

# Endogenous Returns to Scale\*

Alexandr Kopytov

University of Rochester

Mathieu Taschereau-Dumouchel

Cornell University

Zebang Xu

Cornell University

January 30, 2026

## Abstract

We develop a general equilibrium model in which firms choose how scalable their production technologies are. More scalable technologies make it easier for firms to expand output but are less effective at small scale. In equilibrium, more productive firms adopt more scalable technologies and grow disproportionately large. As a result, the tail of the size distribution becomes thicker and, as resources reallocate to the most productive producers, GDP increases. Over the long-run, as aggregate productivity rises, firms adopt more scalable technologies, which lowers input prices, leading to further increases in scalability. Through this supply-chain amplification process, endogenous returns to scale raise the growth rate of GDP. A calibrated version of the model shows that these effects are quantitatively significant. We also document support for the model's predictions in firm-level data.

**JEL Classifications:** E23, D24, D57, O40, L11

---

\*We thank Guangbin Hong, Ron Kaniel, Narayana Kocherlakota, Ezra Oberfield, Christian Opp, Tommaso Porzio, Markus Poschke and participants at various seminars and conferences for helpful suggestions.

# 1 Introduction

At the turn of the twentieth century, the automobile was a luxury good, carefully assembled by skilled craftsmen. State-of-the-art models like the 1909 Cadillac Model Thirty were already built with interchangeable parts, yet their production remained artisanal and their price prohibitively high. Just a few years later, Henry Ford’s Highland Park plant was producing a Model T every 93 minutes. By adopting the moving assembly line, Ford reorganized production to increase its returns to scale. The benefit was a large decline in production costs that transformed the automobile into a mass-market product.

Ford’s story is not unique. In his seminal work, Chandler (1990) argued that the adoption of processes, organizational structures, and technologies designed to achieve greater economies of scale was a key driver of modern economic growth. This historical perspective suggests that scalability is not a fixed constraint but a dimension of technology that firms actively manage. Building on this insight, we explore the idea that returns to scale are endogenous equilibrium objects driven by incentives, and study its implications for the firm size distribution, the response of the economy to shocks, and long-run growth.

To do so, we develop a multi-sector general equilibrium model with endogenous returns to scale. Within each sector, a continuum of firms with heterogeneous productivity produce a common good using labor and intermediate inputs. Importantly, firms in our setup are free to choose their returns to scale subject to a technological trade-off: achieving larger scale economies comes at a cost in terms of productivity. As a result, while more scalable technologies make it easier for firms to expand output, they are less effective at small scale.

The existence of a trade-off between scale and productivity is well documented. Chandler (1977) describes how the historical shift from small, artisanal producers to large, integrated enterprises depended on costly investment in a new “technology of organization.” Achieving and managing large scale required the creation of professional hierarchies and complex administrative systems. While this “visible hand” of management enabled firms to coordinate large-scale production, it introduced administrative overhead and reduced the operational flexibility that smaller firms enjoyed.

Our analysis of the model begins with the individual firm’s decisions in partial equilibrium. Following McKenzie (1959), we interpret decreasing returns to scale as arising from a fixed entrepreneurial factor. A firm’s choice of returns to scale therefore reflects a trade-off between this constrained in-house factor and the variable bundle

of labor and intermediate inputs. We show that the optimal degree of scalability is reached when the marginal productivity loss from expanding scalability exactly offsets the cost savings from relying less on the fixed factor. This condition dictates how firms adapt to their environment: any change that pushes the firm to expand puts pressure on the fixed factor and encourages the adoption of a technology with higher returns to scale. Consequently, higher productivity, higher output prices, or cheaper intermediate inputs all induce the firm to adopt a more scalable, input-intensive production function.

This mechanism generates a “double blessing” for the most productive firms. Their intrinsic productivity (the first blessing) naturally leads to larger size and higher profits. This expansion, in turn, tightens the constraint imposed by the fixed entrepreneurial factor, creating an incentive to adopt more scalable technologies (the second blessing), which leads to a further increase in size. This disproportionate growth creates superstar firms and a thick Pareto tail in the firm-size distribution.

Despite heterogeneity in returns to scale, we show that firms in a sector can be aggregated in a tractable way. Because of free entry, production exhibits constant returns to scale at the sector level, with firm-level scalability decisions affecting the importance of labor and intermediate inputs in sectoral production. Returns-to-scale decisions also manifest themselves in sectoral productivity. We show that by allowing high-productivity firms to grow larger, endogenous returns to scale increases sectoral productivity and, through that channel, the level of GDP.

The model admits a unique and efficient competitive equilibrium, which can be characterized as the solution to a social planner’s problem. We use this characterization to study the determinants of returns to scale in general equilibrium. We find that any shock that lowers the relative cost of intermediate inputs induces firms to adopt more scalable technologies. For instance, a productivity improvement in an upstream sector reduces input costs for downstream customers, encouraging them to increase their scalability. This shift toward greater scalability implies a heavier reliance on intermediate inputs, which raises the Domar weights of upstream suppliers. Consequently, sectors experiencing productivity gains become more central to the economy, amplifying their impact on GDP. Through this channel, endogenous returns to scale magnifies the benefits of positive shocks. Symmetrically, the same mechanism dampens the adverse impact of negative shocks, as firms substitute away from intermediate inputs and reduce the importance of the affected sectors.

Endogenous returns to scale also matter for long-run growth. With recurrent productivity improvements, firms continuously adopt more scalable technologies. This

leads to an increase in Domar weights, making subsequent productivity gains even more impactful. As this process unfolds, the economy enters an acceleration phase in which the growth rate of GDP rises over time, eventually converging to a long-run rate strictly higher than in a fixed-technology economy. In our model, growth is therefore driven by the interaction between exogenous innovation and endogenous scaling decisions. A back-of-the-envelope calculation suggests that this mechanism can have a significant impact on long-run growth.

To study how policy interventions and market frictions affect returns-to-scale decisions, we extend our baseline model to include wedges, such as sales taxes or tariffs on intermediate inputs. We find that such distortions, by incentivizing firms to shrink, lead to the adoption of inefficiently low returns to scale. In this distorted equilibrium, productivity shocks have a first-order effect on welfare by altering the economy's returns to scale. A positive productivity shock, for instance, not only increases output directly but also acts as a corrective force: by lowering input costs, it encourages firms to increase their scalability, moving the economy's production technology closer to the efficient benchmark.

We use detailed data covering the near-universe of firms in Spain to test the core predictions of the model. Consistent with our theory, we document a strong positive correlation between firm productivity, size, and returns to scale, both in the cross-section and within firms over time. We also find evidence supporting the model's input-cost mechanism. By exploiting variation in import tariffs, we show that firms more exposed to costlier intermediate inputs tend to reduce their returns to scale, in line with the model's prediction. More broadly, cross-country patterns corroborate these mechanisms at the aggregate level. Countries where returns to scale are more responsive to productivity—indicating a stronger endogenous scalability mechanism—exhibit higher income per capita. This suggests that endogenous scalability decisions may play a role in long-run economic development.

Finally, to quantify the importance of our mechanism, we calibrate the model to the Spanish economy. We find that endogenous returns to scale are a first-order determinant of economic performance. Eliminating the ability of firms to adjust their scalability reduces the level of GDP by nearly 12% and lowers the long-run growth rate of the economy by 0.8 percentage points. Crucially, these gains are driven by the capacity of high-productivity firms to adopt more scalable technologies. To illustrate this, we examine the impact of size-dependent distortions that are particularly detrimental to large firms. Constructing wedges as in Hsieh and Klenow (2009), we confirm that larger firms face higher effective distortions in the data. We find that

removing these wedges yields welfare gains that are more than twice as large in our model compared to a standard framework. This highlights that policies burdening large firms are particularly costly when they stifle the adoption of highly scalable technologies.

## Literature review

Early work emphasizes the importance of changing returns to scale for economic outcomes. Kuznets (1973) argues that the rise of large-scale firms reflected an adaptive process through which firms learned to coordinate production and distribution across expanding markets. Chandler (1977) documents how managerial hierarchies and integrated production systems enabled firms to realize “economies of scale and scope.” These classic accounts emphasized the importance of changes in scalability for growth. Our work formalizes that idea in a general equilibrium framework.<sup>1</sup>

Since we focus on the aggregate impact of endogenous scalability, we adopt a holistic approach and do not take a stance on the underlying margins firms use to adjust their returns to scale (several are likely at work). In contrast, some recent studies have focused on specific mechanisms. Argente et al. (2025) propose a model of multi-product firms in which standardization increases a firm’s returns to scale. Like us, they find that endogenous scalability leads to fat-tailed firm-size distributions. Engbom et al. (2025) build a model in which entrepreneurs can hire white-collar workers to complete administrative tasks, thereby increasing their returns to scale. They find that the scarcity of skilled labor in developing countries limits this reorganization, and that increasing the supply of skills can explain two-thirds of the shift toward large firms observed during development. Also in the development literature, Gottlieb et al. (2025) propose a model in which firms can choose between high- and low-returns-to-scale technologies. They use the model to explain empirical patterns related to the effect of skill endowments on the firm size distribution.

Smirnyagin (2023) proposes a business cycle model with financial frictions in which firms can choose between two returns to scale levels. Focusing on long-run patterns, Lashkari et al. (2024) document that the decline in IT prices led to an increase in returns to scale in France. To explore the implications of this finding, they propose a model with non-homothetic production in which returns to scale can vary with input factors. Hubmer et al. (2025) use administrative data from Canada and the United States to document that larger firms operate technologies with higher returns

---

<sup>1</sup>Our work also relates to classic models of firm heterogeneity and dynamics (Lucas, 1978; Hopenhayn, 1992). This literature typically assumes exogenous returns to scale.

to scale—a finding that is consistent with our model. They explore the implications of these patterns in an entrepreneurial model with fixed heterogeneous returns to scale and financial frictions.<sup>2</sup>

A distinguishing feature of our work is that we study the impact of input-output linkages on scalability. In doing so, we identify a novel channel through which adjustments in returns to scale propagate through supply chains, reshaping Domar weights throughout the economy. This channel has important implications for the aggregate impact of endogenous scalability and is essential for our long-run growth results.

Our work also relates to different strands of the production network literature (Long and Plosser, 1983; Acemoglu et al., 2012). As in Baqaee and Farhi (2019b), Hulten’s (1978) theorem only provides a first-order approximation to the economy’s response to shocks in our model. We also build on previous work that studies the role of wedges in network economies (Jones, 2011; Baqaee and Farhi, 2019a; Liu, 2019; Bigio and La’O, 2020). Finally, we relate to a literature on production networks that treats production functions as endogenous (Oberfield, 2018; Acemoglu and Azar, 2020; Kopytov et al., 2024a; Kopytov et al., 2024b). We share with this literature the assumption that firms have control over their production technologies. Unlike this literature, we endogenize returns to scale, which yields novel predictions for the firm size distribution and has important aggregate implications.

## 2 A model of endogenous returns to scale

We introduce endogenous returns to scale into a multisector economy with input–output linkages. Each sector produces a differentiated good that can be used for final consumption and as an intermediate input. Within each sector, there is a continuum of firms with heterogeneous productivity. A representative household supplies labor, owns the firms, and consumes a bundle of sectoral goods. Crucially, firms optimally choose their returns to scale as part of their production decisions. Consequently, changes in the environment affect firm-level scalability and, through that channel, macroeconomic aggregates.

### 2.1 Production technology

There are  $N$  goods, each produced by a different sector. Each sector  $i$  consists of a continuum of competitive firms whose mass  $M_i$  is determined by a free-entry condition. Upon paying  $\kappa_i > 0$  units of labor to enter, a firm  $l$  draws a random

---

<sup>2</sup>In many models, firms must pay a fixed cost to operate, and this fixed cost therefore influences the *average* returns to scale of the firm. In contrast, our setup focuses on *marginal* returns to scale, which capture how a marginal increase in size affects the marginal cost of production.

productivity level  $\varepsilon_{il} \sim \text{iid } \mathcal{N}(\mu_i, \sigma_i^2)$  from a normal distribution. The firm can then produce using a Cobb–Douglas technology and, crucially, choose how scalable that technology is. Specifically, if it selects returns to scale  $0 < \eta_{il} < 1$ , firm  $l$ 's output is

$$Q_{il} = F_i(L_{il}, X_{il}, \eta_{il}) := e^{\varepsilon_{il}} A_i(\eta_{il}) \zeta_i(\eta_{il}) \left( L_{il}^{1-\sum_{j=1}^N \alpha_{ij}} \prod_{j=1}^N X_{ij,l}^{\alpha_{ij}} \right)^{\eta_{il}}, \quad (1)$$

where  $L_{il}$  is labor,  $X_{il} = (X_{i1,l}, \dots, X_{iN,l})$  is a vector of intermediate inputs, and  $\zeta_i(\eta_{il})$  is a normalization term to simplify subsequent expressions.<sup>3</sup> The parameter  $\alpha_{ij} \geq 0$  is the intensity of intermediate input  $j$  (with  $1 - \sum_j \alpha_{ij} > 0$ ). We assume that operating technologies with higher returns to scale is costly, so that the productivity shifter  $A_i(\eta_{il})$  is strictly decreasing in  $\eta_{il}$ . This captures the idea that achieving greater scalability often requires more complex processes, incurs significant coordination and communication costs, and demands more managerial attention (Chandler, 1977). We impose that  $A_i$  is smooth, strictly log-concave and that  $A_i(\eta_{il}) \rightarrow 0$  as  $\eta_{il} \rightarrow 1$ .

## 2.2 Firm problem

To explore what drives a firm's returns-to-scale decision, we first analyze its problem in partial equilibrium. A firm  $l$  in sector  $i$  simultaneously chooses its returns to scale  $\eta_{il}$  and its inputs to maximize profits:

$$\Pi_{il} := \max_{\eta_{il}, L_{il}, X_{il}} P_i F_i(L_{il}, X_{il}, \eta_{il}) - W L_{il} - \sum_{j=1}^N P_j X_{ij,l}. \quad (2)$$

where  $W$  is the wage and  $P_j$  is the price of good  $j$ . Solving this problem, we derive the firm's marginal cost of production as a function of its output  $Q_{il}$  and its chosen technology  $\eta_{il}$ .

**Lemma 1.** *The firm's marginal cost of production  $\lambda_{il}$  is given by*

$$\lambda_{il} := \frac{1}{e^{\varepsilon_{il}} A_i(\eta_{il})} H_i^{\eta_{il}} \Pi_{il}^{1-\eta_{il}}, \quad (3)$$

where  $H_i := W^{1-\sum_{j=1}^N \alpha_{ij}} \prod_{j=1}^N P_j^{\alpha_{ij}}$  is the price of the variable input bundle in sector  $i$ , and  $\Pi_{il}$  is profits,<sup>4</sup>

$$\Pi_{il} = (1 - \eta_{il}) \lambda_{il} Q_{il}. \quad (4)$$

As usual, the firm's marginal cost  $\lambda_{il}$  decreases with productivity and increases

---

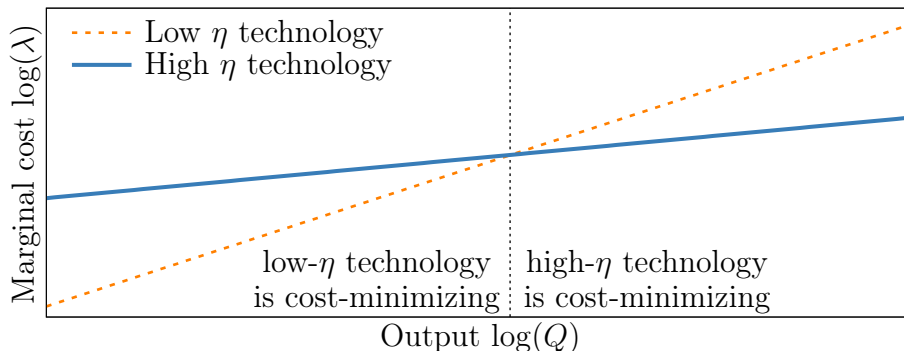
<sup>3</sup>We set  $(\zeta_i(\eta))^{-1} := \left( (1 - \sum_j \alpha_{ij}) \eta \right)^{\eta(1-\sum_j \alpha_{ij})} \prod_j (\eta \alpha_{ij})^{\eta \alpha_{ij}} (1 - \eta)^{1-\eta}$  to simplify the unit cost expression below. Here, this term can be subsumed in  $A_i(\eta_{il})$ .

