Scale Economies and Aggregate Productivity

Joel Kariel^{*} Anthony Savagar[†]

December 14, 2023

Abstract

We present a theory of rising scale economies and stagnating productivity in a model of heterogeneous firms with imperfectly competitive product markets and firm selection. Our theory shows that scale economies arising from fixed costs versus returns to scale differ in their effect on aggregate productivity. Using UK data, we estimate a long-run increase in fixed costs and returns to scale. Our model implies that this should have significantly increased aggregate productivity, both through stronger selection of high-technical-efficiency firms and better allocation of resources across firms. However, increasing markups can offset the productivity gain. Higher markups cushion low-productivity firms' revenues, allowing them to survive, and constrain firm output, which limits exploitation of scale economies.

JEL: E32, E23, D21, D43, L13.

Keywords: Returns to Scale, Scale Economies, Productivity, Market Structures, Firm Dynamics, Fixed Costs, Marginal Costs.

Disclaimer: This work was produced using statistical data from ONS. The use of the ONS statistical data in this work does not imply the endorsement of the ONS in relation to the interpretation or analysis of the statistical data. This work uses research datasets which may not exactly reproduce National Statistics aggregates. We thank seminar participants at Cardiff, CompNet, Nottingham, York, Durham Macro workshop, EUI, Bank of Italy, MWM Clemson, University of Washington, RES 2023, EMF Bern 2023, Birmingham, St Louis Fed, Lancaster, King's, Bath, IFN Stockholm, EEA-ESEM 2022, IAAE 2022, CEF 2022, RES 2022, SNDE 2022, AMEF 2022, SES 2022, MMF 2022, Kent Firm Dynamics Workshop 2021, Exeter Macro Workshop 2022 and Bristol for their helpful comments. We thank the following people for feedback: Jan de Loecker, Jan Eeckhout, Mark Bils, Alex Monge, Max Gillman, Mark Wright, Omar Licandro, Julian Neira, Tom Schmitz, Petr Sedláček, Danial Lashkari, John Morrow, Anthony Priolo, Riccardo Silvestrini and Kunal Sangani.

^{*}Competition and Markets Authority and University of Kent, joel.kariel@cma.gov.uk [†]University of Kent, a.savagar@kent.ac.uk.

This research is funded under ESRC project reference ES/V003364/1.

Recent technological advances, such as cloud computing, can raise scale economies allowing firms to expand at lower cost. But, as these technologies have emerged in economies such as the US and UK, productivity has stagnated. In this paper, we develop a theory to relate firm-level scale economies to aggregate productivity. We show that increases in scale economies should have increased aggregate productivity significantly. However, rising markups offset the productivity gains.

We make three contributions: first, we document rising scale economies from two determinants: higher returns to scale in variable production, which lowers marginal costs, and higher fixed costs. Second, we develop a tractable model to study the effect of these determinants of scale economies on aggregate productivity. Third, we conduct a quantitative exercise to replicate growing scale economies but stagnating productivity in the UK economy.

We develop a heterogeneous firm model with monopolistic competition, fixed costs, returns to scale and endogenous entry. We derive firm scale economies which is the ratio of average cost to marginal cost. Firm scale economies are endogenous. They consist of a fixed-cost component and a returns-to-scale component which dictates marginal cost.¹ The fixed cost and returns to scale determinants of scale economies lead to different aggregate productivity outcomes. Higher fixed costs may increase or decrease aggregate productivity depending on returns to scale in variable inputs at the firm. Higher returns to scale in variable inputs may increase or decrease aggregate productivity depending on the level of the markup.

We decompose aggregate productivity into allocative efficiency and technical efficiency. *Allocative efficiency* depends on the division of aggregate resources between firms. If there are increasing returns to scale, allocative efficiency improves when aggregate resources are concentrated on a small number of producers, since they exploit scale economies. Whereas with decreasing returns the opposite holds. The number of firms does not affect allocative efficiency when there are constant returns. *Technical efficiency* measures the average technology of *active* firms. Technology is an exogenous

¹Returns to scale are *returns to scale in variable inputs*. This measures the slope of the marginal cost curve and is the sum of output elasticities to variable inputs.

productivity characteristic that is revealed to firms upon entry. Given a technology draw, a firm decides to be active or inactive based on a period-by-period fixed cost. Therefore, technical efficiency captures the firm selection channel. That is, where the exogenous productivity distribution is truncated.

Our theoretical results show that rising scale economies, whether through fixed costs or returns to scale, strengthen selection, thus improving average technical efficiency. However, in high-markup environments, the selection channel is weaker. With high markups, selection weakens because small (low technology draw) firms get more revenue for each unit sold, so it is easier to cover fixed costs and survive. Allocative efficiency declines because markups increase the number of firms, constrain output, and therefore limit the exploitation of scale economies. Therefore, ceteris paribus, increases in scale economies should increase productivity.

Our quantitative exercise applies the theoretical insights to UK aggregate productivity. We show that estimated increases in returns to scale in variable production accompanied by estimated increases in markups replicate UK aggregate productivity dynamics well. Rising fixed costs cannot explain the data as well. If markups had not increased, the aggregate productivity of the UK would have been 20% higher through efficiency gains from scale economies.

Our paper abstracts from the specific technologies that may have changed scale economies, other than to characterise them by increasing fixed costs or raising returns to scale in variable production, which lowers marginal costs. Industry studies provide some qualitative insight. Ganapati (2021) shows that information technology reduced marginal costs and increased markups in the wholesale sector. For the manufacturing sector, Bloom, Garicano, Sadun, and Van Reenen (2014) study specific information technologies, such as enterprise resource planning, that increase managers' span of control and, therefore, lower marginal costs. Syverson (2019) hypothesises a shift towards products with lower marginal costs, such as software and pharmaceuticals. These examples seem relevant for a services-dominated economy like the UK. Therefore, our conclusion is that emerging technologies have increased returns to scale, which has decreased marginal costs and enhanced scale economies. These scale economies should translate into productivity gains. However, increasing market power limits the exploitation of scale economies and, in turn, productivity gains.

Related Literature

Recent work by Bilbiie and Melitz (2020), Edmond, Midrigan, and Xu (2021), and Baqaee, Farhi, and Sangani (2023) demonstrates the importance of returns to scale for aggregate analysis. The work is mostly focused on external returns to scale (love of variety) that arise from aggregation. However, Baqaee, Farhi, and Sangani (2023) also note that returns to scale at the firm level magnify aggregate returns to scale. Similarly to our analysis, the effects of scale economies are smaller in efficient (low markup) economies. Our analysis is parametric; we focus on the technical parameters of the production function that cause scale economies, and in turn affect aggregate TFP through technical and allocative efficiency. This is complementary to Baqaee and Farhi (2020) who provide non-parametric aggregation results for economies with scale economies. Both parametric and non-parametric approaches find that the role of allocative efficiency grows as distortions increase. And, we show that this is quantitatively relevant to replicate UK productivity dynamics.

In order to understand the consequenes of rising market power, De Loecker, Eeckhout, and Mongey (2021) present a quantitative model with oligopolistic competition and fixed costs. This allows them to compare the role of technology on the supply-side versus competitive factors on the demand-side. We differ by focusing on analytical results to understand the supply-side mechanisms through which different technologies affect scale economies, and in turn aggregate productivity. Our demand-side is restricted to monopolistic competition for tractability. Collectively, our papers advance the idea that to reconcile changing technologies on the supply side, market power must increase on the demand side.

Recent research in endogenous growth theory shows that changing technologies

affect firm cost structures, which in turn explains stagnating growth. Scale economies are not the direct focus, but they are implicit in the arguments. De Ridder (2019) models intangible inputs as reducing marginal costs and raising fixed costs. Aghion, Bergeaud, Boppart, Klenow, and Li (2019) model a fixed cost that increases with the number of product lines, but as technology improves, the fixed cost becomes less sensitive to the number of products.

Our model is a neoclassical growth model with heterogeneous firms based on Hopenhayn and Rogerson (1993), Restuccia and Rogerson (2008), and Barseghyan and DiCecio (2016). The model is similar to two-factor closed economy versions of Melitz (2003) and Ghironi and Melitz (2005). We include firm production with fixed costs and returns to scale similar to models by J. Kim (2004), Atkeson and P. J. Kehoe (2005), Bartelsman, Haltiwanger, and Scarpetta (2013), and D. Kim (2021).

Several recent articles provide estimates of returns to scale in the US economy. Gao and Kehrig (2021) estimate slightly decreasing returns to scale in US manufacturing firms. Using similar US data, Ruzic and Ho (2019) find a decline in returns to scale from 1982 to 2007. Using Compustat data, Chiavari (2022) documents rising returns to scale through production function estimation, and De Loecker, Eeckhout, and Unger (2020, Figure 7) documents increasing overhead cost shares as evidence of rising scale economies. Baqaee, Farhi, and Sangani (2023) also document economies of scale in US firms. Lashkari, Bauer, and Boussard (2019) find cost elasticity below one for French corporations, which implies economies of scale. For the UK economy, Oulton (1996), Harris and Lau (1998), and Girma and Görg (2002) document constant or slightly decreasing returns to scale for manufacturing firms.

1 Scale Economies

In this section, we define some concepts that are subject to ambiguity across fields.

Internal vs. External Returns to Scale: Our interest is internal returns to scale, not external returns to scale that arise from aggregation. Internal returns to scale and

scale economies arise within the firm from the production technology or fixed costs. External returns to scale are gains in aggregate output from changing aggregate inputs. They arise from grouping firms together.²

Scale Economies: Scale economies describe the response of firm costs to output changes. They are measured by the inverse cost elasticity, which is the average cost to marginal cost ratio.³

Returns to scale: Returns to scale are a property of the production technology. They are captured by the degree of homogeneity of the production function. On the cost side, this parameter represents the slope of a firm's marginal cost curve.⁴ For homothetic production functions, the scale elasticity of the cost function equals the returns to the scale of the production function.⁵ Fixed costs lead to non-homothetic production functions which break this relationship.

Imprecision over the terms scale economies and returns to scale extends beyond semantics. Erroneous conclusions and calibrations occur when the AC/MC ratio is estimated but is interpreted as the production function returns to scale.⁶

1.1 Graphical Intuition of Scale Economies

To aid understanding throughout the paper, it is helpful to present the cost curve scenarios of the production functions we consider. We define scale economies as the inverse cost elasticity, which is the ratio of average cost to marginal cost. With firm

²On the demand-side, with a consumption aggregator, the analogous concept is love-of-variety. Other terms used are 'thick markets' (Caballero and Lyons 1992), Ethier effects (Ethier 1982), and agglomeration effects (Krugman 1991).

³This definition of scale economies is common in industrial organization textbooks (Panzar 1989; Church and Ware 2000; Davis and Garcés 2009), recent examples are (Syverson 2019; Conlon, Miller, Otgon, and Yao 2023). It is sometimes recognised in macroeconomics, for example (Rotemberg and Woodford 1993; Basu 2008; Baqaee, Farhi, and Sangani 2023; Lashkari, Bauer, and Boussard 2019).

⁴Occasionally, researchers recognise this parameter as 'span of control' since it is mathematically analogous to the span of control parameter in Lucas (1978). In that context, it captures diminishing returns in managerial span of control. Hopenhayn (2014) analyses the equivalence with returns to scale.

⁵Silberberg and Suen (2000, Ch. 8) present traditional proofs.

⁶Basu (2008) discusses this in detail. Since homothetic production functions are common in macroeconomics, the term returns to scale is often used universally even in the presence of fixed costs.

output *y*, we have:

$$S(y) \equiv \left(\frac{\partial \mathcal{C}}{\partial y}\frac{y}{\mathcal{C}}\right)^{-1} = \frac{AC(y)}{MC(y)}$$

where AC $\equiv C/y$ and MC $\equiv \partial C/\partial y$. There are economies of scale if S(y) > 1; constant scale economies if S(y) = 1; and diseconomies of scale if S(y) < 1. Figure 1 presents a firm with a U-shaped average cost curve due to increasing marginal costs and fixed cost.⁷ At the intersection of average and marginal cost, a firm has constant scale economies. To the left there are economies of scale. To the right there are diseconomies of scale. Therefore, the S(y) curve shows that size and scale economies are negatively related at the firm level.⁸



Figure 1: Fixed Cost with Increasing MC, U-Shaped AC Curve

Profits, Markups and Scale Economies: Scale economies can be represented directly from the profit definition. This yields an expression based on market structure, namely markups and profits. Scale economies can also be written in terms of technical properties of the production function, namely fixed costs and the homogeneity parameter. This will depend on the production function and can be derived from the cost function or the production function.⁹ Consider the definition of profits as revenue minus

⁷In the appendix we present plots considering the three main cases that arise in our theory: a fixed cost with increasing, constant or decreasing marginal cost.

