ENTRY VS. RENTS:
AGGREGATION WITH ECONOMIES OF SCALE

David Baqaee
Emmanuel Farhi

Entry vs. Rents: Aggregation with Economies of Scale
David Baqaee and Emmanuel Farhi
NBER Working Paper No. 27140
May 2020, Revised July 2021
JEL No. E0,E1,E3,O0,O11,O21,O25,O3,O4,O41

## ABSTRACT

This paper characterizes the response of aggregate output to micro shocks in dis-aggregated economies with entry, internal or external returns to scale, input-output linkages, and distortions. We decompose changes in output into changes in technical and allocative efficiency, and show that the latter depend on changes in rents and quasi-rents across markets. In addition, we characterize the social costs of distortions. We prove that while first-best industrial policy is network-independent, second-best policy does depend on the network, and boosts upstream industries with high returns to scale. As an example, we quantify the impact of misallocation caused by markups on aggregate efficiency in the US. In our preferred specification, losses are 40% of GDP whereas if we abstract from endogenous entry they are only 20%. Our baseline is sensitive not only to the presence of entry, but also to the specifics of how entry is modeled, in ways that our social-costs-of-distortions formulas clarify.

David Baqaee
Department of Economics
University of California at Los Angeles
Bunche Hall
Los Angeles, CA 90095
and CEPR
and also NBER
baqaee@econ.ucla.edu

Emmanuel Farhi
Harvard University
*NA user is deceased

# 1 Introduction

One of the major challenges of macroeconomics is the aggregation problem. The problem of translating microeconomic disturbances into macroeconomic consequences. Aggregation results like those of Hulten (1978) and Baqaee and Farhi (2019a) provide some structure through which the mapping from microeconomic primitives to macroeconomic outcomes can be understood. In this paper, we generalize these results to environments with non-constant returns to scale and product entry and exit. We decompose aggregate effects into technical and allocative efficiency changes, characterize the efficiency losses from misallocation, and derive optimal first- and second-best industrial policies.

Our analysis is relatively general, and allows for scale effects both due to downward-sloping demand, as in Dixit and Stiglitz (1977), or upward-sloping supply curves, as in Lucas (1978). We also allow for an arbitrary pattern of distorting wedges and technological heterogeneity. We allow for within and cross-industry heterogeneity, as well as unrestricted input-output linkages in both production and entry costs. We characterize how aggregate output and aggregate productivity respond to changes in technology and changes in wedges (which we take as primitives).

We decompose changes in output into changes in technical and allocative efficiency. Technical efficiency measures the direct impact of technology shocks, holding fixed the allocation of resources, and allocative efficiency measures the indirect effect of shocks due to the reallocation of resources.[1]

We show that changes in technical efficiency are given by a weighted sum of microeconomic technology shocks. The weight on each technology shock depends on expenditure shares and can be thought of as a cost-based or distortion-adjusted Domar weight or sales share. The intuition for this cost-based Domar weight is similar to the logic of Hulten (1978) or Liu (2017). This term captures the direct benefits of the technology shock on the production of the consumption bundle taking into account direct and indirect linkages but holding the allocation of resources fixed.

If the equilibrium is efficient, the change in aggregate output is equal to the change in technical efficiency. This is because, when the equilibrium is efficient, reallocation effects can be ignored to a first order. Hence, even in models with non-convexities and product creation and destruction, the logic of Hulten (1978) continues to apply as long as the equilibrium is efficient.

However, in economies with non-convexities and product entry and exit, efficiency

---

[1]There are different notions of changes in allocative efficiency. In this paper, we define them as changes in output due to reallocations of resources. See Baqaee and Farhi (2019a) for a detailed discussion.

is rarely attained. Once we stray from efficiency, we show that changes in allocative efficiency can play a theoretically and quantitatively important role in determining the aggregate consequences of disturbances. These reallocation effects depend on which markets expand and shrink, and on whether these adjustments in market sizes occur through changes in the size of existing producers or through changes in the number of producers.

We show that the resulting changes in allocative efficiency can be summarized by changes in *rents* and *quasi-rents*. Here, rent is income earned by firms after variable costs have been deducted from revenues (variable profit). Proprietors earn rents because of upward sloping marginal cost curves (diminishing returns giving rise to Ricardian rents) and because of deviations from marginal-cost pricing (e.g. markups or taxes giving rise to monopoly rents).[2] Quasi-rent, on the other hand, is a rent earned that is dissipated by entry costs. We show that changes in rents and quasi-rents capture reallocation effects in equilibrium.[3]

To model entry, we assume that there are potentially different types of entrants, and when entrants pays the entry cost, they draw a mix of different production technologies with some probability. This treatment is relatively general and nests the cases of directed technical change, in which entrants choose their technology from a menu (each technology has its own entry condition), and undirected innovation (entrants are assigned technologies randomly). We show that entry by new producers, and the quasi-rents associated with that entry, can be represented using linear projections. In particular, we show that the change in quasi-rents associated with some market is the projection of changes in rents on a vector space representing the entry technology.

Intuitively, suppose that rents in some market $i$ increase. Entrants respond to this by attempting to enter. However, entrants may not be able to perfectly direct their entry into market $i$. It could be that when entrants pay the entry cost, they only end up with $i$'s technology with some probability. We show that the amount of entry, and the quasi-rents associated with that entry, are given by a linear projection. Therefore, in a least-squares sense, new entry minimizes new rents claimed by existing producers.

Whereas the projection measures changes in quasi-rents, the residual from this projection is related to pure rents and measures the inability of entry to respond to changes in rents . A positive residual for the $i$th market means that rents in market $i$ have increased

---

[2] Here, monopoly rents also includes all the revenues collected by distortionary wedges (since other distorting wedges, say taxes, can be represented as markups).

[3] This generalizes the intuition in Baqaee and Farhi (2019a) that reallocation effects are summarized by changes in factor income shares. From the lens of this paper, income earned by factors are rents and since there is no entry or exit in that paper, quasi-rents are always equal to zero.

but entry is unable to keep up with the increase in profit opportunities. If new entrants can perfectly direct their entry decisions, then this residual is always zero. This projection (quasi-rents) and the residual (pure rents) summarize reallocation effects in general equilibrium.

There is a folk wisdom in the literature that modelling entry via diminishing marginal product, in the spirit of Lucas (1978), Hopenhayn (1992), or Restuccia and Rogerson (2008), is in some sense isomorphic to modelling entry via diminishing marginal utility, in the spirit of Dixit and Stiglitz (1977), Melitz (2003), or Hsieh and Klenow (2009). By allowing for both possibilities, we show that this intuition is very fragile and fails outside of very simple single-sector models. In particular, for the first class, the relevant sufficient statistic for reallocation is pure rent (the residual from the projection described above). On the other hand, for second class, the relevant sufficient statistic for reallocation is quasi-rent (the projection itself).

We complete this perspective by characterizing how changes in rents and quasi-rents are determined in equilibrium as functions of the microeconomic primitives and the shocks. That is, elasticities of substitution, expenditure shares, and wedges. This characterization, which depends on backward and forward propagation through supply chains, also determines how every price and quantity responds to a shock in equilibrium. This characterization fully pins down the model's positive properties as a function of primitives.

We use these positive results to study some normative questions. We characterize optimal industrial policy as well as the gains from implementing it. We generalize the influential insights of Harberger (1954) to economies with non-convexities and entry and exit. In particular, we show that the social cost of inefficiencies is, up to a second-order approximation, equal to the sales-weighted sum of a series of Harberger triangles. Some of these triangles are associated with production and some are associated with entry.

Using our positive results, we can characterize these Harberger triangles in terms of microeconomic primitives. In doing so, we overturn a common intuition, valid in models without entry, that the social cost of misallocation is monotone in the elasticities of substitution. Whereas a high elasticity of substitution increases the size of Harberger triangles associated with variable production, a low elasticity of substitution increases the size of Harberger triangles associated with entry. This results in non-monotonocity of losses with respect to elasticities of substitution.

While first-best policy is network-independent, and fixes distortions market by market, we also consider some second-best policies. In particular, we consider how a marginal entry or production subsidy affects output starting at the decentralized equilibrium. Unlike

first-best policies, the effect of second-best policies are network-dependent. In particular, for economies with external increasing returns to scale, we rationalize and revise Hirschman (1958)'s influential argument that policy should encourage expansion in sectors with the most forward and backward linkages, and we give precise formal definitions for these concepts. We show that the optimal marginal intervention aims to boost the sales of sectors that are relatively upstream and have strong scale economies.

We provide an example application by quantifying the social costs of markups using micro data for the U.S. We decompose the losses into losses arising from misallocation of resources in production (due to dispersion in markups) and misallocation of resources in entry (due to the suboptimal level of average markups). To use a concrete example, without entry, we find that markups estimated by a production-function approach à la De Loecker et al. (2019) reduce aggregate productivity by around 20%.[4] Accounting for entry can double these losses. Similarly, we also show that the gains from marginal tax interventions are potentially large.

Of course, there is a great deal of uncertainty about the specific number one attaches to these exercises. For us, the goal is to provide some back-of-the-envelope sense of the order of magnitudes and to shed light on what sufficient statistics determine these numbers. In particular, our aim is to show how assumptions on the production structure, including the strength of external economies, the extent to which entry is targeted, the type of resources used for entry, and the view one takes on the presence of entry restrictions affect the gains from policy and the losses from misallocation. While our results show that these features are critical theoretically and quantitatively, little is know about them empirically, and more empirical work is needed to bridge theory and measurement.

The structure of the paper is as follows. In Section 2, we set up the general model and define the equilibrium. In Section 3, we prove conditions under which the equilibrium is efficient and derive comparative statics for the efficient case. In Section 4, we specialize the model and introduce notation necessary to analyze inefficient equilibria. In Section 5, we provide and discuss the aggregation formula for how shocks affect aggregate output. Section 6 contains backward and forward propagation equations that determine how rents respond to shocks as a function of primitives. In Section 7, we analyze the normative properties of the economy, including first- and second-best optimal policy and the social costs of distortions. Section 8 is a quantitative application where we use a calibrated model to compute and dissect the social costs of markups and the benefits of industrial

---

[4]We also use alternative approaches for estimating markups: an alternative implementation of the production-function (PF) approach with different categories of costs, the user-cost approach (UC), and the accounting-profits (AP) approach. Although the numbers depend on the specification, the qualitative message remains the same.

policy in the U.S. using firm-level data on markups.

**Related Literature.** Our results apply to a broad range of influential models. For instance, our framework encompasses and generalizes models of entry like Dixit and Stiglitz (1977) or (a finite-horizon version of) Hopenhayn (1992), the closed economy version of Melitz (2003), and finite-horizon versions of models of endogenous growth with lab-equipment like Romer (1987) and Grossman and Helpman (1991). It also nests multi-sector and production network models like Hulten (1978), Long and Plosser (1983), and much of the subsequent literature like Gabaix (2011), Acemoglu et al. (2012), Jones (2013), Bigio and La'O (2016), and Baqaee and Farhi (2019b), amongst others.[5]

This paper is closely related to Baqaee (2018) and Baqaee and Farhi (2019a) which establish aggregation and propagation results for inefficient production networks with and without entry. Baqaee (2018) considers a tightly-parameterized class of production networks with external economies, entry, and distortions. This paper dispenses with the parametric restrictions, allows for a more sophisticated handling of entry, returns to scale, production functions, and network linkages in both production and entry. Furthermore, unlike Baqaee (2018), this paper also characterizes reallocation, misallocation, and optimal policy. On the other hand, Baqaee and Farhi (2019a) analyze reallocation and misallocation but, unlike this paper, abstract from entry.

This paper also relates to the literature on cross-sectional misallocation and policy interventions, with or without externalities, like Restuccia and Rogerson (2008), Hsieh and Klenow (2009), Epifani and Gancia (2011), Liu (2017), Osotimehin and Popov (2017), Behrens et al. (2016), Bartelme et al. (2019), Boehm and Oberfield (2020), Rubbo (2020), and La'O and Tahbaz-Salehi (2020). Our analysis of the economy's distance to the frontier is also related to Edmond et al. (2018), who analyze the social cost of markups. Our paper contributes to this literature by focusing on how the entry margin interacts with the input-output network to affect the costs of distortions. By showing that even in non-neoclassical economies with entry social losses can be approximated using Harberger triangles, the paper also extends the insights of Harberger (1954) and Harberger (1964).

Another strand of the literature which this paper relates to is the literature studying link-formation in production networks. In contrast to the approach in this paper, this literature takes discreteness of decisions seriously and is often studied with a non-Walrasian equilibrating mechanism. Some examples are Oberfield (2017), Acemoglu and Azar (2020), Acemoglu and Tahbaz-Salehi (2020), Taschereau-Dumouchel (2020), Lim (2017), Kikkawa et al. (2018), Dhyne et al. (2021), and Elliott et al. (2020). We abstract from these

---

[5]See Carvalho and Tahbaz-Salehi (2018) for a review of this literature.

issues in our analysis, assuming that individual firms are infinitesimal and that the mass of entrants and number of links formed adjusts smoothly in response to perturbations of primitives. In exchange for these simplifications, we can provide a fairly general local characterization of the equilibrium.

## 2 General Framework

The model consists of a representative household, a set of producers, and a set of entrants. In this section, we describe the model, and define the equilibrium. A circular flow diagram of the economy is depicted in Figure 1. Each rectangle represents a type of agent in the model. Loosely speaking, entrants buy resources to enter. After paying the entry costs, entrants are (perhaps randomly) assigned to be producers. Producers produce using intermediate materials they purchase from other producers. The representative household owns all resources in the economy and purchases consumption goods using national income. We begin by describing the problem each agent is faced with, starting with the producers.



Figure 1: Circular flow schematic of the economy showing the flow of resources.

### 2.1 Markets and Producers

There is a set of *markets* indexed by $i \in \mathcal{N}$. Each market $i$ is populated by an endogenous mass $M_i$ of identical producers with output

$$y_i = A_i f_i \left( \{x_{ij}\}_{j \in \mathcal{N}} \right),$$

where $f_i$ is a neoclassical production function, $A_i$ is some scalar indexing productivity, $x_{ij}$ is the input quantity of market good $j$ (including primary factors) used by $i$. Each producer

minimizes costs and sets its price $p_i^y$ equal to its marginal cost times an exogenous markup $\mu_i$.

The output good of market $i$ is given by

$$Y_i = F_i(M_i y_i),$$

where the market aggregator $F_i$ may have constant, decreasing, or increasing returns to scale in the producer-level output $y_i$. The price of the market good $P_i^Y$ is equal to the marginal cost of producing $Y_i$ times an exogenous wedge $\mu_i^Y$. Unlike the producer-level markup $\mu_i$, revenues generated by the market-level wedge $\mu_i^Y$ are *not* rebated to the owner of $i$ and instead go directly to the household. This distinction matters because revenues generated by $\mu_i$ incentivize entry, whereas revenues generated by $\mu_i^Y$ do not. The market-level wedge $\mu_i^Y$ therefore acts like an output tax.

To see the versatility of this modeling block, consider the following examples. Let $x$ denote a bundle of inputs and ignore productivity by setting $A_i = 1$ as an argument. Assume that $f_i(x) = x^{1-\gamma_i}$ for some $\gamma_i \in [0,1]$, and $F_i(x) = x^{\frac{1}{\gamma_i}}$. Then $Y_i = (M_i x^{\gamma_i})^{\frac{1}{\gamma_i}}$ captures a CES aggregator with an elasticity of substitution $1/(1-\gamma_i)$ between differentiated varieties produced under constant returns to scale. Suppose instead that $F_i(x) = x$. Then $Y_i = M_i x^{\gamma_i}$ captures a market structure with perfectly substitutable varieties produced under decreasing returns to scale.

