized in the recent literature studying multi-product firms (e.g., [Mayer et al., 2021](#)). We focus on the channel of the cost of quality and document a positive relationship between product scope and the contribution of intra-firm resource reallocation. This result illustrates the importance of intra-firm resource reallocation within multi-product firms due to quality differences as a novel channel affecting overall productivity at the firm level.

In the rest of the paper, [Section 2](#) describes our model with the production and demand functions and the firm’s endogenous decisions on input and output. In [Section 3](#) we develop the estimation methodology. [Section 4](#) describes the data used in the estimation. [Section 5](#) presents the estimation results. [Section 6](#) illustrates the trade-off between productivity and quality, while [Section 7](#) quantitatively assesses the significance of the cost of quality in intra-firm resource reallocation. We conclude in [Section 8](#).

2 Model

This section develops a framework for describing a firm’s static input and output decisions, which underpins our empirical estimation by leveraging the optimization conditions implied by this model.⁶ Although we present a specific functional form in the main text, our identification strategy is not dependent on this assumption. Our methodology can be applied to more general settings, explicitly accounting for product complementarity or substitution in demand and production, as demonstrated in Online Appendix A.

Consider an industry with J firms indexed by $j = 1, 2, \dots, J$. There is a total of N products, indexed by $n = 1, 2, \dots, N$, that firms can choose to produce. The timeline of the decisions is as follows. At the beginning of period t , the set of products that firm j has decided (at the end of the previous period) to produce in this period is Λ_{jt} . Each product $n \in \Lambda_{jt}$ is associated with a level of technical efficiency ω_{jnt} and a level of quality ξ_{jnt} , both of which have been determined and observed by the firm at the end of the previous period. The firm’s capital stock is also determined in the previous period via an investment decision.

The firm’s static decisions in the current period consist of the material input, labor input, and quantities of individual products to maximize its total period profit subject to demand and production functions, after observing the material price, wage rate, and capital stock. At the end of period t , the firm makes dynamic decisions on the capital stock and the set of products to be produced with their levels of product quality and technical efficiency for the following period, after observing the associated adjustment or investment costs.

2.1 Demand

The entire set of products that the firm can choose to produce is divided into N horizontal categories, such as women’s and men’s shoes. For each product category $n \in \{1, 2, \dots, N\}$, the output of each firm is vertically differentiated according to its choice of quality level Ξ_{jnt} . This means that, although the demand for each of the N product categories is segmented, there is monopolistic competition across firms that produce vertically differentiated products in the same category. This assumption is also adopted by De Loecker (2011) and Valmari (2022) in modelling the demand functions in the multi-product context.⁷

⁶The firm’s static choices are made conditional on a set of dynamic decisions, including output quality, technical efficiency, product scope, and investment. Online Appendix B outlines the firm’s dynamic decisions related to these choices, offering conceptual insights into how these choices are endogenously determined. While we do not estimate the complex dynamic model due to the high dimensionality of the state variables, it serves to clarify the firm’s dynamic decision-making process.

⁷While we extend our model to allow for a flexible, general demand system to account for potential complementarity and substitution on the demand side in Online Appendix A, we maintain a simpler demand function assumption in the main text of the paper.

Specifically, for each product category n , a representative consumer has constant elasticity of substitution (CES) preferences in terms of both the quality and the quantity of the products offered by firms:⁸

$$U_{nt} = \left[\sum_j (\Xi_{jnt}^{\frac{1}{\eta_n-1}} Q_{jnt})^{\frac{\eta_n-1}{\eta_n}} \right]^{\frac{\eta_n}{\eta_n-1}}, \quad (1)$$

where $\eta_n > 1$ is the elasticity of substitution across the varieties offered by the firms. Q_{jnt} is the physical quantity and Ξ_{jnt} is the product quality of firm j in period t , respectively. That is, the consumer values the quality-adjusted quantity of the product, $\Xi_{jnt}^{\frac{1}{\eta_n-1}} Q_{jnt}$, which forms the basis for constructing the quality-adjusted productivity in Section 5.

Given the consumer's total expenditure B_{nt} and the product price P_{jnt} , the consumer's utility maximization problem implies the following demand function for product n from firm j :

$$\ln Q_{jnt} = -\eta_n \ln P_{jnt} + \xi_{jnt} + \phi_{nt} + \psi_{jn} + v_{jt}, \quad (2)$$

where $\xi_{jnt} = \ln \Xi_{jnt}$. Intuitively, a higher quality level shifts the demand curve upwards. Beyond quality, three other components also influence demand. First, $\phi_{nt} = \ln \left(\frac{B_{nt}}{\sum_j \Xi_{jnt} P_{jnt}^{1-\eta_n}} \right)$ is a product-specific expenditure index that depends on macroeconomic conditions captured in B_{nt} such as consumer income and market size in period t . Second, ψ_{jn} represents factors that affect demand at the firm-product level but do not vary over time such as consumers' subjective tastes, brand image related to specific products, number (or variety) of subcategories contained in each product category under our classification and product measurement units (e.g., grams vs. liters).⁹ Finally, v_{jt} captures the demand heterogeneity such as firm effort in marketing that varies by firm and year. For demonstration, we summarize the structural terms that shift the demand function as $\tilde{\xi}_{jnt} = \xi_{jnt} + \phi_{nt} + \psi_{jn} + v_{jt}$. The firm observes $\tilde{\xi}_{jnt}$ for all products before making production decisions.

Remark: Essentially, $\tilde{\xi}_{jnt}$ is a demand shifter, which captures all sorts of demand heterogeneity that influences product demand but is not accounted for by product prices. Empirically, $\tilde{\xi}_{jnt}$ is usually referred to as “perceived product appeal/demand” (e.g., Pozzi and Schivardi, 2016; Orr, 2022; Valmari, 2022; Eslava et al., 2023) or “quality” (e.g., Melitz, 2000; Khandelwal, 2010; Hottman et al., 2016). In our paper, we follow this tradition and

⁸The power of Ξ_{jnt} , $\frac{1}{\eta_n-1}$, is used to simplify the notation to reach a commonly used demand function (2). A large literature that treats demand residual as output quality implicitly shares the same setup (e.g., Melitz, 2000; Khandelwal, 2010; Pozzi and Schivardi, 2016; Valmari, 2022).

⁹Units of measurement can be different across product categories. Consequently, the quantities and prices of different product categories are not readily comparable. In the demand function (2), ψ_{jn} absorb such differences. Similarly, in our empirical analysis in Section 6, we use firm dummies and product dummies to tease out ξ_{jnt} from such differences.

acknowledge that it embodies quality (ξ_{jnt}) as well as non-quality components, such as consumer tastes, brand/firm image, marketing efforts and market size. Yet, our setting with multiple-product firms provides us with a rich set of fixed effects at the product-year (ϕ_{nt}), firm-product (ψ_{jn}), and firm-year (v_{jt}) levels to control for the non-quality component that varies at these levels. For this reason, we define $\chi_{jnt} = \phi_{nt} + \psi_{jn} + v_{jt}$ and refer to χ_{jnt} as a demand shock in this paper. Notably, this advantage is not available in the analysis using firm-level data, and thus it helps to tease out a finer measure of quality (i.e., ξ_{jnt}) from residual demand (i.e., product appeal, $\tilde{\xi}_{jnt} = \xi_{jnt} + \chi_{jnt}$) that is traditionally used as quality.

2.2 Production Technology

We use a transformation function to model the production technology. Specifically, given the set of products to be produced (Λ_{jt}) and associated product appeal ($\tilde{\xi}_{jnt}$, $n \in \Lambda_{jt}$), the firm uses labor (L_{jt}), material (M_{jt}), and capital (K_{jt}) to produce output quantity (Q_{jnt} , $n \in \Lambda_{jt}$) following a constant elasticity of substitution (CES) transformation function:

$$\sum_{n \in \Lambda_{jt}} e^{-\tilde{\omega}_{jnt}} Q_{jnt} = F(L_{jt}, M_{jt}, K_{jt}) \equiv [\alpha_L L_{jt}^\gamma + \alpha_M M_{jt}^\gamma + \alpha_K K_{jt}^\gamma]^{\frac{\rho}{\gamma}}, \quad (3)$$

where $\tilde{\omega}_{jnt}$ is the Hicks neutral, quantity-based productivity (i.e., so-called physical productivity, or TFPQ) of firm j in producing product n in period t . In this paper, we use quantity-based productivity, TFPQ, and productivity interchangeably. $\gamma \equiv \frac{\sigma-1}{\sigma}$ governs the elasticity of substitution across inputs, i.e., labor, material, and capital. ρ is a parameter for the returns to scale in the transformation of inputs into output. α_L , α_M , and α_K are distribution parameters associated with labor, material, and capital, respectively. We normalize their sum to 1.

Remark: A few features of the transformation function are worth noticing. First, the transformation function is compatible with the single-product CES production functions traditionally used in the literature. In the context of multi-product firms, a similar transformation function is adopted by [Cairncross et al. \(2023\)](#), who derive the transformation function (as a general output distance function) from individual product production functions with shared inputs across products. In fact, the transformation function (3) in our setup is a restricted version of the output distance function proposed by [Cairncross et al. \(2023\)](#), which

includes a CES aggregator on the output side with a parameter θ :¹⁰

$$\left[\sum_{n \in \Lambda_{jt}} e^{-\tilde{\omega}_{jnt}} Q_{jnt}^\theta \right]^{\frac{1}{\theta}} = F(L_{jt}, M_{jt}, K_{jt}) \equiv [\alpha_L L_{jt}^\gamma + \alpha_M M_{jt}^\gamma + \alpha_K K_{jt}^\gamma]^{\frac{\rho}{\gamma}}. \quad (4)$$

Parameter θ governs the degree of substitution or complementarity across outputs. Our restriction in the implementation is $\theta = 1$. This effectively assumes that outputs are perfectly substitutable and the rate of substitution between any two products is determined by their relative productivity, which flexibly varies at the firm-product-time level to absorb any potential substitution or complementarity on the production side. Nonetheless, our model is extendable to allow for a flexible θ to be estimated and account for substitution or complementarity as described in Online Appendix A.

Second, although we assume the production parameters of the transformation function are the same across all firms, regardless of the combinations of products they produce, our model is readily extendable to accommodate potential differences in technologies for producing different sets of products. That is, a more flexible specification of the transformation function:

$$\sum_{n \in \Lambda(o)} e^{-\tilde{\omega}_{jnt}} Q_{jnt} = F_o(L_{jt}, M_{jt}, K_{jt}) \equiv [\alpha_{L_o} L_{jt}^{\gamma_o} + \alpha_{M_o} M_{jt}^{\gamma_o} + \alpha_{K_o} K_{jt}^{\gamma_o}]^{\frac{\rho_o}{\gamma_o}}, \quad \forall o \in \mathcal{O}, \quad (5)$$

where \mathcal{O} is the set of all possible permutations of products (including single products) that are observed in the data and $o \in \mathcal{O}$ is one of the permutations. In such a flexible setting, the production parameters, representing input intensity and substitutability, vary by the permutation o . This feature is similar to (and potentially more general than) that of Valmari (2022) who allows the production parameters to vary by product. The estimation strategy described in Section 3 can be directly applied to each permutation $o \in \mathcal{O}$, provided there are sufficient observations of firm-year pairs in each permutation.¹¹ This flexibility of our model is due to the distinct feature of our estimation strategy, which does not rely on time series variation, unlike traditional proxy-based approaches. In contrast, extending existing methods that rely on the evolution of productivity (e.g., Orr, 2022; Dhyne et al., 2022) to such a setting is more challenging, especially when the set of products chosen by firms varies frequently over time.

¹⁰In this restricted setup with $\theta = 1$, parameter $\rho > 1$ implies that there is input sharing in production, as shown Cairncross et al. (2023).

¹¹For example, suppose there are two permutations of products in the data: firm-year observations either produce products 1 and 2 or produce products 2 and 3. We can estimate two transformation functions for each type of firm-year observations, allowing the production parameters to differ. However, in our empirical exercise, the lack of sufficient observations in each permutation refrains us from implementing such a flexible model.

Third, for multi-product firms, the transformation function can be interpreted as the frontier of joint production of all products, Q_{jnt} , $n \in \Lambda_{jt}$. This interpretation has three implications: (i) different products are manufactured with the same set of inputs (i.e., labor, material, and capital); (ii) the inputs can be costlessly transferred across different products within the firm; (iii) producing more of one product means producing less of another product, holding inputs fixed. These implications are consistent with the modelling assumptions used by [Dhyne et al. \(2022\)](#), [Orr \(2022\)](#), and [Valmari \(2022\)](#).

Finally, our transformation function allows for shared inputs or the joint utilization of inputs across different products, which may contribute to economies of scope in the spirit of [Panzar and Willig \(1977, 1981\)](#). Input allocation within a firm is not explicitly modelled in our framework. This methodology is in contrast to the existing methods that rely on imputing the intra-firm (exclusive) allocation of inputs and thus abstract away from imperfectly divisible inputs with properties of a public good within a firm.

2.3 Productivity

A key element of our model is the quantity-based productivity $\tilde{\omega}_{jnt}$ in (3), which varies by firm, product, and period. While we do not impose restrictions on $\tilde{\omega}_{jnt}$ to estimate the parameters in (3), in this subsection we discuss the potential components and evolution of $\tilde{\omega}_{jnt}$ to highlight the key differences compared with the assumptions in the existing literature.

Departing from the literature, we unpack productivity into two components:

$$\tilde{\omega}_{jnt} = \omega_{jnt} - h(\xi_{jnt}), \tag{6}$$

where ω_{jnt} is technical efficiency and $h(\xi_{jnt})$ is a function of product quality ξ_{jnt} . It is crucial to model $h(\xi_{jnt})$ as a part of quantity-based productivity because varieties of the same product category produced by different firms can be vertically differentiated by quality. Producing one unit of the high-quality product may require more production procedures (e.g., longer refinements in the steel industry in [Li et al., 2023](#)), better (or more specialized, exclusive) machinery, higher-quality (or more) intermediate materials, higher standards of quality control (e.g., lower septic infections rate in the healthcare industry [Grieco and McDevitt, 2017](#)), and extra dedicated workers (e.g., promoting quality or demand rather than production as discussed by [Bond et al., 2021](#)). In turn, this leads to a lower quantity of output, holding the inputs fixed, and thus it implies an increase in the *marginal cost* of production (or equivalently a lower productivity). Thus, we refer to $h(\xi_{jnt})$ as the **cost of quality**.¹²

¹²Note that the term cost of quality in this paper refers only to the impact of quality on the marginal cost

As a result, differences in quantity-based productivity can be due to not only technical efficiency but also the cost of quality. Theoretically, explicitly modelling the cost of quality $h(\xi_{jnt})$ as a component of productivity allows for a trade-off between product quantity and quality, conditional on inputs. From an empirical perspective, this also implies that comparisons of quantity-based productivity across firms and over time require control for quality differences. Accordingly, we deal with $\tilde{\omega}_{jnt}$ as a whole rather than its components (ω_{jnt} and $h(\xi_{jnt})$) when estimating the model in Section 3. That is, our estimation method does not rely on how output quality is chosen. We explore the trade-off between (quantity-based) productivity and output quality in Section 6 after they are estimated.

While our empirical model does not rely on the evolution of productivity, we include it to explore the relationship between productivity and quality *after* estimating the empirical model in Section 6 and to facilitate the modelling of dynamic decisions in Online Appendix B. Specifically, we model the evolution of ω_{jnt} as a flexible, endogenous Markov process:

$$\omega_{jnt} = g_n(\boldsymbol{\omega}_{jt-1}, x_{jt-1}) + \epsilon_{jnt}, \quad \forall n = 1, 2, \dots, N, \quad (7)$$

where $g_n(\cdot)$ is a function specific to product category n , ϵ_{jnt} is an innovation term, and x_{jt-1} a set of firm-level decisions implemented in $t - 1$ (such as investment in research and development as emphasized by Doraszelski and Jaumandreu, 2013) that influences the future path of technical efficiency. Importantly, $\boldsymbol{\omega}_{jt} = (\omega_{jt1}, \omega_{jt2}, \dots, \omega_{jnt})$ is a vector of firm-product level technical efficiency of *all* products of firm j in period t . That is, the evolution process of the technical efficiency of one product can be influenced by the previous levels of technical efficiency of other products due to, for instance, intra-firm technology spillovers. The firm observes the realization of $\boldsymbol{\omega}_{jt}$ before making the production decisions specified in Section 2.4.

Remark: Our modelling of the evolution processes is different from that of the literature in three aspects. First, we model the evolution of the underlying technical efficiency rather than quantity-based productivity as in the literature. When quality is an endogenous choice made by the firm and has an impact on quantity-based productivity, quantity-based productivity may no longer evolve in an auto-regressive way, even if the underlying technical efficiency is auto-regressive. An example explicitly used to explore the relationship between productivity and quality is illustrated later by (26) in Section 6.

Second, we allow the evolution processes to be interdependent across products. From a computational perspective, adopting and estimating such flexible evolution processes would add a significant computational burden to the existing proxy-based approach in dealing

of production, rather than the overall cost of quality (including research cost for new products with higher quality, which is more dynamic in nature, or the installation cost of new equipment to produce higher quality products, which are usually one-time fixed costs).

with firms producing many products. Fortunately, our estimation methodology utilizes the first-order conditions of profit maximization to map observable firm input, output, and price choices to unobservable productivity, without relying on the evolution processes, as will become clear in Section 3. This feature is in contrast to the existing estimation methods (e.g., Orr, 2022; Valmari, 2022), which rely on the evolution assumption of productivity and thus exclude a flexible interdependency of productivity among different products.

Third, the literature usually only models the evolution processes of manufactured products due to data and computational limitations. But this approach potentially suffers from an endogenous selection problem because firms only manufacture products when they are profitable. This problem could be severe if the product turnover (i.e., adding and dropping products) is frequent. An appropriate approach is to model the evolution processes of *all* products. But this imposes a challenge in dealing with the latent variables that determine product selection. Our estimation methodology saves us from the data and computational challenges, because it does not rely on the productivity evolution processes.

2.4 Inputs and Outputs Decisions

At the beginning of period t , the firm observes the vector of state variables, which includes the product scope, capital stock, intermediate input price, wage rate, technical efficiency, and product quality of all the products. We summarize the state variables in $s_{jt} = (\Lambda_{jt}, \omega_{jt}, \xi_{jt}, K_{jt}, P_{Mjt}, P_{Ljt}, \chi_{jt})$, where ω_{jt} , ξ_{jt} and χ_{jt} are the vectors of technical efficiency, product quality and demand shocks of *all* the products of firm j in period t , respectively. Note that the observation of technical efficiency and product quality implies that the firm also knows productivity, $\tilde{\omega}_{jt}$, because the firm knows the trade-off (6). P_{Mjt} and P_{Ljt} are the firm-level material price and the wage rate, respectively. Importantly, both of them can be different across firms and vary over time.