<sup>4</sup>An Online Supplement is available at the authors' websites. All proofs are in Supplement C.

with input prices. The wage and the price of intermediate inputs, in particular, affect  $\lambda_{il}$  through the variable input bundle price  $H_i$ . Crucially, the firm’s profit  $\Pi_{il}$  also shows up as an input price in (3). To understand why, recall that we can interpret any decreasing-returns production function as a constant-returns technology with an additional fixed entrepreneurial factor in unit supply (McKenzie, 1959). Under this interpretation, profits  $\Pi_{il}$  are simply the payment to that input. As output  $Q_{il}$  increases, the pressure on this fixed factor rises, increasing its shadow cost  $\Pi_{il}$  and, in turn, driving up the marginal cost  $\lambda_{il}$ .

Returns to scale  $\eta_{il}$  play a dual role in shaping marginal costs. Higher returns to scale lower baseline productivity through  $A_i(\eta_{il})$ , but they also increase the weight of the variable bundle (labor and intermediate) relative to the fixed entrepreneurial factor. Through this second channel,  $\eta_{il}$  determines how steeply the marginal cost  $\lambda_{il}$  rises with output  $Q_{il}$ . Figure 1 illustrates this trade-off by plotting  $\lambda_{il}$  as a function of  $Q_{il}$  for a high and a low value of  $\eta_{il}$ . The high- $\eta$  technology offers greater scalability and thus a flatter marginal cost curve, allowing the firm to increase its size with only a small increase in its marginal cost. This makes it particularly effective for large firms. However, because it incurs a large productivity penalty through  $A_i$ , this technology is inefficient at small scales. In contrast, the low- $\eta$  technology benefits from high productivity  $A_i(\eta_{il})$ , making it the preferred choice for small-scale production.

Figure 1: The trade-off between baseline productivity and returns to scale.



This leads to the key sorting mechanism of our model. While adopting a technology with higher returns to scale is costly in terms of baseline productivity, firms that choose these technologies are, in equilibrium, more productive overall. This is because only firms with a sufficiently high idiosyncratic productivity draw  $\varepsilon_{il}$  find it optimal to operate at the large scale necessary to make a high- $\eta$  technology worthwhile. In Section 6, we will show that this positive correlation between productivity and returns to scale is supported by the data.

Finally, note that profit maximization implies that the firm selects output  $Q_{il}$  so that its marginal cost  $\lambda_{il}$  equals the price of its good  $P_i$ .

### 2.3 Choosing returns to scale

Building on these insights, we can characterize the optimal returns-to-scale decision of the firm.

**Lemma 2.** *At an interior solution, the firm chooses its returns to scale  $\eta_{il} \in (0, 1)$  according to*

$$\frac{da_i(\eta_{il})}{d\eta_{il}} = \log H_i - \log \Pi_{il}, \quad (5)$$

where  $a_i(\eta_{il}) := \log A_i(\eta_{il})$ .

Equation (5) describes the key trade-off behind the returns-to-scale choice. It is better understood as the derivative of the log of  $\lambda_{il}$ , given by (3), with respect to  $\eta_{il}$ . When increasing  $\eta_{il}$  at the margin, the firm shifts its input mix away from the fixed entrepreneurial factor, whose shadow cost is  $\Pi_{il}$ , toward the variable input bundle, whose cost is  $H_i$ . The right-hand side of (5) captures the change in cost associated with that shift. The firm balances that change in cost with any loss in TFP associated with the higher returns to scale, as reflected by the left-hand side of (5).

From (5), we can determine how  $\eta_{il}$  responds to changes in the economic environment. The top two panels of Figure 2 illustrate the forces involved. Since  $a_i$  is concave, its derivative is decreasing. Therefore, any change that increases profits  $\Pi_{il}$ —such as a higher output price  $P_i$  or a better productivity draw  $\varepsilon_{il}$ —makes the fixed factor more expensive, pushing the firm to adopt a higher  $\eta_{il}$ . Conversely, an increase in the variable input cost  $H_i$ , as it incentivizes the firm to rely more on its fixed factor, lowers the optimal  $\eta_{il}$ . The following lemma formalizes this intuition.

**Lemma 3.** *At an interior solution, the returns-to-scale parameter  $\eta_{il}$  satisfies<sup>5</sup>*

$$\frac{d\eta_{il}}{d\varepsilon_{il}} = \frac{d\eta_{il}}{d \log P_i} = - \left[ (1 - \eta_{il}) \frac{d^2 a_i}{d\eta_{il}^2} \right]^{-1} > 0, \quad \text{and} \quad \frac{d\eta_{il}}{d \log H_i} = \left[ (1 - \eta_{il}) \frac{d^2 a_i}{d\eta_{il}^2} \right]^{-1} < 0.$$

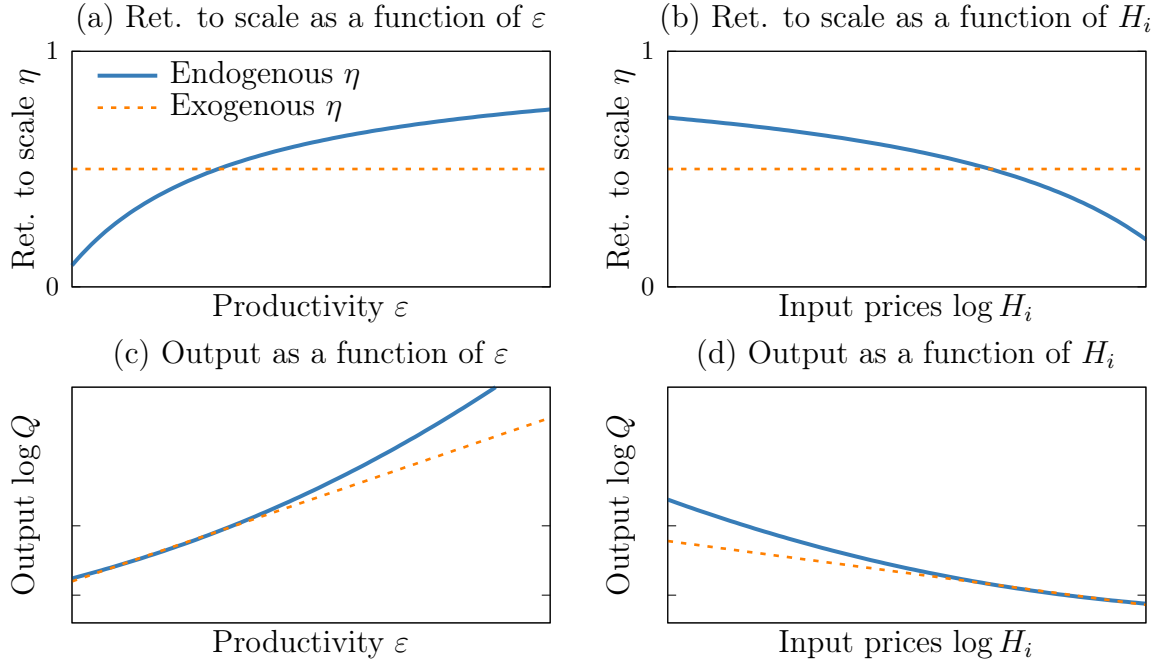
This result highlights that the elasticities of the returns to scale with respect to prices and productivity depend on  $\eta_{il}$  itself and on the concavity of  $a_i$ .

Endogenous scalability also has crucial implications for the output  $Q_{il}$  of the firm.

---

<sup>5</sup>When increasing  $P_i$ , we keep the price of the variable input bundle constant to distinguish the two channels that affect  $\eta_{il}$ .

Figure 2: Impact of productivity  $\varepsilon_{il}$  and input prices  $H_i$  on the firm



**Lemma 4.** *At an interior solution, the elasticity of output  $Q_{il}$  with respect to productivity  $\varepsilon_{il}$  is given by*

$$\frac{d \log Q_{il}}{d \varepsilon_{il}} = \underbrace{\frac{1}{1 - \eta_{il}}}_{\text{Fixed } \eta \text{ effect}} + \underbrace{\frac{1}{1 - \eta_{il}} \frac{d \eta_{il}}{d \varepsilon_{il}}}_{\text{Flexible } \eta \text{ effect}} > 0.$$

*In addition, the elasticities of output  $Q_{il}$  with respect to prices are given by*

$$\frac{d \log Q_{il}}{d \log P_i} = \underbrace{\frac{\eta_{il}}{1 - \eta_{il}}}_{\text{Fixed } \eta \text{ effect}} + \underbrace{\frac{1}{1 - \eta_{il}} \frac{d \eta_{il}}{d \log P_i}}_{\text{Flexible } \eta \text{ effect}} > 0, \text{ and } \frac{d \log Q_{il}}{d \log H_i} = \underbrace{-\frac{\eta_{il}}{1 - \eta_{il}}}_{\text{Fixed } \eta \text{ effect}} + \underbrace{\frac{1}{1 - \eta_{il}} \frac{d \eta_{il}}{d \log H_i}}_{\text{Flexible } \eta \text{ effect}} < 0.$$

*Furthermore, the impact of a change in  $\varepsilon_{il}$ ,  $\log P_i$  or  $\log H_i$  on  $\log Q_{il}$  is amplified because of the endogenous response of  $\eta_{il}$ .*

With fixed returns to scale, productivity and prices affect output  $Q_{il}$  through standard channels, captured by the first terms in the expressions of Lemma 4. Higher productivity  $\varepsilon_{il}$ , for instance, allows the firm to produce larger quantities before its marginal cost reaches the price  $P_i$  of its output. The magnitude of this response depends on returns to scale: a high- $\eta$  firm is more sensitive to productivity and prices than a low- $\eta$  firm.

In addition to this fixed- $\eta$  mechanism, Lemma 4 reveals an additional mechanism at work when returns to scale are endogenous. Following an increase in productivity

$\varepsilon_{il}$ , the firm not only expands to exploit its lower marginal cost but also increases its returns to scale to better accommodate the higher production volume. This amplification mechanism creates a superstar effect, causing high-productivity firms to grow disproportionately large. A similar mechanism operates in response to price changes.

The bottom two panels of Figure 2 illustrate these forces. With exogenous returns to scale (dashed orange lines),  $\log Q_{il}$  varies linearly with  $\varepsilon_{il}$  and  $\log H_i$ , as in standard models. In contrast, with endogenous returns to scale (blue lines), the response is convex: productivity and input prices have an outsized impact on output.

This amplification mechanism has important implications for the firm-size distribution. To explore them transparently, it helps to specialize the returns-to-scale cost function  $a_i$ .<sup>6</sup>

**Assumption 1.** *The TFP shifter function  $A_i$  takes the form*

$$a_i(\eta_{il}) = -\frac{\gamma_i}{1 - \eta_{il}}, \quad (6)$$

where the parameter  $\gamma_i > \sigma_i^2/2$  governs the productivity cost of raising  $\eta_{il}$  in sector  $i$ .

We also define the *effective productivity dispersion*  $\varphi_i := \sigma_i^2/(2\gamma_i)$  as a measure of productivity dispersion in sector  $i$  relative to the cost of adjusting returns to scale. Since  $\gamma_i > \sigma_i^2/2$  by Assumption 1, we have that  $0 < \varphi_i < 1$ .<sup>7</sup> This parameter will play an important role in our analysis.

With that assumption, we can describe the impact of endogenous returns to scale on the tail of the firm-size distribution.

**Proposition 1.** *Suppose that Assumption 1 holds. Without endogenous returns to scale, the distribution of  $Q_{il}$  in sector  $i$  is log-normal. With endogenous returns to scale, the right tail of the distribution of  $Q_{il}$  behaves like a Pareto distribution with tail index  $1/\varphi_i$ , in the sense that*

$$\log(\mathbb{P}(Q_{il} > q)) \sim -\frac{1}{\varphi_i} \log q, \text{ as } q \rightarrow \infty.$$

In the absence of endogenous scalability, all firms within a sector operate with identical returns to scale. Consequently, the distribution of firm output mirrors that of

---

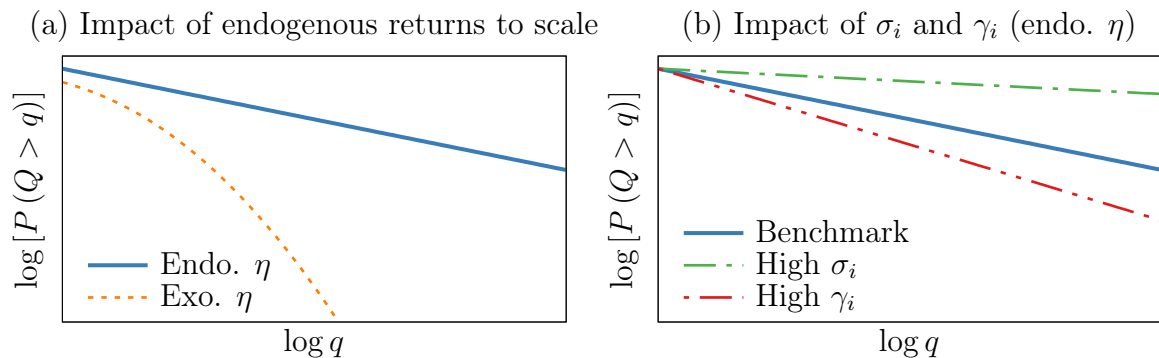
<sup>6</sup>Assumption 1 imposes that  $A_i(\eta_{il})$  satisfies an Inada condition as  $\eta_{il} \rightarrow 1$ , but not as  $\eta_{il} \rightarrow 0$ . Since productivity shocks  $\varepsilon_{il}$  are unbounded, firms with very low  $\varepsilon_{il}$  might choose  $\eta_{il} \notin (0, 1)$ . In Online Supplement D.1, we explore the model in which the distribution of  $\varepsilon_{il}$  is truncated such that  $\eta_{il} \in (0, 1)$  for all firms. We show that aggregate quantities converge to their main-text counterparts as the mass of firms picking  $\eta_{il} \notin (0, 1)$  shrinks. The mass of such firms is indeed small under our calibration.

<sup>7</sup>Without this constraint, returns to scale increase so rapidly with  $\varepsilon_{il}$  that output can be infinite.

the underlying productivity distribution and is log-normal. In contrast, when returns to scale are endogenous, the most productive firms choose higher returns to scale. This stretches the right tail of the distribution, making it thick and Pareto-like. Proposition 1 shows that the tail's thickness depends on the effective productivity dispersion  $\varphi_i$ . When productivity shocks are highly dispersed (large  $\sigma_i^2$ ) or scalability is cheap (low  $\gamma_i$ ), the firm-size distribution is thicker. Figure 3 illustrates these mechanisms.

The thick tail of the firm size distribution is well-documented empirically (Axtell, 2001). Our model generates this property endogenously from fundamental productivity shocks that are not themselves fat-tailed, with superstar firms emerging from the decisions of high-productivity producers to become more scalable.

Figure 3: Tail of the distribution of firm-level output  $Q_{it}$



## 2.4 Household

A representative household owns the firms, supplies  $\bar{L} > 0$  units of labor inelastically, and consumes a bundle  $Y := \prod_{i=1}^N (\beta_i^{-1} C_i)^{\beta_i}$  of the different consumption goods, where  $\sum_{i=1}^N \beta_i = 1$ . Since  $Y$  measures aggregate value added in this economy, we refer to it as (real) GDP.

The household maximizes  $Y$  subject to the budget constraint<sup>8</sup>

$$\sum_{i=1}^N P_i C_i \leq W \bar{L}. \quad (7)$$

Because of the free-entry condition, all profits from the firms are dissipated through entry costs, and the household's only income comes from labor.

Maximization by the household implies that spending on good  $i$  amounts to a fraction  $\beta_i$  of total expenditure, so that  $P_i C_i = \beta_i \bar{P} Y$ , where  $\bar{P} := \prod_{i=1}^N P_i^{\beta_i}$  is the

<sup>8</sup>Because the model is static and there is no aggregate uncertainty, the household could instead maximize a strictly increasing function of  $Y$  without affecting the results.

ideal price index which we adopt as the numeraire. Consequently, nominal and real GDP are equal, and the household's budget constraint simplifies to  $Y = W\bar{L}$ .

## 2.5 Equilibrium conditions

For any firm-level quantity  $B_{il}$ , we denote by  $B_i = \int_0^{M_i} B_{il} dl$  the sum of that quantity across all firms in sector  $i$ . We also use brackets  $\{B_{il}\}$  to denote the set of that quantity over all sectors and firms.

We define an equilibrium as an allocation in which the optimality conditions of the firms and the household hold simultaneously, and all markets clear.