⁸In the appendix we present a graphical explanation of scale economies from the production side.

⁹In this paper we will show this for labor denominated fixed costs beginning with the production function. Savagar (2021) shows it for output-denominated fixed costs beginning with the cost function.

costs

$$Profit = Price \times Output - Cost = Revenue - Cost.$$

Divide by revenue, define AC=Cost/Output, and multiply by MC/MC, yields:

$$\frac{AC}{MC} = \frac{Price}{Marginal Cost} \left(1 - \frac{Profit}{Revenue}\right).$$

This shows that a firm's scale economies are its markup multiplied by its profit share remainder (*i.e.* total cost share).¹⁰ A firm that makes zero-profits has scale economies equal to its markup.¹¹ And, a firm with positive profits will have lower scale economies than the zero-profit firm. Higher scale economies imply higher markups lower profit shares.

2 **Empirical Motivation**

We are motivated by the presence of rising scale economies at the firm level, while aggregate measures of productivity are stagnating.

2.1 **Productivity**

Figure 2 shows UK aggregate TFP growth over time. Aggregate productivity growth increases until 2007 but then declines and stagnates. This captures the UK 'productivity puzzle' (Barnett, Batten, Chiu, Franklin, and Sebastia-Barriel 2014; Goodridge, Haskel, and Wallis 2016).

¹⁰The total cost share is the sum of the variable cost share and the fixed cost share.

¹¹This result was used in earlier empirical work on returns to scale, when profits in the US economy were close to zero (Basu and Fernald 1997).



TFP growth (aggregate) is from the Penn World Table 10.01 (Feenstra, Inklaar, and Timmer 2015), accessed from FRED: Total Factor Productivity at Constant National Prices for United Kingdom (RTF-PNAGBA632NRUG).

2.2 Returns to Scale in Variable Inputs

To measure returns to scale, we estimate firm-level production functions on UK data from the Annual Respondents Database (ARDx). The data contains approximately 50,000 firms each year, 11 million workers, and two-thirds of gross value added. Firms report a range of production data, including gross output, value added, labor, materials, and investment.¹²

We assume that each firm *j* has the following Cobb-Douglas production function

$$y_{jt} = A_{jt} k_{jt}^{\beta_k} \ell_{jt}^{\beta_\ell}$$

where y_t , k_{jt} , ℓ_{jt} are firm value-added (or gross output) and inputs of capital and labour.¹³ A_{jt} is a measure of firm-level technical efficiency which we do not observe. Our aim is to estimate the β_k and β_ℓ parameters which represent output elasticities. The sum of

¹²In the appendix, we provide details about the data, data cleaning, deflation, capital construction, SIC code matching, and summary statistics.

¹³Output is represented by value-added or gross output depending on the estimation methodology.

these output elasticities is returns to scale in variable inputs.

Production function estimation suffers from omitted variable bias. The bias occurs because the input variables are correlated with the unobserved firm-level technology term. We use several production function estimation methodologies which are designed to address this problem. Further details are available in Olley and Pakes (1996), Levinsohn and Petrin (2003), Ackerberg, Caves, and Frazer (2015), and Gandhi, Navarro, and Rivers (2020). Since we estimate Cobb-Douglas production functions, we obtain a single, time-invariant, coefficient for each input in the production function. We divide the data into sub-periods to estimate changes over time.

Figure 3 shows average returns to scale across firms in the UK over time using the estimation methodology of Gandhi, Navarro, and Rivers (2020). There is a rising trend in returns to scale, from weakly decreasing to above unity. In the appendix, we provide estimates at the industry level and for alternative estimation methodologies. All the results imply rising returns to scale.



Figure 3: UK RTS, 2001 - 2014

RTS are the sum of firm-level coefficients from a Cobb-Douglas, gross-output, production function estimated with the methodology of Gandhi, Navarro, and Rivers (2020). To obtain time-varying estimates of RTS, we estimate production functions over rolling windows.

2.3 Fixed Cost Share in Revenue

An alternative contributor to firm scale economies is the fixed cost share. In Figure 4, we use the administration expense share in revenue as a proxy for a companies' fixed cost share. The figure shows rising fixed cost shares which is consistent with rising scale at the firm level. Administration expenses in UK company accounts are the costs incurred by a company that are not directly related to the production, manufacture or sale of goods or services. In the Appendix we discuss the data in greater detail and provide examples of administrative costs.

Figure 4: Median Fixed Cost Share in Sales, Source: BvD FAME



The plot shows the median 'Administration Expenses' share in 'Turnover' for UK firms.

3 Model

The household side of the model follows a neoclassical growth setup. The production side of the economy has firm entry and exit, monopolistic competition, and production functions that have different sources of scale economies. There are two stages to the firm problem. First, a firm decides whether to pay a fixed, output-denominated, entry cost based on the expected profits they would receive from optimal production decisions. Second, given a firm has entered, it makes optimal production decisions. Upon entering, the firm receives a productivity draw at which point it decides whether to produce or not, and if so how much to produce. The decision to produce or not is based on whether producing output will generate enough revenue to cover a fixed, period-by-period, labour-denominated, overhead cost. At the end of the period all firms exit exogenously.

3.1 Households

A representative household maximizes lifetime utility subject to a budget constraint

$$\max_{\{C_t, K_{t+1}\}_{t=0}^{\infty}} \sum_{t=0}^{\infty} \beta^t \frac{C_t^{1-\sigma} - 1}{1 - \sigma}, \quad \beta \in (0, 1),$$

s.t. $C_t + I_t = r_t K_t + w_t L^s + \Pi_t + T_t$ (1)

$$I_t = K_{t+1} - (1 - \delta)K_t.$$
 (2)

Households own all firms in the economy and receive profit Π_t . T_t is a lump-sum transfer from the government which will equal to the entry fees paid by firms. Households supply a fixed amount of labour that is not time-varying, we normalize this to one:

$$L^{\rm s} = 1. \tag{3}$$

Households own the capital stock and rent it to firms at a rental rate r_t , hence the capital investment decision is part of the household problem. The household optimization problem satisfies the following condition

$$\left(\frac{C_{t+1}}{C_t}\right)^{\sigma} = \beta(r_{t+1} + (1-\delta)).$$
(4)

plus a transversality condition and the resource constraint.

3.2 Firms

3.2.1 Final goods producer

The final goods aggregator is

$$Y_{t} = N_{t} \left[\frac{1}{N_{t}} \int_{0}^{N_{t}} y_{t}(t)^{\frac{1}{\mu}} dt \right]^{\mu}.$$
(5)

There are N_t intermediate producers on the interval $t \in (0, N_t)$. The parameter $\mu \ge 1$ captures product substitutability.¹⁴ The aggregator has constant returns to scale.¹⁵

The maximization problem of the final goods producer is

$$\Pi_t^F = \max_{y_t(\iota)} \quad Y_t - \int_0^{N_t} p_t(\iota) y_t(\iota) d\iota$$
(6)

s.t.
$$Y_t = N_t \left[\frac{1}{N_t} \int_0^{N_t} y_t(\iota)^{\frac{1}{\mu}} d\iota \right]^{\mu}$$
 (7)

The firm is infinitesimal so firm level output does not affect Y_t . The first-order condition with respect to $y_t(t)$ gives the inverse-demand for a firm

$$p_t(t) = \left(\frac{N_t y_t(t)}{Y_t}\right)^{\frac{1-\mu}{\mu}}.$$
(8)

3.2.2 Intermediate goods producer

The timeline for the intermediate goods producer is as follows. The firm pays cost κ to enter. It receives a draw $j \in (0, 1)$ from an i.i.d uniform distribution which translates to productivity A(j). It then decides whether to produce which incurs a fixed overhead cost. If the firm does not produce it remains inactive which we refer to as endogenous exit. All firms, active and inactive, exit at the end of one period.

¹⁴Perfectly substitutable products $\mu = 1$ are admissible when intermediate producers have a fixed cost and increasing marginal cost ($\phi > 0$ and $\nu \in (0, 1)$). This is the case of perfect competition where profit maximizing intermediate producers take price as given. Under perfect competition all firms produce at the minimum on their average cost curves with perfectly-elastic, horizontal, demand curves.

¹⁵A typical CES production function would have the pre-multiplying term as N_t^{μ} , such that is cancels with the $1/N_t$ inside the square brackets. However, this creates increasing scale economies in aggregation. Since our interest is scale economies at the firm level, we remove this additional source of scale in aggregation.

The production function for a firm with productivity *j* is given by

$$y_t(j) = A(j) \left[k_t(j)^{\alpha} \ell_t(j)^{1-\alpha} \right]^{\nu}.$$
(9)

The parameter $0 < \alpha < 1$ captures the capital cost in total variable cost. The parameter $\nu > 0$ captures returns to scale in variable inputs. This represents the slope of the marginal cost curve or returns to scale in variable inputs. There are decreasing returns in variable production when $\nu \in (0, 1)$, constant returns when $\nu = 1$, and increasing returns when $\nu > 1$. As $\nu : 0 \rightarrow 1$ the marginal cost curve flattens which raises returns to scale (increases returns to scale in variable inputs), when $\nu = 1$ the marginal cost curve is flat, and as $\nu : 1 \rightarrow \infty$ the marginal cost curve is increasingly downward sloping.¹⁶ The labour employed to produce output is:

$$\ell_t(j) = \ell_t^{\text{tot}}(j) - \phi, \tag{10}$$

where $\ell_t^{\text{tot}}(j)$ represents the total labour employed by the firm, and ϕ is an overhead cost. Both ϕ and ν determine scale at the firm level.

The firm solves the following profit maximization problem:

$$\max_{k_t(j),\ell_t(j)} p_t(j) y_t(j) - r_t k_t(j) - w_t(\ell_t(j) + \phi)$$
(11)

subject to the production function (9) and inverse demand function (8). The optimality conditions imply constant factor shares in revenue:

$$\frac{r_t k_t(j)}{p_t(j) y_t(j)} = \frac{\nu}{\mu} \alpha \tag{12}$$

$$\frac{w_t \ell_t(j)}{p_t(j) y_t(j)} = \frac{\nu}{\mu} (1 - \alpha).$$
(13)

For the second-order conditions on profit maximization to hold, a necessary condition is: $v < \mu$. We present the first- and second-order conditions in Appendix D.1. Addi-

¹⁶We show that downward sloping MC curve must be shallower than the downward sloping demand curve to ensure a profit-maximizing equilibrium where MR = MC exists.

tionally, we assume $\alpha \nu < 1.^{17}$ Therefore, we assume the following upper-bound on returns to scale in variable inputs.

Assumption 1. Increasing returns in variables inputs are limited as follows:

$$\nu < \min\left\{\frac{1}{\alpha}, \mu\right\}. \tag{14}$$

This always holds with decreasing returns in variable inputs since $\nu < 1$. A higher markup and a lower capital cost share in variable costs allow for greater returns to scale in variable inputs. Empirically, the markup constraint is more likely to prevail. For example, markups of 1.25 and 1.5 are large but plausible, and capital shares in variable costs of 0.25 and 0.5, which are also plausible, give constraints of 4 and 2 which are higher than the markup constraints.

From the factor market equilibrium conditions, the ratio $v/\mu = (w\ell + rk)/py$ is variable cost share in revenue. The remaining share, $1 - (v/\mu)$, is the profit plus fixed cost share in revenue. Additionally, $\alpha = rk/(w\ell + rk)$ and $1 - \alpha = w\ell/(w\ell + rk)$ are the share of capital and production labour in variable costs. Also, $\alpha v = \mu(rk/py)$ is the capital share in revenue scaled by the markup.

3.2.3 Ratio of firm size

Firm output, revenue and inputs are proportional to productivity to the power of a constant $y(j)^{\frac{1}{\mu}}$, p(j)y(j), k(j), $\ell(j) \propto A(j)^{\frac{1}{\mu-\nu}}$. Consequently, for a given distribution of A(j) across firms, changes in μ and ν affect the distribution of labour, capital, revenue and output across firms.