The firm's objective is to maximize its total profit from all products in period t after observing its state, by optimally choosing the quantity of material (M_{jt}), the quantity of labor (L_{jt}), and the quantities of all the products to be produced ($\mathbf{Q}_{jt} = \{Q_{jnt}\}, n \in \Lambda_{jt}$). Specifically, the period (static) profit is:

$$\begin{aligned} \pi(s_{jt}) &= \max_{\mathbf{Q}_{jt}, M_{jt}, L_{jt}} \sum_{n \in \Lambda_{jt}} P_{jnt} Q_{jnt} - P_{Mjt} M_{jt} - P_{Ljt} L_{jt} \\ \text{subject to:} & \quad (2) \text{ and } (3). \end{aligned} \tag{8}$$

Remark: In commonly available data, while P_{Ljt} is usually observable to researchers as the wage rate, P_{Mjt} is rarely recorded at the firm level. As documented by Atalay (2014)

using US Census Bureau data, P_{Mjt} can be significantly heterogeneous across firms due to geography, bargaining power and access to the input market, suppliers’ marginal costs, etc. It is well understood that such input price heterogeneity should be controlled for in the production function estimation to avoid bias (i.e., input price bias as emphasized in [Ornaghi, 2006](#); [De Loecker et al., 2016](#); [Grieco et al., 2016](#)). Recent developments in the estimation of multi-product production functions usually assume the availability of P_{Mjt} (or a firm-level index of it, e.g., [Orr, 2022](#); [Valmari, 2022](#)). In contrast, our method is tailored to accommodate common situations where input prices vary at the firm level but are unobservable to researchers. In particular, we maintain the assumption of the literature that P_{Mjt} varies at the firm level (as opposed to the firm-product level) because we model the production as a transformation function (rather than an individual production plant for each product).¹³ We control for P_{Mjt} following the insights of [Grieco et al. \(2016, 2022\)](#), [Harrigan et al. \(2021\)](#), and [Li and Zhang \(2022\)](#), as will be shown in Section 3. Consequently, our empirical method for estimating multi-product production functions offers broader applicability in commonly available datasets compared to existing methods.

3 Estimation

The estimation method leverages a set of implications from the model that can be used to estimate productivity and quality at the firm-product-period level. The method is built upon the insights of [Grieco et al. \(2016, 2022\)](#), [Harrigan et al. \(2021\)](#) and [Li and Zhang \(2022\)](#), who utilize the first-order conditions of static profit maximization to control for unobservable variables in the production function estimation, but it is extended to the multi-product setting where within-firm allocation of inputs is unobserved. Specifically, while researchers do not observe key variables such as productivity and quality, the firm observes them before making optimal production decisions. Thus, the idea is to invert the implications from the profit maximization problem to establish a unique one-to-one mapping from observable production decisions to variables that are unobservable to researchers and control for them in the estimation of the transformation function. Crucially, our model always admits such a mapping regardless of the number of products.

Compared with the existing methods in the literature, our method has several important innovations and advantages, as summarized by Table 1. First, our method models the production technology flexibly as a transformation function and not as a collection of single-product production functions ([De Loecker et al., 2016](#); [Orr, 2022](#); [Valmari, 2022](#); [Chen and Liao, 2022](#)). This saves us from potentially restrictive assumptions regarding how firms

¹³This assumption holds if the input can be costlessly transferred across product lines within the firm, as assumed by [Orr \(2022\)](#) and [Valmari \(2022\)](#).

allocate inputs to produce different products. This is especially important in the presence of shared inputs that serve as public goods within firms. In this regard, [Dhyne et al. \(2022\)](#)'s model is the most similar to ours. Second, our model offers the advantage of scalability as it does not require proxies for product-level productivity and rather relies on static optimization conditions that naturally increase with the number of products. This advantage allows for the analysis of industries with a large number of products without relying on assumptions to aggregate products. Third, our method does not rely on productivity evolution processes. This enables researchers to explore the productivity evolution *after the estimation*, contrary to the existing methods which rely on productivity evolution *for the estimation*. More broadly, this advantage is useful in applications to explore complex (e.g., interdependent) productivity dynamics or when product turnover is frequent and endogenously depends on latent variables (e.g., in the context of exported products). Fourth, our method is designed to deal with the bias caused by unobserved material prices, like [De Loecker et al. \(2016\)](#), but we employ the variation of labor and material expenditure ratio (conditional on the wage rate) to identify material prices. This is particularly useful when material prices are heterogeneous across firms and over time but are unobservable to researchers. Fifth, we assume a demand system, as in [Orr \(2022\)](#), [Valmari \(2022\)](#), and [Chen and Liao \(2022\)](#), but our estimation of demand functions leverages the within-firm revenue relationship implied by profit maximization to estimate demand elasticities with commonly available firm-level IVs. This alleviates the need for firm-product level IVs in the demand estimation that are rarely available.

Table 1: Comparison to existing estimation methods

	Production system	Firm-product productivity	Proxy free	Evolution free	Material price unobservable	Demand system
DGKP	Product				✓	
Orr	Product	✓				✓
Valmari	Product	✓				✓
CL	Product	✓				✓
DPSW	Transformation	✓				
Us	Transformation	✓	✓	✓	✓	✓

Notes: [1] DGKP refers to [De Loecker et al. \(2016\)](#), Orr refers to [Orr \(2022\)](#), Valmari refers to [Valmari \(2022\)](#), CL refers to [Chen and Liao \(2022\)](#), and DPSW refers to [Dhyne et al. \(2022\)](#). [2] In terms of the assumed demand system, our implementation adopts a CES demand function while in Online Appendix A we adopt a general demand function. In comparison, [Valmari \(2022\)](#) and [Chen and Liao \(2022\)](#) also assume a CES demand function while [Orr \(2022\)](#) uses a more flexible Logit model.

This section is organized as follows. In Section 3.1, we first describe how the static profit maximization conditions lead to one-to-one mapping between the observed data and variables

that are unobservable to researchers. In Section 3.2, we derive the estimating equations using the mapping established in Section 3.1 and describe our estimation strategy in detail.

3.1 From Observables to Unobservables: a One-to-one Mapping

We start the description of the estimation strategy by clarifying the observable and unobserved variables in the estimation procedure. We observe capital stock K_{jt} , labor input L_{jt} , labor expenditure E_{Ljt} , material expenditure E_{Mjt} , and quantity Q_{jnt} and price P_{jnt} of each product $n \in \Lambda_{jt}$. We do not observe P_{Mjt} (or M_{jt}), and $\tilde{\xi}_{jnt}$ and $\tilde{\omega}_{jnt}$ for $n \in \Lambda_{jt}$. Our goal is to estimate these unobserved variables together with the production and demand function parameters. Next, we describe the mapping between observed data and unobservables on the basis of the firm's profit maximization.

Note that the firm observes the state s_{jt} (in particular, $\tilde{\omega}_{jnt}$ and $\tilde{\xi}_{jnt}$ for all $n \in \Lambda_{jt}$, P_{Ljt} and P_{Mjt}) as described in Section 2.4 and optimally chooses quantities of inputs and outputs subject to the demand and production functions. The Lagrange function implied by the static profit maximization problem (8) is:

$$\begin{aligned} \mathcal{L}_{jt} = & \sum_{n \in \Lambda_{jt}} (Q_{jnt})^{1-\frac{1}{\eta_n}} e^{\frac{\tilde{\xi}_{jnt}}{\eta_n}} - P_{Ljt}L_{jt} - P_{Mjt}M_{jt} \\ & - \lambda_{jt} \left\{ \left[\sum_{n \in \Lambda_{jt}} e^{-\tilde{\omega}_{jnt}} Q_{jnt} \right] - F(L_{jt}, M_{jt}, K_{jt}) \right\}. \end{aligned} \quad (9)$$

The first-order conditions with respect to labor and material inputs are, respectively:

$$\frac{\partial \mathcal{L}_{jt}}{\partial L_{jt}} = -P_{Ljt} + \lambda_{jt} \frac{\partial F(L_{jt}, M_{jt}, K_{jt})}{\partial L_{jt}} = 0, \quad (10)$$

$$\frac{\partial \mathcal{L}_{jt}}{\partial M_{jt}} = -P_{Mjt} + \lambda_{jt} \frac{\partial F(L_{jt}, M_{jt}, K_{jt})}{\partial M_{jt}} = 0. \quad (11)$$

The first-order condition with respect to each product quantity Q_{jnt} , $n \in \Lambda_{jt}$, is:

$$\frac{\partial \mathcal{L}}{\partial Q_{jnt}} = \frac{\eta_n - 1}{\eta_n} P_{jnt} - \lambda_{jt} e^{-\tilde{\omega}_{jnt}} = 0, \quad (12)$$

where we have used $P_{jnt} = (Q_{jnt})^{-\frac{1}{\eta_n}} e^{\frac{\tilde{\xi}_{jnt}}{\eta_n}}$ according to the demand function (2). The implication of (12), $P_{jnt} = \frac{\eta_n}{\eta_n - 1} \lambda_{jt} e^{-\tilde{\omega}_{jnt}}$, is intuitive: the price is the product of the markup ($\frac{\eta_n}{\eta_n - 1}$) and the marginal cost ($\lambda_{jt} e^{-\tilde{\omega}_{jnt}}$). Within a firm, the marginal cost of a given product differs only due to productivity $\tilde{\omega}_{jnt}$, although the marginal cost also varies across firms due

to λ_{jt} . This is a direct result of the costless input transferability assumption of the production transformation function and profit maximization. Therefore, conditional on a firm, the variation in product prices identifies the productivity difference across products within the firm (after accounting for the markup).

From the perspective of researchers, we do not observe $\tilde{\xi}_{jnt}$, $\tilde{\omega}_{jnt}$ and P_{Mjt} . Nonetheless, we observe the optimal choices which are made based on them by the firm. Thus, utilizing the optimization conditions allows us to recover the unobserved state variables as functions of the observable variables. Specifically, our strategy is to recover $\tilde{\xi}_{jnt}$, $\tilde{\omega}_{jnt}$ and P_{Mjt} as functions of parameters and observable variables including capital stock K_{jt} , labor input L_{jt} , labor expenditure E_{Ljt} , material expenditure E_{Mjt} , quantity Q_{jnt} and price P_{jnt} of each product $n \in \Lambda_{jt}$.

First, we write $\tilde{\xi}_{jnt}$ as a function of observed price and quantity according to the demand function (2):

$$\tilde{\xi}_{jnt} = \ln Q_{jnt} + \eta_n \ln P_{jnt}. \quad (13)$$

Once η_n is estimated, we can recover $\tilde{\xi}_{jnt}$ as above.

Second, we write P_{Mjt} as a function of observable variables. Taking the ratio of equations (10) and (11) and utilizing the expenditure identities (i.e., $E_{Ljt} = L_{jt}P_{Ljt}$ and $E_{Mjt} = M_{jt}P_{Mjt}$), we have:

$$M_{jt} = \left[\frac{\alpha_L E_{Mjt}}{\alpha_M E_{Ljt}} \right]^{\frac{1}{\gamma}} L_{jt}. \quad (14)$$

This implies that material quantity can be recovered from observable variables up to unknown parameters $(\alpha_L, \alpha_M, \gamma)$. Thus, P_{Mjt} is naturally derived by substituting (14) in the expenditure identity (i.e., $E_{Mjt} = M_{jt}P_{Mjt}$):

$$P_{Mjt} = \left[\frac{\alpha_M}{\alpha_L} \right]^{\frac{1}{\gamma}} \left[\frac{E_{Mjt}}{E_{Ljt}} \right]^{1-\frac{1}{\gamma}} P_{Ljt}. \quad (15)$$

In the same spirit of [Grieco et al. \(2016\)](#), the identification of P_{Mjt} comes from the variation of labor and material expenditure ratio (conditional on wage rate), which is implied by the optimality condition under non-Hicks neutrality of the material price in the transformation function.

The third step is to recover $\tilde{\omega}_{jnt}$ for $n \in \Lambda_{jt}$. Specifically, by substituting (14) into (10), we can solve for λ_{jt} as:

$$\lambda_{jt} = \frac{E_{Ljt}}{\rho \alpha_L L_{jt}^\gamma} \left[\alpha_L L_{jt}^\gamma \left(1 + \frac{E_{Mjt}}{E_{Ljt}} \right) + \alpha_K K_{jt}^\gamma \right]^{1-\frac{\rho}{\gamma}}. \quad (16)$$

Then, we substitute (16) into (12) to get:

$$e^{\tilde{\omega}_{jnt}} = \frac{\eta_n}{(\eta_n - 1)P_{jnt}} \frac{E_{Ljt}}{\rho\alpha_L L_{jt}^\gamma} \underbrace{\left[\alpha_L L_{jt}^\gamma \left(1 + \frac{E_{Mjt}}{E_{Ljt}} \right) + \alpha_K K_{jt}^\gamma \right]}_{\lambda_{jt}}^{1-\frac{\rho}{\gamma}}. \quad (17)$$

Noticeably, there are two major components in (17) that identify firm-product-period specific productivity, $\tilde{\omega}_{jnt}$. The first is a firm-level component, λ_{jt} as in (16). This component is the analog of single-product-firm productivity modelled by Grieco et al. (2016) (see their equation (7)). This (unobserved) productivity component is identified from the (unobserved) material price because productivity is Hicks-neutral while the material price is not in our framework.¹⁴ That is, a change in the material price causes a change in the (observable) labor-material expenditure ratio, but a productivity change does not. The second major component, which varies by firm and by product, consists of P_{jnt} and η_n . Intuitively, the variation in product prices helps identify the differences in productivity across products within the same firm, conditional on the elasticity of demand. That is, firms with higher (quantity-based) productivity pass the cost-saving to the product prices (as in Foster et al., 2008). Consequently, the product with a lower price has higher productivity compared with another product manufactured by the same firm, after controlling for the markup (implied by the elasticity of demand). In sum, our identification of $\tilde{\omega}_{jnt}$ uses the variations both at the firm level and at the firm-product level.

Remark: The proxy-based methodology, originated from Olley and Pakes (1996) along with a long list of methodological papers, uses observable variables (such as capital investment and material input) to control for productivity when estimating production functions. Extending the proxy-based approach to the multiple-product context requires valid proxies, which have to admit a one-to-one mapping between the proxies and firm-product level productivity. This is a challenging assumption in the context of a large number of products due to the high dimension of the problem. More importantly, the number of proxies has to increase with the number of products (as recognized by Dhyne et al., 2022), making the extension even more challenging without additional assumptions. The recent development in methods (i.e., Chen and Liao, 2022; Orr, 2022; Valmari, 2022) circumvents this challenge by using production functions of individual products as proxy functions directly, after imputing intra-firm input allocation from firm optimization conditions. This approach assumes that there is no transitory error in production (which is explicitly modelled and dealt with by Olley

¹⁴If both productivity and the material price are Hicks-neutral in the production function, as in Cobb-Douglas production functions, then this identification strategy fails. However, in this case, the labor-material expenditure ratio would be a constant under the optimality condition, which is not supported by the data.

and Pakes, 1996) and that the persistent error (as in the traditional notion of productivity) evolves independently according to a Markov process.

In contrast, our methodology uses first-order conditions to construct an explicit one-to-one mapping for productivity (up to the parameters to be estimated). This does not only guarantee the existence and uniqueness of the mapping from observable data to unobservable heterogeneity, but also lends us a significant advantage in dealing with scenarios where firms produce a large number of products, because the number of first-order conditions naturally increases with the number of products. In addition, this methodology of recovering unobservable heterogeneity (instead of imputing input allocation) saves us from estimating productivity evolution processes as a part of the production estimation, which can dramatically complicate the existing methods in the literature, especially when there are endogenous, frequent entry and exit of products or the evolution processes of productivity are interdependent. More broadly, this feature allows our method to be widely applied to analyzing the impact of policy shocks on productivity, which would have to be otherwise considered as a part of the evolution processes (as emphasized by Chen et al., 2021) and further complicate the estimation process using the existing methods.

3.2 Estimating Equations and Strategy

In the previous subsection, we have explicitly constructed a one-to-one mapping from observable variables to the unobserved $\tilde{\xi}_{jnt}$, $\tilde{\omega}_{jnt}$, and P_{Mjt} (or M_{jt} equivalently) up to a set of parameters to be estimated. This mapping is the key to developing the estimating equations, which we derive in this subsection. Next, we describe in detail the strategy to estimate the key parameters of the model.

We assume that there is a firm-level measurement error (or unexpected shock) in revenue: $R_{jnt} = P_{jnt}Q_{jnt}e^{u_{jt}}$, where u_{jt} is a mean-zero independent and identically distributed shock.¹⁵

¹⁵An alternative way to view this unexpected shock u_{jt} is to consider it as a transitory shock (in addition to productivity) to the transformation function (3) in the estimation: $\sum_{n \in \Lambda_{jt}} e^{-\tilde{\omega}_{jnt}} Q_{jnt} = [\alpha_L L_{jt}^\gamma + \alpha_M M_{jt}^\gamma + \alpha_K K_{jt}^\gamma]^\frac{p}{\gamma} e^{u_{jt}}$. The distinction between u_{jt} and productivity ($\tilde{\omega}_{jnt}$) is that the firm observes productivity when making decisions, causing it to be correlated with input choices, whereas u_{jt} is not observed by the firm (only *ex post*) and is thus uncorrelated with input choices. Importantly, because u_{jt} is only observed *ex post*, it does not affect the production decisions (i.e., the first-order conditions) and only becomes an additive error in the estimating equation (18). A full derivation with such an *ex post* error term in the revenue is described in Online Appendix A. Nonetheless, such a shock, realized *ex post*, will be reflected a part of observed P_{jt} or Q_{jt} and enter the estimated quality (13) or productivity (17) in an additive way. Consequently, we include a firm-time specific dummy as a control variable in studying the relationship between quality and productivity in Section 6.

We use this definition and substitute (14) and (17) into (3) to obtain the estimating equation:¹⁶

$$\ln \left[\sum_{n \in \Lambda_{jt}} \frac{(\eta_n - 1)\rho}{\eta_n} R_{jnt} \right] = \ln \left[E_{M_{jt}} + E_{L_{jt}} \left(1 + \frac{\alpha_K}{\alpha_L} \left(\frac{K_{jt}}{L_{jt}} \right)^\gamma \right) \right] + u_{jt}. \quad (18)$$

This equation is the multi-product version of the estimating equation proposed by Grieco et al. (2016) (see their equation 8), who assume that each firm produces a single product.¹⁷ In the context of multi-product firms, the individual product revenues are adjusted by the reciprocal of their corresponding markups.¹⁸ This equation is extendable to a general demand system which explicitly allows for complementarity or substitutability across products, as shown in Online Appendix A. Although u_{jt} does not affect production decisions, it does appear as a part of observed product revenues. A higher shock implies a higher realized revenue R_{jnt} . Thus, u_{jt} and R_{jnt} are correlated. This also implies that we need to estimate the model via the generalized method of moments (GMM).

Nonetheless, estimating all the parameters using (18) alone faces two challenges. First, ρ is not separately identified from demand elasticities in (18). In fact, only a combination of η_n and ρ (i.e., $\frac{(\eta_n - 1)\rho}{\eta_n}$) is identified by (18).¹⁹ Second, (18) requires (at least) the same number of instrumental variables as the number of products to identify $\frac{(\eta_n - 1)\rho}{\eta_n}$ of each product, because all product revenues are correlated with u_{jt} .

To address the two challenges at the same time, we explore the relationship between the revenues of any two products implied by the firm's static maximization problem, taking into account that the markets for different products are segmented. Notably, η_n influences the sales of *individual* products, while ρ represents the returns to scale of the production transformation function and affects the overall sales of *all* products. Thus, the firm's optimal decision on trading off the sales of different products within the firm helps identify η_n from ρ . In other words, the variation in the sales of a product relative to another product contains information on how the elasticities of the two products differ. This addresses the first challenge. Meanwhile, the identified relationship between elasticities reduces the number

¹⁶The detailed algebra for a more general model is demonstrated in Online Appendix A. Also, as shown in Online Appendix A, in the more general function form (4) where output demonstrates complementarity, (18) stays the same.

