

How Costly Are Markups?*

Chris Edmond[†] Virgiliu Midrigan[‡] Daniel Yi Xu[§]

First draft: July 2018. This draft: September 2022

Abstract

We study the welfare costs of markups in a dynamic model with heterogeneous firms and endogenously variable markups. Our general framework encompasses a range of popular market structures. We provide aggregation results showing how the macro implications of micro-level markup heterogeneity can be summarized by a few key statistics. We calibrate our model to match US Census of Manufactures firm-level data and find that the welfare costs of markups can be large. We decompose the costs of markups into three channels: (i) an aggregate markup that acts like a uniform output tax, (ii) misallocation of factors of production, and (iii) inefficient entry. Across all specifications, we find that the aggregate markup and misallocation channels account for the bulk of the costs of markups and that the entry channel is much less important. Subsidizing entry is not an effective tool in our model. While an increase in competition reduces incumbents' markups, it also reallocates market shares towards larger incumbent firms and the net effect is that the aggregate markup changes little.

Keywords: competition, concentration, misallocation, firm dynamics.

JEL classifications: D4, E2, L1, O4.

*We thank Greg Kaplan and five anonymous referees for valuable comments and suggestions. We are also particularly grateful for insightful feedback from our discussants Salomé Baslandze, Ariel Burstein, Jan Eeckhout, and Gino Gancia, and from Gauti Eggertsson, Emmanuel Farhi, Oleg Itskhoki, Pete Klenow and Iván Werning. We also thank participants at the Fall 2018 NBER EFG meeting, 2018 Bank of Italy-CEPR-EIEF conference on firm dynamics and economic growth and seminar participants at Columbia University, Duke University, ETH Zurich, the FRB New York, FRB St Louis, the Graduate Institute of Geneva, Keio University, the LSE, MIT, Northwestern, Notre Dame, NYU, Rochester University, UC Berkeley, USC, University of Adelaide, University of Chicago, University of Michigan, and University of Queensland for their comments. Edmond thanks the Australian Research Council for financial support under grant DP-150101857.

[†]University of Melbourne, cedmond@unimelb.edu.au.

[‡]New York University and NBER, virgiliu.midrigan@nyu.edu.

[§]Duke University and NBER, daniel.xu@duke.edu.

1 Introduction

How large are the welfare costs of product market distortions? What kinds of policies can best overcome these distortions? We answer these questions using a dynamic model with heterogeneous firms and endogenously variable markups. In our model, markups distort allocations through three channels. First, the *aggregate markup* acts like a uniform tax on all firms. Second, there is cross-sectional markup dispersion because larger firms face less competition and so charge higher markups. This markup dispersion gives rise to *misallocation* of factors of production. Third, there is *inefficient entry*. Our goal is to quantify these three channels using US micro data and to evaluate policies aimed at reducing the costs of markups.

Our focus in this paper is normative: we quantify the welfare costs of markups through the lens of a dynamic model. But the specific endogenous markup mechanism we study is consistent with key facts stressed in the recent empirical literature. In our model, within a given sector, more productive firms are, in equilibrium, larger and face less elastic demand and so charge higher markups than less productive firms. Shocks that allow more productive firms to grow at the expense of less productive firms will be associated with an increase in the aggregate markup and a decline in the aggregate labor share. In this sense, our model is consistent with the reallocation of production from firms with relatively high measured labor shares to firms with relatively low measured labor shares (Autor, Dorn, Katz, Patterson and Van Reenen, 2020; Kehrig and Vincent, 2021) and the observation that firms with high markups have been getting larger, driving up the aggregate markup (Baqaee and Farhi, 2020).

Our general framework encompasses a range of popular market structures including (i) *monopolistic competition* with Kimball (1995) demand or symmetric translog demand as in Feenstra (2003), and (ii) *oligopolistic competition* with nested-CES demand as in Atkeson and Burstein (2008) and Edmond, Midrigan and Xu (2015). We consider settings where firms can differ in both productivity and quality and provide aggregation results showing that the macro implications of micro-level markup heterogeneity can be summarized by a few key statistics. One such result is that the aggregate markup, the ‘wedge’ in aggregate employment and investment decisions, is given by the *cost-weighted* average of firm-level markups.¹ By contrast, the empirical literature on the macro implications of markup heterogeneity typically reports the *sales-weighted* average of firm-level markups.² We show that the sales-weighted average is the cost-weighted average plus a term reflecting the *variance* of markups. In this sense the sales-weighted average overstates the aggregate markup by including a term that reflects misallocation rather than the level of markups per se. Importantly, these aggregation results hold independent of the market structure details.

¹Or the sales-weighted *harmonic* average, as in Edmond, Midrigan and Xu (2015) and Grassi (2017).

²For example, for the US economy De Loecker, Eeckhout and Unger (2020) estimate a sharply increasing sales-weighted average markup rising from about 1.2 in 1980 to about 1.6 in 2016. By contrast the cost-weighted average is lower and has risen by less, from about 1.1 to about 1.25. The difference reflects the increase in cross-sectional markup dispersion. We discuss these measures at length in Appendix A.

Regardless of market structure, we find that markups distort allocations through the three channels mentioned at the outset: the aggregate markup, misallocation due to markup dispersion, and inefficient entry. We show that the efficient allocation can be implemented by a specific nonlinear schedule of direct subsidies with two components, a *uniform component* that subsidizes all firms and that can be used to eliminate the aggregate markup, and a *size-dependent* component that jointly eliminates misallocation and the entry distortion.

We quantify the welfare costs of markups by asking how much the representative consumer would benefit if the economy transitioned from an initial steady state with markup distortions to the efficient steady state. Because eliminating the markup distortions entails a large increase in the capital stock, taking into account the cost of building up the capital stock is critical to correctly assess the welfare gains from such policies. We calibrate the initial steady state using US Census of Manufactures firm-level data from 1972 to 2012 to match levels of sales concentration and the firm-level relationship between markups and market shares observed in 6-digit NAICS sectors, controlling for firm fixed effects and 6-digit NAICS sector-year effects to control for other persistent sources of firm and sector heterogeneity.

Our calibration strategy makes use of the fact that, though the precise mapping depends on market structure, all versions of our model imply a simple firm-level relationship between markups and market shares. We use the estimated parameter values from this relationship to calculate firm-level markups in the model and calculate the welfare costs of these markups. That is, we do not feed into the model separately estimated firm-level markups. We prefer to use the markups implied by our model for two reasons. First, in the Census data we only observe firm-level revenues, not prices and quantities separately. Absent firm-level quantities we cannot disentangle markup *levels* from output elasticities in production (see [Bond, Hashemi, Kaplan and Zoch, 2021](#); [De Ridder, Grassi and Morzenti, 2022](#), for extensive discussion). Second, we would be cautious to interpret such estimates as ‘true markups’ even if output elasticities were accurately estimated, since such estimates potentially confound the true markup with other distortionary ‘wedges’ — e.g., implicit or explicit input or revenue taxes, factor-adjustment costs, or price rigidities, etc. For these reasons the Census data we use leads a relatively wide range of empirically plausible markup levels. Given this, we report the welfare costs of markups for a wide range of values for the aggregate markup, recalibrating the model each time.

We find that the welfare costs of markups can be large. It turns out that the welfare costs are not just increasing in the level of the aggregate markup we target, they are increasing and *convex*. Because of this convexity, for some parameterizations of the model we find very large welfare costs of markups, as high as 50% in consumption-equivalent terms. Overall we also find that the costs tend to be lower if we assume monopolistic competition but are much higher if we assume oligopolistic competition.

We then turn to quantifying the relative importance of the three channels by which

markups reduce welfare in our model. Across all specifications, we find that the aggregate markup and misallocation channels account for the bulk of the costs of markups and that the entry channel is much less important. That said, the relative importance of the aggregate markup and misallocation channels vary depending on the market structure and target for the aggregate markup. For example, the Kimball specification implies that the *share* of the total costs accounted for by the aggregate markup increases from $1/2$ to $3/4$ as we increase the aggregate markup from 1.05 to 1.35. The balance of the costs are almost entirely due to misallocation, the losses from the entry distortion are negligible.

Although the losses from misallocation in our model can be sizeable, accounting for value-added TFP losses of around 2% to 6%, depending on the specification, they are small relative to standard estimates in the literature (Restuccia and Rogerson, 2008; Hsieh and Klenow, 2009). This is because we measure misallocation using the dispersion in marginal revenue products implied by the endogenous markup distribution in our model, i.e., that relatively small share of the dispersion in marginal revenue products systematically related to market shares. We *do not* attribute all variation in observed marginal revenue products to markups.

In representative firm models, subsidizing entry (or reducing barriers to entry) so as to increase competition is a powerful tool for reducing the aggregate markup and hence reducing the costs of markups (Bilbiie, Ghironi and Melitz, 2008, 2019). By contrast we find that, with heterogeneous firms, subsidizing entry is *not* a powerful tool. For all our specifications, we find that even large increases in the number of firms have small effects on the aggregate markup.³ To understand this, recall that the aggregate markup is a cost-weighted average of firm-level markups. An increase in the number of firms has two effects on this weighted average. The direct effect is a reduction in the markup of each firm, due to a reduction in each firm’s market share. But there is also an important compositional effect: small firms face more elastic demand and are more vulnerable to competition from entrants; large firms face less elastic demand and are less vulnerable. So when there is an increase in the number of firms, small, low markup firms contract by more than large, high markup firms and the resulting reallocation keeps the aggregate markup almost unchanged, despite the reduction in firm-level markups. In all our specifications, this offsetting compositional effect is almost as large as the direct effect so overall the aggregate markup falls by a small amount.⁴

The different specifications we consider each have their own strengths and weaknesses. The model with Kimball demand is more flexible than the model with symmetric translog demand and is better able to match our calibration targets. But the model with translog

³There are however standard love-of-variety gains from increasing the number of firms.

⁴These offsetting direct and compositional effects are reminiscent of results in the trade literature, e.g., Bernard, Eaton, Jensen and Kortum (2003) and especially Arkolakis, Costinot, Donaldson and Rodríguez-Clare (2019). We derive analogous results for Kimball and translog demand but unlike in their analysis, we do not assume from the outset that the ‘choke price’ in either demand system is binding, since this is an equilibrium outcome. For the translog case, we also provide closed-form solutions for the aggregate markup and the cutoff productivity that pins down the cross-sectional distributions of markups and market shares.

demand is more tractable than Kimball demand and leads to sharp analytic results. Both monopolistic competition models are simple computationally. The oligopoly model is computationally challenging but has richer empirical content. Though our aggregation results hold regardless of the assumed market structure, the oligopoly model makes a number of predictions that differ from the monopolistic competition models. First, we find larger amounts of markup dispersion and hence larger losses from misallocation in the oligopoly model than in either of the monopolistic competition models. Second, while the monopolistic competition models predict that there are *too few* firms in equilibrium, the oligopoly model predicts that there are *too many*. But since the entry margin is not a quantitatively important source of losses in any specification, this qualitative difference is not important.

Existing results on costs of markups. The starting point for discussion of the welfare costs of markups is [Dixit and Stiglitz \(1977\)](#), though the literature goes back to [Lerner \(1934\)](#). Recent work such as [Zhelobodko, Kokovin, Parenti and Thisse \(2012\)](#), [Dhingra and Morrow \(2019\)](#) and [Behrens, Mion, Murata and Suedekum \(2020\)](#) studies variable markups in static models with heterogeneous firms. By contrast, our model is dynamic. Like us, [Bilbiie, Ghironi and Melitz \(2008, 2019\)](#) study a dynamic model and quantify the costs of markups but they assume a representative firm. We find, however, that firm heterogeneity plays a crucial role in understanding the costs of markups. In our model, markups compensate firms for sunk investments in the creation of a new variety. To the extent that there are positive spillovers from the creation of new varieties, as in the endogenous growth literature, our results may overstate the costs of markups. [Atkeson and Burstein \(2010, 2019\)](#) provide a welfare analysis of innovation policies in firm dynamics models but abstract from variable markups. [Peters \(2020\)](#) studies innovation, firm dynamics, and variable markups but does not evaluate the welfare costs of markups.

Markups and misallocation. In our model markups increase with firm size. This is one form of misallocation in the sense of [Restuccia and Rogerson \(2008\)](#), and [Hsieh and Klenow \(2009\)](#). We find that the gross output productivity losses from this form of misallocation are on the order of 1 to 3%, with value-added productivity losses about double that, on the order of 2 to 6%, reflecting a materials share in gross output just under one-half. We view these numbers as an upper bound on the the gains from size-dependent subsidies since we attribute all of the systematic relationship between firm revenue productivity and firm size to market power, and not to, say, overhead costs as in [Autor, Dorn, Katz, Patterson and Van Reenen \(2020\)](#) and [Bartelsman, Haltiwanger and Scarpetta \(2013\)](#). Because of this we are likely somewhat overstating the true relationship between markups and firm size and overstating the losses from this form of misallocation.

It is important to recognize that we abstract from all other sources of markup variation that may cause misallocation. Firms may operate in different locations or sell different prod-

ucts in different sectors and charge different markups depending on the amount of competition they face in those different markets.⁵ Policies that condition on location or other relevant market details may be able to address these forms of misallocation too. But implementing finely-tuned policies that condition on details of market conditions location-by-location seems challenging in practice. Given this, we restrict our attention to size-dependent markup variation and we find that the value-added productivity gains from eliminating misallocation due to size-dependent markup variation are likely no more than 2 to 6%.

In related work, [Baqaee and Farhi \(2020\)](#) calculate that the value-added aggregate productivity gains from eliminating markups are about 20%, much larger than in our model. They find much larger effects because they feed into their calculation all the variation in estimated markups (as in [De Loecker, Eeckhout and Unger, 2020](#); [Gutiérrez and Phillippon, 2017b](#)) whereas we feed in that component of markups that systematically varies with firm market shares. Because the estimated markups they use are more dispersed than the markups from our model, they find larger effects of markup dispersion on aggregate productivity.

2 Model

There is a representative consumer with preferences over final consumption and labor supply and who owns all the firms. The final good is produced by perfectly competitive firms using inputs from many sectors. Within each sector there are heterogeneous imperfectly competitive firms producing differentiated products using capital, labor and materials. Firms enter by paying a sunk cost in units of labor and then obtain a one-time productivity draw in a randomly allocated sector. Exit is random and there is no aggregate uncertainty. We focus on characterizing the steady state and transitional dynamics after a policy change.

2.1 Setup

A key feature of our analysis is a set of aggregation results that hold regardless of the details of market structure within each sector. We proceed in two steps, first explaining the basic setup and aggregate outcomes that hold independent of market structure within each sector and then turning to the remaining details where market structure matters.

Representative consumer. The representative consumer maximizes

$$\sum_{t=0}^{\infty} \beta^t \left(\log C_t - \psi \frac{L_t^{1+\nu}}{1+\nu} \right) \tag{1}$$

⁵[Rossi-Hansberg, Sarte and Trachter \(2020\)](#) show that while aggregate US product-market concentration has been rising since the early 1990s, concentration in geographically-specific local markets has been falling.

subject to the budget constraint

$$C_t + I_t = W_t L_t + R_t K_t + \Pi_t \quad (2)$$

where C_t denotes consumption of the numeraire final good, $I_t = K_{t+1} - (1 - \delta)K_t$ denotes investment, K_t denotes physical capital, L_t denotes labor supply, W_t the real wage, R_t the rental rate of capital, and Π_t denotes aggregate profits net of the cost of creating new firms.

The representative consumer's labor supply satisfies

$$\psi C_t L_t^\nu = W_t \quad (3)$$

and their investment choice satisfies

$$1 = \beta \frac{C_t}{C_{t+1}} (R_{t+1} + 1 - \delta) \quad (4)$$

Since firms are owned by the representative consumer, they use the one-period discount factor $\beta C_t / C_{t+1}$ to discount future profit flows.

Final good producers. Let Y_t denote *gross output* of the final good. This can be used for consumption C_t , investment I_t , or as materials X_t , so that

$$C_t + I_t + X_t = Y_t \quad (5)$$

The use of the final good as materials gives the model a simple ‘roundabout’ production structure, as in [Jones \(2011\)](#) and [Baqaee and Farhi \(2020\)](#).

The final good Y_t is produced by perfectly competitive firms using inputs $y_t(s)$ from a continuum of sectors

$$Y_t = \left(\int_0^1 y_t(s)^{\frac{\eta-1}{\eta}} ds \right)^{\frac{\eta}{\eta-1}} \quad (6)$$

where $\eta > 1$ is the elasticity of substitution *across* sectors $s \in [0, 1]$. Let $p_t(s)$ denote the price index for sector s . Since the final good is the numeraire, these satisfy

$$1 = \left(\int_0^1 p_t(s)^{1-\eta} ds \right)^{\frac{1}{1-\eta}} \quad (7)$$

Within sectors. Within each sector there are imperfectly competitive firms producing differentiated goods. As discussed extensively below, we consider two market structures: *monopolistic competition* with a continuum of firms $i \in [0, n_t(s)]$ per sector, or *oligopolistic competition* with a finite number of firms $i = 1, \dots, n_t(s)$ per sector. Except where noted, our results below hold for both cases.

Technology. Firms enter by paying a sunk cost κ in units of labor and then obtain a one-time productivity draw $z_i(s) \sim G(z)$ in a random sector s . A firm's *gross output* is then

$$y_{it}(s) = z_i(s) \left(\phi^{\frac{1}{\theta}} v_{it}(s)^{\frac{\theta-1}{\theta}} + (1-\phi)^{\frac{1}{\theta}} x_{it}(s)^{\frac{\theta-1}{\theta}} \right)^{\frac{\theta}{\theta-1}} \quad (8)$$

where $v_{it}(s)$ is the firm's *value-added*, a composite of physical capital and labor

$$v_{it} = k_{it}(s)^\alpha l_{it}(s)^{1-\alpha} \quad (9)$$

We impose a unit elasticity of substitution between capital and labor. The elasticity of substitution between value-added $v_{it}(s)$ and materials $x_{it}(s)$ is given by θ .

Input demands. Taking input prices as given, cost minimization gives the input demands

$$R_t k_{it}(s) = \alpha \left\{ \left(\frac{R_t}{\alpha} \right)^\alpha \left(\frac{W_t}{1-\alpha} \right)^{1-\alpha} \right\} v_{it}(s) \quad (10)$$

$$W_t l_{it}(s) = (1-\alpha) \left\{ \left(\frac{R_t}{\alpha} \right)^\alpha \left(\frac{W_t}{1-\alpha} \right)^{1-\alpha} \right\} v_{it}(s) \quad (11)$$

where the term in braces on the right is the price index for the value-added composite. In turn, demand for the value-added composite and demand for materials are given by

$$v_{it}(s) = \phi \left\{ \frac{\left(\frac{R_t}{\alpha} \right)^\alpha \left(\frac{W_t}{1-\alpha} \right)^{1-\alpha}}{\Omega_t} \right\}^{-\theta} \frac{y_{it}(s)}{z_i(s)} \quad (12)$$

and

$$x_{it}(s) = (1-\phi) \left\{ \frac{1}{\Omega_t} \right\}^{-\theta} \frac{y_{it}(s)}{z_i(s)} \quad (13)$$

where Ω_t is the input price index dual to the technologies in (8) and (9), namely

$$\Omega_t = \left(\phi \left\{ \left(\frac{R_t}{\alpha} \right)^\alpha \left(\frac{W_t}{1-\alpha} \right)^{1-\alpha} \right\}^{1-\theta} + (1-\phi) \right)^{\frac{1}{1-\theta}} \quad (14)$$

where materials have a relative price of 1 since they are in units of the numeraire. Notice that the capital/labor and value-added/materials ratios are common to all firms.