**Primary Factors.** A subset of markets $\mathcal{F} \subset \mathcal{N}$ correspond to *primary factors*. For primary factors $f \in \mathcal{F}$, we assume that $M_f$ is exogenous, the production function $f_f$ has zero returns to scale, and the market aggregator has constant returns to scale $F_f(M_f y_f) = M_f y_f$. In addition, there are no markups/wedges $\mu_f = \mu_f^Y = 1$. In words, there is no entry into the factor market (since $M_f$ is fixed), each producer produces a fixed amount of output (since $f_f$ has zero returns to scale), and producer outputs are aggregated linearly (since $F_f$ is linear). This means that total market output of each factor is also fixed. This allows us to capture endowments of primary factors such as labor, land, or the initial capital stock.

## 2.2 Entrants

There is an infinite supply of potential entrants who are grouped into types indexed by $j \in E$. Entrants pay fixed costs and enter subject to a zero-profit condition.

**Fixed Costs.** To enter, potential type-$j$ entrants pay a fixed cost

$$g_j \left( \left\{ x_{E,ji} \right\}_{i \in \mathcal{N}} \right), \tag{1}$$

where $g_j$ has constant returns, and $x_{E,ji}$ is the input quantity of market good $i$. A simple example is when firms pay entry costs in units of labor if they choose to enter, as in Hopenhayn (1992) or Melitz (2003).

**Entry Technology.** The entry matrix $\zeta$ is an $|E| \times |\mathcal{N} - \mathcal{F}|$ positive-valued matrix. Type-$j$ entrants who pay the fixed cost are randomly assigned, according to $\zeta(j, i)$, the ability to produce in market $i \in \mathcal{N} - \mathcal{F}$. Without loss of generality, assume that the rows of $\zeta$ are linearly independent.[6] A simple example is that there is only one type of entrant and technology is assigned randomly, as in Hopenhayn (1992) or Melitz (2003). We denote by $M_{E,j}$ the endogenous mass of type-$j$ entrants who pay the entry cost.

If there is no way to enter market $i \in \mathcal{N}$, which occurs when $\zeta(j, i) = 0$ for all $j \in E$, then we allow for an exogenous mass $M_i$ of incumbents to operate in market $j$ without having to enter.

We refer to markets where entry is not possible as *uncontested markets* and denote their collection by $\mathcal{N}^u$. We also sometimes simply call them *incumbents*, since each of these markets operate like a representative incumbent. We refer to markets where entry is possible as *contested markets* and denote their collection by $\mathcal{N}^c$. Note that since we are flexible in the way we define and combine markets, we can capture a situation where incumbents and entrants compete by having them operate in different markets that are highly-substitutable with one another.

**Sunk and Overhead Costs.** The entry matrix $\zeta$ can capture sunk and overhead costs simultaneously. To capture sunk costs, suppose that $\zeta(j, i)$ has positive support for a range of different $i$'s. In this case, once the entry cost $j$ has been paid, the entrant will always choose to operate all of its technologies since the entry cost is sunk. At the other extreme, suppose that $\zeta(j, i) = 1$ for one specific $i$ and zero otherwise. In this case, entrant $j$ will only choose to pay the cost if operating technology $i$ is worth paying the fixed cost. In other words, the fixed cost is not sunk.[7]

---

[6]If the rows of $\zeta$ are not linearly independent, then some entry types are redundant (can be replicated by playing a mixed entry strategy).

[7]We can also consider intermediate situations in which entrant $j$ pays a sunk cost and draws a mixture of zero-returns technologies $j'$. Other entrants $j''$ can purchase the output of $j'$ and combine it with another fixed cost to enter with certainty into producing $i$. This structure mimics the entry decision in standard models such as Hopenhayn (1992) and Melitz (2003) where potential entrants first pay a sunk

**Zero-Profit Conditions.** The zero-profit condition for type-$j$ entrants equates expected profits post-entry with the costs of entry

$$\sum_i \frac{\zeta(j,i)M_{E,j}}{M_i}\lambda_{\pi,i} = M_{E,j}\sum_{k\in\mathcal{N}} P_k^Y x_{E,jk} = \lambda_{E,i},$$

where

$$\lambda_{\pi,i} = M_i p_i^y y_i - M_i \sum_{j\in\mathcal{N}} P_j^Y x_{ij}$$

is the total *rent* or *variable profit* (we use the two terms interchangeably) earned by all the producers of market $i$. The left-hand side of the zero-profit condition is the expected total rent earned by type-$j$ entrants and the right-hand side is the total cost of entry. This condition ensures that the rents earned by type-$j$ entrants are *quasi-rents* rather than *pure rents*.

## 2.3   Households

There is a representative household whose preferences are given by a homothetic utility function over market goods

$$Y = \mathcal{D}\left(\{C_i\}_{i\in\mathcal{N}}\right).$$

To avoid corners, we require that $Y \leq 0$ whenever $C_i = 0$ for any $i \in \mathcal{N}$. The budget constraint of this representative household requires total final expenditure to equal total income defined as revenues net of expenditures

$$\sum_{i\in\mathcal{N}} P_i^Y C_i = \sum_{i\in\mathcal{N}} P_i^Y Y_i - \sum_{j\in\mathcal{N}} P_j^Y x_{ij} - \sum_{j\in E} M_{E,j} \sum_{k\in\mathcal{N}} P_k^Y x_{E,jk}.$$

Note that payments to primary factors are included as the revenues of zero-returns-to-scale incumbents in markets $\mathcal{F} \subset \mathcal{N}$.

---

cost and then decide whether or not to pay an additional overhead cost before operating. The difference between our treatment of overhead costs and that in Hopenhayn (1992) and Melitz (2003) is that we assume divisibility and that they assume non-divisibility. We could capture non-divisibility by letting $g_j(\{x_{E,ji}\}_{i\in\mathcal{N}})$ have variable (possibly increasing) returns to scale (for example, by making it a step function). This would not affect Theorems 1 or 2. See Appendix B for more details and for a more general formalization which sidesteps these issues.

## 2.4 Resource Constraints

The resource constraint for market good $i \in \mathcal{N}$ is

$$Y_i = C_i + \sum_{j \in \mathcal{N}} M_j x_{ji} + \sum_{j \in E} M_{E,j} x_{E,ji},$$

in words, the total supply of good $i$ is equal to demand by households, producers (as intermediate inputs), and entrants (as fixed costs). The mass of producers in a contested market $i \in \mathcal{N}^c$ is the sum of the share of entrants $j \in E$ that obtained technology $i$:

$$M_i = \sum_{j \in E} \zeta(j, i) M_{E,j}.$$

The mass of producers $M_i$ in uncontested markets $i \in \mathcal{N}^u$ is exogenous.

## 2.5 Equilibrium

The decentralized equilibrium is an allocation of resources and collection of prices which clears markets and solves each agents' decision problem.

**Definition 1.** A *decentralized equilibrium* is a collection of prices $\{P_i^Y, p_i^y\}$ and quantities $\{C_i, Y_i, y_i, x_{ij}, x_{E,ij}, M_{E,j}, M_i\}$, such that given productivities $\{A_i\}$ and markups/wedges $\{\mu_i, \mu_i^Y\}$: (i) the representative household maximizes utility; (ii) each price is equal to marginal cost times the markup; (iii) entrants earn zero profits; (iv) prices clear all markets.

In this paper, we derive comparative statics with respect to changes in technologies $A_i$ and changes in wedges $\mu_i$ and $\mu_i^Y$. Since these are reduced-form wedges, many types of distortions like taxes, financial frictions, or nominal rigidities, are nested as special cases. For instance, to capture a financial friction on $i$'s ability to purchase inputs, add a fictitious incumbent producer to the model who buys inputs on behalf of $i$. An output wedge on this fictitious producer can then implement the same allocation as a financial friction on $i$.

Of course, many wedges, like variable markups or nominal rigidities, are themselves endogenous. Since given the changes in wedges, aggregation is the same, in this paper, we do not commit to a specific theory of wedge determination and instead ask how changing wedges affects output. Some questions, like the economy's distance from the Pareto efficient frontier, can be answered without endogenizing the wedges. Endogenizing wedges, for example variable markups or nominal rigidities, requires additional assumptions but would be a relatively straightforward extension.

Going forward, the primary object of interest is the response of aggregate output $d \log Y$ (or welfare) to productivity shocks and shocks to wedges.[8] Since the supply of primary factors is fixed, changes in aggregate output coincide with changes in aggregate productivity $d \log TFP$.[9]

## 2.6 Noteworthy Special Cases

At this level of abstraction, the model nests many general equilibrium models with entry, including models where goods are perfectly substitutable and firms have diminishing returns, as well as models where goods are imperfectly substitutable and firms have constant marginal cost. For example, it nests models of industry dynamics like (a finite-horizon version of) Hopenhayn (1992), the closed-economy version of Melitz (2003), models with product variety like Dixit and Stiglitz (1977), Krugman (1979), (finite-horizon versions of) growth models with lab-equipment, like Romer (1987) and Grossman and Helpman (1991), and models of production networks without entry like Acemoglu et al. (2012), Baqaee and Farhi (2019a), and Bigio and La'O (2016) or with entry like Baqaee (2018).

# 3 Marginal-Cost-Pricing Benchmark

In this section, we consider the marginal-cost pricing benchmark defined as follows

**Definition 2.** A *marginal-cost pricing equilibrium* is a decentralized equilibrium where $\mu_i = \mu_i^Y = 1$ for all $i \in \mathcal{N}$.

Theorem 1, below, shows that the marginal cost-pricing benchmark is efficient. Therefore, it generalizes the first welfare theorem to an environment with fixed and sunk costs of operation and entry. From a normative perspective, it clarifies how the optimal allocation can be implemented using linear taxes, and we use this implementation in Section 7 when we approximate the decentralized economy's distance from the Pareto-efficient frontier.

**Theorem 1** (First Welfare Theorem). *The marginal-cost pricing equilibrium is Pareto-efficient.*

This normative theorem is also important from a positive perspective since it ensures that the response of aggregate output to shocks can easily be obtained by applying the envelope theorem. Theorem 2 uses this insight to derive comparative statics.

---

[8]Since we allow for productivity shocks to producers, we can capture productivity shocks to the entry or overhead costs of operation by adding fictitious incumbents who buy inputs to be used for entry and shock their productivity. Finally, using the Arrow-Debreu trick of indexing commodities by dates and states of the world, we can capture dynamic stochastic models.

[9]We abstract away from the well-understood issues related to the treatment of new goods in the measurement of aggregate output. Therefore, in this paper, real GDP coincides with welfare.

**Theorem 2** (Comparative Statics under Efficiency). *In the marginal-cost pricing equilibrium, the response of aggregate output to a Hicks-neutral productivity shock $d \log A_i$ is given by*

$$\frac{d \log Y}{d \log A_i} = \frac{M_i p_i^y y_i}{GDP},$$

*which is the total sales of market i as a share of GDP. Similarly, the response of aggregate output to an entry productivity shock $d \log \zeta(j, i)$ is given by*

$$\frac{d \log Y}{d \log \zeta(j, i)} = \frac{\lambda_{\pi,i} \zeta(i, j) M_{E,j}}{GDP},$$

*which is the rents earned by type-j entrants from producing in market i as a share of GDP.*

Theorem 2 is an envelope theorem which extends Hulten (1978) to economies with selection, fixed costs, increasing returns, and an extensive margin of product creation and destruction. In particular, it shows that, for marginal-cost-pricing equilibria, simple and readily observable sufficient statistics like the sales or profit shares summarize the macroeconomic impact of microeconomic disturbances in general equilibrium.[10,11] In the next section, we derive comparative statics for shocks when the economy is inefficient.

# 4 Framework for Inefficient Equilibrium

Comparative statics in efficient models are easy to derive because, following the logic of the envelope theorem, reallocation effects can be ignored. Comparative statics in inefficient models are harder to obtain since reallocation effects can no longer be ignored. In this section, we impose some additional assumptions on the model, and introduce input-output notation, to analyze inefficient equilibria.

## 4.1 Additional Assumptions

To emphasize our mechanisms of interest, we specialize the general framework. Assume that each good is *either* produced as differentiated varieties using a CES aggregator, as in

---

[10]For the proof in Appendix B, we also allow for non-divisible overhead costs.

[11]Extending Theorem 2 to cover biased technical change, for example factor-augmenting shocks, or shocks to the entry or overhead costs of operation is trivial. To model these shocks, say a shock to $i$'s ability to use input $k$, simply introduce a new producer who buys from $k$ and sells to $i$. A Hicks-neutral shock to this new producer is the same as a biased shock in the original model. This trick allows us to restrict attention to Hicks-neutral shocks without loss of generality.

Dixit and Stiglitz (1977) or Melitz (2003), or as perfect substitutes with decreasing returns, as in Hopenhayn (1992).[12,13]

**Assumption 1** (Iso-elasticity). For each $i \in \mathcal{N} - \mathcal{F}$, either $i$ is a CES aggregate of imperfectly substitutable varieties

$$Y_i = \left( M_i y_i^{\gamma_i} \right)^{\frac{1}{\gamma_i}}, \qquad y_i = A_i f_i \left( \{x_{ij}\}_{j \in \mathcal{N}} \right), \tag{2}$$

or $i$ is a linear aggregate of perfectly substitutable varieties

$$Y_i = (M_i y_i), \qquad y_i = A_i f_i \left( \{x_{ij}\}_{j \in \mathcal{N}} \right)_i^{\varepsilon_i}, \tag{3}$$

where are $\gamma_i, \varepsilon_i \in [0, 1]$ and $f_i$ has constant returns to scale.

The parameters $\varepsilon_i$ and $\gamma_i$ control internal and external returns to scale (on the margin) respectively.[14] Since goods that are produced according to (2) have increasing external returns (due to love-of-variety), we refer to them as *IRS goods*. On the other hand, since goods according to (3) have decreasing internal returns, we refer to them as *DRS goods*. Denote the set of IRS goods by $\mathcal{N}^{IRS}$ and the set of DRS goods by $\mathcal{N}^{DRS}$.[15]

The next assumption rules out corners in $M_i$ by ensuring that markups are not so low that producer $i$ always makes negative profits. The rent (or variable profit) of market $i$ is

$$\lambda_{\pi,i} = \frac{P_i Y_i}{GDP} \pi_i, \quad \text{with} \quad \pi_i = \left( 1 - \frac{\varepsilon_i}{\mu_i} \right) \mathbf{1}_{i \in \mathcal{N}^{DRS}} + \left( 1 - \frac{1}{\mu_i} \right) \mathbf{1}_{i \in \mathcal{N}^{IRS}}. \tag{4}$$

Here $\pi_i$ is the share of market $i$'s sales that are claimed as profits. The profit margin $\pi_i$ consists of the rents due to market power and the rents due to diminishing returns.

**Assumption 2** (Positive profits). If $i$ is a DRS good, then $\mu_i > \varepsilon_i$, and if $i$ is an IRS good, then $\mu_i > 1$.

---

[12]Appendix A discusses the relationship between these two assumptions.

[13]Appendix F relaxes Assumption 1 and extends our results to the case where internal diseconomies are non-isoelastic (allowing for variable returns to scale at the producer level) and external economies are non-isoelastic (along the lines of Kimball, 1995).

