

When do employers share? Rent sharing, monopsony and minimum wages

Ihsaan Bassier* and Joshua Budlender†

November 27, 2025

[Click here for latest version](#)

When firm productivity or product demand rises, workers typically share in the gains through higher wages or expanded employment. We show that for firms under monopsony with a binding minimum wage, this link from firm gains to worker outcomes breaks sharply. Revenue-productivity improvements raise revenues but not wages or employment: firms simply maintain the minimum wage and absorb the gains into higher wage markdowns. We find compelling evidence for these predictions using South African administrative data, based on a cross-sectional kink design as well as within-firm responses to internal and shift-share trade shocks. These results reveal a previously overlooked monopsonistic margin—productivity-induced markdown adjustment—and we show using a structural model that this substantially diminishes the intended returns of policies such as employment subsidies.

Keywords: Monopsony, Rent-sharing, Minimum wage, Firm productivity.

JEL codes: D22, J31, J38, J42, O33.

Author ordering is alphabetical; the authors contributed equally to this work. For comments and suggestions we owe particular thanks to Arindrajit Dube, as well as David Berger, Daniele Girardi, Leonard Goff, C. Friedrich Kreuser, Attila Lindner, Alan Manning, Suresh Naidu and Matt Woerman. We are grateful for comments from participants of several conferences (ASSA 2024, ESSA 2023), workshops (UCL-IFS 2025, Princeton Global Minimum Wages 2024, SALDRU December 2024, Pretoria-Stellenbosch PhD 2024, UNU-WIDER Work-in-Progress 2023), and working groups (University of Massachusetts Amherst PhD Labor-Development 2019, Columbia University PhD 2020-2021). We also thank the NT-SDF and SA-TIED teams for their assistance, particularly Ayanda Hlatshwayo, Murray Leibbrandt, Michael Kilumelume, Carol Newman and John Rand. We gratefully acknowledge funding from UNU-WIDER, and from the research initiative “Structural Transformation and Economic Growth” (STEG), contract reference STEG.LOA.2206.Budlender. A previous version of this paper was posted in 2023 under the title “Rent sharing, wage floors and development.”

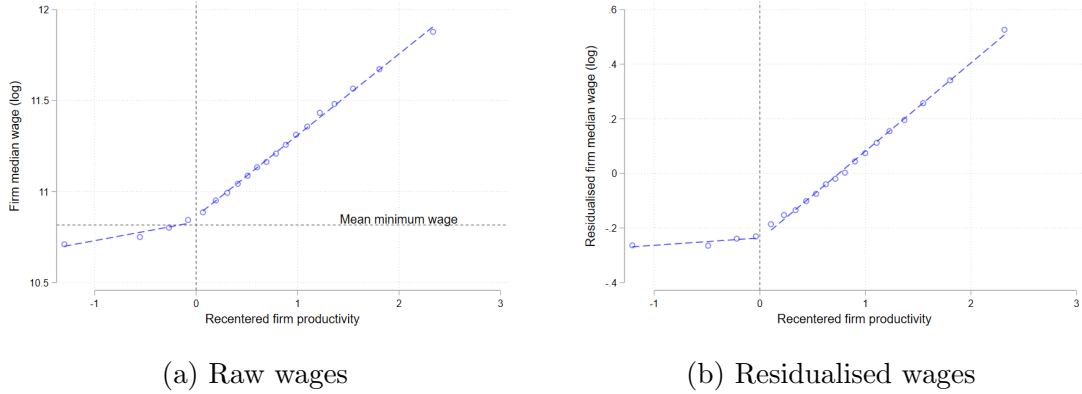
*Economics, University of Surrey; Centre for Economic Performance, London School of Economics & Political Science; Southern Africa Labour and Development Research Unit (SALDRU), University of Cape Town (UCT). Email: i.bassier@surrey.ac.uk

†SALDRU, UCT. Email: joshua.budlender@uct.ac.za

1 Introduction

We begin with an intuitive stylized fact: a large segment of firms with low but varying productivity pay around the minimum wage, with a concomitantly low cross-sectional rent-sharing (or wage pass-through) elasticity. This contrasts with higher productivity firms, which do appear to share rents. Figure 1 shows this using South African administrative data for raw firm median wages (panel (a)), and even more starkly when controlling for industry-location fixed effects and worker quality (panel (b)). Below a particular productivity threshold, firm median wages do not increase with productivity. Above that threshold, there is a kink in the wage-productivity curve and the expected pass-through dynamics commence. Once one knows to look for it, this kink is surprisingly common: it is evident in Card et al.’s (2016) normalization method for the gender wage gap in Portugal, and also in Palladino et al.’s (2025) application of the same normalization method for nine European countries. Why are firms not sharing rents in this region and why does rent-sharing behaviour change at this threshold? What does this imply for other firm outcomes such as employment and profits?

Figure 1: A kink in the wage–productivity curve



Notes: Figure shows the firm-level wage–productivity curve, using raw median wages at the firm (Panel (a)) and these wages residualised on industry–locations and controlling for worker quality (Panel (b)). The x-axis is firm-specific total factor productivity, recentered around an estimated threshold (vertical dashed line) that divides supply-constrained firms (where the minimum wage binds) from those unconstrained by the minimum wage. Productivity estimation and the recentering procedure are described in Section 4. The horizontal dashed line is the mean minimum wage across firms (firms may be subject to different minima; see Section 3).

We take off from two literatures. The first explains firm rent sharing as an outcome of monopsonistic labour markets (e.g. Card et al. 2018; Kline et al. 2019; Lamadon et al. 2022). The second seeks to understand the direct employ-

ment or welfare effects of minimum wages under monopsony ([Dickens et al. 1999](#); [Engbom and Moser 2022](#); [Berger et al. 2025](#)). We show that incorporating the labour market constraint of a binding minimum wage into a very general and simple monopsony model explains the stylised fact we started with, and has striking further implications for firm behaviour. Specifically, we predict that just below the kink-point in the wage-productivity curve, minimum-wage bound firms will absorb favorable revenue-productivity shocks into increased markdowns, instead of increasing wages or employment, and firms will only increase wages and employment with productivity once markdowns have reached optimal monopsony levels. We test these predictions in South African administrative data using a cross-sectional kink design as well as within-firm responses to internal and shift-share trade shocks, and find strong evidence consistent with these predictions. We show that this mechanism is likely to be empirically important for a range of countries and policies where minimum wages are substantially binding.

In the model, the kink-point in the wage-productivity curve defines a productivity threshold dividing minimum wage-constrained and -unconstrained firms. The intuition for just-constrained firms absorbing revenue-productivity increases into markdowns starts with noting that the minimum wage statically redistributes rents from firms to workers at these low productivity constrained firms. These firms therefore receive a lower markdown than they would in the absence of the minimum wage. Higher productivity unconstrained firms, in contrast, optimally markdown wages and receive excess rents. There is consequently a productivity region in which just-constrained firms absorb revenue-productivity growth as increased rents until they have reclaimed the markdown associated with the unconstrained productivity level. Firms above this productivity threshold “share” revenue-productivity gains with workers in the form of higher wages, because in a monopsony model this is the prerequisite for firm expansion.

Observationally, this rent-sharing pattern is not limited to monopsony models, and we show that these striking wage and markdown patterns come out clearly from other models of imperfect competition such as Diamond-Mortensen-Pissarides or DMP ([Mortensen and Pissarides 1999](#)). However, we also establish a prediction more peculiar to monopsony models, which is that such just-constrained firms also do not increase employment in response to productivity increases, while their wage markdowns rise. The intuition here is that such firms cannot attract additional workers without increasing the wage, but they will not increase it above the minimum wage until their markdown has been restored to the optimal unconstrained monopsony level.

The key empirical predictions of the model, then, are differential employment and profit responses to revenue-productivity shocks along the firm productivity distribution, with a break at the point where wage pass-through commences.¹ After estimating firm-specific productivity using [Akerberg et al. \(2015\)](#), we test this using a kink design focusing on the productivity threshold identified on the wage-productivity curve. The empirical patterns fit the model predictions remarkably well: profit-share steeply increases with productivity below the threshold, and then is more constant above it; and employment does not change much with productivity below the threshold, but increases strongly above it.

We then complement this kink design by considering how firms on either side of the productivity threshold *respond* to revenue-productivity shocks, replicating and extending leading approaches in the literature. Following [Lamadon et al. \(2022\)](#), we use an “internal instrument” method, which entails constructing a stacked event study where firm-specific treatment is defined as an unusually large observed increase in firm value-added. We also use an “external instrument” based on a shift-share variable we construct from firm-specific shares of destination-country exports and imports (the “share”) and destination-country GDP movements (the “shift”), similar to [Garin and Silvério \(2023\)](#). Our results strongly support our theoretical predictions: compared to responses in the unconstrained region, in the constrained region the wage response (the rent-sharing elasticity) is approximately 30% lower, the profit share response is almost 3 times higher, and the employment response is approximately 25% lower. Estimates by productivity bin support the prediction that the break in response size occurs around the productivity threshold. The estimates are remarkably similar using the external and internal instruments given how different the sources of variation are, which builds confidence in a causal interpretation. We conduct a variety of robustness tests.

Our primary contribution is to the rent-sharing and monopsony literatures. Theoretically and empirically, we identify a novel region where firms bound by minimum wages react to positive revenue-productivity shocks with muted wage and employment responses, and instead increase their profit share. This speaks to the broader economic question of the extent to which workers benefit from firm growth, and we highlight that this process is not automatic for a potentially large subset of firms. A related contribution is that, as noted above, these patterns support the existence of monopsonistic mechanisms in the labour market.

¹Because markdowns are unobserved, in our empirical specification we approximate markdowns with a measure of the “profit share”, as one minus the labor share as it is defined in [Gouin-Bonenfant \(2022\)](#). In our structural model, we formally show that the model-implied predictions for the markdown and profit share are very similar.

The inconsistency in rent sharing also highlights an interesting distinction between two related concepts in the literature (Bell et al. 2024; Kline et al. 2019; Risch 2024): the average rent-sharing level (wage over marginal revenue product), versus its marginal elasticity or pass-through (change in wage with respect to revenue-productivity). In many cases, these two concepts do closely align, as is implicitly assumed in much of the literature. However, our model and results highlight that when the minimum wage increases, the rent-sharing level rises, yet the pass-through declines. A similar distinction holds for employment: for firms just-constrained by minimum wages the level of employment is higher than otherwise, even as the employment response to productivity increases is low.

A secondary contribution is that our estimated firm-specific labour-supply and pass-through elasticities add to very few such estimates for developing countries, where higher labour surpluses and frictions may increase monopsony power (Armangué-Jubert et al. 2025; Bassier 2023). More generally, a popular approach in the literature estimating labor supply elasticities is to estimate wage and employment responses to a revenue-productivity shock (Kroft et al., 2025; Kline, 2025). We note such studies should take care to estimate these elasticities for unconstrained firms, as when shocks take place across firms constrained, for example, by minimum wages, then the estimate will not identify the labour supply elasticity.

Our findings are likely to be applicable to a broad range of countries with high minimum wages. The mechanism we identify may be particularly important for countries seeking to develop through increasing firm productivity while protecting workers with high minimum wages (e.g. Brazil and South Africa), as well as low-income countries with relatively high wage floors associated with subsistence or efficiency wages (Breza et al. 2021; Muralidharan et al. 2023). We caution that a potentially non-trivial fraction of firms may absorb such productivity gains entirely as profits, undermining this development agenda of inclusive growth (Lewis 1954; Verhoogen 2023). To illustrate this, we estimate a structural model using our reduced form estimates fitted to the data. We find that the mechanism we identify substantially diminishes the returns of popular industrial and labour market policies which upgrade firm productivity or subsidize employment.

The core ideas of the model are introduced in Section 2, and the data and context are discussed in Section 3. The evidence using the cross-sectional kink design is presented in Section 4, then the within-firm shock evidence in Section 5. Section 6 structurally estimates and simulates policy effects, and Section 7 discusses alternative models of the labour market. Section 8 concludes.

2 Theoretical predictions

2.1 Simple model of monopsonistic firms with a minimum wage

Standard argument. We briefly recapitulate a simple model of firm responses to minimum wages under monopsony, following [Dickens et al. \(1999\)](#) and [Manning \(2003, pp. 338-345\)](#).

As is characteristic of monopsony models, we assume a firm-facing labour supply curve $w_i = \varepsilon n_i$, where lower case letters denote logs, w_i is the firm wage, n_i is firm employment, and there are many firms. The firm-facing labour supply elasticity $1/\varepsilon$ is constant across firms and is finite. Such an upwards sloping labor supply curve implies a marginal cost of labor greater than the wage for firm i :²

$$\text{mcl}_i = \ln(1 + \varepsilon) + w_i = \ln(1 + \varepsilon) + \varepsilon n_i. \quad (1)$$

The marginal revenue product of labour of firm i , is a simple downwards sloping labour demand curve:

$$\text{mrpl}_i = a_i - \eta n_i, \quad (2)$$

where a_i is a demand or productivity shifter. The elasticity of the labour demand curve under perfect competition would be $1/\eta$. This can be motivated by a production function such as $Y_i = \frac{1}{1-\eta} A_i N_i^{1-\eta}$, where additional factors such as capital can be log-additively included without loss of generality.

This model setup represents a very simple and general monopsonistic form, and remains agnostic as to the source of monopsony power (e.g. search frictions, amenities or concentration). As such it nests different popular approaches to modeling monopsonistic competition ([Azar and Marinescu, 2024](#); [Card et al., 2018](#); [Manning, 2003](#)). Setting marginal product equal to marginal cost, the unconstrained employment and wage for firm i are:

$$n_i^* = \frac{1}{\varepsilon + \eta} (a_i - \ln(1 + \varepsilon)) \quad (3)$$

$$w_i^* = \varepsilon n_i = \frac{\varepsilon}{\varepsilon + \eta} (a_i - \ln(1 + \varepsilon)). \quad (4)$$

When a minimum wage is introduced, a firm finds itself in one of three quali-

²The intuition for this is that for a monopsonist to hire an additional worker they must increase the wage, which also applies to the wages of already-employed workers. $\frac{\partial WL}{\partial L} = \frac{\partial W}{\partial L} L + W = \varepsilon W + W$.

tatively distinct regimes, depending on its productivity a_i . These three regimes are depicted in Figure 2 panel (a) as the three marginal revenue product of labour (MRPL) curves, corresponding to varying a_i in equation 2. The labour supply (LS) and marginal cost of labour (MCL) curves are as described above per equation 1.