¹⁷More broadly, (18), without logarithms, is also similar to the estimating equations used by Das et al. (2007), Aw et al. (2011), and Li (2018) with data on the firm's total variable cost to estimate demand elasticities in multiple markets.

¹⁸If the elasticities (markups) are the same, then the estimating equation is the same as in Grieco et al. (2016). We also allow for the returns to scale parameter, ρ , to be estimated, while Grieco et al. (2016) assume it to be one.

¹⁹The identification of $\frac{(\eta_n - 1)\rho}{\eta_n}$ relies on the normalization condition that the mean of unexpected revenue shock is zero: $E(u_{jt}) = 0$. If $E(u_{jt}) \neq 0$, then one cannot identify whether a higher revenue is from a larger return to scale or a higher unexpected shock.

of parameters to be estimated in (18). Consequently, the number of instrumental variables required to estimate the rest of the parameters does not increase with the number of products. This addresses the second challenge.

To implement this idea, we define a reference product. In principle, the reference product can be any product. However, the reference product may not be produced by all firms. Thus, from an empirical point of view, we use the product that is manufactured by most firms in the industry, to maximize the number of observations one could use in the estimation.²⁰ Without loss of generality, we denote the reference product as product 1. For any firm j , taking the ratio of (12) of the reference product and that of another product n and using $R_{jnt} = P_{jnt}Q_{jnt}$, we obtain:²¹

$$\ln(R_{jt1}) = c_n + \frac{\eta_1 - 1}{\eta_n - 1} \ln(R_{jnt}) + \mu_{jnt}, \quad n = 2, \dots, N, \quad (19)$$

where

$$c_n = (1 - \eta_1) \ln \left[\frac{\eta_1}{\eta_1 - 1} \frac{\eta_n - 1}{\eta_n} \right]$$

and

$$\mu_{jnt} = (\eta_1 - 1) \left[\underbrace{\left(\tilde{\omega}_{jt1} + \frac{1}{\eta_1 - 1} \tilde{\xi}_{jt1} \right) - \left(\tilde{\omega}_{jnt} + \frac{1}{\eta_n - 1} \tilde{\xi}_{jnt} \right)}_{\text{difference in quality-adjusted productivity}} + \underbrace{\frac{\eta_1 - \eta_n}{(\eta_1 - 1)(\eta_n - 1)} u_{jt}}_{\text{measurement error component}} \right].$$

The latter, μ_{jnt} , contains the *difference* of the capability (or quality-adjusted productivity, $\tilde{\omega} + \frac{1}{\eta-1} \tilde{\xi}$, as will be formally defined in Section 5) of producing a product relative to that of the reference product and composition of the unexpected shock. This equation predicts that the (logarithmic) revenues of two products are linearly related conditional on the *difference* of production capability. In particular, firm-level inputs are not a part of the equation explicitly. This equation is similar to the estimating equation developed by Grieco et al. (2022), who explore the relationship of revenues of two markets (domestic sales and exports).²²

²⁰Across the three industries in our empirical exercise, the percentage of firm-year pairs that produce the reference product ranges from 62% in footwear to 72% in printing and 88% in pharmaceutical.

²¹In the more general function form (4) where output demonstrates complementarity, this equation (19) becomes $\ln(R_{jt1}) = c_n + \frac{1-\theta}{1-\theta\frac{\eta_n-1}{\eta_1-1}} \ln(R_{jnt}) + \mu_{jnt}$, $n = 2, \dots, N$, where $c_n = \frac{1-\theta}{1-\theta\frac{\eta_1}{\eta_1-1}} \ln \left[\frac{\eta_1}{\eta_1-1} \frac{\eta_n-1}{\eta_n} \right]$. However, it can be shown that the additional parameter θ cannot be separately identified from ρ using this equation and (18). In Online Appendix A, we extend our model to identify θ using additional restrictions, which requires estimating the demand system (and thus demand elasticities) directly using appropriate instrumental variables. Overall, the advantage of imposing $\theta = 1$ is to allow us to estimate demand elasticities jointly from (18) and (19) without estimating the demand function directly.

²²One difference is that Grieco et al. (2022) model the error term as an unexpected shock because the

Intuitively, because the demand for each product is segmented in our setting, as discussed in Sections 2.1 and 4, the relative revenue of one product over another product in the same firm depends on their own demand elasticities (conditional on their relative levels of productivity and quality, measured as μ_{jnt}) rather than on complementarity or substitution between them. As a result, the variation of one revenue *relative* to another in (19) provides the identification of the ratio, $\frac{\eta_1-1}{\eta_n-1}$ for $n = 2, 3, \dots, N$. In contrast, the variation of revenue *levels* in (18) identifies $\frac{(\eta_n-1)\rho}{\eta_n}$, $n = 1, 2, \dots, N$. That is, the returns to scale parameter affects the sales of all products but not the relative relationship of sales between different products, while demand elasticities affect both the level and the relative relationship of sales of different products. As a result, ρ and η_n , $n = 1, 2, \dots, N$, are separately identified as long as there are at least two products with different demand elasticities in the industry. The model is over-identified when there are more than two products produced by the firms in the industry. More precisely, the elasticities and returns to scale parameter can be identified as long as there is a firm that manufactures two products with different demand elasticities for a number of periods, which is a very mild assumption.

To estimate (19), we treat μ_{jnt} as an error term. We allow the mean of μ_{jnt} to vary by product and year and use a set of flexible product-year dummies as controls (which also absorb c_n). μ_{jnt} is likely correlated with R_{jnt} – the revenue of product n is lower if the capability of producing n is lower than that of the reference product. We use a set of IVs to address the endogeneity issue. In our implementation, the IV set consists of a constant and the logarithm of the wage rate (P_{Ljt}), the capital stock (K_{jt}), and the ratio of material expenditure to labor (E_{Mjt}/L_{jt} , as a proxy for material prices after conditional on wage rate).²³ Grieco et al. (2022) uses a similar set of firm-level IVs to estimate an equation analogous to (19) in a two-product scenario. The same insight carries over in our context. These firm-level variables influence the *level* of revenue (i.e., R_{jnt}), but they are uncorrelated with the *difference* of capability (i.e., μ_{jnt}) between two products. For example, conditional on everything else, a higher level of capital stock potentially leads to higher revenues of a given product, but it is not necessarily associated with the production capability of one product being larger than that of another product within the same firm. Thus, we use these firm-level variables as IVs for all product pairs in (19).²⁴

productivity and quality of the domestic and export products are assumed to be the same and thus cancel out.

²³To see this, note that (15) is equivalent to $P_{Mjt} = \left[\frac{\alpha_M}{\alpha_L}\right]^{\frac{1}{\gamma}} \left[\frac{E_{Mjt}}{L_{jt}}\right]^{1-\frac{1}{\gamma}} P_{Ljt}^{\frac{1}{\gamma}}$. Taking logarithm, we obtain $\ln(P_{Mjt}) = \frac{1}{\gamma} \ln\left[\frac{\alpha_M}{\alpha_L}\right] + (1 - \frac{1}{\gamma}) \ln\left[\frac{E_{Mjt}}{L_{jt}}\right] + \frac{1}{\gamma} \ln(P_{Ljt})$. Because we include the logarithm of the wage rate, $\ln(P_{Ljt})$, in the IV set, using $\ln\left[\frac{E_{Mjt}}{L_{jt}}\right]$ is equivalent to using $\ln(P_{Mjt})$ in this setting, although P_{Mjt} is not observable. Our result is quantitatively similar if the expenditure ratio of material and labor is used as an IV.

²⁴The model is over-identified if there is more than one IV. For example, if there are 2 IVs, then there are

The validity of these IVs relies on the condition that the production of a product is not systematically more intensive in the use of a specific input (e.g., capital) than other products and that the wage rate and input price are not systematically correlated with the capability *differences* between products. We use Monte Carlo exercises to demonstrate the performance of our approach and IVs under this condition in Online Appendix D. In our structural framework, given the assumption of transformation function and costless transferability of inputs across products, this condition is satisfied.²⁵

Remark: In cases where this condition does not hold, alternative solutions are available. First, the demand elasticities can be, in principle, identified by the demand function (2) using the variation of prices and quantities. The primary challenge in directly estimating (2) is the availability of IVs – ideally at the firm-product-time level – that are uncorrelated with product quality. Commonly used cost shifters are often unsuitable because firms producing high-quality products typically use higher-cost, high-quality inputs. However, if appropriate IVs are available for directly estimating the demand function (2), it eliminates the need to estimate (19), allowing the main equation (18) to be estimated using a Nonlinear Least Squares estimator, rather than GMM. For example, Orr (2022) estimate a demand system directly by employing sophisticated IVs that exploit variations in product sets and input price growth across firms in different output markets that use similar inputs.

Second, if one is willing to assume constant return to scale (i.e., $\rho = 1$), then the demand elasticities can be identified using (18) alone, without relying on the strategy involving (19). In fact, with the constant return to scale assumption, (18) degenerates to the estimating equations used by Das et al. (2007), Aw et al. (2011), and Li (2018). These papers utilize the relationship between the total variable cost (as our counterpart of the right-hand side of (18)) and export revenues (as our counterpart of the left-hand side of (18)) of the same firm to estimate demand elasticities in multiple export markets.

We denote the estimated relationship between elasticities as $\hat{b}_n = \frac{\eta_n - 1}{\eta_n - 1}$, $n = 2, \dots, N$, and, naturally, $\hat{b}_1 = 1$ by definition. Thus, $\eta_n = \frac{1}{\hat{b}_n}(\eta_1 - 1) + 1$. Substitute it as η_n in (18) and

$2(N - 1)$ moment equations that can be formed to identify $(N - 1)$ coefficients (i.e., $\frac{\eta_n - 1}{\eta_n - 1}$, $n = 2, \dots, N$).

²⁵To examine this assumption empirically, we check whether the IVs are correlated with either the within-firm product shares or the ratio of log sales of a product over that of the baseline product as alternative measures of relative production capability. Specifically, we regress each IV on either the interactions between product fixed effects and within-firm revenue shares (including firm and year fixed effects) or the interactions between product fixed effects and the ratio of log sales of a given product over that of the baseline product (including firm and year fixed effects). We find that at least 85% of coefficients (i.e., products) are not significant at the 1% level in these tests for our IVs.

solve for u_{jt} to construct moment conditions for the GMM estimation:

$$u_{jt} = \ln \rho + \ln \left[\sum_{n \in \Lambda_{jt}} \frac{\eta_1 - 1}{\eta_1 - 1 + \hat{b}_n} R_{jnt} \right] - \ln \left[E_{M_{jt}} + E_{L_{jt}} \left(1 + \frac{\alpha_K}{\alpha_L} \left(\frac{K_{jt}}{L_{jt}} \right)^\gamma \right) \right]. \quad (20)$$

There are only four parameters, $\beta \equiv (\rho, \eta_1, \frac{\alpha_K}{\alpha_L}, \gamma)$, to be estimated. This means that the number of instrumental variables required does not increase with the number of products. In particular, firm-level input choices can serve as valid IVs because they are not correlated with the unexpected shock u_{jt} . In the implementation, we use $Z_{jt} = (1, E_{M_{jt}}, E_{L_{jt}}, L_{jt}, K_{jt}/L_{jt})$ as IVs. Our results are robust to a set of alternative firm-level IVs.

Equation (20) can only identify $\frac{\alpha_K}{\alpha_L}$ rather than α_L , α_M , and α_K separately. As shown by Grieco et al. (2016), the full set of $(\alpha_L, \alpha_M, \alpha_K)$ can be identified with two constraints naturally implied by the model. The first constraint is a normalization of distribution parameters in the CES production function: $\alpha_L + \alpha_M + \alpha_K = 1$. The second constraint equalizes the ratio of geometric means of labor expenditure (\bar{E}_L) and material expenditure (\bar{E}_M) to the ratio of distribution parameters in the CES production function. That is, $\frac{\alpha_M}{\alpha_L} = \frac{\bar{E}_M}{\bar{E}_L}$. This constraint results from taking the geometric mean of (14), which is implied by the first-order conditions of labor and material quantities, (10) and (11), of all firms.²⁶

As a result, β can be estimated as:

$$\begin{aligned} \hat{\beta} &= \operatorname{argmin}_{\beta} \left[\frac{1}{\mathbb{N}} \sum_{j,t} u_{jt} Z_{jt} \right]' W \left[\frac{1}{\mathbb{N}} \sum_{j,t} u_{jt} Z_{jt} \right], \quad (21) \\ \text{subject to:} \quad &\alpha_L + \alpha_M + \alpha_K = 1 \quad \text{and} \quad \frac{\alpha_M}{\alpha_L} = \frac{\bar{E}_M}{\bar{E}_L}, \end{aligned}$$

where W is a weight matrix and \mathbb{N} is the number of firm-time observations.

As a summary of the full estimating approach, the first step is to estimate $\hat{b}_n = \frac{\eta_1 - 1}{\eta_n - 1}$, $n = 2, \dots, N$ via Two-Stage Least Squares (2SLS) using the relationship imposed by the within-firm relative sales in (19). The second step is to estimate $(\hat{\rho}, \hat{\eta}_1, \hat{\alpha}_L, \hat{\alpha}_M, \hat{\alpha}_K, \hat{\gamma})$ using (20) via GMM. With these estimates, the demand elasticities can be recovered as $\hat{\eta}_n = \frac{1}{\hat{b}_n}(\hat{\eta}_1 - 1) + 1$. After that, we compute $\tilde{\xi}_{jnt}$, $\tilde{\omega}_{jnt}$, and $P_{M_{jt}}$ via (13), (17), and (15), respectively. We demonstrate that our method is able to recover the true parameter values in the Monte Carlo

²⁶As shown by Grieco et al. (2016), this constraint holds conditional on a normalization of the CES production function. Thus, we follow the same procedure to normalize the inputs using their corresponding industry-level geometric means as in the literature (e.g., Klump and de La Grandville, 2000; León-Ledesma et al., 2010). Nonetheless, to ease our notation, we directly denote the normalized input variables as (L_{jt}, M_{jt}, K_{jt}) . As a result, the ratio of the geometric means of material and labor is $\frac{\bar{M}}{\bar{L}} = 1$, which implies $\frac{\alpha_M}{\alpha_L} = \frac{\bar{E}_M}{\bar{E}_L}$, by taking the geometric mean of (14) across firms.

exercises shown in Table A7 of Online Appendix D.

4 Data

We estimate our model using firm-level Mexican manufacturing data, collected by the *Instituto Nacional de Estadística y Geografía* (National Institute of Statistics and Geography, INEGI henceforth) and covering the period 1994-2007. We use two datasets: the *Encuesta Industrial Anual* (Annual Industrial Survey, EIA henceforth), the main annual survey covering the manufacturing sector, and the *Encuesta Industrial Mensual* (Monthly Industrial Survey, EIM henceforth), a monthly survey that monitors short-term trends related to employment and output.²⁷ These datasets are particularly useful for our analysis because they provide quantity and sales information at the firm-product level.

Next, we describe in more detail these two surveys and the variables we extract from them.²⁸ The EIA contains information on 6867 firms in 1994, but this number decreases over time due to attrition. It covers roughly 85 percent of all manufacturing output value based on information from the industrial census, but it excludes assembly plants, i.e., “maquiladoras”. The EIA includes variables related to output indicators, inputs, and investment. These data make it possible to calculate the value of intermediate inputs and physical capital stock based on information on investment and the perpetual inventory method. The EIM runs in parallel with the EIA and covers the same firms. The EIM contains information on the number of workers and their wage bills so that the average wage at the firm level can be calculated. The EIM also contains output-related variables, in particular values and quantities of sales at the product level, so that an implicit average unit price can be calculated.²⁹

Firms are classified by INEGI into one of the classes of activity based on their principal product. A class of activity is the most disaggregated level of industrial classification and is defined at six digits according to the 1994 *Clasificación Mexicana de Actividades y Productos* (Mexican System of Classification for Activities and Products, CMAP henceforth). Firms report information product by product based on their industries and a list of products provided by INEGI.

In this paper, we focus on three specific classes of activities: manufacturing of footwear, mainly of leather (class 324001, footwear in short); printing and binding (class 342003, printing in short); and manufacturing of pharmaceutical products (class 352100, pharmaceuticals in

²⁷The unit of observation in both surveys is a plant rather than a firm and the sample includes all plants with more than 100 employees as well as a sample of smaller plants. For simplicity and in line with the literature, we will use the term “firm” to refer to a plant.

²⁸More information on the EIA and EIM can be found in Caselli et al. (2017) and Caselli (2018).

²⁹All nominal variables are deflated using the consumer price index. To facilitate comparison, we normalize average industry output prices to 1. Initial capital stock and investment are deflated using industry-level price indices.

short). These three industries were chosen because each industry is made up of more than 500 firm-year pairs, a number of observations large enough for our estimation strategy. More importantly, multi-product firms are particularly prevalent in these industries – 65% of firms in these three industries are multi-product producers and such firms account for 86% of total revenues and produce on average 6.7 products per year.³⁰ They also represent a diverse set of manufacturing industries with clear concepts/characteristics of product quality: for example, advanced design and assembly that provide superior comfort and durability in the footwear industry; acid-free paper and durable binding in the printing industry; potent active ingredients and degrading-preventing packaging in the pharmaceutical industry.

For the purpose of the production function estimation in Section 5, all products with fewer than 100 observations are aggregated together in a residual product category.³¹ The prices and quantities of the aggregated residual product category are estimated following [Diewert et al. \(2009\)](#) and [Caselli \(2018\)](#). While this aggregation is required to estimate the demand elasticity of substitution for each product based on a large enough number of observations, it only implies that the demand elasticity of substitution is by assumption equal across all products included in the residual product category within an industry. In addition, this aggregation involves a relatively small share of products: the main (i.e., not aggregated) products account for between 74% and 92% of observations and 82% and 90% of revenue across the three industries. Accordingly, the descriptive statistics and patterns demonstrated in this section are reported based on the aggregated categories, which is the data used in the estimation in Section 5.

There are a few patterns worth noticing. First, multi-product production is an essential feature of the firms in our sample. We demonstrate this point by using an index that is analogous to the traditional Herfindahl–Hirschman Index (HHI). Specifically, we construct an analog index of HHI as the sum of the squared shares of sales within a firm. A higher HHI index means a higher level of concentration of sales within a firm.³² The index is naturally equal to one for single-product producers. For firms with a larger product scope, HHI decreases sharply becoming close to 0.3 for firm-year pairs producing 5 products and close to 0.2 for firm-year pairs producing 10 or more products.³³ These values imply that

³⁰Tables [A1](#), [A2](#) and [A3](#) in the Online Appendix show how detailed the product-level information is by reporting the list of products with at least 100 observations for each of the three chosen industries.

³¹The residual product category is defined as “Others” (product code 99) in Tables [A1](#), [A2](#) and [A3](#) in the Online Appendix.

³²In Figure [A1](#) in the Online Appendix, we aggregate the firm-level index with weights equal to the firms’ total revenues, by firm-year pairs’ product scope.

³³These values indeed show some degree of concentration of sales within firms. For example, if a firm produces 5 products with equal sales, the index would be 0.2. The fact that the index is close to 0.3 implies that there exists an uneven distribution of sales. We explore this heterogeneity using quality and productivity within firms in Section 6.

producers are genuine multi-product firms – they do not concentrate production entirely on their top products, and all products, albeit to different degrees, are important for firms’ total revenues.³⁴ Thus, multi-product firms need to be treated and modelled as such and they cannot be simplified as single-product producers.