[14]We normalize $A_i$ to ensure that $Y_i$ is unit-elastic in the productivity shock. For comparison, note that we did not impose this unit-elasticity normalization in Section 2.

[15]To map IRS goods into the framework in Section 2, note that a CES aggregate is isomorphic to an industry aggregator with increasing returns $F_i(x) = x^{1/\gamma_i}$ and a production function with decreasing returns $f_i(x) = x^{\gamma_i}$, where $x$ is the bundle of inputs. Since the aggregators do not generate any revenue (but have curvature), they must set price equal to average costs. This implies that, for a CES model, the outer aggregator is charging an implicit markup $\mu_i^Y = 1/\gamma_i$, and the inner aggregator is charging an implicit markdown $\mu_i^y = \gamma_i$. We come back to this in Section 7, in particulare footnote 26, when we discuss optimal policy and misallocation.

## 4.2 Notation and Other Preliminaries

We normalize nominal GDP to one throughout, so all prices are quoted in the nominal GDP numeraire and all sales, revenues, expenditures, and costs are expressed as shares of GDP.

We also represent the final demand function $Y = \mathcal{D}(C_1, \ldots, C_N)$ as the first producer in $\mathcal{N}$. In other words, we represent real GDP as the output of some incumbent producer standing in for the household. To emphasize the unique role the household plays in the economy, we index it by the number 0, to remind the reader that the zero-th "producer" is the household.

All the objects introduced below are defined at the initial equilibrium (around which we provide first-order and second-order approximations). We normalize the mass of entrants $M_{E,j}$ to one at the initial equilibrium.

**The Normalized Entry Matrix.** Define the $|E| \times |\mathcal{N}|$ normalized entry matrix $\tilde{\zeta}$ by

$$\tilde{\zeta}(j, i) = \frac{\zeta(j, i) M_{E,j}}{\sum_{k \in E} \zeta(k, i) M_{E,k}}$$

whenever market $i$ is contested and zero otherwise. This matrix gives the fraction of producers in market $i$ who are type-$j$ entrants (if $i$ is contested, then the $i$th column of $\tilde{\zeta}$ sums to one).

**IO Matrices.** We introduce the *forward* and *backward* Input-Output (IO) matrices $\Omega^{\mathrm{F}}$ and $\Omega^{\mathrm{B}}$ and their accompanying Leontief inverses $\Psi^{\mathrm{F}}$ and $\Psi^{\mathrm{B}}$. Intuitively, the forward matrix captures the transmission of prices from upstream suppliers to their downstream customers (forward linkages), whereas the backward matrix encodes the transmission of sales from downstream customers to their upstream suppliers (backward linkages).

Let $\Omega^V$ be the $|\mathcal{N}| \times |\mathcal{N}|$ matrix whose $ij$th element is equal to $i$'s variable expenditures on inputs from $j$ as a share of revenues

$$\Omega_{ij}^V \equiv \frac{M_i P_j^Y x_{ij}}{P_i^Y Y_i}.$$

Let $\Omega^E$ be the $|E| \times |\mathcal{N}|$ matrix whose $ij$th element is equal to entrant $i$'s expenditures on

15

inputs from $j$ as a share of the total entry costs

$$\Omega_{ij}^{E} \equiv \frac{P_{j}^{Y} x_{E,ij}}{\sum_{k \in \mathcal{N}} P_{k}^{Y} x_{E,ik}}.$$

**Backward IO Matrix.** Let $\pi$ be the $|\mathcal{N}| \times |\mathcal{N}|$ diagonal matrix of profit shares defined in (4). The *backward* IO matrix combines variable and fixed expenditures

$$\Omega^{B} = \Omega^{V} + \pi \tilde{\zeta}' \Omega^{E}.$$

Its $ij$th element $\Omega_{ji}^{B}$ is the fraction of the revenues of $j$ directly paid out to $i$ for variable production and entry. The associated backward Leontief inverse is

$$\Psi^{B} = \left(I - \Omega^{B}\right)^{-1} = I + \Omega^{B} + \left(\Omega^{B}\right)^{2} + \cdots.$$

Its $ij$th element $\Psi_{ij}^{B}$ is the fraction of the revenues of $i$ directly and indirectly (through the network) paid out to $j$ for variable production and entry.[16]

**Forward IO Matrix.** Let $\mathcal{E}$ be the $|\mathcal{N}| \times |\mathcal{N}|$ diagonal matrix whose $i$th diagonal element is equal to $\mathcal{E}_{ii} = (1 - \varepsilon_i)$ if $i \in \mathcal{N}^{DRS}$ or $\mathcal{E}_{ii} = (1/\gamma_i - 1)$ if $i \in \mathcal{N}^{IRS}$. Intuitively, $\mathcal{E}_{ii}$ measures how an increase in entry in $i$ affects the price of good $i$. The forward IO matrix is defined by

$$\Omega^{F} = \mu \Omega^{V} + \mathcal{E} \tilde{\zeta}' \Omega^{E},$$

where $\mu$ is a diagonal matrix of markups. The $ij$th element $\Omega_{ij}^{F}$ is the fraction of the cost of $i$ directly attributable to the price of $j$ through variable production and entry. The associated forward Leontief inverse is

$$\Psi^{F} = \left(I - \Omega^{F}\right)^{-1} = I + \Omega^{F} + \left(\Omega^{F}\right)^{2} + \cdots.$$

---

[16] The sales of $j$ can be broken down into to its sales to the the different $i$'s according to $\lambda_{j}^{B} = \sum_{i} \lambda_{i}^{B} \Omega_{ij}^{B}$. By implication, the $ij$th element of the backward IO matrix therefore encodes the elasticity of the sales of $j$ to the sales of $i$, so that $\Omega_{ij}^{B} = \partial \log \lambda_{j}^{B} / \partial \log \lambda_{i}^{B}$, where the partial derivative holds $\Omega^{B}$ and other sales $\lambda^{B}$ constant. The $ij$th element of the backward Leontief inverse therefore encodes the elasticity of the sales of $j$ to the sales of $i$, so that $\Psi_{ij}^{B} = \partial \log \lambda_{j}^{B} / \partial \log \lambda_{i}^{B}$, where the partial derivative holds $\Omega^{B}$ constant but accounts for changes in sales $\lambda^{B}$. As we shall see, this is equivalent to holding relative prices constant, since when relative prices are constant, $\Omega^{B}$ is also held constant.

Its $ij$th element $\Omega_{ij}^F$ is the fraction of the cost of $i$ directly and indirectly (through the network) attributable to the price of $j$ through variable production and entry.[17]

**Backward and Forward Domar Weights.** Following Domar (1961), the *Domar weight* of market $i$ is

$$\lambda_i^B = \frac{P_i^Y Y_i}{GDP} = P_i^Y Y_i,$$

where the last equality follows from our choice of numeraire. Theorem 1 implies that for the efficient benchmark, Domar weights are key sufficient statistics.

As a matter of accounting the Domar weight of $i$ coincides with its *backward Domar weight* defined as the $i$th element of the zero-th row of the backward Leontief inverse

$$\lambda_i^B = \sum_j \Omega_{0j}^B \Psi_{ji}^B = \Psi_{0i}^B.$$

This captures the household's exposure to $i$ via backward linkages or equivalently $i$'s centrality in demand.

The *forward Domar weight* of product $i$ is the $i$th element of the zero-th row of the forward Leontief inverse

$$\lambda_i^F = \Psi_{0i}^F = \sum_j \Omega_{0j}^F \Psi_{ji}^F.$$

This captures the household's exposures to $i$ via forward linkages or equivalently $i$'s centrality in supply.[18]

In the efficient marginal-cost pricing benchmark, the forward and backward Domar weights of market $i$ coincide $\lambda_i^B = \lambda_i^F$, so that the supply centrality (forward Domar weight) of the market is equal to its demand centrality (backward Domar weight), and both are equal to its sales share. By contrast, with inefficiencies, in general, the backward and forward Domar weights of market $i$ differ $\lambda_i^B \neq \lambda_i^F$ and their ratio $\lambda_i^F / \lambda_i^B$ measures the wedge between the supply and demand centralities of the market, or equivalently the reduction in the size of the market caused by the cumulated distortions in its downstream supply chain.

---

[17]By Shepard's lemma, the $ij$th element of the forward IO matrix encodes the elasticity of the price of $i$ to the price of $j$, so that $\partial \log P_i^Y / \partial \log P_j^Y = \Omega_{ij}^F$, where the partial derivative indicates that sales and shocks as well as other prices are held constant. By repeated applications of Shepard's lemma, the $ij$th element of the forward Leontief therefore encodes the elasticity of the price of $j$ to the price of $i$, so that $\Psi_{ij}^F = \partial \log P_i^Y / \partial \log P_j^Y$, where the partial derivative now indicates that sales and shocks are held constant but that other prices are allowed to vary.

[18]The backward and forward Domar weight generalize the revenue- and cost-based Domar weights in Baqaee and Farhi (2019a).

# 5 Aggregation

We now generalize Theorem 2 to inefficient economies. We provide our comparative statics in two steps. First, in this section, we provide an aggregation equation that gives the response to shocks of aggregate output as a function of changes in sales, rents, and quasi-rents. Second, in the next section, we provide propagation equations which give changes in sales (or rents) and quasi-rents, as a function of microeconomic primitives. Putting the two steps together yields our result. The shocks that we consider are shocks to all productivities and markups/wedges which we write in vector form as $(d \log A, d \log \mu)$.[19]

## 5.1 The Aggregation Equation

Let $\lambda_\pi$ be the $|\mathcal{N}| \times |\mathcal{N}|$ diagonal matrix of rents, and let $d \log \lambda_\pi$ be the $|\mathcal{N}| \times 1$ vector of changes in rents.[20] Define the (rent-weighted) projection of $d \log \lambda_\pi$ on the entry matrix $\tilde{\zeta}$ by

$$\widehat{d \log \lambda_\pi} = \tilde{\zeta}'(\tilde{\zeta}\lambda_\pi\tilde{\zeta}')^{-1}\tilde{\zeta}\lambda_\pi \, d \log \lambda_\pi. \tag{5}$$

**Lemma 1.** *Denote ith component of this projection by $\widehat{d \log \lambda_{\pi,i}}$. Let $\lambda_{E,j}$ be the expenditures on entry (quasi-rents) by type-$j$ entrants. In equilibrium,*

$$\widehat{d \log \lambda_{\pi,i}} = \sum_{j \in E} \tilde{\zeta}(j,i) d \log \lambda_{E,j}.$$

In words, $\widehat{d \log \lambda_{\pi,i}}$ is the expected log change in the *quasi-rent* associated with market $i$, that is, the expected change in the rents $d \log \lambda_E$ earned by entrants who entered into market $i$ (we discuss this in more detail in Section 5.2).

**Theorem 3** (Comparative Statics with Inefficiencies). *The response of aggregate output to shocks $(d \log A, d \log \mu)$ is given by*

$$d \log Y = \sum_{i \in \mathcal{N}} \lambda_i^F d \log A_i - \sum_{i \in \mathcal{N}^{IRS}} \lambda_i^F d \log \mu_i - \sum_{i \in \mathcal{N}^{DRS}} \lambda_i^F \varepsilon_i d \log \mu_i \tag{6}$$

---

[19] An output wedge on $i$ not rebated back to the proprietor, in our notation $\mu_i^Y$, can be captured by adding a fictitious incumbent middleman who buys $i$'s output and sells to the rest of the economy. A markup on this fictitious middleman is isomorphic to an output wedge on $i$. Therefore, comparative statics in $\mu$ encompass both output wedges and markups. In Appendix B, which contains the proofs, we explicitly distinguish between markups, $\mu_i$, and output wedges $\mu_i^Y$.

[20] We are purposefully defining $\lambda_\pi$ as an $|\mathcal{N}| \times |\mathcal{N}|$ diagonal matrix and $d \log \lambda_\pi$ be the $|\mathcal{N}| \times 1$ vector in order to streamline the matrix expressions for projections below. Throughout the paper, in order to lighten the notation, we often use the same symbol to denote vectors and their counterparts as diagonalized matrices.

$$- \sum_{i \in \mathcal{N}^{DRS}} \lambda_i^F (1 - \varepsilon_i) \left( d \log \lambda_i^B - \widehat{d \log \lambda}_{\pi,i} \right) + \sum_{i \in \mathcal{N}^{IRS}} \lambda_i^F \left( \frac{1}{\gamma_i} - 1 \right) \widehat{d \log \lambda}_{\pi,i}.$$

Theorem 3 is a key result of the paper, and we spend the rest of this section unpacking its intuition and working through some examples. The terms $\sum_{i \in \mathcal{N}} \lambda_i^F d \log A_i$ in (6) capture changes in output caused by changes in technology, holding fixed the allocation of resources. These are changes in *technical efficiency*, holding fixed the allocation of resources. When the equilibrium is efficient, as in Theorem 2, these are the only terms that matter. The remaining terms in (6) are changes in output caused by changes in the allocation of resources, holding fixed technology. We refer these to changes in *allocative efficiency*, caused by reallocations.[21]

Since nominal GDP is normalized to one, changes in real GDP are the negative of changes in the household price index $d \log Y = - d \log P_0^Y$. Therefore, one way to understand (6) is to think through how shocks affect the price of the consumer price index. Focus on the first line, which captures changes in consumer prices when sales and quasi-rents are held constant. The first term captures the direct effect of productivity shocks, which are weighted by their forward Domar weights. The second and third terms capture the effect of an increase in markups on consumer prices, which are weighted by the forward Domar weight of the bundle of inputs ($\lambda_i^F$ for IRS goods and $\lambda_i^F \varepsilon_i$ for DRS goods).

The second line accounts for changes in sales and quasi-rents. Intuitively, the first term on the second line captures how for each DRS market $i$, changes in the scale of operation of individual producers affect the price of the market good because of decreasing internal returns to scale. The second term on the second line captures how, for each IRS market $i$, changes in entry affect the price of the market good by stimulating external economies. In both cases, what matters is then how, for each market $i$, the change in the price of the good affects the price of final-demand.

We can also rewrite (6) in terms of changes in rents and quasi-rents like we do in the introduction. To do this, we use the fact that $d \log \lambda_\pi = d \log \lambda^B + d \log \pi$:

$$d \log Y = \sum_i \lambda_i^F d \log A_i - \sum_{i \in \mathcal{N}^{IRS}} \lambda_i^F d \log \mu_i - \sum_{i \in \mathcal{N}^{DRS}} \lambda_i^F \left[ \frac{1 - \varepsilon_i}{\pi_i} - 1 \right] d \log \mu_i$$

$$- \sum_{i \in \mathcal{N}^{DRS}} \lambda_i^F (1 - \varepsilon_i) \left( d \log \lambda_{\pi,i} - \widehat{d \log \lambda}_{\pi,i} \right) + \sum_{i \in \mathcal{N}^{IRS}} \lambda_i^F \left( \frac{1}{\gamma_i} - 1 \right) \widehat{d \log \lambda}_{\pi,i}, \qquad (7)$$

where the projection $\widehat{d \log \lambda}_\pi$ captures changes in quasi-rents (profits dissipated by entry

---

[21]See Appendix C for a formal discussion.

costs) and the residual $d \log \lambda_\pi - \widehat{d \log \lambda_\pi}$ captures changes in the difference between rents and quasi-rents (profits not dissipated by entry costs).