If the minimum wage w_m is not binding, i.e. $w_m \leq w_i^*$, equations 3 and 4 hold and $w_i = w_i^*$ and $n_i = n_i^*$. These are *unconstrained* firms, shown in the first regime (MRPL1) in Figure 2 panel (a), and these firms are not directly affected by the minimum wage.

When the unconstrained wage is lower than the minimum wage ($w_m > w_i^*$), firms must pay a wage equal to the minimum wage ($w_i = w_m$). These are *constrained* firms. Within the constrained firms, there are two regimes. Firms for which $w_m > w_i^*$ but which have their marginal revenue product of labor *above* the minimum wage are *supply-constrained*, shown as regime two (MRPL2) in Figure 2 panel (a). We have the well-known result that for such firms, employment *increases* as a result of the minimum wage, because they attract additional workers at this higher wage, and accept all workers supplied at that wage, $n_i = (1/\varepsilon) w_m$. To allow for a more general framing beyond monopsonistic firms (see Section 7), we sometimes also refer to these firms as *just-constrained* firms.

The other regime of constrained firms are those with $w_m > \text{mrpl}_i$, denoted *demand-constrained* and shown as regime three (MRPL3) in Figure 2 panel (a). These firms *reduce* employment in response to the minimum wage until $\text{mrpl}_i = w_m$. Firms must still pay the minimum wage, but the new employment level is now governed by the firm labour demand constraint, so that $n_i = (1/\eta) (a_i - w_m)$.

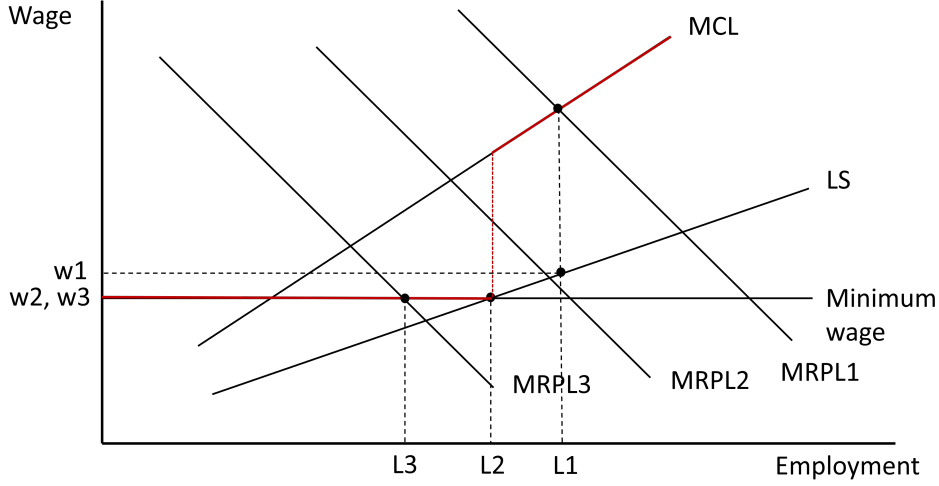
The more general model, presented in Appendix B.1, incorporates the average market-level wage as a determinant of aggregate labor supply, and allows for a firm-specific labour supply shifter (e.g. disamenities) as well as the firm-specific demand-shifter above (a_i)—but the qualitative regimes are exactly the same.

Novel insight. We first note that the introduction of a minimum wage changes the “latent” MCL curve to a new “effective” MCL curve indicated by the discontinuous red line shown in Figure 2 panel (a). Since wages cannot be below the minimum wage, and the cost of hiring an additional worker is simply the minimum wage paid to that worker, the marginal cost of labour when the LS curve is below the minimum wage is simply the minimum wage itself. Of particular interest to us is the discontinuity in the effective MCL at L2 on the employment axis, where firms switch from being minimum wage-constrained to unconstrained.

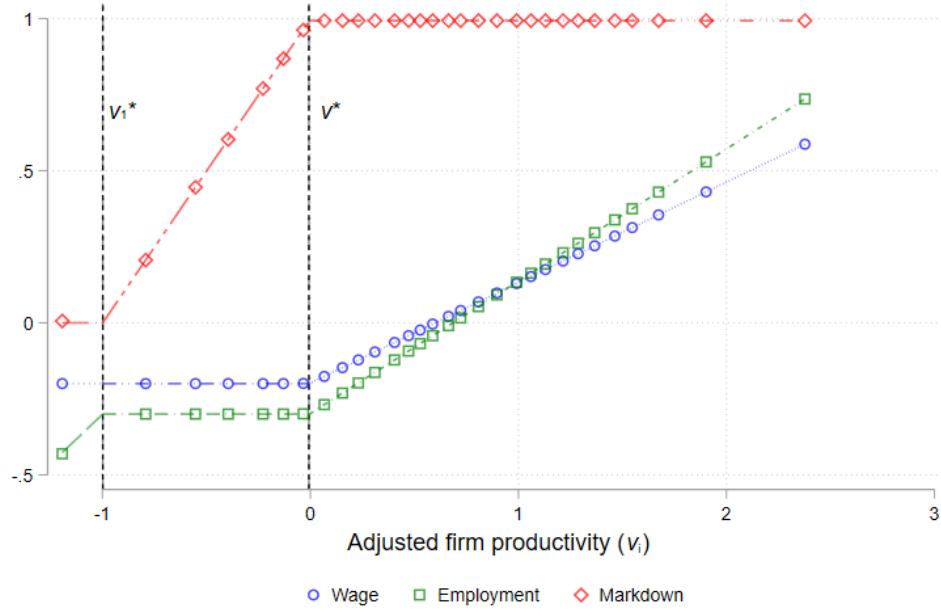
Firm-specific employment (on the x-axis) is set where the MRPL curve intersects with the effective MCL curve, while the wage is marked down from this point

Figure 2: Model of monopsonistic firms of differing productivity with a minimum wage

(a) Three revenue-productivity regimes in the presence of a minimum wage



(b) Simulation: Wage, employment and markdowns against productivity



Notes: Panel (a) shows the standard graphical depiction of monopsonistic firms facing minimum wages, adapted from [Manning \(2003, p. 343\)](#), for firms of three different productivity levels (MRPL1, MRPL2, MRPL3). The red line shows the “effective” MCL in the presence of the minimum wage, which is discontinuous at L2. Panel (b) shows a simulation of the simple model in Section 2.2 with each firm’s “adjusted productivity” on the x-axis. Outcomes are the firm wage (blue), firm employment (green), and markdown (red). All outcome values are in logs. Wages and employment are normalized relative to their average values in the perfectly competitive case with no minimum wage. The dotted vertical lines indicate the boundaries of the different regions: demand-constrained firms are in the left-most region (2.6% firms), supply-constrained firms in the centre (20.4%), and unconstrained firms in the right-most region (77%). See Section 2.2 for simulation details.

(along the y-axis) to the level on the labour supply curve. For firms with productivity MRPL2, this MRPL-MCL intersection occurs in the region of the effective MCL discontinuity. These are the supply-constrained firms. It is easy to see that for local shifts in MRPL2, the intersection with the effective MCL curve remains at the discontinuity point, and subsequently that these shifts in productivity do not change firm employment. They also do not change the wage, which is marked down to the minimum wage level. Instead, local shifts in MRPL in this region are reflected as changes in the size of the markdown from the marginal revenue productivity of labour to the (minimum) wage. The intuition is that firms with a range of MRPL intersecting the MCL at its discontinuity cannot attract additional workers without increasing the wage above the minimum wage (because the LS curve is now above the minimum wage for $n_i > L2$), but as long as their unconstrained wage (read on the labour supply curve in the region $n_i \leq L2$) is below the minimum ($w_m > w_i^*$), there is no incentive for them to do so. Instead the additional productivity per worker is reflected in increased markdowns. The reason there is a discontinuously large cost associated with increasing the wage above the minimum is because this wage increase would also apply to their existing workers currently paid at the minimum wage.

Our main insight, then, is that for supply-constrained firms there is a range of productivity increases (decreases) which do not change the firm’s wage or employment, and which are instead reflected in increases (decreases) in the markdown.

2.2 Simulations

We demonstrate these patterns by simulating the model above in Figure 2 panel (b), focusing on a fixed minimum wage with productivity on the x-axis and firm wages, employment and markdowns on the y-axis.³ Following Manning (2003), it is useful to focus on an “adjusted productivity” term we denote v_i , which determines the regime of a firm:

$$v_i = \frac{\varepsilon a_i}{\eta + \varepsilon} \quad (5)$$

This is just the variable component of equation 4, where firm productivity a_i is adjusted by the elasticities of firm labour supply and demand.

³We impose that MRPL shifters a_i follow a standard normal distribution, we set the market labour supply elasticity to 1.2 and the firm-facing labour supply elasticity to 1.3 (these are based on the cross-sectional patterns, see section 4). The simulations are based on 1,000 observations, each representing a firm. Wages and employment are normalized by comparison to the average wage and employment under the perfectly competitive case—that is, no monopsony nor minimum wage. The minimum wage is set at -0.2 log units. We trim the 1% tails of productivity.

Firms with v_i above some threshold v^* will have $w_i^* \geq w_m$ and will be unconstrained, firms with v_i below v^* but above another threshold v_1^* will be supply-constrained, and firms with v_i below v_1^* will be demand-constrained.⁴ Depending on the value of v_i , then, Figure 2 shows the three regimes:

1. Unconstrained (i.e. higher productivity, MRPL1): $v_i \geq v^*$ The right-most region, after the second vertical line (indicating v^*), delineates firms whose optimal monopsony wages are above the minimum wage, and so are not affected directly by the minimum wage. Wages and employment increase in productivity (equations 3 and 4), wages are marked down relative to MRPL as in the standard monopsony optimization, and the markdown level is constant.
2. Supply-constrained or “just-constrained” (i.e mid/lower productivity, MRPL2): $v^* > v_i \geq v_1^*$ The middle region between the vertical lines is our subject of interest, where the optimal wage is just below the minimum wage. Firms in this region keep wages fixed at the minimum, and do not increase employment as productivity increases. Instead, increased productivity is absorbed in higher markdowns, until the markdown is at the optimal level for an unconstrained monopsonist (at v^*).
3. Demand-constrained (i.e. very low productivity, MRPL3): $v_i \leq v_1^*$ The left-most region, before the first vertical line (indicating v_1^*), shows firms constrained to set wages equal the minimum wage, but with too low productivity to employ all of the workers available at that wage. These firms employ more workers as their productivity increases. MRPL is equal to the minimum wage and there is no markdown.

In Figure 2 panel (b), 20.4% of firms are supply-constrained and 2.6% are demand-constrained, which compare favorably with our empirical results discussed in Section 4. Appendix Figure A1 illustrates that in cases with less monopsony (ε smaller) and/or lower minimum wages, the supply-constrained region shrinks and moves down the productivity distribution, capturing fewer firms. The qualitative patterns are the same in the more detailed model, for example Appendix Figure A2 includes firm amenities with alternative assumptions regarding its correlation with firm productivity.

⁴We provide expressions for v^* and v_1^* in Appendix B.1. In the fuller model discussed in Appendix B.1, the term v_i also includes the firm-specific amenities shifter, which we abstract from for our purposes here.

2.3 Model predictions

What are the precise patterns predicted by the model above? In practical terms, even if the mechanisms outlined above are important, the empirical patterns will diverge from the simulations due to the influence of unmodeled factors and measurement error.

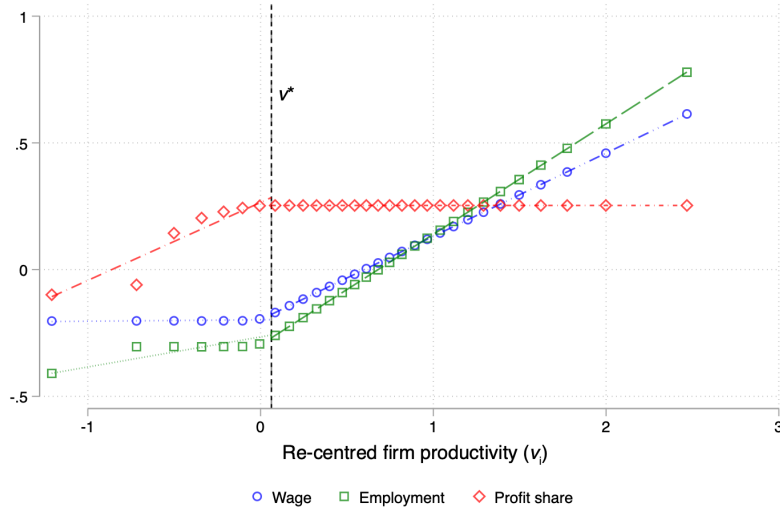
Perhaps the most important among these departures from the simulations is that we focus only on the region around v^* , the kink in the wage-productivity curve that divides the constrained and unconstrained firms. We provide more details in the next section, but a central point is to note that our institutional context involves about 40 different minimum wage systems, with correspondingly different threshold points. In our empirical exercise we pool these different regimes and re-center around the same adjusted productivity threshold to ensure enough statistical power to test our predicted patterns. We have to choose to recenter around either v^* or v_1^* , because there is little reason to expect a similar productivity range between these thresholds across minimum wage regimes. The v^* threshold is the more natural point as it contains the wage kink, and is the distinguishing point between our firms of interest (the supply-constrained firms) from the bulk of firms in the data, which will be unconstrained. The number of demand-constrained firms that we observe is likely to be small, because firms with very low productivity draws will be unobserved or under-represented in our actual firm data due to informality, fixed costs and endogenous exit (Olley and Pakes 1996; De Loecker and Syverson 2021).

Another point of divergence between the model predictions and empirical observation is that our key variables will be measured with error. Firm productivity will be estimated, introducing noise into the simulated kink-point v^* . The markdown is also unobserved, and so we use the gross profit share as a proxy. This is defined as gross profits over gross profits plus the firm wagebill, and so is equivalent to one minus the labor share as it is defined in Gouin-Bonenfant (2022). We (imprecisely) also refer to this as the “capital share”, borrowing from the macroeconomics literature which divides income into labor and capital income. We return to the correspondence between the markdown and profit share in our structural model in Section 6.

How do these divergences affect the model-predicted patterns? Figure 3 simulates a scenario much closer to our empirical context, where the simulation above (Figure 2 panel (b)) is repeated for 40 different labour markets (i.e. differing minimum wages and labour supply elasticities), and then firms are re-centered around

the labour market-specific wage-kink v^* . The key implication is that we do not observe a clean kink point within the constrained region that divides demand-constrained and supply-constrained firms. Yet the wage-kink (v^* , dividing just-constrained and unconstrained firms) is still very apparent and we still have differential slopes on either side of this kink point for wages, employment and profit share. The constrained region does have some positive slope in employment, because it contains both demand- and supply-constrained firms, but the key point is that it is much flatter than the unconstrained slope.