The inverse demand condition and factor price equilibrium conditions imply that for any two firms, *i* and *j*, their relative revenue and input choices are proportional to

¹⁷This assumption is not required for profit maximization to hold. Imperfect competition ensures that firm-level revenue is concave in inputs, even if output is not concave in inputs. That is, marginal revenue products are decreasing in their respective inputs, even if marginal products are not. Specifically, $0 < \alpha v < 1$ ensures firm-level output is concave in capital, and aggregate output is concave in aggregate capital and not decreasing in aggregate labour.

their relative (scaled) productivity:

$$\frac{p_t(j)y_t(j)}{p_t(i)y_t(i)} = \frac{k_t(j)}{k_t(i)} = \frac{\ell_t(j)}{\ell_t(i)} = \left(\frac{A(j)}{A(i)}\right)^{\frac{1}{\mu-\nu}}, \quad \forall i, j.$$
(15)

Additionally, if we use equation (8) to substitute out p_t , we can write:

$$\frac{y_t(j)}{y_t(\iota)} = \left(\frac{A(j)}{A(\iota)}\right)^{\frac{\mu}{\mu-\nu}}.$$
(16)

3.2.4 Zero-profit firm

We assume there is a threshold productivity draw $J_t \in (0, 1)$ characterised by zero profits, which yields threshold technology \underline{A}_t . If a firm receives a productivity draw below the threshold productivity level they would make negative profits from production. Consequently, they prefer to produce zero and make zero profits. Therefore we define profits and characterise the threshold productivity as follows:

$$\pi_t(j) = p_t(j)y_t(j) - r_tk_t(j) - w_t(\ell_t(j) + \phi)$$
(17)

$$\pi_t(J_t) = 0. \tag{18}$$

A helpful reduced-form expression for profits combines the profit condition with equilibrium factor prices, with the zero-profit condition and with the ratio of revenues to scaled productivity:

$$\pi_t(j) = \phi w_t \left[\left(\frac{A(j)}{\underline{A}_t} \right)^{\frac{1}{\mu - \nu}} - 1 \right].$$
(19)

3.2.5 Free Entry

All firms die after one period. A firm only produces if it makes positive profits, hence firm value is given by

$$v_t(j) = \max\{\pi_t(j), 0\}.$$
 (20)

We assume a free entry condition which implies that the unconditional expected value from entering equals to the entry cost κ :

$$\mathbb{E}[v_t(j)] = \kappa. \tag{21}$$

The cost of entry κ is denominated in consumption units and will be rebated to households in a lump-sum. Combining (20) and (21) with our reduced-form profit expression (19) yields:

$$\phi w_t (1 - J_t) \left[\left(\frac{\hat{A}_t}{\underline{A}_t} \right)^{\frac{1}{\mu - \nu}} - 1 \right] = \kappa.$$
(22)

This shows that profits from being active multiplied by the probability of being active $1-J_t$ equals the entry cost. We have defined the power mean of technology, conditional on being active, as

$$\hat{A}(J_t) \equiv \mathbb{E}\left[A(j)^{\frac{1}{\mu-\nu}} \middle| j > J_t\right]^{\mu-\nu} = \left[\frac{1}{1-J_t} \int_{J_t}^1 A(j)^{\frac{1}{\mu-\nu}} dj\right]^{\mu-\nu}.$$
(23)

The power mean is a weighted average of firm-level productivity.¹⁸

3.3 Entry

Operating firms N_t are the subset of firms who decide to produce once receiving their productivity draw. Entrants E_t are all firms who pay the entry cost.

$$N_t = \int_0^{N_t} dt = E_t \int_{J_t}^1 dt = E_t (1 - J_t).$$
(24)

We can interpret the productivity cut-off J_t as the probability of exit and $1 - J_t$ as the probability of surviving.

¹⁸The term $\hat{A}(J_t)$ generalizes Melitz (eq. 7 2003, p. 1700) and Colciago and Silvestrini (eq. 31 2022, p. 10). This term is equivalent to these papers if $\nu = 1$ and the markup is expressed in terms of elasticities of substitution between goods, for example $\mu = \theta/(\theta - 1)$ where θ is the elasticity parameter. Notably, with $\nu \neq 1$, we *cannot* represent the power mean of technology \hat{A} as an output-weighted harmonic average of unscaled technology draws.

3.4 Aggregation

To get aggregate output and aggregate inputs, we use that the index of operating firms $(0, N_t)$ is equivalent to the measure of entering firms E_t constrained over the region of operation $(J_t, 1)$.

3.4.1 Aggregate Factor Inputs

Aggregate labour is comprised of production labour and non-production labour

$$K_{t} = \int_{0}^{N_{t}} k_{t}(i) di = E_{t} \int_{J_{t}}^{1} k_{t}(j) dj$$
(25)

$$L_{t} = \int_{0}^{N_{t}} \left[\ell_{t}(i) + \phi \right] di = E_{t} \int_{J_{t}}^{1} \left[\ell_{t}(j) + \phi \right] dj.$$
(26)

We define u_t as the fraction of aggregate labour that goes to production

$$u_{t} \equiv \frac{E_{t} \int_{J_{t}}^{1} \ell(j) dj}{L_{t}} = \frac{\int_{0}^{N_{t}} \ell_{t}(i) di}{L_{t}}$$
(27)

$$1 - u_t = \frac{E_t (1 - J_t)\phi}{L_t} = \frac{N_t \phi}{L_t}.$$
 (28)

3.4.2 Aggregate Output

We can express aggregate output as:

$$Y_{t} = N_{t}^{1-\nu} \hat{A}_{t} \left[K_{t}^{\alpha} \left(u_{t} L_{t} \right)^{1-\alpha} \right]^{\nu}$$
(29)

We present a derivation of this result in the appendix. The result shows that aggregate output is the sum of N_t firms, which are homogeneous, *i.e.* dividing aggregate resources evenly $\left[(K_t/N_t)^{\alpha} (u_t L_t/N_t)^{1-\alpha} \right]^{\nu}$, and each endowed with a power mean of technology \hat{A}_t . In equation (29) there are constant returns in capital, production labour and number of firms. That is, if $K_t, u_t L_t, N_t$ change by a fixed proportion, then aggregate output will change by this fixed proportion. If N_t is treated as a fixed factor, then external returns to scale in aggregate capital and production labour is given by ν .

3.4.3 Aggregate Factor Market Equilibrium

The wage, rental rate on capital and zero-profit condition are

$$r_t = \alpha \frac{\nu}{\mu} \frac{Y_t}{K_t} \tag{30}$$

$$w_t = (1 - \alpha) \frac{\nu}{\mu} \frac{Y_t}{u_t L_t} \tag{31}$$

$$\frac{w_t}{Y_t} \frac{N_t \phi}{L_t} = \left(1 - \frac{\nu}{\mu}\right) \left(\frac{\underline{A}_t}{\hat{A}_t}\right)^{\frac{1}{\mu - \nu}}$$
(32)

3.5 Government Budget Constraint and Resource Constraints

The resource constraint is

$$Y_t = C_t + I_t. aga{33}$$

The government rebates entry fees to households. The government budget constraint equates taxes to government expenditure

$$T_t = E_t \kappa. \tag{34}$$

Profits and labour markets clear:

$$\Pi_t = \Pi_t^F \tag{35}$$

$$L_t = L^s. aga{36}$$

Aggregate profits received by the household from owning firms equate to profits earned by the final goods producer. The profits are zero in equilibrium. Labour demanded by the firm equates to labour supplied by the household which is normalised to 1.

3.6 Equilibrium Definition

An equilibrium is a sequence of prices $\{r_t, w_t\}_{t=0}^{\infty}$; firm capital and labour demands $\{\ell_t(j), k_t(j)\}_{t=0}^{\infty}$; firms' operating decisions to be active or inactive, measures of entry and active firms $\{E_t, N_t\}_{t=0}^{\infty}$; consumption and capital $\{C_t, K_{t+1}\}_{t=0}^{\infty}$, such that

- 1. households choose *C* and *K* optimally by solving problem (1);
- 2. firms compete under monopolistic competition and decide optimally whether to produce or remain inactive, and demand factors according to (11);
- 3. the free entry condition holds (21);
- 4. markets clear for aggregate labour (26), aggregate capital (25), goods market (33), labour market (36) and aggregate profits (35);
- 5. the government budget constraint is satisfied (34).

3.7 Model Characteristics

Before imposing a Pareto distribution on the technology draws A(j), we can make some general characterisations from the model equilibrium conditions.

3.7.1 Aggregate Labour Utilized for Production

The level of aggregate labour utilized for production u_t is a function of J_t only. In equation (32), use $N_t\phi/L_t = 1 - u_t$. Hence, there are two equations, (31) and (32), determining the wage as a function of u and J. If we equate these two wage equations, we get the level of utilization as a function of J:

$$u_t = \left[1 + \frac{1}{1 - \alpha} \left(\frac{\mu}{\nu} - 1\right) \left(\frac{\underline{A}_t}{\hat{A}_t}\right)^{\frac{1}{\mu - \nu}}\right]^{-1}.$$

In turn, by equation (31) this implies that the aggregate labour share $w_t L_t/Y_t$ is only a function of J_t .

3.7.2 Aggregate Productivity

We can rearrange aggregate output into Cobb-Douglas form, where we use $N_t = (1 - u_t)L_t/\phi$ for $\phi > 0$, which gives:

$$Y_t = \text{TFP}_t K_t^{\alpha \nu} L_t^{1 - \alpha \nu}$$
(37)

where,
$$\text{TFP}_t \equiv \left(\frac{N_t}{L_t}\right)^{1-\nu} \left(1 - \frac{N_t \phi}{L_t}\right)^{(1-\alpha)\nu} \hat{A}_t$$
 (38)

$$= \left(\frac{1-u_t}{\phi}\right)^{1-\nu} u_t^{(1-\alpha)\nu} \hat{A}_t \tag{39}$$

Aggregate output exhibits constant (external) returns to scale in aggregate capital and aggregate labour when firms are treated as a fixed factor.¹⁹ Aggregate total factor productivity (TFP) measures aggregate output that is not accounted for by aggregate capital and aggregate labour. It includes labour utilized for production u_t and overheads $1 - u_t$, as well as two sources of returns to scale, ϕ and ν , and the capital share in variable costs α . TFP is not the Solow residual because the exponents of aggregate capital and labour do not correspond to aggregate factor shares.²⁰ It is helpful to decompose TFP into allocative efficiency and technical efficiency:

$$TFP_t = \underbrace{\Omega_t}_{\text{allocative technical}} \times \underbrace{\hat{A}_t}_{\text{technical}}.$$
(40)

We define \hat{A}_t as technical efficiency, and we define allocative efficiency as:

$$\Omega_{t} \equiv \underbrace{\left(\frac{N_{t}}{L_{t}}\right)^{1-\nu}}_{\text{Scale effect}} \times \underbrace{\left(1 - \frac{N_{t}\phi}{L_{t}}\right)^{(1-\alpha)\nu}}_{\text{Resource duplication}}.$$

¹⁹To see this consider the sum of coefficients in the log differenced equation:

$$d\ln Y_t = d\ln \mathrm{TFP}_t + \alpha \nu d\ln K_t + (1 - \alpha \nu) d\ln L_t.$$

²⁰The term $\alpha \nu$ is the aggregate capital share in output multiplied by the markup $\alpha \nu = \mu \times rK/Y$.

Allocative efficiency captures the negative effect of more firms duplicating fixed costs, and the scale effect of dividing aggregate labour among more firms, which will depend on returns to scale in variable production $\nu \ge 1.^{21}$ Technical efficiency is the generalised mean, conditional on being active, of exogenously drawn technology, and hence it is determined by selection. Under Pareto distributed A(j), technical efficiency will be a linear function of the threshold productivity level <u>A</u>.

3.7.3 Scale Economies

What are scale economies in this model? The parameters ν and ϕ are both sources of scale economies in the model. Scale economies are measured as the ratio of average cost to marginal cost (the inverse cost elasticity) which is an endogenous object. In this section, we show this from the production side by summing output elasticities. The same result can be shown from the cost function.²²

From equations (9) and (10), the response of firm output to a change in each variable input is constant. Consequently, returns to scale in variable production is constant:

$$\frac{\partial \ln y_t(j)}{\partial \ln k_t(j)} = \nu \alpha, \quad \frac{\partial \ln y_t(j)}{\partial \ln \ell_t(j)} = \nu(1-\alpha), \quad \frac{\partial \ln y_t(j)}{\partial \ln k_t(j)} + \frac{\partial \ln y_t(j)}{\partial \ln \ell_t(j)} = \nu.$$

The effect of a change in total labour input is decreasing in firm size:²³

$$\frac{\partial \ln y_t(j)}{\partial \ln \ell_t^{\text{tot}}(j)} = \nu(1-\alpha) \left(1 + \frac{\phi}{\ell_t(j)} \right) = \nu(1-\alpha) + (\mu-\nu) \left(\frac{\underline{A}_t}{A(j)} \right)^{\frac{1}{\mu-\nu}} \quad \in (\nu(1-\alpha), \mu-\alpha\nu)$$

²³For the second equality, we use the zero-profit condition
$$\left(1 - \frac{\nu}{\mu}\right)p_t(j)y_t(j) = w_t\phi\left(\frac{A(j)}{\underline{A}_t}\right)^{\mu-\nu}$$
 combined with labour demand $\frac{w_t}{p_t(j)y_t(j)} = \frac{\nu(1-\alpha)}{\mu}\frac{1}{\ell(j)}$ to yield $\nu(1-\alpha)\frac{\phi}{\ell_t(j)} = (\mu-\nu)\left(\frac{\underline{A}_t}{\overline{A}(j)}\right)^{\frac{1}{\mu-\nu}}$.