## 5.2 The Role of Entry

The following lemma helps clarify the role played by $\widehat{d \log \lambda_\pi}$, and its relation to quasi-rents and entry.

**Lemma 2.** *In equilibrium, the change in the mass of producers in each market is given by*

$$d \log M = \widehat{d \log \lambda_\pi} - \tilde{\zeta}'(\tilde{\zeta} \lambda_\pi \tilde{\zeta}')^{-1} \lambda_E \, d \log P_E, \qquad (8)$$

*where $d \log M$ is the $|\mathcal{N}| \times 1$ vector of changes in masses of producers, $\lambda_E$ is the $|E| \times |E|$ diagonal matrix of quasi-rents (expenditures on entry), and $d \log P_E$ is the $|E| \times 1$ vector of changes in entry prices.*[22] *Furthermore,*

Holding fixed entry costs ($d \log P_E = 0$), Lemma 2 shows that entry responds to changes in rents across the economy: entry changes to match the changes in rents $d \log \lambda_\pi$ to the extent possible. The normalized entry matrix $\tilde{\zeta}$ acts like the data matrix in a regression, and the response of the entrants to a change in rents is the linear projection of the changes in rents $d \log \lambda_\pi$ onto the space spanned by $\tilde{\zeta}$. Therefore, new entry acts to minimize the new rents going to existing producers.

Intuitively, the zero-profit condition is a moment condition (expected excess profits must equal zero). When we linearize this moment condition, it becomes a linear moment condition. Hence, changes in entry take the form of a linear projection to a first-order. When there are no wedges, this linear moment condition (expected excess profits must equal zero) also coincides with a social planner's first-order conditions. This is because, when prices are equal to marginal cost, excess profits are equal to the marginal social value of entry.

As shown by Lemma 1, The $i$th component of this projection, denoted by $\widehat{d \log \lambda_{\pi,i}}$, measures the changes in quasi-rents associated with market $i$. In other words, it is the change in the amount of resources spent by those entrants who go on to become producers of type $i$.

Holding fixed entry prices ($d \log P_E = 0$), if there are as many entrant types as there are markets $|E| = |\mathcal{N} - \mathcal{F}|$, then a change in profits in a given market maps, one for one, into a

---

[22]The entry price $P_{E,j}$ of the $j$th entrant is the marginal cost associated with the production function in equation (1).

change in the mass of entrants in that market. We call this situation fully directed entry, because in this case, changes in rents are captured entirely by new entrants as quasi-rents.

**Definition 3.** Entry is *fully-directed* if there are as many entrant types as there are markets $|E| = |\mathcal{N} - \mathcal{F}|$.

If there are fewer entrant types than markets $|E| < |\mathcal{N} - \mathcal{F}|$, entry into a particular product type may be restricted, or even impossible. When entry into a product type $i$ is impossible, $\zeta(j, i) = 0$ for every $j \in E$, product $i$ is either not produced, or if it is produced, then it is produced by incumbents. In this case, increases in $i$'s rents $d \log \lambda_{\pi,i}$ will not affect entry into $i$ at all, since $\widehat{d \log \lambda}_{\pi,i} = 0$.

## 5.3 Useful Special Cases

To build more intuition, consider three special cases of Theorem 3: (i) (CRS) all goods are produced with constant-returns to scale so that for every $j \in \mathcal{N} - \mathcal{F}$ either $\varepsilon_j = 1$ or $\gamma_j = 1$; (ii) (DRS) all goods are produced with decreasing returns $\mathcal{N} - \mathcal{F} = \mathcal{N}^{DRS}$; and (IRS) all goods are produced with increasing returns $\mathcal{N} - \mathcal{F} = \mathcal{N}^{IRS}$. For simplicity, for all of these examples, assume there is only one primary factor and that all other markets are contested. Consider a univariate productivity shock $d \log A_i$ (holding constant other productivities, wedges, and markups).[23]

**Productivity shocks with CRS.** When $\varepsilon_j = 1$ or $\gamma_j = 1$ for all $j \in \mathcal{N} - \mathcal{F}$, Theorem 3 reduces to

$$d \log Y = \lambda_i^{\mathrm{F}} d \log A_i. \tag{9}$$

In this model, only the direct technology shock matters and reallocations are irrelevant. The reason is because of free-entry. When $\varepsilon_j = 1$ and $\gamma_j = 1$, this means that entry is socially wasteful, and entry always adjusts to ensure zero profits. This means that entry absorbs or exudes resources in such a way that there are no changes in allocative efficiency, even though there are reallocations and the economy is inefficient. If there was no free-entry, then the behavior of output would be substantially more complicated, since we would then have to account for how technology shocks reallocate resources across producers (as in Baqaee and Farhi, 2019a). This example shows how dramatically free entry alters the behavior of output, even if entry itself is socially wasteful.

---

[23]The intuition for a shock to markups/wedges is similar, but for brevity, we relegate this discussion to Appendix K.

**Productivity shocks with DRS.** Suppose that $\varepsilon_j < 1$ for all $j \in \mathcal{N} - \mathcal{F}$, Theorem 3 then becomes

$$d \log Y = \lambda_i^F d \log A_i - \sum_{j \in \mathcal{N}-\mathcal{F}} \lambda_j^F(1 - \varepsilon_j)\left(d \log \lambda_{\pi,j} - \widehat{d \log \lambda}_{\pi,j}\right).$$

There are decreasing internal economies but no increasing external economies. If for some market $j \in \mathcal{N}$, entry cannot keep up with rents so that $d \log \lambda_{\pi,j} - \widehat{d \log \lambda}_{\pi,j} > 0$, then individual producers in this market scale up and run into diminishing returns. As a result, the prices of their producer-specific factors increase. This reallocation contributes to reducing aggregate output in proportion to the forward Domar weight $\lambda_j^F(1 - \varepsilon_j)$ of these specific fixed factors. The total effect of reallocations is obtained by summing over all markets. Reallocations lead to a more efficient use of resources when they reduce the scarcity of fixed factors by making them cheaper.

Such improvements in allocative efficiency cannot occur when the economy is efficient. In this case factor prices cannot go down on balance across markets. To see this, note that, when there are no markups, $\lambda_{\pi,j} = \lambda_j^F(1 - \varepsilon_j)$ and so, the reallocation terms become

$$\sum_{j \in \mathcal{N}} \lambda_j^F(1 - \varepsilon_j)(d \log \lambda_{\pi,j} - \widehat{d \log \lambda}_{\pi,j}) = \sum_{j \in \mathcal{N}} \lambda_{\pi,j}(d \log \lambda_{\pi,j} - \widehat{d \log \lambda}_{\pi,j}) = 0,$$

because the weighted sum of residuals must be zero.

When there is directed entry ($|E| = |N|$), this expression simplifies further to just

$$d \log Y = \lambda_i^F d \log A_i,$$

so that there are no changes in allocative efficiency. This is similar to (9). Intuitively, in this case, changes in the prices of market goods are determined independently from changes in their sales because changes in sales are accommodated entirely through changes in entry. In other words, even though the equilibrium may be inefficient, reallocations happen entirely on the extensive margin of entry and exit and offset one another.

**Productivity shocks with IRS.** When all non-primary factor markets are IRS, Theorem 3 becomes

$$d \log Y = \lambda_i^F d \log A_i + \sum_{j \in \mathcal{N}-\mathcal{F}} \lambda_j^F\left(\frac{1}{\gamma_j} - 1\right)\widehat{d \log \lambda}_{\pi,j}.$$

If in some market $j$, quasi-rents increase so that $\widehat{d \log \lambda}_{\pi,j} > 0$, then entry in the market increases and triggers external economies from love of variety. This reduces the negative

22

price of the associated specific fixed factor. This reallocation contributes to increasing aggregate output in proportion to the forward Domar weight $\lambda_i^F(1/\gamma_i - 1)$ of these specific fixed factors. The total effect of reallocations is obtained by summing over all markets.

# 6  Propagation

Theorem 3 in the previous section gives changes in aggregate output as a function of changes in rents and quasi-rents. In this section, we complete the theory by deriving propagation equations for the changes in sales (or rents) and quasi-rents. We do this in two steps: forward and backward propagation. In Section 6.1, we characterize the propagation of shocks through forward linkages: how changes in prices feed forward from suppliers to consumers. In Section 6.2, we characterize the propagation of shocks through backward linkages: how changes in sales feed backward from consumers to their suppliers. Together, they pin down changes in sales, rents, and quasi-rents, as well as all other disaggregated variables such as prices and quantities. We consider some worked out examples in Section 6.3.

## 6.1  Propagation Through Forward Linkages

We start by describing the response of prices to shocks.

**Proposition 1** (Forward Propagation). *In response to shocks* $(\mathrm{d}\log A, \mathrm{d}\log \mu)$, *changes in prices are given by*

$$d\log P_i^Y = -\sum_{j\in\mathcal{N}} \Psi_{ij}^F d\log A_j + \sum_{j\in\mathcal{N}^{IRS}} \Psi_{ij}^F \, \mathrm{d}\log \mu_j + \sum_{j\in\mathcal{N}^{DRS}} \Psi_{ij}^F \varepsilon_j \, \mathrm{d}\log \mu_j$$

$$+ \sum_{j\in\mathcal{N}^{DRS}} \Psi_{ij}^F \left(1 - \varepsilon_j\right)\left(d\log \lambda_j^B - \widehat{\mathrm{d}\log \lambda}_{\pi,j}\right) - \sum_{j\in\mathcal{N}^{IRS}} \Psi_{ij}^F \left(\frac{1}{\gamma_j} - 1\right)\widehat{\mathrm{d}\log \lambda}_{\pi,j}$$

Proposition 1 is similar to Theorem 3. In fact, since nominal GDP is normalized to one, changes in real output are just the negative of the changes in the consumer price index $\mathrm{d}\log Y = -\mathrm{d}\log P_0$. Therefore, Proposition 1 can be specialized to yield Theorem 3 by setting $i$ to be the price of the final consumption good 0. Therefore, the intuition for Proposition 1 is similar to the one for Theorem 3.[24]

---

[24]An interesting special case of Proposition 1 is when every good is DRS, there is only one primary factor,

## 6.2 Propagation Through Backward Linkages

For simplicity, we assume that all production and entry functions in the economy $f_i$ and $g_i$ are CES production functions. Without loss of generality (by relabelling the input-output network), we can assume that each CES production function $i$ has a single elasticity of substitution $\theta_i$ associated with it.[25] We make this assumption for clarity, not tractability, and Appendix D generalizes our results to non-CES production functions. We also assume that there are perfectly competitive incumbents that produce the goods required for paying the entry cost, and these perfectly competitive incumbents are added to the input-output network as additional "producers" whose outputs are only used by entrants.

To state our results, we use the *input-output covariance operator*:

$$Cov_m(X, \Psi^B_{(i)}) = \sum_{k \in \mathcal{N}} (1 - \pi_m)^{-1} \Omega^V_{mk} X_k \Psi^B_{ki} - \left( \sum_{k \in \mathcal{N}} (1 - \pi_m)^{-1} \Omega^V_{mk} \Psi^B_{ki} \right) \left( \sum_{k \in \mathcal{N}} (1 - \pi_m)^{-1} \Omega^V_{mk} X_k \right),$$

where $\Psi^B_{(i)}$ is the $i$th column of the backward Leontief inverse $\Psi^B$. This is the covariance between the vector $X$ and the $i$th column of $\Psi^B$, using the $m$th row of $(1 - \pi)^{-1} \Omega^V$ as the probability distribution. We call it a covariance since $\sum_{k \in \mathcal{N}} (1 - \pi_m)^{-1} \Omega^V_{mk} = 1$ for $m \in \mathcal{N} - \mathcal{F}$.

**Proposition 2** (Backward Propagation). *In response to shocks* $(d \log A, d \log \mu)$, *changes in sales are given by*

$$d\lambda^B_i = -\sum_{m \in \mathcal{N}} \lambda^B_m \sum_{k \in \mathcal{N}} \left[ \Omega^V_{mk} - (1 - \pi_m) \sum_{j \in E} \tilde{\zeta}_{jm} \Omega^E_{jk} \right] \Psi^B_{ki} d \log \mu_m$$

$$- \sum_{m \in \mathcal{N}} \lambda^B_m (1 - \pi_m)(\theta_m - 1) Cov_m \left( d \log P^Y, \Psi^B_{(i)} \right).$$

*Proportional changes in sales can be deduced using* $d \log \lambda^B_i = d \lambda^B_i / \lambda^B_i$.

The first term is the effect of changes in markups/wedges on the demand for $i$ holding

---

and entry is fully-directed. In this case, the change in prices simplifies to

$$d \log P^Y_i = -\sum_{j \in \mathcal{N}} \Psi^F_{ij} d \log A_j + \sum_{j \in \mathcal{N}^{DRS}} \Psi^F_{ij} \left( 1 - \frac{(1 - \varepsilon_j)}{\pi_j} \right) d \log \mu_j.$$

In other words, the change in relative prices does not depend on final demand or on the elasticities of substitution in production. This is reminiscent of the no-substitution theorem (Georgescu-Roegen,1951; Samuelson, 1951). However, it holds under different assumptions: in particular, unlike the classic no-substitution theorem, one does not need to assume constant returns to scale nor perfect competition.

[25]See the discussion of standard-form economies in Baqaee and Farhi (2019b) for more information.

fixed relative prices. The term in square brackets is how an increase in $m$'s markup $d \log \mu_m > 0$ affects spending on some input $k$. On the one hand, a higher markup reduces $m$'s variable spending on input $k$ by $\lambda^B_m \Omega^V_{mk}$. On the other hand, a higher markup increases entry, and this increases spending on $k$ by entrant $j$ by $\lambda^B_m (1 - \pi_m) \tilde{\zeta}_{jm} \Omega^E_{jk}$. These two effects in turn change spending on $i$ in proportion to the exposure $\Psi^B_{ki}$ of $k$ to $i$.

The term on the second line captures the effect of substitutions on the intensive margin. Changes in relative prices $d \log P^Y$ caused by the shocks lead individual producers in every market $m \in \mathcal{N}$ to shift their expenditures on their inputs. If $\theta_m > 1$, then $m$'s inputs are gross substitutes. Hence, $m$ substitutes its expenditures towards those inputs that have become relatively cheaper. If those inputs intensively rely on $i$, then $(\theta_m - 1) Cov_m (d \log P, \Psi^B_{(i)})$ is negative. Hence, substitution by $m$ changes $i$'s sales in proportion to $-\lambda^B_m (1 - \pi_m)(\theta_m - 1) Cov_m (d \log P, \Psi^B_{(i)})$.

We continue by describing the responses of rents and quasi-rents to shocks.

**Lemma 3** (Profit Shares). *In response to shocks* $(d \log A, d \log \mu)$, *changes in sales, rents, and profit margins are related through*

$$d \log \lambda_{\pi,i} = d \log \lambda^B_i + d \log \pi_i, \quad \text{where} \quad d \log \pi_i = \frac{1 - \pi_i}{\pi_i} d \log \mu_i.$$

Hence, changes in rents in each sector are driven, either by changes in sales $d \log \lambda^B$ or changes in profit margins $d \log \pi$. Hence, given changes in sales $d \log \lambda^B_i$ it is easy to obtain changes in rents $d \log \lambda_{\pi,i}$ and quasi-rents $\widehat{d \log \lambda}_{\pi,i}$ from Lemma 3 and equation (5).