Figure 3: Pooled simulation demonstrating empirical predictions



Notes: The figure shows results from pooled model simulations, where the simulation of Section 2.2 (i.e. Figure 2 panel (b)) is run separately for 40 different labour markets with randomly varying labour supply elasticities and minimum wages. Industry estimates are then pooled after re-centering their adjusted productivity v_i around the wage-kink threshold v^* . Outcome values are in logs relative to their average values in the perfectly competitive case with no minimum wage. The vertical line shows the wage-kink threshold and divides constrained firms to the left (21% of firms) and unconstrained firms to the right (79% of firms). See Section 2.3 for more details..

Finally, many variables omitted from the model are empirically correlated with productivity and affect the outcomes of concern. An important such example, to which we return later, is firms' skill composition: in principle the model can be applied to firm by occupation cells, but it still abstracts from within-cell wage and quality dispersion.⁵ For these reasons our tests rely on *differential* slopes around the kink-point, much like, for example, the literature on UI benefit floors using regression kink designs (Ganong and Jäger 2018).

⁵A correlation between worker skill and firm productivity will create positive correlations between firm productivity and wages which are separate from the monopsonistic and rent-sharing dynamics, our subject of interest here (Manning, 2003).

In particular, we test for the following patterns:

1. **Rent-sharing or wage-productivity curve:** The two key predictions regarding the co-movement of wages and productivity are (a) there exists a kink-point \hat{v}^* on the adjusted productivity axis at which the linear wage-productivity slope changes, corresponding to the model point v^* , and (b) the linear slope to the left of this point is smaller than the slope to the right, which is positive.
2. **Employment-productivity curve:** At the same kink-point \hat{v}^* , the linear slope to the left is lower than the slope to the right, which is positive.
3. **Profit share-productivity curve:** At the same kink-point \hat{v}^* , the linear slope to the left is positive and higher than the slope to the right.

The next sections take these testable hypotheses to the data, first in section 4 using a cross-sectional kink design and then in section 5 using within-firm shocks.⁶

3 Data and Institutional Context

3.1 Data

For our empirical analysis we use the administrative National Treasury-South African Revenue Services (NT-SARS) tax data held at the National Treasury Secure Data Facility (NT-SDF) in Pretoria. This is restricted-access data which can only be accessed in person at the NT-SDF for approved projects. The data we use consist of annual firm balance sheet information from Company Income Tax returns (“ITR14” forms) which include information such as sales, costs, profits and industry; matched worker-level annual payroll data (“IRP5” forms) which can be used to construct firm-level employment, (approximate) monthly wages, and firm geographic location; and linked firm-level monthly customs data, which provides transaction value and country of origin and destination for imports and exports respectively. The data constitute a panel covering the universe of formal-sector firms in South Africa, and while each dataset covers different periods they all reliably cover at least the period from the 2010 to 2019 tax years (approximately

⁶Rather than looking at variation in responses along the firm productivity distribution, as below, one may think that the minimum wage bite provides the appropriate variation; Appendix D.2 shows why this is not the case. Note that because we consider variation in firm-wide productivity, we do not predict labour-labour substitution as there may be in the case of minimum wage shocks.

the 2009-2018 calendar years).⁷ For the trade shock analysis discussed in Section 5.1.2 we use GDP data from the World Development Indicators.

3.2 Minimum wage institutions

Prior to January 2019, a multilayered wage legislation system operated in South Africa, where minimum wages were set by the government for selected broad industry-locations (“Sectoral Determinations”, SDs), or by publicly-recognized Bargaining Councils (BCs) consisting of employers and employees at the sub-industry-location level. Minimum wages can vary substantially by these sectors, and we therefore examine firms separately by their BC or SD.⁸

BCs cover industry-regions, and are constituted by trade union and employer representatives who negotiate industry-region minimum wages. Supplementary establishment-level wages can then also be negotiated above these minima, allowing for rent sharing. This is a setup common to a variety of European countries (Bhuller et al. 2022), but unlike in some of these countries BC agreements are routinely extended to include non-unionized workers. We identify BC firms in the SARS-NT data by matching firms according to their industry and location, using the Bassier (2022) dataset of BC agreements. There are 39 private sector BCs; after restricting for key missing variables we identify 30 in the data, which cover approximately 26% of the (formal private sector) workers in our sample. This dataset also provides a minimum wage associated with each BC for each year, but it is highly approximate: we must take the lowest BC-specified wage to be the minimum, even though BC agreements typically specify multiple occupation-specific wages, because occupations are not observed in the NT-SARS data.

SDs are government-set wage minima (and conditions of employment) for sectors not fully covered by BCs, often because they are understood as “hard to organise”. There are 11 SDs, 8 of which set minimum wages for formal sector workers. SDs are defined more expansively than BCs, and sometimes overlap with BCs; in these cases the BC minimum wages apply. While SDs, like BCs, may set occupation- and location-specific wages, there is usually less heterogeneity in minimum wages than in BCs. We identify SD firms in the NT-SARS data by matching their industry and location to a dataset we create from promulgated government regulations. We identify the 8 formal sector SDs, which exclusively cover about 32% of (formal

⁷See Appendix C for additional details concerning our use of this NT-SARS data.

⁸A national minimum wage (NMW) was introduced in January 2019 which supercedes a small number of these minimum wages. Only the last 2 months of our observed period overlap with the period of the national minimum wage, so we ignore this given that any dynamics in these months are likely to be irrelevant for our results in our empirical designs.

private sector) workers in our sample. These are predominantly workers at the lower-end of the wage distribution, unlike BCs which have coverage concentrated in the upper half of the wage distribution (Bassier 2022). We also include minimum wages from these regulations, but these are approximate for the same reason as the BC minima.

4 Evidence from a cross-sectional kink design

Implicit in the model above is a measure of firm productivity. We begin by estimating this productivity, and then test for the patterns in Section 2.3 by identifying and using kinks in the cross-section along the firm productivity axis.

4.1 Production function estimation

In estimating firm-specific productivity we draw from a substantial Industrial Organization literature concerned with production function estimation (Olley and Pakes 1996; Levinsohn and Petrin 2003; Akerberg et al. 2015; De Loecker and Syverson 2021). Recognising issues with OLS estimation of productivity such as simultaneity/transmission bias and selection/survival bias, we estimate productivity using the proxy variable/control function method of Akerberg et al. (2015) (ACF) with materials as the proxy variable, probably the leading approach in the literature (De Loecker and Syverson 2021; Yeh et al. 2022). Cognisant of the Gandhi et al. (2020) critique of attempts to estimate gross output production functions using proxy variable methods, we specify a value-added production function with a flexible translog form:

$$y_{it} = \beta_l l_{it} + \beta_{ll} l_{it}^2 + \beta_k k_{it} + \beta_{kk} k_{it}^2 + \beta_{lk} l_{it} k_{it} + \omega_{it} + \varepsilon_{it} \quad (6)$$

where y_{it} is value-added for firm i in period t , l_{it} is firm employment and k_{it} is firm capital stock, all in logs, while $\omega_{it} + \varepsilon_{it}$ is the productivity residual made up of productivity shocks which are observed or predictable for the firm at time t (ω_{it}) and those which are not (ε_{it}).

We show in Appendix Figures A6 and A7 that our main results are robust to a variety of alternative methods of estimating production functions: ACF with a Cobb-Douglas functional form, the Olley and Pakes (1996) method, the Levinsohn and Petrin (2003) method, the ACF correction applied after Olley and Pakes (1996) rather than Levinsohn and Petrin (2003) estimation, and the ACF method

estimated separately by various industry categories. Our results are unaffected by controlling for each firm’s share of the labour market wagebill during production function estimation.⁹ They are also robust to excluding observations around the threshold estimated below, addressing measurement error in the “running variable” following [Dong and Kolesár \(2023\)](#).

4.2 Empirical strategy: cross-sectional kink design

Overview. Our first step is to identify a kink (“knot”) (i.e. a threshold of a piecewise function) in the wage-productivity curve. The pattern we look for is as follows: Wages (measured as firm medians) are a piecewise linear continuous function of estimated productivity, defined over two intervals, and containing a discontinuity in its derivative (the knot) at the boundary between the intervals. We identify such a knot in the observed distribution by running two OLS regressions, one to the left and another to the right, for each point in the productivity distribution, and selecting the point or threshold which maximizes the R-squared.¹⁰ This is a candidate for the productivity threshold v^* in the model where firms move from being supply-constrained to unconstrained.

Given such a wage knot denoted \hat{v}^* , we test the predictions from [Section 2.3](#). That is, at the same productivity threshold, there are significantly different slopes on either side (this is what produces the kinks) in the cross-sectional firm wage (by construction), markdown and employment curves along the productivity axis, and these differential slopes are of the predicted sign. We separately regress each of these variables on productivity to the left and right of \hat{v}^* , which allows us to examine whether there is indeed a change in the slope as we expect.

Estimation details. As noted in [Section 2.3](#), due to the different minimum wages which operate in each BC/SD, the above knot-finding exercise is implemented separately for firms in each BC and SD. We estimate productivity for each firm in “pre-period” windows, and then the knot-finding exercise above is implemented only for the years after this pre-period.¹¹

⁹This is an input-market analogy of the [De Loecker et al. \(2016\)](#) approach to adjusting ACF for output market price-setting power, in case unobserved period t productivity shocks also affect labour demand via period t wage changes. This is likely a second-order margin in any case, or can also be resolved with a structural assumption that firms cannot instantaneously adjust wages (this is consistent with our dynamic results in [Section 5](#)).

¹⁰This procedure is analogous to that used by [Card et al. \(2016\)](#) to identify a similar kink in the distribution of AKM firm wage premia against firm log value added.

¹¹See [Appendix C](#) for a description of how we construct these pre-periods. We define firm productivity as the average of firm-year-specific productivity in the pre-period, winsorized at the 1% tails. [De Loecker and Syverson \(2021\)](#) note, in a related context, that this averaging may

With the BC- and SD-specific kink-points and cross-sectional patterns having been separately estimated, we then pool the results across all the different BCs and SDs. As discussed above, in order to account for the different minimum wages and other market-level characteristics of each BC and SD, which will necessarily lead to different v^* wage-kink productivity thresholds, we re-center productivity in each BC and SD around the estimated wage-kink productivity threshold \hat{v}^* in that BC/SD, so that re-centered productivity above 0 indicates an unconstrained firm and below 0 indicates a constrained firm. We then pool all BCs and SDs for testing differential slopes and for plotting.¹²

Though we do not focus on the demand-constrained region, we do try to isolate it from the supply-constrained region by identifying a kink on the employment-productivity curve to the left of the wage threshold (i.e. \hat{v}_1^*) for each BC/SD and in the pooled aggregate. In practice, this point is not well-identified, for the reasons discussed in section 2.3.

Identification and interpretation. Our tests rely on statistically significant breaks in slopes around the wage-kink threshold, in the spirit of the regression kink design (Card et al. 2015; Ganong and Jäger 2018). Porter and Yu (2015) show that when the discontinuity point is unobserved and has to be estimated, treatment effects and standard errors are asymptotically valid. A prominent application is Card et al. (2008), who estimate tipping point thresholds under a Schelling model of racial segregation separately by city, and similarly recenter and pool observations. Our tests are more demanding since, in addition, we test whether employment and profit share break at the same threshold as the wage.

A threat to the validity of these tests would be if the breaks in the wage-, employment- and profit share-productivity relationships were caused by discontinuous *breaks* in unobserved confounding variables exactly where minimum wages bind. This seems unlikely, given that the productivity threshold is estimated separately for each BC/SD and then recentered, meaning that such a confounding variable would also need to systematically break at each of these different BC/SD-specific productivity thresholds. To check this, we show that the observed model-exogenous characteristics of constrained and unconstrained firms are very similar around the wage-kink threshold (Appendix Table A2). We also show the kinks occur across various percentiles of wages, employment and profit share (Appendix

reduce misspecification error.

¹²In a few cases our knot-finding algorithm does not identify a plausible interior wage-kink \hat{v}^* and instead identifies a kink at extreme wage values; in order to exclude these cases we trim the estimated wage-kinks \hat{v}^* at the 1st and 99th percentile of the pooled firm distribution. However these cases are few and our results are robust to whether or not we trim here.

Figure A4).

Of course, it is not difficult to think of factors that generally co-vary with firm productivity. As discussed in Manning (2003) regarding the “employer size-wage effect”, plausible candidates are unobserved worker quality and amenities. Empirically, we find more highly skilled workers are indeed more likely to work in more productive, larger firms, creating a positive correlation between wages and employment which is not due to the shape of the labour supply curve.¹³ The assumption then of our cross-sectional kink design is that the confounding effect is somewhat constant across the distribution of firms, or varying but just not discontinuous at exactly the wage-kink threshold. This would mean our kink design provides a credible test for our main theoretical predictions, but the estimated slopes are not directly informative about the underlying model parameters.

4.3 Results from the cross-sectional kink design

Figure 4 shows the results from the pooled kink design, with no controls in panel (a) and with controls in panel (b).¹⁴ The vertical dotted line indicates the productivity threshold corresponding to the estimated wage kink productivity threshold \hat{v}^* .

The pooled aggregate case clearly exhibits the predicted features of the model. To start, the wage curve (blue, reproduced in Figure 1) has a relatively flat slope to the left and is more upwards-sloping to the right of the wage-kink at \hat{v}^* . Though there is a break in this relationship at the wage-kink \hat{v}^* by construction, the result that the signs and differential magnitudes of these slopes match the predictions is not by construction. Moreover, the flat portion corresponds very closely to the level of the average minimum wage (indicated by the horizontal dotted line in panel (a)), despite the minimum wage not being used in the estimation of the kink point. This strongly supports the estimated kink point \hat{v}^* as identifying the productivity threshold v^* below which firms are constrained by the minimum wage. Further supporting this, firms with greater minimum wage bite (as proxied by workers’

¹³We estimate implied cross-sectional elasticities with respect to value added in Appendix Table A1, instrumenting value-added with the recentered productivity value. As expected, the implied (and confounded) pass-through or rent-sharing elasticity of 0.438 (0.0061) is substantially higher than most identified estimates of this parameter in prior literature, while the implied firm-facing labour supply elasticity, 0.648 (0.022) is substantially lower (see Section 5 for discussion of these prior estimates).