Table 2: Descriptive statistics: prevailing multi-product firms

Variable	Footwear	Printing	Pharmaceutical
Product scope, all firms	1.289 (0.627)	3.708 (3.752)	6.858 (3.740)
Product scope, MPFs only	2.388 (0.602)	5.891 (3.841)	7.925 (3.024)
Share of MPFs	0.208	0.554	0.846
Revenue share of MPFs	0.389	0.599	0.940
Total number of products	4	14	16
Total number of firms	72	83	82
Average number of firms per product-year	21	19	43
Number of firm-year pairs	707	831	928

Notes: The table reports the means and standard deviations (in parenthesis) for each variable by industry. Product scope is the number of products manufactured by firm. MPFs refers to multi-product firms only.

The importance of multiple-product production is also present in all the industries of our analysis, albeit with some degrees of variation, as shown in Table 2.³⁵ The percentage of multi-product firms ranges from 21% in the footwear industry to 55% in printing and 85% in pharmaceuticals and they account for an even larger share of revenues (from 39% in the footwear industry to 94% in pharmaceuticals). The average product scope is larger in printing and pharmaceuticals (respectively, 5.9 and 7.9 for multi-product firms) than in the footwear industry (2.4). These differences in average product scope are in line with the number of product categories available in each industry, which ranges from 4 in footwear to 16 in pharmaceuticals.

Second, the status of being a multi-product firm is quite persistent, and so is the product scope. In particular, using a simple autoregressive process of the number of products produced by each firm, we measure the persistence coefficients are 0.87, 0.95, and 0.98 in the three industries, respectively.³⁶ Thus, multi-product firms unequivocally dominate manufacturing

³⁴To confirm that firms rely heavily on all products for their total sales, Online Appendix Table A4 shows the average within-firm product shares by product scope. For instance, for firms producing 5 or more products, the share of products other than the top product is 0.557 and the share of products with rank 5 and beyond is 0.146, on average.

³⁵Additional descriptive statistics are available in Table A5 in the Online Appendix.

³⁶The entry of new products and the exit of old products only account for 3.8 and 4.4 percent of the observations, respectively.

production in our data and their within-firm adjustment across products is more salient than the extensive margin adjustment in changing the number of products.

These patterns imply that both within-firm and across-firm heterogeneity is important. On the one hand, there exist persistent characteristics at the firm level that determine the performance across firms. On the other hand, intra-firm heterogeneity and product scope play a significant role in shaping these characteristics within firms. These implications are in line with the specification for productivity (17), which contains a common component at the firm level to capture the differences across firms as well as an individual component varying at the firm-product level to explain the variation of performance within a firm.

Finally, the sample reflects patterns consistent with the model’s demand assumption. On average, about 19 to 43 firms are competing in the market for any given product in any given year. The majority of the firms do not command a dominant share of the market – the median (traditionally defined) HHI at the product-year level ranges between 0.11 in the pharmaceutical industry and 0.26 in the printing industry. More importantly, given the level of product disaggregation, the markets for different products (e.g., women’s shoes vs. men’s shoes in the footwear industry) are reasonably assumed as segmented. For each product, firms’ outputs are vertically differentiated as evidenced by the large dispersion in prices.³⁷ Overall, these patterns support abstracting from demand cannibalization across products manufactured by the same firm and assuming that firms face monopolistic competition within each product category.

5 Estimation Results

In this section, we apply the empirical model to the data and estimate the production and demand function parameters by industry, which then allows us to compute firm-product level productivity and quality. Notably, our approach employs a novel method, and despite this novelty, the resulting structural parameter estimates align closely with existing literature. Moreover, the productivity and quality measures derived from these estimates exhibit economically meaningful properties. Because our empirical analysis relies on estimated variables, we employ bootstrapping with 100 samples to compute all standard errors presented in the subsequent tables, ensuring robustness and accuracy in our findings.

Table 3 presents the production function parameters. α_M is significantly larger than α_L and α_K , consistent with the intensive use of intermediate material input across all industries. α_K in the pharmaceutical industry is the highest among the three industries, reflecting the importance of capital in this industry. Parameter σ , which is the elasticity of substitution

³⁷For example, the interquartile range of prices in logarithm is 1.4 (i.e., a 400% difference) within a product category, on average, across the three industries.

across inputs, i.e., labor, material, and capital, is greater than one across all industries. This is different from those in the classical literature which does not control for heterogeneous material prices. But it is largely consistent with the estimates in [Grieco et al. \(2016, 2022\)](#), [Harrigan et al. \(2021\)](#), and [Li and Zhang \(2022\)](#) based on a similar method but using different datasets from Colombia (ranging from 1.4 to 2.6), France (ranging from 1.1 to 2.0), and China (ranging from 1.2 to 2.7), respectively. It is also close to the average estimate (around 1.4) of the elasticity of substitution among Chinese industries by [Berkowitz et al. \(2017\)](#) using a different method. Finally, the returns to scale parameter ρ of the three industries is larger than one, but it is not significantly different from one, implying that the production is close to constant returns to scale in these industries, except in the case of the footwear industry.

Table 3: Production function estimates

Parameter	Footwear	Printing	Pharmaceutical
α_L	0.202 (0.013)	0.229 (0.017)	0.227 (0.021)
α_M	0.774 (0.035)	0.673 (0.032)	0.595 (0.063)
α_K	0.023 (0.044)	0.099 (0.041)	0.178 (0.079)
σ	1.518 (0.568)	1.244 (0.171)	1.168 (0.241)
ρ	1.227 (0.107)	1.078 (0.095)	1.002 (0.127)

Note: Bootstrapped standard errors clustered at the firm level and stratified by industry and scope are shown in parentheses (100 repetitions).

Table 4 presents the estimates of the demand elasticities of substitution of different products in the three industries. These estimates generally fall within a similar range as those found in the existing literature (e.g., see [Roberts et al. \(2018\)](#); [Grieco et al. \(2016\)](#); [Dubois and Lasio \(2018\)](#)). Our approach is in contrast to the literature, which often relies on direct estimation of the demand function while assuming time-invariant product quality and/or using firm-level instrumental variables, such as capital stock. By leveraging the multi-product context, as described in (19), we capitalize on the advantage of utilizing firm-level IVs that may be potentially correlated with the *level* of quality but are less likely to be correlated with the *difference* in production capabilities of any two products within the firm.³⁸

³⁸When we estimate the demand function (2) directly using the same firm-level IVs, the estimated demand elasticities are significantly biased towards zero: the mean elasticities are -0.005, 1.941, and -0.395 for the footwear, printing, and pharmaceutical industries, respectively.

Table 4: Demand function estimates

Parameter	Footwear	Printing	Pharmaceutical
η_1	2.823 (0.539)	4.523 (1.583)	3.688 (1.275)
η_2	2.455 (0.540)	8.661 (2.683)	3.037 (1.616)
η_3	3.588 (0.699)	4.432 (1.451)	4.209 (2.145)
η_4	3.250 (0.713)	7.321 (2.204)	3.999 (2.000)
η_5		4.448 (1.596)	4.010 (2.057)
η_6		4.769 (1.913)	2.712 (0.904)
η_7		5.140 (1.704)	3.544 (1.620)
η_8		6.157 (2.325)	3.210 (1.352)
η_9		7.139 (2.202)	3.133 (1.761)
η_{10}		4.838 (1.490)	3.263 (1.408)
η_{11}		6.682 (1.845)	3.418 (1.934)
η_{12}		5.588 (1.669)	3.047 (1.027)
η_{13}		4.279 (2.009)	4.713 (2.058)
η_{14}		5.379 (1.416)	7.279 (2.462)
η_{15}			2.431 (1.937)
η_{16}			2.809 (1.654)

Note: Bootstrapped standard errors clustered at the firm level and stratified by industry and scope are shown in parentheses (100 repetitions).

The variations in demand elasticities across products, as documented above, lead to differences in markups at the firm-year level. These markups can be calculated as the weighted average of product markups considering their respective shares within firms.³⁹

³⁹Across the three industries, the average markup at the firm-year level is 1.40 with a standard deviation

Estimated dispersion in markups is smaller than the estimate reported by [De Loecker and Warzynski \(2012\)](#). This is because our variation of firm-year-level markups only captures the heterogeneous revenue shares and sets of products (as well as their associated markups) manufactured by different firms. Despite this narrower focus, the dispersion of markups at the firm-year level remains significant.

After all model parameters are estimated, we compute the firm-product-time varying output quality and productivity according to (13) and (17) in logarithm, respectively.⁴⁰ Nonetheless, these two measures are not directly comparable across and within firms. This is because the varieties (in the same product category) are of different quality levels and the unit of measurement across different products can be also different (e.g., grams vs. liters). However, the quality-adjusted output is readily comparable across firms and products, as shown by [Melitz \(2000\)](#), [Orr \(2022\)](#), and [Li et al. \(2023\)](#). Thus, we follow the literature to construct a combined measure that takes both quality and productivity into account. In our context, given the setup of quality-adjusted output in (1), we define a quality-adjusted productivity (ATFP) measure as⁴¹

$$\text{ATFP}_{njt} = \tilde{\omega}_{njt} + \frac{1}{\eta_n - 1} \tilde{\xi}_{njt}. \quad (22)$$

As expected, ATFP reflects significant dispersion across firms even within a specific product category.⁴² The mean interquartile range within a product is about 2.8 (calculated across all products in the three industries), which is similar in magnitude to that of revenue productivity documented by [Grieco et al. \(2022\)](#) in the Chinese paint industry. Regarding the components of ATFP, the interquartile range of $\tilde{\omega}_{jnt}$ within a product has a mean of 2.8, while the interquartile range of $\frac{1}{\eta_n - 1} \tilde{\xi}_{njt}$ within a product has a mean of 1.8.⁴³

Overall, our estimation results reflect reasonable parameter estimates and productivity and quality measures at the firm-product level. In the following sections, we turn to use these measures to explore the role of productivity and quality in shaping intra-firm performance heterogeneity.

of 0.14. The interquartile range (using the logarithm of the markups) is 0.16.

⁴⁰We also compute firm-level intermediate input prices according to (15). We find that there is significant heterogeneity in intermediate prices, as documented by [Ornaghi \(2006\)](#) using observed intermediate price data.

⁴¹This measure is similar to the conventionally defined revenue-based productivity (a.k.a., TFPR).

⁴²The distributions of ATFP by product, as well as the distributions of its components, $\tilde{\omega}_{njt}$ and $\tilde{\xi}_{njt}$, are reported in Figures [A2](#), [A3](#) and [A4](#), respectively.

⁴³The interquartile range of $\tilde{\omega}_{jnt}$ is slightly larger than that of ATFP because the two components of ATFP, productivity, and quality, are negatively related, as will be clear in [Section 6](#).

6 Intra-firm Heterogeneity: Productivity and Quality

With the structural parameters reasonably estimated and the rich distributions of productivity and quality revealing significant heterogeneity even within narrowly defined product lines, we now turn to the pivotal question: what new insights about multi-product firms emerge from our analysis? We focus on how firm-product-level heterogeneity in productivity and quality influences the relative performance of different products within a firm and, importantly, the relationship between these attributes.

The literature traditionally emphasizes the role of productivity in explaining the growth and performance of firms and industries (e.g., Jovanovic, 1982; Hopenhayn, 1992; Ericson and Pakes, 1995; Melitz, 2003). Recently, a growing literature shows that demand is also important for firm turnover and growth (e.g., Foster et al., 2008; Pozzi and Schivardi, 2016; Kumar and Zhang, 2019). However, this strand of the literature usually focuses on across-firm analysis using firm-level data. Taking into account the joint nature of production in multi-product firms, our estimation method allows us to uncover rich, flexible dimensions of heterogeneity within firms and explore the role of productivity and demand at the firm-product level.

Specifically, we estimate the following regression equation to explore the relationship between the within-firm product rank of sales and firm-product-level productivity and quality:⁴⁴

$$\text{Log product rank}_{jnt} = \alpha_{\tilde{\omega}} \tilde{\omega}_{jnt} + \alpha_{\tilde{\xi}} \tilde{\xi}_{jnt} + d_{jn} + d_{jt} + d_{nt} + \epsilon_{jnt}, \quad (23)$$

where the product rank (in logarithm) is defined based on the sales of products within firm-year pairs.⁴⁵ The rank of the top product (i.e., with the largest sales) is 1. An increase in rank indicates a product further away from the core competency of a firm. We include d_{jn} , d_{jt} , and d_{nt} as firm-product, firm-year, and product-year fixed effects, respectively, to capture different characteristics other than productivity and quality that vary at these levels and influence the product rank.

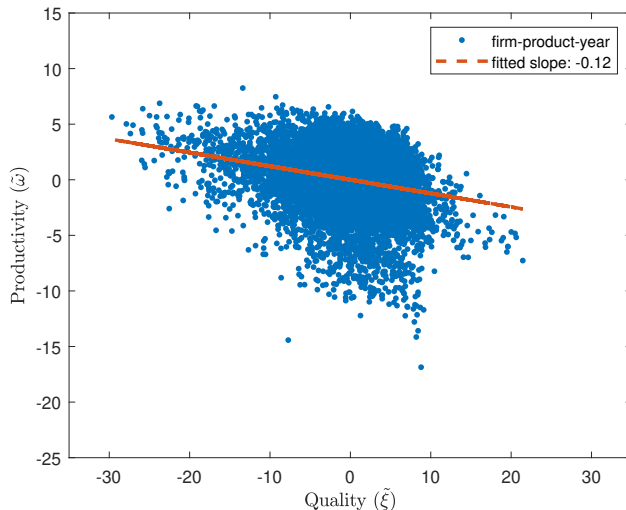
We find that products closer to firms' core competence (i.e., with a lower rank value) have both higher productivity and quality. An increase of 1 percent in productivity and quality moves the rank of the product up by 0.602 percent and 0.170 percent, respectively,

⁴⁴In an unreported result, a similar regression equation is estimated using the growth of sales (instead of the level of sales) within a firm to measure the product rank. The regression result shows a similar pattern.

⁴⁵Equation (23) examines the relationship between product rank on the one hand and productivity and quality on the other by using product rank as the dependent variable. The purpose of the regression is to study directly the importance of productivity and quality for differences in sales across products within firms. In fact, regression equation (23) is predicted by the implication of our model. To see this, summing (12) over products produced by the firm and using (3), we obtain: $\ln s_{jnt} = (\eta_n - 1)\tilde{\omega}_{jnt} + \tilde{\xi}_{jnt} + (\eta_n - 1)\ln\left(\frac{\eta_n - 1}{\eta_n}\right) + \frac{1}{\eta_n - 1}\ln\left(\frac{F(L_{jt}, M_{jt}, K_{jt})}{(\sum_{n \in \Lambda_{jt}} \frac{(\eta_n - 1)}{\eta_n} R_{jnt})(\sum_{n \in \Lambda_{jt}} R_{jnt})^{\eta_n - 1}}\right)$, where s_{jnt} is the within-firm share of product n in period t .

conditional on all other factors. This is consistent with the literature that has theoretically postulated cost (i.e., productivity) or demand (i.e., quality) as key determinants of such within-firm variation in sales (e.g., [Berman et al., 2012](#); [Chatterjee et al., 2013](#); [Mayer et al., 2014, 2021](#); [Eckel et al., 2015](#); [Arkolakis et al., 2021](#)). Our results provide empirical support for both of these hypotheses.

Figure 1: The relationship between productivity and quality



As both productivity and quality influence the intra-firm performance of a product, it is natural to ask whether and how these different dimensions of within-firm heterogeneity are related. As a starting point, Figure 1 presents the raw relationship between our two key estimated measures of heterogeneity, i.e., (quantity-based) productivity ($\tilde{\omega}_{jnt}$) and quality ($\tilde{\xi}_{jnt}$).⁴⁶ This raw correlation is negative, consistent with the emerging literature (e.g., [Grieco and McDevitt, 2017](#); [Orr, 2022](#); [Li et al., 2023](#)) highlighting the trade-off between these two dimensions of firm heterogeneity. This empirical pattern suggests that producing higher-quality products increases the marginal cost of production by decreasing output quantity per unit of inputs, which in turn reduces (quantity-based) productivity. This is a relationship between productivity and quality that we allow for in (6) of the model but do not impose in our structural estimation.

To formally quantify the trade-off between the two dimensions within firms, we propose

⁴⁶When we tease out the fixed effects at firm-product, firm-year, and product-year levels from $\tilde{\xi}_{jnt}$ to obtain a finer measure of quality (i.e., ξ_{jnt}) as defined in Section 2, the correlation is also negative. The firm-product fixed effects may contain parts of quality that only vary at the firm-product level. The correlation is robustly negative when we include the firm-product fixed effects as a part of the quality measure.

to estimate a linear version of (6):⁴⁷

$$\tilde{\omega}_{jnt} = \omega_{jnt} - \gamma_{\xi} \tilde{\xi}_{jnt}, \quad (24)$$

where $\gamma_{\xi} \tilde{\xi}_{jnt}$ is interpreted as the cost (in terms of lowering productivity or raising marginal cost) of increasing quality, holding inputs fixed. γ_{ξ} measures the elasticity of productivity with respect to the change in quality. We refer to it as the **cost responsiveness** of quality.

Nonetheless, estimating the cost responsiveness of quality (γ_{ξ}) is challenging. Technical efficiency (ω_{jnt}) is correlated with the quality choice (ξ_{jnt}) if firms choose to produce different quality products based on technical efficiency. To address this challenge, we exploit a simplified version of the evolution of ω_{jnt} following equation (7):

$$\omega_{jnt} = g_1 \omega_{jnt-1} + d_{jt} + d_{nt} + \epsilon_{jnt}, \quad (25)$$

where d_{nt} and d_{jt} are product-year and firm-year fixed effects, and ϵ_{jnt} is an i.i.d. innovation shock.

Replacing technical efficiency in (25) by that in (24) gives:

$$\tilde{\omega}_{jnt} = g_1 \tilde{\omega}_{jnt-1} - \gamma_{\xi} \tilde{\xi}_{jnt} + g_1 \gamma_{\xi} \tilde{\xi}_{jnt-1} + d_{jt} + d_{nt} + \epsilon_{jnt}. \quad (26)$$

Although all variables (except ϵ_{jnt}) are already estimated from our structural model, the innovation shock ϵ_{jnt} can be correlated with contemporaneous quality choice $\tilde{\xi}_{jnt}$. To address such an endogeneity problem, we estimate (26) via GMM using a set of instrumental variables that includes the average productivity and average quality of products that are produced by other firms in period $t - 2$. These variables are uncorrelated with the innovation term ϵ_{jnt} which is an i.i.d. shock.

The estimation results are presented in Table 5. As expected, technical efficiency is highly persistent. More importantly, we find a negative trade-off between productivity and quality at the firm-product level across various specifications of (26). Column (1) reports the estimated coefficients of (26) including only product-year fixed effects, while the estimation in Column (2) also includes firm-year fixed effects. The comparison of the coefficient estimates in Columns (1) and (2) suggests that it is important to control for unobserved firm-year

⁴⁷Our estimated measure of quality, $\tilde{\xi}_{jnt}$, is derived as the residual from the demand function (2). Consequently, its variation across products, firms, and over time may be influenced by factors such as demand conditions (e.g., macroeconomic conditions and market size), firm-brand image, product measurement units (e.g., grams vs. liters), and firm-time-specific measurement errors, as discussed in Section 2.1. To isolate the actual impact of quality (ξ_{jnt}) from these potential confounders and control for unobserved product and firm characteristics, we include product-year and firm-year fixed effects in the analysis.

fixed effects to minimize the potential selection bias despite the use of valid instrumental variables. According to Column (2), a 1-percent increase in quality lowers productivity (and thus increases marginal cost) by 0.234 percent, holding all other variables fixed. Such an estimate of cost responsiveness of quality is consistent with other analyses using different approaches in various industries and countries.⁴⁸ Importantly, this also suggests that it is necessary to control for quality differences as a part of the evolution in the existing estimation methods (e.g., Orr, 2022; Valmari, 2022) that utilize the evolution process of productivity in the estimation of production parameters.