**Combining Forward and Backward Propagation.** Proposition 1 can be plugged into Proposition 2 to give a linear system in sales shares $d \log \lambda^B$. The solution pins down changes in sales shares in every market, and these can then be plugged back into Theorem 3 for welfare changes. In Section 6.3, we apply these propositions to specific examples, and show how the model allows for a very rich set of behavior.

## 6.3 Illustrative Examples

In this section, we consider a multi-sector economy with homogenous firms within sectors, targeted free entry in all sectors with entry costs paid in units of labor. We focus on a shock to sector-level productivities. We explain how different views on returns to scale (IRS vs. DRS) lead to very different responses of aggregate output. For brevity, we keep

the examples in this section very simple. A variety of additional worked-out examples, including ones with input-output linkages, are presented in Appendix G.

We use the following notation throughout. Given three vectors $U$, $V$, and $W$ with $\sum_k U_k = 1$, we write $\mathbb{E}_U(V) = \sum_k U_k V_k$ and $Cov_U(V, W) = \sum_k U_k(V_k W_k) - (\sum_k U_k V_k)(\sum_k U_k W_k)$. We also sometimes use overlines to signal initial values when there is an ambiguity , but we drop them when there is none: for example, we alternatively write $\overline{\lambda}_i^{\mathrm{B}}$ or $\lambda_i^{\mathrm{B}}$ depending on the context.

**CRS without Entry.** We start by considering an economy without entry, and then consider the cases with entry. Consider an economy where aggregate output

$$Y = \left( \sum_k Y_k^{\frac{\theta_0 - 1}{\theta_0}} \right)^{\frac{\theta_0 - 1}{\theta_0}}$$

is a CES aggregate of differentiated inputs indexed by $k$ with an elasticity of substitution $\theta_0$. Each sector $k$'s output

$$Y_k = \left( M_k y_k^{\frac{\theta_k - 1}{\theta_k}} \right)^{\frac{\theta_k}{\theta_k - 1}}$$

is itself a CES aggregate of some mass $M_k$ of differentiated varieties with an elasticity of substitution $\theta_k > \min\{\theta_0, 1\}$. Each variety in sector $k$ is produced from labor with constant returns and productivity $A_k$ by a single firm and sold at a markup $\mu_k > 1$ over marginal cost

$$y_k = A_k l_k, \quad p_k^y = \mu_k mc_k.$$

We submit this economy to a vector of sector-level productivity shocks $d \log A$. To apply our formulas, we use the fact that in this example, the backward Domar weight and the forward Domar weight of each sector are equal to each other $\lambda_k^{\mathrm{B}} = \lambda_k^{\mathrm{F}}$. Applying Theorem 3, the change in aggregate output is given by

$$d \log Y = \mathbb{E}_{\lambda^{\mathrm{B}}} (d \log A) - d \log \lambda_L^{\mathrm{B}},$$

where $d \log \lambda_L^{\mathrm{B}}$ is the change in labor's share of income. The first term is the change in technical efficiency, holding fixed the allocation of resources, and the second term is the change caused by reallocations. Applying Proposition 2, the second term is

$$d \log \lambda_L^{\mathrm{B}} = (\theta_0 - 1) \frac{1}{\lambda_L^{\mathrm{B}}} Cov_{\lambda^{\mathrm{B}}} \left( \frac{1}{\mu}, d \log A \right).$$

26

To understand the intuition, suppose that sectors are substitutes ($\theta_0 > 1$) and that the shock disproportionately increases the productivity of sectors with high markups. Since the shock disproportionately increase the productivity of high-markup firms ($Cov_{\lambda^B}(d \log A, 1/\mu) < 0$), and since goods are substitutes ($\theta_0 > 1$), the shock reallocates labor towards high-markup firms and reduces the labor share (rents earned by labor). This reallocation improves allocative efficiency, because high-markup firms were too small to begin with from a social perspective, and boosts aggregate output.

**IRS with Free Entry à la Dixit and Stiglitz (1977).**   Consider the same model as above, but now suppose that there is targeted entry into every sector $k$, with potential entrants choosing which sector to enter into after paying a fixed cost in units of labor. To apply our formulas, we use the fact that in this example, we also use the fact that the backward Domar weight (income share) of labor is one $\lambda_L^B = 1$ because of free entry. Denote by $\theta$ the vector of within-sector elasticities of substitution. From Theorem 3, changes in aggregate output are now given by

$$d \log Y = \mathbb{E}_{\lambda^B}(d \log A) + \sum_k \frac{1}{\theta_k - 1} \lambda_k^B d \log\left(\lambda_k^B\left(1 - \frac{1}{\mu_k^q}\right)\right),$$

where the first term is the direct technology effect and the second term is the reallocation effect. Note that the reallocation effect is very different to what it was without free entry. From Proposition 2, the second term is equal to

$$\sum_k \frac{1}{\theta_k - 1} \lambda_k^B d \log\left(\lambda_k^B\left(1 - \frac{1}{\mu_k^q}\right)\right) = \frac{Cov_{\lambda^B}\left(\frac{\theta-1}{\theta-\theta_0}, d \log A\right)}{\mathbb{E}_{\lambda^B}\left(\frac{\theta-1}{\theta-\theta_0}\right)}.$$

To understand the intuition, suppose that sectors are substitutes ($\theta_0 > 1$) and that the shock disproportionately increases the productivity of sectors with high external economies (low elasticities of substitution). Then $Cov_{\lambda^B}((\theta - 1)/(\theta - \theta_0), d \log A) > 0$ and so the shock leads to improvements in allocative efficiency. Intuitively, the shock triggers beneficial reallocations of labor towards sectors with high external economies which were too small to begin with from a social perspective. These forces operate in reverse when sectors are complements with $\theta_0 < 1$.

   The correlation of productivity shocks and markups, which was key in the economy without entry is now irrelevant. This is because now labor reallocations happen purely on the extensive margin via changes in entry in the different sectors, while the intensive margin remains unchanged as individual producers in the different sectors keep operating

at the same scale. Instead, the key is now the correlation of productivity shocks and returns to scale.

**DRS with Free Entry à la Hopenhayn (1992).** We now show that changes in aggregate output are very different under DRS. Consider the same multi-sector model but now assume that each sector $k$'s output

$$Y_k = M_k y_k$$

is a linear aggregate of an endogenous mass $M_k$ of undifferentiated varieties. Each variety in sector $k$ is produced from labor with decreasing returns $\varepsilon_k$ and productivity $A_k$ by a single firm and sold at a markup $\mu_k$ over marginal cost

$$y_k = A_k l_k^{\varepsilon_k}, \quad p_k^y = \mu_k mc_k.$$

Like in the IRS case, we are again interested in how a long-run steady-state with targeted free entry in all sectors responds to a vector of sector-level productivity shocks $d \log A$.

Changes in aggregate output are given by

$$d \log Y = \mathbb{E}_{\lambda^B}(d \log A).$$

Changes in technical efficiency are captured by the same Hulten-like term as in the IRS case. By contrast, there are no longer any changes in allocative efficiency. This occurs despite the fact that there are equilibrium reallocations and that the model is inefficient. Basically, the adjustment in the sizes of the different sectors happen entirely on the extensive margin via changes in entry. Individual producers in the different sectors keep operating at the same scale so that there is no change on the intensive margin. Since in addition there are no external economies, the price of each good at the sectoral level is the same as the price of the good for individual firms. Reallocations therefore do not are therefore neutral on efficiency grounds. This example clarifies that the IRS and DRS models are, in general, very different.

# 7 Optimal Policy and Misallocation

In this section, we turn to optimal policy. We describe first-best policies when instruments are unrestricted. Then we characterize the gains from implementing optimal policy (i.e. the distance to the efficient frontier). Finally, we consider second-best policies when only limited instruments are available.

## 7.1 First-Best Policy

Theorem 1 implies that the first best is attained when $\mu_i^Y = \mu_i = 1$ for all $i \in \mathcal{N}$. Consider a planner with access to unrestricted linear taxes. Then applying Theorem 1 to the class of economies that satisfy Assumption 1 results in the following.

**Corollary 1** (First-Best). *The decentralized equilibrium is efficient if in each market $i \in \mathcal{N}$, the planner introduces output taxes $\tau_i^Y$ and $\tau_i^y$ on industry-level and firm-level output satisfying*

$$\tau_i^Y = \mathbf{1}_{\{i \in \mathcal{N}^{DRS}\}} + \gamma_i \mathbf{1}_{\{i \in \mathcal{N}^{IRS}\}}, \qquad \tau_i^y = \frac{1}{\mu_i} \mathbf{1}_{\{i \in \mathcal{N}^{DRS}\}} + \frac{1}{\gamma_i \mu_i} \mathbf{1}_{\{i \in \mathcal{N}^{IRS}\}},$$

*where revenues collected by the industry-level tax are paid to the household, and revenues collected by the firm-level tax are paid to producers of the good.*

Corollary 1 shows that first-best policy is independent of the input-output network. The policy intervention in each market depends only on the markups/wedges and the returns to scale in that market. In particular, for DRS markets, optimal policy ensures marginal cost pricing. For IRS markets, optimal policy sets the markup charged by each producer equal to $1/\gamma_i$, which is the Dixit and Stiglitz (1977) markup, to incentivize entry and subsidises output by $\gamma_i$ to offset markup.[26]

## 7.2 Social Costs of Distortions

In this section we characterize the gains from optimal policy, which coincide with the social costs of distortions, the distance from the efficient frontier, or the amount of misallocation. We show that even with non-neoclassical ingredients like entry, non-convexities, and external economies, the distance to the frontier can be approximated via a Domar-weighted sum of Harberger triangles associated with variable production and entry. We then specialize the result and work through a series of examples to emphasize the importance of accounting for entry.

For any equilibrium variable $X$, we denote by $\mathrm{d}\log X$ the log-deviation of $X$ from its value at the efficient allocation, which can also be thought of as the change in $X$

---

[26]At first glance, it may seem that Corollary 1 contradicts the marginal cost pricing result in Theorem 1, since for IRS markets, the firms do not seem to be marginal cost pricing. To see that Corollary 1 is a consequence of Theorem 1, one needs to recognize that an IRS industry with production function $Y_i = \left(M_i y_i^{\gamma_i}\right)^{1/\gamma_i}$ is equivalent to $Y_i = F_i(M_i f_i(y_i))$, where $F_i(x) = x^{1/\gamma_i}$ and $f_i(x) = x_i^\gamma$ are as in Section 3. Since aggregators $F_i$ and $f_i$ have curvature but generate no income, this means that they implicitly charge price equal to average cost (as opposed to marginal cost). Therefore, to restore efficiency, the planner must introduce taxes that undo these implicit markups and markdowns.

caused by the deviations of $\mathrm{d}\log\tau_i$ and $\mathrm{d}\log\tau_i^Y$ of the firm-level and industry-level output wedges from their efficient values in Corollary 1. We provide a second-order approximation in these deviations $(\mathrm{d}\log\tau, \mathrm{d}\log\tau^Y)$ of the associated aggregate efficiency loss $\mathcal{L} = -(1/2)\,\mathrm{d}^2\log Y.$[27]

**Proposition 3** (Deadweight-Loss). *As long as $\varepsilon_i, \gamma_i < 1$, the efficiency loss can be approximated, up to second-order approximation, as*

$$\mathcal{L} \approx \frac{1}{2}\sum_{i\in\mathcal{N}}\lambda_i^B\,\mathrm{d}\log y_i\,\mathrm{d}\log\left(\tau_i\tau_i^Y\right) + \frac{1}{2}\sum_{i\in\mathcal{N}^{IRS}}\lambda_i^B\frac{1}{\gamma_i}\,\mathrm{d}\log M_i\,\mathrm{d}\log\tau_i^Y + \frac{1}{2}\sum_{i\in\mathcal{N}^{DRS}}\lambda_i^B\,\mathrm{d}\log M_i\,\mathrm{d}\log\tau_i^Y.$$

Hence, the social cost of distortions is, up to a second-order approximation, a Domar-weighted sum of Harberger triangles associated with variable production and entry. In conjunction with the forward and backward propagation equations in Propositions 1 and 2, we can rewrite these loss functions in terms of microeconomic primitives (the input-output matrix, the elasticities of substitution, and returns to scale).[28] We relegate this general formula to Appendix B, and focus on a few prominent examples obtained by considering a special class of models with a sectoral structure.

## 7.3 Sectoral Models.

To generate examples, we will use *sectoral* models defined by the following conditions:

1.  every goods market $i \in \mathcal{N}-\mathcal{F}$ is assigned to a unique sector $\mathcal{I}$, with common returns to scale so that its output matters only through sectoral output. Sectoral output is

$$Y_{\mathcal{I}} = \sum_{i\in\mathcal{I}} y_i^{\varepsilon_{\mathcal{I}}}, \qquad \text{or} \qquad Y_{\mathcal{I}} = \left(\sum_{i\in\mathcal{I}} y_i^{\gamma_{\mathcal{I}}}\right)^{\frac{1}{\gamma_{\mathcal{I}}}},$$

    depending on whether $\mathcal{I}$ is DRS or IRS;

2.  individual producers $i$ in sector $\mathcal{I}$ have the same constant-returns production function upto a productivity shifter $y_i = A_i f_{\mathcal{I}}(\{x_{i\mathcal{J}}\})$, where $x_{i\mathcal{J}}$ indicates that inputs are sectoral aggregates;

---

[27]By Corollary 2, around the efficient point, the first-order loss is zero as long as $\varepsilon_i, \gamma_i < 1$. If either $\varepsilon_i = 1$ or $\gamma_i = 1$, and there is entry into $i$, then the losses are first-order, and we must instead use Theorem 3.

[28]To do this, note that $\mathrm{d}\log Y_i = \mathrm{d}\log\lambda_i^B - \mathrm{d}\log P_i$, where Proposition 1 gives $\mathrm{d}\log P_i$ and Proposition 2 gives $\mathrm{d}\log\lambda_i^B$. Next, observe that $\mathrm{d}\log Y_i = \mathrm{d}\log y_i + 1/\gamma_i\,\mathrm{d}\log M_i$ if $i$ is IRS and $\mathrm{d}\log Y_i = \mathrm{d}\log y_i + \mathrm{d}\log M_i$ if $i$ is DRS. Finally, note that $\mathrm{d}\log M$ is given by Lemma 2. Putting this altogether will allow us to write Proposition 3 in terms of primitives.

3. there is one type of entrant for each sector $\mathcal{I}$, and entrants are randomly assigned to $i \in \mathcal{I}$ according to some fixed distribution;

4. individual producers $i$ in sector $\mathcal{I}$ charge different markups $\tau_i^y$ but share common industry-level wedges $\tau_i^Y = \tau_{\mathcal{I}}^Y$.

Sectoral models, common in the literature, are worth singling out because their within-sector heterogeneity can be aggregated. For sectoral models, we can break the problem of computing the distance to the frontier into two blocks, within and across sectors. See Appendix E for detailed derivations.

Throughout the following examples, we define the sales share of sector $\mathcal{I}$ to be $\lambda_{\mathcal{I}}^B = \sum_{j \in \mathcal{I}} \lambda_j^B$, and producer $i$'s share of sector $\mathcal{I}$ to be $\lambda_i^{\mathcal{I},B} = \lambda_i^B / \lambda_{\mathcal{I}}^B \mathbf{1}_{\{i \in \mathcal{I}\}}$. We will denote by $\mathbb{E}_{\lambda^{\mathcal{I},B}}(\text{d}\log \tau^y)$ and $Var_{\lambda^{\mathcal{I},B}}(\text{d}\log \tau^y)$ the within-sector weighted expectations and variances of changes in markups/wedges $\text{d}\log \tau_i^y$ of producers $i \in \mathcal{I}$ with weights $\lambda_i^{\mathcal{I},B}$.