¹⁴Specifically we control for detailed industry-by-location cell (2-digit industry by District Council), average AKM worker fixed effect at the firm (Abowd et al. 1999), and poaching ratio at the firm (E-E hires to all hires). Like with productivity, the average AKM worker effect at the firm and the poaching ratio are estimated in a “pre-period” which is not used when calculating the cross-sectional relationship.

wages or the Kaitz index) are correlated with being classified as constrained, as expected (see Appendix D).

Even more strikingly, Figure 4 depicts the predicted kink patterns for firm employment (green) and the profit share (red), despite these variables not being used for kink estimation and therefore not having any by-construction relationship to the kink-point. Specifically, for firm employment (green), the slope is near-flat to the left of the wage-kink threshold \hat{v}^* , and then sharply increases to the right. Just as strikingly, the slope of the profit share (red) is increasing to the left of the wage-kink threshold, and then changes to be near-flat to the right of the threshold. The plot is remarkably similar to the simulated model prediction plot in Figure 3, and even more so for Figure 4 panel (b) which includes controls.

Table 1 reports that the differences in the slopes around the wage-kink threshold in Figure 4 are statistically significant at the 1% level, for all three curves. We drop the 5% of firms above and below the estimated threshold to address possible measurement error in \hat{v}^* (a “donut” specification (Dong and Kolesár, 2023)), but this makes essentially no difference to the results. Because there is no by-construction relationship between firm size or the profit share and the wage-kink estimation routine, we regard these observed patterns as a compelling test and validation of the model predictions. 24% of firms in our sample are found in the constrained region, suggesting that the mechanism we discuss affects about 1 in 4 firms in South Africa.

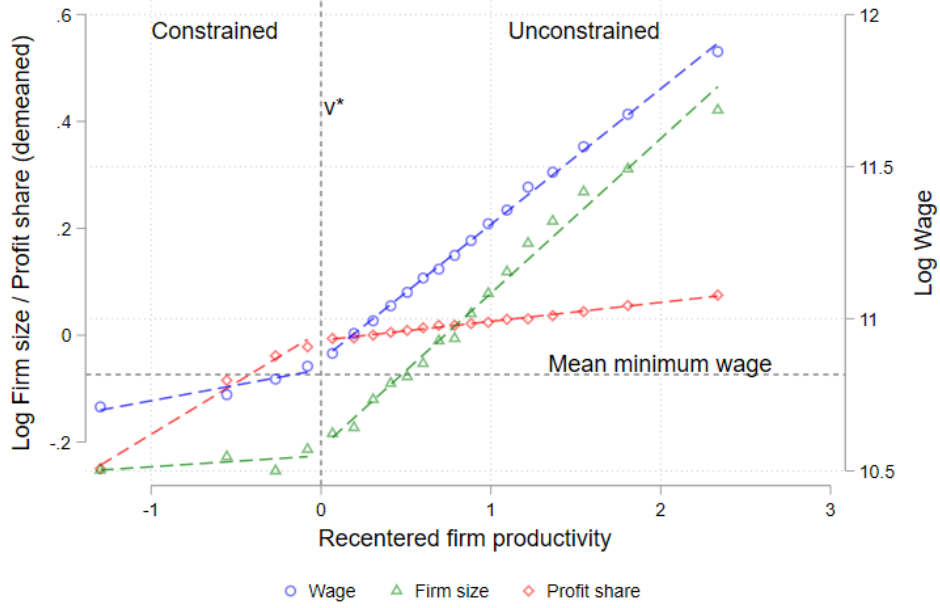
While the same patterns hold for most of the individual BCs and SDs (Appendix Figure A3), they are sometimes noisy or simply not evident for particular BCs or SDs. We view this as unsurprising, given that we are testing a strong prediction, which in any case may not be detectable given the unavoidably approximate nature of our productivity estimation routine and varying BC/SD sample size. The pooled figure above, however, is robust to several checks: the kink patterns occur for various percentiles of wages, employment and profits at \hat{v}^* , i.e. not just the medians or associated with a particular outcome level (Appendix Figure A4); results are similar when the sample is restricted to only one “event” per firm rather than the stacked event structure; and we also check robustness to the number of productivity bins, trimming, and re-estimating \hat{v}^* based on the pooled sample.

5 Responses to within-firm shocks

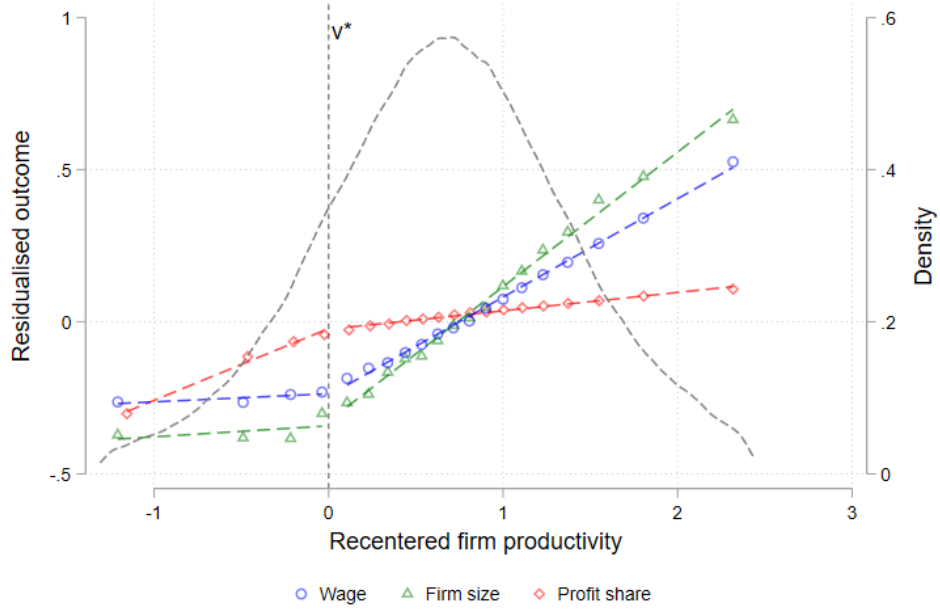
We complement the evidence from the cross-sectional kink design by looking at heterogeneity in within-firm dynamic *responses* to revenue-productivity increases,

Figure 4: Cross-sectional kink design: Empirical results

(a) Raw outcomes



(b) Residualized outcomes



Notes: The figures show firm median wage (blue), employment (green), and profit share (red) (all in logs) by 20 recentered firm productivity bins. Firm productivity is estimated using the ACF method, and recentered around the wage-kink \hat{v}^* (vertical dashed line), which is estimated separately for each minimum wage regime. The wage-kink \hat{v}^* divides minimum wage-constrained (to the left) and unconstrained firms (to the right). Panel (a) shows the raw outcomes, where the horizontal line is the average minimum wage across firms. Panel (b) shows the outcomes residualized on the average AKM worker fixed at the firm, the firm poaching ratio, and industry (2-digit) by region (district council) labour market dummies. It also shows the firm distribution density with a black dashed line. See Table 1 for slope magnitudes and Section 4 for estimation details.

Table 1: Differential slopes from cross-sectional kink design, firm outcomes on re-centered productivity

	No controls			With controls		
	Constrained	Unconstrained	Difference	Constrained	Unconstrained	Difference
Wage	0.085*** (0.0067)	0.443*** (0.0118)	0.358*** (0.0136)	0.032*** (0.0067)	0.325*** (0.0069)	0.293*** (0.0096)
Employment	0.009 (0.0189)	0.290*** (0.0247)	0.281*** (0.0311)	0.068*** (0.0161)	0.466*** (0.0188)	0.398*** (0.0248)
Profit-share	0.207*** (0.0084)	0.036*** (0.0029)	-0.171*** (0.0089)	0.254*** (0.0096)	0.055*** (0.0020)	-0.199*** (0.0098)
N	237763	978665		179035	807201	

Notes: The table shows the slopes from the kink design in Figure 4, separately for constrained and unconstrained firms, as well as the difference in slopes between the constrained and unconstrained firms. These are the slopes of the outcome variable (wage, employment, or profit share, all in logs) on the recentered productivity term. Controls refer to the average AKM worker fixed at the firm, the firm poaching ratio, and industry (2-digit) by region (district council) labour market dummies. 5% of firms above and below the threshold \hat{v}^* are dropped to create a “donut” estimate. Standard errors are shown in parentheses and clustered at labour market by event; differences are calculated using the Delta method. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

for firms along the productivity distribution. This design resolves biases in the cross-sectional magnitudes discussed earlier (e.g. firm skill composition) and should produce identified estimates of the rent-sharing and firm-specific labour supply elasticities comparable to the rest of the literature, as well the employment- and profit share-revenue productivity elasticities. To this end, we follow a large number of papers which use shocks to firm value added to identify firm rent sharing and employment responses (e.g. [Amodio and De Roux 2022](#); [Kroft et al. 2025](#)). In particular, we follow the approach in [Lamadon et al. \(2022\)](#) (LMS) with both an “internal” and “external” instrument for value added. Also following LMS, our wage measures focus on workers who stay at the same firm over the full estimation period. This also helps address potentially confounding worker quality and other compositional issues such as tenure effects, which could otherwise bias our estimates.

5.1 Empirical strategy: within-firm shock design

Following the pooling procedure used in our kink design, we use the estimated re-centered productivity variable and associated threshold \hat{v}^* (see Section 4.2) to test for heterogeneous responses to the left (constrained firms) compared to the

right of the threshold (unconstrained firms). We implement the same robustness checks as in Section 4.2 to address possible measurement error in the estimation of firm productivity (the running variable) and the v^* threshold.

The relevant testable hypotheses stated at the end of Section 2 were differential slopes on either side of the kink point for each of the outcome variables with respect to productivity. In the context of firm responses to shocks to marginal revenue productivity (such as demand shocks), these hypotheses translate into differential marginal responses, i.e. higher responses to the right of \hat{v}^* for wages and employment, and lower responses to the right of \hat{v}^* for profit share. We test this directly by estimating the heterogeneous marginal responses to value added shocks along the firm productivity distribution, using two approaches described below.

Recall that the key test is *differential* responses around the threshold \hat{v}^* , rather than 0 wage or employment response in the constrained region, because practical considerations mean the theoretically-implied 0 response is unlikely to be identified in our empirical exercise. One reason is simply noise due to measurement and sampling error: our threshold between constrained and unconstrained firms is an estimated quantity rather than an exact delineation, based on estimated firm-specific productivity, so there will be some unconstrained firms in the constrained region (around the threshold), and vice versa. Another more systematic issue is that constrained firms receiving a value added shock may cross over into the unconstrained region, thus exhibiting (attenuated) behavior associated with the unconstrained region. Indeed, a back-of-envelope calculation suggests that approximately 20% of constrained firms move into the unconstrained region as a result of the internal shock.¹⁵ Thirdly, minimum-wage compliance in South Africa is far from perfect (Bhorat et al. 2012), which will be another reason that firms identified in the constrained region will actually be unconstrained and exhibit unconstrained responses. Lastly, the unconstrained region contains demand-constrained as well as supply-constrained firms, and demand-constrained firms will have employment responses to productivity shocks, and no profit share response, like firms in the unconstrained region. Fortunately, these issues would *attenuate* the differential responses between the constrained and unconstrained regions, making our predic-

¹⁵To calculate this we run a firm-level cross-sectional regression of the re-centered productivity measure on the log of value-added (with various controls), finding a coefficient of about 0.3. Multiplying this by the size of the value-added shock in the constrained region (0.34 log points), we find that the value-added shock increases constrained firms' revenue-productivity measure by about 0.10 log points. Empirically, 20% of constrained firms are within 0.10 re-centered productivity log points of the threshold, and so approximately 1 in 5 constrained firms transition into the unconstrained region as a result of the TFP shock.

tion of a differential response a *more* demanding test of the model predictions.

5.1.1 Internal instrument

The core of the LMS “internal” instrument method is a firm-level event study analysis where treatment is defined as an above-median increase in firm value-added between periods -1 and 0, with some additional specification and variable- and sample-definition restrictions. LMS focus on the effects of these binary value-added shocks on earnings. We extend this firstly by also examining effects on employment and the capital share, but most importantly by examining heterogeneous responses along the firm productivity distribution. A related exercise is undertaken by [de Frahan et al. \(2022\)](#), who also extend LMS to examine heterogeneous effects on employment as well as earnings along the firm size distribution.¹⁶

The sample and event definitions are important to the LMS reduced form analysis, and we closely follow them. While we outline these restrictions in detail in Appendix C.3, the core of the approach is to construct events which use a balanced panel of firms which have a minimum number of worker “stayers” who remain employed at the same firm for the period of the event. For our baseline specification we keep firms which have at least 2 stayers.¹⁷ We stack four 6-year events (2010-2017, 2011-2018, 2012-2019, and 2013-2020 tax years), where we treat the first three periods as the pre-period and the latter three as the post. Each event is a balanced panel.

For the internal LMS instrument approach, treatment is defined as an above-median increase in firm value-added between periods -1 and 0 for each event, where the median increase is weighted by firm size. As in LMS, period -2 is used as the omitted reference period to allow for some mean reversion dynamics in period -1, while period -3 is used to assess pre-period parallel trend violations. For the same reason, periods 1 and 2 are considered the post periods of interest, rather than period 0 (results are essentially unchanged if we use only period 2 instead).

¹⁶While not a focus of this paper, we note that we have been able to replicate the [de Frahan et al. \(2022\)](#) findings in our data, finding qualitatively similar results.

¹⁷This is our only notable divergence from LMS, who use a 10-stayers minimum as their baseline. A 10-stayers minimum is overly restrictive when it comes to the South African firm-size distribution and a labour market context defined by high churn. We use a 2-stayer minimum to be as nonrestrictive as possible but to mitigate measurement error in one-stayer firms where the one stayer may be an owner or otherwise unrepresentative of general employer/employee dynamics. In Appendix Figure A8 we show that our main results are not sensitive to the number of stayers.