**Definition 1.** An *equilibrium* is a set of prices  $(P^*, W^*)$ , a choice of returns to scale  $\{\eta_{il}^*\}$ , a tuple of quantities  $\{C_i^*, L_{il}^*, X_{il}^*, Q_{il}^*\}$ , and a mass of firms  $M^*$  in each sector such that

1. (Optimal returns-to-scale choice) For each  $i \in \{1, \dots, N\}$  and  $l \in [0, M_i^*]$ , the returns-to-scale decision  $\eta_{il}^*$  solves (2) given prices  $(P^*, W^*)$ .
2. (Optimal input choice) For each  $i \in \{1, \dots, N\}$  and  $l \in [0, M_i^*]$ , factor demands  $L_{il}^*$  and  $X_{il}^*$  solve (2) given prices  $(P^*, W^*)$ .
3. (Consumer optimization) The consumption vector  $C^*$  maximizes GDP  $Y$  subject to (7) given prices  $(P^*, W^*)$ .
4. (Free entry) For each  $i \in \{1, \dots, N\}$ , the expected profit of a potential entrant in sector  $i$  solves

$$E_i [\Pi_i(\varepsilon_{il}, P^*, W^*)] = \kappa_i W^*, \quad (8)$$

where  $\Pi_{il}$  is given by (4), and where the expectation  $E_i$  is taken over  $\varepsilon_{il}$ .

5. (Market clearing) For each  $i \in \{1, \dots, N\}$ ,

$$C_i^* + \sum_{j=1}^N X_{ji}^* = Q_i^* = \int_0^{M_i^*} F_i(L_{il}^*, X_{il}^*, \eta_{il}^*) dl, \text{ and } \sum_{i=1}^N L_i^* + \sum_{i=1}^N M_i^* \kappa_i = \bar{L}. \quad (9)$$

Conditions 2 to 5 are standard and imply that the household and the firms maximize their objective functions, that all markets clear, and that the free-entry condition holds. Condition 1 states that firms pick their returns to scale to maximize profits.

## 3 Aggregation

In this section, we aggregate the economy and derive equations for equilibrium prices and GDP. While most of our firm-level results hold under general  $A_i$ 's, we need

to impose additional restrictions to aggregate the economy in a tractable way. We therefore assume that Assumption 1 holds from now on. Under that assumption, we can derive a tractable mapping between a firm's productivity  $\varepsilon_{il}$  and its returns to scale  $\eta_{il}$  in equilibrium. Indeed, (5) implies that

$$\frac{1}{1 - \eta_{il}} = \frac{1}{2\gamma_i} (\varepsilon_{il} + \log P_i - \log H_i). \quad (10)$$

In the remainder of this section, we take advantage of (10) by first aggregating firms in a sector and then for the whole economy.

### 3.1 Sectoral aggregation

To aggregate the economy, we define the Domar weight of a production unit (a firm or a sector) as the share of its sales in nominal GDP. For a firm  $l$  in sector  $i$  and for sector  $i$  as a whole, those are given by

$$\omega_{il} := \frac{P_i Q_{il}}{\bar{P}Y} \quad \text{and} \quad \omega_i := \frac{P_i Q_i}{\bar{P}Y}.$$

We also introduce the *effective returns to scale*  $\hat{\eta}_i$ , defined as the sales-weighted average of firm-level returns to scale in sector  $i$ :

$$\hat{\eta}_i := \int_0^{M_i} \frac{P_i Q_{il}}{P_i Q_i} \eta_{il} dl. \quad (11)$$

This quantity will play an important role in our analysis. One can show, for instance, that the *sectoral* input cost shares depend on  $\hat{\eta}_i$ :

$$\frac{WL_i}{P_i Q_i} = \hat{\eta}_i \left( 1 - \sum_{j=1}^N \alpha_{ij} \right), \quad \frac{P_j X_{ij}}{P_i Q_i} = \hat{\eta}_i \alpha_{ij}, \quad \text{and} \quad \frac{\Pi_i}{P_i Q_i} = 1 - \hat{\eta}_i.$$

In addition, we can characterize the returns to scale of any firm in sector  $i$  using  $\hat{\eta}_i$ .

**Lemma 5.** *The returns to scale  $\eta_{il}$  of firm  $l$  in sector  $i$  is given by*

$$\frac{1}{1 - \eta_{il}} = \frac{1 - \varphi_i}{1 - \hat{\eta}_i} + \frac{\varepsilon_{il} - \mu_i}{2\gamma_i}. \quad (12)$$

*Furthermore, the moments of the firm-level returns-to-scale distribution are given by*

$$E_i \left[ \frac{1}{1 - \eta_{il}} \right] = \frac{1 - \varphi_i}{1 - \hat{\eta}_i}, \quad V_i \left[ \frac{1}{1 - \eta_{il}} \right] = \frac{\varphi_i}{2\gamma_i}, \quad \text{and} \quad \text{Cov}_i \left[ \frac{1}{1 - \eta_{il}}, \varepsilon_{il} \right] = \varphi_i > 0. \quad (13)$$

Equation (12) links a firm's own returns to scale  $\eta_{il}$  to its productivity  $\varepsilon_{il}$  and the effective sectoral returns to scale  $\hat{\eta}_i$ . For the firm with the median productivity

( $\varepsilon_{il} = \mu_i$ ), this equation simplifies to

$$\hat{\eta}_i = \eta_i(\mu_i) + \varphi_i(1 - \eta_i(\mu_i)). \quad (14)$$

Since  $\varphi_i > 0$ , it follows that the effective returns to scale  $\hat{\eta}_i$  of a sector is larger than that of its median firm. This is because high-productivity firms, which have higher returns to scale and are larger (Lemmas 3 and 4), are weighted more heavily in the calculation of  $\hat{\eta}_i$ .

Equation (14) also shows that the gap between  $\hat{\eta}_i$  and  $\eta_i(\mu_i)$  increases with  $\varphi_i = \sigma_i^2 / (2\gamma_i)$ . Intuitively, greater productivity dispersion  $\sigma_i^2$  implies that there are relatively more high-productivity firms. Lower adjustment costs  $\gamma_i$  also allow these high-productivity firms to adopt more scalable technologies and grow more aggressively, leading to a higher sectoral returns to scale  $\hat{\eta}_i$ .

The second part of Lemma 5 describes the cross-sectional moments of returns to scale within a sector. The first moment shows again that  $\varphi_i$  controls the gap between the median and effective returns to scale. The second moment shows that a higher productivity dispersion  $\sigma_i^2$  and a lower adjustment cost  $\gamma_i$  both contribute to greater cross-sectional dispersion in  $\eta_{il}$ . The third moment confirms that high  $\varepsilon_{il}$  firms choose higher  $\eta_{il}$ . As the endogenous returns-to-scale mechanism shuts down ( $\varphi_i \rightarrow 0$ ), the covariance between productivity and returns to scale goes to zero. Later on, we will rely on that covariance to measure the strength of the mechanism in the data.

Aggregating firms within a sector using the free-entry condition (8) yields the following results.<sup>9</sup>

**Proposition 2.** *The marginal cost of sector  $i$  is given by*

$$\lambda_i = \frac{1}{Z_i(\hat{\eta}_i)} W^{1-\hat{\eta}_i} \sum_{j=1}^N \alpha_{ij} \prod_{j=1}^N P_j^{\hat{\eta}_i \alpha_{ij}}, \quad (15)$$

where sectoral total factor productivity  $Z_i(\hat{\eta}_i)$  is defined as

$$\log Z_i(\hat{\eta}_i) := \underbrace{\mu_i + a_i(\hat{\eta}_i) + \frac{\sigma_i^2}{2} \frac{1}{1-\hat{\eta}_i}}_{\text{Exogenous returns to scale}} + \underbrace{\frac{1}{2}(1-\hat{\eta}_i) \log\left(\frac{1}{1-\varphi_i}\right)}_{\text{Superstar effect}} - \underbrace{(1-\hat{\eta}_i) \log \kappa_i}_{\text{Entry cost}}. \quad (16)$$

<sup>9</sup>Since all firms in a sector face the same output price, they have the same marginal cost through profit maximization. We therefore define the marginal cost  $\lambda_i$  of sector  $i$  as the marginal cost of any firm in that sector. Equivalently, since, as we show later, the economy is efficient, one can write the cost minimization problem of the sector and find the same expression.

Furthermore, the effective returns to scale  $\hat{\eta}_i$  is given by

$$\frac{1}{1 - \hat{\eta}_i} = \frac{1}{2\gamma_i(1 - \varphi_i)} (\mu_i + \log P_i - \log H_i). \quad (17)$$

The sectoral marginal cost  $\lambda_i$  takes the standard form associated with a Cobb–Douglas production function with two factors: labor, which is used for production and entry costs, and intermediate inputs. As the cost shares of those inputs sum to one, the sector as a whole exhibits *constant* returns to scale. Intuitively, free entry acts as an adjustment margin: while individual producers operate under decreasing returns, the entry of new firms allows the sector to expand to achieve constant returns.

The effective returns to scale  $\hat{\eta}_i$  is the main driver of the sectoral cost shares. Specifically, the sectoral labor share is the sum of labor used for entry (share  $1 - \hat{\eta}_i$ ) and labor used for production (share  $\hat{\eta}_i \left(1 - \sum_j \alpha_{ij}\right)$ ). These sum to a total labor share of  $1 - \hat{\eta}_i \sum_j \alpha_{ij}$ , which is decreasing in  $\hat{\eta}_i$ . Intuitively, a higher  $\hat{\eta}_i$  implies that firms are larger and more scalable, so fewer entrants are required to produce a given output. This economizes on fixed entry costs—which are paid in labor—and shifts the input mix toward the variable bundle, which includes intermediate goods. Consequently, as scalability increases, the sector becomes less labor-intensive.

Equation (16) characterizes the sector’s total factor productivity  $Z_i(\hat{\eta}_i)$ . As expected, sectoral productivity depends on the mean firm-level productivity  $\mu_i$  and the productivity cost  $a_i(\hat{\eta}_i)$  associated with the effective returns to scale. However, firm heterogeneity also plays a role. The third term in (16) captures a standard *dispersion effect*: as more productive firms grow larger, they receive a larger share of input factors, raising sectoral productivity. These first three terms would also appear in an exogenous returns-to-scale model in which all firms share a common  $\eta_{il} = \hat{\eta}_i$ . The fourth term in (16), however, captures a novel amplification channel from endogenous returns to scale. In our model, high-productivity firms not only produce more but also choose more scalable technologies, which allows them to grow even larger. This *superstar effect* amplifies their contribution to sectoral productivity beyond the standard dispersion effect. Finally, the last term in (16) captures the role of entry costs. A higher entry cost  $\kappa_i$  diverts labor away from production, lowering sectoral productivity, with the magnitude of this loss determined by the importance of the fixed factor,  $1 - \hat{\eta}_i$ .

Proposition 2 also provides an expression for the sector’s effective returns to scale. This expression is analogous to the one determining firm-level returns to scale, given by (10), but it includes the adjustment term  $\varphi_i$  to account for firm heterogeneity. This adjustment reflects that larger, more productive firms have a disproportionate

impact on the aggregate measure of returns to scale  $\hat{\eta}_i$ .

### 3.2 Prices and GDP

We can now aggregate the economy to derive expressions for prices and GDP. To do so, we define the sectoral Leontief inverse matrix  $\mathcal{L} := (I - \text{diag}(\hat{\eta})\alpha)^{-1}$ , where  $\text{diag}(\hat{\eta})$  is the diagonal matrix with the effective returns-to-scale vector  $\hat{\eta}$  on the main diagonal. An element  $\mathcal{L}_{ij}$  of this matrix captures the importance of sector  $j$  in the production of good  $i$ , taking into account direct and indirect connections through the production network. For example,  $\mathcal{L}_{ij}$  is large if sector  $i$  spends a large share of its costs on inputs from  $j$  (i.e.,  $\hat{\eta}_i\alpha_{ij}$  is large), or if  $i$  relies on another sector  $k$  that, in turn, relies heavily on  $j$ , and so on.

We show in Online Supplement C.1 that the Leontief inverse can be used to write the sectoral Domar weights as

$$\omega_i := \frac{P_i Q_i}{\bar{P} Y} = \beta^\top \mathcal{L} \mathbf{1}_i, \quad (18)$$

where  $\mathbf{1}_i$  is the  $i$ th standard basis vector. As usual in network economies, the Domar weight  $\omega_i$  provides a measure of the importance of sector  $i$ . A sector  $i$  has a large Domar weight if its output is heavily demanded, either directly by the household (high  $\beta_i$ ), or indirectly by other sectors favored by the household (high  $\beta_j \mathcal{L}_{ji}$ ).

Sectoral returns to scale  $\hat{\eta}$  play an important role in shaping the production network. Intuitively, when a downstream producer increases its returns to scale, it effectively shifts its input mix toward intermediate goods, thereby increasing its demand for upstream suppliers. This strengthens the input-output linkages and raises the Domar weights of those suppliers. As we will show, this mechanism has important implications for the impact of endogenous returns to scale on the macroeconomy.

We can now characterize equilibrium prices and GDP.

**Proposition 3.** *The equilibrium price vector  $P = (P_1, \dots, P_N)$  satisfies*

$$\log \frac{P}{\bar{W}} = -\mathcal{L}(\hat{\eta}) z(\hat{\eta}), \quad (19)$$

where  $z(\hat{\eta}) = (\log Z_1(\hat{\eta}_1), \dots, \log Z_N(\hat{\eta}_N))$  is the vector of log sectoral productivities (16). Furthermore, equilibrium log GDP  $y := \log Y$  is given by

$$y(\hat{\eta}) = \underbrace{[\omega(\hat{\eta})]^\top z(\hat{\eta})}_{\text{Aggregate productivity}} + \underbrace{\log \bar{L}}_{\text{Labor endowment}} \quad (20)$$

In equilibrium, prices must equal marginal production costs ( $P_i = \lambda_i$ ). This condition, combined with Proposition 2, allows us to solve for the vector of sectoral

prices as (19). Intuitively, the price of good  $i$  is low if its key suppliers—both direct and indirect, as captured by the  $i$ -th row of  $\mathcal{L}$ —have high productivity  $z$ .

Equation (20) shows that the contribution of a sector’s productivity  $z_i$  to GDP is proportional to its Domar weight  $\omega_i$ , as in standard network economies. A key feature of our model, however, is that both  $\omega$  and  $z$  depend on the endogenous effective returns to scale  $\hat{\eta}$ . We explore below the role played by  $\hat{\eta}$  in shaping GDP.

### 3.3 Equilibrium existence, uniqueness and efficiency

The preceding analysis describes key equilibrium objects, such as prices and GDP, as functions of the vector of effective returns to scale  $\hat{\eta}$ . To solve for  $\hat{\eta}$  itself and characterize how it responds to changes in the environment, it is convenient to rely on the problem of a social planner. Since there is a single representative household in the economy, the planner’s problem is to maximize that household’s utility (GDP) subject to the physical constraints of the environment. The following result characterizes that problem and its relation to the set of equilibria.

**Proposition 4.** *There exists a unique equilibrium, and it is efficient. Furthermore, the equilibrium vector of effective returns to scale  $\hat{\eta}$  maximizes GDP  $y(\hat{\eta})$ .*

The proof of this proposition establishes an equivalence result between the set of equilibria and the set of efficient allocations. It further shows that since there exists a unique efficient allocation, there is also a unique equilibrium. We can then use the first-order conditions of the planner to find the equilibrium  $\hat{\eta}$ . With that object in hand, the returns to scale of all the firms can be recovered using (12).

## 4 Forces shaping returns-to-scale decisions

In this section, we study how changes in the environment affect returns to scale in equilibrium. To do so, we rely on the fact that the equilibrium is efficient, and that the effective returns-to-scale vector  $\hat{\eta}$  maximizes GDP. The first-order condition associated with that problem is<sup>10</sup>

$$\underbrace{\omega_i \alpha_i^\top \mathcal{L} z}_{\text{Network adjustment } d\omega^\top/d\hat{\eta}_i} + \omega_i \underbrace{\left[ \frac{da_i}{d\hat{\eta}_i} + \frac{\sigma_i^2}{2} \frac{1}{(1 - \hat{\eta}_i)^2} - \frac{1}{2} \log \left( \frac{1}{1 - \varphi_i} \right) + \log \kappa_i \right]}_{\text{Productivity adjustment } dz_i/d\hat{\eta}_i} = 0. \quad (21)$$

A marginal increase in  $\hat{\eta}_i$  has two effects on the economy. First, it makes intermediate inputs more important production factors. As a result, the network becomes more connected and Domar weights increase. Through this channel, captured by the first

<sup>10</sup>See the proof of Proposition 6 for a derivation.

term in (21), the productivities  $z$  of the sectors that are upstream of  $i$  have a larger impact on GDP.