Marginal cost. These factor demands imply that a firm's *marginal cost* is given by

$$\frac{\Omega_t}{z_i(s)} \quad (15)$$

Profits and markups. A firm's profits are then given by

$$\pi_{it}(s) = p_{it}(s)y_{it}(s) - \frac{\Omega_t}{z_i(s)} y_{it}(s) \quad (16)$$

Firms maximize profits subject to the demand system they face, which depends on the market structure details. At the optimum a firm's price can be written as a markup $\mu_{it}(s)$ over marginal cost

$$p_{it}(s) = \mu_{it}(s) \frac{\Omega_t}{z_i(s)}, \quad \mu_{it}(s) = \frac{\sigma_{it}(s)}{\sigma_{it}(s) - 1} \quad (17)$$

where $\sigma_{it}(s)$ denotes the (endogenous) *demand elasticity* facing firm i . Different demand systems imply different determinants of $\sigma_{it}(s)$ as discussed below. Profits can then be written in terms of markups and sales

$$\pi_{it}(s) = \left(1 - \frac{1}{\mu_{it}(s)}\right) p_{it}(s)y_{it}(s) \quad (18)$$

Labor shares. Combining a firm's labor demand from (11)-(12) with markup pricing (17), a firm's labor share can be written

$$\frac{W_t l_{it}(s)}{p_{it}(s)y_{it}(s)} = \frac{(1 - \alpha)\zeta_t}{\mu_{it}(s)} \quad (19)$$

where ζ_t denotes the elasticity of output with respect to value-added

$$\zeta_t := \frac{\frac{\phi}{1-\phi} \left\{ \left(\frac{R_t}{\alpha}\right)^\alpha \left(\frac{W_t}{1-\alpha}\right)^{1-\alpha} \right\}^{1-\theta}}{1 + \frac{\phi}{1-\phi} \left\{ \left(\frac{R_t}{\alpha}\right)^\alpha \left(\frac{W_t}{1-\alpha}\right)^{1-\alpha} \right\}^{1-\theta}} \quad (20)$$

This elasticity is common to all firms but in general varies over time. All cross-sectional variation in labor shares is due to cross-sectional variation in markups $\mu_{it}(s)$.

We next briefly outline how the distribution of markups $\mu_{it}(s)$ affects productivity within and across sectors. We focus on aggregation results that obtain independent of within-sector market structure.

Aggregate productivity. Let $k_t(s)$, $l_t(s)$, and $x_t(s)$ denote sector-level capital, labor and materials. These are the integrals (or sums) of $k_{it}(s)$, $l_{it}(s)$ and $x_{it}(s)$ over i within s . We can then write the gross output of sector s as

$$y_t(s) = z_t(s)F(k_t(s), l_t(s), x_t(s)) \quad (21)$$

where

$$F(k, l, x) = \left(\phi^{\frac{1}{\theta}} \left(k^\alpha l^{1-\alpha} \right)^{\frac{\theta-1}{\theta}} + (1 - \phi)^{\frac{1}{\theta}} x^{\frac{\theta-1}{\theta}} \right)^{\frac{\theta}{\theta-1}} \quad (22)$$

and where sector-level productivity satisfies

$$z_t(s) = \left(\int_0^{n_t(s)} \frac{q_{it}(s)}{z_i(s)} di \right)^{-1} \quad (23)$$

where $q_{it}(s) := y_{it}(s)/y_t(s)$ denotes the relative size of firm i in sector s . The only difference having a finite number of firms makes is that the integral should be replaced by a finite sum.

Likewise, let K_t , \tilde{L}_t , and X_t denote aggregate capital, labor *used in production*, and materials. These are the integrals of $k_t(s)$, $l_t(s)$ and $x_t(s)$ over $s \in [0, 1]$. We then have aggregate gross output $Y_t = Z_t F(K_t, \tilde{L}_t, X_t)$ where aggregate productivity is given in the same way as sector productivity

$$Z_t = \left(\int_0^1 \frac{q_t(s)}{z_t(s)} ds \right)^{-1} \quad (24)$$

where $q_t(s) := y_t(s)/Y_t$ denotes the relative size of sector s .

Thus sector-level productivity $z_t(s)$ is a firm-size-weighted harmonic average of firm-level productivity $z_i(s)$ and aggregate productivity Z_t is a sector-size-weighted harmonic average of sector-level productivity. Sector-level productivity and aggregate productivity are affected by markups $\mu_{it}(s)$ through the effects of markups on the distribution of firm-size $q_{it}(s)$ within sectors and the distribution of sector-size $q_t(s)$ across sectors.

Aggregate markup. Let $\mu_t(s)$ denote the sector-level markup, implicitly defined by the sector-level labor share

$$\frac{W_t l_t(s)}{p_t(s) y_t(s)} = \frac{(1 - \alpha) \zeta_t}{\mu_t(s)} \quad (25)$$

Combining the sector-level labor share with its firm-level counterpart (19) we can write the sales-share of firm i in sector s as

$$\frac{p_{it}(s) y_{it}(s)}{p_t(s) y_t(s)} = \frac{\mu_{it}(s)}{\mu_t(s)} \times \frac{l_{it}(s)}{l_t(s)} \quad (26)$$

Integrating both sides, the sector-level markup can be written *either* as an employment-weighted arithmetic average or a sales-weighted harmonic average of firm-level markups, as in [Edmond, Midrigan and Xu \(2015\)](#),

$$\mu_t(s) = \int_0^{n_t(s)} \mu_{it}(s) \frac{l_{it}(s)}{l_t(s)} di = \left(\int_0^{n_t(s)} \frac{1}{\mu_{it}(s)} \frac{p_{it}(s) y_{it}(s)}{p_t(s) y_t(s)} di \right)^{-1} \quad (27)$$

where, again, the only difference having a finite number of firms makes is that the integral should be replaced by a finite sum. From either of these and the expression for sector-level

productivity $z_t(s)$ we see that the sector-level markup satisfies $p_t(s) = \mu_t(s)\Omega_t/z_t(s)$, i.e., the sector price index can be expressed as the sector-level markup over marginal cost.

Likewise, let \mathcal{M}_t denote the aggregate, economy-wide markup. Following the same steps, this can be written either as an employment-weighted arithmetic average or a sales-weighted harmonic average of sector-level markups

$$\mathcal{M}_t = \int_0^1 \mu_t(s) \frac{l_t(s)}{\tilde{L}_t} ds = \left(\int_0^1 \frac{1}{\mu_t(s)} \frac{p_t(s)y_t(s)}{Y_t} ds \right)^{-1} \quad (28)$$

The aggregate markup satisfies $1 = \mathcal{M}_t\Omega_t/Z_t$, i.e., the aggregate price level (normalized to one) is the aggregate markup over aggregate marginal cost. We discuss these and related measures of average markups in more detail in [Appendix A](#).

Markup dispersion and productivity. To see how markup dispersion affects productivity, observe from (6) that sector size $q_t(s) = y_t(s)/Y_t$ satisfies $q_t(s) = p_t(s)^{-\eta}$ and since $p_t(s) = \mu_t(s)\Omega_t/z_t(s)$ and $1 = \mathcal{M}_t\Omega_t/Z_t$ we can write

$$q_t(s) = \left(\frac{\mu_t(s)}{\mathcal{M}_t} \frac{Z_t}{z_t(s)} \right)^{-\eta} \quad (29)$$

Plugging this into our expressions for aggregate productivity and solving for Z_t we obtain

$$Z_t = \left(\int_0^1 \left(\frac{\mu_t(s)}{\mathcal{M}_t} \right)^{-\eta} z_t(s)^{\eta-1} ds \right)^{\frac{1}{\eta-1}} \quad (30)$$

In turn, sector-productivity $z_t(s)$ and markups $\mu_t(s)$ depend on the distribution of firm-level productivity $z_i(s)$ and markups $\mu_{it}(s)$ within sector s — but the details of this layer of aggregation *do* depend on the within-sector market structure.

2.2 Role of market structure

In this section we explain how the details of within-sector market structure matter. First, the market structure matters for determining the relative size distribution $q_{it}(s) = y_{it}(s)/y_t(s)$ within each sector s . That said, *taking* $n_t(s)$ *as given*, we can cover a range of popular specifications in a unified way, as explained below. Second, and more substantively, the market structure matters for the entry problem that determines $n_t(s)$. The entry problem is simple with monopolistic competition but more involved with oligopolistic competition.⁶

⁶In our model with oligopoly, the *potential* number of firms per sector $n_t(s)$ is endogenous. This problem is challenging because potential entrants anticipate their impact on a sector, and the distribution of sectoral configurations is a very high-dimensional object. By contrast in [Atkeson and Burstein \(2008\)](#), [Edmond, Midrigan and Xu \(2015\)](#), and [De Loecker, Eeckhout and Mongey \(2021\)](#), the potential number of firms is static and exogenous, with firms simply deciding whether to operate or not.

Relative size distribution. Taking $n_t(s)$ as given, the relative size distribution $q_{it}(s)$ within sector s is pinned down by the static markup-pricing condition (17). To cover alternative specifications in a unified way, we write this as

$$f(q) = \frac{\sigma(q)}{\sigma(q) - 1} \frac{A_t(s)}{z_i(s)}, \quad A_t(s) := \frac{\Omega_t}{p_t(s)d_t(s)} \quad (31)$$

where the function $f(q)$ is proportional to the inverse demand curve, $\sigma(q)$ is the associated demand elasticity, with markup $\mu(q) = \sigma(q)/(\sigma(q) - 1)$, and where $p_t(s)$ is the price index for sector s and $d_t(s)$ is a demand index that depends on the market structure.⁷ Let $q(z; A)$ denote the solution to $p(q) = \mu(q)A/z$ for arbitrary $A > 0$. We then pick the specific value of A that satisfies the within-sector aggregator. For example:

- (i) **MONOPOLISTIC COMPETITION WITH KIMBALL DEMAND.** Let sector s consist of a mass $n_t(s) > 0$ firms and let sector output be given implicitly by the Kimball aggregator

$$\int_0^{n_t(s)} \Upsilon\left(\frac{y_{it}(s)}{y_t(s)}\right) di = 1 \quad (32)$$

where $\Upsilon(q)$ is strictly increasing and strictly concave. For this specification inverse demand $f(q)$ and the demand elasticity $\sigma(q)$ are given by

$$f(q) = \Upsilon'(q) \quad \text{and} \quad \sigma(q) = -\frac{\Upsilon'(q)}{\Upsilon''(q)q} \quad (33)$$

The associated demand index $d_t(s)$ is given by

$$d_t(s) = \left(\int_0^{n_t(s)} \Upsilon'(q_{it}(s))q_{it}(s) di \right)^{-1} \quad (34)$$

The scalar $A_t(s) := \Omega_t/p_t(s)d_t(s)$ is then pinned down by satisfying the Kimball aggregator and thus depends on the mass of firms $n_t(s)$.

- (ii) **OLIGOPOLISTIC COMPETITION WITH CES DEMAND.** Let sector s consist of a finite $n_t(s) \in \mathbb{N}$ firms and let sector output be given by the CES aggregator

$$\sum_{i=1}^{n_t(s)} \Upsilon\left(\frac{y_{it}(s)}{y_t(s)}\right) = 1, \quad \Upsilon(q) = q^{\frac{\gamma-1}{\gamma}} \quad (35)$$

where $\gamma > \eta > 1$ denotes the elasticity of substitution within sector s . Relative to the Kimball specification we have a finite number of firms, hence genuine strategic

⁷In this notation, a firm of size $q_{it}(s)$ has price $p_{it}(s) = f(q_{it}(s)) \times p_t(s)d_t(s)$.

interactions, but restrict the kernel of the aggregator $\Upsilon(q)$ to be a power function. For this specification inverse demand $f(q)$ is given by

$$f(q) = \Upsilon'(q) = \frac{\gamma - 1}{\gamma} q^{-\frac{1}{\gamma}} \quad (36)$$

while the demand index is simply

$$d_t(s) = \left(\sum_{i=1}^{n_t(s)} \Upsilon'(q_{it}(s)) q_{it}(s) \right)^{-1} = \frac{\gamma}{\gamma - 1} \quad (37)$$

Depending on whether competition is in quantities or prices, the demand elasticity facing a firm of size q is given by

$$\sigma(q) = \begin{cases} \left(\frac{1}{\eta} q^{\frac{\gamma-1}{\gamma}} + \frac{1}{\gamma} (1 - q^{\frac{\gamma-1}{\gamma}}) \right)^{-1} & \text{[Cournot competition]} \\ \eta q^{\frac{\gamma-1}{\gamma}} + \gamma (1 - q^{\frac{\gamma-1}{\gamma}}) & \text{[Bertrand competition]} \end{cases} \quad (38)$$

where $q^{\frac{\gamma-1}{\gamma}}$ is the sales share of a firm of size q , equal to the kernel of the aggregator $\Upsilon(q)$ in the CES case but not in general. The scalar $A_t(s) := \Omega_t/p_t(s)d_t(s)$ is pinned down by satisfying the CES aggregator and thus depends on $n_t(s)$.

With the relative size distribution $q_{it}(s)$ solved for in this way, we then know the distribution of markups $\mu_{it}(s) = \mu(q_{it}(s))$ and hence can compute sector-level productivity $z_t(s)$ and markups $\mu_t(s)$ and then aggregate productivity Z_t and the aggregate markup \mathcal{M}_t .

Entry and exit. Firms enter by paying a sunk cost κ in units of labor and then obtain a one-time productivity draw $z_i(s) \sim G(z)$ in a randomly allocated sector $s \in [0, 1]$. Let $N_t = \int_0^1 n_t(s) ds$ denote the aggregate mass of firms and let $M_t = \int_0^1 m_t(s) ds$ denote the aggregate mass of entrants. With a continuum of sectors, entry per sector $m_t(s)$ is IID Poisson with rate parameter M_t .⁸ Firms operate in their sector, obtaining a stream of profits $\pi_{it}(s)$, until they are hit with an IID exit shock, which happens with probability φ per period. For each sector s we then have

$$n_{t+1}(s) = (1 - \varphi)n_t(s) + m_t(s) \quad (39)$$

and hence the aggregate mass of firms evolves according to $N_{t+1} = (1 - \varphi)N_t + M_t$.

⁸With a finite number of sectors S , entry per sector $m_t(s)$ would be IID Binomial with number of trials $M_t S$ and success per trial $1/S$. Taking $S \rightarrow \infty$ this converges to a Poisson with rate parameter M_t .

Free entry condition. Now consider the decision problem of a potential entrant. In all versions of our model, entry occurs to the point at which ex ante expected discounted profits are offset by the sunk cost

$$\kappa W_t \geq \beta \sum_{j=1}^{\infty} (\beta(1-\varphi))^{j-1} \frac{C_t}{C_{t+j}} \int_0^1 \bar{\pi}_{t+j}(s) ds \quad (40)$$

with strict equality whenever $M_t > 0$ and where $\bar{\pi}_t(s)$ denotes expected profits conditional on operating in sector s . Where these market structures differ is in how these expected profits are calculated. Under *monopolistic competition*, with a continuum $[0, n_t(s)]$ of firms per sector, the entry of any individual firm i has no effect on sector-level variables. But under *oligopolistic competition*, with a finite $n_t(s) \in \mathbb{N}$ firms per sector, the entry of a new firm has non-negligible effects on post-entry sector-level variables. Specifically:

- (i) **MONOPOLISTIC COMPETITION.** Let $\pi_t(z_i, s) := \pi_{it}(s)$ denote the ex post profits of an individual firm with productivity draw z_i in sector s . In the monopolistic competition case, the expected profits conditional on operating in sector s are equal to the average profits of the incumbent firms in that sector

$$\bar{\pi}_t(s) = \int \pi_t(z_i, s) dG(z_i) \quad (41)$$

- (ii) **OLIGOPOLISTIC COMPETITION.** Let $\mathbf{z}(s)$ denote a sector-specific vector

$$\mathbf{z}(s) = (z_1(s), z_2(s), \dots, z_{n_t(s)}(s)) \quad (42)$$

of $n_t(s)$ independent draws from $G(z)$. Let $\pi_t(z_i, \mathbf{z}(s))$ denote the ex post profits of an individual firm with productivity z_i in a sector with $n_t(s)$ other firms with productivities $\mathbf{z}(s)$. The free-entry condition is again given by equation (40) but now the expected profits conditional on operating in sector s are given by

$$\bar{\pi}_t(s) = \iint \pi_t(z_i, \mathbf{z}(s)) dG_{n_t(s)}(\mathbf{z}(s)) dG(z_i) \quad (43)$$

where $G_{n_t(s)}(\mathbf{z}(s)) = G(z_1) \times G(z_2) \times \dots \times G(z_{n_t(s)}(s))$ denotes the joint distribution of the vector $\mathbf{z}(s)$. In the oligopolistic competition case, the *expected* profits from entering sector s are no longer equal to the *average* profits of those that do operate in sector s . There are two reasons for this. First, even if firms were identical, an entrant of non-negligible size would reduce the market shares of incumbents, tending to decrease expected profits. Second, sectors are heterogeneous, even two sectors with the same $n_t(s)$ will have different samples $\mathbf{z}(s)$, and, given this heterogeneity, Jensen's inequality can push expected profits *above* average profits.

2.3 Equilibrium

Given an initial mass of firms $n_0(s)$ per sector and an aggregate capital stock K_0 , an *equilibrium* is (i) a sequence of firm prices $p_{it}(s)$ and allocations $y_{it}(s)$, $k_{it}(s)$, $l_{it}(s)$, $x_{it}(s)$ and (ii) aggregate gross output Y_t , consumption C_t , investment I_t , materials X_t , labor L_t , wage rate W_t , rental rate R_t , and mass of entrants M_t such that firms and consumers optimize and the labor, capital and goods markets all clear. In particular

$$L_t = \iint l_{it}(s) di ds + \kappa M_t \quad (44)$$

$$K_t = \iint k_{it}(s) di ds \quad (45)$$

$$X_t = \iint x_{it}(s) di ds \quad (46)$$

(or the equivalent finite sums over i in the case of oligopolistic competition). Note that κM_t denotes labor used in the entry of new firms.

Solving the model. We discuss the solution method in [Appendix D](#). The key to solving the model is to recognize that aggregate markups \mathcal{M}_t , aggregate productivity Z_t and aggregate expected profits $\bar{\Pi}_t := \int_0^1 \bar{\pi}_t(s) ds$, are given by *time-invariant* functions of the aggregate mass of firms N_t , independent of all other aggregate variables, say $\mathcal{M}_t = \mathcal{M}(N_t)$, $Z_t = Z(N_t)$, and $\bar{\Pi}_t = \Pi(N_t)$. These functions summarize all the implications of market structure for aggregate outcomes. We solve the model by interpolating these functions and then use the remaining conditions, i.e., the production functions, input choices, optimality conditions of the representative consumer, and our aggregation results to simultaneously determine $Y_t, C_t, I_t, X_t, L_t, W_t, R_t, M_t$ given the state variables N_t and K_t .

3 Efficient allocation

In this section we derive the efficient allocation in our economy by considering the problem of a benevolent planner who faces the same technological and resource constraints as in the decentralized economy. Comparing the efficient allocation chosen by the planner to the decentralized allocation reveals three channels through which markups distort outcomes in the decentralized economy: (i) the aggregate markup acts like a uniform output tax, (ii) markup dispersion gives rise to misallocation of factors of production, and (iii) markups distort the entry margin.

3.1 Planner's problem

The planner chooses how many varieties to create, how to allocate inputs, consumption, investment, and employment so as to maximize the representative consumer's utility taking

as given the resource constraints for capital, labor and goods and the production functions for individual varieties. To facilitate comparisons with the decentralized equilibrium, the planner *cannot direct* the creation of new varieties towards specific sectors. We use asterisks to denote variables in the planner's problem.

The planner's problem has two parts: (i) a static allocation problem that determines aggregate productivity, and (ii) a dynamic problem that determines aggregate investment in new varieties, aggregate investment in physical capital, and aggregate employment. The link between the two parts is that the aggregate productivity solving the static allocation problem is a function of the stock of varieties, $Z_t^* = Z(N_t^*)$, which the planner internalizes when choosing how many varieties to create.

Dynamic problem. Starting with the dynamic problem, just as in the decentralized problem, we can use the resource constraints for capital, labor and goods and the production functions for individual varieties to derive the aggregate production function (22). We can then write the the planner's problem as maximizing

$$\sum_{t=0}^{\infty} \beta^t \left(\log C_t^* - \psi \frac{(\tilde{L}_t^* + \kappa(N_{t+1}^* - (1 - \varphi)N_t^*))^{1+\nu}}{1 + \nu} \right) \quad (47)$$

subject to the resource constraint for goods,

$$C_t^* + K_{t+1}^* + X_t^* = Z(N_t^*)F(K_t^*, \tilde{L}_t^*, X_t^*) + (1 - \delta)K_t^* \quad (48)$$

taking as given the function $Z(N_t^*)$ implied by the static allocation problem. The initial conditions for this problem are the mass of varieties N_0 and capital stock K_0 .

The planner's optimality conditions for consumption, investment, and employment are standard. The shadow wage is equated to the marginal product of labor

$$\psi C_t^* L_t^{*\nu} = Z_t^* F_{L,t}^* \quad (49)$$

while the marginal product of capital satisfies

$$1 = \beta \frac{C_t^*}{C_{t+1}^*} \left(Z_{t+1}^* F_{K,t+1}^* + 1 - \delta \right) \quad (50)$$

and the marginal product of materials is simply $Z_t^* F_{X,t}^* = 1$. Comparing these conditions with their decentralized counterparts, we see that the aggregate markup \mathcal{M}_t acts like a uniform output tax, reducing the overall scale of production and hence reducing the use of all inputs relative to the planner's problem.