The aggregate reduced-form event study regression is:

$$y_{i,t,e} = \lambda_{i,e} + \gamma_{t,m(i),e} + \sum_{s=-3, s \neq -2}^2 \beta_s \times \mathbb{1}[t = s] \times D_{i,e} + \varepsilon_{i,t,e} \quad (7)$$

where $y_{i,t,e}$ is the log of the outcome for firm i at time t in labor market m for event e , $\lambda_{i,e}$ is the firm-event fixed effect, $\gamma_{t,m(i),e}$ is a time-varying market-event fixed effect, $D_{i,e}$ is the treatment variable and the β_s are the coefficients of interest, relative to period -2. Standard errors are clustered at the market-event level. The market controls $\gamma_{t,m(i),e}$ constitute 81 labor markets (interacted with event and year), made up of province and 1-digit industry.¹⁸

Since the binary shock may lead to differently-sized shocks to value-added along the distribution, our main specification normalizes the size of the shock by using the binary treatment as an instrument for changes in firm value-added.¹⁹ This also allows us to estimate full elasticities, comparable to other estimates in the literature, rather than reduced-form semi-elasticities. Specifically, we run long-differences regressions of the log change in the outcome between the post-period (periods 1 and 2) and period -2 on the equivalent change in value-added, with the change in value-added instrumented by the binary treatment. We include fixed effects analogous to those in equation 7 and again cluster at the market-event level. The reduced-form results are relevant mainly insofar as they allow for examining the shape of the event study response, to assess issues such as parallel pre-trends, but the reduced-form results are in any case similar to our main results, statistically significant (levels and differences between constrained versus unconstrained), and consistent with the model predictions.

5.1.2 External instrument

While the “internal” approach is a greater focus of their paper, LMS also implement a supplementary “external” instrument specification, where they use firms receiving a procurement contract to define treatment. We use a different external instrument, relying on shift-share trade shocks similar to [Garin and Silvério \(2023\)](#) and [Bassier and Manning \(2025\)](#).

¹⁸Appendix Figure A9 shows results when different industry and geography variables are used, including using the 2-digit industry and the “district” geography variable, which creates 2600 labor market interactions. Our main results are essentially unchanged.

¹⁹This accounts for systematic differences in the binary treatment in each of the heterogeneity regions. In our baseline specification the average value-added shock is 0.34 log points in the constrained region versus 0.29 log points in the unconstrained region.

We use the South African administrative customs dataset to define *firm-specific* shift-share trade instruments based on the country composition of firm exports and imports, combined with movements in destination- and origin-country GDP, respectively. Specifically, for the export instrument, we define firm i 's shock in period t as,

$$D_{i,t}^{exports-IV} = \sum_d \alpha_{i,d}^{export-share} \frac{\psi_{d,t}^{shift}}{\bar{\psi}_d} \quad (8)$$

where their export exposure “share” $\alpha_{i,d}^{export-share}$ is the proportion of their exports made up of exports to country d across the sample period (i.e. a firm-specific constant), and the “shock” $\psi_{d,t}^{shift}$ is the GDP of country d in that year t normalized by the mid-period GDP of that country $\bar{\psi}_d$. Such foreign GDP movements constitute positive revenue-productivity shocks for exporters, since those countries demand more goods. We define an analogous import shock $D_{i,t}^{imports-IV}$, which is defined the same as in Equation 8 except that the term $\alpha_{i,d}^{export-share}$ is replaced with $\alpha_{i,d}^{import-share}$, which is the share of firm i 's imports from country d .

A notable difference with the internal instrument is that the shift-share treatment can only be defined for trading firms, and this decreases the sample size considerably, especially for low-productivity constrained firms. This is unsurprising; it is well-known that exporters are generally higher-productivity firms (Verhoogen 2023). The scale of the reduction is however quite dramatic: while the estimation sample for the main internal instrument specification is made up of 6,167 distinct constrained firms and 34,120 distinct unconstrained firms, when this is restricted to trading firms for the external instrument approach this reduces the equivalent specification's sample to 956 distinct constrained firms and 12,017 distinct unconstrained firms. While the overall number of firms in the sample reduces by just over one third, this is particularly concentrated among the constrained firms, which drop from being about 18% of the sample in the internal instrument case to being about 8% of the sample in the external instrument setting. This dramatically reduces statistical power, rendering us unable to present an event-study or estimate results by fine-grained productivity bins.²⁰ Instead, we define only two bins (constrained firms versus unconstrained firms) and estimate a fixed effects regression using the full sample period (no split into pre- and post-periods).²¹

²⁰Indeed, if we implement the LMS *internal* instrument approach on the trade shock sample, we find that while magnitudes are qualitatively in line with our model predictions and similar to our main results in Table 2, the estimates are much more imprecise.

²¹In Appendix Table A5 we report results for a specification where the export and import shares are estimated only over the pre-periods used in the internal instrument approach, and GDP shocks are only considered in the equivalent post-periods. While standard errors are

The reduced-form specification is very similar to 7 above, but without the event study terms:

$$y_{i,t,e} = \lambda_{i,e} + \gamma_{t,m(i),e} + \beta_{export} \times D_{i,t}^{exports-IV} + \beta_{import} \times D_{i,t}^{imports-IV} + \varepsilon_{i,t,e} \quad (9)$$

As for the internal instrument approach, for our main results we use these shocks (in this case $D_{i,t}^{exports-IV}$ and $D_{i,t}^{imports-IV}$) to instrument for firm value added.

5.2 Results from within-firm shocks

5.2.1 Internal IV results

Figure 5 and Table 2 show our main results from the internal instrument specification, with firm median wage (blue), firm profit share (red), firm full-time equivalent employment (green) and firm hires (purple) responses as elasticities with respect to the induced changes in value-added. To estimate responses along the firm productivity distribution, we essentially estimate the instrumented version of equation 7 separately by productivity bin. The key differentiation is binary: constrained versus unconstrained, which means separate regressions for firms below versus above the re-centered productivity threshold (\hat{v}^* , shown with a vertical dashed line in Figure 5). These estimates are shown as horizontal dashed lines with their associated 95% confidence intervals in Figure 5, and in separate columns in Table 2. We also present estimates by 10 approximately equally-sized productivity bins to get a sense of the shape of response along the distribution. The construction of these bins is discussed in Appendix C, as well as other minor sample and estimation decisions (we show that these decisions are not consequential; our results are robust to alternative specifications).

The results show very clearly the pattern predicted by the theory, with highly statistically significant differences between constrained and unconstrained firms in their wage, employment and profit-share responses, which observably break around the estimated wage-kink \hat{v}^* . Wages and employment responses are lower for constrained firms and higher for unconstrained firms, while the converse is true for the profit share.

While not one of our core theoretical predictions, in panel (b) of Figure 5 and Table 2 we also show effects on hires using the same method. There is an even starker contrast between the constrained and unconstrained firms. It is arguably

much wider and magnitudes vary, the qualitative conclusions are the same as in the baseline specification results we discuss in the following sections.

at exactly the hires margin of employment that the core model mechanism of section 2 operates, as firms tend to dynamically adjust on the hires margin while the other margin, separations, is less in their control.

The magnitudes reported in Table 2 panel (a) imply that in the constrained region, relative to the unconstrained region, the wage elasticity (i.e. rent sharing) is 29% lower, the employment elasticity is 25% lower, the hires elasticity is 41% lower, and the profit share elasticity is 175% higher. All these differences are highly statistically significant and consistent with the model predictions.

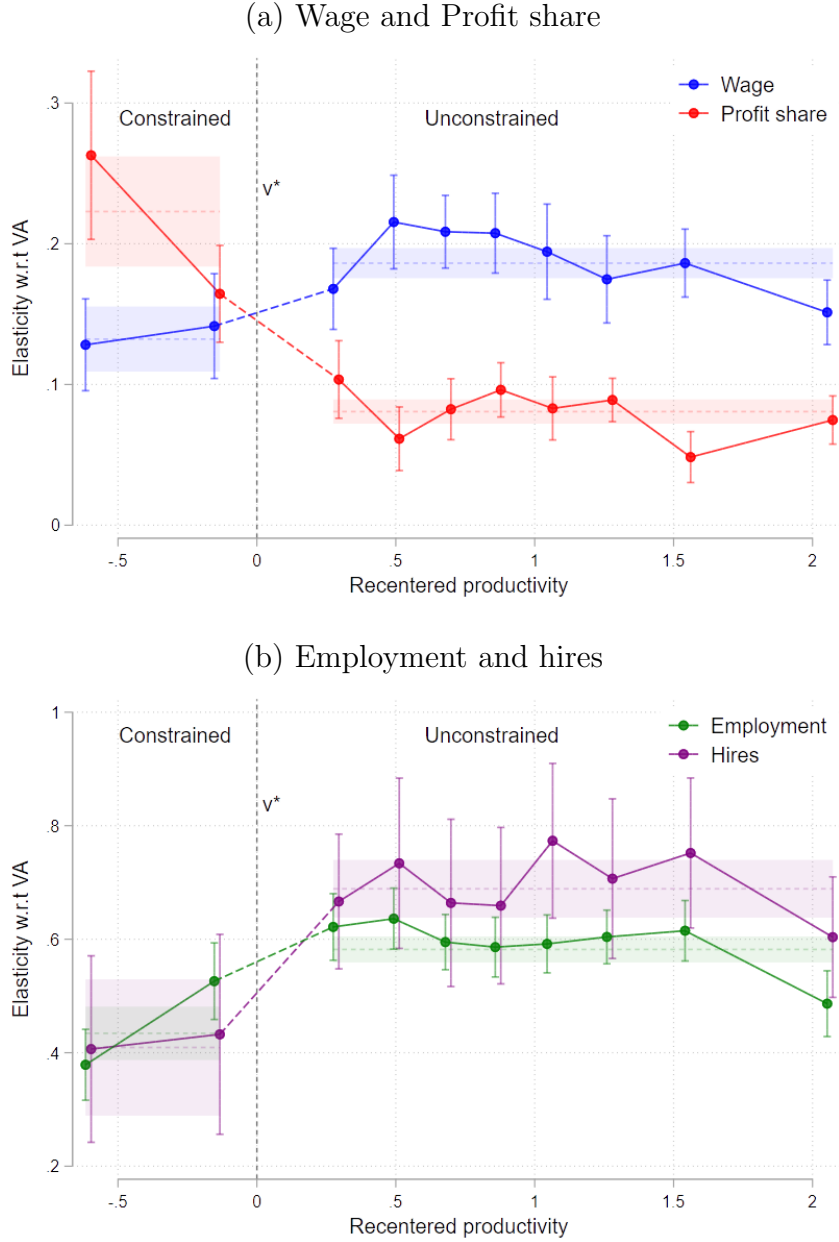
In order to check pre-trends and the robustness of the main results in Figure 5 and Table 2, the reduced-form event-study results for the internal IV are shown in Appendix Figure A5 and Appendix Table A3. Panel (a) of Figure A5 shows the aggregate response across the sample and panel (b) the responses along the recentered productivity distribution. The aggregate reduced form response in Figure A5 panel (a) is very similar for the comparable variables in LMS and de Frahan et al. (2022), which also suggest some mean reversion in VA and other dynamics in period -1, but which are small relative to the size of the post-period effects (and seems to dissipate between period 0 and period 1 in any case). The heterogeneous reduced form result in panel (b) is notable for two main reasons: 1) the reduced form response patterns (semi-elasticities) match the pattern in Figure A5 (and Appendix Table A3 shows that differences between constrained versus unconstrained semi-elasticities are also statistically significant); and 2) Pre-trend coefficients for period -3 (lighter-shade estimates with dashed lines) are generally not statistically significantly different from zero, and always economically much smaller than the post-period effects.

Our results are robust to a number of potential concerns. Appendix Figures A6 and A7 show robustness to alternative methods of estimating the underlying production functions, Appendix Figure A8 shows robustness to the number of stayers, and Appendix Figure A9 shows robustness to the choice of time-varying labour market fixed effects. Appendix Figure A10 Panels (a)-(d) show that the results are robust to allowing for misclassification around the estimated \hat{v}^* threshold, while Panels (e) and (f) show that the patterns are not driven by firms in the tails of the recentered productivity distribution.

5.2.2 External IV results

The second super-column of Table 2 presents the results from the external IV, the shift-share trade shocks. Recall that the source of variation is completely different

Figure 5: Internal IV main results (full elasticities), by recentered productivity bin



Notes: Figure shows main IV results for the internal instrument specification, estimated by productivity bin. Panel (a) shows the wage and profit share responses, while Panel (b) shows the employment and hires responses. Blue is the median wage of firm stayers (incumbents), red is for firm profit share, green is firm employment, and purple is firm hires. Full elasticities are shown, estimated by regressing the pre-post change in the outcome (in logs) on the pre-post change in log value-added, with the change in value-added instrumented by the binary treatment variable, as discussed in Section 5. The horizontal axis is firm productivity (estimated using the ACF method) recentered around the estimated productivity wage-kink \hat{v}^* (see Section 4). Recentered productivity equals zero at the wage-kink, shown with a vertical dashed line, and this line divides minimum wage-constrained (to the left) and -unconstrained (to the right) firms. Ten approximately equally-sized productivity bins (deciles) are created. The solid lines and points show the average treatment effect across post-periods 1 and 2 by bin. 95% confidence intervals are shown with vertical bars. The horizontal dashed lines with attendant shaded regions (95% confidence intervals) show applicable post-period treatment effects estimated across the productivity bins, pooled separately below and above the wage-kink value \hat{v}^* .

Table 2: Firm responses to productivity shocks, constrained versus unconstrained firms

	Internal IV (large VA change)			External IV (trade shock)		
	Constrained (1)	Unconstrained (2)	Difference	Constrained (3)	Unconstrained (4)	Difference
Panel (a)						
Rent sharing	0.132*** (0.0117)	0.186*** (0.0054)	0.054*** (0.0129)	0.129*** (0.0284)	0.188*** (0.0104)	0.059* (0.0302)
Employment	0.434*** (0.0238)	0.582*** (0.0116)	0.148*** (0.0265)	0.470*** (0.0887)	0.635*** (0.0247)	0.164* (0.0921)
Profit-share	0.223*** (0.0198)	0.081*** (0.0044)	-0.142*** (0.0203)	0.228*** (0.0502)	0.080*** (0.0092)	-0.148*** (0.0511)
F-stat	449.9	5956.7		40.6	716.5	
N firms	6167	34120		956	12017	
Obs	13596	98234		13242	224244	
Panel (b)						
Hires	0.409*** (0.0609)	0.689*** (0.0258)	0.280*** (0.0662)	0.271* (0.1549)	0.685*** (0.0594)	0.414** (0.1659)
F-stat	299.4	5013.0		32.6	583.6	
N firms	4884	29587		943	11893	
Obs	9929	77244		11367	191588	

Notes: The table shows main results for the responses of wages (rent sharing), firm employment, profit-share and firm hires (all in logs) to firm value added (VA) shocks, by constrained versus unconstrained status, as well as the difference in responses between constrained and unconstrained. Firm VA (in logs) is instrumented using two approaches: the left super-column shows estimates from the internal IV, i.e. above-median increases in firm VA between event periods -1 and 0. Only the post-period effects are reported. The right super-column shows estimates from the external IV, i.e. the shift-share trade shocks. See section 5 for sample and specification details. Note that the external shock uses a larger pooled sample to account for sample size limitations (as evidenced by the count of individual firms). Standard errors are shown in parentheses and clustered at labor market by event; for differences these are calculated using the Delta method. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

to the internal IV, based on GDP changes in the export destinations and import origin countries, of only the sample of trading firms. Yet the resulting elasticities are remarkably similar, providing reassurance that since two very different instruments give very similar estimates, then there is likely to be little bias. The column showing the differences in the two slopes shows that in the constrained region the elasticity for rent sharing is 0.059 log points lower (0.054 for the internal IV), for employment is 0.164 log points lower (versus 0.148), for hires is 0.414 (versus 0.280), and for profit share is 0.148 higher (versus 0.142). As noted above, the shift-share strategy has less power than the internal IV, with the demanding test for elasticity differences across regions not quite significant at the 5% level for wages ($p = 0.051$) and employment ($p = 0.074$), but for hires and profit share still statistically significant at the 5% level.