Furthermore, increasing  $\hat{\eta}_i$  affects the productivity  $z_i$  of sector  $i$  itself, as captured by the term in square brackets in (21). The first term there captures the additional productivity cost  $a_i$  of selecting higher returns to scale  $\hat{\eta}_i$ . In addition, when returns to scale are larger, production can move more easily to the most productive firms within the sector. This amplified dispersion effect is captured by the second term between brackets. Higher returns to scale also affect the importance of returns-to-scale *heterogeneity* within a sector. When  $\hat{\eta}_i$  is close to 1, there is less room for the most productive firms to increase their returns to scale and, through that channel, reach a larger size. This effect is captured by the third term between brackets. Finally, a higher  $\hat{\eta}_i$  means that firms can scale up more freely and thus produce more. Consequently, fewer firms enter, and sectoral productivity benefits from a reduction in total entry costs. This effectively leads to an increase in sectoral productivity, as the last term in (21) shows.

## 4.1 Sectoral productivity

We now analyze how changes in the environment affect equilibrium returns to scale. Two prices play a key role in this process. The first is the price  $H_i$  of the variable input bundle. Indeed, when inputs are expensive (high  $W$  or high  $P$ ) individual firms move toward technologies with lower returns to scale. But the wage  $W$  also plays an additional role at the sector level. When labor is expensive, entry is costly, and few firms enter. This allows incumbents to expand and encourages the adoption of higher returns to scale. Combining these two channels, we find that the ratio  $H_i/W$  is key in determining how returns to scale evolve at the sector level.<sup>11</sup>

To capture how changes in the environment affect this ratio, we define the input-price sensitivity matrix  $\mathcal{K}$  with typical element

$$\mathcal{K}_{ij} := \frac{\partial}{\partial z_j} \log \left( \frac{H_i}{W} \right) \leq 0,$$

where the partial derivative holds  $\hat{\eta}$  fixed. The matrix  $\mathcal{K}$  summarizes how productivity shocks propagate through the network to affect the input cost ratio. Using the pricing equation (19), we can show that  $\mathcal{K} = -\alpha\mathcal{L}$ . Since higher productivity  $z_j$  lowers the production cost of all of  $j$ 's downstream customers, the elements of  $\mathcal{K}$  are non-positive, with  $\mathcal{K}_{ij} < 0$  whenever sector  $j$  is a direct or indirect upstream supplier to sector  $i$

---

<sup>11</sup>One can show that  $H_i/W$  is the quantity that is raised to the power  $\hat{\eta}_i$  in the marginal cost expression (15), which explains its importance for scalability decisions.

(i.e.,  $\mathcal{L}_{ij} > 0$ ). The matrix  $\mathcal{K}$  plays an important role in our analysis and depends crucially on the input-output structure of the economy. Without network connections,  $\mathcal{K} = 0$  and several of the mechanisms that we explore below disappear.

We also introduce a measure of how responsive returns to scale are to changes in the environment. In the planner's first-order condition (21), the term  $dz_i/d\hat{\eta}_i$  represents the marginal productivity benefit of increasing scalability  $\hat{\eta}_i$ . Consequently, the sensitivity of the optimal  $\hat{\eta}_i$  to changes in the environment depends on the curvature of this function. We therefore define

$$\Psi_i := \frac{d^2 z_i}{d\hat{\eta}_i^2} = (1 - \varphi_i) \frac{d^2 a_i}{d\hat{\eta}_i^2} < 0,$$

where the second equality follows directly from (16). The inverse  $\Psi_i^{-1}$  captures how elastic returns to scale are in sector  $i$ , with a large  $|\Psi_i^{-1}|$  implying that returns to scale are flexible and respond strongly to changes in the environment.

From these definitions, we can characterize how  $\mu_j$  and  $\sigma_j^2$  affect returns to scale.

**Lemma 6.** *An increase in average productivity  $\mu_j$  increases returns to scale in all downstream sectors, such that*

$$\frac{d\hat{\eta}_i}{d\mu_j} = \Psi_i^{-1} \mathcal{K}_{ij} \geq 0. \quad (22)$$

Furthermore, the impact of productivity dispersion  $\sigma_j^2$  on  $\hat{\eta}_i$  with  $i \neq j$  is given by<sup>12</sup>

$$\frac{d\hat{\eta}_i}{d\sigma_j^2} = \Psi_i^{-1} \mathcal{K}_{ij} \frac{\partial z_j}{\partial \sigma_j^2} \geq 0, \quad (23)$$

where  $\frac{\partial z_j}{\partial \sigma_j^2} = \frac{1}{2(1-\hat{\eta}_j)} + \frac{1-\hat{\eta}_j}{4\gamma_j(1-\varphi_j)} > 0$ . In particular,  $d\hat{\eta}_i/d\sigma_j^2 \geq 0$  for  $i \neq j$ .

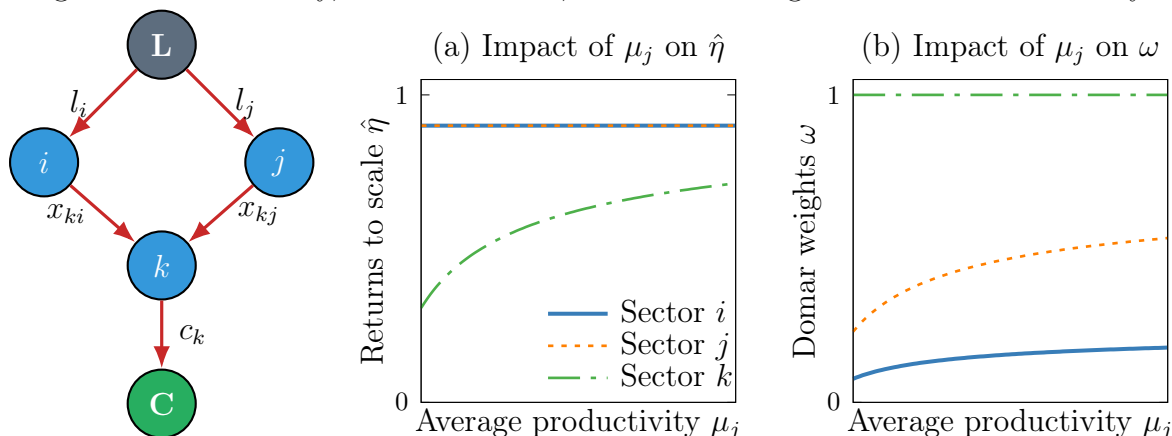
Consider (22) first. An increase in  $\mu_j$  makes firms in sector  $j$  more productive, which lowers the price  $P_j$  through competition. If sector  $i$  is a downstream customer of  $j$  ( $\mathcal{L}_{ij} > 0$ ), this lowers the price of its variable input bundle by an amount proportional to  $|\mathcal{K}_{ij}|$ . This in turn pushes firms in sector  $i$  to increase their returns to scale to take advantage of the cheaper intermediate inputs (Proposition 2). The magnitude of this response depends on how elastic  $\hat{\eta}_i$  is, as given by  $\Psi_i^{-1}$ . Both  $\sigma_i^2$  and  $\gamma_i$  influence this elasticity through  $\varphi_i$ . Specifically, if  $\varphi_i$  is large, which is the case if productivity is dispersed (high  $\sigma_i^2$ ) or scalability is cheap (low  $\gamma_i$ ), the response is stronger. Intuitively, in such sectors, there is a larger mass of high-productivity firms able to aggressively scale up in response to cheaper inputs.

---

<sup>12</sup>In Supplement C.11, we provide a richer version of Lemma 6 that also covers the case  $i = j$ .

**Example.** Consider the economy depicted in the left panel of Figure 4. Since sector  $k$  is downstream from  $j$ , an increase in  $\mu_j$  reduces the price of  $k$ 's input bundle. In response, firms in sector  $k$  increase their returns to scale  $\hat{\eta}_k$ , as panel (a) shows. In contrast, since sector  $i$  is not downstream from  $j$ , input prices in sector  $i$  are unchanged and so are its returns to scale  $\hat{\eta}_i$ .

Figure 4: Productivity, returns to scale, and Domar weights in a vertical economy



The second part of Lemma 6 shows that an increase in  $\sigma_j^2$  affects returns to scale similarly to a productivity shock  $\mu_j$ . Indeed, a higher  $\sigma_j^2$  raises sectoral productivity  $z_j$  through two channels. First, it strengthens the standard dispersion mechanism: even with fixed returns to scale, a higher variance  $\sigma_j^2$  reallocates market share to more productive firms. Second, it amplifies the superstar effect: with more variation in productivity, the most productive firms adopt even higher returns to scale, further boosting sectoral efficiency. This increase in  $z_j$  lowers output prices, triggering a downstream increase in returns to scale  $\hat{\eta}_i$  analogous to the response to a  $\mu_j$  shock.

Lemma 6 focuses on changes in the productivity process, but it is straightforward to derive similar results for changes in entry costs  $\kappa_j$  and the cost of scalability  $\gamma_j$ . We provide these results in Online Supplement D.2. In a nutshell, an increase in the entry cost  $\kappa_j$  in sector  $j$  always reduces the effective returns to scale of any other downstream sector  $i \neq j$  since it leads to an increase in the price of good  $j$ . An increase in the cost of scalability  $\gamma_j$  naturally leads to a reduction in returns to scale in sector  $j$ . The price of good  $j$  increases as a result, leading to lower returns to scale in other downstream sectors as well.

## 4.2 Implications for Domar weights

The comparative statics results derived so far describe how the equilibrium returns to scale  $\hat{\eta}$  respond to the environment. These movements in  $\hat{\eta}$  are important, in part,

because they directly alter the sectoral Domar weights  $\omega$  and, therefore, GDP.

**Lemma 7.** *The impact of a parameter  $\chi \in \{\mu_j, \sigma_j^2, \kappa_j, \gamma_j\}$  on Domar weights is*

$$\frac{d\omega_i}{d\chi} = - \sum_{k=1}^N \mathcal{K}_{ki} \omega_k \frac{d\hat{\eta}_k}{d\chi}. \quad (24)$$

*Proof.* Follows directly from differentiating the Domar weight expression (18).  $\square$

This result shows that if a shock to  $\chi$  leads sector  $k$  to increase its returns to scale ( $\hat{\eta}_k \uparrow$ ), the Domar weights of all sectors *upstream* of  $k$  increase as well. To illustrate this mechanism, consider again the example in Figure 4.

**Example.** Recall that an increase in  $\mu_j$  raises productivity in sector  $j$ , lowering  $P_j$  through competition. Since  $k$  is a downstream customer of  $j$ , firms in  $k$  respond by increasing their returns to scale  $\hat{\eta}_k$  to take advantage of the cheaper input (Lemma 6). Lemma 7 shows that this shift raises the Domar weights of all of  $k$ 's suppliers (panel (b) in Figure 4). Intuitively, as sector  $k$  scales up, it becomes more input-intensive, increasing its demand for upstream goods. Consequently, the sales and the Domar weights of sectors  $i$  and  $j$  rise, reflecting their increased centrality to aggregate production. Notably, sector  $i$ 's importance grows even though its own returns to scale are unaffected by the shock.

This example highlights a key feature of our model: while productivity shocks propagate *downstream* to affect returns to scale, adjustments in returns to scale propagate *upstream* to reshape the network structure and alter Domar weights. As we show next, this interaction affects how GDP responds to changes in the environment.

## 5 Endogenous returns to scale and GDP

Endogenous scalability has important implications for GDP and welfare. It raises the level of GDP, amplifies the impact of beneficial shocks while dampening that of adverse shocks, and increases the long-run growth rate of the economy. We explore these mechanisms in this section.

Throughout our analysis, we compare the equilibrium of our model to a counterfactual economy in which returns to scale are exogenously fixed.

**Definition 2** (Fixed returns-to-scale economy). Let the equilibrium effective returns-to-scale vector in the baseline model be  $\hat{\eta}$ . The *fixed returns-to-scale economy* is an otherwise identical economy in which the returns to scale of all firms are exogenously set to  $\eta_{il} = \hat{\eta}_i$  for all  $i$  and  $l$ . All other quantities are chosen optimally. Furthermore, the returns to scale  $\{\eta_{il}\}$  are fixed and do not respond to changes in the environment.

Endogenous returns to scale do two main things in our model: They create dispersion in  $\eta_{il}$  within a sector, and they allow  $\eta_{il}$  to respond to changes in parameters. The fixed returns-to-scale economy, by construction, shuts down both of these channels. By comparing our baseline model to this counterfactual, we can therefore isolate the full impact of endogenous scalability on economic outcomes.<sup>13</sup> Note that by construction, sectoral Domar weights are the same in both economies. In what follows, we use  $\tilde{\cdot}$  to denote quantities in the fixed returns-to-scale economy.

## 5.1 Endogenous returns to scale and the level of GDP

We first examine the impact of endogenous returns to scale on GDP by comparing its level in the baseline and the fixed returns-to-scale economies. Since both economies share the same Domar weights, any difference in GDP must arise from differences in sectoral productivity. Comparing that quantity in the baseline model ( $Z$ ) with its counterpart in the fixed returns-to-scale economy ( $\tilde{Z}$ ), we find that<sup>14</sup>

$$\log Z_i(\hat{\eta}_i) - \log \tilde{Z}_i(\hat{\eta}_i) := \frac{1}{2}(1 - \hat{\eta}_i) \log \left( \frac{1}{1 - \varphi_i} \right) > 0. \quad (25)$$

Thus, sectoral productivity is always larger in the model with endogenous returns to scale. Intuitively, when returns to scale are fixed, high- $\varepsilon_{il}$  firms are no longer able to adjust their scalability to take advantage of their high productivity. This limits how much they produce, and sectoral productivity falls as a result. The superstar effect is completely shut down in that case. Equation (25) shows that the difference between the two economies is particularly pronounced when the effective dispersion  $\varphi_i$  is large and the effective returns to scale  $\hat{\eta}_i$  is low. In those circumstances, highly productive firms in the baseline model can deviate strongly from  $\hat{\eta}_i$  and thus contribute more to sectoral productivity.

The next result characterizes the aggregate impact of endogenous returns to scale.

**Proposition 5.** *The difference in log GDP between the baseline model and the fixed returns-to-scale economy is given by*

$$y - \tilde{y} = \sum_{i=1}^N \omega_i \frac{1}{2} (1 - \hat{\eta}_i) \log \left( \frac{1}{1 - \varphi_i} \right) > 0. \quad (26)$$

The gain in GDP from endogenous returns to scale is simply the Domar-weighted gain in sectoral productivity. We see from (26) that  $y - \tilde{y}$  is particularly large if the sectors in which high-productivity firms can raise their returns to scale more easily

<sup>13</sup>In Section 7, we disentangle the impact of these two channels in the calibrated economy.

<sup>14</sup>See the proof of Proposition 5 for the derivation.

(high  $\varphi_i$  and low  $\hat{\eta}_i$ ) are also important suppliers (high  $\omega_i$ ). In the calibrated economy of Section 7, we will see that the impact of endogenous returns to scale on the level of GDP can be sizable.

## 5.2 How GDP responds to changes in the environment

In addition to its impact on the level of GDP, endogenous returns to scale also affect how GDP responds to changes in the environment.