Planner's choice of varieties. Now consider the planner's choice of varieties N_{t+1}^* . Letting $W_t^* = \psi C_t^* L_t^{*\nu}$ denote the shadow wage, we can write the first order condition

$$\kappa W_t^* = \beta \frac{C_t^*}{C_{t+1}^*} (1 - \varphi) \kappa W_{t+1}^* + \beta \frac{C_t^*}{C_{t+1}^*} \left(\frac{dZ_{t+1}^*}{dN_{t+1}^*} \frac{N_{t+1}^*}{Z_{t+1}^*} \right) \frac{Y_{t+1}^*}{N_{t+1}^*} \quad (51)$$

Iterating forward this gives

$$\kappa W_t^* = \beta \sum_{j=1}^{\infty} (\beta(1 - \varphi))^{j-1} \frac{C_t^*}{C_{t+j}^*} \left(\frac{dZ_{t+j}^*}{dN_{t+j}^*} \frac{N_{t+j}^*}{Z_{t+j}^*} \right) \frac{Y_{t+j}^*}{N_{t+j}^*} \quad (52)$$

This is the planner's counterpart to the free-entry condition in the decentralized problem. In the decentralized problem, a firm's incentive to enter is given by its expected discounted profits, which depend on its markup and sales. By contrast, the planner's incentive to create new varieties depends on the elasticity of aggregate productivity with respect to the mass of firms — and this depends on the solution to the static allocation problem.

Static allocation problem. Now consider the problem of maximizing aggregate productivity Z_t^* taking as given $n_t(s)$. The allocation of activity across sectors $q_t^*(s) = y_t^*(s)/Y_t^*$ is given by $q_t^*(s) = (z_t^*(s)/Z_t^*)^\eta$ so that in terms of sector-level productivity, aggregate productivity is $Z_t^* = (\int_0^1 z_t^*(s)^{\eta-1} ds)^{1/(\eta-1)}$, i.e., as in (30) but with no dispersion in sector-level markups. In turn, the allocation of activity within sectors $q_{it}^*(s) = y_{it}^*(s)/y_t^*(s)$ is given by

$$\Upsilon'(q_{it}^*(s)) d_t^*(s) = \frac{z_t^*(s)}{z_i(s)} \quad (53)$$

where $d_t^*(s)$ is the planner's demand index, the counterpart of (34) or (37). In other words, at the optimum the planner's shadow value of a variety is simply the planner's marginal cost of producing it. This optimality condition holds for both our monopolistic competition model with Kimball demand and our oligopolistic competition model with CES demand. As in the decentralized problem, the scalar $z_t^*(s)/d_t^*(s)$ is pinned down by satisfying the within-sector aggregator. In our oligopolistic competition model with CES demand this gives sector-level productivity $z_t^*(s) = (\sum_{i=1}^{n_t(s)} z_{it}^*(s)^{\gamma-1})^{1/(\gamma-1)}$ with constant demand index $d_t^*(s) = \frac{\gamma}{\gamma-1}$.

There is *misallocation* in the sense of Hsieh and Klenow (2009) whenever there is variation in marginal revenue products across firms, i.e., when the equilibrium $q_{it}(s)$ does not coincide with the planner's $q_{it}^*(s)$. This happens whenever markups $\mu_{it}(s)$ vary across firms.

Value of an additional variety. Now consider the value to the planner of an additional variety. Abstracting from any integer constraints on $n_t(s)$, an application of the envelope theorem gives

$$\frac{dZ_t^*}{dn_t(s)} \frac{n_t(s)}{Z_t^*} = (d_t^*(s) - 1) q_t^*(s) \frac{Z_t^*}{z_t^*(s)} \quad (54)$$

To interpret this condition, we use the planner’s demand index to write

$$d_t^*(s) - 1 = \int_0^{n_t(s)} (\epsilon_{it}^*(s) - 1) p_{it}^*(s) q_{it}^*(s) di \quad (55)$$

(or the equivalent finite sum in the case of oligopolistic competition), where we define

$$\epsilon_{it}^*(s) := \frac{\Upsilon(q_{it}^*(s))}{\Upsilon'(q_{it}^*(s))q_{it}^*(s)}, \quad \text{and} \quad p_{it}^*(s) := \Upsilon'(q_{it}^*(s))d_t^*(s) \quad (56)$$

The term $\epsilon_{it}^*(s)$ is the *inverse elasticity* of the within-sector aggregator $\Upsilon(q)$ evaluated at the planner’s allocation for a particular variety $q_{it}^*(s)$. The term $p_{it}^*(s)$ is the social value of an additional unit of that variety, i.e., the planner’s counterpart to the market price.

Comparing the free-entry condition in the decentralized equilibrium to the planner’s entry condition, we recover an important insight of [Bilbiie, Ghironi and Melitz \(2008, 2019\)](#), [Zhelobodko, Kokovin, Parenti and Thisse \(2012\)](#) and [Dhingra and Morrow \(2019\)](#), namely that the planner’s incentives to create new varieties are determined by the inverse elasticity $\epsilon_{it}^*(s)$ of the aggregator while the incentives for new firms to enter are determined by their markups $\mu_{it}(s)$. Whether there is too much or too little entry compared to the planner’s allocation is in general ambiguous and depends on precise details of the parameterization.

To summarize, variable markups distort outcomes in the decentralized economy through three channels: (i) the aggregate markup \mathcal{M}_t acts like a uniform output tax, (ii) markup dispersion $\mu_{it}(s)$ gives rise to misallocation of factors of production, and (iii) markups distort the entry margin.

4 Quantifying the model

In this section we outline our parameterization and calibration strategy and our model’s implications for the cross-sectional distribution of markups. We then calculate the aggregate productivity losses due to misallocation.

4.1 Benchmark parameterization

Kimball demand. To this point we have stressed aggregation results that hold regardless of the details of market structure within each sector. But to quantify the model we need to take a stand on demand and market structure. For our benchmark model we assume *monopolistic competition with Kimball demand*, as in [\(32\)](#) above. In particular, we assume the Kimball aggregator has the functional form introduced by [Klenow and Willis \(2016\)](#).

This specification implies that inverse demand curves are given by⁹

$$\Upsilon'(q) = \frac{\bar{\sigma} - 1}{\bar{\sigma}} \exp\left(\frac{1 - q^{\varepsilon/\bar{\sigma}}}{\varepsilon}\right), \quad \bar{\sigma} > 1 \quad (57)$$

which in turn implies that the demand elasticity $\sigma(q)$ is log-linear in relative size

$$\sigma(q) := -\frac{\Upsilon'(q)}{\Upsilon''(q)q} = \bar{\sigma} q^{-\varepsilon/\bar{\sigma}} \quad (58)$$

The parameter $\varepsilon/\bar{\sigma}$ is the elasticity of the demand elasticity with respect to relative size and is often known as the *super-elasticity*. If $\varepsilon = 0$ we have the constant demand elasticity $\sigma(q) = \bar{\sigma}$. If $\varepsilon > 0$, relatively large firms will face less elastic demand and charge high markups. If $\varepsilon < 0$, relatively large firms will face more elastic demand and charge low markups.

Productivity distribution. For parsimony and as is standard in the literature we assume that the distribution of productivity $G(z)$ is Pareto with tail parameter ξ .

Calibration strategy. We assign values to a number of conventional macro parameters that are held constant through all our quantitative exercises. We calibrate the parameters of the demand system and the productivity distribution to match facts on the amount of sales concentration and the relationship between markups and market shares *within* sectors.¹⁰

Assigned parameters. We assume that a period is one year and set the discount factor $\beta = 0.96$ and depreciation rate $\delta = 0.06$. We set the exit rate to $\varphi = 0.04$ to match the employment share of exiting firms, as in Boar and Midrigan (2020). We set the elasticity of value-added to capital $\alpha = 1/3$ and set the elasticity of substitution between value-added and materials to $\theta = 0.5$, both conventional values. Preferences (1) are homothetic and consistent with balanced growth. We set the inverse of the Frisch elasticity of labor supply to $\nu = 1$. We normalize the disutility from labor supply ψ and the entry cost κ to achieve a steady-state output of $Y = 1$ and a steady-state total mass of firms $N = 1$ for our benchmark economy. We report these parameter choices in Panel A of Table 1.

⁹The aggregator $\Upsilon(q)$ itself is given by

$$\Upsilon(q) = 1 + (\bar{\sigma} - 1) \exp\left(\frac{1}{\varepsilon}\right) \varepsilon^{\frac{\bar{\sigma}}{\varepsilon} - 1} \left[\Gamma\left(\frac{\bar{\sigma}}{\varepsilon}, \frac{1}{\varepsilon}\right) - \Gamma\left(\frac{\bar{\sigma}}{\varepsilon}, \frac{q^{\varepsilon/\bar{\sigma}}}{\varepsilon}\right) \right]$$

where $\Gamma(s, x) := \int_x^\infty t^{s-1} e^{-t} dt$ denotes the upper incomplete Gamma function.

¹⁰Our benchmark model with monopolistic competition features identical sectors so there is no variation in outcomes between sectors. In Section 6 we consider an alternative model with oligopolistic competition which features both within- and between-sector variation in concentration. We calibrate our oligopoly model to match within-sector concentration and the sector-level relationship between markups and market shares.

Table 1: Parameterization

Panel A: Assigned Parameters

β	discount factor	0.96
δ	depreciation rate	0.06
φ	exit rate	0.04
α	elasticity of value-added to capital	1/3
ν	elasticity of labor supply	1
θ	elasticity of substitution between value-added and materials	0.5

Panel B: Calibrated Parameters

<i>calibration targets</i>		<i>data</i>				
\mathcal{M}	aggregate markup	1.1 ~ 1.4	1.05	1.15	1.25	1.35
	top 5% sales share	0.57	0.57	0.57	0.57	0.57
	materials share	0.45	0.45	0.45	0.45	0.45
\hat{b}	regression coefficient	0.16	0.16	0.16	0.16	0.16
<i>parameter values</i>						
ξ	Pareto tail		20.70	6.84	4.07	2.89
$\bar{\sigma}$	demand elasticity		29.10	10.86	7.21	5.66
$\varepsilon/\bar{\sigma}$	super-elasticity		0.16	0.16	0.16	0.16
ϕ	weight on value-added		0.51	0.43	0.33	0.21

Panel A reports assigned parameters held constant through all our quantitative exercises. Panel B reports calibrated parameters for our benchmark model with monopolistic competition and Kimball demand. We report four cases corresponding to alternative targets for the level of the aggregate markup, $\mathcal{M} = 1.05, 1.15, 1.25$ and 1.35 , over the range of \mathcal{M} implied by the US Census of Manufactures from 1972 to 2012, as discussed in [Appendix B](#). For each \mathcal{M} we calibrate the Pareto tail ξ , demand elasticity $\bar{\sigma}$, super-elasticity $\varepsilon/\bar{\sigma}$ and weight on value-added ϕ to match the targets shown in Panel B. For each model we choose the super-elasticity $\varepsilon/\bar{\sigma}$ so that the slope coefficient b from equation (59) in the model matches the estimated slope coefficient \hat{b} . See the text for more details.

Calibrated parameters. The level and dispersion of markups in our benchmark model depend crucially on three underlying parameters: (i) the Pareto tail parameter ξ , (ii) the super-elasticity $\varepsilon/\bar{\sigma}$ that determines the sensitivity of a firm’s demand elasticity to its relative size, and (iii) the ‘average’ demand elasticity $\bar{\sigma}$. Intuitively, the Pareto tail parameter ξ is pinned down by the amount of concentration in the distribution of firm size, the super-elasticity $\varepsilon/\bar{\sigma}$ is pinned down by the cross-sectional relationship between markups and market shares, and $\bar{\sigma}$ is pinned down by the overall level of markups. Specifically we target:

- (i) SALES CONCENTRATION. The Pareto tail parameter ξ is pinned down by our target for sales concentration. We target the average sales share of the top 5% of firms (by market share) in 6-digit NAICS sectors. For 2012 US manufacturing, the top 5% of firms on average account for 57% of sales.
- (ii) RELATIONSHIP BETWEEN MARKUPS AND MARKET SHARES. The super-elasticity $\varepsilon/\bar{\sigma}$ is pinned down by the relationship between firm-level markups and market shares in our model.¹¹ As discussed in detail below, in our benchmark model the super-elasticity $\varepsilon/\bar{\sigma}$ corresponds to the slope coefficient b in a regression of (transformed) markups on market shares. We estimate this regression on firm-level data from the US Census of Manufactures 1972 to 2012 and obtain a precisely estimated $\hat{b} = 0.16$. In our benchmark model this slope coefficient *is* the super-elasticity so for our benchmark model we set $\varepsilon/\bar{\sigma} = 0.16$. In other versions of our model with different demand systems we use indirect inference, choosing parameters so that the slope coefficient in the model matches the estimated slope coefficient $\hat{b} = 0.16$.
- (iii) AGGREGATE MARKUP. The average elasticity $\bar{\sigma}$ is pinned down by our target for the aggregate markup \mathcal{M} . As discussed in [Appendix B](#), the aggregate markup we compute in the Census of Manufactures data ranges from about 1.1 to 1.4 depending on the Census year and the specification. The existing literature on markups in the US economy also provides a wide range of estimates for \mathcal{M} .¹² Given this range of estimates, rather than commit to a single target for the aggregate markup, for our benchmark model we recalibrate $\bar{\sigma}$ (jointly, with our other parameters) for \mathcal{M} ranging from 1.05 to 1.45.

Finally, we calibrate the weight ϕ on value-added in the gross-output production function by targeting a materials share of 45% for the US economy in 2012. For each \mathcal{M} we calibrate this parameter jointly with the three key parameters $\xi, \varepsilon/\bar{\sigma}$, and $\bar{\sigma}$ as discussed above.

¹¹This is similar to how we estimated the within-industry relationship between market shares and markups in [Edmond, Midrigan and Xu \(2015\)](#) but adapted to the Kimball demand system used here.

¹²See e.g., [Atkeson, Burstein and Chatzikonstantinou \(2019\)](#), [Barkai \(2020\)](#), [De Loecker, Eeckhout and Unger \(2020\)](#), [Gutiérrez and Phillippon \(2017a,b\)](#), and [Hall \(2018\)](#) etc. [Basu \(2019\)](#) surveys this literature.

Regression specification details. The key to our calibration strategy is the relationship between markups and market shares used to pin down the super-elasticity. To derive this relationship we use the fact that in our model both markups $\mu_{it}(s)$ and market shares $\omega_{it}(s)$ are strictly increasing functions of relative size $q_{it}(s)$. Eliminating $q_{it}(s)$ we can then write markups as a strictly increasing function of market shares. In particular, as shown in [Appendix B](#), in a version of our model with time-invariant firm-specific demand shifters and sector specific Kimball aggregators, the relationship between market shares and markups works out to be

$$\frac{1}{\mu_{it}(s)} + \log \left(1 - \frac{1}{\mu_{it}(s)} \right) = a(s) + a_i(s) + a_t(s) + b(s) \log \omega_{it}(s), \quad b(s) = \frac{\varepsilon(s)}{\bar{\sigma}(s)} \quad (59)$$

where the firm fixed effects $a_i(s)$ control for the time-invariant firm-specific demand shifters and the sector-time fixed effects $a_t(s)$ control for sector-time variation in the Kimball demand index. The transformation on the LHS is strictly increasing in $\mu_{it}(s)$ and independent of other parameters. In this sense the slope coefficient $b(s)$ on the RHS is a measure of the strength of the within-sector relationship between markups and market shares. For our benchmark calibration we take the model at face-value and impose a common slope coefficient $b(s) = b$.¹³ We estimate this regression using data from the US Census of Manufactures from 1972 to 2012. We construct firm-level markups $\mu_{it}(s)$ as discussed below and market shares $\omega_{it}(s)$ within each 6-digit NAICS sector for each Census year. As reported in [Table 2](#), we obtain an estimated slope coefficient $\hat{b} = 0.162$ with standard error 0.002 clustered at the firm level.

Firm-level markups. As discussed in [Appendix B](#), to implement this regression we infer firm-level markups $\mu_{it}(s)$ from the cost-minimization condition¹⁴

$$\mu_{it}(s) = \frac{p_{it}(s)y_{it}(s)}{W_t l_{it}(s)} \times \alpha_t^l(s) \quad (60)$$

Our key assumption is that the elasticity of output with respect to labor $\alpha_t^l(s)$ is *common to all firms within a sector*.¹⁵ Under constant returns to scale,¹⁶ we then have, for each firm

$$\alpha_t^l(s) = \frac{W_t l_{it}(s)}{W_t l_{it}(s) + R_t k_{it}(s) + x_{it}(s)} \quad (61)$$

We estimate this elasticity by averaging [\(61\)](#) over firms within each 6-digit NAICS sector.¹⁷ We allow this elasticity to vary over time by constructing it for each Census year. We then have an estimate of $\alpha_t^l(s)$ that we can plug back into [\(60\)](#) to construct $\mu_{it}(s)$.

¹³We discuss the sensitivity of our results to this common slope coefficient assumption in [Appendix C](#).

¹⁴For multi-establishment firms we construct establishment-level markups $\mu_{eit}(s)$ and then aggregate to firm-level markups $\mu_{it}(s)$ weighting establishments e by their share of the firm's wage bill.

¹⁵For our benchmark model, this elasticity is $\alpha_t^l(s) = (1 - \alpha)\zeta_t$, i.e., the elasticity of output with respect to value-added ζ_t times the elasticity of value-added with respect to labor $(1 - \alpha)$, see [Appendix B](#).

¹⁶Our results are robust to relaxing the assumption of constant returns to scale, see [Appendix C](#).

¹⁷We take this average to reduce the role of measurement error. These calculations also use sector-specific user costs of capital from the NBER-CES and BLS as in [Foster, Grim and Haltiwanger \(2016\)](#).

Table 2: Relationship between Markups and Market Shares

Dependent Variable	$\frac{1}{\mu_{it}(s)} + \log\left(1 - \frac{1}{\mu_{it}(s)}\right)$		
$\log \omega_{it}(s)$	0.063 (0.001)	0.162 (0.002)	0.187 (0.003)
Sector \times Year FE	Y	Y	Y
Firm FE		Y	Y
Firm Age			Y
R^2	0.084	0.531	0.540
Observations	609,000	369,000	315,000

Firm-level markups $\mu_{it}(s)$ constructed from the US Census of Manufactures from 1972 to 2012, as discussed in the text. Market shares $\omega_{it}(s)$ of firm i within each 6-digit NAICS sector s . We include sector \times year fixed effects to control for sector-specific shifts in the Kimball demand index $d_t(s)$. Our benchmark specification also includes firm fixed effects to control for any time-invariant firm-specific component of demand. Results robust to including firm age. Standard errors clustered at the firm level.

An alternative to this would be to estimate sector-specific production functions. But recent work by [Bond, Hashemi, Kaplan and Zoch \(2021\)](#) demonstrates that in the presence of variable markups it is not possible to consistently estimate output elasticities when only revenue data is available.¹⁸ Using the simple labor input expenditure share approach also makes our results easier to compare to recent empirical work, such as [Autor, Dorn, Katz, Patterson and Van Reenen \(2020\)](#) and [De Loecker, Eeckhout and Unger \(2020\)](#), that also report such measures.

Other distortions. Our model abstracts from other distortions at either the firm- or sector-level that may drive a wedge between firm revenues and expenditure on labor input. If the relationship between markups and market shares was log-linear, we could use fixed effects to control for persistent firm- or sector-level distortions that confound the measurement of markups in (60). In a robustness exercise, we implement this approach by taking a log-linear approximation to the LHS of (59). See [Appendix C](#) for details.

Model fit. Panel B of [Table 1](#) reports the parameter values that minimize our objective function for four values of the aggregate markup, $\mathcal{M} = 1.05, 1.15, 1.25$ and 1.35 . To match a low level of markups, $\mathcal{M} = 1.05$, while targeting a top 5% sales share of 0.57 requires

¹⁸That said, [De Ridder, Grassi and Morzenti \(2022\)](#) show by simulation that markups estimated using revenue data are systematically related to the true markups in their model. In this sense the revenue-based estimates are informative about markup variation even if not informative about markup levels.