1

²¹In the appendix we show how allocative efficiency is affected by the number of firms, when we assume Pareto distributed A(j). For $\nu < 1$ there is a number of firms which maximizes allocative efficiency, and for $\nu \ge 1$ allocative efficiency always falls in N_t , in which case a well-defined firm size relies on the markup giving sufficiently downward-sloping demand.

²²Savagar (2021) shows this for a model with output denominated fixed costs.

Therefore, scale economies at the firm are decreasing in firm size:

$$S_t(j) \equiv \frac{\partial \ln y_t(j)}{\partial \ln k_t(j)} + \frac{\partial \ln y_t(j)}{\partial \ln \ell_t^{tot}(j)} = \nu \left(1 + (1 - \alpha)\frac{\phi}{\ell_t(j)}\right) = \nu + (\mu - \nu) \left(\frac{\underline{A}_t}{A(j)}\right)^{\frac{1}{\mu - \nu}} \quad \in (\nu, \mu).$$
(41)

To be more precise, a firm's scale economies decrease as production labour rises relative to the labour overhead, or as firm productivity rises relative to the productivity cut-off.

Figure 5 plots (41) for a given <u>A</u>. More productive firms have lower scale economies. The cut-off firm has the highest level of scale equals to the markup, and scale converges on returns to scale in variable inputs ν for high-productivity firms.



Figure 5: Firm-level Scale Economies in Steady-State

Plot shows equation (41) scale of a firm given its productivity draw. In the shaded region firms are inactive and the dashed line shows their hypothetical scale economies if they were to produce. The horizontal lines show the bounds on scale economies of active firms $S(j) \in (\nu, \mu)$. We have assumed A(j) is Pareto distribution and we have set <u>A</u> arbitrarily.

3.8 Model with Pareto Distribution

We assume that the technology variable is Pareto distributed. Given a random variable j drawn from the uniform distribution on the unit interval [0, 1), then the productivity

variable A(j) given by the quantile function:

$$A(j) = \frac{h}{(1-j)^{\frac{1}{\vartheta}}}.$$
 (42)

The parameter $\vartheta > 1$ is the Pareto shape parameter and *h* is the scale parameter, which is the lowest value of technology, corresponding to j = 0. We set h = 1. A thicker-tailed Pareto distribution occurs as $\vartheta \to 1$, which implies a higher density of high-productivity draws and a lower density of low-productivity draws. A thinner-tailed Pareto distribution occurs as $\vartheta \to \infty$ which implies a lower density of high-productivity draws and a higher likelihood of low-productivity draws.

Under Pareto, the power mean of technology is:

$$\hat{A}_{t} = \left(\frac{\vartheta(\mu - \nu)}{\vartheta(\mu - \nu) - 1}\right)^{\mu - \nu} \underline{A}_{t} = \Gamma \underline{A}_{t} \quad \text{where } \Gamma \equiv \left(\frac{\vartheta(\mu - \nu)}{\vartheta(\mu - \nu) - 1}\right)^{\mu - \nu}.$$
(43)

The constant Γ is the unconditional expectation of scaled technology $A(j)^{\frac{1}{\mu-\nu}}$. If the cutoff took its minimum value $\underline{A}_t = 1$, such that all participants were active and there was no selection $J_t = 0$, this represents the average technology that would arise. To ensure that scaled technology $A(j)^{\frac{1}{\mu-\nu}}$ has a finite expectation, we require that the scaled Pareto shape parameter satisfies the following assumption.

$$\vartheta(\mu - \nu) > 1. \tag{44}$$

This limits the degree of fat tails in the technology distribution. The assumption is analogous to the assumption $\vartheta > 1$ for the Pareto distributed technology before it is scaled.

3.8.1 Equilibrium Conditions with Pareto Distribution

Given the constant ratio between the power mean of technology and cut-off technology in equation (43), several equilibrium conditions simplify. Labour utilized for production is constant, aggregate TFP is a linear function of cut-off technology, and wage is a

log-linear function of cut-off technology:

$$u = \left(1 + \frac{\vartheta(\mu - \nu) - 1}{\upsilon \vartheta(1 - \alpha)}\right)^{-1} \qquad 1 - u = \frac{\vartheta(\mu - \nu) - 1}{\vartheta(\mu - \alpha \nu) - 1} \tag{45}$$

$$TFP_t = \Omega \hat{A}_t$$
, where $\Omega \equiv \left(\frac{1-u}{\phi}\right)^{1-\nu} u^{(1-\alpha)\nu}$ and $\hat{A}_t = \Gamma \underline{A}_t$ (46)

$$w_t = \frac{\kappa}{\phi} \left[\vartheta(\mu - \nu) - 1\right] \underline{A}_t^{\vartheta}.$$
(47)

The final equation determines the wage from the free entry condition. The lowest value \underline{A}_t can take is 1 which is the lowest productivity draw corresponding to J = 0. The constant u implies that total production labour is always a fixed portion of aggregate labour as an economy transitions over time.²⁴ Labour utilized for production is invariant to the fixed cost, increasing in returns to scale in variable inputs and decreasing in the markup:

$$\frac{du}{d\phi} = 0 \qquad \frac{du}{d\nu} = \frac{(1-\alpha)\vartheta(\vartheta\mu - 1)}{(\vartheta(\mu - \alpha\nu) - 1)^2} > 0 \qquad \frac{du}{d\mu} = -\frac{(1-\alpha)\vartheta^2\nu}{(\vartheta(\mu - \alpha\nu) - 1)^2} < 0$$
(48)

The constant u implies that the number of active firms is constant

$$N = \frac{1-u}{\phi} = \frac{1}{\phi} \frac{\vartheta(\mu - \nu) - 1}{\vartheta(\mu - \alpha\nu) - 1}.$$
(49)

Therefore, we can characterise the number of active firms as decreasing in the fixed cost and returns to scale in variable inputs, and increasing in the markup:

$$\frac{dN}{d\phi} = -\frac{N}{\phi} < 0 \qquad \frac{dN}{d\nu} = -\frac{1}{\phi}\frac{du}{d\nu} < 0 \qquad \frac{dN}{d\mu} = \frac{1}{\phi}\frac{du}{d\mu} > 0.$$
(50)

Both sources of scale economies ϕ and ν reduce the number of firms as they increase. As the marginal cost curve becomes flatter $\nu < 1$, horizontal $\nu = 1$ and downward sloping $\nu > 1$, optimal firm size (MR=MC) increases, and more total labour goes toward production (i.e. *u* rises). With larger firms the number of firms declines. For

²⁴In the appendix we provide a diagram to illustrate that this occurs because entry E_t and proportion of active J_t firms adjust over time to keep N fixed.

a rise in the fixed cost ϕ , the fraction of production labour in aggregate labour does not change, consequently the number of firms must fall so that the fraction of total fixed costs in labour also remains unchanged. Unlike the decline in *N* from rising *v* or ϕ , an increase in the markup raises the number of firms. This is because with higher markups, due to greater product differentiation (i.e. more steeply downward sloping demand curves), firms restrict their output more. This leads to the well-studied result that there is excessive entry of 'small' firms under monopolistic competition (Dixit and Stiglitz (1977) and Mankiw and Whinston (1986)). Finally, an implication of constant *u* and *N* is that the aggregate labour share $w_t L_t/Y_t$ is constant with a Pareto distribution, where $L_t = 1$ we have

$$\frac{w}{Y} = \frac{1}{\mu} \left(\mu - \alpha \nu - \frac{1}{\vartheta} \right).$$
(51)

The labour share is increasing in the markup, and decreasing in returns to scale in variable inputs and invariant to the fixed cost. The remaining model equations are unchanged:

$$Y_t - C_t = K_{t+1} - (1 - \delta)K_t$$
$$\left(\frac{C_{t+1}}{C_t}\right)^{\sigma} = \beta \left[r_{t+1} + (1 - \delta)\right]$$
$$Y_t = \text{TFP}_t K_t^{\alpha \nu}$$
$$r_t = \frac{\nu}{\mu} \alpha \frac{Y_t}{K_t}$$
$$w_t = \frac{\nu}{\mu} (1 - \alpha) \frac{Y_t}{u}$$

Therefore, we have reduced the model to seven equations in seven variables $C_t, K_t, Y_t, r_t, w_t, TFP_t, \underline{A}_t$, and u is a constant. In the reduced model there is no individual firm heterogeneity j. The model is an economy of homogeneous firms, each endowed with the power mean of technology \hat{A}_t , which captures all heterogeneity. We can further reduce the equilibrium conditions to two dynamic equations in two variables $\{C_t, K_t\}$.

First, if we equate wages and substitute out Y_t , we get \underline{A}_t as a function of K_t :²⁵

$$\underline{\mathbf{A}}_{t} = \Psi K_{t}^{\frac{\alpha \nu}{\vartheta - 1}}, \quad \text{where } \Psi \equiv \left(\frac{\phi}{\kappa \left[\vartheta(\mu - \nu) - 1\right]} (1 - \alpha) \frac{\nu}{\mu} \frac{\Omega \Gamma}{u}\right)^{\frac{1}{\vartheta - 1}}.$$
(52)

In turn, TFP, wage, rental rate and aggregate output are functions of capital:

$$TFP_t = \Omega \Gamma \Psi K_t^{\frac{\alpha \nu}{\vartheta - 1}}$$
(53)

$$w_t = \frac{\kappa}{\phi} \left[\vartheta(\mu - \nu) - 1\right] \Psi^{\vartheta} K_t^{\frac{\alpha \nu \vartheta}{\vartheta - 1}}$$
(54)

$$r_t = \alpha \frac{\nu}{\mu} \Omega \Gamma \Psi K_t^{\frac{\alpha \nu \vartheta}{\vartheta - 1} - 1}$$
(55)

$$Y_t = \Omega \Gamma \Psi K_t^{\frac{\alpha \nu \vartheta}{\vartheta - 1}}.$$
(56)

Finally, substituting in the rental rate and aggregate output into the two dynamic equations reduces to a dynamic system in $\{K_t, C_t\}$:

$$\Omega\Gamma\Psi K_t^{\frac{\alpha\nu\vartheta}{\vartheta-1}} - C_t = K_{t+1} - (1-\delta)K_t$$
(57)

$$\left(\frac{C_{t+1}}{C_t}\right)^{\sigma} = \beta \left[\alpha \frac{\nu}{\mu} \Omega \Gamma \Psi K_{t+1}^{\frac{\alpha \nu \vartheta}{\vartheta - 1} - 1} + (1 - \delta) \right].$$
(58)

Threshold technology, TFP, wage and aggregate output are increasing in aggregate capital. The rental rate is ambiguously related to capital:

$$\frac{d\ln r_t}{d\ln K_t} = \frac{1 - \vartheta(1 - \alpha\nu)}{\vartheta - 1} \gtrless 0 \iff 1 - \vartheta(1 - \alpha\nu) \gtrless 0.$$

 25 Equating wages with Y_t substituted out yields

$$(1-\alpha)\frac{\nu}{\mu}\frac{\Omega\Gamma\underline{A}_{t}K_{t}^{\alpha\nu}}{u}=\frac{\kappa}{\phi}[\vartheta(\mu-\nu)-1]\underline{A}_{t}^{\varphi}.$$

For a given level of capital, \underline{A}_t adjusts such that the wage markets equate. This relationship gives the intuition for why an increase in capital increases selection. We begin with capital as it is a state variable, determined directly in steady state. From the left-hand wage equation, an increase in capital increases the wage given \underline{A}_t on the left held constant. On the right, which represents wage from the free entry condition, \underline{A}_t must increase – since $\vartheta > 1$, increasing \underline{A}_t on the right-hand side has a stronger wage enhancing effect than increasing \underline{A}_t on the left-hand side. To summarise, an increase in K, increase w in the factor market equilibrium, therefore \underline{A}_t must increase to raise wage in the free entry condition (i.e. a higher wage means only more productive firms survive). We can think of this relationship as two wage curves $\ln w = \ln \underline{A}$ and $\ln w = \vartheta \ln \underline{A}$, since $\vartheta > 1$ wage is more sensitive to selection in the free entry condition than in the factor market condition.