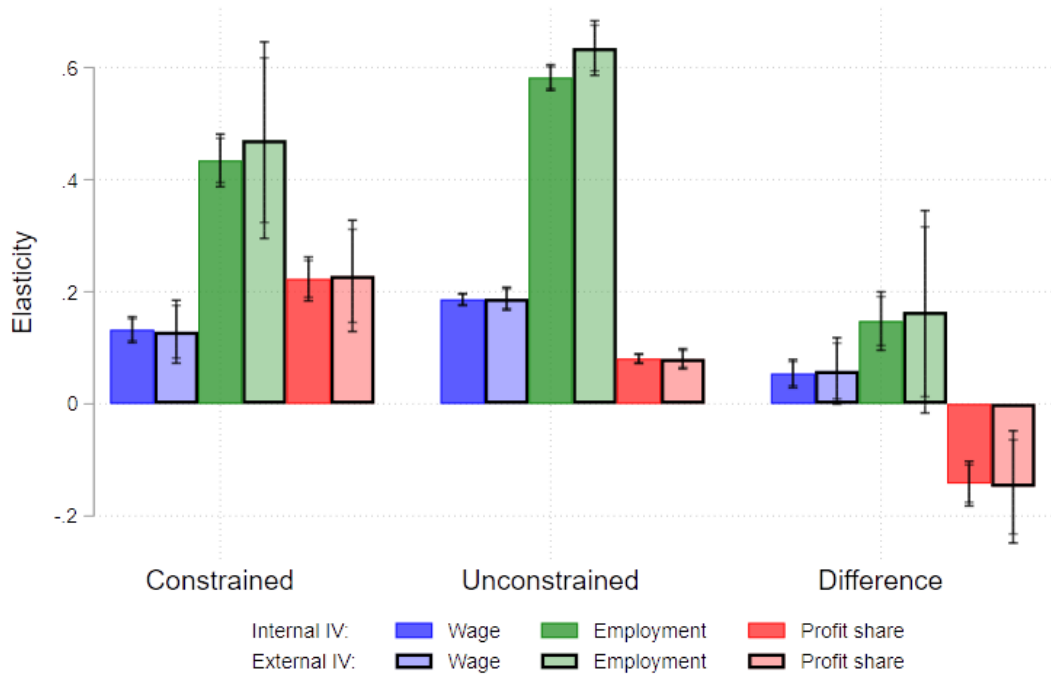
Figure 6 illustrates just how close the estimates are between the internal and external IVs for the main outcomes (Appendix Figure A11 shows the same graph for hires). It is worth highlighting that not only are the differences close, but also the separate estimates for the elasticities in the constrained and unconstrained regions. These region-specific elasticities are all highly significant.

5.2.3 Discussion

The estimates in the unconstrained region are of independent interest since, as discussed earlier, these are interpretable as identified estimates of the corresponding elasticities. Several papers use similar methods in estimating these elasticities, but we highlight here that we are additionally able to isolate the unconstrained region (noting that such elasticities are unidentified in the constrained region). Focusing on the external IV shift-share estimates, the elasticity of the wage with respect to value-added (i.e. the rent-sharing elasticity) is 0.19 (standard error 0.01), and for employment is 0.64 (0.02). The rent-sharing elasticities are very similar to what have been found in the existing literature, for example between 0.14 and 0.16 in Card et al. (2016), and between 0.13 and 0.19 in Lamadon et al. (2022). These are much smaller than the likely biased elasticities implied by naive interpretation of the magnitudes in the cross-sectional kink design, as expected.

The implied firm-facing labour supply elasticity is just the employment elasticity divided by the wage elasticity. Focusing again on the external IV unconstrained sample, our estimate of 3.4 (standard error 0.23) is well within the range of estimates found in the existing literature, for example 2.5 in Amodio and De Roux (2022), 3.8 in Kline et al. (2019), 4.1 in Kroft et al. (2025), and 2.15 in Dal Bó

Figure 6: Comparison of estimates from internal and external IV shocks



Notes: The figure plots the estimates from Table 2, where the internal IV refers to above-median increases in firm VA between event periods -1 and 0 while the external IV refers to the shift-share trade shocks. Wage responses are in blue, employment in green, and profit share in red (all are elasticities). See section 5 for sample and specification details. Vertical bars represent 90% and 95% confidence intervals, where standard errors are clustered at labor market by event and differences these are calculated using the Delta method.

et al. (2013).²² Interestingly, the ratio of differences (rightmost column) between the employment and wage elasticities also implies a fairly similar elasticity of 2.8. If we take the approach as in the regression kink design literature, and think of crossing the productivity threshold v^* as the “treatment”, then these differences in elasticities across the regions can be thought of as a complementary identification strategy for the firm labour supply elasticity focusing on the wage-setting mechanism (activated at the kink point v^*) separately from any other potential mechanisms driving rent sharing and employment responses. Indeed, if labour supply elasticities are heterogeneous, then typical estimates get a weighted average at best while a kink design such as ours pinpoints a LATE at our threshold point v^* (Kline, 2025).

Finally, the profit share elasticity in Table 2 is also of interest especially to the monopsony literature, where there is little evidence on how much of their monopsony power and implied *potential* markdown firms actually exploit (Manning 2021). Using the approach above of using the differenced elasticity as an indication of the wage-setting or monopsony mechanism kicking in between a firm’s constrained and unconstrained state, our estimates imply a markdown due to becoming unconstrained of 0.148 log points or 16%. This is not far from the monopsony markdown on marginal product implied by the firm labour supply elasticity estimated above, equal to 23%.²³ While the correspondence between profit share and marginal markdown is not clear, which means one should be careful about interpreting these magnitudes too closely, we do interpret this as evidence linking monopsony power and profits.

6 Structural model estimation and policy simulations

We now use the results from the previous sections to estimate the model parameters using a structural model. This allows us to draw a tighter link between the theory and empirical results, which in turn clarifies the underlying assumptions and enables us to simulate some striking implications of the model for various policies.

²²It is notable that the implied labour supply elasticity in the constrained region is very similar to the unconstrained elasticity, at 3.6 (1.06). As explained in Section 2, our model implies that the labor supply elasticity is in fact not identified in the constrained region. The similarity may be due to chance, but the striking similarity may be because some unconstrained firms are mis-classified as constrained – see Section 6.2.

²³The monopsony markdown is $1 - \varepsilon/(1 + \varepsilon)$, which using $\varepsilon = 3.4$ yields 23%.

6.1 Estimation of structural model

The simple theoretical model of Section 2 needs to be modified in two ways before we structurally estimate it. Firstly, we noted earlier that the precise magnitudes from the cross-sectional patterns were likely confounded by omitted worker quality. We account for this in our structural model simply as a covariance between firm and worker productivity in the wage-setting equation, as found empirically in several studies of firm and worker effects including in this dataset (e.g., [Bassier, 2023](#); [Card et al., 2016](#); [Engbom and Moser, 2022](#)).

Secondly, while Section 2 considers predictions for the wage markdown, our empirical tests use the gross profit share. As we show below, this is of limited consequence, since the qualitative model predictions are the same in both cases. In the structural model we additionally allow for misspecification in the production function and profit definition: for example, because gross profits may not correctly incorporate the rental cost of capital. Essentially we allow noise in the relationship between observed gross profits and estimated firm productivity, which we capture in the model by adding an unconstrained profits slope parameter.

Appendix E provides details of the estimation, including the model equations and fit. Our procedure begins with the estimate of 3.4 for the firm labour supply elasticity from the within-firm shock design (see Section 5). We then use the cross-sectional data points from Figure 4 Panel (b) along with the modified model to optimize the structural parameters by maximum likelihood estimation (MLE). Table E2 provides the source of identification and values of the 5 estimated structural parameters: the inverse elasticity of labour demand and worker-firm productivity covariance are both pinned down by the slopes from the unconstrained wage and employment curves, and then there are three nuisance parameters.

Appendix figure E1 shows the fit of the structural model against the observed data. The wage and employment curves fit extremely well; indeed, this simple model of a kink, with a linear slope to the right and a constant to the left fits the entire cross section of these outcomes remarkably well. The fit for the gross profit share is worse: the model predicts a constant profit share after the kink, but the data shows it is slightly upwards sloping, and similarly the fit is worse to the left of the kink. We view this as unsurprising, given that profits may be especially prone to misspecification as discussed above, and in fact this shows that there is no by-construction guarantee of the excellent fit we find for the wage and employment curves. A monopsony-based rationale would be that the profit share increases with employment share (and hence productivity) in oligopsonistic

markets (Berger et al., 2022).

6.2 Comparing to results from the within-firm shocks

How does this stack up with the results from the within-firm shock design? Recall that the firm labour supply elasticity from this design is used as an input into the structural model estimates above. The within-firm design also yielded an untargeted reduced-form estimate of the rent-sharing elasticity of 0.19 (see Table 2); the structural model implies a rent-sharing elasticity of 0.14, which is qualitatively similar.

One aspect of the within-firm results that does not match the structural model is the *level* of constrained wage and employment estimates (see Table 2); though significantly lower than the unconstrained estimates, the constrained estimate levels are still positive while the model implies they should be zero.

Several intuitive explanations could account for this. Firstly, in this model we assume minimum wage compliance. Non-compliance, which is widespread in South Africa (Bhorat et al., 2012), would mean many firms to the left of the kink are not in reality constrained and so will increase wages and employment as if they are unconstrained, thus explaining a positive constrained region estimate. Secondly, noting much of our firm sample is part of collective bargaining councils, it may well be that many firms operate under a union-negotiated model; we explore this more later. Thirdly, and more innocuously, if some unconstrained firms are misclassified as constrained (due to specification error in firm productivity estimation or kink-point estimation), this would create the same effect. All three explanations can be modelled in a similar way, and we estimate their potential salience in Appendix E.2. We find that these explanations can rationalize the magnitudes of the constrained wage and employment effects from our within-firm shock results if about 70% of firms are affected by some combination of noncompliance, union bargaining, or measurement error. Such intuitive departures from the simple baseline model reconcile the empirical results and yield a similar estimate of the firm labour supply elasticity of around 3.

6.3 Simulations of rent-sharing and policy implications

Equipped with the parameters from the structural model, we can simulate the effects of policies on our modeled outcomes. Our first simulation illustrates an interesting implication in Figure 7 panel (a), which shows the rent-sharing level (wage over marginal revenue product of labour), elasticity (percentage changes

in wage over MRPL) and pass-through (absolute changes in wage over MRPL) by firm productivity. For lower-productivity firms constrained by the minimum wage, the rent-sharing level is high, since the minimum wage forces firms to pay higher wages than otherwise. However, the elasticity and pass-through are zero: firms absorb productivity increases purely into markdowns and pass nothing on to wages. This matches what we observe empirically (see Appendix Table A4): the rent-sharing *levels* are much higher for constrained firms than unconstrained firms (approximately 0.9 versus 0.3), even as the rent-sharing *elasticities* are higher for unconstrained firms.²⁴

The differential rent-sharing effects mean that productivity increases will have very different effects across the productivity distribution. Policymakers, especially in developing countries, are interested in targeted industrial policy that boosts firm productivity, for example through electrification (Rud, 2012). We illustrate this by simulating the partial equilibrium effects of a 10% rise in firm productivity (see Appendix figure E2). We show that when a greater proportion of workers are in constrained firms, such productivity increases benefit workers much less in terms of wages or employment, and instead boost the profit share. Indeed this may bite precisely where policymakers often target their efforts, i.e. Small, Medium and Micro Enterprises or SMMEs (Grimm and Paffhausen, 2015). Smaller firms tend to have lower productivity and therefore are more constrained, meaning the returns to the policy in wages and employment are low.

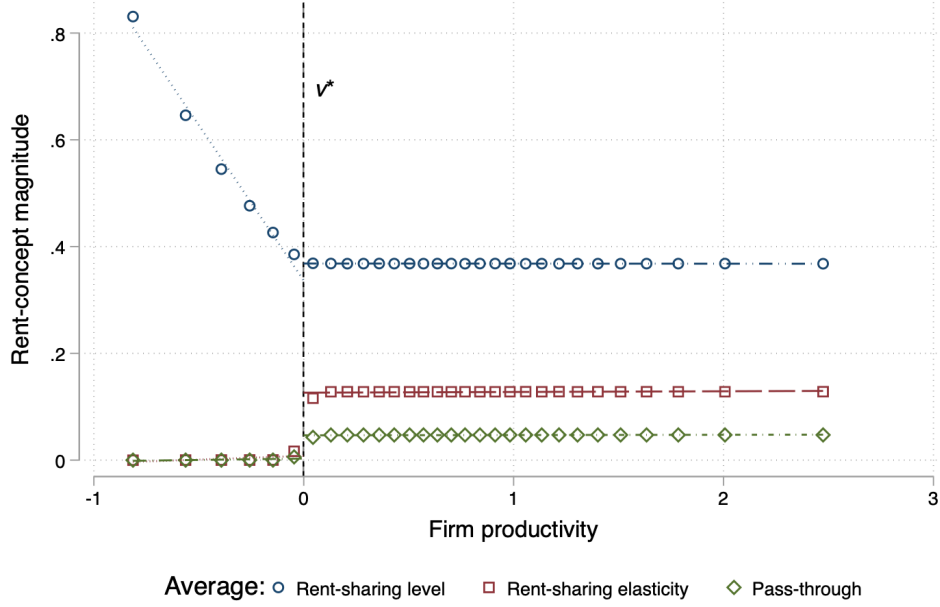
Our second policy simulation is of an employment tax incentive (ETI), the most prominent active labour market policy in South Africa and one that has been adopted in several other countries.²⁵ We simulate a simplified version of this policy as a 25% wage subsidy for each worker. This lowers the cost of labour to the firm, shifting the labour supply curve down, and in our model the constrained productivity thresholds also decrease (shift left). Figure 7 panel (b) shows that for unconstrained firms receiving the subsidy (43% of all firms), the policy works as intended by increasing employment and wages (profits too, but the profit share is constant). However, for a large proportion of firms that are supply-constrained (26%), the subsidy is absorbed purely into profits with no wage or employment increases. The rest of the firms (32%) have some intermediate response. Altogether, only half of the fiscal cost goes towards the wagebill, with the rest going to firm

²⁴These estimates are approximate and do not exactly recover the exact model-consistent measures due to data constraints. For example, the level is calculated as firm wage bill over firm value added, and the rent sharing elasticities are for stayers only.