Table 3: Markup Dispersion and Productivity Losses

<i>cost-weighted distribution of markups</i>				
aggregate markup, \mathcal{M}	1.05	1.15	1.25	1.35
p25 markup	1.04	1.11	1.17	1.23
p50 markup	1.05	1.14	1.23	1.31
p75 markup	1.06	1.18	1.31	1.43
p90 markup	1.07	1.23	1.40	1.58
p99 markup	1.11	1.35	1.63	1.97
<i>aggregate productivity losses, %</i>				
gross output	0.28	0.97	1.83	2.86
value-added	0.61	2.71	6.08	10.73
value-added, $\mathcal{M} = 1$	0.51	1.85	3.63	5.85

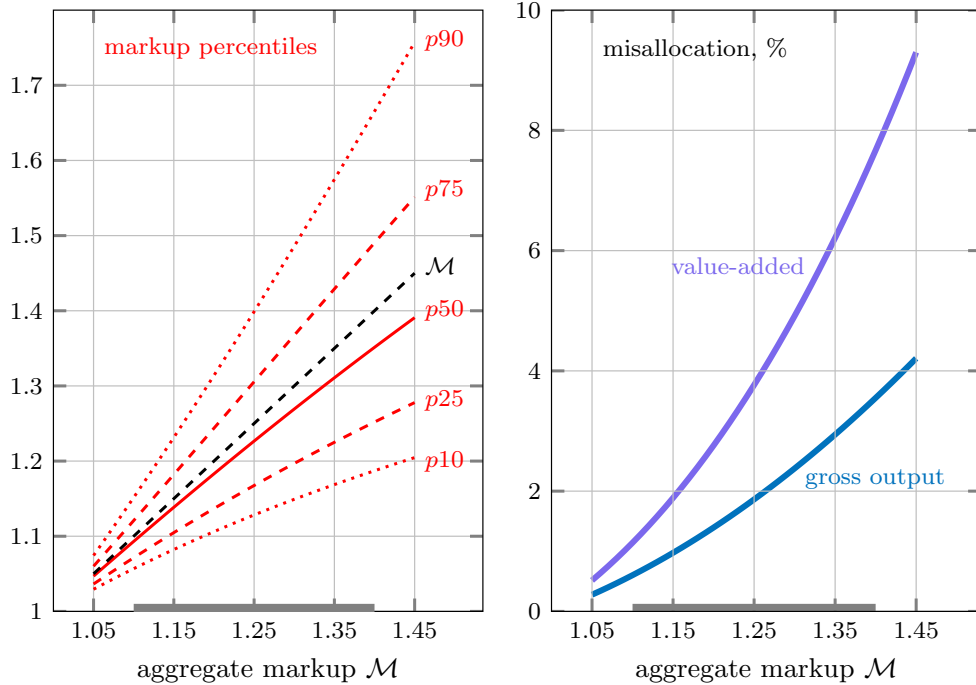
Cost-weighted steady-state distribution of markups and aggregate productivity losses for four calibrations of our benchmark model, corresponding to targets for the aggregate markup $\mathcal{M} = 1.05, 1.15, 1.25$ and 1.35 . Gross output aggregate productivity loss is $(Z - Z^*)/Z^* \times 100$, and similarly for the value-added aggregate productivity loss. To isolate the effect of misallocation on value-added aggregate productivity we also report the value-added aggregate productivity loss with the same amount of markup dispersion but holding $\mathcal{M} = 1$ to eliminate the distortion between value-added and materials, see text for details.

a high average demand elasticity, $\bar{\sigma} = 29.1$, and a thin-tailed productivity distribution, $\xi = 20.7$. To match a high level of markups, $\mathcal{M} = 1.35$, while targeting the same top 5% sales share requires a much lower average demand elasticity, $\bar{\sigma} = 5.66$, and a fatter-tailed productivity distribution $\xi = 2.89$. Though our estimate of $\varepsilon/\bar{\sigma} = 0.162$ is much lower than typically assumed in macro studies that attempt to match the response of prices to changes in monetary policy or exchange rates, it is in line with the micro estimates surveyed by [Klenow and Willis \(2016\)](#). In [Appendix B](#) we find an almost identical super-elasticity $\varepsilon/\bar{\sigma} = 0.16$ best fits the relationship between markups and market shares in the Taiwanese manufacturing firms studied by [Edmond, Midrigan and Xu \(2015\)](#).

4.2 Markups and misallocation

Markup distribution. [Table 3](#) reports the cost-weighted steady-state distribution of markups in our model for the same four values of the aggregate markup. As we target higher levels of the aggregate markup \mathcal{M} the model implies more markup dispersion. This occurs because as we target higher \mathcal{M} , requiring a lower average demand elasticity $\bar{\sigma}$, we need a fatter-tailed productivity distribution to hold the top 5% sales share unchanged. In turn, a fatter-tailed productivity distribution creates more large firms who charge large markups,

Figure 1: Markup Distribution and Misallocation in Benchmark Model



Left panel shows cost-weighted steady-state markup distribution in our benchmark model with monopolistic competition and Kimball demand for a range of targets for the aggregate markup \mathcal{M} . For each \mathcal{M} we recalibrate the Pareto tail ξ , demand elasticity $\bar{\sigma}$, super-elasticity $\varepsilon/\bar{\sigma}$ and weight on value-added ϕ to match the calibration targets in Table 1. Right panel shows implied amounts of misallocation in aggregate gross output and aggregate value-added. To isolate the effect of misallocation on value-added aggregate productivity we report the value-added aggregate productivity loss with the same amount of markup dispersion but holding $\mathcal{M} = 1$ to eliminate the distortion between value-added and materials, see text for details. Shaded interval indicates range of \mathcal{M} implied by the US Census of Manufactures from 1972 to 2012, as discussed in Appendix B.

increasing markup dispersion. We illustrate this in Figure 1 using a fine grid for \mathcal{M} .

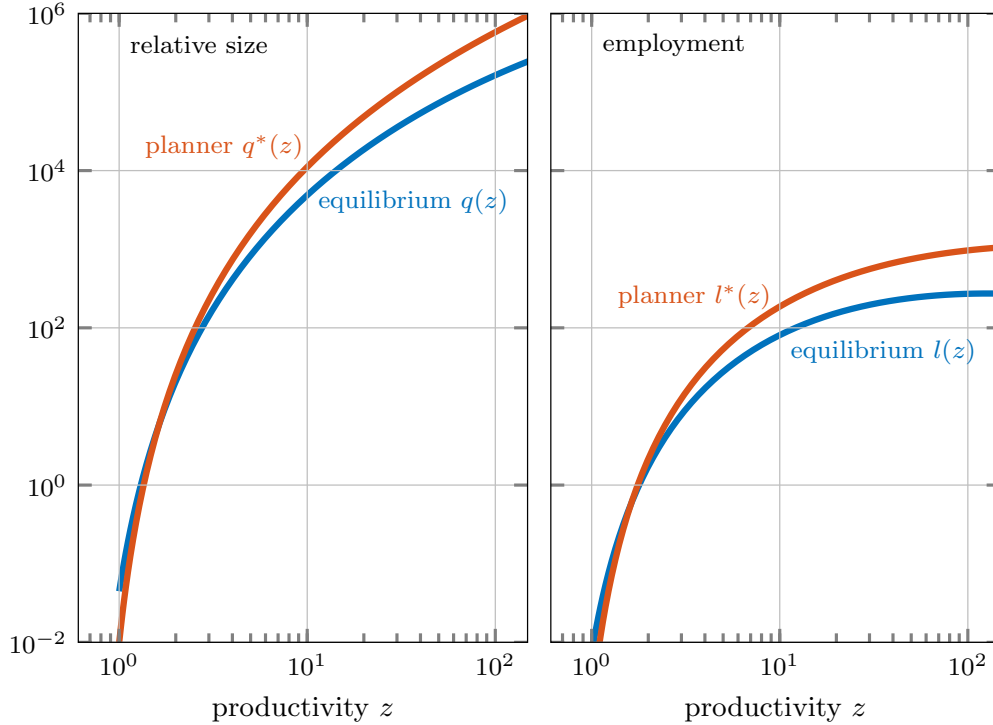
Misallocation. The markup dispersion generated by our model implies that there are aggregate productivity losses due to misallocation. For *gross output* aggregate productivity we compare Z in the the steady state of our benchmark economy to the level of gross output aggregate productivity Z^* that could be achieved by a planner facing the same technology and resource constraints who could reallocate factors of production across producers. As shown in the right panel of Figure 1, for the empirically plausible range of \mathcal{M} , gross output aggregate productivity Z in our benchmark economy is on the order of 1% to 3% below the level of gross output aggregate productivity Z^* that could be achieved by a planner.

We also compute *value-added* aggregate productivity losses. In Appendix G we show that value-added aggregate productivity can be written

$$Z_{\text{value-added}} = \phi^{\frac{1}{\theta-1}} \frac{(1 - (1 - \phi) Z^{\theta-1} \mathcal{M}^{-\theta})}{(1 - (1 - \phi) Z^{\theta-1} \mathcal{M}^{1-\theta})^{\frac{\theta}{\theta-1}}} Z \quad (62)$$

where ϕ is the weight on value-added, θ is the elasticity of substitution between value-

Figure 2: Equilibrium and Planner Allocations



The left panel shows the equilibrium relative size $q(z)$ and the planner's relative size $q^*(z)$ as functions of productivity for our benchmark economy with $\mathcal{M} = 1.15$. The right panel shows the equilibrium employment $l(z)$ and the planner's employment $l^*(z)$ for the same economy. More productive firms have higher markups and produce too little and employ too little compared to the planner's allocation. Less productive firms produce too much and employ too much compared to the planner's allocation. In this figure aggregate employment in the decentralized equilibrium is the same as aggregate employment for the planner. Our measure of *misallocation* is the aggregate output loss implied by the equilibrium allocation relative to the planner's allocation.

added and materials in the gross output production function, and where as above Z is gross output aggregate productivity. While the level of gross output aggregate productivity Z is independent of the level of the aggregate markup \mathcal{M} , depending only on markup dispersion, the level of value-added aggregate productivity *does* depend on the level of \mathcal{M} . This is because the aggregate markup \mathcal{M} directly distorts the choice of materials relative to value-added.

For the planner, value-added aggregate productivity works out to be

$$Z_{\text{value-added}}^* = \phi^{\frac{1}{\theta-1}} \frac{(1 - (1 - \phi) Z^{*\theta-1})}{(1 - (1 - \phi) Z^{*\theta-1})^{\frac{\theta}{\theta-1}}} Z^* \quad (63)$$

where Z^* is the planner's gross output aggregate productivity. In short, markups reduce value-added aggregate productivity relative to the efficient allocation both because markup dispersion reduces Z relative to Z^* and because the aggregate level of markups \mathcal{M} distorts the use of materials relative to value-added. We report these value-added aggregate productivity losses in [Table 3](#). To isolate the role of markup dispersion we also report the value-added productivity losses that would arise if $\mathcal{M} = 1$, as shown in the right panel of [Figure 1](#).

To illustrate the difference in allocations, [Figure 2](#) compares the relative size $q(z)$ and employment $l(z)$ of a firm with productivity z in the decentralized equilibrium to the planner’s counterparts $q^*(z)$ and $l^*(z)$. More productive firms have higher markups and produce and employ too little compared to the planner’s allocation. Less productive firms produce and employ too much compared to the planner’s allocation. Notice that the planner’s allocation is not log-linear in productivity, as it would be with CES demand. The extra concavity reflects strongly diminishing marginal productivity as the relative size q increases. If misallocation losses were calculated assuming a constant demand elasticity $\bar{\sigma}$ rather than variable demand elasticities $\sigma(q) = \bar{\sigma}q^{-\varepsilon/\bar{\sigma}}$ we would find higher misallocation (for a given amount of dispersion in marginal revenue products) because we would overstate the gains from reallocating factors from small, less productive firms to large, more productive firms.

Comparison with Baqaee and Farhi (2020). In related work, [Baqaee and Farhi \(2020\)](#) calculate that the value-added aggregate productivity gains from eliminating all markups are about 20%, about twice as large as the value-added aggregate productivity gains in even the most extreme calibration of our model. Why do they find much larger effects of markup dispersion on productivity? The key point is that they feed into their calculation all the variation in *estimated markups* (e.g., as in [De Loecker, Eeckhout and Unger, 2020](#); [Gutiérrez and Phillippon, 2017b](#)) whereas we feed in that component of markups that systematically varies with firm market shares. In this sense, we use only that part of the cross-sectional variation in markups that is correlated with firm relative size. Because the estimated markups they use are more dispersed than the markups implied by our model, they find larger effects of markup dispersion on aggregate productivity.¹⁹

5 How costly are markups?

We now present our main results on the welfare costs of markups. We first quantify the total welfare costs of markups in our benchmark economy for a range of values for the aggregate markup \mathcal{M} . We then show how the efficient allocation can be implemented by a specific nonlinear schedule of size-dependent subsidies and show how to isolate aspects of this policy to quantify the relative magnitudes of the different markup channels. We also study simple entry subsidies that *indirectly* affect markup distortions through the amount of competition.

We measure the welfare costs of markups by asking how much the representative consumer would benefit from implementing the efficient allocation that eliminates all markup distortions, taking the transitional dynamics into account. We find that the total welfare costs of markups are not only increasing in our target for \mathcal{M} they are increasing and *convex*

¹⁹See [Eslava and Haltiwanger \(2020\)](#) who study the life-cycle of Colombian manufacturing plants and find that markup variation plays only a small role in accounting for variation in average revenue products.

in \mathcal{M} . Because of this, the total welfare costs can be large. For example, for an economy with aggregate markup $\mathcal{M} = 1.15$, implementing the efficient allocation results in a consumption-equivalent welfare gain of about 8.7%, rising to 23.6% for an economy with $\mathcal{M} = 1.25$ and 49.7% for an economy with $\mathcal{M} = 1.35$. We find that a uniform output subsidy that offsets the aggregate markup alone goes a long way towards achieving full efficiency.

5.1 Welfare cost of markups

We first compare the distorted steady state in our decentralized equilibrium to that chosen by a planner, then calculate the welfare gains from implementing the efficient steady state taking the transitional dynamics into account.

Steady state comparisons. The first six columns of [Table 4](#) report the percentage change in consumption C , gross output Y , employment L , mass of firms N , physical capital K , and aggregate productivity Z from the initial distorted steady state to the efficient steady state for each of four values of the aggregate markup \mathcal{M} . The efficient steady state features higher consumption, higher output, and employment. Aggregate productivity is higher, both because of the elimination of misallocation and because of the increase in product variety, i.e., increase in the mass of firms N .²⁰

Welfare gains from implementing efficient allocation. The last column of [Table 4](#) reports the welfare gains for the representative consumer in consumption-equivalent units including the transition, i.e., these take into account the deferred increase in consumption as investment in physical capital and product variety accumulates over time. These dynamics also take into account the time path of employment. We find that if the aggregate markup is low, $\mathcal{M} = 1.05$, the representative consumer needs to be compensated with an additional 1.34% consumption per period in order to be indifferent between the initial distorted steady state and the transition to the efficient steady state. This increases to 8.67% consumption per period if the aggregate markup is $\mathcal{M} = 1.15$ and to 49.66% consumption per period if the aggregate markup is $\mathcal{M} = 1.35$. The welfare gains are higher when we target higher \mathcal{M} . Indeed the gains are *convex* in \mathcal{M} . As we target higher \mathcal{M} for the benchmark economy, both the level of markups and the amount of markup dispersion increase. We illustrate this convexity in [Figure 3](#) using a fine grid for \mathcal{M} with the upper bound extended to 1.45.

5.2 Implementing the efficient allocation

We now show how the efficient allocation can be implemented by a specific nonlinear schedule of size-dependent subsidies. This policy removes the aggregate markup distortion, removes

²⁰We discuss the effects of variety on aggregate productivity in more detail in [Appendix I](#).

Table 4: Implications of Alternative Policies, Benchmark Model

		steady state comparisons, %						
		Y	C	L	N	K	Z	welfare, %
$\mathcal{M} = 1.05$	efficient	15.3	11.0	6.0	12.3	24.1	0.9	1.34
	uniform subsidy	13.7	9.1	5.7	3.4	22.0	0.2	0.65
	size-dependent subsidy	1.4	1.7	0.3	8.1	1.7	0.7	0.71
	entry subsidy	1.1	1.3	0.4	11.3	1.4	0.6	0.06
$\mathcal{M} = 1.15$	efficient	59.6	44.5	18.0	20.1	100.4	4.1	8.67
	uniform subsidy	51.8	35.8	17.0	9.5	88.5	1.5	5.90
	size-dependent subsidy	5.3	6.2	1.0	8.3	6.6	2.3	2.87
	entry subsidy	6.3	7.4	2.4	20.0	8.1	3.0	0.56
$\mathcal{M} = 1.25$	efficient	134.0	102.3	30.1	26.7	246.2	8.9	23.64
	uniform subsidy	112.7	79.7	28.2	15.0	208.2	3.9	17.36
	size-dependent subsidy	10.8	12.5	1.8	8.1	13.6	4.1	6.26
	entry subsidy	17.4	20.3	6.0	29.1	23.0	7.4	1.98
$\mathcal{M} = 1.35$	efficient	263.2	203.4	42.1	32.2	540.3	15.1	49.66
	uniform subsidy	213.5	152.9	39.2	19.8	435.1	7.4	37.41
	size-dependent subsidy	18.4	20.9	2.7	7.6	23.6	6.0	11.32
	entry subsidy	38.6	44.9	11.9	39.0	52.6	14.0	5.11

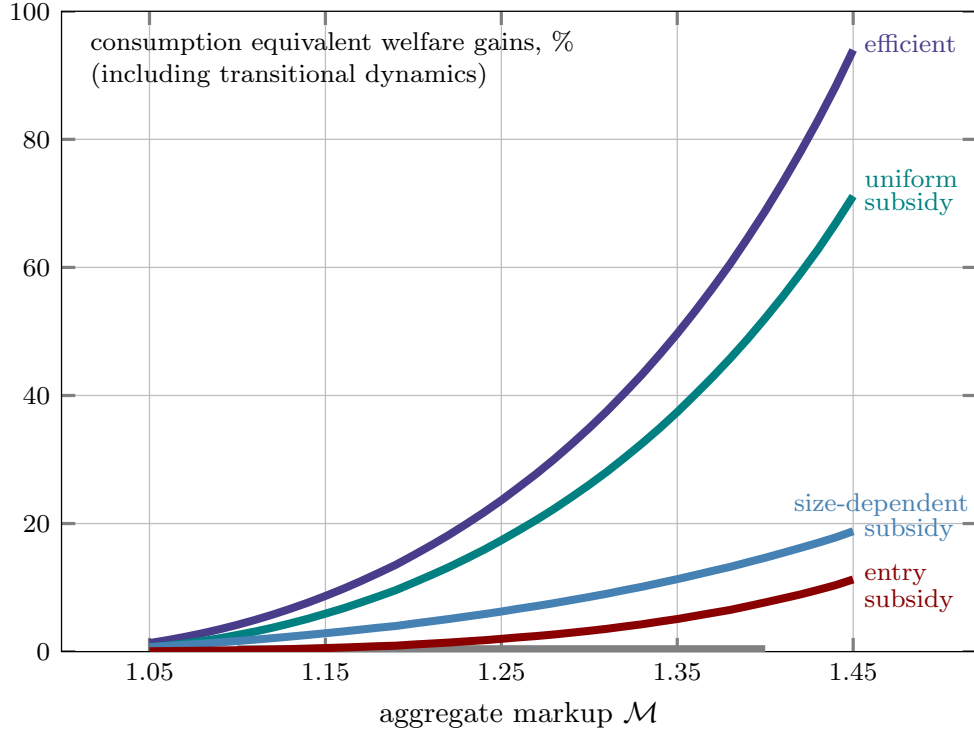
The first six columns report the percentage change from the initial distorted steady state to the new steady state. The last column reports the consumption equivalent welfare gains (including transitional dynamics). For each \mathcal{M} we recalibrate the Pareto tail ξ , demand elasticity $\bar{\sigma}$, super-elasticity $\varepsilon/\bar{\sigma}$ and weight on value-added ϕ . The alternative policies are (i): the *efficient allocation*, where all markups are removed, (ii) a *uniform subsidy* that eliminates the aggregate markup, (iii) *size-dependent subsidies* that eliminate misallocation and the entry distortion, and (iv) the uniform *entry subsidy* that leads to the largest welfare gain.

markup dispersion (and hence misallocation), and removes the entry distortion. We then show how to isolate different aspects of this policy to quantify the relative magnitudes of the different markup channels. This policy is financed by lump-sum taxes on the representative consumer. We view these calculations as a device for isolating the role of each distortion. The actual consequences of such a policy would of course be much more complex in economies with heterogeneous consumers and other frictions (see [Boar and Midrigan, 2020](#), for example).