**Proposition 6.** *In equilibrium, the following holds.*

1. *An increase in average productivity  $\mu_j$  raises GDP:*

$$\frac{dy}{d\mu_j} = \frac{\partial y}{\partial \mu_j} = \omega_j > 0. \quad (27)$$

2. *An increase in productivity dispersion  $\sigma_j^2$  raises GDP:*

$$\frac{dy}{d\sigma_j^2} = \frac{\partial y}{\partial \sigma_j^2} = \omega_j \left( \frac{1}{2} \frac{1}{1 - \hat{\eta}_j} + \frac{1 - \hat{\eta}_j}{4\gamma_j} \frac{1}{1 - \varphi_j} \right) > 0. \quad (28)$$

3. *An increase in entry cost  $\kappa_j$  lowers GDP:*

$$\frac{dy}{d \log \kappa_j} = \frac{\partial y}{\partial \log \kappa_j} = -\omega_j (1 - \hat{\eta}_j) < 0.$$

4. *An increase in the returns-to-scale productivity cost  $\gamma_j$  lowers GDP:*

$$\frac{dy}{d\gamma_j} = \frac{\partial y}{\partial \gamma_j} = -\omega_j \left( \frac{1}{1 - \hat{\eta}_j} + \frac{1 - \hat{\eta}_j}{2\gamma_j} \frac{\varphi_j}{1 - \varphi_j} \right) < 0. \quad (29)$$

*In these expressions, the partial derivatives indicate that returns to scale  $\{\eta_{il}\}$  are taken as fixed.*

*Proof.* The result follows directly from the envelope theorem. □

Since the equilibrium  $\hat{\eta}$  maximizes GDP, any *marginal* adjustment in returns to scale must have no impact on GDP. This implies that GDP responds to marginal changes in the environment *as if* returns to scale were fixed. Consequently, Hulten's (1978) theorem applies: the first-order impact of a productivity shock  $d\mu_j$  is simply the Domar weight  $\omega_j$  of the affected sector. Note, however, that Domar weights themselves are endogenous in our model and depend on the incentives shaping returns to scale. Similarly, the impact of a change in entry costs  $d \log \kappa_j$  is determined by its direct effect on sector  $j$ 's productivity  $z_j$ , captured by  $1 - \hat{\eta}_j$ , weighted by that sector's importance,  $\omega_j$ .

The impacts of  $\sigma_j^2$  and  $\gamma_j$  operate in a similar way, but here endogenous returns to scale features more prominently. Consider first an increase in  $\sigma_j^2$ . This has two positive effects on sector  $j$ 's productivity  $z_j$ . First, the higher variance in  $\varepsilon_{il}$  implies a larger mass of very productive firms, with positive consequences for GDP. This dispersion effect is captured by the term  $\frac{1}{2} \frac{\omega_j}{1-\hat{\eta}_j}$  in (28) and would be at work even without within-sector dispersion in  $\eta_{il}$  (i.e., in the fixed returns-to-scale economy). Second, increasing  $\sigma_j^2$  interacts with the superstar effect. Since high- $\varepsilon_{il}$  firms already have high returns to scale, the increase in dispersion has a disproportionate impact on them. This effect is captured by the remaining term in (28). Overall, increasing  $\sigma_j^2$  always has a positive effect on GDP, and endogenous returns to scale makes that effect larger, even to a first order, as it creates dispersion in within sector firm-level returns to scale.

Conversely, an increase in the cost of adjusting returns to scale  $\gamma_j$  has two adverse effects on GDP. The first effect is mechanical: a higher  $\gamma_j$  directly increases the average productivity cost  $-a_j(\hat{\eta})$  in sector  $j$ , which lowers GDP. This effect is captured by the first term in (29). Second,  $\gamma_j$  interacts with the superstar effect. Since the biggest producers have the highest returns to scale, they suffer particularly strongly from an increase in  $\gamma_j$ . Indeed, recall that  $a_j(\eta_{jl}) = -\gamma_j/(1-\eta_{jl})$  such that for  $\eta_{jl} \approx 1$ , a marginal increase in  $\gamma_j$  has a particularly severe impact on productivity. This effect is captured by the second term in (29). Proposition 6 shows that increasing  $\gamma_j$  always hurts  $y$ , and all the more so when the superstar effect is stronger.

### 5.3 Second-order impact of productivity shocks

In our baseline model, the equilibrium is efficient, which implies that GDP responds to infinitesimal productivity shocks as if returns to scale were held fixed. For larger shocks, however, the Domar weights themselves respond, leading to a further impact on GDP.

**Proposition 7.** *The response of log GDP  $y$  to a shock  $\Delta\mu_i$  is given by*

$$\Delta y = \omega_i \Delta\mu_i + \frac{1}{2} \frac{d\omega_i}{d\mu_i} (\Delta\mu_i)^2 + o((\Delta\mu_i)^2). \quad (30)$$

*Furthermore, the second-order term is non-negative,  $d\omega_i/d\mu_i \geq 0$  and given by (24).*

Equation (30) provides a second-order approximation of the GDP response to a shock  $\Delta\mu_i$ . While the first-order term is the standard Domar-weight effect from Hulten's theorem, the second-order term captures a novel channel driven by the endogenous adjustment in returns to scale. That channel operates through the response

of Domar weights to the shocks.

The intuition is straightforward. Recall that a positive productivity shock in sector  $i$  propagates downstream, lowering input costs for  $i$ 's customer sectors. These sectors, in turn, are incentivized to increase their own returns to scale to capitalize on the cheaper inputs. This shift towards greater scalability makes the production network more reliant on sector  $i$ , which increases its Domar weight. As a result, the derivative  $d\omega_i/d\mu_i$  is positive. This implies that the second-order term is always positive, adding *convexity* to the GDP response. In other words, endogenous returns to scale amplify the impact of positive productivity shocks and dampen the adverse impact of negative shocks.

To identify the sources of this amplification and dampening, we can use Lemma 6 to express the second-order coefficient as

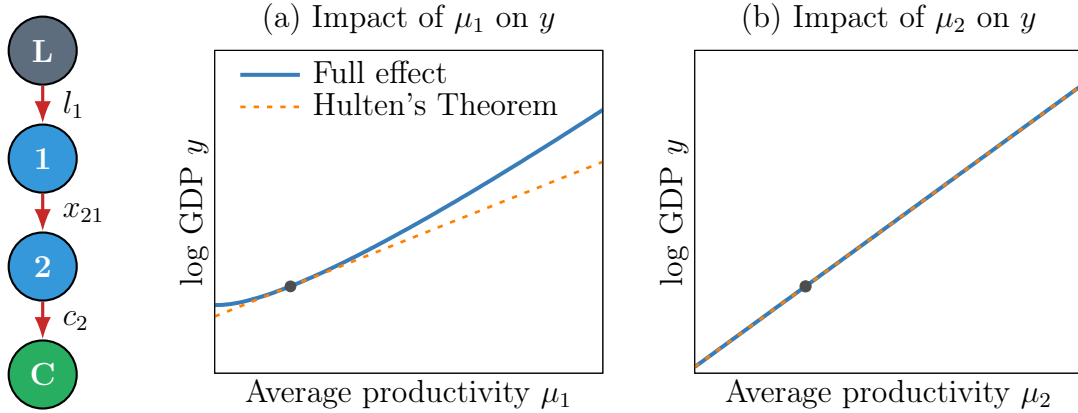
$$\frac{d\omega_i}{d\mu_i} = - \sum_{k=1}^N \omega_k \Psi_k^{-1} \mathcal{K}_{ki}^2 \geq 0. \quad (31)$$

The second-order effect is therefore stronger when the shock hits a sector  $i$  that is a key supplier (large  $|\mathcal{K}_{ki}|$ ) to sectors that are large (high  $\omega_k$ ) and have elastic returns to scale (low  $|\Psi_k|$ , meaning a low cost of adjusting scalability). Those sectors more strongly increase their cost share of good  $i$  after the increase in  $\mu_i$ , contributing to a larger increase in  $\omega_i$ . Crucially, it is the square of  $\mathcal{K}_{ki}$  that shows up in (31), implying that the convex response of GDP is disproportionately stronger for shocks to the most important suppliers.

Proposition 7 implies that whether shocks hit sectors that are upstream or downstream in the supply chain matters for their impact on GDP. The following example illustrates the mechanism.

**Example.** Consider the vertical economy depicted in Figure 5. An upstream sector (Sector 1) sells its entire output to a downstream sector (Sector 2), which in turn sells to the final consumer. Panel (a) shows the impact of a productivity shock  $\mu_1$  to the *upstream* sector. The solid line, representing the response of GDP, is clearly convex and lies above the linear, first-order approximation from Hulten's theorem (dashed line). This illustrates the amplification effect: as  $\mu_1$  increases, the price of good 1 falls, inducing firms in the downstream Sector 2 to increase their returns to scale to capitalize on cheaper inputs. This change in production processes makes Sector 1 a more important supplier, increasing its Domar weight and thus magnifying the aggregate benefit of its higher productivity. Panel (b) shows a starkly different result for a productivity shock  $\mu_2$  to the *downstream* sector. In this case, the full GDP

Figure 5: The impact of sectoral productivity on GDP in a vertical economy



response is linear and coincides with the Hulten’s theorem approximation. Because Sector 2 is at the bottom of the supply chain, a fall in its price provides no benefit to any other sectors. There are no downstream customers to re-optimize their scalability choices, and thus no structural amplification. This example shows a key implication of our model: the macroeconomic impact of a productivity shock depends crucially on a sector’s position in the production network.

## 5.4 Extension: The role of wedges

So far, our analysis has focused on an efficient equilibrium where the envelope theorem holds, meaning that adjustments in  $\hat{\eta}$  have only second-order effects on GDP. We now show that in the presence of frictions, markups or other distortions, this is no longer the case, and changes in returns to scale can have first-order effects on GDP.

We consider a setup with general wedges in Online Supplement D.5, but to illustrate the forces at work transparently, we focus here on distortionary sales wedges. Specifically, we assume that firm  $l$  in sector  $i$  retains only a fraction  $1 - \tau_i^S$  of its revenue, where  $\tau_i^S \in [0, 1)$  is the wedge. From the firm’s perspective, these wedges are equivalent to a reduction in productivity, leading to an inefficient adjustment in scalability. We assume that the proceeds from the wedges are fully rebated lump-sum to the household, ensuring that they act as pure distortions with no direct loss in resources.

Sales wedges directly affect scalability choices, as the following lemma shows.

**Lemma 8.** *An increase in  $\tau_j^S$  decreases the returns to scale in all downstream sectors:*

$$\frac{d\hat{\eta}_i}{d\tau_j^S} = -\frac{1}{1 - \tau_j^S} \Psi_i^{-1} \mathcal{K}_{ij} \leq 0. \quad (32)$$

Intuitively, the sales wedge acts like a markup of  $(1 - \tau_j^S)^{-1}$ , creating a gap between marginal costs and market prices. As  $\tau_j^S$  rises, output prices increase as well, making intermediate inputs more expensive for downstream firms. Facing higher variable input costs, these firms optimally substitute away from input-intensive technologies by reducing their returns to scale. Since firms do not internalize that the revenue from wedges is rebated to the household, this adjustment leads to an equilibrium with inefficiently low returns to scale.

Wedges also have important implications for the behavior of GDP. Propositions 6 and 7 show that without wedges, changes in returns to scale only have second-order effects on GDP. This is no longer the case when wedges are present.

**Proposition 8.** *In the presence of sales wedges, the impact of a parameter  $\chi \in \{\mu_j, \sigma_j, \kappa_j, \gamma_j\}$  on GDP is given by*

$$\frac{dy}{d\chi} = \underbrace{\frac{\partial y}{\partial \chi}}_{\text{Direct effect}} + \underbrace{\sum_{i=1}^N \frac{\partial y}{\partial \hat{\eta}_i} \frac{d\hat{\eta}_i}{d\chi}}_{\text{Structural change effect}}, \quad (33)$$

where  $\partial y / \partial \chi$  is given by Proposition 6,  $d\hat{\eta}_i / d\chi$  is given by Lemmas 6, 9 and 10, and  $\partial y / \partial \hat{\eta}_i \geq 0$ .

This proposition shows that the total effect of a shock is now the sum of a direct effect (the standard Hulten-like term) and a new scalability effect. Since sales wedges push  $\hat{\eta}$  to be inefficiently low, any shock that incentivizes firms to increase their returns to scale ( $d\hat{\eta}/d\chi > 0$ ) now generates an additional, first-order welfare gain.

Consider, for example, a positive productivity shock  $d\mu_j > 0$ . As in the efficient case, this shock directly raises GDP through the Hulten term. However, the lower input costs that it generates also induce firms to increase their  $\hat{\eta}$ . Because the economy started from an inefficiently low level of returns to scale, this structural adjustment is no longer a second-order refinement but a first-order improvement. This implies that  $dy/d\mu_j \geq \omega_j$ . Productivity shocks have a larger impact on GDP in this distorted economy because they not only improve productivity but also partially correct the pre-existing distorted scalability structure.

The nature of this new first-order effect depends on the type of distortion affecting the economy. While sales wedges lead to inefficiently low returns to scale, other distortions can have the opposite effect. For example, as we show in Online Supplement D.5, a corporate profit tax effectively raises the cost of entry, which incentivizes incumbent firms to choose inefficiently high returns to scale. In such an economy,

a productivity shock that further raises  $\hat{\eta}$  could actually be welfare-reducing at the margin.

## 5.5 Implications for growth

Endogenous scalability also propels long-run economic growth. To illustrate this mechanism transparently, we consider a continuous-time setup where the equilibrium of a single-sector version of our baseline model is repeated at each instant. Growth is driven by a constant rate of productivity improvement,  $d\mu/dt = g_\mu > 0$ .

We first characterize the growth rate of returns to scale in this environment.

**Corollary 1.** *The growth of effective returns to scale  $\hat{\eta}$  is given by*

$$\frac{d\hat{\eta}}{dt} = \Psi^{-1} \mathcal{K} g_\mu > 0. \quad (34)$$

Furthermore, as  $t \rightarrow \infty$ , effective returns to scale  $\hat{\eta}$  converges to 1.

The constant improvements in sectoral productivity  $\mu$  result in cheaper intermediate inputs, which push firms to increase their returns to scale, so that  $d\hat{\eta}/dt > 0$ .<sup>15</sup> The strength of that mechanism relies on how fast productivity improves ( $g_\mu$ ), but also on the importance of intermediate inputs in production, as captured by the term  $-\alpha/(1 - \hat{\eta}\alpha) = \mathcal{K}$ . When  $\alpha$  is large, intermediates are more important and  $\hat{\eta}$  grows faster, all else equal. In addition,  $d\hat{\eta}/dt$  depends on the elasticity of  $\hat{\eta}$  in response to a change in input prices. As in the comparative statics exercise of Section 4, this effect is captured by  $1/\Psi$ . An economy that is more flexible and faces stronger incentives will more rapidly reconfigure itself toward a more scalable structure.

The evolution of the economy's returns to scale has implications for the growth rate of GDP.

**Proposition 9.** *The growth rate of GDP is given by*

$$\frac{dy}{dt} = \frac{g_\mu}{1 - \alpha} \times \left( 1 - \frac{1}{\sqrt{\frac{1}{\gamma} \frac{1 - \alpha}{\alpha} \frac{g_\mu}{1 - \varphi} t + K}} \right) > 0, \quad (35)$$

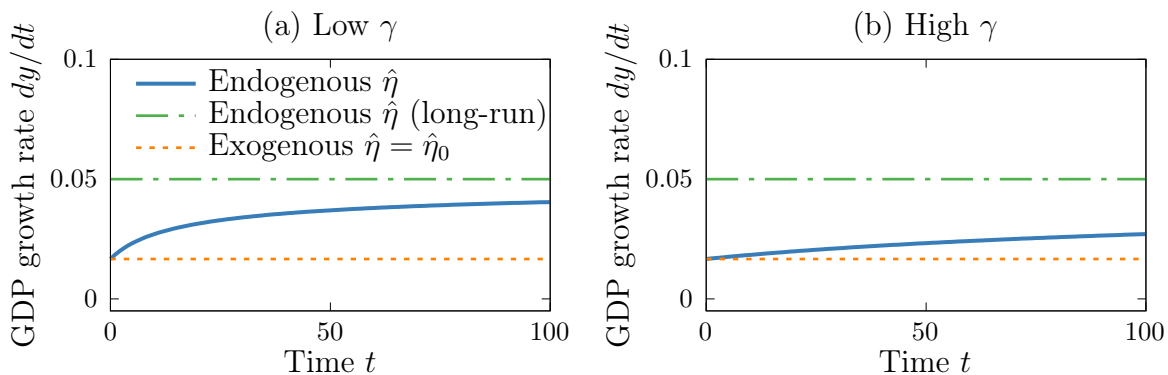
where  $K > 0$  is a time-invariant term given in the proof of the proposition.

This proposition provides a closed-form solution for the growth rate of GDP along its entire transition path. The growth rate  $dy/dt$  can be understood as the product

---

<sup>15</sup>Several empirical studies document rising returns to scale over time. De Loecker et al. (2020) estimate that firm-level returns to scale have increased over the last few decades in the United States. Chiavari and Goraya (2025) show that this finding holds even after accounting for intangible capital. Lashkari et al. (2024) provide evidence that the decline in IT prices led to an increase in returns to scale in France.