### 7.3.1 Sectoral DRS Example.

For sectoral models, we can provide a straightforward characterization of the loss function with DRS. We proceed under the additional assumptions that there is only one primary factor, that entry paid in that factor, and that there are no deviations of output wedges from their efficient benchmarks $\text{d}\log \tau_{\mathcal{I}}^Y = 0$.

**Proposition 4** (Deadweight-Loss in DRS Economy with Entry)**.** *Consider a sectoral model where every sector is DRS, there is only one primary factor, entry is paid in units of the factor, and there are no deviations of output wedges from their efficient benchmarks* $\text{d}\log \tau_{\mathcal{I}}^Y = 0$*. To a second order, the loss function is given by*

$$\mathcal{L} = \frac{1}{2} \sum_{\mathcal{I}} \lambda_{\mathcal{I}} \frac{\varepsilon_{\mathcal{I}}}{1 - \varepsilon_{\mathcal{I}}} Var_{\lambda^{\mathcal{I},B}}(\text{d}\log \tau^y) + \frac{1}{2} \sum_{\mathcal{I}} \lambda_{\mathcal{I}} \frac{\varepsilon_{\mathcal{I}}}{1 - \varepsilon_{\mathcal{I}}} (\mathbb{E}_{\lambda^{\mathcal{I},B}}(\text{d}\log \tau^y))^2.$$

Because there are no output wedges, we know from Proposition 3 that there are no Harberger triangles associated with entry and only Harberger triangles associated with variable production. Of course, this does not mean that the entry margin is irrelevant, but changes in entry only matter through the impact on variable production.

The first term in the loss function captures misallocation arising from distortions in relative producer sizes driven by dispersed markups/wedges within sectors. The second term captures misallocation arising from distortions in the average size of firms, or equivalently, distortions driven by an inappropriate average levels of markups within sectors. The losses increase with the returns to scale: they go to zero in the zero-returns

31

to scale limit where $\varepsilon_I$ goes to one, and they go to infinity in the constant-returns limit where $\varepsilon_I$ goes to zero.

Proposition 4 is surprising if one is familiar with the misallocation literature. Normally, elasticities of substitution are key pieces of information. In this case, this information is not relevant. Intuitively, this is because the wedges do not cause cross-sectoral misallocation. In this model, changes in sectoral markups do not change relative sectoral prices to a first-order. An increase in a sector's markup increases the prices of producers in that sector but reduces their scale. At the efficient point, these effects cancel exactly. Therefore, sectoral prices do not change to a first-order, which means that to a first-order, the elasticity of substitution across sectors is not relevant for how quantities adjust. Since Harberger triangles are products of first-order changes in quantities and first-order changes in the markups, the cross-sectoral elasticity of substitution is irrelevant to a second-order.

### 7.3.2 Sectoral IRS Examples.

Consider a sectoral model with IRS sectors. We denote by $\theta_I = 1/(1 - \gamma_I)$ the elasticity of substitution associated with every sector $I$. In other words, we can think of $\theta_I$ as the elasticity of substitution across different varieties $i$ in sector $I$.

The behavior of sectoral IRS models is substantially more complicated than that of sectoral DRS models. In particular, whereas cross-sector elasticities of substitution are irrelevant under DRS, they are relevant under IRS. Rather than providing the complicated general formula, we instead focus on some simple examples to give intuition. In each case, seemingly small changes in the assumptions about the nature of entry make the welfare costs of distortions quite different.

**One-Sector Economy.** We start with a one-sector model heterogenous-firm economy with IRS and free entry. Aggregate output is given by

$$Y = \left( \sum_i y_i^{\frac{\theta-1}{\theta}} \right)^{\frac{\theta}{\theta-1}}.$$

Each good $i$ is produced from labor with constant returns and productivity $A_i$. The aggregate efficiency loss is, to a second-order, given by

$$\mathcal{L} = \frac{1}{2} \theta Var_{\lambda^B} \left( d \log \tau^y \right) + \frac{1}{2} \theta \mathbb{E}_{\lambda^B} \left( d \log \tau^y \right)^2.$$

The first term captures misallocation on the intensive margin and comes from the fact that high-markup firms are too small and low-markup firms are too big. This term depends on the elasticity of substitution and on the dispersion of markups, and is standard in the literature (see e.g. Hsieh and Klenow, 2009; Baqaee and Farhi, 2019a). The second term captures misallocation on the extensive margin and comes from the fact that there is too much or too little entry. This term depends on the elasticity of substitution and on the level of wedges, and is new to the literature.

**Multi-Sector Economy.** Now, consider a multi-sector economy where output is an aggregate of different sectors $Y_I$ indexed by $I$ with elasticity of substitution $\theta_0$

$$Y = \left( \sum_I Y_I^{\frac{\theta_0-1}{\theta_0}} \right)^{\frac{\theta_0}{\theta_0-1}},$$

where each industry's output is itself a CES aggregate of product varieties in industry $k$

$$Y_I = \left( \sum_{i \in I} M_i y_i^{\frac{\theta_I-1}{\theta_I}} \right)^{\frac{\theta_I}{\theta_I-1}}.$$

Here, $M_i$ is the mass of varieties of type $i$ in sector $I$. We assume that the within-sector elasticity of substitution $\theta_I$ is greater than the between-sector one $\theta_0$. Each good in each industry is produced linearly using labor $y_i = A_i l_i$. Labor is the only primary factor, there is free entry into each sector, and entry costs are paid in units of labor.

Applying Proposition 3, the aggregate efficiency loss is

$$\mathcal{L} = \frac{1}{2} \sum_I \lambda_I \theta_I \left[ Var_{\lambda^{I,B}} (d \log \tau^y) + \mathbb{E}_{\lambda^{I,B}} (d \log \tau^y) \right]^2.$$

Note that while the elasticities of substitution within sectors $\theta_I$ matter, the elasticity of substitution in consumption across sectors $\theta_0$ does not. This is because, at the efficient marginal-cost pricing equilibrium, changes in markups distort the allocation of resources within a sector between the extensive and intensive margins but these distortions have offsetting effects on the price of the sector good. Basically, there is only misallocation within sectors but no misallocation across sectors.

By contrast, with no entry and instead an exogenous mass $M_i$ of incumbents in each

market, Proposition 3 implies the aggregate efficiency loss function is

$$\mathcal{L} = \frac{1}{2}\left[\theta_0 Var_{\lambda^B}\left(\mathbb{E}_{\lambda^{I,B}}(d\log\tau^y)\right) + \sum_I \theta_I Var_{\lambda^{I,B}}\left(d\log\tau^y\right)\right],$$

where the first variance is the weighted variance of sectoral markups with weights given by sectoral sales shares $\lambda_I^B$, and the last set of variances are within-sector variances weighted by within-sector sales shares $\lambda_i^{I,B}$. Whereas $\theta_I$ controls losses from within-sector dispersion, the cross-sectoral elasticity of substitution $\theta_0$ controls the extent to which sectoral misallocation reduces output. Since there is no entry, the level of wedges is no longer relevant, only their dispersion.

This examples illustrates that allowing for entry changes which elasticities of substitution are relevant for misallocation.

**Roundabout Economy.**   We finish with an economy that has intermediate inputs. Consider a roundabout-entry economy with one sector populated by homogenous firms and free entry with entry costs requiring the use of both labor and goods. Aggregate output

$$Y = Y_1 - x_{21}$$

is the output of sector 1 not used for entry. Gross output is given by

$$Y_1 = \left(M_1 y_1^{\frac{\theta_1-1}{\theta_1}}\right)^{\frac{\theta_1}{\theta_1-1}},$$

where the representative producer has a production function $y_1 = A_1 l_1$ that transforms labor into goods linearly.

Suppose there is free entry and potential entrants pay a fixed entry cost that relies on both goods and factors. In particular, the entry good is produced from labor and products and sold at marginal cost

$$Y_2 = (l_2)^{\Omega_{2L}^V} (x_{21})^{1-\Omega_{2L}^V}.$$

When entry only uses labor $\Omega_{2L}^V = 1$, applying Proposition 3, the aggregate efficiency loss from markups is

$$\mathcal{L} = \frac{1}{2}\theta_1(d\log\tau_1^y)^2.$$

The loss is increasing in the elasticity of substitution across products $\theta_1$ since the love-of-variety effect is declining in $\theta_1$, and goes to zero as $\theta_1$ goes to infinity. In this limit, entry

is socially wasteful and the losses from any amount of entry are first-order (which is why the second-order approximation explodes).

Next, suppose that entry uses only products $\Omega_{2L}^V = 1$. Applying Proposition 3, the aggregate efficiency loss from markups is

$$\mathcal{L} = \frac{1}{2} \frac{(\theta_1 - 1)^3}{(\theta_1 - 2)^2} (d \log \tau_1^y)^2.$$

Once again, the losses goes to infinity as $\theta_1$ goes to infinity and for similar reasons. However, the loss is no longer increasing in $\theta_1$, but is instead U-shaped, and also goes to infinity as $\theta_1$ goes to 2 from above, since love of variety becomes so strong that output becomes linear in the mass of entrants. This example breaks the long-standing intuition in the misallocation literature that efficiency losses are increasing in the elasticity of substitution.[29]

This example illustrates how changing the input-output structure of entry can transform the losses from misallocation.

## 7.4 Second-Best Policy

Whereas the first-best policy is network-independent, second-best policies do depend on the details of the network. This section provides bang-for-buck formulas to compare the merits of different marginal interventions. These formulas revive and revise the informal policy recommendations of Hirschman (1958), who argued in favor of encouraging sectors with increasing returns that had the most backward and forward linkages. The analysis reveals the extent to which details matter: effective policy depends crucially on the nature of the intervention, the shape of the production network, and the strength of scale economies.

The marginal bang-for-buck analysis we perform is similar to the one in Liu (2017), however, we focus on an economy with increasing returns to scale due to love-of-variety. To simplify the analysis, we assume that there is only one primary factor which we call labor. We also assume that entry is possible in all markets, or, in other words, that all markets are contested. We consider marginal interventions at the no-intervention equilibrium. We investigate markup regulation and entry subsidization, which can loosely be thought of as capturing respectively competition and industrial policy. These two types of interventions neatly illustrate two very different ways in which forward and backward

---

[29]In Appendix G.4.1, we show that this U-shaped pattern also arises with input-output linkages in variable production rather than entry.

linkages can matter.

**Markup Regulation.** To start with, consider a budget-neutral intervention reducing the markups $d \log \mu_i < 0$ of the producers of market $i$. This can be achieved by placing a subsidy on $i$ and taxing owners of $i$ to fund the subsidy. Applying Theorem 3, the response of aggregate output, normalized by the revenues $-\lambda_i^B d \log \mu_i > 0$ transfered away from the producers by the associated implicit subsidy, is

$$-\frac{1}{\lambda_i^B} \frac{d \log Y}{d \log \mu_i} = \frac{\lambda_i^F}{\lambda_i^B} \frac{1}{\gamma_i} \left( \frac{\mu_i - 1/\gamma_i}{\mu_i - 1} \right) + \sum_{j \in \mathcal{N} - \mathcal{F}} \lambda_j^F \left( \frac{1}{\gamma_j} - 1 \right) \left( -\frac{1}{\lambda_i^B} \frac{d \widehat{\log \lambda_j^B}}{d \log \mu_i} \right).$$

The first term is the direct effect of the markup reduction, holding sales constant. It captures two opposing effects on the price of market good $i$ and in turn on final-demand prices. On the one hand, the policy reduces the price of each individual producer in market $i$, making the good cheaper for the household. On the other hand, the policy also dis-incentivizes entry into market $i$, which increases the effective price of $i$ due to reduced variety. Overall, whether the sign of the direct effect is positive or negative depends on whether there is too little or too much entry in market $i$ to begin with, which in turn depends on whether the initial markup $\mu_i$ is lower or higher than the infra-marginal surplus created by new varieties $1/\gamma_i$.

Under Dixit and Stiglitz (1977) monopolistic competition, the direct effects exactly cancel $\mu_i = 1/\gamma_i$, leaving us the second term

$$-\frac{1}{\lambda_i^B} \frac{d \log Y}{d \log \mu_i} = \sum_{j \in \mathcal{N} - \mathcal{F}} \lambda_j^F \left( \frac{1}{\gamma_j} - 1 \right) \left( -\frac{1}{\lambda_i^B} \frac{d \widehat{\log \lambda_j^B}}{d \log \mu_i} \right). \tag{10}$$

The bang-for-buck impact of the intervention is measured by a simple sufficient statistic: a forward-weighted sum across markets $j$ of the changes in backward-linkages interacted with increasing returns to scale $1/\gamma_j - 1$.[30]

**Entry Subsidies.** Now consider marginal entry subsidies to type-$i$ entrants at the no-intervention equilibrium. Without loss of generality, we treat the entry production function $g_i(x_{E,ij})$ of $i$ as though it were operated by an incumbent producer assembling the resources needed to enter and selling them at marginal cost $\mu_{E,i} = 1$. We capture en-

---

[30]Since the sum of backward linkages $\sum_i \lambda_i^B$ is a natural measure of intermediate input use. There is a sense in which the best marginal intervention acts to encourage intermediate input usage.

36

try subsidies by assuming that the government levies a subsidy on the output of these producers.

Denote the backward and forward Domar weights of these "entry-good producers" by $\lambda_{E,i}^{\mathrm{B}}$ and $\lambda_{E,i}^{\mathrm{F}}$. The backward Domar weight is equal to the quasi-rents of type $i$ entrants, or by the zero profit condition, the profits earned by type-$i$ entrants $\lambda_{E,i}^{\mathrm{B}} = \sum_{j \in \mathcal{N}-\mathcal{F}} \tilde{\zeta}_{ij} \lambda_{\pi,j}$. The forward Domar weight captures the impact of type-$i$ entry on final-demand prices $\lambda_{E,i}^{\mathrm{F}} = \sum_{j \in \mathcal{N}-\mathcal{F}} \tilde{\zeta}_{ji} \lambda_j^{\mathrm{F}} (1/\gamma_j - 1)$.

Introducing an entry subsidy on type-$i$ entrants is equivalent to reducing the markup $d \log \mu_{E,i} < 0$ of the producer of entry good $i$. At the no-intervention equilibrium the budgetary impact is just $-\lambda_{E,i}^{\mathrm{B}} d \log \mu_{E,i}^{Y} > 0$. The response of aggregate output, normalized by its budgetary impact to allow bang-for-buck comparisons, is

$$ -\frac{1}{\lambda_{E,i}^{\mathrm{B}}} \frac{d \log Y}{d \log \mu_{E,i}^{Y}} = \left( \frac{\lambda_{E,i}^{\mathrm{F}}}{\lambda_{E,i}^{\mathrm{B}}} - \frac{\lambda_L^{\mathrm{F}}}{\lambda_L^{\mathrm{B}}} \right) + \sum_{j \in \mathcal{N}-\mathcal{F}} \lambda_j^{\mathrm{F}} \left( \frac{1}{\gamma_j} - 1 \right) \left( -\frac{1}{\lambda_{E,i}^{\mathrm{B}}} \frac{\widehat{d \log \lambda_j^{\mathrm{B}}}}{d \log \mu_{E,i}^{Y}} \right), \tag{11} $$

where, at the no-intervention equilibrium, the sales share of labor $\lambda_L^{\mathrm{B}} = 1$ since all markets are contested (there are no pure profits), but $\lambda_L^{\mathrm{F}} \neq 1$ in general since there are inefficiencies.