Direct policy intervention to remove markup distortions. In the decentralized equilibrium, the profits of a firm with productivity z facing Kimball demand can be written

$$\pi_t(z) = \left[\Upsilon'(q_t(z))q_t(z)D_t - \frac{\Omega_t}{z} q_t(z) \right] Y_t \quad (64)$$

Figure 3: Welfare Gains from Alternative Policies



Consumption equivalent welfare gains (including transitional dynamics), from the initial distorted steady state to the new steady state for a range of targets for the aggregate markup \mathcal{M} . For each \mathcal{M} we recalibrate the Pareto tail ξ , demand elasticity $\bar{\sigma}$, super-elasticity $\varepsilon/\bar{\sigma}$ and weight on value-added ϕ . The alternative policies are (i): the *efficient allocation*, where all markups are removed, (ii) a *uniform subsidy* that eliminates the aggregate markup, (iii) *size-dependent subsidies* that eliminate misallocation and the entry distortion, and (iv) the uniform *entry subsidy* that leads to the largest welfare gain. Shaded interval indicates range of \mathcal{M} implied by the US Census of Manufactures from 1972 to 2012, as discussed in [Appendix B](#).

where D_t denotes the Kimball demand index from (34) above.²¹ Now suppose that firms are paid a size-dependent subsidy $T_t(q)$ given by

$$T_t(q) = \left[\Upsilon(q) - \Upsilon'(q)q \right] D_t Y_t \quad (65)$$

This policy takes away revenues in proportion to $\Upsilon'(q)q$ and returns revenues in proportion to $\Upsilon(q)$ which will then induce firms to price at marginal cost. In particular, given the subsidy $T_t(q)$, a firm has net profits $\hat{\pi}_t(z) := \pi_t(z) + T_t(q_t(z))$ which simplifies to

$$\hat{\pi}_t(z) = \left[\Upsilon(q_t(z)) D_t - \frac{\Omega_t}{z} q_t(z) \right] Y_t \quad (66)$$

This leads to the optimal price

$$p_t(z) = \Upsilon'(q_t(z)) D_t = \frac{\Omega_t}{z} \quad (67)$$

²¹In our benchmark economy, sectors $s \in [0, 1]$ are ex post identical and we have $d_t(s) = D_t$, $y_t(s) = Y_t$, $p_t(s) = 1$, $z_t(s) = Z_t$, $n_t(s) = N_t$ etc.

In other words, this policy induces firms to price at marginal cost with firm-level wedge $\mu_t(z) = 1$. Hence the aggregate wedge is also $\mathcal{M}_t = 1$. Given this, net profits are equal to the transfer $\hat{\pi}_t(z) = T_t(q_t(z))$ and so the free entry condition becomes

$$\begin{aligned} \kappa W_t &= \beta \sum_{j=1}^{\infty} (\beta(1-\varphi))^{j-1} \frac{C_t}{C_{t+j}} \int \left[\Upsilon(q_{t+j}(z)) - \Upsilon'(q_{t+j}(z))q_{t+j}(z) \right] D_{t+j} Y_{t+j} dG(z) \\ &= \beta \sum_{j=1}^{\infty} (\beta(1-\varphi))^{j-1} \frac{C_t}{C_{t+j}} (D_{t+j} - 1) \frac{Y_{t+j}}{N_{t+j}} \end{aligned} \quad (68)$$

where the second line follows using the definitions of the Kimball aggregator (32) and its demand index (34). To see how the free-entry condition under this policy compares to the planner's entry condition, use (54) to write the planner's elasticity of aggregate productivity with respect to new varieties

$$\frac{dZ_t^*}{dN_t^*} \frac{N_t^*}{Z_t^*} = (D_t^* - 1) \quad (69)$$

Plugging this elasticity into the planner's entry condition (52) we see that the free-entry condition under the policy $T_t(q)$ coincides with the planner's entry condition, i.e., this policy also eliminates the entry distortion. We next show how to use a generalization of this policy to isolate and quantify the relative importance of each channel.

5.3 Decomposing the implementation.

The nonlinear schedule $T_t(q)$ directly implements the efficient allocation. To study each channel in isolation, it is helpful to generalize this to

$$T_t(q) = \left[a_0 \Upsilon(q) + a_1 \Upsilon'(q)q \right] D_t Y_t \quad (70)$$

We can then recover the main cases of interest by setting the policy parameters a_0, a_1 appropriately. There are three main cases of interest: (i) setting $a_0 = 1$ and $a_1 = -1$ implements the *efficient allocation* as discussed above, (ii) setting $a_0 = 0$ and $a_1 = \chi > 0$ implements a *uniform subsidy* that leaves the dispersion in marginal revenue products unchanged but drives the aggregate wedge down to $\mathcal{M}/(1+\chi)$, while (iii) setting $a_0 = 1/(1+\chi)$ and $a_1 = -1$ implements *size-dependent subsidies* that eliminate the dispersion in marginal revenue products while leaving an aggregate wedge equal to $1+\chi$.

Uniform subsidy. Setting $a_0 = 0, a_1 = \chi$ implements a *uniform subsidy* giving net profits

$$\hat{\pi}_t(z) = \left[(1+\chi)\Upsilon'(q_t(z))q_t(z)D_t - \frac{\Omega_t}{z} q_t(z) \right] Y_t \quad (71)$$

which leads firms to set the price

$$p_t(z) = \frac{\mu_t(z)}{1+\chi} \frac{\Omega_t}{z} \quad (72)$$

where $\mu_t(z)$ is the benchmark markup of a firm with productivity z . This subsidy induces firms to produce more and to use more of each input, driving the wedge between price and marginal cost down to $\mu_t(z)/(1 + \chi)$ and driving the aggregate wedge in the optimality conditions of the representative firm down to $\mathcal{M}_t/(1 + \chi)$. Thus by setting $\chi = \mathcal{M} - 1$ for the initial distorted steady state we can put in motion a transition to a new steady state where the aggregate wedge has been eliminated. But this uniform subsidy has no effect on relative markups and so leaves steady state misallocation unchanged. This subsidy affects the entry condition but generally leaves it distorted.

Table 4 reports the effect of introducing the uniform subsidy on steady state outcomes for four levels of the aggregate markup \mathcal{M} . Figure 3 reports the effect on welfare, including the transitional dynamics, for a fine grid of \mathcal{M} . For all levels of \mathcal{M} , the uniform subsidy accounts for a large share of the potential welfare gains. For example, if the aggregate markup is low, $\mathcal{M} = 1.05$, the uniform subsidy increases gross output by 13.7%, consumption by 9.1%, and employment by 5.7%. These increases are only slightly smaller than those from implementing the efficient allocation. If the aggregate markup is higher, the uniform subsidy delivers larger increases because the economy is more distorted to begin with. The uniform subsidy delivers less of an increase to aggregate productivity Z and the mass of firms N because these reflect the continued presence of misallocation and a distorted entry margin. Notice that as we increase \mathcal{M} , not only are the welfare gains from the uniform subsidy larger, they are also larger as a share of the total gains. For example, if $\mathcal{M} = 1.05$ the uniform subsidy accounts for about one-half of the total welfare gains (0.65% out of 1.34%), rising to nearly three-quarters of the total welfare gains if $\mathcal{M} = 1.35$ (37.41% out of 49.66%).

Size-dependent subsidies. Setting $a_0 = 1/(1+\chi)$ and $a_1 = -1$ implements *size-dependent subsidies* that drives the wedge between price and marginal cost down to $\mu_t(z)/(1+\chi) = 1$ for each firm but leaves the aggregate wedge in the optimality conditions of the representative firm equal to $1 + \chi$. Thus by setting $\chi = \mathcal{M} - 1$ for the initial distorted steady state we can put in motion a transition to a new steady state where the the aggregate wedge remains \mathcal{M} but where the marginal revenue product of factors are equated across firms, i.e., a new steady state where there is no misallocation, and where the entry distortion is partly offset.

Table 4 shows that such subsidies have a more modest impact than the uniform subsidy. If the aggregate markup is low, $\mathcal{M} = 1.05$, these size-dependent subsidies increase gross output by 1.4%, consumption by 1.7%, and employment by 0.3%, noticeably less than the impact of the uniform subsidy. Where these policies have more success is on aggregate productivity Z which now increases by 0.7% when misallocation is eliminated as opposed to the 0.2% gain from the uniform subsidy driven by love-of-variety effects. If the aggregate markup is higher, the amount of markup dispersion in the benchmark economy is larger and so the level

of misallocation is also higher. In terms of the *share* of the total gains, the size-dependent subsidies account for about one-half if $\mathcal{M} = 1.05$ (0.71% out of 1.34%), falling to about one-quarter if $\mathcal{M} = 1.35$ (11.32% out of 49.66%).

The direct intervention $T_t(q)$ eliminates all markup distortions when both the uniform subsidy component and the size-dependent component are switched on. If only one or other of these components is switched on, entry generally remains distorted as well. We next evaluate the extent to which indirect interventions in the product market, such as those which encourage entry and competition, can reduce markup distortions.

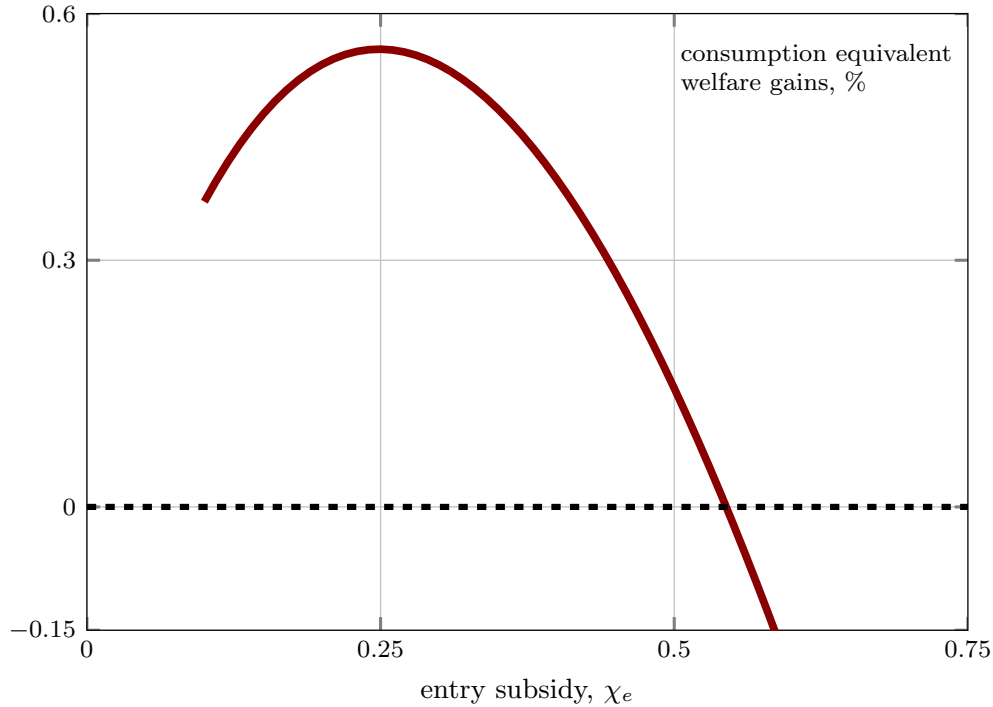
5.4 Subsidizing entry

A policy intervention like $T_t(q)$ reduces markup distortions *directly*, i.e., markups act like a tax on production so subsidizing production reduces the distortion. We now contrast such direct policies with a more *indirect* policy for reducing markup distortions — subsidizing entry, to increase the amount of competition.

Optimal entry subsidy. Consider the introduction of uniform entry subsidy χ_e that reduces the sunk entry cost from κ to $\kappa/(1+\chi_e)$. In [Table 4](#) we report the impact of the *optimal* entry subsidy that delivers the largest total welfare gain. If the aggregate markup is low, $\mathcal{M} = 1.05$, we find that the optimal entry subsidy delivers a relatively large 11.3% increase in the mass of firms N but has a more modest effect on economic activity, increasing gross output by 1.1%, consumption by 1.3%, and employment by 0.4%. Aggregate productivity increases by 0.6%, reflecting the increase in variety. But these increases in activity do not lead to substantial welfare gains, due to the cost of creating new varieties incurred during the transition. The gains from the optimal entry subsidy are 0.06%, about one-twentieth of the total gains available (0.06% out of 1.34%). If the aggregate markup is higher, say $\mathcal{M} = 1.15$, the optimal entry subsidy delivers a 20% increase in the mass of firms N but still entry only accounts for just over one-twentieth of the total gains (0.56% out of 8.67%). If $\mathcal{M} = 1.35$, the optimal entry subsidy delivers a 39% increase in the mass of firms N but still entry accounts for only about one-tenth of the total gains (5.11% out of 49.66%).

Why are the gains from subsidizing entry so low? The gains from entry are low because increasing the number of firms has tiny effects on both the aggregate markup and on misallocation. In this sense, subsidizing entry is too blunt a tool to deal with product market distortions. For example, if the benchmark economy has $\mathcal{M} = 1.05$ the optimal entry subsidy delivers an 11.3% increase in the mass of firms N but the aggregate markup falls by only about 0.02% to $\mathcal{M} = 1.0498$. Similarly if the benchmark economy has $\mathcal{M} = 1.15$, the optimal entry subsidy delivers a 20% increase in the mass of firms N but the aggregate markup hardly

Figure 4: Optimal Entry Subsidy, $\mathcal{M} = 1.15$



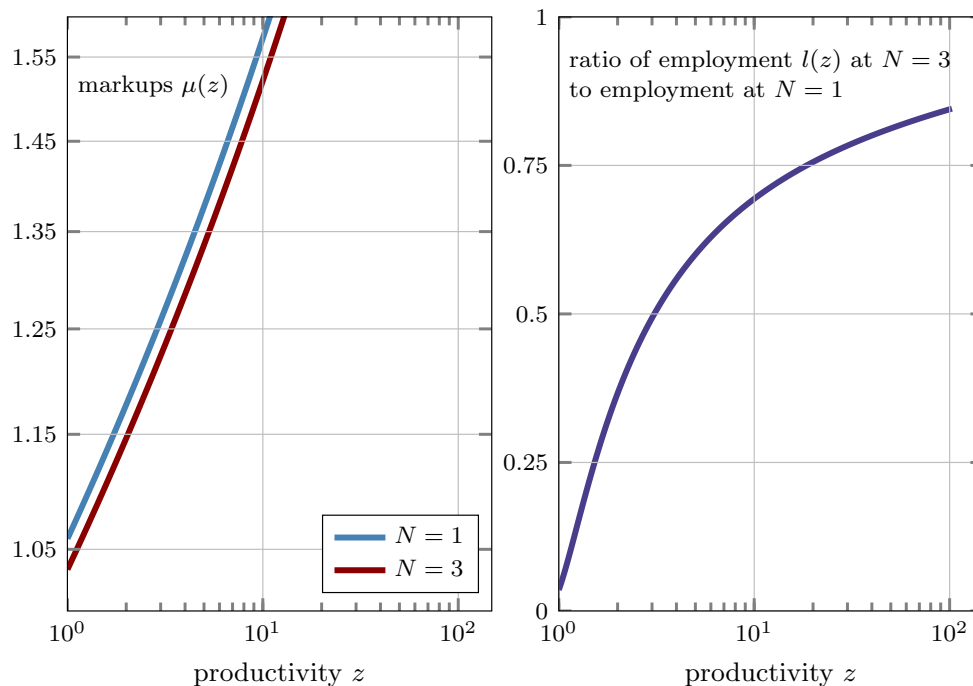
Consumption equivalent welfare gains (including transitional dynamics) as a function of the entry subsidy χ_e for aggregate markup $\mathcal{M} = 1.15$. The welfare gains for entry subsidies reported in Table 4 are for the *optimal* entry subsidies, i.e., for the peak of such curves for each \mathcal{M} . In our benchmark calibration there is insufficient entry in the initial distorted steady state so the optimal entry subsidy is positive. But entry subsidies that are too large lead to welfare losses.

changes, falling to $\mathcal{M} = 1.149$. Entry subsidies do deliver increases in aggregate productivity, but these are due to love-of-variety effects, not due to a reduction in misallocation.

The result that *more competition* does not decrease the aggregate markup may appear counterintuitive but is, in fact, a robust result in a large class of models in the international trade literature which have shown that the removal of trade costs (which subjects domestic producers to more competition) leaves the markup distribution unchanged.²² To understand this result, recall that the aggregate markup is a cost-weighted average of firm-level markups. An increase in the number of firms has two effects on this weighted average. The direct effect is a reduction in the relative size q and hence a reduction in the markups $\mu(q)$ of each firm. But there is also an important compositional effect. Recall that in our model, small firms face more elastic demand. This makes them more vulnerable to competition from entrants. By contrast large firms face less elastic demand and are less vulnerable to competition from entrants. An entry subsidy that increases the number of firms causes small, low markup

²²See Bernard, Eaton, Jensen and Kortum (2003) and Arkolakis, Costinot, Donaldson and Rodríguez-Clare (2019) who show that the markup distribution is invariant to changes in trade costs in models where variable markups arise due to limit pricing and monopolistic competition with non-CES demand, respectively.

Figure 5: Effect of Entry Subsidy on Markups, $\mathcal{M} = 1.15$



The left panel shows steady-state markups $\mu(z)$ for an economy with mass of firms $N = 1$ and an entry subsidy chosen to triple the mass of firms to $N = 3$. The right panel shows the ratio of employment $l(z)$ at $N = 3$ to employment at $N = 1$. Small, low markup firms contract by more than large, high markup firms so that high markup firms get relatively more weight in the aggregate markup calculation. Because of this, the aggregate markup hardly changes. In this example, the aggregate markup barely changes, from $\mathcal{M} = 1.150$ to $\mathcal{M} = 1.146$, even though the mass of firms triples.

firms to contract by more than large, high markup firms and the resulting reallocation means high markup firms get relatively more weight in the aggregate markup calculation. In our model, this offsetting compositional effect is almost exactly as large as the direct effect so that overall the aggregate markup falls by a negligible amount. We develop this argument more formally in [Appendix F](#).

We illustrate the two offsetting effects in [Figure 5](#). For visual clarity, we consider an extreme parameterization in which we make the entry subsidy large enough to *triple* the number of firms. Notice in the left panel that markups fall for all firms when the number of firms increases. But the right panel shows that the largest, most productive firms shrink by much less than the smallest, least productive firms. We show below that similar results are obtained with other market structures.

5.5 Monopolistic competition extensions

We now consider two variations on our benchmark model: (i) where we retain Kimball demand but where firm heterogeneity arises from differences in *quality* (demand shifters) rather

than differences in productivity, and (ii) where we replace Kimball demand with symmetric *translog demand*. For both these variations we retain the assumption of monopolistic competition. We present results for our model with oligopolistic competition and a finite number of firms per sector in the following section.

5.5.1 Heterogeneity in quality

In our benchmark model, markups are pinned down entirely by market shares. We now consider an extension where differences in quality imply differences in demand schedules across firms, breaking the tight link between markups and market shares in our benchmark.

Setup. Let $z \sim G(z)$ denote the *quality* of a firm's product and write the Kimball aggregator

$$N_t \int z \Upsilon\left(\frac{y_t(z)}{Y_t}\right) dG(z) = 1 \quad (73)$$

Following the same steps as in our benchmark model, as shown in [Appendix E](#), this leads to a relationship between markups and market shares of the form

$$\frac{1}{\mu_t(z)} + \log\left(1 - \frac{1}{\mu_t(z)}\right) = a + b \log \omega_t(z) - b \log z, \quad b = \frac{\varepsilon}{\bar{\sigma}} \quad (74)$$

Unlike our benchmark model, cross-sectional variation in market shares is no longer a sufficient statistic for the effect of variation in z . In our benchmark, we interpreted the estimated \hat{b} as a direct estimate of $\varepsilon/\bar{\sigma}$. But in this extension, since the market share is negatively correlated with the empirically *unobserved* quality z , the linear regression coefficient is no longer a consistent estimate of $\varepsilon/\bar{\sigma}$. In recalibrating the model, we use indirect inference to pin down $\varepsilon/\bar{\sigma}$, increasing the value of $\varepsilon/\bar{\sigma}$ until the coefficient in the model b equals its counterpart in the data, $\hat{b} = 0.162$, jointly with our other calibration targets.