Figure 6: Endogenous returns to scale accelerate growth



of two terms: a *potential long-run growth rate*,  $g_\mu / (1 - \alpha)$ , and a convergence factor (the term in parentheses) that starts below one and asymptotically approaches it as  $t \rightarrow \infty$ . This means that the economy's growth rate is always increasing, accelerating towards its long-run potential, which it reaches asymptotically.

The speed of this convergence is governed by the terms inside the square root. The economy adapts its structure and accelerates more quickly when the forces driving reorganization are stronger. Specifically, the transition is faster when: 1) the pace of innovation ( $g_\mu$ ) is high, providing a strong and persistent incentive to adapt, and 2) the economy is structurally flexible, meaning the net cost of adjusting scale is low (low  $\gamma$  and high  $\varphi$ ).

Figure 6 provides an example of the dynamics of the growth rate of GDP. The left panel shows that from its initial condition  $\hat{\eta}(0) = \hat{\eta}_0$ , growth increases before converging to its long-run level  $g_\mu / (1 - \alpha)$ . The right-panel depicts the same economy but with a higher  $\gamma$ . In this case, increasing returns to scale is more costly, and the progression to the long-run growth rate is slower.<sup>16</sup>

Finally, we can compare the growth rate of GDP in the model with endogenous returns to scale to its counterpart in the fixed returns-to-scale economy.

**Corollary 2.** *For any  $t > 0$ , GDP grows faster in the economy with endogenous returns to scale. In the limit as  $t \rightarrow \infty$ , the long-run growth rates satisfy*

$$\lim_{t \rightarrow \infty} \frac{dy}{dt} = \frac{1}{1 - \alpha} g_\mu > \frac{1}{1 - \hat{\eta}_0 \alpha} g_\mu = \lim_{t \rightarrow \infty} \frac{d\tilde{y}}{dt},$$

where  $\tilde{y}$  is log GDP in the fixed returns-to-scale economy, and where  $\hat{\eta}_0$  is effective

<sup>16</sup>While growth is always accelerating ( $d^2y/dt^2 > 0$ ), the rate of acceleration is in general modest. One can show that the second derivative of log GDP scales with  $g_\mu^2$ . Thus, for realistic calibrations of annual productivity growth (e.g.,  $g_\mu \approx 1\%$ ), the acceleration is a gradual, slow-moving process that might be hard to notice over short horizons.

*returns to scale in the baseline economy at  $t = 0$ .*

This result captures the key growth implication of our model: an economy that can adapt its returns to scale grows faster. In the fixed-scale economy, the benefits of technological progress are constrained by a static production structure. In our model, in contrast, as sectoral productivity increases, the economy adopts more scalable technologies and becomes more interconnected. This increases Domar weights and magnifies the benefit of the higher productivity. This effect grows larger over time, with an increasing gap between the growth rates of the two economies. Our model therefore suggests that the macroeconomic consequences of scaling decisions, like Ford’s introduction of the moving assembly line, are not a one-off level effect, but a persistent force that reshapes the economy’s long-run growth trajectory.

The quantitative implications of endogenous returns to scale for growth can be substantial. For instance, using parameters calibrated to the Spanish economy (where  $\alpha \approx 0.67$  and  $\hat{\eta}_0 \approx 0.83$ , see the next sections), a 1% annual rate of underlying technological progress ( $g_\mu$ ) would translate into a 2.2% growth rate in the fixed- $\eta$  economy. With endogenous returns to scale, however, the long-run growth rate converges to 3.0%. Compounded over time, this difference of 0.8 percentage points can amount to a substantial welfare gain.

## 6 Empirical evidence

Our theoretical result describes how returns to scale respond to the economic environment. In this section, we use detailed firm-level data from Spain to provide empirical evidence for these predictions at the firm, sector, and aggregate levels. Consistent with the model, we show a robust positive correlation between productivity and returns to scale in the cross-section of firms. We further use panel data and within-firm variation to demonstrate that returns to scale respond to incentives: firms actively increase scalability as their productivity grows and reduce it when facing higher input costs. At the sector level, the theory predicts that industries with stronger endogenous scalability should exhibit fatter firm-size tails. We find strong support for this mechanism in the data. Finally, our theory implies that endogenous scalability is beneficial for GDP. To test this prediction, we extend our analysis to 24 countries, and show that the strength of the mechanism is indeed a predictor of long-run economic development.

## 6.1 Data

Our primary source of firm-level data is Moody’s Orbis Historical database, which covers the near-universe of Spanish firms between 1995 and 2019. This dataset provides detailed information on sales, labor costs, capital stocks, and material costs. After cleaning, our sample comprises 9,754,405 firm-year observations.<sup>17</sup> We use these data to estimate firm-level returns to scale using a production function approach. We briefly describe our estimation procedure below and provide additional details in Online Supplement A.1.

We complement our analysis with firm-level data from 23 additional countries. We use data from Orbis for 21 European countries with good coverage of the key variables needed for production function estimation. For developing economies, we rely on China’s National Bureau of Statistics (NBS) firm-level database and India’s Annual Survey of Industries (ASI). Both datasets are censuses of manufacturing firms above specific size thresholds. The details of data cleaning and variable construction are provided in Online Supplement A.7.

## 6.2 Estimating returns to scale and productivity

Our theory predicts that within a sector, firms of similar size should have similar returns to scale. We take advantage of this prediction to estimate returns to scale across the firm-size distribution. Specifically, for each sector  $i$  and year  $t$ , we group firms into 10 deciles based on their 7-year moving average of firm-level log sales (years  $t - 3$  to  $t + 3$ ). This construction smooths out short-run fluctuations and measurement error and thus yields a more reliable measure of a firm’s position in the size distribution.

In our baseline approach, we assume that all firms in sector  $i$ , year  $t$  and decile  $d_t$  share the same Cobb-Douglas production technology  $Q_{ilt} = \tilde{A}_{ilt} K_{ilt}^{\beta_{i,d_t(l),t}^K} L_{ilt}^{\beta_{i,d_t(l),t}^L} M_{ilt}^{\beta_{i,d_t(l),t}^M}$ , where  $K_{ilt}$ ,  $L_{ilt}$  and  $M_{ilt}$  denote capital, labor and materials, respectively. We then estimate the output elasticities for each cell  $(i, t, d_t)$  using the Blundell and Bond (2000) IV-GMM estimator on a 7-year rolling-window sample. The estimated returns to scale  $\eta_{ilt}$  for a firm  $l$  in sector  $i$  and year  $t$  is therefore given by the sum of these elasticities:  $\eta_{ilt} = \hat{\beta}_{i,d_t(l),t}^K + \hat{\beta}_{i,d_t(l),t}^L + \hat{\beta}_{i,d_t(l),t}^M$ . We then use the estimated returns to scale for the years 1997-2019 in our empirical analysis and in the calibration of Section 7.<sup>18</sup>

---

<sup>17</sup>We deflate all nominal variables using the Spanish GDP deflators and drop any firm-year observation whose average revenue product for any input (fixed assets, wage bills, or material costs) lies above the 99th percentile or below the 1st percentile of that year’s distribution.

<sup>18</sup>We do not estimate the elasticities for 1995 and 1996 as there are few observations in each

We use the Blundell–Bond estimator as our baseline because it can estimate the gross output production function by leveraging moment conditions that exploit input persistence, without requiring rigid assumptions on input timing. In Online Supplement A.4, we confirm the robustness of our results using a wide range of alternative strategies. These include: (i) alternative production function estimators, including Olley and Pakes (1996) and Levinsohn and Petrin (2003); (ii) controls for market power to alleviate the measurement errors in output using the Akerberg et al. (2015) estimator; and (iii) grouping firms by rolling sales percentiles or by contemporaneous rather than moving-average sales. Our main empirical patterns are robust across all of these alternative estimation strategies.

Finally, to compare productivity in the cross-section despite heterogeneous production technologies, we follow best practices in the development accounting literature (Caves et al., 1982a; Feenstra et al., 2015) and construct a comparable measure,  $\hat{z}_{ilt}$ , using a Törnqvist productivity index. Specifically, we define a hypothetical “average firm” for each sector-year with mean (log) sales, (log) inputs, and output elasticities. The measure  $\hat{z}_{ilt}$  then compares productivity between firm  $l$  and the average firm by looking at how much more output one produces relative to the other, adjusting for differences in technology and input use. We formally define this productivity index and provide more detail in Online Supplement A.3.<sup>19</sup>

## 6.3 Firm-level evidence

### 6.3.1 Cross-sectional patterns: returns to scale, productivity and size

A key prediction of our theory is that larger, more productive firms operate technologies with higher returns to scale. Figure 7 confirms this pattern in the data. Panel (a) plots returns to scale against firms’ sales percentiles within a sector-year. It shows that the largest firms (100th percentile) operate with returns to scale of around 0.91, compared to 0.79 for the smallest firms (1st percentile).<sup>20</sup> Furthermore, panel (b) shows that higher productivity is directly associated with higher returns to scale, which is a distinctive feature of our endogenous returns-to-scale mechanism. In addi-

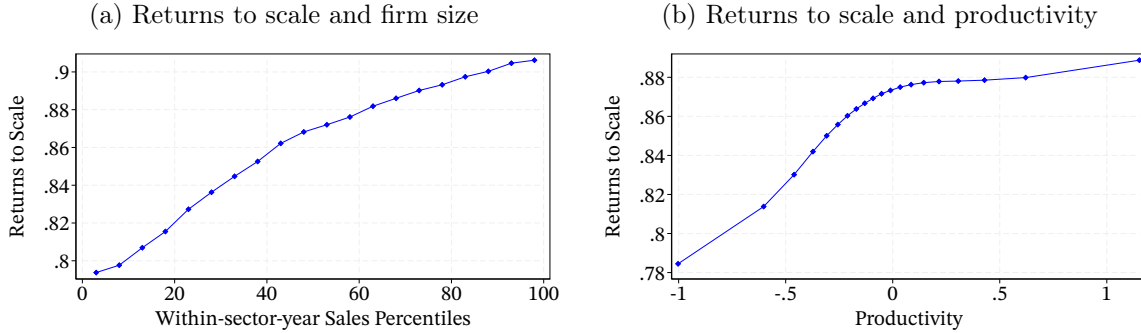
---

sector-year-decade cell in those years. Our 1997–2019 sample covers 9,424,952 firm-year observations.

<sup>19</sup>In our calibrated economy of Section 7, the within-sector correlation between this Törnqvist index and  $\varepsilon_{il}$  is above 0.99. The same number for  $\varepsilon_{il} + \alpha_i (\eta_{il}) + \log \zeta_{il}$  is about 0.92. The Törnqvist index therefore seems to be a good proxy for productivity for reduced-form estimates.

<sup>20</sup>Hubmer et al. (2025) document similar patterns in Canadian and US manufacturing firms. Gao and Kehrig (2025) and McAdam et al. (2024) report that industries with larger average firm size also have higher returns to scale in the United States and in European countries, respectively. Using production data to infer firm- and industry-level returns to scale has a long tradition since Hall (1990) and Basu and Fernald (1997).

Figure 7: Returns to scale, productivity, and firm size in the cross-section



Notes: Panel (a) presents a binned scatter plot of firm-level returns to scale against within-sector-year sales percentiles. Panel (b) presents a binned scatter plot of firm-level returns to scale against productivity, controlling for sector-year fixed effects; the unconditional mean of returns to scale is added back for interpretability. Sectors are defined as the 62 sectors in the Input–Output table of the Annual Spanish National Accounts. Both panels are constructed using a sample of Spanish firms from Orbis. See Supplement A.1 for details on variable construction and sample selection.

tion, the empirical relationship between returns to scale and productivity is concave. This supports our assumption of a rising marginal cost of scalability (concave  $a_i$ ) and qualitatively matches the theoretical prediction illustrated in Figure 2.

### 6.3.2 Panel evidence on endogenous returns to scale

The cross-sectional evidence presented in the last section shows that, on *average*, larger and more productive firms are more scalable. But the theory also predicts that the returns to scale of *individual firms* should also respond to changes in the economic environment. Indeed, Lemma 3 establishes that beneficial shocks, whether higher productivity or lower input prices, induce firms to adopt more scalable technologies. Since these same shocks also drive firm expansion, the model implies that returns to scale and sales should co-move positively within a firm over time.

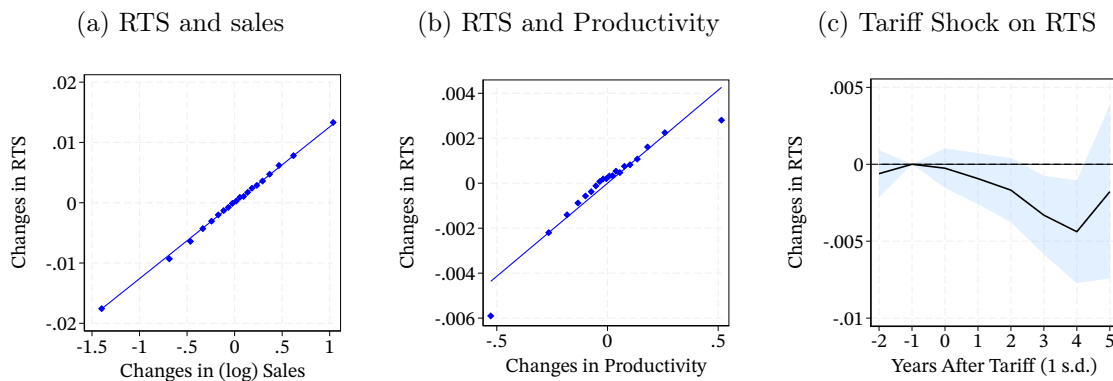
Figure 8 tests these predictions by plotting the *within-firm* variation in returns to scale against changes in sales and productivity, controlling for firm and sector-year fixed effects.<sup>21</sup> Both panels show a strong positive relationship: when a firm experiences an increase in sales or productivity, its returns to scale tend to rise. Quantitatively, a 100% increase in sales is associated with a 0.013 increase in returns to scale, while an analogous increase in productivity raises returns to scale by 0.008. Taken together, these within-firm results provide panel evidence for our mechanism: rather than operating fixed production technologies, firms seem to adopt higher returns to scale as they grow.

We next examine the response of returns to scale to changes in input costs. In

<sup>21</sup>To track within-firm productivity over time, we use a within-firm Törnqvist productivity index that nets out share-weighted input growth from output growth. See Supplement A.3 for details.

the model, an increase in the cost of the variable input bundle, for example, due to higher tariffs on imported intermediates, reduces returns to scale.<sup>22</sup> To test this prediction, we exploit variation in import tariffs, which differentially affect input costs across sectors and over time. We construct a sector-year measure of exposure to tariff-induced changes in input prices,  $\log T_{it}$ , combining data from the OECD multi-country input–output tables and the Global Tariff Project by Teti (2024) (see Online Supplement A.6 for details).

Figure 8: Within-Firm variation in RTS: Productivity, Size, and Tariff Shocks



Notes: Panel (a) presents a binned scatter plot of firm-level returns to scale against log sales. Panel (b) presents a binned scatter plot of firm-level returns to scale against productivity. Firm fixed effects and sector-year fixed effects are controlled for in both panels. Panel (c) plots the estimated dynamic response coefficients  $\hat{\beta}_h$  for horizons  $h = -2, \dots, 5$  to a standardized tariff shock of (36). 90% confidence intervals are constructed using standard errors two-way clustered by firm and industry-year. Industries are defined according to the OECD Input–Output table. See the main text for details on variable construction and sample selection.

We then estimate the dynamic impact of these shocks on returns to scale using panel local projections for horizon years  $h = -2, \dots, 5$ :

$$\eta_{il,t+h} - \eta_{il,t-1} = \beta_h \log T_{it} + \gamma_{lh} + \gamma_{th} + \varepsilon_{ilth}, \quad (36)$$

including firm ( $\gamma_{lh}$ ) and year ( $\gamma_{th}$ ) fixed effects. Panel (c) of Figure 8 plots the estimated dynamic response  $\hat{\beta}_h$  to a standardized tariff shock. Consistent with the theory, firms that are more exposed to tariff-induced cost increases experience larger declines in returns to scale after the shock. A one-standard-deviation shock is associated with a decline in returns to scale of up to 0.004. This adjustment happens progressively, suggesting that changing returns to scale might take time.

Taken together, this evidence is consistent with the key premise of our model that firms adjust their returns to scale in response to both productivity and input-cost changes. This supports our interpretation of returns to scale as an endogenous choice margin rather than a fixed feature of the production function.

<sup>22</sup>We formally analyze the model with a tax on intermediate inputs in Supplement D.5.