The bang-for-buck impact of the intervention depends on two simple sufficient statistics corresponding to the two terms in (11). The second term is exactly the same as that for measuring the bang-for-buck impact of markup regulations in equation (10) and it has the same intuition.

By contrast, the first term is specific to entry subsidies. It depends on the difference between two ratios of forward to the backward Domar weights: that of entry $\lambda_{E,i}^{\mathrm{F}}/\lambda_{E,i}^{\mathrm{B}}$ where the intervention takes place, and that of labor $\lambda_L^{\mathrm{F}}/\lambda_L^{\mathrm{B}}$.[31] The ratio of forward to backward Domar weights $i$ measures the reduction in the size $i$ caused by cumulated wedges downstream. Hence, the first term boils down to a comparison of the cumulated distortions downstream from entry good $i$ compared to labor. Holding sales constant, the entry subsidy stimulates entry by type-$i$ entrants, which reduces final demand prices; but also absorbs more resources into entry, which increases the real price of labor (the labor share) and in turn raises final-demand prices.

In Appendix L, we apply these formulas to a Cobb-Douglas economy, showing that ceteris paribus, markup policies should lower markups in industries that are downstream and entry subsidies should subsidize entry into industries that are upstream.

---

[31] The intuition for this first term is related to Liu (2017), who studies marginal interventions around the decentralized equilibrium of a production network economy with constant returns.

# 8 Quantitative Illustration

In this section, we illustrate the social cost of distortions, or equivalently the gains from optimal policy using a stylized model. We also compute the social bang for a marginal buck of competition or industrial policy. We calibrate the model to fit U.S. data and provide a brief account of how we proceed; the details of how we map the model to data are in Appendix H.

## 8.1 Description of Quantitative Model

The quantitative model has a sectoral structure with heterogenous firms within sectors and one primary factor capturing value-added. We merge firm-level data from Compustat with industry-level data from the BEA. We use annual input-output tables from the BEA with 66 industries, and assign each firm in the our Compustat sample to a BEA industry. In the data, we observe industry-level sales shares for industries $\mathcal{I}$; input-output entries for industries $\mathcal{I}$ and $\mathcal{J}$; the sales shares of the Compustat firms $i$ in industry $\mathcal{I}$; and the markup $\mu_i$ of Compustat firm $i$.

We adopt the baseline estimates of De Loecker et al. (2019) to obtain firm-level markups using a production function estimation approach. In Appendix I, we perform robustness checks by recomputing our results using three alternative methods for estimating markups: an alternative implementation of the production function estimation approach with different categories of costs (including SG&A in variable costs, as in Traina, 2018), and alternative approaches that compute markups by netting out the cost of capital from gross surplus. Although the numbers depend on the specific approach, the qualitative message that accounting for entry and returns to scale is very important remains the same.

The model has a nested CES structure where each firm $i$ in industry $\mathcal{I}$ has a CES production function combining value-added and intermediate inputs with an elasticity of substitution $\theta_1$. The intermediate input component is itself a CES aggregator of inputs from other industries with an elasticity of substitution $\theta_2$. Finally, we have the within-sector elasticities $\varepsilon_{\mathcal{I}}$ or $\gamma_{\mathcal{I}}$ depending on whether we assume the industry is DRS or IRS.

Drawing on estimates from Atalay (2017), Herrendorf et al. (2013), and Boehm et al. (2014), we set the elasticity of substitution across sectors in consumption to be $\theta_0 = 0.9$, between value-added and intermediates to be $\theta_1 = 0.5$, and across sectors in intermediates to be $\theta_2 = 0.2$. Our results are not particularly sensitive to these choices.

We use the same within-sector elasticities for all sectors: $\varepsilon_{\mathcal{I}} = \varepsilon$ and $\gamma_{\mathcal{I}} = \gamma$ and consider two scenarios: (1) every sector is assumed to be IRS with scale elasticity $\gamma$; (2) every sector is assumed to be DRS with scale elasticity $\varepsilon$. In either case, we consider two

different scale elasticities, in the DRS case, we set $\varepsilon = 0.875$ or $\varepsilon = 0.75$. In the IRS case, we set $\gamma = 0.875$ or $\gamma = 0.75$, which corresponds to a within-industry elasticity of substitution of 8 or 4 respectively.

Finally, we experiment with different ways of modeling entry: no entry, entry using primary factors, and entry using primary factors and goods (in the same way as variable production). The model without entry can be thought of as a short-run model and the model with entry as a long-run model.

## 8.2 Social Costs of Distortions

We solve the model nonlinearly and compute the efficiency loss from misallocation. We report the numbers as the percentage gain in welfare achieved by implementing optimal policy starting from the decentralized equilibrium outcome. The results are in Table 1 for different combinations of assumptions regarding entry and returns to scale. Across the board, the benchmark calibration shows that the losses from inefficiency are higher (roughly double) when we allow entry than when we do not, refuting the notion that endogenizing entry necessarily reduces the social cost of markups.

**Decomposing the Results.** For each calibration, Table 1 decomposes the sources of the distance to the frontier. The "Level only" row eliminates the dispersion of markups within each sector by setting all markups within each sector equal to the harmonic average of markups in that sector. The "Dispersion only" row rescales the level of markups in the data so that their harmonic average within each sector is equal to the efficient level (so sectoral markups are equal to the Dixit and Stiglitz (1977) markups when we adopt the IRS benchmark and equal to one when we adopt the DRS benchmark) but keeps their dispersion constant.

When there is no entry, almost the entirety of the losses are explained by the dispersion effect. The losses due to the dispersion effect are due to misallocation across firms within sectors, and are large because markups are very dispersed within sectors and because the relevant elasticities within sectors are large. The losses due to the level effect, when there is no entry, are entirely due to misallocation across sectors, and are small because markups are not so dispersed across sectors and because the cross-sectoral elasticities of substitution are low.

When there is entry, the level effect becomes comparable to the dispersion effect. The losses due to the level effect now also reflect misallocation between entry and variable production within sectors, and these losses are large because markups are in general too

| IRS, $\gamma = 0.875$ | No Entry | Entry Uses Factors | Entry uses Goods and Factors |
|---|---|---|---|
| Level only | 4.6% | 14% | 10% |
| Dispersion only | 30% | 30% | 30% |
| Benchmark | 36% | 50% | 41% |
| IRS, $\gamma = 0.75$ | | | |
| Level only | 4.6% | 17% | 20% |
| Dispersion only | 22% | 23% | 20% |
| Benchmark | 19% | 32% | 37% |
| DRS, $\varepsilon = 0.875$ | | | |
| Level only | 1.5% | 7.8% | 7.6% |
| Dispersion only | 23% | 23% | 23% |
| Benchmark | 26% | 35% | 32% |
| DRS, $\varepsilon = 0.75$ | | | |
| Level only | 0.8% | 9.5% | 10% |
| Dispersion only | 9.2% | 9.2% | 9.2% |
| Benchmark | 9.6% | 19% | 20% |

Table 1: Efficiency losses from misallocation. Firm-level returns to scale $\gamma = 0.875$ under IRS corresponds to elasticity of substitution across firms within sectors equal to 8, whereas $\gamma = 0.75$ corresponds to elasticity of substitution equal to 4.

large and because the relevant elasticities of substitution are large.

Whether entry only uses primary factors or also intermediates has ambiguous effects. Depending on the scale elasticities, the relative size of the gains can go either way. When the entry margin is more important ($\varepsilon$ and $\gamma$ are lower), the gains tend to be higher when entry also uses intermediates, consistent with the roundabout example in Section 7.3.2.

The efficiency losses are different in the IRS benchmark than in the DRS benchmark. This is because quantities are less elastic in the DRS economy and entry distortions are less costly. To understand the latter point, it is useful to think about the limit where $\varepsilon$ and $\gamma$ go to zero, which corresponds to a within-sector across-firm elasticity of one under IRS and a firm-level return to scale of zero under DRS. In this limit, under IRS, the efficiency losses become infinite because love of variety becomes extreme and so do the distortions in entry, as can be seen in Proposition 3. By contrast, under DRS, the efficiency losses go to zero as made clear by Propositions 3 and 4.

**Comparison to Simplified Models.** Our analysis contends that careful modelling of the details of the production network and the entry technology is qualitatively important.

To illustrate this quantitatively, in Table 2, we compare the results of the benchmark model to simplified versions of the model that employ some commonly used shortcuts: ignoring intermediate goods in production or entry (assuming no input-output); using a single-sector economy but allowing for intermediates (roundabout economy); ignoring firm-level heterogeneity within sectors (no firm heterogeneity). We discuss each of these strawmen in turn.

| IRS | No Entry | Entry Uses Factors | Entry uses Goods/Factors |
|---|---|---|---|
| Benchmark | 36% | 50% | 40% |
| No Input-Output | 16% | 20% | – |
| Roundabout | 139% | 182% | 133% |
| Homogeneous Firms | 4.6% | 14% | 10% |
| DRS | | | |
| Benchmark | 26% | 35% | 32% |
| No Input-Output | 13% | 18% | – |
| Roundabout | 91% | 123% | 108% |
| Homogeneous Firms | 1.0% | 7.8% | 7.6% |

Table 2: Efficiency losses from misallocation when different disaggregated aspects of the economy are trivialized. We use firm-level returns to scale $\varepsilon = 0.875$ under DRS, and $\gamma = 0.875$ under IRS. For IRS, this corresponds to an elasticity of substitution across firms within industries of 8.

The "No Input-Output" economy assumes away intermediates, and calibrates the size of each industry to be equal to its value-added share. Without entry, this economy undershoots the benchmark model for reasons discussed by Jones (2011) or Baqaee and Farhi (2019a). The undershooting becomes even more extreme once we allow for entry, underscoring even more strongly the need to model input-output linkages.

The "Roundabout" economy assumes that all firms in the economy belong to a single sector. The output of this sector is used both as the consumption good and as an intermediate input into production. This is a commonly used shortcut for incorporating intermediate inputs into a model. The one-sector roundabout economy overshoots the benchmark by a large amount. This is to be expected since the roundabout economy aggregates all firms in the economy into a single sector. This means cross-sectoral dispersions in markups (which are less costly than within-sectoral dispersions) are treated as if they are within-sectors. Intuitively, dispersed markups now distort input choices across producers by more, since firms in two different industries are treated as if they are highly substitutable.

Finally, the "Homogeneous Firms" economy assumes that all firms in a sector are identical, with the same productivity shifter and the same markup equal to the sectoral markup. The homogeneous sectors economy undershoots the benchmark by a large amount because even though it accounts for cross-sectoral distortions, it abstracts away from within-sector misallocation.
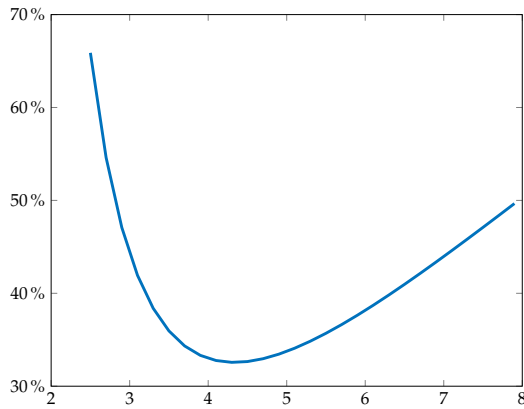
All in all, the sensitivity of these numbers underscores the quantitative importance of modelling and measuring the details as best we can.

**Role of the Elasticity of Substitution Across Firms Within Sectors.** In many models of misallocation without entry, for example (e.g. Hsieh and Klenow, 2009; Baqaee and Farhi, 2019a), the distance to the frontier increases with the elasticity across firms within sectors. As discussed in Section 7.2, this intuition fails when there is entry and markets are of the IRS type.
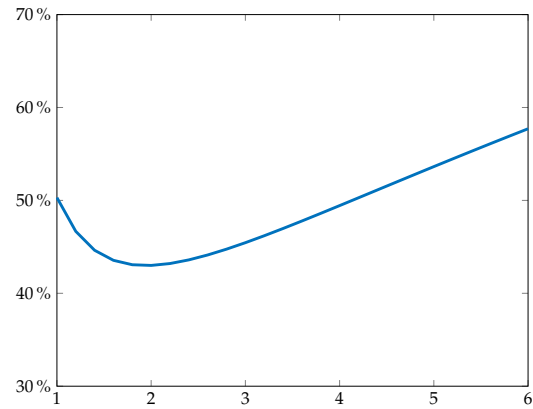
Figure 2a shows that for the IRS benchmark, the distance to the frontier is U-shaped as a function of the within-sector elasticity of substitution $1/(1-\gamma)$. For instance, the losses are 50% when $1/(1-\gamma) = 8$. This number falls to 32% when the elasticity is lowered to 4, before rising to close to 65% when the elasticity is lowered further to 2.5. This is consistent with the theoretical discussion in the last example of Section 7.2. Intuitively, with non-trivial input-output linkages, a lower elasticity reduces misallocation along the intensive margin, but magnifies misallocation along the extensive margin. In the limit where the elasticity goes to one ($\gamma$ goes to zero), misallocation along the extensive margin becomes infinitely costly.

**Role of Barriers to Entry.** In our benchmark specifications with entry, we assume that all rents are quasi-rents rather than pure rents. That is, free-entry holds in every sector. However, it is plausible that, even in the long run, profits are not entirely offset by the costs of entry. For example, it may be that resources spent on entry are less than profits due to barriers to entry from regulations or due to anti-competitive strategic deterrence. We capture these barriers to entry in reduced form by introducing an entry tax/wedge.

Figure 2b displays the implied distance to the frontier as a function of the view that one takes on the size of entry barriers in the data, where the size of entry barriers are measured by the size of the implicit entry tax/wedge (a value of one means that there are no barriers to entry). Perhaps surprisingly, the efficiency losses are non-monotonic in the size of entry barriers. Intuitively, whether barriers to entry increase or decrease the estimated distance to the frontier depends on whether there is too little or too much entry in the equilibrium with no entry barriers. Our estimated markups are relatively high, which implies that if

42

| (a) Elasticity of substitution | (b) Entry wedges |

Figure 2: Efficiency losses for the benchmark IRS model when entry uses factors as a function of the within-industry elasticity of substitution and entry wedges for the benchmark IRS model.

there is free-entry, then there is too much entry in the equilibrium. As a result, if one takes the view that there are entry barriers in the data (so that there is less entry than implied by profits), then one is lead to a lower estimate of the distance to the frontier up (up to some point, after which, entry becomes inefficiently too low).

## 8.3   Bang for Buck of Marginal Policy Interventions

We end this section by considering the effect of a marginal policy intervention in the decentralized equilibrium. Figure 3 shows the bang-for-buck elasticity of aggregate output with respect to a marginal entry subsidy (a form of industrial policy) or markup reduction (a form of competition policy) in different industries. The elasticity is scaled by the revenues associated with the intervention, as in Section 7.4, to make the magnitudes comparable.