Results. For brevity we focus on the case of $\mathcal{M} = 1.15$. As shown in [Appendix E](#), the quality model fits the data just as well as our benchmark. The most important difference is that the super-elasticity needs to be substantially higher, $\varepsilon/\bar{\sigma} = 0.304$ as opposed to 0.162.²³ Given the substantially higher super-elasticity, $\varepsilon/\bar{\sigma} = 0.304$, the quality model implies more markup dispersion, especially in the upper tail. This leads to larger losses from misallocation. Because of this, the total welfare costs are larger than in our benchmark and the gains from size-dependent policies that eliminate misallocation and the entry distortion are both larger in absolute terms and larger as a share of the total than in our benchmark. That said, we continue to find that a uniform output subsidy alone can go more than half way to achieving full efficiency. As in our benchmark, the gains from the optimal entry subsidy are still much smaller than the gains from other policies.

²³This higher super-elasticity is almost exactly what we find in an alternative parameterization where we infer the super-elasticity from a log-linear approximation to (59). In this alternative log-linear specification the quality effect would be absorbed by firm fixed effects, see [Appendix C](#) for details.

5.5.2 Translog demand

We now consider a version of our model where we replace Kimball demand with symmetric *translog* demand as in [Feenstra \(2003\)](#). For this version of the model we revert to our benchmark setting where firm heterogeneity arises from differences in productivity.

Setup. Let the technology for final good producers be given by a symmetric translog expenditure (cost) function which we write

$$\begin{aligned} \log(P_t Y_t) &= \log Y_t + \frac{1}{2\bar{\sigma}N_t} + \int \log p_t(z) dG(z) \\ &+ \frac{\bar{\sigma}N_t}{2} \left(\left(\int \log p_t(z) dG(z) \right)^2 - \int \log p_t(z)^2 dG(z) \right) \end{aligned} \quad (75)$$

Markups and market shares. As shown in [Appendix E](#), the symmetric translog specification implies that the markup $\mu_t(z)$ of a firm with productivity z solves the static condition

$$\mu + \log \mu = 1 + \log \left(\frac{z}{z_t^*} \right), \quad z > z_t^* \quad (76)$$

where z_t^* is an endogenous productivity cutoff such that firms with $z < z_t^*$ have zero sales. Moreover the translog specification implies that there is a *linear* relationship between markups and market shares

$$\mu_t(z) = 1 + \frac{1}{\bar{\sigma}} \omega_t(z) \quad (77)$$

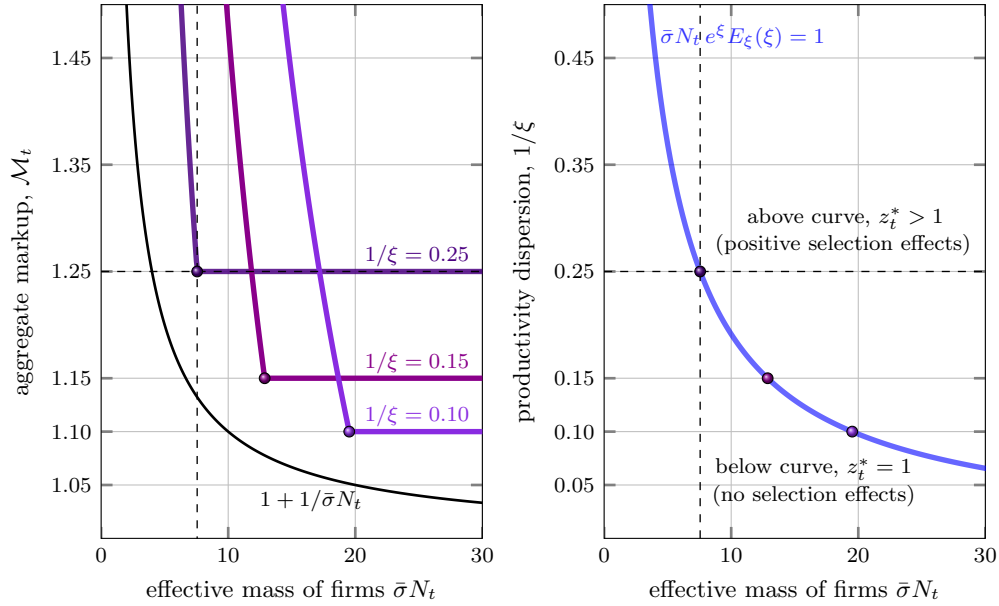
As in our benchmark model, firms with higher market shares have higher markups. With translog demand, the strength of this relationship is governed by $1/\bar{\sigma}$. The productivity cutoff z_t^* is the *only* aggregate variable that matters for the cross-sectional distribution of markups — and hence the only aggregate variable that matters for the the cross-sectional distributions of market shares $\omega_t(z)$.

To this point, our characterization of the translog model has restated standard results in the trade literature, familiar from [Feenstra \(2003\)](#), [Rodríguez-Lopez \(2011\)](#), and [Arkolakis, Costinot, Donaldson and Rodríguez-Clare \(2010, 2019\)](#) among others. We next show that given a Pareto distribution of firm-level productivity $G(z)$ we can solve explicitly for the cutoff productivity z_t^* and then aggregate markup \mathcal{M}_t . Though closely related to these existing papers, to the best of our knowledge, the following results are novel and may be of some independent interest to researchers working with translog demand and Pareto distributions.

Solving for the cutoff z_t^* . As shown in [Appendix F](#), the cutoff z_t^* is given by

$$z_t^* = \max \left[1, \bar{\sigma}N_t e^{\xi} E_{\xi}(\xi) \right]^{1/\xi} \quad (78)$$

Figure 6: Aggregate Markup with Translog Demand



Left panel shows the solution for the aggregate markup \mathcal{M}_t with translog demand as a function of the effective mass of firms $\bar{\sigma}N_t$ for various levels of the Pareto tail ξ . Right panel shows how the parameter space is partitioned into regions where there are positive selection effects, $z_t^* > 1$, or no selection effects, $z_t^* = 1$. Whenever there are positive selection effects, the aggregate markup is constant at $\mathcal{M}_t = 1 + 1/\xi$. By contrast with identical firms, the aggregate markup would be given by $1 + 1/\bar{\sigma}N_t$ as shown. With firm heterogeneity, the aggregate markup is decreasing in $\bar{\sigma}N_t$ only if there are no selection effects, $z_t^* = 1$.

where $E_n(x) := \int_1^\infty t^{-n} e^{-xt} dt$ denotes the generalized exponential integral. Since the mass of firms N_t is a state variable (is predetermined), this determines z_t^* and from (76) we then know the entire distribution of markups, market shares and relative prices given N_t . The constant $e^\xi E_\xi(\xi)$ depends only on the Pareto tail parameter $\xi > 1$ and is strictly decreasing in ξ , i.e., increasing in productivity dispersion $1/\xi$. If either the ‘effective’ mass of firms $\bar{\sigma}N_t$ is sufficiently low or productivity dispersion $1/\xi$ is sufficiently low we have $z_t^* = 1$, meaning that there are no selection effects and all firms operate. But if either $\bar{\sigma}N_t$ is sufficiently high or productivity dispersion $1/\xi$ is sufficiently high we have $z_t^* > 1$, meaning that there are positive selection effects. Intuitively, when demand is more elastic, or when the mass of firms is larger, or when productivity is more dispersed, there is more competitive pressure and selection effects are stronger, increasing the cutoff z_t^* . The right panel of Figure 6 illustrates, showing how the locus $\bar{\sigma}N_t e^\xi E_\xi(\xi) = 1$ partitions the parameter space into the regions where $z_t^* = 1$ (below the curve) and $z_t^* > 1$ (above the curve).

Solving for the aggregate markup \mathcal{M}_t . Our assumption that $G(z)$ is Pareto also implies a simple solution for \mathcal{M}_t . As shown in Appendix F, using the fact that the aggregate markup can be written as a harmonic weighted average of firm-level markups, the linear relationship between market shares and markups (77), the static markup condition (76), and our solution

for the cutoff productivity z_t^* , the aggregate markup \mathcal{M}_t is given by

$$\mathcal{M}_t = \left(1 + \frac{1}{\xi}\right) \times \left(\max \left[1, \bar{\sigma} N_t e^\xi E_\xi(\xi) \right]\right)^{-1} \quad (79)$$

Since the mass of firms N_t is a state variable, this determines \mathcal{M}_t . Now observe from (78) that if $\bar{\sigma} N_t e^\xi E_\xi(\xi) \leq 1$, implying $z_t^* = 1$, then the aggregate markup is strictly decreasing in N_t with an elasticity of -1 . But whenever $\bar{\sigma} N_t e^\xi E_\xi(\xi) > 1$, i.e., whenever there are positive selection effects, $z_t^* > 1$, then the aggregate markup is *constant* at the specific value

$$\mathcal{M}_t = 1 + \frac{1}{\xi}, \quad \text{whenever } z_t^* > 1 \quad (80)$$

So whenever there are positive selection effects, $z_t^* > 1$, e.g., $\bar{\sigma}$ or productivity dispersion $1/\xi$ is sufficiently high, then increases in the mass of firms N_t have *no effect* on the aggregate markup \mathcal{M}_t . Instead, increases in N_t are absorbed by increases in the cutoff z_t^* , i.e., by stronger selection effects. This analytic result reinforces the lesson from our benchmark model with Kimball demand where we found numerically that the aggregate markup is extremely insensitive to changes in N_t .²⁴ The reason is the same: whenever $z_t^* > 1$, an increase in N_t increases z_t^* thereby directly reducing all firm-level markups $\mu_t(z)$ according to (76). But low markup firms contract by more than large, high markup firms and the resulting reallocation means high markup firms get relatively more weight in the aggregate markup calculation. In the translog case, so long as parameters are such that $z_t^* > 1$, this offsetting compositional effect is exactly as large as the direct effect so that overall the aggregate markup is unchanged.

Role of heterogeneity. Firm heterogeneity is essential to this result. If by contrast all firms were identical, as in say Bilbiie, Ghironi and Melitz (2008, 2019), each firm would have market share $1/N_t$ and the aggregate markup would be $\mathcal{M}_t = 1 + 1/\bar{\sigma} N_t$ and would always be decreasing in N_t . In the representative firm setting, there is the direct effect of an increase in N_t on firm-level markups but this effect is the same for all firms so there is no offsetting compositional effect. In this sense, accounting for the role of firm heterogeneity is crucial for understanding the welfare effects of changes in the mass of firms N_t .²⁵

²⁴As discussed in Appendix F, the model with Kimball demand is qualitatively similar to translog demand in that for Kimball demand the aggregate markup \mathcal{M}_t is also invariant to N_t if there are positive selection effects. But in our benchmark calibration of the Kimball model, there are no selection effects and changes in N_t do change \mathcal{M}_t albeit by negligible amounts.

²⁵Rodriguez-Lopez (2011) derives a related result, solving for the *average* markup $\int \mu_t(z) dG(z)$ with translog demand and Pareto productivity and shows that this depends only on the Pareto tail ξ . Also related, Arkolakis, Costinot, Donaldson and Rodriguez-Clare (2019) show that with translog demand and Pareto productivity the univariate distribution of markups $\text{Prob}[\mu' \leq \mu]$ depends only on the Pareto tail ξ . Our key analytic contribution is to explicitly compute the *aggregate* markup, the sales-weighted harmonic average $\mathcal{M}_t = (N_t \int (\omega_t(z)/\mu_t(z)) dG(z))^{-1}$, which, as we have stressed throughout, is the key wedge in the optimality conditions of the representative firm.

Quantitative results. As discussed in [Appendix E](#), the translog model does less well in reproducing our calibration targets. As with the quality differences model, the translog model implies considerably more markup dispersion, especially in the upper tail. This leads to larger losses from misallocation relative to our benchmark model. Because of the larger amount of misallocation in the initial distorted steady state, the total welfare costs are larger than in our benchmark and the gains from size-dependent policies that eliminate misallocation and the entry distortion are both larger in absolute terms and larger as a share of the total than in our benchmark. Again we find that the gains from the optimal entry subsidy are much, much smaller than the gains from other policies.

The extensions consider in this section show that, overall, our benchmark results are robust to different monopolistically competitive setups. But one might reasonably suspect that this has more to do with the assumption of monopolistic competition than the specific aggregator we use. Perhaps a fundamentally different market structure will lead to much larger losses from markups? To assess this, we now turn to an alternative model featuring *oligopolistic competition* with genuine strategic interactions between firms.

6 Oligopolistic competition

How much does the assumed market structure matter? To assess this, we now present calculations based on an alternative model featuring oligopolistic competition rather than monopolistic competition as used in our benchmark. Our aggregation results hold regardless of the market structure, but we will see that the oligopoly model has richer empirical content and makes a number of predictions that differ from the monopolistic competition benchmark. In particular, we find larger amounts of misallocation and hence larger gains from size-dependent subsidies than in our benchmark model.

Setup. Let there be $n_t(s) \in \mathbb{N}$ firms per sector with IID productivity draws $z_i(s) \sim G(z)$. Let the within-sector aggregator be $\Upsilon(q) = q^{\frac{\gamma-1}{\gamma}}$ for $\gamma > \eta > 1$ so that the model has the nested-CES structure used by [Atkeson and Burstein \(2008\)](#) and [Edmond, Midrigan and Xu \(2015\)](#). For our quantitative work we assume *Cournot competition* so that, as in [\(38\)](#) above, the demand elasticity of a firm is given by the sales-weighted harmonic average

$$\sigma_{it}(s) = \left(\frac{1}{\eta} \omega_{it}(s) + \frac{1}{\gamma} (1 - \omega_{it}(s)) \right)^{-1} \quad (81)$$

where $\omega_{it}(s) = q_{it}(s)^{\frac{\gamma-1}{\gamma}}$ denotes the market share of firm i in sector s .²⁶ As stressed at length above, this oligopoly model is encompassed by our general framework except that for

²⁶In this oligopoly model, sectors are ex post heterogeneous so we put back dependence on s in the notation.

the free-entry condition (40) expected profits are given by (43). In practice however, solving this oligopoly model with a forward-looking free-entry condition endogenously determining the number of firms is challenging.²⁷ This is because there are many firms that each have non-negligible effects on sector-level outcomes — outcomes in any given sector are a function of the vector $\mathbf{z}(s) = (z_1(s), z_2(s), \dots, z_{n_t(s)}(s))$ of productivities. Because of the finite number of firms, sectors are heterogeneous and we cannot invoke the law of large numbers to compute expected profits. And because $n_t(s)$ is typically large, we need to compute high-dimensional integrals with respect to the joint distribution $G_{n_t(s)}(\mathbf{z}(s))$ of $\mathbf{z}(s)$. In principle, the heterogeneity across sectors creates incentives for firms to *direct* entry towards more profitable sectors. But to simplify the problem computationally, we assume that entry is *random*, that firms can not direct entry in this way. We discuss these issues in more detail in [Appendix D](#).

Relationship between markups and market shares. This nested-CES specification implies that the inverse markup is linear decreasing in the market share

$$\frac{1}{\mu_{it}(s)} = \left(1 - \frac{1}{\gamma}\right) - \left(\frac{1}{\eta} - \frac{1}{\gamma}\right) \omega_{it}(s) \quad (82)$$

As in our benchmark model, firms with higher market shares have higher markups. Here, the strength of this relationship is governed by the gap between the between-sector elasticity of substitution η and the within-sector elasticity of substitution $\gamma > \eta$. Multiplying both sides of (82) by $\omega_{it}(s)$ and summing over all firms i within sector s gives

$$\frac{1}{\mu_t(s)} = \left(1 - \frac{1}{\gamma}\right) - \left(\frac{1}{\eta} - \frac{1}{\gamma}\right) \sum_{i=1}^{n_t(s)} \omega_{it}(s)^2 \quad (83)$$

The model predicts a linear decreasing relationship between the sector-level inverse markup $1/\mu_t(s)$ and the sector’s *Herfindahl-Hirschman index* (HHI) of sales concentration. From (25), the sector-level labor share is proportional to the inverse markup, $W_{it}(s)/p_t(s)y_t(s) = (1-\alpha)\zeta_t/\mu_t(s)$. Motivated by this, in calibrating the oligopoly model we use indirect inference to pin down the gap between γ and η , choosing parameters so that our model reproduces the $\hat{b} = -0.21$ slope coefficient in a regression of the change over time of sector-level labor shares on the change in sector-level HHIs, as in [Autor, Dorn, Katz, Patterson and Van Reenen \(2020\)](#), jointly with our other calibration targets.

Calibration. The oligopoly model features both within- and between-sector variation in concentration. We calibrate the oligopoly model targeting measures of concentration within 4-digit sectors in the 2012 US Census of Manufactures as reported by [Autor, Dorn, Katz, Patterson and Van Reenen \(2020\)](#). In particular, we target their top 4 sales share (CR4)

²⁷Other applications of this oligopoly setup, e.g., [Atkeson and Burstein \(2008\)](#), [Edmond, Midrigan and Xu \(2015\)](#), and [De Loecker, Eeckhout and Monge \(2021\)](#), treat the number of potential producers as *exogenous*.

Table 5: Parameterization, Oligopoly

<i>calibration targets</i>		<i>data</i>				
\mathcal{M}	aggregate markup	1.1 ~ 1.4	1.05	1.15	1.25	1.35
CR4	top 4 sales share	0.43	0.37	0.43	0.43	0.43
CR20	top 20 sales share	0.72	0.76	0.72	0.72	0.72
	materials share	0.45	0.45	0.45	0.45	0.45
\hat{b}	regression coefficient	-0.21	-0.21	-0.21	-0.21	-0.21
<i>parameter values</i>						
ξ	Pareto tail		28.08	8.51	5.15	3.72
γ	elasticity of substitution within sectors		59.69	12.76	7.16	5.21
η	elasticity of substitution between sectors		1.62	1.35	1.15	0.99
N	average number of firms per sector		415	359	143	112
ϕ	weight on value-added		0.70	0.58	0.46	0.30

The calibrated parameters for our oligopoly model. We calibrate the Pareto tail ξ , the within- and between-sector elasticities of substitution γ and η , the sunk entry cost κ , and weight on value-added ϕ to match the targets shown. In practice, we choose the average number of firms N per sector and back out the sunk cost κ that rationalizes N . The cross-sectional regression is of the change over time in sector-level labor shares on the change in sector-level HHIs, as discussed in the text. All other parameters are assigned as in Panel A of [Table 1](#).

of 0.43 and top 20 sales share (CR20) of 0.73. We also target the slope in a regression of the change over time in sector-level labor shares (inverse markups) on the change in sector-level HHIs of $\hat{b} = -0.21$, i.e., we also target the sectoral relationship between markups and concentration.²⁸ As in our benchmark model we target a materials share of 0.45 and consider a range of targets for the aggregate markup \mathcal{M} . Intuitively, the two measures of sales concentration pin down the Pareto tail ξ , which controls the amount of productivity dispersion, and the sunk entry cost κ . The aggregate markup then pins down γ , while the slope coefficient pins down the gap between γ and η . As shown in [Table 5](#), the oligopoly model hits all our calibration targets except when the target for the aggregate markup is low, $\mathcal{M} = 1.05$. For low levels of the aggregate markup, the oligopoly model struggles to reproduce the top 4 sales concentration in the data. For any given \mathcal{M} , the oligopoly model requires less productivity dispersion than the benchmark model with Kimball demand and monopolistic competition. For example, with $\mathcal{M} = 1.15$ the oligopoly model requires Pareto tail $\xi = 8.51$ as opposed to $\xi = 6.84$ in the benchmark model. On average, there is a relatively large number of firms per sector, $N = 359$, but most of these firms are very small.

²⁸The CR4 and CR20 are reported in Panel A of Figure 4 while the regression coefficient \hat{b} is from Table 2 baseline column 3 in [Autor, Dorn, Katz, Patterson and Van Reenen \(2020\)](#).

Table 6: Markup Dispersion and Productivity Losses, Oligopoly

<i>cost-weighted distribution of markups</i>								
aggregate markup, \mathcal{M}	1.05		1.15		1.25		1.35	
	$\mu_t(s)$	$\mu_{it}(s)$	$\mu_t(s)$	$\mu_{it}(s)$	$\mu_t(s)$	$\mu_{it}(s)$	$\mu_t(s)$	$\mu_{it}(s)$
p25 markup	1.04	1.02	1.12	1.09	1.21	1.17	1.29	1.25
p50 markup	1.05	1.04	1.14	1.11	1.23	1.19	1.32	1.27
p75 markup	1.06	1.07	1.16	1.17	1.27	1.26	1.37	1.36
p90 markup	1.07	1.10	1.21	1.27	1.33	1.41	1.45	1.54
p99 markup	1.11	1.20	1.35	1.57	1.57	1.90	1.78	2.24
<i>aggregate productivity losses, %</i>								
gross output	2.99		3.19		3.32		3.81	
value-added	5.55		6.85		8.89		12.52	
value-added, $\mathcal{M} = 1$	5.45		6.02		6.51		7.74	

Cost-weighted steady state distribution of firm-level markups $\mu_{it}(s)$ and sector-level markups $\mu_t(s)$ and the implied aggregate productivity losses for our oligopoly model. Gross output aggregate productivity loss is $(Z - Z^*)/Z^* \times 100$, and similarly for the value-added aggregate productivity loss. To isolate the effect of misallocation on value-added aggregate productivity we also report the value-added aggregate productivity loss with the same amount of markup dispersion but holding $\mathcal{M} = 1$ to eliminate the distortion between value-added and materials, see text for details.