To explore a potential source of the cost of quality, we allow the cost responsiveness to vary by product age. Here, product age is defined for each product-firm pair as the number of years since the product variety first appeared in the sample. Specifically, we extend (24) by introducing an interaction term between quality and the logarithm of product age, alongside the logarithm of product age itself.⁴⁹ The result is presented in Column (3) of Table 5. The negative coefficient for $\gamma_{a\xi}$ indicates that the trade-off between productivity and quality diminishes as a firm continues to produce a specific product over time. This suggests that firms with extensive experience in manufacturing a particular product develop better production management capabilities, allowing them to achieve higher quality without sacrificing efficiency. Based on the estimate in Column (3), a straightforward calculation shows that five years of experience in product manufacturing results in approximately a 6.7 percentage point reduction in the impact of quality on productivity, translating to a 26.2 percent decrease in the overall effect of quality on productivity.⁵⁰

Finally, while we adopt a simple evolution model for technical efficiency – commonly represented in the literature as an AR(1) process (25) – our estimation approach offers a unique

⁴⁸For example, Jaumandreu and Yin (2014) find strong negative correlations (ranging from -0.99 to -0.59, by industry) between their measures of cost advantage and demand advantage of exporters in the Chinese manufacturing industries. Grieco and McDevitt (2017) show that reducing a healthcare center’s quality standards can increase its patient load, and they document a quality-quantity (number of patients) trade-off elasticity of -0.2 in the dialysis industry in the United States. Atkin et al. (2019) find that firms that make lower quality rugs produce more quickly among rug-makers in Egypt, demonstrating a reverse correlation between quantity productivity and quality productivity with an elasticity of -0.40. Orr (2022) estimates firm-product level measures of productivity and “product appeal” from the Indian machinery manufacturing industry and finds a negative correlation of about -0.28 between them. Using an objective output quality measure, Li et al. (2023) find that about half of the benefit created by quality is offset by the cost of producing the quality in the Chinese steel-making industry. Forlani et al. (2023) document an even stronger negative correlation (about -0.9) between demand and quantity-based productivity at the firm level in various Belgian industries, suggesting a trade-off between the quality of a firm’s products and their production cost.

⁴⁹That is, (24) becomes: $\tilde{\omega}_{jnt} = \omega_{jnt} - (\gamma_{\xi} + \gamma_{a\xi}age_{jnt})\tilde{\xi}_{jnt} + \gamma_a age_{jnt}$, where age_{jnt} is the logarithm of product age. Consequently, the estimating equation (26) becomes $\tilde{\omega}_{jnt} = g_1\tilde{\omega}_{jnt-1} - \gamma_{\xi}\tilde{\xi}_{jnt} + g_1\gamma_{\xi}\tilde{\xi}_{jnt-1} + \gamma_a age_{jnt} - g_1\gamma_a age_{jnt-1} - \gamma_{a\xi}\tilde{\xi}_{jnt}age_{jnt} + g_1\gamma_{a\xi}\tilde{\xi}_{jnt-1}age_{jnt-1} + d_{jt} + d_{nt} + \epsilon_{jnt}$.

⁵⁰The calculation of the level of the impact is: $0.086 \times (\log(5+1) - \log(1)) = 0.067$. Relative to the overall impact of quality on productivity, the calculation is: $\frac{0.086 \times (\log(5+1) - \log(1))}{0.255} = 0.262$.

Table 5: Cost of quality

Dep. var.: Productivity	(1)	(2)	(3)	(4)
g_1	0.867*** (0.028)	0.980*** (0.045)	0.888*** (0.042)	0.818*** (0.064)
γ_ξ	0.160*** (0.057)	0.234*** (0.074)	0.255*** (0.087)	0.151*** (0.056)
γ_a			-2.584 (1.903)	
$\gamma_{a\xi}$			-0.086*** (0.032)	
g_s				0.049 (0.073)
Product-Year FE	yes	yes	yes	yes
Firm-Year FE	no	yes	yes	no
Observations	7122	7122	7122	7122

Note: The dependent variable is quantity-based productivity at the firm-product-year level. The coefficients are estimated via GMM. The instrument set includes lagged productivity, lagged quality, twice-lagged quality, twice-lagged average productivity of the same product produced by other firms and twice-lagged average quality of the same product produced by other firms in all specifications. The instrument set in Column (3) also includes log product age, lagged log product age, the interaction between lagged quality and lagged log product age and the interaction between twice-lagged quality and twice-lagged log product age. The instrument set in Column (4) also includes lagged productivity and lagged quality of the top ranked product of each firm-year pair. Bootstrapped standard errors clustered at the firm level and stratified by industry and scope are shown in parentheses (100 repetitions). *** $p < 0.01$, ** $p < 0.05$.

advantage over existing methods. Specifically, it allows for the straightforward exploration of more complex (e.g., interdependent) structures in the evolution of firm-product-level technical efficiency. This advantage stems from our approach to estimating production parameters and productivity, which does not rely on the time-series relationship of productivity in the estimation process unlike some of the existing methods (e.g., [Orr, 2022](#); [Valmari, 2022](#)). Consequently, this allows us to examine the time-series dynamics of productivity *after* the estimation. As a demonstration of such an advantage, we investigate potential spillovers of technical efficiency from a firm's top-ranked product to its other products. That is, we extend (25) by adding a term, $g_s \omega_{jt-1}^*$, to the right-hand side. Here, ω_{jt-1}^* represents the technical efficiency of the top-ranked product within firm j in period $t - 1$, and the coefficient g_s measures the strength of the spillover effect on the technical efficiency of product n within the

same firm.⁵¹ We estimate this specification report the result in Column (4) of Table 5. We find positive (although not statistically significant) spillover effects from the top-ranked product within a firm. Importantly, the cost-responsiveness coefficient, γ_ξ , remains qualitatively similar to the other specifications.

Overall, the results obtained above demonstrate a robust negative relationship between quality and quantity-based productivity. However, when considering quality and quality-adjusted productivity (ATFP) which takes into account both the costs and benefits of quality as indicated by its definition in (22), a significantly positive relationship emerges. Across all three industries, the correlation coefficient between ATFP and quality at the firm-product level is 0.44. This positive relationship is intuitive. While the cost of quality tends to lower ATFP as quality increases, the benefits of quality contribute to a positive association with ATFP. The dominance of the latter force results in an overall positive relationship between ATFP and quality. This finding aligns with previous analyses that emphasize firms with high production capability choose to produce high-quality output endogenously (e.g., Verhoogen, 2008; Kugler and Verhoogen, 2009, 2012; Feenstra and Romalis, 2014; Hottman et al., 2016; Fan et al., 2018). Our results not only highlight the positive sorting within firms but also indicate that it is conditional upon acknowledging both the increasing cost and benefit of producing higher-quality products. This observation is consistent with the findings of Li et al. (2023), who utilize a firm-level objective quality measure from the Chinese steel industry.

In sum, our analysis on productivity and quality highlights the significance of considering the cost of quality and the relationship between these variables at the firm-product level. A notable implication arises from the relationship: reducing the cost of quality (e.g., through long experience in production) not only contributes directly to an increase in the ATFP of a firm but also indirectly stimulates growth through intra-firm resource reallocation towards the production of higher-quality products, which subsequently enhances the firm's ATFP further. In the following section, we shift our focus to evaluating the cost of quality and study the role of product scope in firm growth through intra-firm resource reallocation resulting from a reduction in the cost of quality.

7 How Costly is Quality?

The results regarding the cost of quality are meaningful because they imply that a reduction in the cost responsiveness of quality can lead to growth in ATFP. Intuitively, conditional on the underlying technical efficiency (i.e., ω) and product quality, a reduction of the cost

⁵¹As a result, the estimating equation (26) becomes $\tilde{\omega}_{jnt} = g_1\tilde{\omega}_{jnt-1} - \gamma_\xi\tilde{\xi}_{jnt} + g_1\gamma_\xi\tilde{\xi}_{jnt-1} + g_s\tilde{\omega}_{jt-1}^* + g_s\gamma_\xi\tilde{\xi}_{jt-1}^* + d_{nt} + \epsilon_{jnt}$, where $\tilde{\omega}_{jt-1}^*$ and $\tilde{\xi}_{jt-1}^*$ are the productivity and quality of the top-ranked product, respectively. The firm-year fixed effects are not included because the top-ranked productivity and quality vary at the firm-year level.

responsiveness of quality (i.e., γ_ξ) means a direct increase in quantity-based productivity (i.e., $\tilde{\omega}$) according to (24) and, thus, a corresponding increase in ATFP as defined in (22). More importantly, the impact on higher-quality products is larger for a given reduction in the cost responsiveness of quality. Thus, in the short run, multi-product firms can endogenously reallocate resources towards high-quality and high-productivity products, which consequently improves ATFP at the firm level.⁵²

We focus on the short-term effects of reducing the cost responsiveness of quality while keeping quality choices fixed.⁵³ We emphasize the role of product scope in driving productivity gains through resource reallocation within firms.

To explore this, we conduct a counterfactual exercise by reducing the cost responsiveness of quality and comparing the resulting ATFP at the firm level with the baseline scenario (i.e., without a reduction in the cost of quality). The improvement in ATFP is then decomposed into a direct increase due to the reduced cost of quality and the gains due to the intra-firm reallocation of resources.

Specifically, in the counterfactual scenario, we reduce the cost responsiveness of quality (γ_ξ) by 1 percent for all firm-product pairs. This leads to a direct improvement in quantity-based productivity: $\tilde{\omega}'_{jnt} = \tilde{\omega}_{jnt} + 0.01 \times \gamma_\xi \tilde{\xi}_{jnt}$, where $\tilde{\omega}'_{jnt}$ is the counterfactual productivity and $\tilde{\omega}_{jnt}$ and $\tilde{\xi}_{jnt}$ are the baseline quantity-based productivity and quality, respectively. Here γ_ξ denotes the estimated cost responsiveness of quality specific to each industry and reported in Column (3) of Table 5. The direct improvement in quantity-based productivity, $0.01 \times \gamma_\xi \tilde{\xi}_{jnt}$, drives an increase in ATFP at the firm-product level according to (22).⁵⁴

More interestingly, there is an indirect improvement in firm-level ATFP due to intra-firm resource reallocation across products for multi-product firms. To see this mechanism, note that the 1-percent decline in γ_ξ leads to a differential improvement in the counterfactual productivity across products within firms, depending on the baseline quality level ($\tilde{\xi}_{jnt}$). For a product with higher quality, the resulting productivity improvement due to the reduction in the cost responsiveness of quality is larger. As a result, firms can react to the differential productivity improvement by re-optimizing their intra-firm allocation of inputs and outputs. Because ATFP and quality are positively related as documented in Section 6, multi-product

⁵²In addition, considering that product quality is endogenously chosen by firms based on productivity as emphasized in the literature (e.g., Verhoogen, 2008; Kugler and Verhoogen, 2009, 2012; Feenstra and Romalis, 2014; Fan et al., 2018), a cost of quality reduction implies an incentive for quality upgrading, thus increasing ATFP even further in the long run. Our static empirical model does not capture the long-term endogenous reaction of quality choices.

⁵³As a result, our evaluation of the cost of quality should be seen as a lower bound of the actual impact on firm performance.

⁵⁴Throughout the analysis, we treat all the dynamic decisions (i.e., product quality, scope, and investment) described in Online Appendix B as fixed.

firms tend to reallocate more production resources to products with higher ATFP and higher quality. Consequently, this reallocation leads to an indirect improvement in firm-level ATFP.

Both the direct and indirect improvements contribute to the increase in firm-level ATFP. To understand their magnitude and relative importance, we aggregate firm-level ATFP from firm-product-level ATFP using sales as weights. We apply the within-industry across-firm decomposition proposed by [Olley and Pakes \(1996\)](#) to compute the intra-firm decomposition. That is, for each firm j in period t ,

$$\text{ATFP}_{jt} = \overline{\text{ATFP}}_{jt} + \sum_{n \in \Lambda_{jt}} (s_{jnt} - \bar{s}_{jt})(\text{ATFP}_{jnt} - \overline{\text{ATFP}}_{jt}), \quad (27)$$

where $\overline{\text{ATFP}}_{jt}$ is the simple average of the exponent of quality-adjusted productivity, ATFP_{jnt} , across products produced by the same firm. s_{jnt} is the within-firm sales share of product n by firm j in period t . \bar{s}_{jt} is the simple average of the sales shares (that is, the inverse of the product scope). Intuitively, an increase in firm-level ATFP can be caused by an increase in ATFP of all products as well as a reallocation of resources towards more productive products. Accordingly, intra-firm resource reallocation is defined as the difference of the covariance term (the second term on the right-hand side) in (27) between the counterfactual scenario and the baseline scenario. To obtain the overall improvement in ATFP at the industrial level, we aggregate firm-level ATFP improvement using firms' total sales as weights. The relative contribution of intra-firm resource reallocation to the firm-level ATFP improvement is aggregated to the industry level in the same way.

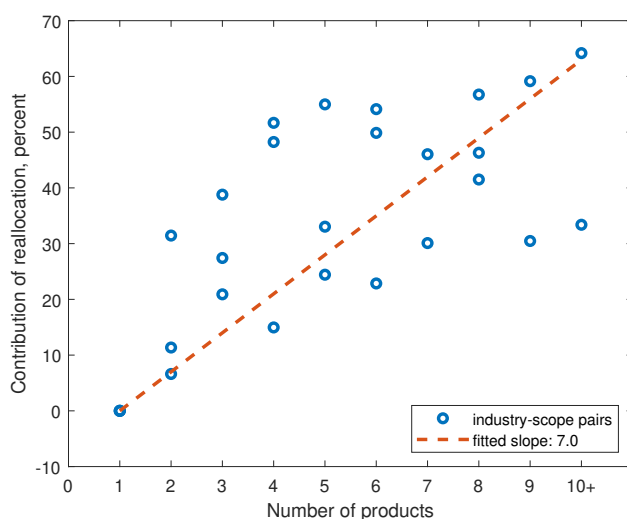
Table 6: Impact of 1-percent reduction in cost responsiveness of quality on ATFP

Industry	All firms				MPFs
	All	Footwear	Printing	Pharmaceutical	All
Total improvement, percent	2.836 (0.097)	1.191 (0.375)	3.512 (0.833)	2.896 (0.109)	2.918 (0.103)
Intra-firm reallocation, percent	0.854 (0.251)	0.120 (0.043)	0.569 (0.177)	0.966 (0.286)	0.997 (0.286)
percentage relative to total	30.1 (4.0)	10.0 (2.2)	16.2 (3.8)	33.4 (5.1)	34.2 (4.6)

Note: The improvement in ATFP at the industry level is measured in percentage and calculated as the weighted average of the improvements in ATFP at the firm-year level with firms' total sales in the baseline scenario as weights. MPFs refers to multi-product firms only. Bootstrapped standard errors clustered at the firm level and stratified by industry and scope are shown in parentheses (100 repetitions).

Table 6 reports the overall improvement in firm-level ATFP as well as the contribution from the intra-firm resource reallocation of multi-product firms in the three industries. A 1-percent decline in the cost responsiveness of quality leads to an improvement in ATFP by approximately 1.2, 3.5, and 2.9 percent for the footwear, printing, and pharmaceutical industries, respectively. This is a sizable magnitude. More importantly, the contribution of the within-firm resource reallocation accounts for roughly 10 percent to 33 percent of the overall improvement in ATFP across the three industries. This is essentially a lower bound of the contribution because the calculation is based on all firms including the single-product firms that experience, by definition, zero within-firm reallocation. When focusing on multi-product firms only, the contribution is on average approximately 34 percent across the three industries. This result establishes the economic significance of the cost of quality within multi-product firms as a channel impacting overall quality-adjusted productivity.

Figure 2: Contribution of within-firm resource reallocation to ATFP growth



Notes: All firms producing more than 10 products are clustered in the “10+” group.

A large literature on resource reallocation focuses on across-firm analysis and shows that much of the aggregate productivity growth is attributable to the resource reallocation towards more productive firms (e.g., Baily et al., 1992; Bartelsman and Doms, 2000; Baily et al., 2001; Aw et al., 2001; Foster et al., 2006, 2008; Syverson, 2011; Collard-Wexler and De Loecker, 2015). Complementary to the literature, our firm-product-level analysis shows that the contribution of within-firm resource reallocation is also sizable. Interestingly, compared to the footwear industry, the relatively higher intra-firm contribution in the printing and pharmaceutical industries is consistent with the relatively larger number of products in these industries. Indeed, as shown in Table 2, firms in the printing and pharmaceutical industries produce 3.7 and 6.9 products on average, respectively, while firms in the footwear industry

produce 1.3 products. Intuitively, a larger product scope allows for a greater potential to reallocate resources across products.

To unpack such a heterogeneous pattern, we group firms by the number of products produced and compute the sales-weighted average contribution of intra-firm resource reallocation to firm-level ATFP improvement (due to the reduction in the cost of quality). This computation is conducted for each industry. We plot the relationship between product scope and the contribution of intra-firm reallocation (in percentage) in Figure 2.⁵⁵ Each dot represents the average contribution of within-firm reallocation by product scope and industry. The dashed line represents the fitted line obtained from a simple OLS regression of within-firm reallocation against product scope. The upward-sloping fitted line establishes that, on average, the role of within-firm reallocation increases in firms with a larger scope with more room for within-firm adjustment. The slope of the fitted line suggests that producing one more product can increase the contribution of within-firm reallocation in improving ATFP by 7 percent. In sum, our results highlight that multi-product firms with larger scope experience larger productivity gains when the cost of quality is lower. This reveals a new mechanism for enhancing the performance of multi-product firms.

8 Conclusion

Multi-product firms account for a significant share of our economy. Yet, the traditional firm-level analysis in the literature masks the intra-firm heterogeneity. In this paper, we propose a novel method to estimate firm-product-level productivity and quality along with demand and transformation function parameters. Compared with the existing methods in the literature, our method does not impose assumptions on how inputs are allocated across the production of different products within firms, nor does it restrict how productivity evolves over time. This flexibility allows researchers to explore complex productivity dynamics after estimation. Importantly, the method can be easily scaled up to estimate production functions with a large number of products, without relying on the availability of productivity proxies. Finally, the method accounts for heterogeneous intermediate input prices that are usually unobservable to researchers and lead to biased estimation results when ignored.

We apply our method to three major industries in the Mexican manufacturing sector. We find that both quality (demand) and productivity play significant roles in explaining intra-firm revenue heterogeneity. However, firms face a trade-off between upgrading quality and productivity. After taking both the costs and benefits of quality into account, quality-adjusted productivity shows a strong positive intra-firm correlation with quality.