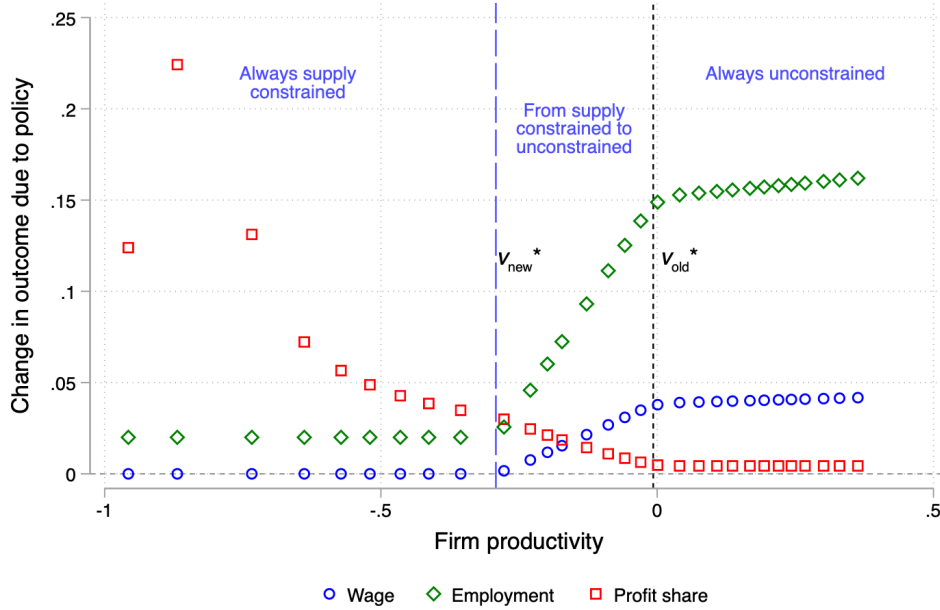
²⁵Bordos et al. (2015) list 26 such policies across 22 different countries as varied as the USA, Brazil, Hungary, Tunisia, Turkey and Sweden.

Figure 7: Simulations based on structural estimation

(a) Rent-sharing level, elasticity and pass-through



(b) Policy simulation: Employment Tax Incentive



Notes: This figure presents simulations based on the structural estimation described in section 6. Panel (a) shows the model-implied rent-sharing level compared to the elasticity and pass-through at different levels of firm productivity. For minimum wage-bound firms, the rent-sharing level is high while the elasticity and pass-through are zero. Panel (b) shows the effect of a simplified employment tax incentive, which shifts the labour supply curve down, and v^* to the left. The old productivity threshold v^* is shown as a black vertical line, while the new policy-induced threshold is marked in blue. While many firms will increase employment as intended by the policy (“always unconstrained”), many other firms will instead absorb the subsidy purely into profits (“always supply-constrained”), with variations in-between. For simplicity of presentation, we trim the small number of demand constrained firms, though these are part of the full simulation and results.

profits. In addition to other reasons given by studies of South Africa’s ETI, this may help explain why such studies struggle to detect significant positive effects on employment ([Ranchhod and Finn, 2016](#); [Ebrahim and Pirttilä, 2025](#)).

7 Theoretical implications

In this section, we discuss the theoretical implications of this evidence in light of the monopsony model in [Section 2](#) and some alternative prominent models of imperfect competition in labour markets. We begin by outlining the baseline versions of these alternative models, restricting our attention to such models which explain rent sharing, and then compare these with the baseline monopsony model and the evidence. We show that these other models do share the prediction that profits increase more in the just-constrained region, but in these models employment still increases in this region, unlike the monopsony model and our empirical evidence. However, we do not claim that one cannot make ad-hoc extensions to these models to accommodate the empirical findings. Instead, we interpret this baseline comparison as reason to think the profit results apply to a range of theoretical models with rent sharing, but that the employment result additionally supports the relevance of a core and interesting mechanism of monopsony.

7.1 Baseline Diamond-Mortensen-Pissarides (DMP) model

The Diamond-Mortensen-Pissarides (DMP) model is the workhorse model of bilateral worker-firm bargaining ([Mortensen and Pissarides 1999](#); [Cahuc et al. 2014](#)). [Appendix B.2](#) derives our baseline version, with simulation details. The key features are that workers and firms are matched through a function depending on the tightness of the labour market (vacancies over unemployment), and that firms face a non-zero cost for posting vacancies. DMP delivers rent sharing because when a worker and firm match, they split the surplus which depends on productivity. Baseline simplifications include that there is a linear production function, no on-the-job search, workers are homogeneous, and the firm discount rate and worker reservation wage are set to zero.

[Figure 8](#) panel (a) presents a simulation of this baseline DMP model. The region on the right shows firms unconstrained by the minimum wage: as one may expect given that surplus increases with productivity, wages and profits also increase with productivity. Vacancies increase too, since firms gain more profits

from each additional match. Appendix Figure B1 shows the proportional split is near-constant in this region.

In the region on the left, where productivity is below the minimum wage, firms post no vacancies as these would be purely loss-making. Like in the monopsony model, it is the middle region which is most of interest. In this region where productivity is just above the minimum wage, firms post vacancies and are constrained to pay the minimum wage, though they would prefer to pay lower wages. The higher minimum wage means a disproportionate surplus goes to wages, and so – as in the monopsony case – productivity increases are absorbed into profits. Figure B1 shows the profit *share* rises throughout this region until it reaches its steady-state share.

However, unlike in the monopsony case, vacancies and employment *increase* in this middle region.²⁶ The intuition is that, as vacancies become more profitable, firms increase vacancies at a greater rate, before reaching a steady-state slope in the unconstrained region.

Overall, in addition to delivering rent sharing and employment growth with productivity when unconstrained, DMP also features a rapid increase in the mark-down or profit share of just-constrained firms. These features are consistent with the evidence we present in sections 4 and 5. However, contrary to our evidence, the rate of increase in hires and employment is higher in the just-constrained region compared to the unconstrained region.

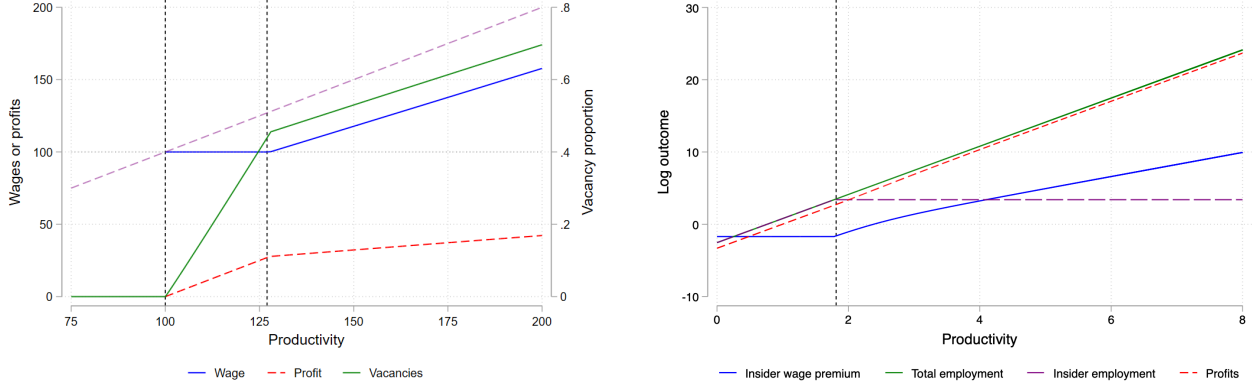
A natural question is whether an extension to on-the-job search, as in Postel-Vinay and Robin (2002), could rationalize an employment kink. Since such an extension is non-trivial (Flinn, 2006), we simply note that the intuition for DMP above still holds strongly in this case. Monopsonistic firms only have the wage margin to adjust employment, while firms that bargain individually additionally have a vacancy-posting decision. As profits increase in the region of no rent sharing, while monopsonistic firms are therefore constrained to the same employment, higher profits incentivise individually-bargaining firms to post more vacancies, thereby increasing employment.

7.2 Baseline union bargaining model

Another prominent set of models which delivers rent sharing involves unions bargaining with firms. We focus on a baseline version of the insider-outsider model

²⁶Given exogenous separations, vacancies are proportional to employment N . This can be seen from $R = sN$, where in steady state recruits R are equal to separations, and the separation rate s is constant.

Figure 8: Simulations of alternative models: Firm outcomes with productivity



(a) Diamond-Mortensen-Pissarides (DMP)

(b) Union bargaining

Notes: Panel (a) shows the DMP model (Mortensen and Pissarides 1999); see Appendix B.2 for details of the model and simulation. The minimum wage is set at 100, and the wage is missing where there are no vacancies. The purple dashed line is the 45-degree line representing productivity. The left vertical line indicates where productivity equals the minimum wage (vacancies first profitable), and the right line indicates where the minimum wage is no longer binding (analogous to v^* in Section 2). Vacancies are proportional to firm employment. Panel (b) shows the union bargaining model with insiders and outsiders; see Appendix B.3 for details. The minimum wage is 1.5, with the vertical line indicating where maximum insider employment is reached (analogous to v^* in Section 2).

(Cahuc et al. 2014), with theory and simulation details given in Appendix B.3. “Insiders”, represented by the union, bargain for a wage premium above the minimum wage. A key assumption is that the firm can hire “outsiders” at the minimum wage, conditional on insiders being retained.²⁷ Insiders who are fired are still given the premium as severance. This means that the marginal cost of labor to the firm is just the minimum wage, since the additional insider cost is fixed.

The insider-outsider model thus delivers a weaker form of rent sharing in that only insiders, not all workers, see wages increase with productivity. Empirically, this is consistent with many studies which find higher pass-through to incumbent workers (Cho and Krueger, 2022; Garin and Silv rio, 2023; Kline et al., 2019). Figure 8 panel (b) shows in the region to the right that the insider wage premium increases with productivity, as does total employment and profits. In the region to the left, where the firm retains less than the number of insiders, the union prefers to ensure members are employed and so only maintains a minimal premium until all insiders are retained. Profits and total employment increase at a steady rate across the wage-premium kink-point, however, providing a contrast to both the

²⁷One can think of this as imposing a limit on the size of the union, or as local changes in employment, where over a longer time horizon the number of insiders may change.

DMP and monopsony models.

7.3 Discussion of models and evidence

The previous discussion shows how our key variables of interest respond in other rent-sharing models. We show there is a wage-kink point in at least two prominent alternative models, the DMP and insider-outsider models. The insider-outsider model demonstrates that such a wage-kink need not necessarily generate a corresponding kink in profit share at this point; on the other hand, DMP shows that this kink in profit share is not unique to monopsony models, with similar intuition of recovering the optimal unconstrained level of profits. However, DMP implies *higher* employment growth in this region of just-constrained firms (and the insider-outsider model has no differential employment prediction), and so the differentially *lower* employment growth of the monopsony models and empirical results is not as easy to explain with other models.

In our view, it is likely that a mixture of these, and other, models applies to firms even within the same labour market. [Kline \(2025\)](#) presents a model in which bargaining and monopsony-like mechanisms operate in the same firm. In [Section 6.2](#), we discussed how our simple model mechanism may only apply to a part of the labour market, with the rest explained by noncompliance or, here, non-monopsonistic models such as in DMP or unions. [Appendix section E.2](#) shows that the parameter estimates are again very similar in such a mixture model. We therefore would not like to claim that the differentially lower employment response is evidence *against* the other models, but rather that it suggests the relevance of the core mechanism of monopsony described in [Section 2](#).

8 Conclusion

This paper starts with the stylized fact that in the presence of binding minimum wages, lower-productivity “constrained” firms tend to pay around the minimum wage, and as such, wages do not increase much with revenue-productivity for these firms in the cross-section. We show that per a very general model of monopsonistic firms, a firm whose preferred wage lies just below a minimum wage will not raise its own wage (i.e. no pass-through) nor expand employment when it experiences a revenue-productivity increase. Instead, the firm maintains the minimum wage and absorbs the additional revenue into a higher markdown. This behavior persists until the firm’s revenue-productivity grows enough such that the minimum wage

is no longer binding, and normal monopsonistic rent-sharing behavior – raising wages to attract more labour – resumes.

We test this prediction in South African administrative data, finding strong support for the theoretical predictions. It is particularly striking that we find support for our results using three different sources of variation: cross-sectional variation in productivity, within-firm changes in value-added using the [Lamadon et al. \(2022\)](#) approach, and within firm variation driven by trade shocks in the spirit of [Garin and Silvério \(2023\)](#).

The productivity region we identify is substantively important, with about a quarter of formal firms constrained in our data. Using a structural model, we estimate the effects of prominent firm- and labour-market policies which upgrade productivity or subsidize employment, and show that the mechanism we identify substantively diminishes the positive effects of these policies. We discuss the extent to which our mechanism and results are compatible with other models of the labour market, such as individual or union wage bargaining.

While our mechanism applies across labour market settings, it may be especially salient in developing countries. The proportion of firms affected by the supply-constrained region is larger when minimum wages are relatively high, as in South Africa and other developing countries.²⁸ The size of the supply-constrained region also increases in our model when there is more monopsony power, and the limited existing evidence (including our own) indicates that monopsony power may in fact be higher in developing countries.²⁹ The mechanism we identify may thus weaken a developmental path predicated on firms sharing the gains of productivity growth in the form of higher wages and expanded employment.³⁰ The mechanism may also help explain part of the the common developing-country complaint of “stalled” development or “jobless growth” ([Kannan and Raveendran 2009](#); [Sanyal 2014](#)).