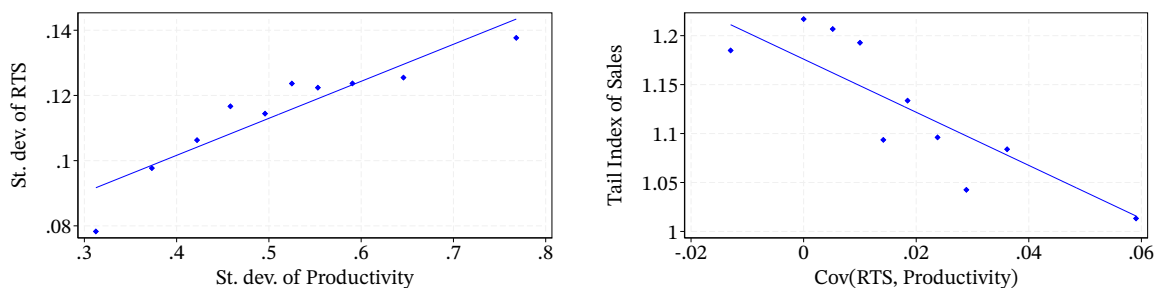
## 6.4 Sectoral evidence

Our model also has implications for cross-industry variation. For instance, because high-productivity firms adopt more scalable technologies to leverage their productivity advantage, sectors with greater productivity dispersion should also exhibit greater dispersion in returns to scale (Lemma 5). The first panel of Figure 9 confirms this prediction in the Spanish data. It shows a strong positive relationship between these moments: sectors in the top decile of productivity dispersion have a standard deviation of returns to scale that is 0.06 higher than those in the bottom decile. This pattern supports the model’s “double blessing” mechanism, which gives rise to superstar firms through the adoption of high-scalability technologies.

This mechanism also implies that the shape of the firm-size distribution should vary systematically across sectors. Specifically, Proposition 1 shows that sectors where adopting high scalability is easier (high  $\varphi_i$ ) should have a thicker right tail of the sales distribution. Since  $\varphi_i$  is not directly observable, we construct a proxy for the strength of the mechanism: the within-sector covariance between returns to scale and the productivity index,  $\text{Cov}(\eta_{it}, \varepsilon_{it} + a(\eta_{it}))$ . Theoretically, this covariance should be strictly positive when returns to scale are endogenous ( $\varphi_i > 0$ ), but zero if returns to scale are fixed.<sup>23</sup> Using this measure, panel (b) of Figure 9 confirms the model prediction: sectors where firms can more easily adjust their returns to scale exhibit significantly thicker tails in their sales distribution.<sup>24</sup>

Figure 9: Dispersion in returns to scale, productivity and market concentration

(a) Dispersion in returns to scale and productivity (b) Tail indices of sales and endogenous scalability



Notes: Panel (a) presents a binned scatter plot of the standard deviation of firm-level returns to scale against the standard deviation of firm-level productivity, for all sector–year observations. Panel (b) presents a binned scatter plot of the tail index of firm sales against the covariance between firm-level returns to scale and productivity, for all sector–year observations with at least 50 firms. Year fixed effects are controlled for in both panels and unconditional means are added back for interpretability. See Supplement A.1 for details on variable construction and sample selection.

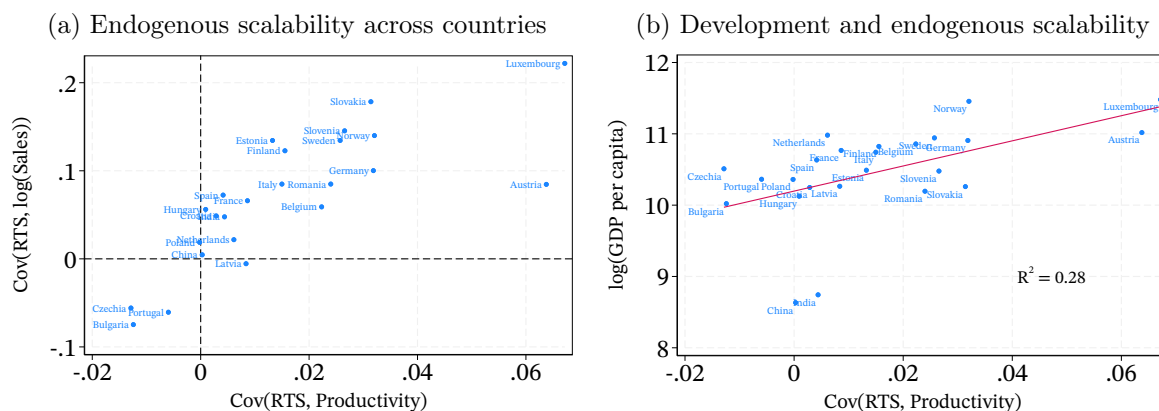
<sup>23</sup>One can show that in the model this quantity is increasing in  $\varphi_i$ .

<sup>24</sup>We use the tail index estimator of Gabaix and Ibragimov (2011). See Supplement A.5 for details.

## 6.5 Cross-country development and endogenous scalability

We conclude this section by providing cross-country evidence for the importance of endogenous returns to scale in our sample of 24 countries. To ensure comparability, our main cross-country analysis focuses on manufacturing firms only. For each country, we select the seven-year window with the largest number of firm-year observations between 2001 and 2018.<sup>25</sup> Within this window, we estimate size-decile-specific production functions using the same Blundell–Bond procedure described above to construct firm-level measures of returns to scale and productivity.

Figure 10: Cross-country evidence on productivity, sales, and returns to scale



Notes: Panel (a) plots, for each country, the covariance between firm-level returns to scale and log sales against the covariance between firm-level returns to scale and productivity; dashed lines denote zero covariances. Panel (b) plots the seven-year-average of log GDP per capita against the covariance between firm-level returns to scale and productivity; the solid line reports the fitted linear relationship and the figure reports the associated  $R^2$ . Each marker corresponds to a country. See Online Supplement A.7 for details on variable construction and sample selection.

Figure 10 summarizes our results. Panel (a) reports the covariance between returns to scale and log sales, as well as between returns to scale and productivity. Mirroring the Spanish evidence in Section 6.3.1, we find that in most countries, larger and more productive firms systematically operate with higher returns to scale. This suggests that endogenous scalability is at work across a broad set of countries.

Panel (b) relates the strength of the endogenous returns-to-scale mechanism to economic development. As before, we measure the intensity of endogenous scalability using the firm-level covariance between productivity and returns to scale. Plotting this measure against log GDP per capita reveals a clear relationship: countries with stronger endogenous scalability are systematically richer.<sup>26</sup> This is consistent with

<sup>25</sup>For several countries, the quality of the data is uneven over the sample period. We therefore restrict our analysis to the time window with the most observations to improve the quality of our estimates. We then apply identical data cleaning and estimation procedures to all country samples.

<sup>26</sup>India and China are two outliers but since those countries are at an earlier stage of development, their lower GDP per capita level might be explained by their distance to the technology frontier.

the model prediction that endogenous scalability increases the level of GDP. Quantitatively, variation in this measure accounts for approximately one quarter of the cross-country variation in GDP per capita. This suggests that the capacity of productive firms to adopt more scalable technologies might be an important driver of long-run economic development.

## 7 Calibration to the Spanish economy

To evaluate the quantitative importance of endogenous returns to scale, we calibrate the model to the Spanish economy. We rely on the detailed firm-level data introduced in the previous section to discipline the model parameters. Our calibration strategy is summarized here, and more details are in Online Supplement B.

### 7.1 Calibration strategy

We calibrate the household preference vector  $\beta$  and the input-share matrix  $\alpha$  directly to match the sectoral final consumption shares and intermediate expenditure shares observed in the 2010 Spanish Input-Output tables.<sup>27</sup> The firm-level returns to scale  $\eta_{il}$  are taken from our estimates of Section 6. We then compute each sector’s effective returns to scale  $\hat{\eta}_i$  as the sales-weighted average of  $\eta_{il}$ . We find substantial heterogeneity in  $\hat{\eta}_i$  across sectors, ranging from 0.54 to 0.98, with a mean and median of 0.83 and 0.82, respectively. Figure 15 in Online Supplement B.4 shows  $\hat{\eta}_i$  for all sectors.

The productivity dispersion  $\sigma_i$  and the cost of scalability  $\gamma_i$  are jointly identified by targeting within-sector moments that are informative about scalability choices. As we show in Online Supplement B, the model implies a mapping from the pair  $(\sigma_i, \gamma_i)$  to the cross-sectional dispersion of firm-level returns to scale  $\eta_{il}$  and profits  $\Pi_{il}$ . Accordingly, we choose  $\sigma_i$  and  $\gamma_i$  for each sector to match the interquartile ranges of these two variables in the data.<sup>28</sup> Online Supplement B.4 reports the calibrated values (Figure 14) and shows that the model matches the targeted moments well (Figure 13). By construction, the calibrated model also matches the empirical effective returns to scale  $\hat{\eta}$  perfectly.<sup>29</sup>

---

<sup>27</sup>We calibrate the model to the 2010 Spanish economy because it is close to the midpoint of our sample period and provides detailed input–output tables that can be cleanly matched to the Orbis firm-level data at the NACE 2-digit level. Over time, changes in the national-accounts base complicate consistent harmonization of detailed input–output tables across years.

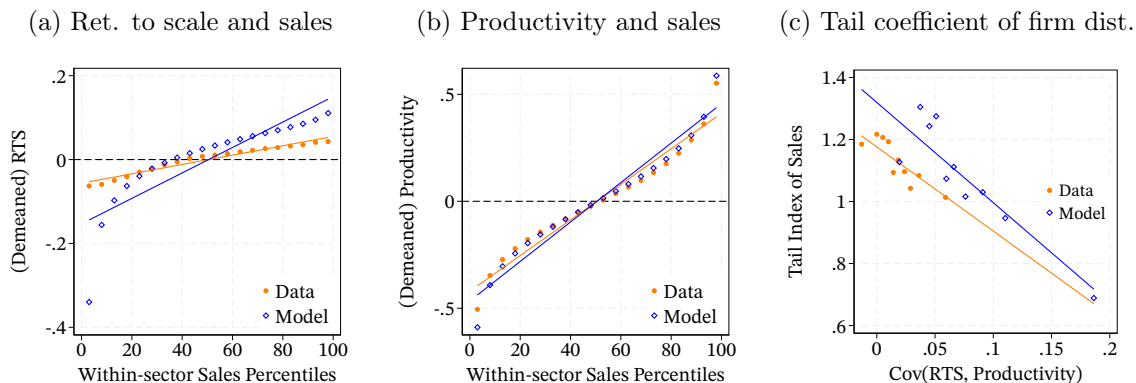
<sup>28</sup>To be consistent with the model, we compute the profits of firm  $l$  in sector  $i$  as  $(1 - \eta_{il}) P_i Q_{il}$ . We target interquartile ranges instead of variances to reduce the impact of outliers.

<sup>29</sup>As we study *changes* to the environment, the results presented below are independent of the average productivity  $\mu$  and the entry cost  $\kappa$ , and so there is no need to specify their values. Given all the other parameters, we can always find a combination of  $\mu$  and  $\kappa$  to match  $\hat{\eta}$  perfectly.

To validate the calibration, we test the model’s ability to reproduce untargeted moments of the firm distribution. Recall that our mechanism implies that larger firms operate technologies with higher returns to scale. Panel (a) of Figure 11 compares the relationship between these quantities in the model and the data. Although this moment was not targeted, the model replicates the empirical pattern reasonably well. Two differences, however, are worth noting. First, the fit is less precise for very small firms. But since these producers account for a negligible share of aggregate output, they have little influence on the counterfactual exercises that follow. Second, the relationship between sales and returns to scale is steeper in the model than in the data. This is likely because firms in the model adjust their returns to scale instantly in response to productivity shocks, whereas real-world adjustments are subject to frictions and delays.

Panel (b) of Figure 11 shows that the model closely matches the empirical link between productivity (measured using the Törnqvist index) and sales. This is reassuring, as the link between these variables in the model is driven by the endogenous returns-to-scale parameter  $\eta_{il}$ .<sup>30</sup> Finally, panel (c) shows that the tail of the firm-size distribution is thicker in sectors where the endogenous scalability mechanism is stronger. Once again, the model matches the data reasonably well.

Figure 11: RTS, sales, productivity, and tail coefficients in the model and in the data



Notes: *Data* correspond to Spanish firms in Orbis; *Model* corresponds to simulated firm-level outcomes from the calibrated model. Simulated firm observations are reweighted at the sector level so that the sector composition matches the data. Panels (a) and (b) report binscatter plots of returns to scale and productivity against sales percentiles. The y-axis variables are residualized by sector–year fixed effects in the *Data* series and sector fixed effects in the *Model* series. Panel (c) plots the tail index of sales against the within-sector-year covariance of returns to scale and productivity: the *Data* series uses sector–year statistics (computed only for sector–years with at least 50 firms) and includes year fixed effects, whereas the *Model* series uses sector-level statistics computed from simulated data and includes no fixed effects.

<sup>30</sup>Productivity and returns to scale are also positively correlated in the calibrated model, as they are in the data.

## 7.2 Contribution of endogenous returns to scale to GDP

Using our calibrated model, we first evaluate the importance of endogenous returns to scale for the level of GDP. To do so, we compare our calibrated baseline economy to the “fixed returns-to-scale” counterfactual where each firm’s returns to scale is exogenously fixed at its sector’s average,  $\eta_{il} = \hat{\eta}_i$  (Definition 2). As shown in Proposition 5, the difference in log GDPs between those two economies is given by (26). Using our calibrated parameters, we find that this gain is 11.7%, so that allowing high-productivity firms to choose more scalable technologies has a substantial effect on GDP. Figure 16 in Online Supplement B.4 decomposes this aggregate gain, showing the contribution of each sector.

Next, we evaluate the importance of our mechanism for long-run macroeconomic growth. We conduct an experiment in which we increase the mean productivity  $\mu_i$  of all sectors by one percentage point every year. We then compare the response of GDP in our full model to two counterfactual benchmarks. The purpose of this exercise is to decompose the full effect of endogenous scalability into two components: the gain from the initial *static dispersion* in returns to scale, and the additional gain from allowing these returns to scale to *dynamically adjust* over time. Our first benchmark is the fixed returns-to-scale economy where  $\eta_{il}(t) = \hat{\eta}_i(0)$  for all firms, shutting down both channels. Our second benchmark is a *dispersed return-to-scale* economy, where  $\eta_{il}$  is fixed at each firm’s initial level  $\eta_{il}(t) = \eta_{il}(0)$ , thus featuring the initial dispersion but not the dynamic adjustment.<sup>31</sup>

Panel (a) of Figure 12 shows the evolution of GDP over time in each economy, relative to the fixed returns-to-scale benchmark. The initial gap at  $t = 0$  reflects the level effect from (26). The solid blue line captures the full effect of endogenous returns to scale on GDP. We see that after one hundred years, GDP has grown an extra 5.5% in the full model. Panel (c) illustrates the underlying mechanisms. As  $\mu$  increases, the price of intermediate inputs falls, which incentivizes firms to adopt more scalable technologies. This increase in returns to scale leads to higher Domar weights, which makes the increase in productivity more impactful.<sup>32</sup>

Figure 12 allows us to decompose the 5.5% gain in GDP. It shows that the dispersed returns-to-scale economy grew by an extra 2.8% over a century compared to the fixed benchmark. This implies that the static reallocation channel and the dy-

---

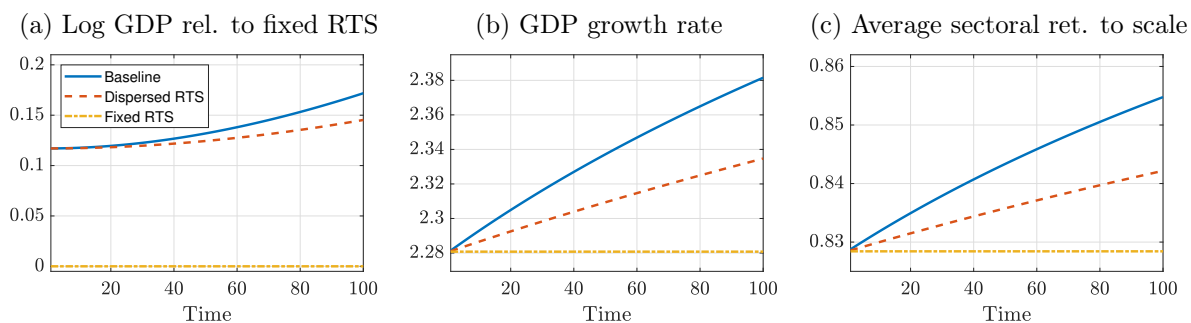
<sup>31</sup>We formally analyze this economy in Online Supplement D.7.