To understand how costly quality is for productivity growth and intra-firm resource

⁵⁵The relationship is similar when the figure is plotted by industry.

allocation, we conduct a counterfactual exercise where we reduce the cost responsiveness of quality by 1 percent. The reduction leads to substantial productivity gains, especially for multi-product firms. A sizable portion of the productivity gain of multi-product firms is due to the within-firm reallocation of resources towards more-productive and higher-quality products. In particular, we show that a larger product scope allows more room for intra-firm resource reallocation, leading to a higher productivity gain when there is a reduction in the cost of quality. This result establishes the quantitative significance of intra-firm resource reallocation in enhancing the performance of multi-product firms that dominate manufacturing production. This channel, thus, has strong potential implications for aggregate productivity growth.

References

- Akerberg, D. A., K. Caves, and G. Frazer (2015, November). Identification properties of recent production function estimators. *Econometrica* 83(6), 2411–2451.
- Arkolakis, C., S. Ganapati, and M.-A. Muendler (2021). The extensive margin of exporting products: A firm-level analysis. *American Economic Journal: Macroeconomics* 13(4), 182–245.
- Atalay, E. (2014, June). Materials prices and productivity. *Journal of the European Economic Association* 12(3), 575–611.
- Atkin, D., A. K. Khandelwal, and A. Osman (2019). Measuring productivity: Lessons from tailored surveys and productivity benchmarking. In *AEA Papers and Proceedings*, Volume 109, pp. 444–49.
- Aw, B. Y., X. Chen, and M. J. Roberts (2001). Firm-level evidence on productivity differentials and turnover in taiwanese manufacturing. *Journal of Development Economics* 66(1), 51–86.
- Aw, B. Y., M. Roberts, and D. Y. Xu (2011). R&d investment, exporting, and productivity dynamics. *American Economic Review* 101, 1312–1344.
- Baily, M. N., E. J. Bartelsman, and J. Haltiwanger (2001). Labor productivity: structural change and cyclical dynamics. *Review of Economics and Statistics* 83(3), 420–433.
- Baily, M. N., C. Hulten, D. Campbell, et al. (1992). Productivity dynamics in manufacturing plants. *Brookings Papers on Economic Activity* 23(1992 Microeconomics), 187–267.
- Bartelsman, E. J. and M. Doms (2000). Understanding productivity: Lessons from longitudinal microdata. *Journal of Economic Literature* 38(3), 569–594.
- Berkowitz, D., H. Ma, and S. Nishioka (2017, 10). Recasting the Iron Rice Bowl: The Evolution of China’s State Owned Enterprises. *The Review of Economics and Statistics* 99(4), 735–747.
- Berman, N., P. Martin, and T. Mayer (2012). How do different exporters react to exchange rate changes? *The Quarterly Journal of Economics* 127(1), 437–492.
- Berry, S. (1994, Summer). Estimating discrete-choice models of product differentiation. *RAND Journal of Economics* 25(2), 242–262.
- Berry, S., J. Levinsohn, and A. Pakes (1995, July). Automobile prices in market equilibrium. *Econometrica* 63(4), 841–890.
- 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* 121, 1–14.
- Cairncross, J., P. Morrow, S. Orr, and S. Rachapalli (2023). Multi-product markups.
- Caselli, M. (2018, August). Do all imports matter for productivity? Intermediate inputs vs capital goods. *Economia Politica: Journal of Analytical and Institutional Economics* 35(2), 285–311.
- Caselli, M., A. Chatterjee, and A. Woodland (2017). Multi-product Exporters, Variable Markups and Exchange Rate Fluctuations. *Canadian Journal of Economics* 50(4), 1130–1160.
- Caselli, M., L. Nesta, and S. Schiavo (2021). Imports and labour market imperfections: Firm-level evidence from France. *European Economic Review* 131, 103632.
- Chatterjee, A., R. Dix-Carneiro, and J. Vichyanond (2013). Multi-product firms and exchange rate fluctuations. *American Economic Journal: Economic Policy* 5(2), 77–110.
- Chen, Y., M. Igami, M. Sawada, and M. Xiao (2021). Privatization and productivity in china. *The RAND Journal of Economics* 52(4), 884–916.
- Chen, Z. and M. Liao (2022). Production Function Estimation for Multi-Product Firms. Working Paper 3968432, SSRN.
- Collard-Wexler, A. and J. De Loecker (2015). Reallocation and technology: Evidence from the us steel industry. *American Economic Review* 105(1), 131–71.
- Das, S., M. J. Roberts, and J. R. Tybout (2007, May). Market entry costs, producer heterogeneity and export dynamics. *Econometrica* 75(3), 837–873.
- De Loecker, J. (2011). Product differentiation, multi-product firms and estimating the impact of trade liberalization on productivity. *Econometrica* Vol. 79, No. 5, pp. 1407–1451.
- De Loecker, J., P. K. Goldberg, A. K. Khandelwal, and N. Pavcnik (2016). Prices, markups, and trade reform. *Econometrica* 84(2), 445–510.
- De Loecker, J. and F. Warzynski (2012). Markups and firm-level export status. *American Economic Review* 102(6), 2437–71.
- Demirer, M. (2022). Production Function Estimation with Factor-Augmenting Technology: An Application to Markups. mimeo.
- Dhyne, E., A. Petrin, V. Smeets, and F. Warzynski (2022). Theory for extending single-product production function estimation to multi-product settings. Nber working paper, National Bureau of Economic Research.
- Diewert, E., K. J. Fox, and L. Ivancic (2009). Scanner Data, Time Aggregation and the Construction of Price Indexes. UBC Discussion Paper 09-09, Department of Economics, University of British Columbia.
- Doraszelski, U. and J. Jaumandreu (2013). R&d and productivity: Estimating endogenous productivity. *Review of Economic Studies* 80, 1338 – 1383.
- Dubois, P. and L. Lasio (2018). Identifying industry margins with price constraints: Structural estimation on pharmaceuticals. *American Economic Review* 108(12), 3685–3724.
- Eckel, C., L. Iacovone, B. Javorcik, and J. P. Neary (2015). Multi-product firms at home and away: Cost-versus quality-based competence. *Journal of International Economics* 95(2), 216–232.
- Ericson, R. and A. Pakes (1995, January). Markov-perfect industry dynamics: A framework for empirical work. *Review of Economic Studies* 62(1), 53–82.
- Eslava, M., J. Haltiwanger, and N. Urdaneta (2023). The Size and Life-Cycle Growth of Plants: The Role of Productivity, Demand, and Wedges. *The Review of Economic*

Studies forthcoming.

- Fan, H., Y. A. Li, and S. R. Yeaple (2018). On the relationship between quality and productivity: Evidence from china's accession to the wto. *Journal of International Economics* 110, 28–49.
- Feenstra, R. C. and J. Romalis (2014, May). International prices and endogenous quality. *The Quarterly Journal of Economics* 129(2), 477–527.
- Forlani, E., R. Martin, G. Mion, and M. Muùls (2023, 04). Unraveling Firms: Demand, Productivity and Markups Heterogeneity. *The Economic Journal forthcoming*.
- Foster, L., J. Haltiwanger, and C. J. Krizan (2006). Market selection, reallocation, and restructuring in the us retail trade sector in the 1990s. *The Review of Economics and Statistics* 88(4), 748–758.
- Foster, L., J. Haltiwanger, and C. Syverson (2008). Reallocation, firm turnover, and efficiency: Selection on productivity or profitability? *American Economic Review* 98(1), 394–425.
- Gandhi, A., S. Navarro, and D. A. Rivers (2020). On the identification of gross output production functions. *Journal of Political Economy* 128(8), 2973–3016.
- Grieco, P., S. Li, and H. Zhang (2016, May). Production function estimation with unobserved input price dispersion. *International Economic Review* 57(2), 665–690.
- Grieco, P., S. Li, and H. Zhang (2022). Input Prices, Productivity and Trade Dynamics: Long-run Effects of Liberalization on Chinese Paint Manufacturers. *The RAND Journal of Economics* 53(3), 516–560.
- Grieco, P. L. and R. C. McDevitt (2017). Productivity and quality in health care: Evidence from the dialysis industry. *The Review of Economic Studies* 84(3), 1071–1105.
- Harrigan, J., A. Reshef, and F. Toubal (2021, February). Techies, Trade, and Skill-Biased Productivity. CEPR Discussion Papers 15815, C.E.P.R. Discussion Papers.
- Hopenhayn, H. A. (1992, September). Entry, Exit, and Firm Dynamics in Long Run Equilibrium. *Econometrica* 60(5), 1127–1150.
- Hottman, C. J., S. J. Redding, and D. E. Weinstein (2016). Quantifying the sources of firm heterogeneity. *The Quarterly Journal of Economics* 131(3), 1291–1364.
- Jaumandreu, J. and H. Yin (2014). Cost and product advantages: A firm-level model for the chinese exports and industry growth. working paper, Boston College.
- Jovanovic, B. (1982, May). Selection and the Evolution of Industry. *Econometrica* 50(3), 649–670.
- Khandelwal, A. K. (2010). The long and short (of) quality ladders. *Review of Economic Studies* 77, 1450–1476.
- Kirov, I. and J. Traina (2023). Labor Market Power and Technological Change in US Manufacturing. mimeo.
- Klump, R. and O. de La Grandville (2000). Economic growth and the elasticity of substitution: Two theorems and some suggestions. *American Economic Review* 90(1), 282–291.
- Kugler, M. and E. Verhoogen (2009). Plants and imported inputs: New facts and an interpretation. *American Economic Review* 99(2), 501–07.
- Kugler, M. and E. Verhoogen (2012). Prices, plant size, and product quality. *Review of Economic Studies* 79(1), 307–339.
- Kumar, P. and H. Zhang (2019). Productivity or unexpected demand shocks: What determines firms' investment and exit decisions? *International Economic Review* 60(1), 303–327.
- León-Ledesma, M. A., P. McAdam, and A. Willman (2010). Identifying the elasticity of

- substitution with biased technical change. *American Economic Review* 100(4)(4), 1330–1357.
- Levinsohn, J. and A. Petrin (2003). Estimating production functions using inputs to control for unobservables. *The Review of Economic Studies* 70(2), 317–341.
- Li, J., S. Li, and H. Zhang (2023). Output Quality, Productivity, and Demand Advantage: Evidence from the Chinese Steel Industry. Unsw working paper, University of New South Wales.
- Li, S. (2018). A structural model of productivity, uncertain demand, and export dynamics. *Journal of International Economics* 115, 1–15.
- Li, S. and H. Zhang (2022, February). Does External Monitoring from the Government Improve the Performance of State-Owned Enterprises? *The Economic Journal* 132(642), 675–708.
- Mayer, T., M. J. Melitz, and G. I. Ottaviano (2014). Market size, competition, and the product mix of exporters. *American Economic Review* 104(2), 495–536.
- Mayer, T., M. J. Melitz, and G. I. Ottaviano (2021). Product mix and firm productivity responses to trade competition. *Review of Economics and Statistics* 103(5), 874–891.
- Melitz, M. J. (2000). Estimating firm-level productivity in differentiated product industries. unpublished paper.
- Melitz, M. J. (2003, November). The impact of trade on intra-industry reallocations and aggregate industry productivity. *Econometrica* 71(6), 1695–1725.
- Morlacco, M. (2020). Market Power in Input Markets: Theory and Evidence from French Manufacturing. mimeo.
- Olley, G. S. and A. Pakes (1996). The dynamics of productivity in the telecommunications equipment industry. *Econometrica* 64(6), 1263–1297.
- Ornaghi, C. (2006). Assessing the effects of measurement errors on the estimation of production functions. *Journal of Applied Econometrics* 21(6), 879–891.
- Orr, S. (2022). Within-firm productivity dispersion: Estimates and implications. *Journal of Political Economy* 130(11), 000–000.
- Panzar, J. C. and R. D. Willig (1977). Economies of scale in multi-output production. *The Quarterly Journal of Economics* 91(3), 481–493.
- Panzar, J. C. and R. D. Willig (1981). Economies of scope. *The American Economic Review* 71(2), 268–272.
- Pozzi, A. and F. Schivardi (2016, July). Demand or productivity: What determines firm growth? *RAND Journal of Economics* 47(3), 608–630.
- Raval, D. (2023, 01). Testing the Production Approach to Markup Estimation. *The Review of Economic Studies* 90(5), 2592–2611.
- Roberts, M., D. Y. Xu, X. Fan, and S. Zhang (2018, 11). The Role of Firm Factors in Demand, Cost, and Export Market Selection for Chinese Footwear Producers. *Review of Economic Studies* 85(4), 2429–2461.
- Syverson, C. (2011). What determines productivity? *Journal of Economic Literature* 49(2), 326–65.
- Valmari, N. (2022). Estimating production functions of multiproduct firms. *Review of Economic Studies* 130(11), 000–000.
- Verhoogen, E. A. (2008). Trade, quality upgrading, and wage inequality in the mexican manufacturing sector. *The Quarterly Journal of Economics* 123(2), 489–530.

Online Appendix

A Extension to More General Demand and Transformation Functions

In the main text of the paper, we have assumed specific forms for the demand and production transformation functions. In this appendix, we demonstrate how to extend the estimation method to accommodate general demand and production transformation functions, whose parameters can be identified and estimated using appropriate data. Specifically, we generalize the method in two directions. On the demand side, we generalize the demand for a product to account for influences from the sales of other products within the same firm (i.e., cannibalization) as well as from products of competing firms (i.e., competition). On the production side, we extend the linear aggregator of output to a nonlinear form, allowing for flexible substitution or complementarity across products produced by the same firm. In addition, the model is readily extendable to allow for flexible transformation functions whose parameters vary by product permutation, as described in Section 2.2.

A.1 Demand and Transformation Functions

We start from describing the demand function. The demand for product n of firm j in period t is modelled as a general inverse demand function:

$$P_{jnt} = P_{jnt}(\mathbf{Q}_{jt}, \mathbf{Q}_{-jt}; \boldsymbol{\xi}_t), \quad (\text{A.1})$$

where P_{jnt} is the product price. Importantly, $\mathbf{Q}_{jt} = \{Q_{jnt}\}$, $n \in \Lambda_{jt}$ is a vector of quantities of the products produced by firm j in period t ; $\mathbf{Q}_{-jt} = \{Q_{kt}\}$, $k \neq j$ is a vector of quantities of the products produced by the competitors of firm j in period t ; $\boldsymbol{\xi}_t$ is a vector of quality levels of products produced by firm j and its competitors.

As an identification condition, we assume that the demand system admits a unique solution for the quality levels given observable price and quality outcomes. That is,

$$\boldsymbol{\xi}_{jnt} = P_{jnt}^{-1}(\mathbf{P}_t, \mathbf{Q}_t), \quad (\text{A.2})$$

where \mathbf{P}_t and \mathbf{Q}_t are the vectors of prices and qualities of all products and firms, respectively. This assumption is satisfied by a wide set of demand functions, such as CES and Logit. In particular, in a CES demand system, this equation degenerates to (13).

To ease the exposure of the model, we denote the price elasticity of demand of any product

n of firm j with respect to another product k as:

$$\frac{\partial Q_{jnt}}{\partial P_{jtk}} \frac{P_{jtk}}{Q_{jnt}} = -\eta_{jtnk}. \quad (\text{A.3})$$

Note that we have slightly abused the notation (in order to simplify it) because product k can be either a product of the same firm j (i.e., cannibalization) or a product of any other firm $j' \neq j$ (i.e., competition). That is, the elasticity defined in (A.3) is fully flexible and varies by firm, product and time. In the paper, where a standard constant-elasticity demand function is assumed, $\eta_{jtnk} = \eta_n$ if $n = k$ and $\eta_{jtnk} = 0$ if $n \neq k$.

Now we turn to the production side. As in the paper, we use a transformation function to model the production technology. Specifically, given the set of products to be produced (Λ_{jt}) and associated product appeal ($\tilde{\xi}_{jnt}$, $n \in \Lambda_{jt}$), the firm uses labor (L_{jt}), material (M_{jt}), and capital (K_{jt}) to produce output quantity (Q_{jnt} , $n \in \Lambda_{jt}$) following a constant elasticity of substitution (CES) transformation function as the input aggregator. However, instead of assuming a linear output aggregator as in the paper, we use a CES aggregator as in [Cairncross et al. \(2023\)](#), who derive the transformation function from individual product production functions and show that the output aggregator should be in a CES format if there are shared inputs across products within the firm.

Formally, the transformation production function is modelled as:

$$G(\mathbf{Q}_{jt}) = F(L_{jt}, M_{jt}, K_{jt}), \quad (\text{A.4})$$

where \mathbf{Q}_{jt} is the vector of output quantities, and the output aggregator is

$$G(\mathbf{Q}_{jt}) \equiv \left[\sum_{n \in \Lambda_{jt}} e^{-\tilde{\omega}_{jnt}} Q_{jnt}^\theta \right]^{\frac{1}{\theta}}, \quad \theta \leq 1, \quad (\text{A.5})$$

and the input aggregator is⁵⁶

$$F(L_{jt}, M_{jt}, K_{jt}) \equiv [\alpha_L L_{jt}^\gamma + \alpha_M M_{jt}^\gamma + \alpha_K K_{jt}^\gamma]^{\frac{\rho}{\gamma}}, \quad (\text{A.6})$$

While the interpretation of the variables and parameters remain the same as in Section 2.2, this setup introduces a new parameter θ to the production model. If $\theta = 1$, then the model

⁵⁶Generalizing the CES input aggregator to more flexible functional forms, such as translog, is possible but more challenging because there are more production parameters to be identified and estimated. If one is willing to sacrifice the advantage of not relying on productivity evolution, a specific approach of such an extension for single-product firms is provided by [Grieco et al. \(2016\)](#).

degenerates to our model in the main body of the paper.

A.2 The Unique One-to-One Mapping from Observables to Unobservables

As in the paper, the estimating equation is derived from the first-order conditions of firm profit maximization. Specifically, the firm's objective is to maximize its total profits from all products in period t after observing its state, by optimally choosing the quantity of material (M_{jt}), the quantity of labor (L_{jt}), and the quantities of all the products to be produced ($\mathbf{Q}_{jt} = \{Q_{jnt}\}, n \in \Lambda_{jt}$):

$$\begin{aligned} \pi(s_{jt}) &= \max_{\mathbf{Q}_{jt}, M_{jt}, L_{jt}} \sum_{n \in \Lambda_{jt}} P_{jnt} Q_{jnt} - P_{M_{jt}} M_{jt} - P_{L_{jt}} L_{jt} \\ \text{subject to:} & \quad (\text{A.1}) \text{ and } (\text{A.4}). \end{aligned} \quad (\text{A.7})$$

The Lagrange function implied by the profit maximization problem is:

$$\begin{aligned} \mathcal{L}_{jt} &= \sum_{n \in \Lambda_{jt}} P_{jnt}(\mathbf{Q}_{jt}, \mathbf{Q}_{-jt}; \boldsymbol{\xi}_t) Q_{jnt} - P_{L_{jt}} L_{jt} - P_{M_{jt}} M_{jt} \\ &\quad - \lambda_{jt} \left\{ G(\mathbf{Q}_{jt}) - F(L_{jt}, M_{jt}, K_{jt}) \right\}. \end{aligned} \quad (\text{A.8})$$

The first-order conditions with respect to labor and material inputs are, respectively:

$$\frac{\partial \mathcal{L}_{jt}}{\partial L_{jt}} = -P_{L_{jt}} + \lambda_{jt} \frac{\partial F(L_{jt}, M_{jt}, K_{jt})}{\partial L_{jt}} = 0, \quad (\text{A.9})$$

$$\frac{\partial \mathcal{L}_{jt}}{\partial M_{jt}} = -P_{M_{jt}} + \lambda_{jt} \frac{\partial F(L_{jt}, M_{jt}, K_{jt})}{\partial M_{jt}} = 0. \quad (\text{A.10})$$

While the above two first-order conditions regarding inputs are the same as in the main text, the first-order condition with respect to each product quantity Q_{jnt} , $n \in \Lambda_{jt}$, is generalized:

$$\frac{\partial \mathcal{L}}{\partial Q_{jnt}} = \sum_{k \in \Lambda_{jt}} \frac{\partial P_{jtk}(\mathbf{Q}_{jt}, \mathbf{Q}_{-jt}; \boldsymbol{\xi}_t)}{\partial Q_{jnt}} Q_{jtk} + P_{jnt} - \lambda_{jt} \frac{\partial G(\mathbf{Q}_{jt})}{\partial Q_{jnt}} = 0, \quad (\text{A.11})$$

where $\frac{\partial G(\mathbf{Q}_{jt})}{\partial Q_{jnt}} = e^{-\tilde{\omega}_{jnt}} Q_{jnt}^{\theta-1} \left[\sum_{n \in \Lambda_{jt}} e^{-\tilde{\omega}_{jnt}} Q_{jnt}^\theta \right]^{\frac{1}{\theta}-1}$. Intuitively, $\lambda_{jt} \frac{\partial G(\mathbf{Q}_{jt})}{\partial Q_{jnt}}$ is the marginal cost of Q_{jnt} , which we denote it as $mc(Q_{jnt})$.