For this exercise, we focus on the IRS case where $\gamma = 0.875$. We consider two alternative calibrations: one where we set markups equal to their CES monopolistic values, and one where we set markups equal to their estimated values. We begin by discussing the case where all markups are set equal to their CES Dixit and Stiglitz (1977) values. Then we discuss the case where markups are equal to their estimated values in the data. In both cases, we abstract from endogenous changes in markups in response to the policy.[32]

The monopolistic-markups calibration is a useful starting point for understanding the

---

[32]Here, we assume that the policy maker can directly change the wedges. As pointed out by Gupta (2020), in practice, a linear tax may not be able to achieve this since firm-level wedges may respond to the policy instrument.

(a) Markup reduction for CES markups

(b) Entry subsidies for CES markups

(c) Markup reduction for estim. markups
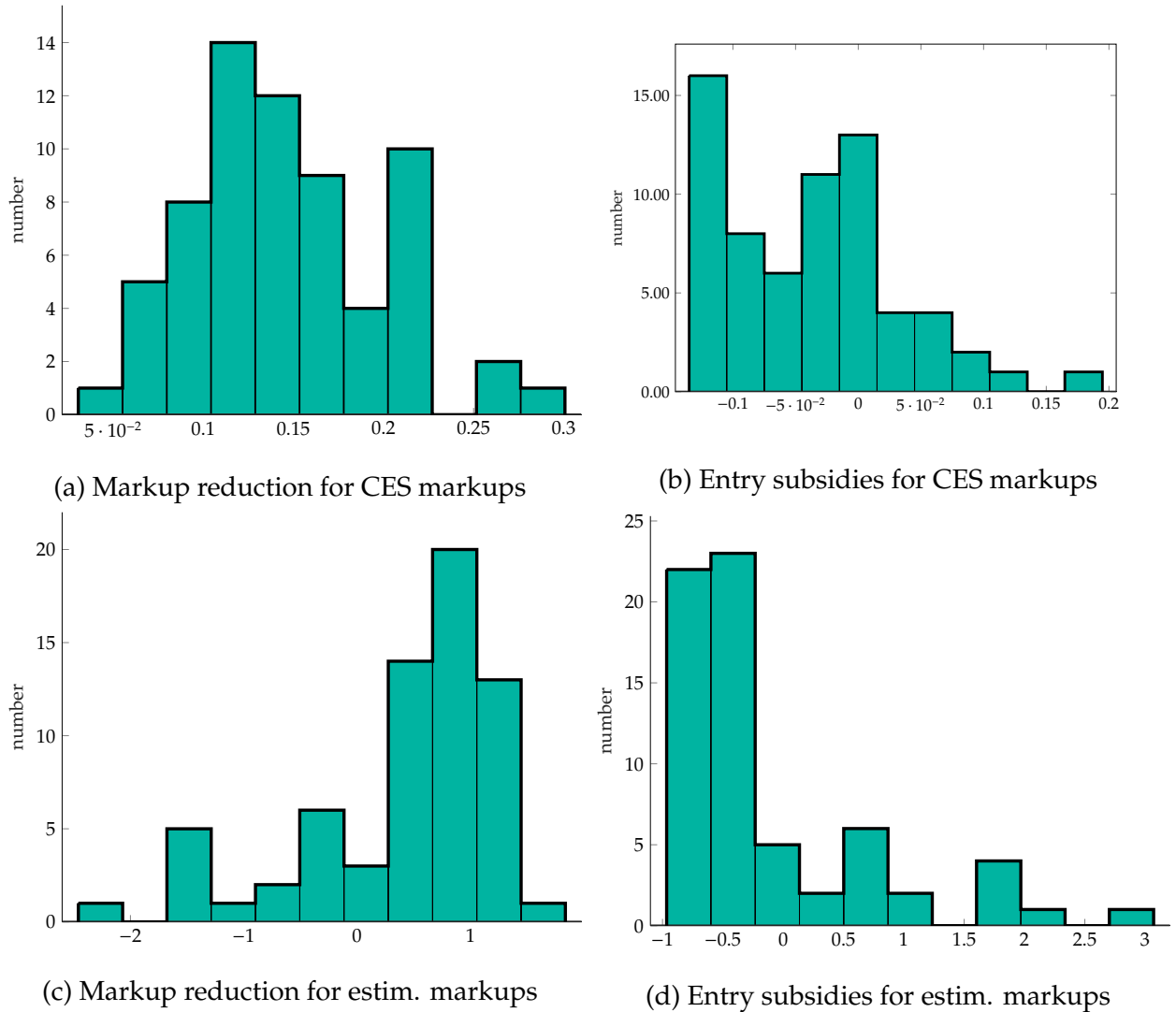
(d) Entry subsidies for estim. markups

Figure 3: The elasticity of output with respect to reductions in markups or an entry subsidy to different sectors normalized by the cost of the intervention. The top row uses CES markups, whereas the bottom row uses estimated PF markups.

results, since by setting markups to be the same in every sector, it helps isolate the role played by the input-output network on its own. In this case, markup reductions, plotted in Figure 3a, are always beneficial. Because we have imposed the same love-of-variety parameter in all sectors, the greatest bang-for-buck comes from reducing markups for those sectors with more complex supply chains, namely manufacturing industries like motor vehicles, metals, and plastics. Intuitively, reducing markups in these sectors allows more entry into their supply chains. The smallest gains come from those industries with the simplest supply chains, mostly service industries like housing or legal services but also primary industries like oil extraction or forestry. For entry subsidies, plotted in Figure 3b, the biggest gains, on the other hand, come from subsidizing those industries which are upstream in complex supply chains, namely primary industries like forestry, oil, and

mining, whereas subsidizing entry into relatively downstream industries, like nursing, hospitals, or social assistance, is actually harmful.[33]

When we move to the estimated markups, plotted in Figures 3c and 3d, the shape of the input-output network is not the only determinant of the relative ranking of different industries, as now we must also consider whether each sector's markups are too high or too low on average relative to its external economies. Since, for simplicity, we have imposed the same love-of-variety effect in all sectors but we have estimated markups for each sector, we do not read too much into the exact relative ranking of the different industries.

However, these figures are still useful because they show that as we move farther away from the efficient frontier, which we do when we go from monopolistic markups to estimated markups, the potency of second-best policies increases dramatically. To see this, note that the elasticities in the top row are an order of magnitude smaller than the elasticities in the bottom row of Figure 3.

But the larger effect sizes are a mixed blessing. Once we are far away from the frontier, the scope for policy having unintended consequences also increases. Although there appear to be many free lunches available to policy makers, interventions can equally have large negative as well as positive effects. In other words, as implied by the theory of the second-best, interventions that seem sensible in isolation, like reducing markups, can reduce output once we are deep inside the frontier.

# 9   Conclusion

Traditional theories of aggregation, by relying on aggregate envelope theorems, imply that the aggregate production function can be treated like a black-box machine whose contents are irrelevant to a first order. Aggregate productivity changes are simply the sales-weighted averages of the exogenous microeconomic productivity shocks. Under this view, these exogenous changes in aggregate productivity are responsible for a large fraction of both the cycle and the trend in aggregate output.

For inefficient economies, this first-order approach is untenable. In a disaggregated economy, where many different margins can be misallocated, total factor productivity is endogenous and affected, to a first order, by reallocation effects. Furthermore, unlike exogenous productivity shocks, which are likely to be gradual and positive, reallocation effects can be abrupt and have either sign. This paper shows that these reallocation effects

---

[33]For more intuition about this, in Appendix L, we work through a Cobb-Douglas example.

can be very potent in the presence of non-convexities and entry, and provides a framework for studying them.

# References

Acemoglu, D. and P. D. Azar (2020). Endogenous production networks. *Econometrica 88*(1), 33–82.

Acemoglu, D., V. M. Carvalho, A. Ozdaglar, and A. Tahbaz-Salehi (2012). The network origins of aggregate fluctuations. *Econometrica 80*(5), 1977–2016.

Acemoglu, D. and A. Tahbaz-Salehi (2020). Firms, failures, and fluctuations: the macroeconomics of supply chain disruptions. Technical report, National Bureau of Economic Research.

Ackerberg, D. A., K. Caves, and G. Frazer (2015). Identification properties of recent production function estimators. *Econometrica 83*(6), 2411–2451.

Atalay, E. (2017). How important are sectoral shocks? *American Economic Journal: Macroeconomics (Forthcoming)*.

Autor, D., D. Dorn, L. Katz, C. Patterson, and J. Van Reenen (2017). The fall of the labor share and the rise of superstar firms.

Baqaee, D. R. (2018). Cascading failures in production networks. *Econometrica 86*(5), 1819–1838.

Baqaee, D. R. and E. Farhi (2019a). Productivity and misallocation in general equilibrium. *The Quarterly Journal of Economics 135*(1), 105–163.

Baqaee, D. R. and E. Farhi (2019b). The Macroeconomic Impact of Microeconomic Shocks: Beyond Hulten's Theorem. *Econometrica 87*(4), 1155–1203.

Bartelme, D. G., A. Costinot, D. Donaldson, and A. Rodriguez-Clare (2019, August). The Textbook Case for Industrial Policy: Theory Meets Data. NBER Working Papers 26193, National Bureau of Economic Research, Inc.

Behrens, K., G. Mion, Y. Murata, and J. Suedekum (2016). Distorted monopolistic competition.

Bigio, S. and J. La'O (2016). Financial frictions in production networks. Technical report.

Boehm, C., A. Flaaen, and N. Pandalai-Nayar (2014). Complementarities in multinational production and business cycle dynamics. Technical report, Working paper, University of Michigan.

Boehm, J. and E. Oberfield (2020). Misallocation in the market for inputs: Enforcement and the organization of production. *The Quarterly Journal of Economics 135*(4), 2007–2058.

Carvalho, V. M. and A. Tahbaz-Salehi (2018). Production networks: A primer.

Ciccone, A. and K. Matsuyama (1996). Start-up costs and pecuniary externalities as barriers to economic development. *Journal of Development Economics 49*(1), 33–59.

Claus, J. and J. Thomas (2001). Equity premia as low as three percent? evidence from analysts' earnings forecasts for domestic and international stock markets. *The Journal of Finance 56*(5), 1629–1666.

De Loecker, J., J. Eeckhout, and G. Unger (2019). The rise of market power and the macroeconomic implications. Technical report.

Dhyne, E., A. K. Kikkawa, M. Mogstad, and F. Tintelnot (2021). Trade and domestic production networks. *The Review of Economic Studies 88*(2), 643–668.

Dixit, A. K. and J. E. Stiglitz (1977). Monopolistic competition and optimum product diversity. *The American Economic Review*, 297–308.

Domar, E. D. (1961). On the measurement of technological change. *The Economic Journal 71*(284), 709–729.

Edmond, C., V. Midrigan, and D. Y. Xu (2018). How costly are markups? Technical report, National Bureau of Economic Research.

Elliott, M., B. Golub, and M. V. Leduc (2020). Supply network formation and fragility. *Available at SSRN 3525459*.

Epifani, P. and G. Gancia (2011). Trade, markup heterogeneity and misallocations. *Journal of International Economics 83*(1), 1–13.

Gabaix, X. (2011). The granular origins of aggregate fluctuations. *Econometrica 79*(3), 733–772.

Georgescu-Roegen, N. (1951). Some properties of a generalized leontief model. *Activity Analysis of Allocation and Production. John Wiley & Sons, New York*, 165–173.

Grossman, G. M. and E. Helpman (1991). *Innovation and growth in the global economy*. MIT press.

Gupta, A. (2020). Firm heterogeneity, demand for quality and prices: Evidence from india. Technical report.

Gutiérrez, G. and T. Philippon (2016). Investment-less growth: An empirical investigation. Technical report, National Bureau of Economic Research.

Harberger, A. C. (1954). Monopoly and resource allocation. In *American Economic Association, Papers and Proceedings*, Volume 44, pp. 77–87.

Harberger, A. C. (1964). The measurement of waste. *The American Economic Review 54*(3), 58–76.

Herrendorf, B., R. Rogerson, and A. Valentinyi (2013). Two perspectives on preferences and structural transformation. *American Economic Review 103*(7), 2752–89.

Hirschman, A. O. (1958). *The strategy of economic development*, Volume 58. Yale University Press New Haven.

Hopenhayn, H. A. (1992). Entry, exit, and firm dynamics in long run equilibrium. *Econometrica*, 1127–1150.

Hsieh, C.-T. and P. J. Klenow (2009). Misallocation and manufacturing TFP in China and India. *The quarterly journal of economics 124*(4), 1403–1448.

Hulten, C. R. (1978). Growth accounting with intermediate inputs. *The Review of Economic Studies*, 511–518.

Jones, C. I. (2011). Intermediate goods and weak links in the theory of economic development. *American Economic Journal: Macroeconomics*, 1–28.

Jones, C. I. (2013). Input-Output economics. In *Advances in Economics and Econometrics: Tenth World Congress*, Volume 2, pp. 419. Cambridge University Press.

Kikkawa, A. K., G. Magerman, E. Dhyne, et al. (2018). Imperfect competition in firm-to-firm trade. Technical report.

Kimball, M. S. (1995). The quantitative analytics of the basic neomonetarist model. *Journal of Money, Credit, and Banking 27*(4).

Krugman, P. R. (1979). Increasing returns, monopolistic competition, and international trade. *Journal of international Economics 9*(4), 469–479.

La'O, J. and A. Tahbaz-Salehi (2020). Optimal monetary policy in production networks. Technical report, National Bureau of Economic Research.

Levinsohn, J. and A. Petrin (2003). Estimating production functions using inputs to control for unobservables. *The review of economic studies 70*(2), 317–341.

Lim, K. (2017). Firm-to-rm trade in sticky production networks.

Liu, E. (2017). Industrial policies and economic development. Technical report.

Long, J. B. and C. I. Plosser (1983). Real business cycles. *The Journal of Political Economy*, 39–69.

Lucas, R. E. (1978). On the size distribution of business firms. *The Bell Journal of Economics*, 508–523.

McKenzie, L. W. (1959). On the existence of general equilibrium for a competitive market. *Econometrica: journal of the Econometric Society*, 54–71.

Melitz, M. J. (2003). The impact of trade on intra-industry reallocations and aggregate industry productivity. *Econometrica 71*(6), 1695–1725.

Oberfield, E. (2017). A theory of input-output architecture. Technical report.

Olley, G. S. and A. Pakes (1996). The dynamics of productivity in the telecommunications equipment industry. *Econometrica: Journal of the Econometric Society*, 1263–1297.

Osotimehin, S. and L. Popov (2017). Misallocation and intersectoral linkages. Technical report, Mimeo.

Restuccia, D. and R. Rogerson (2008). Policy distortions and aggregate productivity with heterogeneous establishments. *Review of Economic dynamics 11*(4), 707–720.

Romer, P. M. (1987). Growth based on increasing returns due to specialization. *The American Economic Review 77*(2), 56–62.

Rubbo, E. (2020). Networks, phillips curves and monetary policy. Technical report, mimeo, Harvard University.

Samuelson, P. A. (1951). Abstract of a Theorem Concerning Substitutability in Open Leontief Models. In T. Koopmans (Ed.), *Activity Analysis of Production and Allocation*, New York. Wiley.

Taschereau-Dumouchel, M. (2020). Cascades and fluctuations in an economy with an endogenous production network. *Available at SSRN 3115854*.

Traina, J. (2018). Is aggregate market power increasing? production trends using financial statements. Technical report.

Vincent, N. and M. Kehrig (2017). Growing productivity without growing wages: The micro-level anatomy of the aggregate labor share decline. In *2017 Meeting Papers*, Number 739. Society for Economic Dynamics.