Results. Table 6 reports the cost-weighted steady-state distribution of firm-level markups $\mu_{it}(s)$ and sector-level markups $\mu_t(s)$ for four levels of the aggregate markup \mathcal{M} . The distribution of sector-level markups alone is as dispersed as the unconditional markup distribution in our benchmark model with monopolistic competition calibrated to the same aggregate markup \mathcal{M} (for which sectors are identical).²⁹ For brevity we focus on the case of $\mathcal{M} = 1.15$. The unconditional distribution of markups in the oligopoly model is considerably more dispersed than in our benchmark, especially in the upper tail. The gross output losses from misallocation are 3.19%, up from 0.97% in the benchmark.

As reported in Table 7, in many respects the oligopoly model implies similar long-run changes in economic activity as the benchmark model calibrated to the same \mathcal{M} . For example, for $\mathcal{M} = 1.15$ our benchmark model implies an output increase of 59.6%, consumption increase of 44.5%, and employment increase of 18.0%. For $\mathcal{M} = 1.15$ the oligopoly model implies an output increase of 55.9%, consumption increase of 39.9%, and employment increase of 14.9%. One notable difference however is that in our oligopoly model the initial distorted

²⁹This amount of dispersion in sector-level markups is however less costly, because of the low elasticity of substitution η between sectors. The amount of markup dispersion within sectors is more important.

Table 7: Implications of Alternative Policies, Oligopoly

	steady state comparisons, %						welfare, %
	Y	C	L	N	K	Z	
<i>oligopoly, $\mathcal{M} = 1.05$</i>							
efficient	20.1	15.9	4.8	-29.5	30.3	1.8	8.71
uniform subsidy	14.2	9.4	6.0	3.7	23.1	0.1	0.58
size-dependent subsidy	5.3	6.1	-1.1	-31.4	6.1	1.7	7.97
entry subsidy	-2.0	-2.4	-1.1	-28.9	-2.5	-1.2	0.50
<i>oligopoly, $\mathcal{M} = 1.15$</i>							
efficient	55.9	39.9	14.9	-10.5	94.0	1.9	14.66
uniform subsidy	50.5	34.4	17.0	9.8	86.6	1.1	5.14
size-dependent subsidy	4.1	4.4	-1.8	-18.4	4.5	0.8	8.70
entry subsidy	-2.1	-2.5	-1.0	-8.7	-2.7	-1.1	0.12
<i>oligopoly, $\mathcal{M} = 1.25$</i>							
efficient	112.6	79.0	25.3	-1.7	206.6	3.1	26.76
uniform subsidy	108.4	75.2	28.4	7.9	201.1	3.0	15.20
size-dependent subsidy	3.0	3.0	-2.4	-14.3	3.1	0.2	9.34
entry subsidy	0.9	1.0	0.4	1.9	1.2	0.4	0.01
<i>oligopoly, $\mathcal{M} = 1.35$</i>							
efficient	204.0	142.7	35.3	3.6	412.1	5.0	48.63
uniform subsidy	201.8	141.0	39.6	20.3	411.9	5.6	32.38
size-dependent subsidy	2.8	2.7	-3.1	-12.7	2.7	-0.1	11.28
entry subsidy	8.3	9.4	3.4	11.1	11.2	3.2	0.44

The first six columns report the percentage change from the initial distorted steady state with to the new steady state. The last column reports the consumption equivalent welfare gains (including transitional dynamics). The alternative policies are (i) the *efficient allocation*, where all markups are removed, (ii) a *uniform subsidy* that eliminates the aggregate markup, (iii) *size-dependent subsidies* that eliminate misallocation and the entry distortion, and (iv) the optimal *entry subsidy* (or tax).

steady state often features *too many* firms, for $\mathcal{M} = 1.15$ the efficient steady state involves reducing the average number of firms N by about 10.5%. In any case, because of the larger amount of misallocation, the oligopoly model implies substantially larger costs of markups, 14.66% in consumption-equivalent terms, up from 8.67% for our benchmark. The gains from size-dependent subsidies that eliminate misallocation and the entry distortion are 8.70% for the oligopoly model, up from 2.87% for our benchmark. The gains from a uniform subsidy that eliminates the aggregate markup distortion are similar to our benchmark, 5.14% down slightly from 5.90%, but are correspondingly a smaller share of the total. Again, the gains from the optimal entry subsidy are much, much smaller than the gains from other policies.³⁰

There are two important caveats regarding these results. First, in the oligopoly model, subsidies to eliminate misallocation would have to be both sector- and size-dependent, as opposed to just size-dependent as they are in our benchmark model with monopolistic competition. Second, the losses from misallocation may be lower if entry could be directed to specific sectors. It remains an open question and an important direction for future research to assess how much misallocation would be reduced if firms could direct entry.

7 Conclusion

We study the welfare costs of product market distortions in a dynamic model with heterogeneous firms and endogenously variable markups. Our model encompasses several popular market structures and we provide aggregation results showing how the macro implications of micro-level markup heterogeneity can be summarized by a few key statistics. We calibrate our model to match levels of sales concentration and the firm-level relationship between labor shares and market shares observed in 6-digit US Census of Manufactures data. We find that the welfare costs of markups can be large. Depending on the market structure and assumed level of the aggregate markup, the representative consumer can gain as much as 50% in consumption-equivalent terms if all markup distortions are eliminated, once transitional dynamics are taken into account.

In our model markups reduce welfare because the aggregate markup distortion acts like a uniform output tax, reducing employment and investment by all firms, because markup variation across firms causes misallocation of factors of production, and because there is an inefficient rate of entry due to the misalignment between private and social incentives to create new firms. Across all specifications, we robustly find that the aggregate markup and misallocation channels account for the bulk of the costs of markups and that the entry channel is much less important.

Although we focus on the normative side of our model, our results also have clear empirical

³⁰Since the initial steady state has too many firms, the optimal entry subsidy is a *tax*.

implications. One simple but important finding is that the overall level of markups is best measured as a *cost-weighted* average of firm-level markups. This is the relevant ‘wedge’ in aggregate employment and investment decisions. By contrast a *sales-weighted* average of firm-level markups, as used in the empirical literature, overstates the rise in the overall level of market power. In addition, our results provide two reasons to be skeptical of explanations for the simultaneous rise in concentration and markups that focus on increasing barriers to entry. First, in our model increasing barriers to entry *reduce concentration*, because the resulting lack of competition makes it easier for small firms to survive. Second, in our model changes in entry have negligible effects on the overall level of markups because entry is associated with a reallocation of production towards high productivity, high markup firms.

To keep our model tractable enough that we can aggregate cross-sectional outcomes and study transitional dynamics for a broad range of alternative market structures, we have abstracted from a number of considerations that might play an important role in the development of a more complete account of the macroeconomic implications of product market distortions. First, while markups in our model are a return to sunk investments, there are no positive spillovers from such investment to the stock of knowledge in the economy at large and hence no implications for endogenous growth. But as emphasized by [Atkeson, Burstein and Chatzikonstantinou \(2019\)](#), in the endogenous growth models they survey, a higher markup acts like a uniform subsidy to innovation and is welfare-improving, the quantitative details depending sensitively on the specification of the technology for research. In principle, these effects could be large. That said, in endogenous growth models with variable markups, such as [Peters \(2020\)](#), the interactions between entry, aggregate innovation and misallocation are more subtle with the overall effects on growth ambiguous. An important challenge for future work in this area is to provide detailed evidence on technologies for research and the magnitudes of spillovers that can be used to refine such models to help quantify the relative importance of these growth effects and the level effects of markups emphasized in this paper.

Second, we have made the assumption, standard in the literature, that the underlying sources of firm size differences are fundamental differences in productivity or quality. Because of this, large firms with high markups represent a lost opportunity — they should be even larger, not smaller, but charge lower prices. But if large firms are large not because they are more productive or because their products are higher quality but instead because they receive special tax breaks, or have political connections that help them evade antitrust actions or other forms of regulation, then such firms may well be too large, not too small. Another important challenge for future work in this area is to build models that blend political connections, as in [Akcigit, Baslandze and Lotti \(2018\)](#), with endogenous product market distortions so that we can quantitatively evaluate size-dependent policy interventions when both fundamental and non-fundamental sources of firm size are operative.

Finally, to keep the analysis focused, we have abstracted from distortionary tax wedges

and frictions in factor markets (e.g., monopsony power) that affect aggregate employment and capital accumulation. For standard second best reasons, such distortions may either amplify or mitigate the costs of product market distortions. Quantifying the interactions between these different types of distortions also seems a natural topic for ongoing research.

Appendix

A Cost-weighted vs. sales-weighted average markups

In this appendix we derive an exact relationship between a *cost-weighted* average markup \mathcal{M} and a *sales-weighted* average markup $\tilde{\mathcal{M}}$. The key result is

$$\boxed{\frac{\tilde{\mathcal{M}} - \mathcal{M}}{\mathcal{M}} = \text{Var}[\hat{\mu}_i]}$$

where $\text{Var}[\hat{\mu}_i]$ is a measure of the cross-sectional dispersion in the idiosyncratic component in markups, $\hat{\mu}_i := \mu_i/\mathcal{M}$. This derivation makes no assumptions about demand or market structure but makes one key assumption about technology, specifically, that all firms within a given industry have the same *cost elasticity*.

Notation. Consider an industry with $i = 1, 2, \dots, n$ firms. Let p_i, y_i, μ_i and c_i denote respectively a firm's price, output, markup, and *total variable costs*.

Cost elasticity assumption. Let $\vartheta > 0$ denote a firm's *cost elasticity*, that is, the elasticity of total variable costs with respect to output

$$\vartheta := \frac{\partial \log c}{\partial \log y} = \frac{\partial c}{\partial y} \frac{y}{c} = \frac{\text{marginal cost}}{\text{average cost}} \quad (\text{A1})$$

Our key assumption is that the cost elasticity ϑ is common to all firms within a given industry, $\vartheta_i = \vartheta$. In other words, all firms within a given industry have the same *returns to scale*, but this may be either increasing, constant, or decreasing at the industry level. Marginal costs are then given by $\vartheta c_i/y_i$. Importantly we do not put any restrictions on marginal costs, these can vary arbitrarily across firms within the industry.

Aggregate markup. Given the assumption that all firms within a given industry have the same cost elasticity ϑ , it is straightforward to show that the industry aggregate markup, that is, the ratio of industry price to industry marginal cost, is given by a cost-weighted average of firm-level markups (equivalently, a sales-weighted *harmonic* average). Following the same steps as in the main text, since prices p_i are a markup μ_i over marginal cost $\vartheta c_i/y_i$ we have revenues $p_i y_i = \vartheta \mu_i c_i$ so if we are to write \mathcal{M} as the 'wedge' between industry revenue $PY := \sum_i p_i y_i$ and industry costs $\vartheta \sum_i c_i$ (i.e., so that \mathcal{M} is the ratio of the industry price level to industry marginal costs), then

$$\mathcal{M} = \sum_{i=1}^n \mu_i \omega_i, \quad \omega_i := \frac{c_i}{\sum_i c_i} \quad (\text{A2})$$

where in slight abuse of notation we now use ω_i to denote the *cost-weights*. Notice that this derivation makes no assumptions about the demand system or market structure that generates the markups μ_i .

Relationship between cost-weighted and sales-weighted averages. By contrast, the applied literature on markups has emphasized sales-weighted averages, which can be written

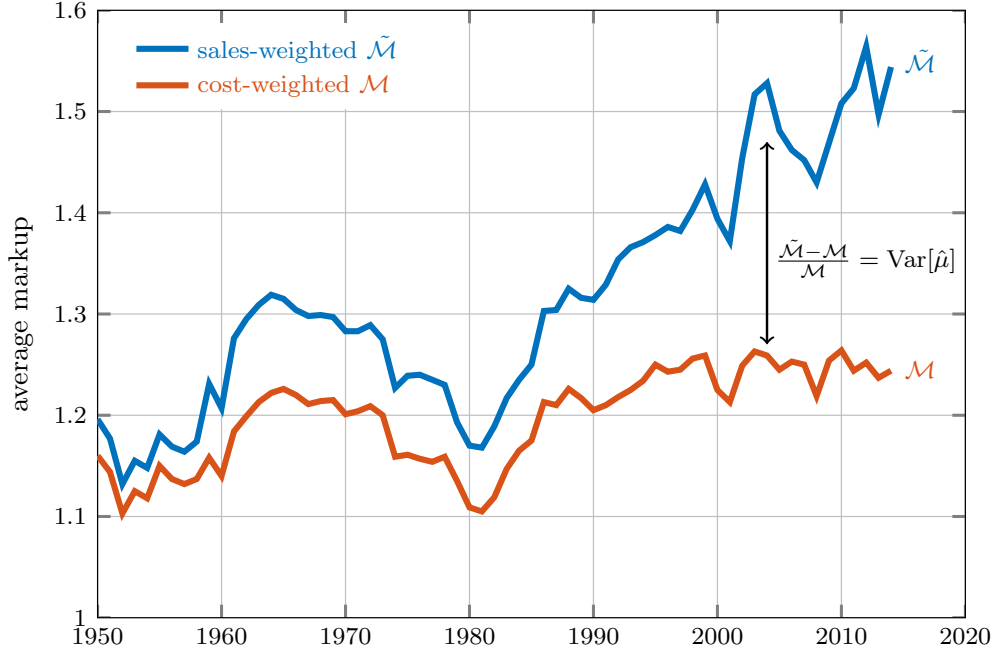
$$\tilde{\mathcal{M}} = \sum_{i=1}^n \mu_i \tilde{\omega}_i, \quad \tilde{\omega}_i := \frac{p_i y_i}{\sum_i p_i y_i} \quad (\text{A3})$$

We will now show that the sales-weighted average $\tilde{\mathcal{M}}$ can be decomposed into the cost-weighted average \mathcal{M} plus a term that reflects the cross-sectional dispersion in markups.

Let $\mathbb{E}[\cdot]$ denote averages with respect to the cost weights so that $\mathcal{M} = \mathbb{E}[\mu_i]$. Then we can write the sales-weighted average as

$$\tilde{\mathcal{M}} = \sum_{i=1}^n \mu_i \tilde{\omega}_i = \sum_{i=1}^n \mu_i \frac{\tilde{\omega}_i}{\omega_i} \omega_i = \mathbb{E}\left[\mu_i \frac{\tilde{\omega}_i}{\omega_i}\right] \quad (\text{A4})$$

Figure A1: Cost-Weighted vs. Sales-Weighted Average Markups, Compustat



The sales-weighted average $\tilde{\mathcal{M}}$ of firm-level markups in Compustat data, as in De Loecker, Eeckhout and Unger (2020), and the cost-weighted average of firm-level markups \mathcal{M} . The former is higher and has increased by a larger amount. The proportional difference between the two averages reflects the cross-sectional dispersion in markups, which has been increasing.

Expanding the expectation of the product into the covariance plus the product of the expectations then gives

$$\begin{aligned}\tilde{\mathcal{M}} &= \mathbb{E}\left[\mu_i \frac{\tilde{\omega}_i}{\omega_i}\right] = \text{Cov}\left[\mu_i, \frac{\tilde{\omega}_i}{\omega_i}\right] + \mathbb{E}[\mu_i] \mathbb{E}\left[\frac{\tilde{\omega}_i}{\omega_i}\right] \\ &= \text{Cov}\left[\mu_i, \frac{\tilde{\omega}_i}{\omega_i}\right] + \mathcal{M}\end{aligned}\tag{A5}$$

since $\mathcal{M} = \mathbb{E}[\mu_i]$ and $\mathbb{E}\left[\frac{\tilde{\omega}_i}{\omega_i}\right] = \sum_i \tilde{\omega}_i = 1$. In short, the *absolute difference* between the sales-weighted and cost-weighted average markups is given by the covariance of the markups μ_i and the relative weights $\tilde{\omega}_i/\omega_i$.

But under the assumption of a common cost elasticity ϑ the relative weights are *proportional to the markups themselves*

$$\frac{\tilde{\omega}_i}{\omega_i} = \frac{p_i y_i}{c_i} \frac{\sum_i c_i}{\sum_i p_i y_i} = \frac{\mu_i \vartheta \frac{c_i}{y_i} y_i}{c_i} \frac{\sum_i c_i}{\sum_i p_i y_i} = \frac{\mu_i}{\mathcal{M}}\tag{A6}$$

where the last equality follows because \mathcal{M} is the ‘wedge’ between industry revenue $\sum_i p_i y_i$ and industry costs $\vartheta \sum_i c_i$. In short, we can write

$$\text{Cov}\left[\mu_i, \frac{\tilde{\omega}_i}{\omega_i}\right] = \text{Cov}\left[\mu_i, \mu_i \frac{1}{\mathcal{M}}\right] = \frac{1}{\mathcal{M}} \text{Var}[\mu_i]\tag{A7}$$

And hence our key decomposition can be written

$$\tilde{\mathcal{M}} = \mathcal{M} + \frac{1}{\mathcal{M}} \text{Var}[\mu_i]\tag{A8}$$

That is, the sales-weighted average can be expressed as the cost-weighted average plus a term that reflects the cross-sectional dispersion in markups.

Multiplicative decomposition. A slightly more intuitive version of this decomposition obtains if we decompose the markups μ_i multiplicatively into the common \mathcal{M} component and an idiosyncratic component $\hat{\mu}_i$ with mean normalized to one

$$\hat{\mu}_i := \mu_i / \mathcal{M} \quad (\text{A9})$$

Then $\text{Var}[\mu_i] = \mathcal{M}^2 \text{Var}[\hat{\mu}_i]$ and we can write

$$\frac{\tilde{\mathcal{M}} - \mathcal{M}}{\mathcal{M}} = \text{Var}[\hat{\mu}_i] \quad (\text{A10})$$

That is, the percentage difference between the sales-weighted average and the cost-weighted average is given by the cross-sectional variance of the idiosyncratic component $\hat{\mu}_i$.

Hence $\tilde{\mathcal{M}} \geq \mathcal{M}$ with equality only if there is no markup dispersion. The statistic $\tilde{\mathcal{M}}$ can rise over time either due to increasing \mathcal{M} or increasing $\text{Var}[\hat{\mu}_i]$ or both. The statistic $\tilde{\mathcal{M}}$ can be rising even if \mathcal{M} is constant. Indeed $\tilde{\mathcal{M}}$ can be rising even if \mathcal{M} is falling if the increase in dispersion $\text{Var}[\hat{\mu}_i]$ is large enough.

Compustat example. To get a quantitative sense of the difference between the cost-weighted average \mathcal{M} and the sales-weighted average $\tilde{\mathcal{M}}$, we compute these statistics using publicly available Compustat data for the US economy. We follow the approach of [De Loecker, Eeckhout and Unger \(2020\)](#) using the ratio of sales to the cost of goods sold, scaled by estimates (at the 2-digit industry level) of the output elasticity of the production function from [Karabarbounis and Neiman \(2019\)](#). We show the results in [Figure A1](#).³¹ Clearly the sales weighted average $\tilde{\mathcal{M}}$ is higher and has risen by substantially more than the cost-weighted average \mathcal{M} . The additional increase in $\tilde{\mathcal{M}}$ reflects the increasing dispersion of markups.

Although researchers may not always have reliable data on total variable costs, under the assumption that all firms within a given industry share the same cost elasticity ϑ , the cost-weighted *arithmetic* average is equivalent to the sales-weighted *harmonic* average, which can of course be computed if the sales-weighted arithmetic average can.

B Census data and markup estimates

We use data from the US Census of Manufactures from 1972 to 2012. We focus on the Census of Manufactures for two reasons: (i) it has higher-quality input data relative to other sectors, such as Services, and (ii) the vast majority of manufacturing goods are easily transportable and not limited to local markets.