To understand this ambiguity, consider that $r_t = \alpha \frac{\nu}{\mu} Y_t / K_t$. Since $Y_t = TFP_t K_t^{\alpha\nu} = \Omega \Gamma \Psi K_t^{\frac{\alpha\nu}{\vartheta-1}} \times K_t^{\alpha\nu}$, therefore $Y_t / K_t = \Omega \Gamma \Psi K_t^{\frac{\alpha\nu}{\vartheta-1}} \times K_t^{\alpha\nu-1}$ where $\alpha \nu - 1 < 0$ by assumption. Heterogeneity ϑ matters through the TFP_t component. If ϑ decreases, Pareto tails become fatter and there is a greater density of high technology draws i.e. more heterogeneity. This strengthens the TFP response to aggregate capital, and consequently aggregate output responds more to aggregate capital, such that aggregate output could increase at an increasing rate in capital. If ϑ increases, Pareto tails become thinner and there is a greater density of low technology draws i.e. no heterogeneity, TFP responds less to capital, and r_t will decrease in K_t . For the remainder of the paper, we impose that the price of capital decreases in the quantity of aggregate capital, therefore:

$$1 - \vartheta(1 - \alpha \nu) < 0. \tag{59}$$

Therefore, we have imposed two assumptions on the Pareto shape parameter. These assumptions restrict the thickness of the Pareto tail, and are summarised in the next assumption.

Assumption 2. The Pareto shape parameter must satisfy

$$\frac{1}{\vartheta} < \min\{\mu - \nu, 1 - \alpha\nu\}.$$
(60)

This ensures that the scaled technology distribution is Pareto distributed with a finite mean and the price of capital *r* is decreasing in aggregate capital *K*. For example, consider a markup of $\mu = 1.2$, a constant marginal cost curve $\nu = 1$, and a capital share in variable costs of $\alpha = 0.25$. Then, we require $\frac{1}{\vartheta} < \min\{0.2, 0.75\}$, so the Pareto shape parameter must satisfy $\vartheta > 5$.

3.8.2 Steady-state with Pareto Distribution

In steady state the system satisfies $K_{t+1} = K_t = K$ and $C_{t+1} = C_t = C$. This yields the following steady-state solution for capital and consumption:

$$K = \left[\frac{\alpha \nu \Omega \Gamma \Psi}{\mu r}\right]^{\frac{\vartheta - 1}{\vartheta(1 - \alpha \nu) - 1}} \tag{61}$$

$$C = K \left(\frac{\mu r}{\alpha \nu} - \delta\right). \tag{62}$$

where $r = \frac{1}{\beta} - (1 - \delta)$. The remaining variables follow by substituting the expression for *K* into the reduced model. In particular, solving for <u>A</u> yields:

$$\underline{\mathbf{A}} = \left[\nu^{\nu} \frac{1}{\mu} \left(\frac{\alpha}{r}\right)^{\alpha\nu} (\phi(1-\alpha))^{\nu(1-\alpha)} \vartheta^{\mu-1} (\mu-\nu)^{\mu-\nu} \frac{1}{\kappa^{1-\alpha\nu}} \frac{1}{(\vartheta(\mu-\nu)-1)^{\mu-\alpha\nu}}\right]^{\frac{1}{\vartheta(1-\alpha\nu)-1}}.$$
 (63)

This captures how the technology cut-off, which represents *selection*, responds to underlying model parameters. An increase in <u>A</u> implies stronger selection and a decrease in <u>A</u> implies weaker selection.²⁶

The threshold technology is log linear in the entry cost κ and overhead costs ϕ . They have opposite effects on threshold technology <u>A</u>. We can understand this through the free entry condition which states that the constant entry cost equals to the expected value of a firm. If κ increases, then the expected value of a firm must rise to be at equality with κ . Given $\vartheta(1 - \alpha \nu) - 1 > 0$, a rise in expected profits will occur if the survival threshold falls, so there is more chance of being active upon entry.

Higher ϕ increases expected profits because break-even firm size increases (this is ²⁶The threshold productivity cannot be lower than the minimum productivity which we have normalized to one ($\underline{A} \ge A_{min} = 1$). Therefore, we require that

$$1 \le \left[\nu^{\nu} \frac{1}{\mu} \left(\frac{\alpha}{r}\right)^{\alpha\nu} (\phi(1-\alpha))^{\nu(1-\alpha)} \vartheta^{\mu-1} (\mu-\nu)^{\mu-\nu} \frac{1}{\kappa^{1-\alpha\nu}} \frac{1}{(\vartheta(\mu-\nu)-1)^{\mu-\alpha\nu}}\right]^{\frac{1}{\vartheta(1-\alpha\nu)-1}}.$$

Consequently, we can constrain the entry cost parameter such that it satisfies

$$\kappa \leq \left[\nu^{\nu} \frac{1}{\mu} \left(\frac{\alpha}{r}\right)^{\alpha\nu} (\phi(1-\alpha))^{\nu(1-\alpha)} \vartheta^{\mu-1} (\mu-\nu)^{\mu-\nu} \frac{1}{(\vartheta(\mu-\nu)-1)^{\mu-\alpha\nu}}\right]^{\frac{1}{1-\alpha\nu}}.$$

If κ satisfies this with equality, then $\underline{A} = 1$ (and J = 0), therefore all entrants are active N = E.

a ceteris paribus statement assuming the threshold does not adjust), then to offset this increase in firm size and to reduce expected profits back to their equality with κ the productivity cut-off must rise.

If ϕ increases, the threshold must rise because only more productive firms generate sufficient revenue to cover the higher overhead costs.

4 Theoretical Analysis

Changes in aggregate productivity occur through an allocation component $d \ln \Omega$ and a technical efficiency component $d \ln \hat{A}$:

$$d\ln TFP = d\ln \Omega + d\ln \hat{A}$$

4.1 The Effect of Entry Cost on Aggregate Productivity

The entry cost κ does not affect allocative efficiency Ω , but affects technical efficiency \hat{A} . If the entry cost increases, then technical efficiency decreases because the threshold technology level falls, thus weakening selection.²⁷ Selection weakens as the entry cost increases because, by the free-entry condition, the expected value of the firm must increase. The expected value increases if the threshold productivity declines.

4.2 The Effect of Fixed Costs on Aggregate Productivity

Changes in fixed costs affect aggregate TFP through an allocation component and a technology component:

$$\frac{d\ln TFP}{d\ln \phi} = \frac{d\ln\Omega}{d\ln\phi} + \frac{d\ln\hat{A}}{d\ln\phi}$$

²⁷Barseghyan and DiCecio (2011) Study this in a perfectly competitive economy, where the entry cost is in terms of output κ/Y . They find empirical evidence that higher entry costs decrease aggregate TFP across countries.

Under Pareto, technical efficiency depends on the technology threshold <u>A</u> only since the constant Γ is invariant to ϕ , therefore:

$$\frac{d\ln\hat{A}}{d\ln\phi} = \frac{d\ln\Gamma}{d\ln\phi} + \frac{d\ln\underline{A}}{d\ln\phi} = \frac{\nu(1-\alpha)}{\vartheta(1-\alpha\nu)-1} > 0.$$

The technology threshold is increasing in the overhead cost if $\vartheta(1 - \alpha \nu) - 1 > 0$. This is the condition for the rental rate *r* to be decreasing in aggregate capital.

The allocation effect depends on the degree of returns to scale in variable production:

$$\frac{d\ln\Omega}{d\ln\phi} = -(1-\nu).$$

The result is independent of the Pareto distribution assumption. We can interpret the allocation effect through the number of firms. Note that $\Omega = \left(\frac{1-u}{\phi}\right)^{1-\nu} u^{(1-\alpha)\nu} = N^{1-\nu}u^{(1-\alpha)\nu}$ and u is independent of ϕ . An increase in ϕ , decreases the number of active firms. With increasing returns ($\nu > 1$), allocative efficiency is improved by having fewer firms, as they benefit more from the increasing returns. On the other hand, with decreasing returns ($\nu < 1$), then having fewer firms is detrimental to allocative efficiency, as the effect of decreasing returns is accentuated. Lastly, with constant returns ($\nu = 1$), the number of firms has no effect on allocative efficiency.

Combining the allocative and technical efficiency effects, shows that the response of aggregate TFP to a change in fixed costs will depend on the level of returns to scale in variable inputs v.

$$\frac{d\ln TFP}{d\ln \phi} = -(1-\nu) + \frac{\nu(1-\alpha)}{\vartheta(1-\alpha\nu) - 1}$$
(64)

Figure 6 simulates equation (64) for different values of ν based on our benchmark calibration (Table 1).



In Figure 7, we decompose the three cases from Figure 6.²⁸ Technical efficiency always rises as the fixed cost increases, while the allocative efficiency component is determined by $\nu \gtrless 1$, as previously discussed.

Figure 7: Effect of $\ln \phi$ on TFP decomposed into \hat{A} and Ω for different ν



²⁸The effect of ϕ on TFP is the same regardless of μ .

4.3 The Effect of Returns to Scale on Aggregate Productivity

The nonlinearity of Equation (63) in ν makes it difficult to obtain a closed-form expression for the influence of ν on TFP. Therefore, we present simulations to illustrate this effect.

Calibration

	Parameter	Value	Target
β	Discount rate	0.96	Real interest rate
δ	Depreciation rate	0.08	Office for National Statistics
ν	Variable RTS	0.99 - 1.05	ABS (authors' estimates)
μ	Markup	1.21 - 1.28	CMA (2022)
α	Capital share	0.25	ABS (authors' calculations)
θ	Pareto shape	10	Hopenhayn (2014)
κ	Entry cost	0.017	Model-implied maximum given range of ν , μ
ϕ	Overhead cost	0.135	Match share inactive firms

Table 1: Parameter Values for Comparative Statics

The model is calibrated as in Table 1. We set the discount factor β to match the average real interest rate of 2.08 percent over the period. To do this, we use the equation for steady-state interest rate $r = \frac{1}{\beta} + 1 - \delta$.²⁹ The depreciation rate δ is determined by a weighted-average from ONS data. Our estimates of the returns to scale ν come from our estimates of the production function using the estimation of Gandhi, Navarro, and Rivers (2020). Markup estimates are from CMA (2022). They use a different dataset and estimation strategy. The markup estimates are consistent with other studies that show rising markups over this time period (ONS 2022; Hwang, Savagar, and Kariel 2022). In the model $\frac{\alpha\nu}{\mu}$ is the capital share in revenue and $\frac{(1-\alpha)\nu}{\mu}$ is the production labour share in revenue. Given our ν and μ estimates, we set $\alpha = 0.25$ to match a capital share of 20%.³⁰

²⁹Data on UK long-term government bond and inflation used to compute the real interest rate from FRED database: IRLTLT01GBM156N and FPCPITOTLZGGBR.

³⁰The ratio ν/μ is the revenue elasticity, which is typically set to 0.85 in US studies (Restuccia and Rogerson 2008; Barseghyan and DiCecio 2011). Hopenhayn (2014) discusses this common calibration.

Our theory imposes restrictions on the Pareto shape paramter ϑ . First, $\vartheta > \frac{1}{\mu-\nu}$ which ensures scaled productivity is Pareto distributed and the first moment exists, and second, $\vartheta > \frac{1}{1-\alpha\nu}$ which ensures aggregate output is concave in aggregate capital, so that the interest rate is decreasing in aggregate capital.³¹ Our calibrated markup minus our estimated returns to scale $\mu - \nu$ is between 0.198 and 0.234 from 2001 - 2014. Therefore, our restrictions imply that we must set $\vartheta > 5$, similar to Hopenhayn (2014) who sets the Pareto shape between 5 and 10. We set ϑ to this upper bound.

The entry cost parameter κ and the fixed cost parameter ϕ must satisfy restrictions such that $J_t \in (0,1)$. Then to calibrate these parameters at an empirically plausible level, we target the $\kappa/\phi w$ ratio. Barseghyan and DiCecio (2011) report a range of values from industry studies. In most industries, the ratio is less than one, so entry costs are less than overhead costs. The average they report is 0.82. Our experiments vary ν , ϕ , μ parameters, so the entry-to-overhead cost ratio will vary as we change these values, but the outcome always remains below 1.