Because the first-order conditions, (A.9) and (A.10), are the same as in the paper, the solution for M_{jt} and λ_{jt} (which are implied by the two first-order conditions) are unchanged

too:

$$M_{jt} = \left[\frac{\alpha_L E_{Mjt}}{\alpha_M E_{Ljt}} \right]^{\frac{1}{\gamma}} L_{jt}. \quad (\text{A.12})$$

and

$$\lambda_{jt} = \frac{E_{Ljt}}{\rho \alpha_L L_{jt}^\gamma} \left[\alpha_L L_{jt}^\gamma \left(1 + \frac{E_{Mjt}}{E_{Ljt}} \right) + \alpha_K K_{jt}^\gamma \right]^{1 - \frac{\rho}{\gamma}}. \quad (\text{A.13})$$

Substituting (A.13) into (A.11) gives the solution for productivity:

$$e^{\tilde{\omega}_{jnt}} = \frac{Q_{jnt}^{\theta-1} G(\mathbf{Q}_{jt})^{1-\theta}}{\sum_{k \in \Lambda_{jt}} \frac{\partial P_{jtk}(\mathbf{Q}_{jt}, \mathbf{Q}_{-jt}; \boldsymbol{\xi}_t)}{\partial Q_{jnt}} Q_{jtk} + P_{jnt}} \underbrace{\frac{E_{Ljt}}{\rho \alpha_L L_{jt}^\gamma} \left[\alpha_L L_{jt}^\gamma \left(1 + \frac{E_{Mjt}}{E_{Ljt}} \right) + \alpha_K K_{jt}^\gamma \right]^{1 - \frac{\rho}{\gamma}}}_{\lambda_{jt}}, \quad (\text{A.14})$$

where $\boldsymbol{\xi}_t$ is given by (A.2). This equation is the analog of (17). If the demand function is CES and $\theta = 1$, then this equation degenerates to (17) exactly. If $\theta \neq 1$, (A.14) suggests that the productivity measure absorbs a product scale effect ($Q_{jnt}^{\theta-1}$) and firm scale effect ($G(\mathbf{Q}_{jt})^{1-\theta}$).

Critically, as in the paper, we have derived a unique one-to-one mapping from observable variables ($\mathbf{Q}_{jt}, \mathbf{P}_{jt}, L_{jt}, P_{Ljt}, E_{Mjt}, K_{jt}$) to unobservable variables ($\boldsymbol{\xi}_{jt}, M_{jt}, \lambda_{jt}, \omega_{jt}$) specified by (A.2), (A.12), (A.13), and (A.14). Thus, these variables can be computed directly once the demand and production parameters are estimated.

A.3 Estimating Demand and Production Parameters

Estimating a flexible demand system like (A.1) is challenging due to unobservable demand factors (e.g., quality) and the endogeneity of prices. Most existing approaches (e.g. Berry, 1994; Berry et al., 1995) rely on a set of valid instrumental variables. Since our focus is on the production transformation function, we assume the existence of a valid set of instrumental variables, allowing us to estimate the demand system (A.1). Consequently, the firm-product-time-specific quality can be recovered via (A.2).

Thus, our primary focus is on estimating the production parameters, conditional on the demand system (A.1) being estimated and quality being recovered. This approach, which involves first estimating the demand system and then the production functions, is adopted in the literature (e.g. Orr, 2022).

A.3.1 Estimating the main production parameters

To derive our main estimating equation, we start by multiplying both sides of the equation implied by (A.11) by Q_{jnt} . Rearranging this equation gives:

$$\sum_{k \in \Lambda_{jt}} \frac{\partial P_{jtk}(\mathbf{Q}_{jt}, \mathbf{Q}_{-jt}, \boldsymbol{\xi}_t)}{\partial Q_{jnt}} \frac{Q_{jnt}}{P_{jtk}} P_{jtk} Q_{jtk} + P_{jnt} Q_{jnt} = \lambda_{jt} G(\mathbf{Q}_{jt}) \frac{e^{-\tilde{\omega}_{jnt}} Q_{jnt}^\theta}{\left[\sum_{n \in \Lambda_{jt}} e^{-\tilde{\omega}_{jnt}} Q_{jnt}^\theta \right]}. \quad (\text{A.15})$$

Use the definition $R_{jnt} = P_{jnt} Q_{jnt} e^{u_{jt}}$, where u_{jt} is a mean-zero i.i.d. shock (i.e., measurement error or unexpected shock), and (A.3), the above equation can be written as:

$$R_{jnt} - \sum_{k \in \Lambda_{jt}} \frac{1}{\eta_{jtnk}} R_{jtk} = \lambda_{jt} G(\mathbf{Q}_{jt}) \frac{e^{-\tilde{\omega}_{jnt}} Q_{jnt}^\theta}{\left[\sum_{n \in \Lambda_{jt}} e^{-\tilde{\omega}_{jnt}} Q_{jnt}^\theta \right]} e^{u_{jt}}. \quad (\text{A.16})$$

Sum the above equation over $n \in \Lambda_{jt}$ to obtain:

$$\begin{aligned} \sum_{n \in \Lambda_{jt}} \left(1 - \sum_{k \in \Lambda_{jt}} \frac{1}{\eta_{jtnk}} \frac{R_{jtk}}{R_{jnt}} \right) R_{jnt} &= \lambda_{jt} G(\mathbf{Q}_{jt}) e^{u_{jt}} \\ &= \lambda_{jt} F(L_{jt}, M_{jt}, K_{jt}) e^{u_{jt}} \\ &= \lambda_{jt} \left[\alpha_L L_{jt}^\gamma + \alpha_M M_{jt}^\gamma + \alpha_K K_{jt}^\gamma \right]^{\frac{\rho}{\gamma}} e^{u_{jt}} \\ &= \frac{E_{Ljt}}{\rho \alpha_L L_{jt}^\gamma} \left[\alpha_L L_{jt}^\gamma \left(1 + \frac{E_{Mjt}}{E_{Ljt}} \right) + \alpha_K K_{jt}^\gamma \right] e^{u_{jt}} \\ &= \frac{1}{\rho} \left[E_{Mjt} + E_{Ljt} \left(1 + \frac{\alpha_K}{\alpha_L} \left(\frac{K_{jt}}{L_{jt}} \right)^\gamma \right) \right] e^{u_{jt}}. \end{aligned} \quad (\text{A.17})$$

The second equality comes from the transformation function and the second last equality is a result of substituting λ_{jt} and M_{jt} by (A.12) and (A.13), respectively.

Take logarithm of the above equation to obtain:

$$\ln \left[\sum_{n \in \Lambda_{jt}} \left(1 - \sum_{k \in \Lambda_{jt}} \frac{1}{\eta_{jtnk}} \frac{R_{jtk}}{R_{jnt}} \right) R_{jnt} \right] = -\ln \rho + \ln \left[E_{Mjt} + E_{Ljt} \left(1 + \frac{\alpha_K}{\alpha_L} \left(\frac{K_{jt}}{L_{jt}} \right)^\gamma \right) \right] + u_{jt}. \quad (\text{A.18})$$

Note that $\left(1 - \sum_{k \in \Lambda_{jt}} \frac{1}{\eta_{jtnk}} \frac{R_{jtk}}{R_{jnt}} \right)$ is the reciprocal of markup of product n of firm j in

period t . To see this, the firm-product-time specific markup is defined as:

$$\begin{aligned}
\mu_{jnt} &\equiv \frac{P_{jnt}}{mc(Q_{jnt})} \\
&= \frac{P_{jnt}}{\sum_{k \in \Lambda_{jt}} \frac{\partial P_{jtk}(\mathbf{Q}_{jt}, \mathbf{Q}_{-jt}; \boldsymbol{\xi}_t)}{\partial Q_{jnt}} Q_{jtk} + P_{jnt}} \\
&= \frac{1}{\sum_{k \in \Lambda_{jt}} \frac{\partial P_{jtk}(\mathbf{Q}_{jt}, \mathbf{Q}_{-jt}; \boldsymbol{\xi}_t)}{\partial Q_{jnt}} \frac{Q_{jtk}}{P_{jnt}} + 1} \\
&= \frac{1}{\sum_{k \in \Lambda_{jt}} \frac{\partial P_{jtk}(\mathbf{Q}_{jt}, \mathbf{Q}_{-jt}; \boldsymbol{\xi}_t)}{\partial Q_{jnt}} \frac{Q_{jnt}}{P_{jtk}} \frac{P_{jtk} Q_{jtk}}{P_{jnt} Q_{jnt}} + 1} \\
&= \frac{1}{1 - \sum_{k \in \Lambda_{jt}} \frac{1}{\eta_{jtnk}} \frac{R_{jtk}}{R_{jnt}}}, \tag{A.19}
\end{aligned}$$

where the second equality comes from the definition of marginal cost, as defined below (A.11), and the last equality comes from the definition of revenue and demand elasticity, as defined in (A.3).

Therefore, (A.18) describes the relationship between revenues (adjusted by reciprocal of markups) and inputs for a *general system system*, conditional on that the firm maximizes profit. It is clear that this equation is analog of (18) in the paper. In particular, in the setup of a CES demand function, $\left(1 - \sum_{k \in \Lambda_{jt}} \frac{1}{\eta_{jtnk}} \frac{R_{jtk}}{R_{jnt}}\right) = \frac{\eta_n - 1}{\eta_n}$, and thus (A.18) degenerates to (18). Note that θ does not appear in this equation. As a result, even if the outputs are indeed complementary in production (i.e., $\theta < 1$), (A.18) is not mis-specified for a model assuming $\theta = 1$ and thus the rest of production parameters are still correctly identified and estimated. Of course, such robustness comes with a challenge – we need an extra equation to identify and estimate θ , which will be provided in Section A.3.2.

Given that the demand system (A.1) is estimated, the left-hand side of (A.18) can be computed. Denote the computed value as

$$\hat{\Upsilon}_{jt} = \left[\sum_{n \in \Lambda_{jt}} \left(1 - \sum_{k \in \Lambda_{jt}} \frac{1}{\hat{\eta}_{jtnk}} \frac{R_{jtk}}{R_{jnt}} \right) R_{jnt} \right].$$

Intuitively, $\hat{\Upsilon}_{jt}$ is the firm-level revenue adjusted by the reciprocals of markups.

The only unobserved shock in (A.18) is the unexpected i.i.d. shock u_{jt} . As a result, (A.18)

can be estimated using a Nonlinear Least Square estimator:⁵⁷

$$\hat{\beta} = \operatorname{argmin}_{\beta} \frac{1}{\mathbb{N}} \sum_{jt} \left\{ \ln \hat{\Upsilon}_{jt} + \ln \rho - \ln \left[E_{M_{jt}} + E_{L_{jt}} \left(1 + \frac{\alpha_K}{\alpha_L} \left(\frac{K_{jt}}{L_{jt}} \right)^\gamma \right) \right] \right\}^2$$

subject to: $\alpha_L + \alpha_M + \alpha_K = 1$ and $\frac{\alpha_M}{\alpha_L} = \frac{\bar{E}_M}{\bar{E}_L}$

where $\beta = (\rho, \alpha_L, \alpha_M, \alpha_K, \gamma)$.

A.3.2 Estimating the remaining production parameter

While (A.18) provides a simple way of estimating the main production parameters, the parameter that governs the elasticity of substitution of the output aggregator, θ , is left unidentified. To estimate θ , we adopt a similar strategy as employed in Section (3.2) of the paper – utilizing the relationship across products with a firm.

Specifically, (A.16) implied by the first-order condition with respect to a product can be expressed as:

$$\left(1 - \sum_{k \in \Lambda_{jt}} \frac{1}{\eta_{jtnk}} \frac{R_{jtk}}{R_{jnt}} \right) R_{jnt} = \lambda_{jt} G(\mathbf{Q}_{jt}) \frac{e^{-\tilde{\omega}_{jnt}} Q_{jnt}^\theta}{\left[\sum_{n \in \Lambda_{jt}} e^{-\tilde{\omega}_{jnt}} Q_{jnt}^\theta \right]} e^{u_{jt}}, \quad \forall n \in \Lambda_{jt}. \quad (\text{A.20})$$

From the the estimated demand system (A.1), compute and define

$$\hat{\Upsilon}_{jnt} = \left[\left(1 - \sum_{k \in \Lambda_{jt}} \frac{1}{\hat{\eta}_{jtnk}} \frac{R_{jtk}}{R_{jnt}} \right) R_{jnt} \right].$$

Intuitively, $\hat{\Upsilon}_{jnt}$ is the firm-product-level revenue adjusted by the reciprocal of product markup.

Without loss of generality, we assume that product 1 is the reference product (that is, the one produced by most firms in the industry). Take the ratio of (A.20) of any other products to the main product. The logarithm of the ratio is:

$$\ln \left(\frac{\hat{\Upsilon}_{jnt}}{\hat{\Upsilon}_{jt1}} \right) = \theta \ln \left(\frac{Q_{jnt}}{Q_{jt1}} \right) + v_{jnt}, \quad (\text{A.21})$$

where $v_{jnt} = \tilde{\omega}_{jt1} - \tilde{\omega}_{jnt}$ is the relative *difference* between the productivity of the two products.

⁵⁷If one is concerned with the error term contained in revenues, which are used in computing the elasticities in $\hat{\Upsilon}_{jt}$, a GMM approach can be implemented to estimate (A.18) using the same set of instrumental variables as proposed in estimating (21) as described in Section 3.2.

Like (19), (A.21) states the relationship between the main product and any other product produced by the same firm. In the setup of a CES demand and $\theta = 1$ as in the paper, it is straightforward to show (by substituting Q_{jnt} by a function of R_{jnt} using the CES demand function) that (A.21) degenerates to (19).

We estimate θ from (A.21) treating v_{jnt} as an error term via a 2-Stage Least Square estimator, using a set of instrumental variables. While firm-product-level instrumental variables would be preferred (if available), firm-level variables such as capital stock (K_{jt}) and wage rate (P_{Ljt}) are sufficient – an advantage due to the use of *relative* differences in revenues and quantities of the two products. A similar argument has been applied to the estimation of (19). Specifically, while the *levels* of these variables are uncorrelated with the *differences* in productivity (i.e., v_{jnt}) between two products, they influence the ratio of the quantities of the two products. For example, conditional on everything else, a higher level of capital stock leads to high quantities of both products, but the product (e.g., the reference product) with less elastic demand expands more than the other, leading to a lower ratio quantity (Q_{jnt}/Q_{jt1}).⁵⁸

A.3.3 Advantages and challenges of estimating the general model

The general model is attractive due to its flexibility. Unlike our CES demand model specified in the paper, the general model incorporates a comprehensive demand system, allowing us to study the cannibalization effects across different products produced by the same firm and to account for markups that vary flexibly at the firm-product-time level. Moreover, once this general demand system is estimated, it simplifies the estimation of the production parameters. In particular, unlike the GMM approach implemented in the paper to estimate the demand elasticity and production parameters, the estimation of the main equation (A.18) can be accomplished using NLLS. Additionally, the general model features a more flexible output aggregator structure, with an additional parameter governing the rate of substitution among different outputs in the transformation function.

Nonetheless, estimating such a general model can be challenging. First, the complexities of estimating a richer structure of the demand system (A.1) are well recognized. One particular challenge arises from the availability of suitable instrumental variables for endogenous product prices. Traditionally used instrumental variables, such as cost shifters, may fail as valid instruments when products are vertically differentiated and firms choose higher-quality, more costly inputs to improve the quality of products. In general, it requires carefully constructed instrumental variables that are orthogonal to quality differences. For example, [Berry et al.](#)

⁵⁸In our Monte Carlo experiment, the correlation between capital stock and (Q_{jnt}/Q_{jt1}) (both in logarithm) is -0.2 to -0.1, depending on the demand elasticity differences. In the data, the correlation is about -0.06.

(1995) utilize the characteristics of other automobiles produced by the firm itself and similar automobiles produced by its rivals. In the context of estimating the production function of multi-product firms, Orr (2022) designs sophisticated instrumental variables to leverage the variation in product sets and material input price growth experienced by firms in other output markets that use similar inputs in order to estimate a flexible demand function. These strong, valid variables are not always available, like in our case.

Second, the estimation of the CES output aggregator, θ , relies on the availability of suitable instrumental variables that are correlated with the quantity ratio (Q_{jnt}/Q_{jt1}). In the approach in Section A.3.2, we propose the use of firm-level variables, such as capital stock and wage rate, as instrumental variables when firm-product-level instruments are not available. However, the relevance of such firm-level instrumental variables depends on the characteristics of demand across products. For example, if the demand elasticities are the same across products, then capital stock will not shift the quantity ratio and thus will not serve as a valid instrumental variable for estimating (A.21) for θ .

The model implemented in the paper is a simplified version of the general model. There, the demand system is a standard, commonly used CES demand function, and the parameter governing the substitution of outputs in the output aggregator is assumed to be one. These restrictions provide the advantage of estimating the demand and production parameters jointly without relying on the availability of suitable instrumental variables required for estimating the general model.

Despite its simplicity, this model retains several key advantages of the general model. On the demand side, although it abstracts away from cannibalization effects within a firm, it does allow for correlation of demand of products produced by the same firm. On the production side, although it assumes a linear output aggregator, the potential complementary effects governed by θ is absorbed by the flexible productivity (A.14), and consequently, production parameters estimated from (A.18) are not biased. In terms of estimating strategy, it eliminates the need for imputing firm-product input shares or imposing productivity evolution processes, while maintaining flexibility in the relationship between productivity and quality. Additionally, it is scalable to accommodate numerous products and can address bias caused by unobserved heterogeneous intermediate input prices. Nonetheless, our methodology can be readily implemented to the general model if suitable instrumental variables for estimating the demand system (A.1) and (A.21) are available.