Framework. We now spell out the assumptions we need to infer firm-level markups from this Census data. Suppose firms face an inverse demand function and let $\sigma_{it}(s)$ and $\mu_{it}(s)$ denote the implied demand elasticity and markup

$$\sigma_{it}(s) := -\frac{\partial \log y_{it}(s)}{\partial \log p_{it}(s)}, \quad \mu_{it}(s) := \frac{\sigma_{it}(s)}{\sigma_{it}(s) - 1} \quad (\text{B1})$$

Suppose firms have production function

$$y_{it}(s) = F_s(k_{it}(s), l_{it}(s), x_{it}(s)) \quad (\text{B2})$$

and let $\alpha_{it}^k(s)$, $\alpha_{it}^l(s)$, $\alpha_{it}^x(s)$ denote the elasticities of output with respect to capital, labor, and materials

$$\alpha_{it}^k(s) := \frac{\partial \log y_{it}(s)}{\partial \log k_{it}(s)}, \quad \alpha_{it}^l(s) := \frac{\partial \log y_{it}(s)}{\partial \log l_{it}(s)}, \quad \alpha_{it}^x(s) := \frac{\partial \log y_{it}(s)}{\partial \log x_{it}(s)} \quad (\text{B3})$$

Taking factor prices as given, suppose $k_{it}(s)$, $l_{it}(s)$, $x_{it}(s)$ are chosen to maximize profits

$$p_{it}(s)y_{it}(s) - R_t k_{it}(s) - W_t l_{it}(s) - x_{it}(s) \quad (\text{B4})$$

³¹See also Figure II, Panel B in [De Loecker, Eeckhout and Unger \(2020\)](#).

subject to the inverse demand curve and production function given above. The key first order conditions for this problem can be written

$$R_t k_{it}(s) = \alpha_{it}^k(s) \frac{p_{it}(s)y_{it}(s)}{\mu_{it}(s)} \quad (\text{B5})$$

$$W_t l_{it}(s) = \alpha_{it}^l(s) \frac{p_{it}(s)y_{it}(s)}{\mu_{it}(s)} \quad (\text{B6})$$

$$x_{it}(s) = \alpha_{it}^x(s) \frac{p_{it}(s)y_{it}(s)}{\mu_{it}(s)} \quad (\text{B7})$$

which implies, for example,

$$\frac{W_t l_{it}(s)}{R_t k_{it}(s) + W_t l_{it}(s) + x_{it}(s)} = \frac{\alpha_{it}^l(s)}{\alpha_{it}^k(s) + \alpha_{it}^l(s) + \alpha_{it}^x(s)} \quad (\text{B8})$$

To infer markups from these conditions using data from the Census of Manufactures we impose two additional assumptions: (i) that each firm i within a given sector s has the same factor elasticities, i.e., for each factor $j = k, l, x$ the elasticities $\alpha_{it}^j(s) = \alpha_t^j(s)$ for all i in s , and (ii) the degree of returns to scale in each sector is the same, $\text{RTS} := \sum_j \alpha_t^j(s)$ for all s . The Census gives us the value of revenue $p_{it}(s)y_{it}(s)$ and the wage bill $W_t l_{it}(s)$ for each firm i in each 6-digit NAICS sector s . Thus if we are equipped with an estimate of the elasticity of output with respect to labor, $\hat{\alpha}_t^l(s)$ our estimated markups are

$$\hat{\mu}_{it}(s) = \frac{p_{it}(s)y_{it}(s)}{W_t l_{it}(s)} \times \hat{\alpha}_t^l(s) \quad (\text{B9})$$

Multi-establishment firms. In practice, we begin with the value of shipments $p_{eit}(s)y_{eit}(s)$ and total salaries/wages $W_t l_{eit}(s)$ for each establishment e of firm i in each 6-digit NAICS sector s . In the case of a single-establishment firm i in sector s , we have

$$\mu_{it}(s) = \mu_{eit}(s) = \frac{p_{eit}(s)y_{eit}(s)}{W_t l_{eit}(s)} \times \alpha_t^l(s) \quad (\text{B10})$$

where $\alpha_t^l(s)$ is the elasticity of output with respect to labor in sector s , as discussed above. For multi-establishment firms we aggregate over the establishments e of firm i to get

$$\mu_{it}(s) = \alpha_t^l(s) \sum_{e \in i} \mu_{eit}(s) \frac{W_t l_{eit}(s)}{\sum_{e' \in i} W_t l_{e'it}(s)} \quad (\text{B11})$$

Output elasticities. The empirical literature has proposed various strategies for recovering the output elasticities $\alpha_t^l(s)$ specific to sector s . In principle, one could estimate sector-specific production functions to recover these elasticities. However, recently [Bond, Hashemi, Kaplan and Zoch \(2021\)](#) have shown that in the presence of variable markups it is not possible to consistently estimate output elasticities when only revenue data is available. Given this, we follow an alternative approach, more in the spirit of growth accounting, where we use the firm's cost minimization conditions to write, for each establishment e and firm i

$$\alpha_t^l(s) = \frac{W_t l_{eit}(s)}{W_t l_{eit}(s) + R_t k_{eit}(s) + x_{eit}(s)} \times \text{RTS} \quad (\text{B12})$$

Because of measurement error at the establishment level, we take averages within sector s for some given RTS. Following [Foster, Grim and Haltiwanger \(2016\)](#), we take the cost-weighted average of labor input expenditure shares of establishments within each sector s . This provides us with an estimate of $\alpha_t^l(s)$ for each 6-digit NAICS sector s in each Census year t . For our benchmark results we assume constant returns to scale, $\text{RTS} = 1$. We discuss the sensitivity of our results to the RTS in [Appendix C](#).

Markup regression specification. The key to our calibration of the benchmark model with Kimball demand is the cross-sectional relationship between markups and market shares within a given sector. To see this relationship precisely, consider a version of our model with sector-specific Kimball aggregators with inverse demand curves of the form

$$p_{it}(s) = \Upsilon'_s(q_{it}(s))\gamma_i(s)d_t(s), \quad \Upsilon'_s(q) = \frac{\bar{\sigma}(s) - 1}{\bar{\sigma}(s)} \exp\left(\frac{1 - q^{\varepsilon(s)/\bar{\sigma}(s)}}{\varepsilon(s)}\right) \quad (\text{B13})$$

where $d_t(s)$ is the Kimball demand index, common to all firms i in sector s . Relative to our benchmark model, this more general setting allows for time-invariant firm-specific demand shifters $\gamma_i(s)$ and sector-specific elasticity parameters $\varepsilon(s), \bar{\sigma}(s)$. Market shares are $\omega_{it}(s) = p_{it}(s)q_{it}(s)$ so the log market share can be written

$$\log \omega_{it}(s) = \log q_{it}(s) + \frac{1 - q_{it}(s)^{\varepsilon(s)/\bar{\sigma}(s)}}{\varepsilon(s)} + \log\left(\gamma_i(s)d_t(s)\frac{\bar{\sigma}(s) - 1}{\bar{\sigma}(s)}\right) \quad (\text{B14})$$

With Kimball demand the markup $\mu_{it}(s)$ is related to relative size $q_{it}(s)$ according to

$$\frac{1}{\mu_{it}(s)} = 1 - \frac{1}{\bar{\sigma}(s)} q_{it}(s)^{\varepsilon(s)/\bar{\sigma}(s)} \quad (\text{B15})$$

We can then eliminate $q_{it}(s)$ between equations (B14)-(B15) and collect terms to get

$$\frac{1}{\mu_{it}(s)} + \log\left(1 - \frac{1}{\mu_{it}(s)}\right) = a(s) + a_i(s) + a_t(s) + b(s) \log \omega_{it}(s)$$

the same as (59) above, with fixed effects

$$a(s) = \frac{\bar{\sigma}(s) - 1}{\bar{\sigma}(s)} - \log \bar{\sigma}(s) - \frac{\varepsilon(s)}{\bar{\sigma}(s)} \log\left(\frac{\bar{\sigma}(s) - 1}{\bar{\sigma}(s)}\right) \quad (\text{B16})$$

$$a_i(s) = -\frac{\varepsilon(s)}{\bar{\sigma}(s)} \log \gamma_i(s) \quad (\text{B17})$$

$$a_t(s) = -\frac{\varepsilon(s)}{\bar{\sigma}(s)} \log d_t(s) \quad (\text{B18})$$

and slope coefficient

$$b(s) = \frac{\varepsilon(s)}{\bar{\sigma}(s)} \quad (\text{B19})$$

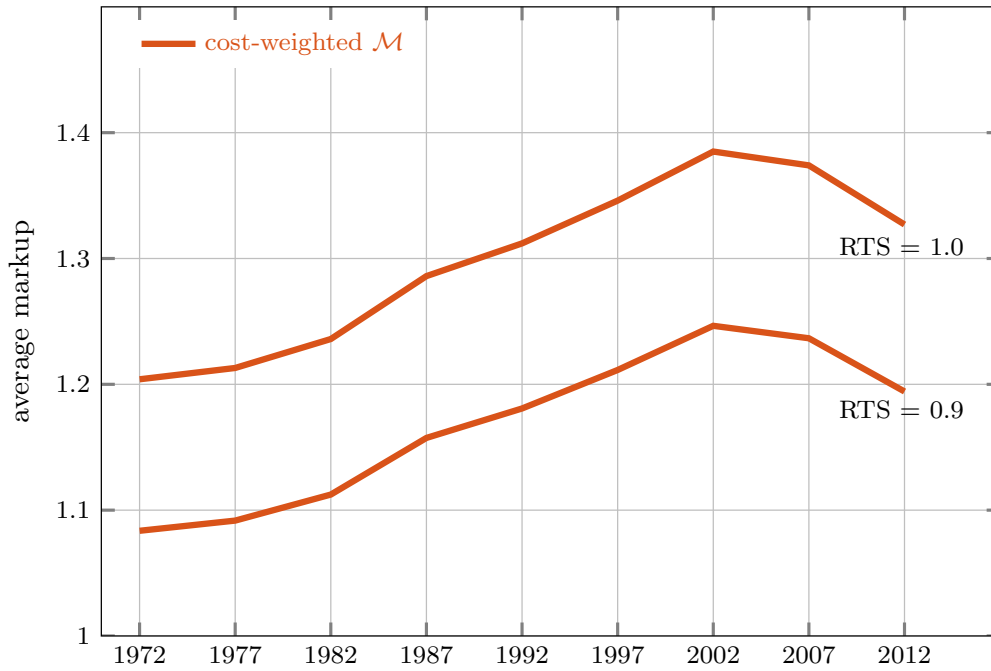
To summarize, the model then tells us that the superelasticity is pinned down by the strength of the covariation between (transformed) markups and market shares after having controlled for firm- and sector-time fixed effects. The firm effects control for time-invariant firm-specific demand, $\gamma_i(s)$. The sector-time effects control for sector-specific implications of shocks that shift the Kimball demand index $d_t(s)$. For our benchmark specification we take the model at face value and impose a common super-elasticity $b(s) = b$ for all sectors s . We discuss alternative estimates that relax the assumption of a common super-elasticity and estimate different $b(s)$ for different subsamples of sectors in [Appendix C](#).

Outliers. We trim outliers by winsorizing establishment-level markups $\mu_{eit}(s)$ at the top and bottom 5% of each Census year.

Interpreting markup estimates. In our view, these markup estimates should be interpreted with some caution, both because of the issue of disentangling markups from output elasticities discussed above and because of the possibility that the firms' cost-minimization problem is misspecified — in which case, estimated markups will confound true markups with any other distortionary ‘wedge’ between prices and marginal cost, e.g., implicit or explicit input or revenue taxes, factor-adjustment costs, or price rigidities, etc.

Still, if one is prepared to take our estimated firm-level markups from the Census at face value, assuming away any other distortions etc, then one can compute the aggregate markup by taking the appropriate

Figure B1: Aggregate Markup from Census Data



Cost-weighted aggregate markup \mathcal{M} computed from the firm-level markups $\mu_{it}(s)$ constructed using micro data from the US Census of Manufactures from 1972 to 2012, as discussed in the text, for different values of the returns to scale (RTS). Our benchmark model assumes constant returns to scale, $RTS = 1$, but our results are robust to lower returns to scale.

weighted average. We report the results of this exercise in [Figure B1](#). This figure shows the evolution of the aggregate markup (cost-weighted average markup) for two cases, constant returns to scale ($RTS = 1.0$) and decreasing returns to scale ($RTS = 0.9$) for each Census year. For $RTS = 1.0$, the aggregate markup ranges from 1.20 in 1972 to a peak of 1.40 in 2002 before declining to about 1.33 in 2012. For $RTS = 0.9$ the aggregate markup is proportionately lower, ranging from 1.08 in 1972, peaking at 1.25 in 2002 before declining to about 1.20 in 2012.

References

- Akcigit, Ufuk, Salomé Baslandze, and Francesca Lotti**, “Connecting to Power: Political Connections, Innovation, and Firm Dynamics,” October 2018. NBER Working Paper 25136.
- Arkolakis, Costas, Arnaud Costinot, Dave Donaldson, and Andrés Rodríguez-Clare**, “Gains From Trade under Monopolistic Competition: A Simple Example with Translog Expenditure Functions and Pareto Distributions of Firm-Level Productivity,” September 2010. Yale University Working Paper.
- , —, —, and —, “The Elusive Pro-Competitive Effects of Trade,” *Review of Economic Studies*, January 2019, *86* (1), 46–80.
- Atkeson, Andrew and Ariel Burstein**, “Pricing-to-Market, Trade Costs, and International Relative Prices,” *American Economic Review*, 2008, *98* (5), 1998–2031.
- and —, “Innovation, Firm Dynamics, and International Trade,” *Journal of Political Economy*, June 2010, *118* (3), 433–484.
- and —, “The Aggregate Implications of Innovation Policy,” *Journal of Political Economy*, 2019, *127* (6), 2625–2683.
- , —, and **Manolis Chatzikonstantinou**, “Transitional Dynamics in Aggregate Models of Innovative Investment,” *Annual Review of Economics*, August 2019, *11*, 273–301.
- Autor, David, David Dorn, Lawrence F. Katz, Christina Patterson, and John Van Reenen**, “The Fall of the Labor Share and the Rise of Superstar Firms,” *Quarterly Journal of Economics*, May 2020, *135* (2), 645–709.
- Baqae, David Rezza and Emmanuel Farhi**, “Productivity and Misallocation in General Equilibrium,” *Quarterly Journal of Economics*, February 2020, *135* (1), 105–163.
- Barkai, Simcha**, “Declining Labor and Capital Shares,” *Journal of Finance*, 2020, *75* (5), 2421–2463.
- Bartelsman, Eric, John Haltiwanger, and Stefano Scarpetta**, “Cross-Country Differences in Productivity: The Role of Allocation and Selection,” *American Economic Review*, 2013, *103* (1), 305–334.
- Basu, Susanto**, “Are Price-Cost Markups Rising in the United States? A Discussion of the Evidence,” *Journal of Economic Perspectives*, 2019, *33* (3), 3–22.
- Behrens, Kristian, Giordano Mion, Yasusada Murata, and Jens Suedekum**, “Quantifying the Gap between Equilibrium and Optimum Under Monopolistic Competition,” *Quarterly Journal of Economics*, November 2020, *134* (4), 2299–2360.
- Bernard, Andrew B., Jonathan Eaton, J. Bradford Jensen, and Samuel Kortum**, “Plants and Productivity in International Trade,” *American Economic Review*, September 2003, *93* (4), 1268–1290.

- Bilbiie, Florin O., Fabio Ghironi, and Marc J. Melitz**, “Monopoly Power and Endogenous Product Variety: Distortions and Remedies,” October 2008. NBER Working Paper 14383.
- , – , and – , “Monopoly Power and Endogenous Product Variety: Distortions and Remedies,” *American Economic Journal: Macroeconomics*, October 2019, 11 (4), 140–174.
- Boar, Corina and Virgiliu Midrigan**, “Markups and Inequality,” November 2020. NBER Working Paper 25952.
- Bond, Steve, Arshia Hashemi, Greg Kaplan, and Piotr Zoch**, “Some Unpleasant Markup Arithmetic: Production Function Elasticities and their Estimation from Production Data,” *Journal of Monetary Economics*, July 2021, 121, 1–14.
- De Loecker, Jan, Jan Eeckhout, and Gabriel Unger**, “The Rise of Market Power and the Macroeconomic Implications,” *Quarterly Journal of Economics*, May 2020, 135 (2), 561–644.
- , – , and **Simon Mongey**, “Quantifying Market Power and Business Dynamism in the Macroeconomy,” April 2021. KU Leuven Working Paper.
- De Ridder, Maarten, Basille Grassi, and Giovanni Morzenti**, “The Hitchhiker’s Guide to Markup Estimation,” June 2022. LSE Working Paper.
- Dhingra, Swati and John Morrow**, “Monopolistic Competition and Optimum Product Diversity Under Firm Heterogeneity,” *Journal of Political Economy*, February 2019, 127 (1), 196–232.
- Dixit, Avinash K. and Joseph E. Stiglitz**, “Monopolistic Competition and Optimum Product Diversity,” *American Economic Review*, 1977, 67 (3), 297–308.
- Edmond, Chris, Virgiliu Midrigan, and Daniel Yi Xu**, “Competition, Markups, and the Gains from International Trade,” *American Economic Review*, October 2015, 105 (10), 3183–3221.
- Eslava, Marcela and John C. Haltiwanger**, “The Size and Life-Cycle Growth of Plants: The Role of Productivity, Demand and Wedges,” May 2020. NBER Working Paper 27184.
- Feenstra, Robert C.**, “A Homothetic Utility Function for Monopolistic Competition Models, Without Constant Price Elasticity,” *Economics Letters*, January 2003, 78 (1), 79–86.
- Foster, Lucia, Cheryl Grim, and John Haltiwanger**, “Reallocation in the Great Recession: Cleansing or Not?,” *Journal of Labor Economics*, January 2016, 34 (S1), S293–S331.
- Grassi, Basile**, “IO in I-O: Size, Industrial Organization, and the Input-Output Network Make a Firm Structurally Important,” December 2017. Bocconi University Working Paper.
- Gutiérrez, Germán and Thomas Philippon**, “Declining Competition and Investment in the US,” November 2017. NYU Stern Working Paper.

- **and** —, “Investment-Less Growth: An Empirical Investigation,” *Brooking Papers on Economic Activity*, Fall 2017, pp. 89–169.
- Hall, Robert E.**, “New Evidence on the Markup of Prices over Marginal Costs and the Role of Mega-Firms in the US Economy,” May 2018. NBER Working Paper.
- Hsieh, Chang-Tai and Peter J. Klenow**, “Misallocation and Manufacturing TFP in China and India,” *Quarterly Journal of Economics*, November 2009, *124* (4), 1403–1448.
- Jones, Charles I.**, “Intermediate Goods and Weak Links in the Theory of Economic Development,” *American Economic Journal: Macroeconomics*, 2011, *3* (2), 1–28.
- Karabarbounis, Loukas and Brent Neiman**, “Accounting for Factorless Income,” in Martin Eichenbaum and Jonathan A. Parker, eds., *NBER Macroeconomics Annual*, University of Chicago Press, June 2019.
- Kehrig, Matthias and Nicolas Vincent**, “The Micro-Level Anatomy of the Aggregate Labor Share Decline,” *Quarterly Journal of Economics*, May 2021, *136* (2), 1031–1087.
- Kimball, Miles S.**, “The Quantitative Analytics of the Basic Neomonetarist Model,” *Journal of Money, Credit, and Banking*, 1995, *27* (4, Part 2), 1241–1277.
- Klenow, Peter J. and Jonathan L. Willis**, “Real Rigidities and Nominal Price Changes,” *Economica*, July 2016, *83*, 443–472.
- Lerner, Abba P.**, “The Concept of Monopoly and the Measurement of Monopoly Power,” *Review of Economic Studies*, 1934, *1* (3), 157–175.
- Peters, Michael**, “Heterogeneous Mark-Ups, Growth and Endogenous Misallocation,” *Econometrica*, September 2020, *88* (5), 2037–2073.
- Restuccia, Diego and Richard Rogerson**, “Policy Distortions and Aggregate Productivity with Heterogeneous Establishments,” *Review of Economic Dynamics*, October 2008, *11* (4), 707–720.
- Rodriguez-Lopez, Jose Antonio**, “Prices and Exchange Rates: A Theory of Disconnect,” *Review of Economic Studies*, July 2011, *78* (3), 1135–1177.
- Rossi-Hansberg, Esteban, Pierre-Daniel Sarte, and Nicholas Trachter**, “Diverging Trends in National and Local Concentration,” in Martin Eichenbaum and Erik Hurst, eds., *NBER Macroeconomics Annual*, University of Chicago Press, 2020.
- Zhelobodko, Evgeny, Sergey Kokovin, Mathieu Parenti, and Jacques-François Thisse**, “Monopolistic Competition: Beyond the Constant Elasticity of Substitution,” *Econometrica*, November 2012, *80* (6), 2765–2784.