Figure 8 shows the effect of ν on aggregate productivity for different values of the markup μ . We observe that aggregate productivity rises unambiguously in ν in a low markup economy, but not when the markup is higher. Both the level and the slope of the relationship is falling in μ .

Our estimates for ν divided by our calibrated markup μ yield a ratio from 0.81 to 0.84 between 2001 and 2014.

³¹The first restriction implies that scaled technology, $A(j)^{\frac{1}{\mu-\nu}} = (1-j)^{-\frac{1}{\vartheta(\mu-\nu)}}$, is Pareto distributed. In some experiments, we take $\nu \to \mu$ from below, and this requires us to raise the value of ϑ . The relevant value for us is the scaled Pareto parameter $\vartheta(\mu - \nu)$, since labour is distributed proportionally to this term.

Figure 8: Effect of variable RTS on ln TFP for different levels of the markup



In TFP for calibrated model for a range of ν and μ .

In Figure 9 we provide a decomposition into technical efficiency and allocative efficiency for each of these markup cases. We observe that the weakening passthrough of returns to scale to TFP occurs because of weakening technical efficiency (i.e. less selection), and worsening allocative efficiency.



Figure 9: Effect of ν on TFP decomposed into \hat{A} and Ω for different μ

As returns to scale v increase, technical efficiency \hat{A} increases. This implies stronger selection of high A(j) firms. However, the effect is weaker as market power increases. Hence, in high-markup economies, there is weaker selection of high productivity firms as returns to scale increase.

Returns to scale ν have a U-shaped relationship with allocative efficiency. This occurs because an increase in ν decreases the number of firms. With decreasing returns ($\nu < 1$), fewer firms harm allocative efficiency. However, with increasing returns ($\nu > 1$), fewer firms improve allocative efficiency. As market power increases, the minimum point moves right, causing a wider range of declining allocative efficiency. This occurs because higher markups increase the number of firms. Hence, the benefits of growing returns to scale for allocative efficiency are counteracted by higher markups, reducing the size of firms and limiting their ability to benefit from increasing returns.

5 Quantitative Application

In Section 4, we examined the impact on aggregate productivity of the parameters of the production function that cause scale economies. We now analyse the quantitative plausibility of scale economies alongside stagnating productivity, which has occured in the US and UK in recent years. We find that changing returns to scale in variable inputs alongside rising markups explains the data well.

5.1 Rising Returns to Scale

We calibrate the parameter ν to our annual estimates from 2001 to 2014, while the parameter μ is set to annual estimates from CMA 2022. We set $\phi = 0.135$ such that the share of inactive firms is empirically plausible in our benchmark calibration.

Figure 10 compares the trends in TFP in the data and our model. It reveals a rise in both series prior to the Financial Crisis, followed by a sharp decline in the data and a more gradual decrease in the model. Fixing the markup to its 2001 value highlights the significant impact of rising returns to scale on aggregate productivity. If market power had remained constant, higher returns to scale would have boosted aggregate productivity by over 20% between 2001 and 2014. However, when we incorporate the simultaneous increase in markups and returns to scale, our estimated productivity trend aligns more closely with observed data.



Figure 10: TFP Growth: Model vs Data

We give the model estimates of μ and estimates of ν and solve in each year for steady-state to obtain the model-implied TFP. The TFP data series is from the Penn World Table 10.01 (Feenstra, Inklaar, and Timmer 2015), accessed from FRED: Total Factor Productivity at Constant National Prices for United Kingdom (RTFPNAGBA632NRUG).

5.2 Rising Overhead Costs

The rise in both ν and μ in the UK explains aggregate productivity growth well. However, our empirical evidence shows that payments to administration costs as a share of sales has increased for the median firm. We consider this data series as a proxy for $w_t \phi/Y_t$ in the model.³²

In Figure 11 we calibrate ϕ to match our estimates of this ratio. The results highlight the opposing response of aggregate TFP conditional on the level of ν that we discussed in our theoretical analysis. Therefore the level of returns to scale in variable production is crucial for the implied effect of changing overhead costs. In our estimates, ν is greater than one, which implies productivity should have risen 10% over the period.

³²Since changing ϕ has general equilibrium effects on w_t and Y_t , increasing this ratio does not necessarily mean ϕ increases each period. This is relevant because our theoretical analysis focuses on changing ϕ , not the ratio. However, in practice for our calibration, ϕ and the ratio move together.



Figure 11: TFP Growth: Model (fixed ν and μ , with variable ϕ) vs Data

We fix μ to its 2001 level and calibrate ϕ to match the overhead share in BvD data. We solve the model steady-state in each year to obtain the model-implied TFP. The TFP data series is from the Penn World Table 10.01 (Feenstra, Inklaar, and Timmer 2015), accessed from FRED: Total Factor Productivity at Constant National Prices for United Kingdom (RTFPNAGBA632NRUG).

In Figure 12, we also re-calibrate μ each year to match CMA (2022) estimates. In this case, aggregate TFP growth underperforms TFP growth in the data, regardless of returns to scale in variable production. Therefore, the markup effect dominates the fixed cost effect and we do not observe opposing dynamics for productivity conditional on $\nu \ge 1$.



Figure 12: TFP Growth: Model (fixed ν , with variable μ and ϕ) vs Data

We give the model estimates of μ and calibrate ϕ to match the overhead share in BvD data. We solve the model steady state in each year to obtain the model-implied TFP. The TFP data series is from the Penn World Table 10.01 (Feenstra, Inklaar, and Timmer 2015), accessed from FRED: Total Factor Productivity at Constant National Prices for United Kingdom (RTFPNAGBA632NRUG).

6 Conclusion

In this paper, we analyse the relationship between firm-level scale economies and aggregate productivity. First, we estimate that returns to scale in the UK have risen since 1998. Then we develop a theory to relate firm-level scale economies to aggregate productivity. We clarify that scale economies can arise through fixed costs or returns to scale in variable inputs. We show that these two sources of scale economies have different implications for aggregate productivity through their effect on firm selection and allocation of resources to fixed costs. Finally, we simulate the model with our estimated series for returns to scale in variable inputs. This shows that, ceteris paribus, higher scale should have raised aggregate productivity significantly. However, rising markups in the UK offset this mechanism, and the combined impacts of higher scale and higher market power can explain stagnant productivity growth in this period.

References

- Ackerberg, Daniel A., Kevin Caves, and Garth Frazer (2015). "Identification Properties of Recent Production Function Estimators". In: *Econometrica* 83.6, pp. 2411–2451.
- Aghion, Philippe, Antonin Bergeaud, Timo Boppart, Peter J. Klenow, and Huiyu Li (Nov. 2019). *A Theory of Falling Growth and Rising Rents*. NBER Working Papers 26448. National Bureau of Economic Research, Inc.
- Asturias, Jose, Sewon Hur, Timothy J. Kehoe, and Kim J. Ruhl (2022). "Firm Entry and Exit and Aggregate Growth". In: *American Economic Journal: Macroeconomics*.
- Atkeson, Andrew and Patrick J Kehoe (2005). "Modeling and measuring organization capital". In: *Journal of political Economy* 113.5, pp. 1026–1053.
- Baqaee, David and Emmanuel Farhi (May 2020). *Entry vs. Rents: Aggregation with Economies of Scale*. Working Paper 27140. National Bureau of Economic Research.
- Baqaee, David, Emmanuel Farhi, and Kunal Sangani (June 2023). "The Darwinian Returns to Scale". In: *The Review of Economic Studies*.
- Barnett, Alina, Sandra Batten, Adrian Chiu, Jeremy Franklin, and Maria Sebastia-Barriel (2014). "The UK productivity puzzle". In: *Bank of England Quarterly Bulletin*, Q2.
- Barseghyan, Levon and Riccardo DiCecio (Oct. 2011). "Entry costs, industry structure, and cross-country income and TFP differences". In: *Journal of Economic Theory* 146.5, pp. 1828–1851.
- (2016). "Externalities, endogenous productivity, and poverty traps". In: *European Economic Review* 85.C, pp. 112–126.
- Bartelsman, Eric, John Haltiwanger, and Stefano Scarpetta (2013). "Cross-country differences in productivity: The role of allocation and selection". In: *American economic review* 103.1, pp. 305–334.
- Basu, Susanto (2008). "Returns to Scale Measurement". In: *The New Palgrave Dictionary of Economics: Volume 1 8*. Ed. by Steven N. Durlauf and Lawrence E. Blume. London: Palgrave Macmillan UK, pp. 5559–5562.

- Basu, Susanto and John Fernald (1997). "Returns to scale in US production: Estimates and implications". In: *Journal of political economy* 105.2, pp. 249–283.
- Bilbiie, Florin O. and Marc J Melitz (Dec. 2020). *Aggregate-Demand Amplification of Supply Disruptions: The Entry-Exit Multiplier*. Working Paper 28258. National Bureau of Economic Research.
- Bloom, Nicholas, Luis Garicano, Raffaella Sadun, and John Van Reenen (2014). "The distinct effects of information technology and communication technology on firm organization". In: *Management Science* 60.12, pp. 2859–2885.
- Caballero, Ricardo J and Richard K Lyons (1992). "External effects in US procyclical productivity". In: *Journal of Monetary Economics* 29.2, pp. 209–225.
- Chiavari, Andrea (Oct. 2022). *Customer Accumulation, Returns to Scale, and Secular Trends*. Tech. rep. Working paper (Oct 2022).
- Church, Jeffrey R. and Roger Ware (2000). Industrial Organization: A Strategic Approach. Irwin McGraw-Hill.
- CMA (Apr. 2022). The State of UK Competition. Tech. rep.
- Colciago, Andrea and Riccardo Silvestrini (2022). "Monetary policy, productivity, and market concentration". In: *European Economic Review* 142, p. 103999.
- Conlon, Christopher, Nathan H. Miller, Tsolmon Otgon, and Yi Yao (May 2023). "Rising Markups, Rising Prices?" In: *AEA Papers and Proceedings* 113, pp. 279–83.
- Davis, P. and E. Garcés (2009). *Quantitative Techniques for Competition and Antitrust Analysis*. Princeton University Press.
- De Loecker, Jan, Jan Eeckhout, and Simon Mongey (2021). *Quantifying market power and business dynamism in the macroeconomy*. Tech. rep. Working paper (Sept 2021).
- De Loecker, Jan, Jan Eeckhout, and Gabriel Unger (Jan. 2020). "The Rise of Market Power and the Macroeconomic Implications*". In: *The Quarterly Journal of Economics* 135.2, pp. 561–644.
- De Ridder, Maarten (Mar. 2019). *Market Power and Innovation in the Intangible Economy*. Discussion Papers 1907. Centre for Macroeconomics (CFM).

- Dixit, Avinash K. and Joseph E. Stiglitz (1977). "Monopolistic Competition and Optimum Product Diversity". English. In: *The American Economic Review* 67.3, pp. 297– 308.
- Edmond, Chris, Virgiliu Midrigan, and Daniel Yi Xu (May 2021). *How Costly Are Markups?* NBER Working Papers 24800. National Bureau of Economic Research, Inc.
- Ethier, Wilfred J (1982). "National and international returns to scale in the modern theory of international trade". In: *The American Economic Review* 72.3, pp. 389–405.
- Feenstra, Robert C., Robert Inklaar, and Marcel P. Timmer (Oct. 2015). "The Next Generation of the Penn World Table". In: *American Economic Review* 105.10, pp. 3150–82.
- Ganapati, Sharat (2021). "The Modern Wholesaler: Global Sourcing, Domestic Distribution, and Scale Economies". In:
- Gandhi, Amit, Salvador Navarro, and David A. Rivers (2020). "On the Identification of Gross Output Production Functions". In: *Journal of Political Economy* 128.8, pp. 2973–3016.
- Gao, Wei and Matthias Kehrig (July 2021). "Returns to Scale, Productivity and Competition: Empirical Evidence from US Manufacturing and Construction Establishments". In: *Working Paper (Jul 2021)*.
- Ghironi, Fabio and Marc J Melitz (2005). "International trade and macroeconomic dynamics with heterogeneous firms". In: *The Quarterly Journal of Economics* 120.3, pp. 865–915.
- Girma, S. and H. Görg (2002). "Foreign Ownership, Returns to Scale and Productivity: Evidence from UK Manufacturing Establishments". In: *CEPR Discussion Paper Series*.
- Goodridge, Peter, Jonathan Haskel, and Gavin Wallis (2016). "Accounting for the UK Productivity Puzzle: A Decomposition and Predictions". In: *Economica*.