B Dynamic Decisions

This section describes the dynamic decisions made by the firm as a completion of the full model. At the end of each period t , the firm chooses the set of products to produce, their associated quality levels, and investment in technical efficiency improvement (e.g., research and development), for the next period ($t + 1$). These decisions are made conditional on the current state $s_{jt} = (\Lambda_{jt}, \omega_{jt}, \xi_{jt}, K_{jt}, P_{Mjt}, P_{Ljt}, \chi_{jt})$ and after observing the adjustment costs of product scope and quality levels. Although the evolution of K_{jt} , P_{Mjt} , P_{Ljt} and χ_{jt} can be endogenous, we remain agnostic on modelling their exact evolution processes because our estimation method focuses on the static decisions and does not rely on how these variables evolve over time. The adjustment costs of product scope capture the costs incurred by the firm to install and arrange new production lines. The adjustment costs of product quality contain the costs of modifying the production procedure and sourcing new suppliers of the material input to meet the new quality levels.

In making decisions regarding product scope, quality levels, and investment, the firm is forward-looking and takes into account the impact of the current decisions on the future paths of the state variables. In particular, the firm knows that the choice of improving the quality of a product for the next period will reduce the associated (quantity-based) productivity in the next period (i.e., due to the cost of quality). As a result, these decisions are dynamic.

Although we do not estimate the complex dynamic model in this paper (due to the considerably high dimension of the state variables),⁵⁹ the model serves the crucial purpose of clarifying the (dynamic) choices made by the firm and their implications when we estimate the static model. In particular, the dynamic model implies that even if the underlying technical efficiency follows a simple AR(1) process, the resulting productivity (i.e., TFPQ) is not an AR(1) process as assumed in the literature. To see this, note that quality ξ_{jt+1n} is endogenously determined by the firm based on the state variable vector s_{jt} , including the technical efficiency of all products (i.e., ω_{jt}). Considering the impact of quality on productivity shown in (6), the productivity of any product n in period $t + 1$ depends on the entire state vector s_{jt} in a highly nonlinear way. Ignoring such interdependent relationships may potentially result in biased estimation. Fortunately, our empirical method does not use any assumptions regarding how technical efficiency and productivity evolve, as discussed in Section 3.

⁵⁹For example, even in the footwear industry with only four products, the dynamic state includes at least 10 continuous variables – 4 variables for technical efficiency, 4 variables for product quality, and 2 for the material and labor prices.

C Additional Tables and Figures

Table A1: Product list, manufacturing of footwear, mainly of leather (class 324001)

Industry	Product description	Code
324001	Cow leather, for men	1
324001	Cow leather, for women	2
324001	Cow leather, for kids	3
324001	Others	99

Table A2: Product list, printing and binding (class 342003)

Industry	Product description	Code
342003	Printing of Calendars and almanacs	5
342003	Folding boxes	6
342003	Labels and prints	13
342003	Brochures and catalogs	14
342003	Continuous forms	15
342003	Accounting, administrative and tax forms	16
342003	Telephone directories	17
342003	Books	18
342003	Journals	19
342003	Checks	21
342003	Commemorative and business cards	23
342003	Commercial flyers	24
342003	Posters	25
342003	Others	99

Table A3: Product list, manufacturing of pharmaceutical products (class 352100)

Industry	Product description	Code
352100	Medicinal products, for human use with a specific action, anti-infectious: Bactericides	11
352100	Antiparasitics	13
352100	Dermatological	15
352100	Other products with specific action not included in other categories	19
352100	Medicinal products for human use for specialties with action on: Circulatory system	21
352100	Digestive system and metabolism	22
352100	Human musculoskeletal system	23
352100	Respiratory system	24
352100	Sensory organs	25
352100	Genitourinary organs, except hormones	26
352100	Blood and hematopoietic organs	27
352100	Central nervous system	28
352100	Hormones	32
352100	Vitamins and Vitamin Compounds	43
352100	Non-therapeutic products	59
352100	Others	99

Table A4: Within-firm product shares by product scope

Product scope	Product rank (by sales level)				
	1	2	3	4	5+
1	1.000				
2	0.783	0.217			
3	0.675	0.238	0.087		
4	0.560	0.283	0.117	0.040	
5+	0.443	0.204	0.124	0.083	0.146

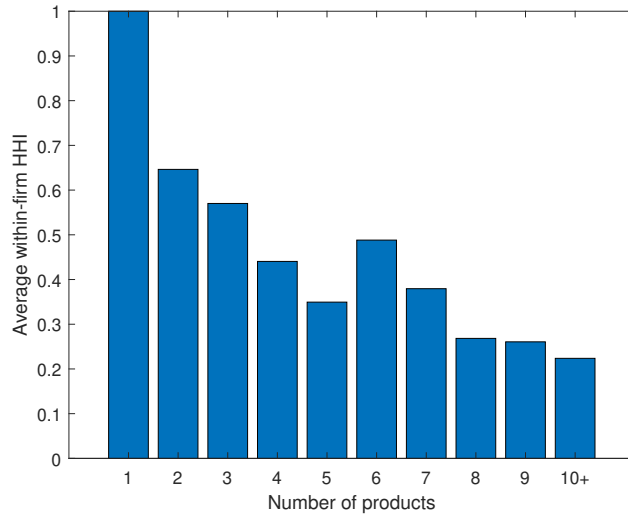
Note: All firm-year pairs producing 5 products or more are clustered in the “5+” group. All products ranked 5 or lower are clustered in the “5+” group.

Table A5: Descriptive statistics

Variable	Footwear	Printing	Pharmaceutical
Revenue per product (R)	64.846 (101.261)	29.142 (73.091)	100.167 (207.205)
Number of workers (L)	236.180 (361.356)	157.704 (153.598)	450.222 (482.926)
Labor expenditure (E_L)	13.785 (28.726)	17.435 (22.124)	88.792 (110.791)
Material expenditure (E_M)	50.446 (77.265)	65.568 (90.175)	263.033 (382.666)
Capital stock (K)	3.603 (8.444)	21.491 (47.314)	22.534 (31.437)

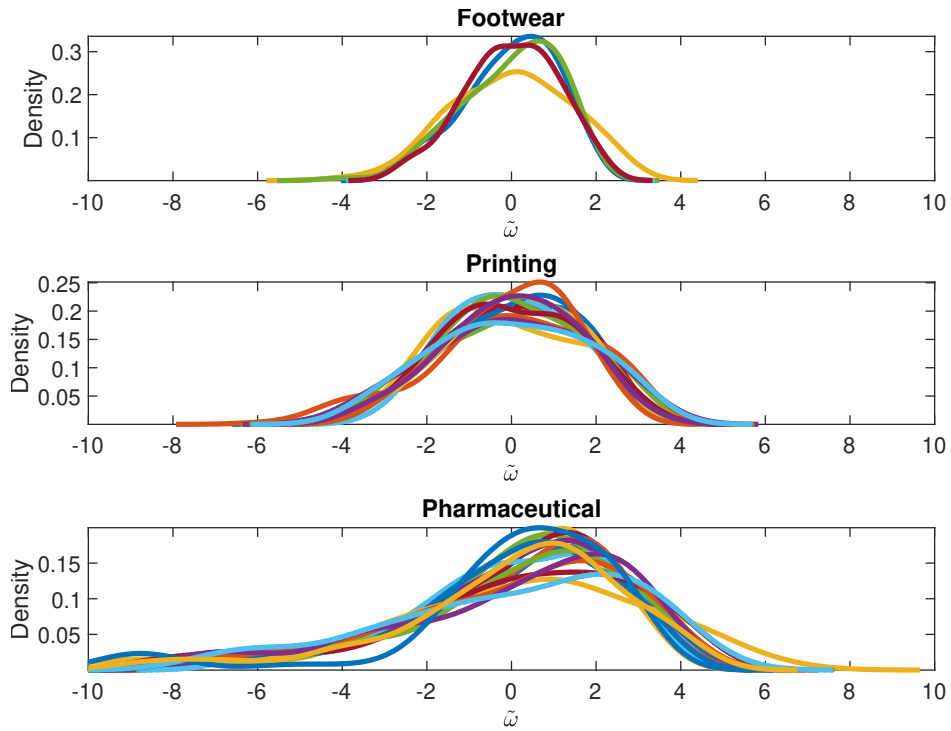
Notes: The table reports the means and standard deviations (in parenthesis) for each variable by industry. R is revenues by product (1 million 2007 Mexican Peso, 1M MXN); L is the number of workers by firm, K is the capital stock by firm (1000 physical units); E_L is the expenditure on labor (wage bill) by firm (1M MXN); E_M is the expenditure on intermediates by firm (1M MXN).

Figure A1: Weighted average within-firm HHI, by number of products



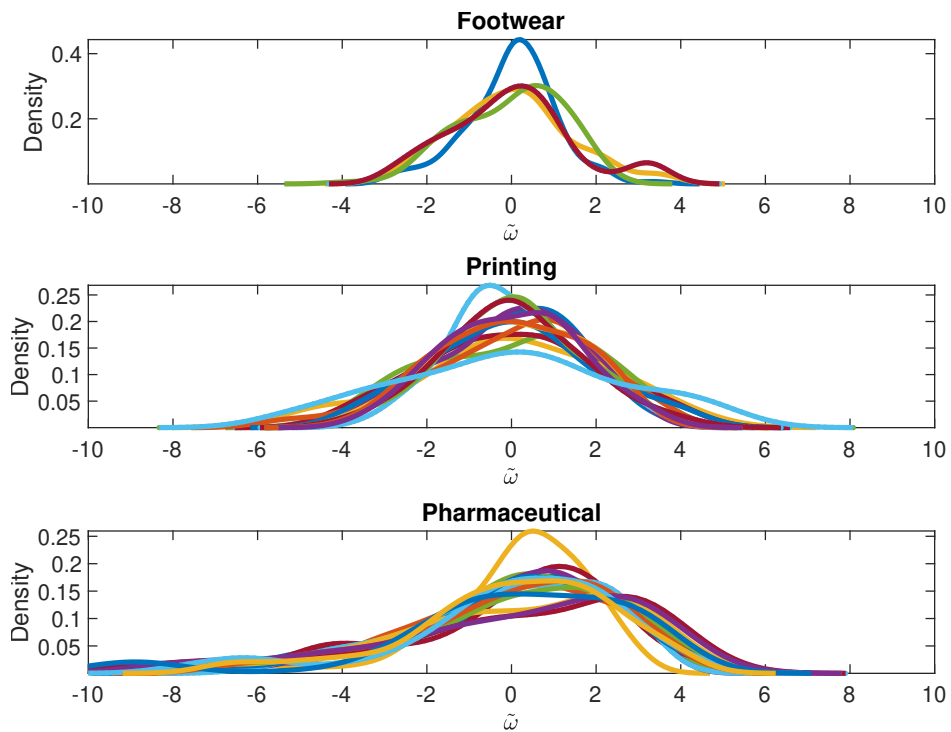
Notes: All firm-year pairs producing 10 products or more are clustered in the “10+” group. The weighted average is calculated using revenues as weights.

Figure A2: Distribution of quality-adjusted productivity, ATFP



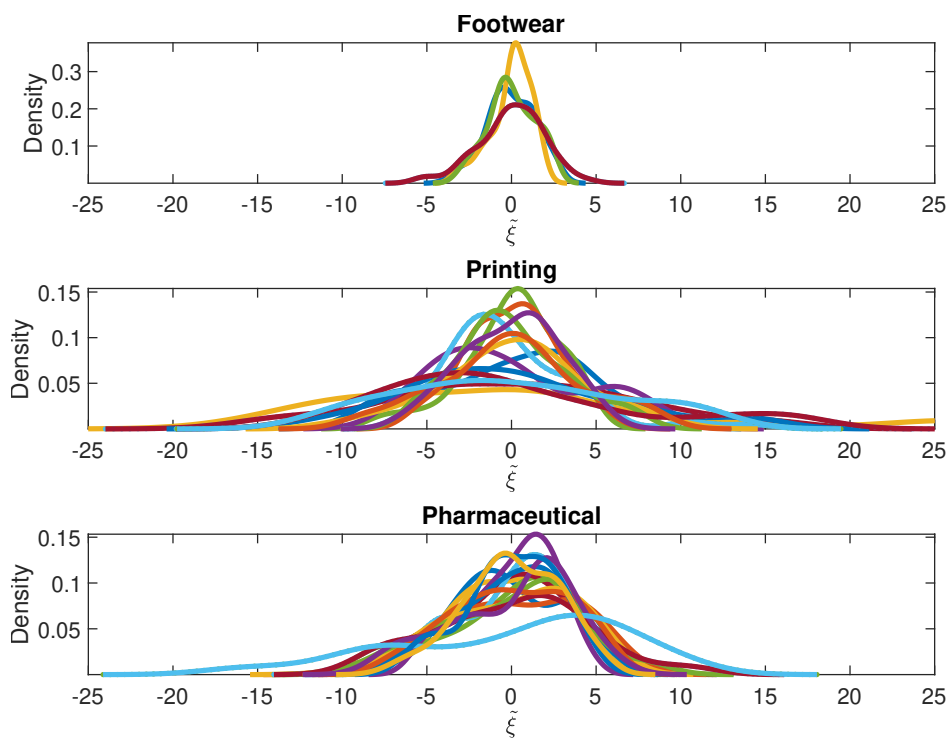
Notes: ATFP is demeaned, and only products with at least 100 observations are included.

Figure A3: Distribution of productivity, $\tilde{\omega}$



Notes: $\tilde{\omega}$ is demeaned, and only products with at least 100 observations are included.

Figure A4: Distribution of quality, $\tilde{\xi}$



Notes: $\tilde{\xi}$ is demeaned, and only products with at least 100 observations are included.

D Monte Carlo Exercises

In this appendix, we present the results of Monte Carlo exercises to demonstrate the performance of our estimation method.

In this Monte Carlo setting, the choice of product sets is exogenous and random. Wage rate, material prices, and capital stock are serially correlated.⁶⁰ The levels of productivity and quality of any given product are not only serially correlated over time but also negative correlated with each other. With this setting, the Monte Carlo exercises consist of N replications of simulated data sets of J firms in T years, given a set of true parameters of interest for 5 products, namely, $(\eta_1, \eta_2, \eta_3, \eta_4, \eta_5, \alpha_L, \alpha_M, \alpha_K, \sigma, \rho)$.

Specifically, in each replication, we simulate productivity $(\tilde{\omega}_{jnt})$ and quality $(\tilde{\xi}_{jnt})$ for each product n , firm j , and time t . We also simulate the wage rate (P_{Ljt}) , the material price (P_{Mjt}) and the capital stock (K_{jt}) for each firm j and time t . All of these variables are serially correlated. In addition, we simulate the negative relationship between productivity and quality as documented in the paper by allowing for a negative correlation r between the shocks in their evolution processes. Specifically, the evolution process of each of these variables for each firm follows an AR(1) process:

$$\begin{aligned}\tilde{\omega}_{jnt} &= g_{0\omega}^n + g_{\omega}^n \tilde{\omega}_{jnt-1} + \varepsilon_{jnt}^{\omega}, & \forall n, \\ \tilde{\xi}_{jnt} &= g_{0\xi}^n + g_{\xi}^n \tilde{\xi}_{jnt-1} + \varepsilon_{jnt}^{\xi}, & \forall n, \\ \ln(P_{Ljt}) &= g_{0\ell} + g_{\ell} \ln(P_{Ljt-1}) + \varepsilon_{jt}^{\ell}, \\ \ln(P_{Mjt}) &= g_{0m} + g_m \ln(P_{Mjt-1}) + \varepsilon_{jt}^m, \\ \ln(K_{jt}) &= g_{0k} + g_k \ln(K_{jt-1}) + \varepsilon_{jt}^k,\end{aligned}$$

where ε is the innovation shock realized in period t , which is assumed to be a normally distributed error term with zero mean and standard deviation $sd(\varepsilon)$. While the shocks in the processes of P_{Ljt} , P_{Mjt} , and K_{jt} are i.i.d., those of $\tilde{\omega}_{jnt}$ and $\tilde{\xi}_{jnt}$ are correlated with a coefficient of r . Although the evolution of the capital stock is exogenous in this setup, the Monte Carlo result is similar if investment (and hence the capital stock) depends on productivity and quality levels.

Given these variables, we use the firm's static profit maximization problem to derive a sequence of optimal choices of labor and material inputs (L_{jt} and M_{jt}), the optimal output quantity (Q_{jnt}) and price (P_{jnt}) for firm j and product n in each period t .

⁶⁰The Monte Carlo result is similar if the evolution of capital stock depends on an investment rule which is a function of capital stock and the levels of productivity and quality.

In this way, we generate a data set of variables for the Monte Carlo experiments. Among them, we use the following variables for the estimation procedure (including the sets of IVs) described in Section 3: $\{Q_{jt1}, \dots, Q_{jt5}, R_{jt1}, \dots, R_{jt5}, K_{jt}, L_{jt}, E_{L_{jt}}, E_{M_{jt}}\}$. The values of the parameters used for the data generation process are reported in Table A6. The mean estimates of the key parameters, together with their corresponding standard errors, are reported in Table A7. Overall, the result shows that our estimation recovers the true parameters of the production and demand functions well.

Table A6: Monte Carlo Parameter Values

Parameter	Description	Value
$\eta_1, \eta_2, \eta_3, \eta_4, \eta_5$	Demand elasticities	7, 6, 5, 4, 3
σ	Elasticity of substitution	2
α_L	Distribution parameter of labor	0.2
α_M	Distribution parameter of material	0.6
α_K	Distribution parameter of capital	0.2
$g_\omega^1, g_\omega^2, g_\omega^3, g_\omega^4, g_\omega^5$	Persistence parameters in productivity evolution	0.75, 0.7, 0.65, 0.6, 0.55
$g_\xi^1, g_\xi^2, g_\xi^3, g_\xi^4, g_\xi^5$	Persistence parameter in quality evolution	0.75, 0.7, 0.65, 0.6, 0.55
g_l	Persistence parameter in wage rate evolution	0.8
g_m	Persistence parameter in material price evolution	0.8
g_k	Persistence parameter in capital evolution	0.8
r	Correlation between productivity and quality shocks	-0.2
$sd(\varepsilon^\omega)$	Standard deviation of productivity shock	0.02
$sd(\varepsilon^\xi)$	Standard deviation of quality shock	0.02
$sd(\varepsilon^\ell)$	Standard deviation of wage rate shock	0.1
$sd(\varepsilon^m)$	Standard deviation of material price shock	0.1
$sd(\varepsilon^k)$	Standard deviation of capital stock shock	0.1
$sd(u)$	Standard deviation of revenue measurement error (u)	0.01
T	Number of periods	15
J	Number of firms	400
N	Number of Monte Carlo replications	300

Table A7: Monte Carlo Estimates of Production and Demand Function Parameters

Parameter	True	Estimate
$\frac{\eta_1-1}{\eta_2-1}$	1.200	1.199 (0.021)
$\frac{\eta_1-1}{\eta_3-1}$	1.500	1.499 (0.027)
$\frac{\eta_1-1}{\eta_4-1}$	2.000	1.999 (0.037)
$\frac{\eta_1-1}{\eta_5-1}$	3.000	3.002 (0.053)
α_L	0.200	0.200 (0.002)
α_M	0.600	0.600 (0.001)
α_K	0.200	0.200 (0.002)
σ	2.000	2.000 (0.010)
ρ	1.100	1.101 (0.009)
η_1	7.000	7.000 (0.350)
η_2	6.000	6.003 (0.254)
η_3	5.000	5.001 (0.204)
η_4	4.000	4.002 (0.157)
η_5	3.000	2.998 (0.102)

Note: The parameter estimates are reported as the mean estimates from the Monte Carlo simulations. Standard errors in parentheses are computed as the standard deviation of the estimates.