<sup>32</sup>Over one hundred years, effective returns to scale in the calibrated model increase by about 0.025, which is within the range found by empirical studies. For instance, De Loecker et al. (2020) find that firm-level returns to scale increased between 0.01 and 0.1 over periods of 40 to 60 years, depending on the estimation method.

dynamic adjustment channel each contribute roughly half of the total effect. The static channel is powerful because, even when individual  $\eta_{il}$  cannot change, the economy-wide productivity boom disproportionately benefits the firms that can scale more easily—that is, those that started with higher  $\eta_{il}$  (Lemma 4). As these firms expand their output disproportionately, their sales shares increase. This compositional shift has two beneficial effects. First, it endogenously raises effective returns to scale  $\hat{\eta}$  and the Domar weights  $\omega$ , which amplifies growth even though no firm changes its technology. Second, since these high- $\eta$  firms are also the most productive, this reallocation of market share directly raises aggregate productivity.

Figure 12 also shows that the evolution of log GDP over time in the full model is convex, with growth constantly accelerating (panel b). This indicates that the effects of endogenous returns to scale become increasingly important over time, with the gap between the three curves in panel (a) increasing indefinitely. In the long-run, as  $t \rightarrow \infty$ , the baseline economy converges to a growth rate of 3.1% under our parametrization. In contrast, in the fixed returns-to-scale economy, where the endogenous returns-to-scale mechanism is shut down, long-run growth is only 2.3%. The large gap between these two numbers suggests that endogenous scalability might play a significant role in shaping long-run economic outcomes.<sup>33</sup>

Figure 12: Endogenous returns to scale and productivity: Implications for GDP



Notes: Panel (a) shows  $y(t) - \bar{y}(t)$  (“Baseline”) and  $y^d(t) - \bar{y}(t)$  (“Dispersed RTS”) as the productivity  $\mu$  of all sectors grows at 1% per year, where  $y$ ,  $\bar{y}$ , and  $y^d$  are log GDPs in the economies with fully flexible returns to scale, fixed returns to scale, and dispersed returns to scale that do not respond to changes in  $\mu$ , respectively. Panel (b) shows the growth rates of log GDP and panel (c) shows average returns to scale  $\sum_{i=1}^N \hat{\eta}_i(t) / N$ .

<sup>33</sup>As  $t \rightarrow \infty$ , the *growth rates* (but not the levels) of GDP in the baseline model and the dispersed returns-to-scale economy converge. In that limit, output in both economies is dominated by a vanishingly small fraction of extremely productive firms operating with near-constant returns to scale. In that case, the gains from the *dynamic channel* of endogenous returns to scale are exhausted, and both economies grow at the same pace of 3.1%.

### 7.3 Distorted economy

One implication of our model is that the ability of high-productivity firms to adopt more scalable technologies is important for the level and growth rate of GDP. Yet, several studies have documented that larger firms often suffer from higher wedges (e.g., Restuccia and Rogerson, 2008). With endogenous returns to scale, those wedges would disproportionately distort the scalability decisions of those superstar firms, with potentially large consequences for welfare. In this section, we provide a simple exercise to quantify this new adverse effect of wedges.<sup>34</sup>

Following Hsieh and Klenow (2009), we compute the sales wedge  $\widehat{\tau}_{il}^S$  for firm  $l$  in sector  $i$  as the ratio of the marginal revenue product of labor to the wage:

$$\frac{1}{1 - \widehat{\tau}_{il}^S} = \frac{\eta_{il} \left(1 - \sum_{j=1}^N \alpha_{ij}\right) P_i Y_{il}}{W L_{il}}.$$

Consistent with the literature, we find these wedges to be large. The average sales-weighted wedge in a sector is 0.40. Furthermore, wedges are positively correlated with firm size: the average within-sector correlation between  $\widehat{\tau}_{il}^S$  and firm sales is 0.27.

To study the impact of these wedges in our calibrated economy, we assume that each firm is subject to a sales wedge  $\tau_{il}^S$  given by

$$\log(1 - \tau_{il}^S) = \log(1 - \tau_i^S) - b_i(\varepsilon_{il} - \mu_i), \quad (37)$$

where  $\log(1 - \tau_i^S)$  is an average sector-wide term and  $-b_i(\varepsilon_{il} - \mu_i)$  captures size-dependent distortions, with  $b_i > 0$  implying that more productive firms faces stronger distortions. As in the case where the wedge corresponds to a tax or a markup, we assume that net revenues collected from  $\tau_{il}^S$  are rebated to the household lump-sum. We calibrate the intercept  $\tau_i^S$  and the slope  $b_i$  to match the average sectoral wedge and the covariance between firm profits and the estimated wedges in the data. We estimate  $b_i > 0$  in nearly all sectors.<sup>35</sup> After fixing  $\tau_i^S$  and  $b_i$ , we recalibrate the remaining model parameters so that the economy continues to match our empirical targets. Further details on this procedure are provided in Online Supplement B.3.

To evaluate the importance of these distortions, we conduct an experiment in which we remove all wedges.<sup>36</sup> The first two columns of Table 1 show that, in the

<sup>34</sup>Hubmer et al. (2025) also look at the impact of wedges in an economy with heterogeneous, but exogenous, returns to scale.

<sup>35</sup>Our estimates of size-dependent distortions ( $b_i$ ) are consistent with the cross-country evidence in Ayerst et al. (2024) and with recent estimates for India in Hsieh et al. (2025).

<sup>36</sup>For some sectors, removing wedges would push  $\varphi_i$  above one. For these sectors, we set  $\varphi_i = 0.99$ . We conduct sensitivity analysis to the value of this threshold in Online Supplement B.5.

baseline model, this leads to a large increase in both returns to scale and GDP. To understand the mechanisms behind this result, the table also reports the outcome of the same experiment in the dispersed and the fixed returns-to-scale economies. In the dispersed return-to-scale economy, where firms choose their returns to scale in the presence of wedges but cannot adjust them once wedges are removed, the impact of eliminating wedges is substantially smaller. In this case, firms at the top of the distribution cannot increase their scalability to fully benefit from the removal of distortions, which limits welfare gains.<sup>37</sup> In the fixed return-to-scale economy, where the endogenous scalability mechanism is shut down entirely, the effect of wedges is even more muted. Because all firms operate with the same returns to scale, the firm-size distribution is more compressed and GDP is produced more evenly across firms. Distorting high-productivity producers is thus less damaging in that economy.

Table 1: Returns to scale and GDP when wedges are removed

	Size-dependent wedges		Flat wedges	
	$\Delta$ RTS	$\Delta$ log GDP	$\Delta$ RTS	$\Delta$ log GDP
Baseline economy	0.067	167%	0.020	62%
Dispersed ret. to scale	0.046	138%	0.010	60%
Fixed ret. to scale	0	70%	0	58%

Notes: Increases in average effective returns to scale,  $\Delta \left[ \sum_{i=1}^N \hat{\eta}_i / N \right]$ , and in log GDP,  $\Delta y$ , due to removal of sales wedges in the baseline economy, and in the economies with dispersed and fixed returns to scale. The “Size-dependent wedges” columns reports the results when sales taxes are correlated with firm productivity. The “Flat wedges” columns reports the results when sales wedges are identical for all firms within a sector.

To further show that wedges that disproportionately affect the top firms are particularly harmful when those firms endogenously choose the scalability of their operation, the last two columns of Table 1 repeat this analysis for an economy with flat wedges ( $b_i = 0$ ). In this setting, productive firms face no additional penalty relative to smaller firms, so removing distortions yields much smaller GDP gains. Moreover, the results are nearly identical across the three model specifications. This confirms that the interaction between endogenous scalability and size-dependent distortions is the primary driver of our results.

In summary, our quantitative analysis suggests that endogenous returns to scale might play a substantial role in shaping both the level and the growth rate of GDP. When highly productive firms are able to adopt more scalable technologies as the technological frontier advances, they can expand disproportionately, with substantial gains for welfare. Taxes or distortions that fall on those firms, however, can disrupt

<sup>37</sup>Effective sectoral returns to scale still increase in this case since high- $\eta_{il}$  firms become larger.

this process. As our results show, size-dependent distortions discourage the adoption of high-scale technologies, stifle the growth of superstar firms, and generate efficiency losses that can exceed those in standard models. Policy interventions that disproportionately burden high-productivity firms may therefore be particularly harmful.

## 8 Conclusion

We develop a theory in which returns to scale are endogenous equilibrium objects driven by incentives. At the micro level, this mechanism gives rise to superstar firms and fat-tailed firm-size distributions. At the macro level, it endows the economy with greater resilience by dampening the impact of adverse shocks while amplifying favorable ones. It also provides an engine for long-run growth, as the economy’s organizational structure co-evolves with its technological frontier. Input-output connections between sectors play a crucial role in these mechanisms.

Several extensions would be worth pursuing. First, introducing capital would change the incentives to increase returns to scale. Since capital can be accumulated, its presence might affect the growth properties of the model. Second, the superstar firms that emerge in our model might, in reality, gain market power, creating a feedback effect that would further increase their incentives to scale. Third, modeling individual margins through which firms achieve scale would allow for a tighter connection to the data and a deeper understanding of the importance of scalability for macroeconomic outcomes.

## References

- Acemoglu, Daron and Pablo D Azar (2020). “Endogenous Production Networks”. *Econometrica* 88.1, pp. 33–82.
- Acemoglu, Daron, Vasco M Carvalho, Asuman Ozdaglar, and Alireza Tahbaz-Salehi (2012). “The Network Origins of Aggregate Fluctuations”. *Econometrica* 80.5, pp. 1977–2016.
- Akerberg, Daniel A, Kevin Caves, and Garth Frazer (2015). “Identification properties of recent production function estimators”. *Econometrica* 83.6, pp. 2411–2451.
- Argente, David, Sara Moreira, Ezra Oberfield, and Venky Venkateswaran (2025). *Scalable Expertise: How Standardization Drives Scale and Scope*. Tech. rep. National Bureau of Economic Research.
- Axtell, Robert L (2001). “Zipf distribution of US firm sizes”. *Science* 293.5536, pp. 1818–1820.
- Ayerst, Stephen, Duc M Nguyen, and Diego Restuccia (2024). *The micro and macro productivity of nations*. Tech. rep. National Bureau of Economic Research.

- Baqae, David Rezza and Emmanuel Farhi (2019a). “Productivity and Misallocation in General Equilibrium”. *Quarterly Journal of Economics* 135.1, pp. 105–163.
- (2019b). “The Macroeconomic Impact of Microeconomic Shocks: Beyond Hulten’s Theorem”. *Econometrica* 87.4, pp. 1155–1203.
- Basu, Susanto and John G Fernald (1997). “Returns to scale in US production: Estimates and implications”. *Journal of political economy* 105.2, pp. 249–283.
- Bigio, Saki and Jennifer La’O (2020). “Distortions in Production Networks”. *Quarterly Journal of Economics* 135.4, pp. 2187–2253.
- Blundell, Richard and Stephen Bond (2000). “GMM estimation with persistent panel data: an application to production functions”. *Econometric reviews* 19.3, pp. 321–340.
- Caves, Douglas W, Laurits R Christensen, and W Erwin Diewert (1982a). “Multilateral comparisons of output, input, and productivity using superlative index numbers”. *Economic Journal* 92.365, pp. 73–86.
- Chandler, Alfred Dupont (1977). *The Visible Hand: The Managerial Revolution in American Business*. Harvard University Press.
- (1990). *Scale and scope: The dynamics of industrial capitalism*. Harvard University Press.
- Chiavari, Andrea and Sampreet Singh Goraya (2025). “The rise of intangible capital and the macroeconomic implications”.
- De Loecker, Jan, Jan Eeckhout, and Gabriel Unger (2020). “The rise of market power and the macroeconomic implications”. *Quarterly Journal of Economics* 135.2, pp. 561–644.
- Engbom, Niklas, Hannes Malmberg, Tommaso Porzio, Federico Rossi, and Todd Schoellman (2025). *Economic Development According to Chandler*. Tech. rep.
- Feenstra, Robert C, Robert Inklaar, and Marcel P Timmer (2015). “The next generation of the Penn World Table”. *American economic review* 105.10, pp. 3150–3182.
- Gabaix, Xavier and Rustam Ibragimov (2011). “Rank-  $1/2$ : a simple way to improve the OLS estimation of tail exponents”. *Journal of Business & Economic Statistics* 29.1, pp. 24–39.
- Gao, Wei and Matthias Kehrig (2025). *Returns to scale, productivity and competition: Empirical evidence from US manufacturing and construction establishments*. Working Paper.
- Gottlieb, Charles, Markus Poschke, and Michael Tueting (2025). “Skill Supply, Firm Size, and Economic Development”. *Paper for World Development Report*.
- Hall, Robert E (1990). “Invariance properties of Solow’s productivity residual”. *Growth, Productivity, Unemployment: Essays to Celebrate Bob Solow’s Birthday*.
- Hopenhayn, Hugo A (1992). “Entry, exit, and firm dynamics in long run equilibrium”. *Econometrica: Journal of the Econometric Society*, pp. 1127–1150.
- Hsieh, Chang-Tai and Peter J Klenow (2009). “Misallocation and manufacturing TFP in China and India”. *Quarterly Journal of Economics* 124.4, pp. 1403–1448.
- Hsieh, Chang-Tai, Peter J. Klenow, Emmanuella Kyei Manu, and Emma Rockall (2025). *Misallocation versus Inequality*. working paper.

- Hubmer, Joachim, Mons Chan, Serdar Ozkan, Sergio Salgado, and Guangbin Hong (2025). *Scalable versus Productive Technologies*. Tech. rep.
- Hulten, C. R. (1978). “Growth Accounting with Intermediate Inputs”. *Review of Economic Studies* 45.3, pp. 511–518.
- Jones, Charles I (2011). “Intermediate Goods and Weak Links in the Theory of Economic Development”. *AEJ: Macroeconomics* 3.2, pp. 1–28.
- Kopytov, Alexandr, Bineet Mishra, Kristoffer Nimark, and Mathieu Taschereau-Dumouchel (2024a). “Endogenous production networks under supply chain uncertainty”. *Econometrica* 92.5, pp. 1621–1659.
- Kopytov, Alexandr, Mathieu Taschereau-Dumouchel, and Zebang Xu (2024b). “The Origin of Risk”. *Available at SSRN*.
- Kuznets, Simon (1973). “Modern economic growth: findings and reflections”. *The American economic review* 63.3, pp. 247–258.
- Lashkari, Danial, Arthur Bauer, and Jocelyn Boussard (2024). “Information technology and returns to scale”. *American Economic Review* 114.6, pp. 1769–1815.
- Levinsohn, James and Amil Petrin (2003). “Estimating production functions using inputs to control for unobservables”. *Review of Economic Studies* 70.2, pp. 317–341.
- Liu, Ernest (2019). “Industrial Policies in Production Networks”. *Quarterly Journal of Economics* 134.4, pp. 1883–1948.
- Long, John B. and Charles I. Plosser (1983). “Real Business Cycles”. *Journal of Political Economy* 91.1, pp. 39–69.
- Lucas, Robert E (1978). “On the size distribution of business firms”. *The Bell Journal of Economics*, pp. 508–523.
- McAdam, Peter, Philipp Meinen, Chris Papageorgiou, and Patrick Schulte (2024). *Returns to scale: New evidence from administrative firm-level data*. Tech. rep. 24/2024.
- McKenzie, Lionel W (1959). “On the existence of general equilibrium for a competitive market”. *Econometrica: journal of the Econometric Society*, pp. 54–71.
- Oberfield, Ezra (2018). “A Theory of Input-Output Architecture”. *Econometrica* 86.2, pp. 559–589.
- Olley, G Steven and Ariel Pakes (1996). “The Dynamics of Productivity in the Telecommunications Equipment Industry”. *Econometrica* 64.6, pp. 1263–1297.
- Restuccia, Diego and Richard Rogerson (2008). “Policy distortions and aggregate productivity with heterogeneous establishments”. *Review of Economic Dynamics* 11.4, pp. 707–720.
- Smirnyagin, Vladimir (2023). “Returns to scale, firm entry, and the business cycle”. *Journal of Monetary Economics* 134, pp. 118–134.
- Teti, Feodora A. (2024). “Missing Tariffs”.