- Harris, Richard and Eunice Lau (Apr. 1998). "Verdoorn's law and increasing returns to scale in the UK regions, 1968–91: some new estimates based on the cointegration approach". In: *Oxford Economic Papers* 50.2, pp. 201–219.
- Hopenhayn, Hugo (2014). "Firms, misallocation, and aggregate productivity: A review". In: *Annual Review of Economics* 6.1, pp. 735–770.
- Hopenhayn, Hugo and Richard Rogerson (1993). "Job Turnover and Policy Evaluation:
 A General Equilibrium Analysis". In: *Journal of Political Economy* 101.5, pp. 915–938.
- Hwang, Kyung In, Anthony Savagar, and Joel Kariel (2022). *Market Power in the UK*. Working Papers.
- Kim, Daisoon (2021). "Economies of scale and international business cycles". In: *Journal of International Economics* 131, p. 103459.
- Kim, Jinill (2004). "What determines aggregate returns to scale?" In: Journal of Economic Dynamics and Control 28.8, pp. 1577–1594.
- Krugman, Paul (1991). "Increasing returns and economic geography". In: *Journal of political economy* 99.3, pp. 483–499.
- Lashkari, Danial, Arthur Bauer, and Jocelyn Boussard (2019). *Information Technology and Returns to Scale*. 2019 Meeting Papers 1380. Society for Economic Dynamics.
- Levinsohn, James and Amil Petrin (2003). "Estimating production functions using inputs to control for unobservables". In: *The Review of Economic Studies* 70.2, pp. 317– 341.
- Lucas, Robert E. (1978). "On the Size Distribution of Business Firms". In: *The Bell Journal of Economics* 9.2, pp. 508–523.
- Luttmer, Erzo G. J. (Aug. 2007). "Selection, Growth, and the Size Distribution of Firms". In: *The Quarterly Journal of Economics* 122.3, pp. 1103–1144.
- Mankiw, N Gregory and Michael D Whinston (1986). "Free entry and social inefficiency". In: *The RAND Journal of Economics*, pp. 48–58.
- Martin, Ralf (2002). *Building the capital stock*. CeRiBA Working Paper. The Centre for Research into Business Activity.

- Melitz, Marc J (Nov. 2003). "The impact of trade on intra-industry reallocations and aggregate industry productivity". In: *Econometrica* 71.6, pp. 1695–1725.
- Olley, G. Steven and Ariel Pakes (1996). "The dynamics of productivity in the telecommunications equipment industry". In: *Econometrica* 64.6, pp. 1263–1297.
- ONS (2022). Estimates of markups, market power and business dynamism from the Annual Business Survey, Great Britain: 1997 to 2019. Tech. rep. Office for National Statistics website.
- Oulton, Nicholas (1996). "Increasing Returns and Externalities in UK Manufacturing: Myth or Reality?" In: *The Journal of Industrial Economics* 44.1, pp. 99–113.
- Panzar, John C. (1989). "Technological determinants of firm and industry structure".In: vol. 1. Handbook of Industrial Organization. Elsevier, pp. 3–59.
- Restuccia, Diego and Richard Rogerson (2008). "Policy distortions and aggregate productivity with heterogeneous establishments". In: *Review of Economic dynamics* 11.4, pp. 707–720.
- Rotemberg, Julio J. and Michael Woodford (Oct. 1993). *Dynamic General Equilibrium Models with Imperfectly Competitive Product Markets*. Working Paper 4502. National Bureau of Economic Research.
- Ruzic, Dimitrije and Sui-Jade Ho (Aug. 2019). "Returns to Scale, Productivity Measurement, and Trends in U.S. Manufacturing Misallocation". In: *INSEAD working paper*.
- Savagar, Anthony (2021). "Measured productivity with endogenous markups and economic profits". In: *Journal of Economic Dynamics and Control* 133, p. 104232.
- Silberberg, E and W C Suen (2000). *The Structure of Economics: A Mathematical Analysis 3rd Edition*. McGraw-Hill.
- Syverson, Chad (2019). "Macroeconomics and market power: Context, implications, and open questions". In: *Journal of Economic Perspectives* 33.3, pp. 23–43.

Appendix

A Graphical Illustration of Scale Economies (Cost based)

It is helpful to consider the three types of cost curve scenarios faced by firms in our model.

Figures 13, 14 and 15 show a firm's cost curves for the case where there is a fixed cost and increasing, constant or decreasing marginal costs. The diagrams show average total cost (ATC), average variable cost (AVC), average fixed cost (AFC) and marginal cost (MC) as firm output varies. Specifically, total cost is the sum variable cost and a fixed cost: TC = VC + FC, and averages are the components when divided by output *y*. The demand curve (p(y)) and marginal revenue (MR) curve ($\frac{d p(y)y}{dy}$) are not shown. We can imagine them as horizontal in the perfectly competitive case and downward sloping with imperfect competition, for example, due to product differentiation. The first case (Figure 13) allows for a perfectly competitive equilibrium when the demand curve is horizontal and firms produce at minimum average cost. The second and third cases (Figure 14 and 15) require imperfect competition. The demand curve must be downward sloping for MR = MC to occur.

Figure 13 illustrates the cost curves of a firm with a fixed cost and increasing marginal cost curve. The firm's marginal cost intersects the average total cost at its minimum. This minimum point is the firm's *minimum efficient scale* (MES) which would arise under perfect competition and at this minimum the firm has constant scale. To the left-hand side of the MES the firm has economies of scale and to the right-hand side the firm has diseconomies of scale.



Figure 13: Fixed Cost with Increasing MC, U-Shaped AC Curve

Figure 14 has a constant marginal cost curve and a fixed cost, so there are globally decreasing returns and ATC=MC in the limit. In this case there must be a downward sloping demand curve for an equilibrium where MR = MC to exist. Any degree of slope in the demand curve is sufficient to give an equilibrium, unlike in the next example example which requires a sufficiently steep demand curve (or a sufficiently shallow decreasing marginal cost).



Figure 14: Fixed Cost with Constant MC, Globally Decreasing Returns

Figure 15 has a decreasing marginal cost and a fixed cost so there are global diseconomies of scale. In this case there must be a downward sloping demand curve for an equilibrium where MR = MC to exist. The demand curve must be steeper than the downward-sloping marginal cost curve to ensure this occurs.

Figure 15: Fixed Cost with Decreasing MC, Globally Decreasing Returns

B Graphical Explanation of Scale Economies (Production based)

Figure 16 illustrates scale economies from the production side. It conveys the counterintuitive idea that small firms have high scale economies, whilst large firms have low scale economies. The graph represents an economy where firm output is produced by production labour. In order to produce there is some overhead labour that is the same for both firms. Total labour is the sum of production labour and overhead labour. The figure shows that a 10% rise in total labour at a firm raises production labour by 100% for the small firm, but only 13% for the large firm. Therefore, a proportional change in inputs has a larger effect on output for the small firm.

Figure 16: Scale Economies for Large and Small Firm

C Pareto Distributed Productivity

We obtain a measure of productivity A(j) from a random draw on the unit interval $j \in [0,1]$ using inverse transform sampling. The Pareto CDF is given by

$$F(A; \vartheta) = 1 - \left(\frac{h}{A}\right)^{\vartheta}; \quad A \ge h > 0 \quad \text{and} \quad \vartheta > 0.$$

If $\mathcal{J} \sim Uniform(0,1]$, then for $j \in \mathcal{J}$, we have

$$1 - \left(\frac{h}{A}\right)^{\vartheta} = j$$

Therefore, the quantile function is

$$A(j) = h(1-j)^{-\frac{1}{\vartheta}}.$$

Typically we set the scale parameter, which is the minimum possible value of A, to h = 1. Calibrations of the shape parameter (tail index) vary, for example $\vartheta = 1.15$ in Barseghyan and DiCecio (2011) and $\vartheta = 1.06$ in Luttmer (2007) and $\vartheta = 6.10$ in Asturias, Hur, T. J. Kehoe, and Ruhl (2022). These estimates are set to match the firm size distribution in terms of employment since in these models A(j) is proportional to

employment, though, as below, scaling can affect this.

Figure 17: Productivity with Pareto Distribution, $h = 1, \vartheta = \{1.06, 1.15\}$. Domain $j \in (0:0.97)$

Figure 18 plots scaled technology $A(j)^{\frac{1}{\mu-\nu}}$ for different calibrations of $\nu = \{0.95, 1.00, 1.05\}$ given fixed values of $\mu = 1.1$ and $\vartheta = 50$. The Pareto shape parameter must be large such that $(\mu - \nu)\vartheta > 1$. The distribution of scaled technology is proportional to the distribution of labour, capital and revenue. We require $(\mu - \nu)\vartheta > 1$ so that the expected value of scaled technology is finite, and consequently the expected value of labour per firm, capital per firm and revenue per firm is not infinite.

We observe that a higher ν leads to a greater scaled technology for any given j draw. Since a higher ν decreases the tail index for scaled technology, it causes a lower density of firms to have low-productivity draws and a greater density of firms to have highproductivity draws. Therefore, it thickens the tail of the probability density function. Since employment, capital, and revenue are proportional to this, it also means the distribution of firms is denser towards large firms in terms of labour, capital and employment, and with a lower density of small firms.

Figure 18: Scaled Technology with Pareto Distribution, $h = 1, \vartheta = 50$ and $\mu = 1.1, \nu = \{0.95, 1.00, 1.05\}$. Domain $j \in (0:0.95)$.

D Additional Model Derivations

D.1 Profit Maximization Problem

First-Order Conditions

We drop time subscripts *t* and firm-specific notation *j*. Fixed parameters are $\{\nu, \mu, \alpha, \phi\}$ and endogenous variables are $\{N, Y, A, k, \ell, r, w\}$. The revenue function is

$$py = N^{\frac{1-\mu}{\mu}}Y^{\frac{\mu-1}{\mu}}y^{\frac{1}{\mu}} = N^{\frac{1-\mu}{\mu}}Y^{\frac{\mu-1}{\mu}}A^{\frac{1}{\mu}}k^{\frac{\alpha\nu}{\mu}}\ell^{\frac{(1-\alpha)\nu}{\mu}}.$$

The variables $\{N, Y, A, w, r\}$ are taken as given by the firm. The firm maximizes revenue less costs:

$$\max_{k,\ell} \quad p(k,\ell)y(k,\ell) - rk - w(\ell + \phi).$$

The first-order conditions of the maximization problem state that the marginal revenue product of labour (MRPL) and marginal revenue product of capital (MRPK) – *i.e.* the revenue derivatives with respect to labour and capital – equal to wage and rental rate at optimal choices:

$$MRPL = \frac{\nu(1-\alpha)}{\mu} \frac{p(k^{*}, \ell^{*})y(k^{*}, \ell^{*})}{\ell^{*}} = w$$
$$MRPK = \frac{\nu\alpha}{\mu} \frac{p(k^{*}, \ell^{*})y(k^{*}, \ell^{*})}{k^{*}} = r.$$

Since $0 < \alpha < 1$, $\mu \ge 1$, $\nu > 0$ the marginal revenue products are positive. Asterisk notation denotes the profit-maximizing levels of capital and labour.

Second-Order Conditions

The second-order conditions for maximization require that, at the optimal point $\{k^*, \ell^*\}$, the objective function is decreasing in capital and labour and the determinant of the Hessian of the objective function is positive. This implies that $MRPL_{\ell} < 0$ and $MRPK_k < 0$ where subscripts denote derivatives. And, $MRPL_{\ell}MRPK_k - MRPL_k^2 > 0$. First note:

$$MRPL_{k} = MRPK_{\ell} = \frac{\nu\alpha}{\mu} \frac{MRPL}{k^{*}} = \frac{\nu(1-\alpha)}{\mu} \frac{MRPK}{\ell^{*}}.$$

Therefore the following conditions must be satisfied:

$$MRPL_{\ell} = \left(\frac{\nu(1-\alpha)}{\mu} - 1\right) \frac{MRPL}{\ell^{*}} < 0$$
$$MRPK_{k} = \left(\frac{\nu\alpha}{\mu} - 1\right) \frac{MRPK}{k^{*}} < 0$$
$$MRPL_{\ell}MRPK_{k} - MRPL_{k}^{2} = \frac{MRPL \times MRPK}{k^{*}\ell^{*}} \left(1 - \frac{\nu}{\mu}\right) > 0$$

These conditions hold if $\nu < \mu$.

D.2 Reduced-form Aggregate Output

We can show that aggregate output reduces to a Cobb-Douglas function of capital and labour scaled by a power mean measure of technology.

$$Y_{t} = N_{t} \left[\frac{1}{N_{t}} \int_{0}^{N_{t}} y_{t}(t)^{\frac{1}{\mu}} dt \right]^{\mu} = N_{t} \left[\frac{E_{t}}{N_{t}} \int_{J_{t}}^{1} y_{t}(j)^{\frac{1}{\mu}} dj \right]^{\mu} = N_{t} \left[\frac{1}{1 - J_{t}} \int_{J_{t}}^{1} y_{t}(j)^{\frac{1}{\mu}} dj \right]^{\mu}$$
(65)

Next, we use the technique of expressing firm-level variables relative to the threshold firm variable, which in turn can be summarised by relative productivity. Here, we rewrite as the ratio of firm output $y_t(j)$ to threshold firm output $y_t(J_t)$, where threshold firm output is a constant over *j*:

$$Y_{t} = N_{t} y_{t}(J_{t}) \left[\frac{1}{1 - J_{t}} \int_{J_{t}}^{1} \left[\frac{y_{t}(j)}{y_{t}(J_{t})} \right]^{\frac{1}{\mu}} dj \right]^{\mu}$$
(66)

Use the result that

$$\left[\frac{y_t(j)}{y_t(J_t)}\right]^{\frac{1}{\mu}} = \frac{p_t(j)y_t(j)}{p_t(J_t)y_t(J_t)} = \left(\frac{A(j)}{\underline{A}_t}\right)^{\frac{1}{\mu-\nu}}$$
(67)

Hence

$$Y_t = N_t y_t(J_t) \left[\frac{1}{1 - J_t} \int_{J_t}^1 \left(\frac{A(j)}{\underline{A}_t} \right)^{\frac{1}{\mu - \nu}} dJ \right]^{\mu} = N_t y_t(J_t) \left(\frac{\hat{A}_t}{\underline{A}_t} \right)^{\frac{\mu}{\mu - \nu}}$$
(68)

This shows that aggregate output depends on the number of active firms, the size of the threshold firm and the ratio of average technology to threshold technology.³³ Substituting in $y_t(J_t) = \underline{A}_t \left[k_t(J_t)^{\alpha} \ell_t(J_t)^{1-\alpha} \right]^{\nu}$ yields:

$$Y_t = N_t \hat{A}_t^{\frac{\mu}{\mu-\nu}} \underline{A}_t^{\frac{-\nu}{\mu-\nu}} \left[k_t (J_t)^{\alpha} \ell_t (J_t)^{1-\alpha} \right]^{\nu}$$
(69)

The next step again applies the technique of representing firm-level variables relative to the threshold-firm. This allows us to replace $k_t(J_t)$ and $\ell_t(J_t)$ in terms of aggregates.

$$K_t = E_t \int_{J_t}^1 k_t(j) \, dj = \frac{N_t}{1 - J_t} \int_{J_t}^1 k_t(j) \, dj = \frac{N_t k_t(J_t)}{1 - J_t} \int_{J_t}^1 \frac{k_t(j)}{k_t(J_t)} \, dj$$
(70)

$$= \frac{N_t k_t(J_t)}{1 - J_t} \int_{J_t}^1 \left(\frac{A_t(j)}{\underline{A}_t}\right)^{\frac{1}{\mu - \nu}} dj = N_t k_t(J_t) \left(\frac{\hat{A}_t}{\underline{A}_t}\right)^{\frac{1}{\mu - \nu}}$$
(71)

$$L_{t} = E_{t} \int_{J_{t}}^{1} \ell_{t}(j) + \phi \, dj = \frac{N_{t}}{1 - J_{t}} \int_{J_{t}}^{1} \ell_{t}(j) + \phi \, dj = \frac{N_{t}\ell_{t}(J_{t})}{1 - J_{t}} \int_{J_{t}}^{1} \frac{\ell_{t}(j)}{\ell_{t}(J_{t})} + \frac{\phi}{\ell_{t}(J_{t})} \, dj \tag{72}$$

$$=\frac{N_t\ell_t(J_t)}{1-J_t}\int_{J_t}^1 \left(\frac{A_t(j)}{\underline{A}_t}\right)^{\frac{1}{\mu-\nu}} + \frac{\phi}{\ell_t(J_t)}\,dj = N_t\ell_t(J_t)\left(\frac{\hat{A}_t}{\underline{A}_t}\right)^{\frac{1}{\mu-\nu}} + N_t\phi \tag{73}$$

³³Gao and Kehrig (2021) present an analogous result for the partial equilibrium case with perfect competition ($\mu = 1$) and no external returns to scale (love of variety).

Therefore we can express threshold firm capital and labour as

$$k(J_t) = \left(\frac{\underline{A}_t}{\hat{A}_t}\right)^{\frac{1}{\mu-\nu}} \frac{K_t}{N_t}$$
(74)

$$\ell(J_t) = \left(\frac{\underline{A}_t}{\hat{A}_t}\right)^{\frac{1}{\mu-\nu}} \frac{u_t L_t}{N_t}, \quad \text{where } u_t \equiv 1 - \frac{N_t \phi}{L_t}.$$
(75)

Finally, substituting these two expressions into our reduced-form expression for output yields:

$$Y_{t} = N_{t}^{1-\nu} \hat{A}_{t} \left[K_{t}^{\alpha} \left(u_{t} L_{t} \right)^{1-\alpha} \right]^{\nu}.$$
(76)

E Fixed Cost Share Data

We use the administration expenses share in turnover to proxy the fixed cost share for UK firms. Figure 4 shows the median administration expenses share in turnover for UK firms from 2004 to 2023.

Administrative Expenses

In UK company accounts, 'Administrative Expenses' are defined as expenses an organization incurs that are not directly related to a specific function such as manufacturing, production, or sales. These expenses can include things like: rent, utilities, insurance, wages and benefits for administrative staff, depreciation on office furniture and equipment, professional fees (e.g., accounting and legal fees), and travel expenses. They are necessary for the day-to-day operation of a business, but they do not directly contribute to the generation of revenue. Expenses related to the generation of revenue fall under cost of goods sold (COGs). Administration expenses are typically reported on a company's income statement, below the cost of goods sold (COGS) line.

FAME data

We use the Bureau van Dijk FAME dataset, a UK version of Orbis, to obtain firm financial information. The dataset records the annual financial statements of all incorporated companies in the UK. Over the entire period, there are 16,426,460 company entries. We restrict our analysis to companies that have at least one entry in administration expenses for any year between 2004 and 2023. The company does not need to be active today; it could have dissolved. This restriction reduces the number of companies to 680,763. The companies removed in this step have no administration expenses recorded over the sample period. This occurs because smaller companies can submit micro-entity accounts which do not include this information. Medium and large companies submit 'full accounts' which do record this information. Due to download restrictions, we take a random sample of 250,000 companies, and we keep this same sample of firms every year. Since a firm only needs to have an administration expense in one year, there will be many blanks in any given year for any given company, either because it is inactive or because administration expenses were not recorded because it is a micro-entity. In the end, there are approximately 50,000 firms each year that have an entry in both administration expenses and turnover.

F ARD Data

We use the Annual Respondents Database (ARD) or the time-series version known as ARDx. The ARD is based on the Annual Business Survey (ABS). The ABS is an annual survey of firms in the UK economy. It is a core ONS product used in the construction of national accounts. The ARD adds information from other business surveys to the ABS data.³⁴ Firms are legally obligated to respond to the survey. The survey forms a firm-level panel that covers all large firms and a representative sample of small firms by geography, size and sector. Large firms are surveyed annually, while small firms are surveyed for a fixed number of years. The ARDx Methodology and ABS Methodology

³⁴Specifically, the ARD brings together the ABS and the Business Register and Employment Survey (BRES), and prior to 2009 it brought together the two parts of the Annual Business Inquiry (ABI).

provide more detail.

F.1 Capital Construction

The Perpetual Inventory Method (PIM) allows the construction of firm-level capital stocks when such data are unavailable, but investment data is present. The method here follows Martin (2002) and Hwang, Savagar, and Kariel (2022). The PIM is constructed using the following equation:

$$K_t = (1 - \delta)K_{t-1} + I_t.$$

 K_t is the capital stock in period t, and I_t is investment in period t. However, to use this method, we need K_0 – the initial capital stock of a company, which is not in this survey. To construct this series, each firm's K_0 is a revenue-weighted share of the industry-level capital stock in the first year that firm appears in the panel. The capital stock is then constructed for all future years with the above equation, with the missing investment data interpolated. The depreciation rate is taken to be 18.195%, which is a weighted average of the ONS depreciation rates for the three different capital categories: Building, Vehicles, Other.

F.2 Deflating

We convert firm gross output and value added into real values using the ONS industry deflators. Material inputs are deflated with the ONS producer price inflation data. The capital stock is deflated with the ONS gross fixed capital formation deflator.

F.3 Cleaning

For the purpose of our production function estimation, we exclude sectors: Agriculture, Public Sector, Finance & Insurance, Education, and Health. Standard Industrial Classification (SIC) 2007 codes: A, K, O, P, Q. These sectors were excluded from the survey after 2012. K,O,P were fully excluded and A,Q had various subsectors excluded. We set out rules for SIC re-coding to ensure compatibility pre- and post-2007, when the classification is changed. For SIC codes post-2007, we divide the number by 1000 to match with pre-2007 codes. To avoid outliers, which may represent recording errors in the surveys, we winsorize firms with the top and bottom 0.1% of factor shares in revenue (M/Y, K/Y, L/Y) in each year. Table 2 contains number of firms at each stage of the data cleaning process, along with the final number of observations for estimation.

	# Firms	
All ARD firm-year obs	854,732	
Drop if no 2-digit sector	852,424	
Drop if < 100 firms in sector	852,331	
Drop sectors A,K,O,P,Q	761,348	
Take logs of regression variables	539,368	
Drop outlier factor shares	527,813	

Table 2: Data Cleaning: Firms Dropped

F.4 Summary Statistics

Table 3 presents aggregate descriptive statistics of the variables used in our regression analysis.

	Mean	SD	p10	p50	p90	No. Obs
Revenue	39,736	675,183	92	1,458	42,797	527,813
Labour	224	2,213	2	20	349	527,813
Capital	7,696	150,007	22	351	7,915	527,813
Materials	29,651	636,176	32	703	26,255	527,813
Materials Share	0.55	-	0.17	0.58	0.87	527,813
Labour Share	0.26	-	0.04	0.23	0.52	527,813
Capital Share	0.27	-	0.06	0.19	0.60	527,813

Table 3: Descriptive Statistics of Regression Variables for Full Sample

Table 4 presents descriptive statistics by broad industry group.

	Mean	SD	p10	p50	p90	No. Obs
Manufacturing						
Revenue	36,005	235,437	336	4,294	58,896	125,737
Labour	192	576	8	54	431	125,737
Capital	10,362	75,776	148	1,498	16,154	125,737
Materials	24,954	178,528	122	2,400	38,999	125,737
Materials Share	0.57	-	0.30	0.58	0.81	125,737
Labour Share	0.28	-	0.11	0.27	0.47	125,737
Construction						
Revenue	17,812	108,789	111	1,414	48,782	51,784
Labour	103	395	2	11	214	51,784
Capital	2,309	41,523	11	104	2,210	51,784
Materials	12,467	89,027	18	343	16,896	51,784
Materials Share	0.51	-	0.17	0.52	0.81	51,784
Labour Share	0.25	-	0.00	0.24	0.49	51,784
Trade, Wholesale, Transport						
Revenue	62,673	1,102,305	111	1,414	48,782	182,814
Labour	256	3,404	2	14	244	182,814
Capital	7,092	103,075	20	245	5,667	182,814
Materials	52,666	1,044,112	61	929	26,219	182,814
Materials Share	0.69	-	0.37	0.74	0.92	182,814
Labour Share	0.16	-	0.02	0.13	0.35	182,814
Services						
Revenue	25,276	284,335	65	728	28,673	179,028
Labour	249	1,627	2	17	403	179,028
Capital	8,821	228,905	20	218	5,435	179,028
Materials	14,417	209,297	15	242	11,263	179,028
Materials Share	0.41	-	0.09	0.38	0.77	179,028
Labour Share	0.34	-	0.06	0.32	0.68	179,028

 Table 4: Descriptive Statistics of Regression Variables by Broad Sector