In terms of policy, our paper prompts a re-evaluation of the welfare effects of minimum wages in monopsonistic settings. That minimum wages can increase efficiency and improve worker welfare under static monopsony is well understood. In response to the introduction of a binding minimum wage, firms in the supply-

²⁸Appendix Figure [A12](#) shows the minimum to median wage, or Kaitz index, by cross-country gross national income. The Kaitz index is about 14 percentage points higher for lower and middle income countries, with a standard error of 5).

²⁹See [Amodio and De Roux \(2022\)](#); [Dal Bó et al. \(2013\)](#); [Naidu et al. \(2016\)](#); [Sharma \(2023\)](#), with firm labour supply elasticities of 1-2.5 in Colombia, India, Mexico, and UAE.

³⁰Correspondingly, other settings will be less affected. For example, [Berger et al. \(2025\)](#) find a negligible supply-constrained region in the United States. We replicate this finding in our simple model from Section 2: if we use the firm-facing labour supply elasticity and minimum wage bite parameters from [Berger et al. \(2025\)](#), the supply-constrained region is only 4% of firms.

constrained region will increase wages and employment, at the expense of reduced firm monopsony rents. However, we show that in a dynamic setting, such worker welfare benefits in the supply-constrained region will erode as firms' revenue productivity increases, because wages and employment remain constant as firms claw back their markdowns. Thus the static minimum wage cannot force such firms to accept lower monopsony rents in perpetuity.

More generally, when labour supply constraints bind (e.g. subsistence or efficiency wages), and firms earn profits below their desired level, there too they may choose to respond to market and productivity expansions by simply absorbing these gains as windfall profits rather than sharing benefits with workers. Ensuring that workers capture the benefits of productivity increases and growth thus requires additional dynamic interventions, such as a policy of raising minimum wages in line with productivity, or more fundamental reforms which reduce monopsony power and empower countervailing institutions.

References

- Abowd, J. M., Kramarz, F., and Margolis, D. N. (1999). High wage workers and high wage firms. *Econometrica*, 67(2):251–333.
- Akerberg, D. A., Caves, K., and Frazer, G. (2015). Identification properties of recent production function estimators. *Econometrica*, 83(6):2411–2451.
- Amodio, F. and De Roux, N. (2022). Labor market power in developing countries: Evidence from Colombian plants. Mimeo, Preprint.
- Armangué-Jubert, T., Guner, N., and Ruggieri, A. (2025). Labor market power and development. *American Economic Review: Insights*, 7(2):177–195.
- Azar, J. and Marinescu, I. (2024). Monopsony power in the labor market. In *Handbook of Labor Economics*, volume 5, pages 761–827. Elsevier.
- Bassier, I. (2022). Collective bargaining and spillovers in local labor markets. CEP Discussion Paper 1895, Centre for Economic Performance.
- Bassier, I. (2023). Firms and inequality when unemployment is high. *Journal of Development Economics*, 161.
- Bassier, I. and Manning, A. (2025). Estimating labour market power: The long and short of it. Discussion paper 2108, Centre for Economic Performance.
- Bell, B., Bukowski, P., and Machin, S. (2024). The decline in rent sharing. *Journal of Labor Economics*, 42(3):683–716.
- Berger, D., Herkenhoff, K., and Mongey, S. (2022). Labor market power. *American Economic Review*, 112(4):1147–1193.
- Berger, D., Herkenhoff, K., and Mongey, S. (2025). Minimum wages, efficiency, and welfare. *Econometrica*, 93(1):265–301.
- Bhorat, H., Kanbur, R., and Mayet, N. (2012). Minimum wage violation in south africa. *International Labour Review*, 151(3):277–287.
- Bhuller, M., Moene, K. O., Mogstad, M., and Vestad, O. L. (2022). Facts and fantasies about wage setting and collective bargaining. *Journal of Economic Perspectives*, 36(4):29–52.
- Bordos, K., Csillag, M., and Scharl, A. (2015). What works in wage subsidies for young people: A review of issues, theory, policies and evidence. ILO Working Paper 994898973402676, International Labour Organization.
- Breza, E., Kaur, S., and Shamdasani, Y. (2021). Labor rationing. *American Economic Review*, 111(10):3184–3224.
- Budlender, J. and Ebrahim, A. (2021). Estimating employment responses to south africa’s employment tax incentive. Technical report, WIDER Working Paper.
- Cahuc, P., Carcillo, S., and Zylberberg, A. (2014). *Labor economics*. MIT press.

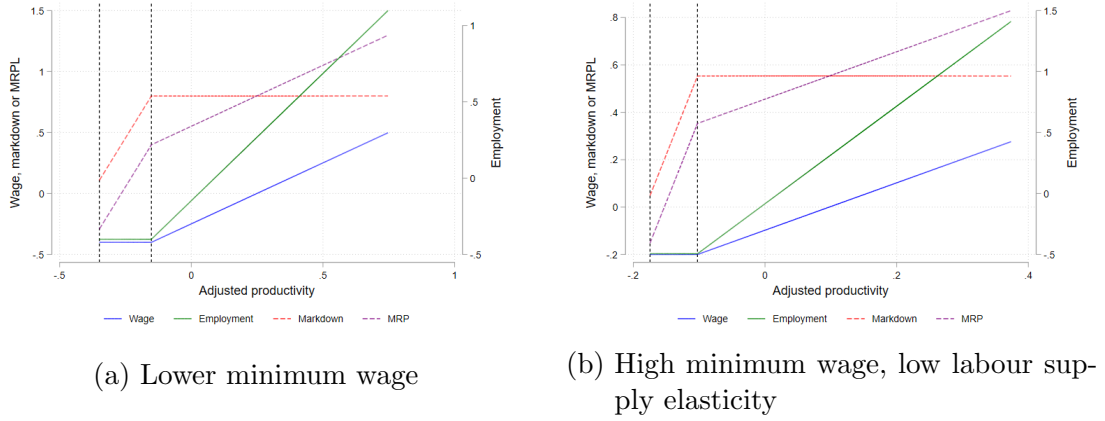
- Card, D., Cardoso, A. R., Heining, J., and Kline, P. (2018). Firms and labor market inequality: Evidence and some theory. *Journal of Labor Economics*, 36(S1):S13–S70.
- Card, D., Cardoso, A. R., and Kline, P. (2016). Bargaining, sorting, and the gender wage gap: Quantifying the impact of firms on the relative pay of women. *The Quarterly Journal of Economics*, 131(2):633–686.
- Card, D., Lee, D. S., Pei, Z., and Weber, A. (2015). Inference on causal effects in a generalized regression kink design. *Econometrica*, 83(6):2453–2483.
- Card, D., Mas, A., and Rothstein, J. (2008). Tipping and the dynamics of segregation. *The Quarterly Journal of Economics*, 123(1):177–218.
- Cengiz, D., Dube, A., Lindner, A., and Zipperer, B. (2019). The effect of minimum wages on low-wage jobs. *The Quarterly Journal of Economics*, 134(3):1405–1454.
- Cho, D. and Krueger, A. B. (2022). Rent sharing within firms. *Journal of Labor Economics*, 40(S1):S17–S38.
- Dal Bó, E., Finan, F., and Rossi, M. A. (2013). Strengthening state capabilities: The role of financial incentives in the call to public service. *The Quarterly Journal of Economics*, 128(3):1169–1218.
- de Frahan, I., Lamadon, T., Meling, T., and Mogstad, M. (2022). Why do larger firms have lower labor shares? NBER Labor Studies Program Meeting, Spring.
- De Loecker, J., Goldberg, P. K., Khandelwal, A. K., and Pavcnik, N. (2016). Prices, markups, and trade reform. *Econometrica*, 84(2):445–510.
- De Loecker, J. and Syverson, C. (2021). An industrial organization perspective on productivity. In *Handbook of industrial organization*, volume 4, pages 141–223. Elsevier.
- Dickens, R., Machin, S., and Manning, A. (1999). The effects of minimum wages on employment: Theory and evidence from Britain. *Journal of Labor Economics*, 17(1):1–22.
- Dong, Y. and Kolesár, M. (2023). When can we ignore measurement error in the running variable? *Journal of Applied Econometrics*, 38(5):735–750.
- Ebrahim, A. and Pirttilä, J. (2025). A policy for the jobless youth in south africa. *Journal of Development Economics*, 172:103394.
- Engbom, N. and Moser, C. (2022). Earnings inequality and the minimum wage: Evidence from brazil. *American Economic Review*, 112(12):3803–3847.
- Flinn, C. J. (2006). Minimum wage effects on labor market outcomes under search, matching, and endogenous contact rates. *Econometrica*, 74(4):1013–1062.
- Gandhi, A., Navarro, S., and Rivers, D. A. (2020). On the identification of gross

- output production functions. *Journal of Political Economy*, 128(8):2973–3016.
- Ganong, P. and Jäger, S. (2018). A permutation test for the regression kink design. *Journal of the American Statistical Association*, 113(522):494–504.
- Garin, A. and Silvério, F. (2023). How responsive are wages to firm-specific changes in labor demand? evidence from idiosyncratic export demand shocks. *Review of Economic Studies*.
- Gouin-Bonenfant, É. (2022). Productivity dispersion, between-firm competition, and the labor share. *Econometrica*, 90(6):2755–2793.
- Grimm, M. and Paffhausen, A. L. (2015). Do interventions targeted at micro-entrepreneurs and small and medium-sized firms create jobs? a systematic review of the evidence for low and middle income countries. *Labour Economics*, 32:67–85.
- Kannan, K. and Raveendran, G. (2009). Growth sans employment: A quarter century of jobless growth in india’s organised manufacturing. *Economic and Political weekly*, pages 80–91.
- Kline, P., Petkova, N., Williams, H., and Zidar, O. (2019). Who profits from patents? rent-sharing at innovative firms. *The Quarterly Journal of Economics*, 134(3):1343–1404.
- Kline, P. M. (2025). Labor market monopsony: Fundamentals and frontiers. Working paper, National Bureau of Economic Research.
- Kroft, K., Luo, Y., Mogstad, M., and Setzler, B. (2025). Imperfect competition and rents in labor and product markets: The case of the construction industry. *American Economic Review*, 115(9):2926–2969.
- Lamadon, T., Mogstad, M., and Setzler, B. (2022). Imperfect competition, compensating differentials, and rent sharing in the us labor market. *American Economic Review*, 112(1):169–212.
- Levinsohn, J. and Petrin, A. (2003). Estimating production functions using inputs to control for unobservables. *The Review of Economic Studies*, 70(2):317–341.
- Lewis, W. A. (1954). Economic development with unlimited supplies of labour. *The manchester school*, 22(2):139–191.
- Manning, A. (1993). Wage bargaining and the phillips curve: the identification and specification of aggregate wage equations. *The Economic Journal*, 103(416):98–118.
- Manning, A. (2003). *Monopsony in motion: Imperfect competition in labor markets*. Princeton University Press, Princeton.
- Manning, A. (2021). Monopsony in labor markets: a review. *ILR Review*, 74(1):3–26.

- Mortensen, D. T. and Pissarides, C. A. (1999). New developments in models of search in the labor market. *Handbook of labor economics*, 3:2567–2627.
- Muralidharan, K., Niehaus, P., and Sukhtankar, S. (2023). General equilibrium effects of (improving) public employment programs: Experimental evidence from india. *Econometrica*, 91(4):1261–1295.
- Naidu, S., Nyarko, Y., and Wang, S.-Y. (2016). Monopsony power in migrant labor markets: evidence from the united arab emirates. *Journal of Political Economy*, 124(6):1735–1792.
- Olley, G. and Pakes, A. (1996). The dynamics of productivity in the telecommunications equipment industry. *Econometrica*, 64(6):1263–1297.
- Palladino, M., Bertheau, A., Hijzen, A., Kunze, A., and the LinkEED team (2025). Firms and the gender wage gap: A comparison of eleven countries. Technical report, Unpublished manuscript.
- Porter, J. and Yu, P. (2015). Regression discontinuity designs with unknown discontinuity points: Testing and estimation. *Journal of Econometrics*, 189(1):132–147.
- Postel-Vinay, F. and Robin, J.-M. (2002). Equilibrium wage dispersion with worker and employer heterogeneity. *Econometrica*, 70(6):2295–2350.
- Ranchhod, V. and Finn, A. (2016). Estimating the short run effects of south africa’s employment tax incentive on youth employment probabilities using a difference-in-differences approach. *South African Journal of Economics*, 84(2):199–216.
- Risch, M. (2024). Does taxing business owners affect employees? evidence from a change in the top marginal tax rate. *The Quarterly Journal of Economics*, 139(1):637–692.
- Rud, J. P. (2012). Electricity provision and industrial development: Evidence from india. *Journal of development Economics*, 97(2):352–367.
- Sanyal, K. (2014). *Rethinking capitalist development: Primitive accumulation, governmentality and post-colonial capitalism*. Routledge India.
- Sharma, G. (2023). Monopsony and gender. Technical report, MIT working paper.
- Verhoogen, E. (2023). Firm-level upgrading in developing countries. *Journal of Economic Literature*, 61(4):1410–1464.
- Yeh, C., Macaluso, C., and Hershbein, B. (2022). Monopsony in the us labor market. *American Economic Review*, 112(7):2099–2138.

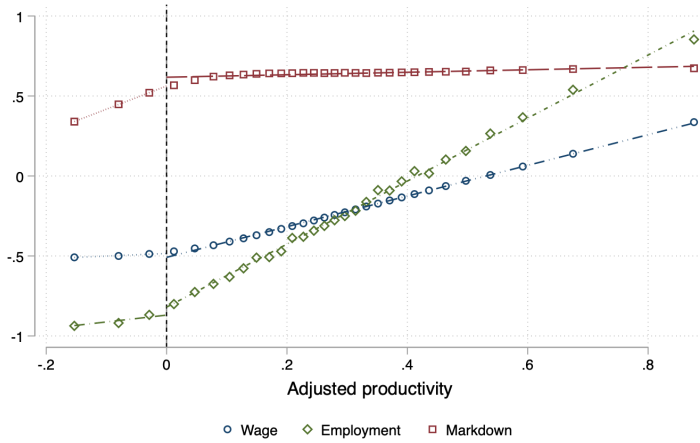
Appendix A: Additional Tables and Figures

Figure A1: Additional simulations (monopsony model)



Notes: Simulations are as in Figure 2 Panel (b), except the minimum wage is lower (−0.5 log points) in Panel (a) and the firm-facing labour supply elasticity is set to 4 in Panel (b). The (supply plus demand) constrained regions are smaller than the baseline in both cases, covering 8.9% and 3.2% of firms respectively.

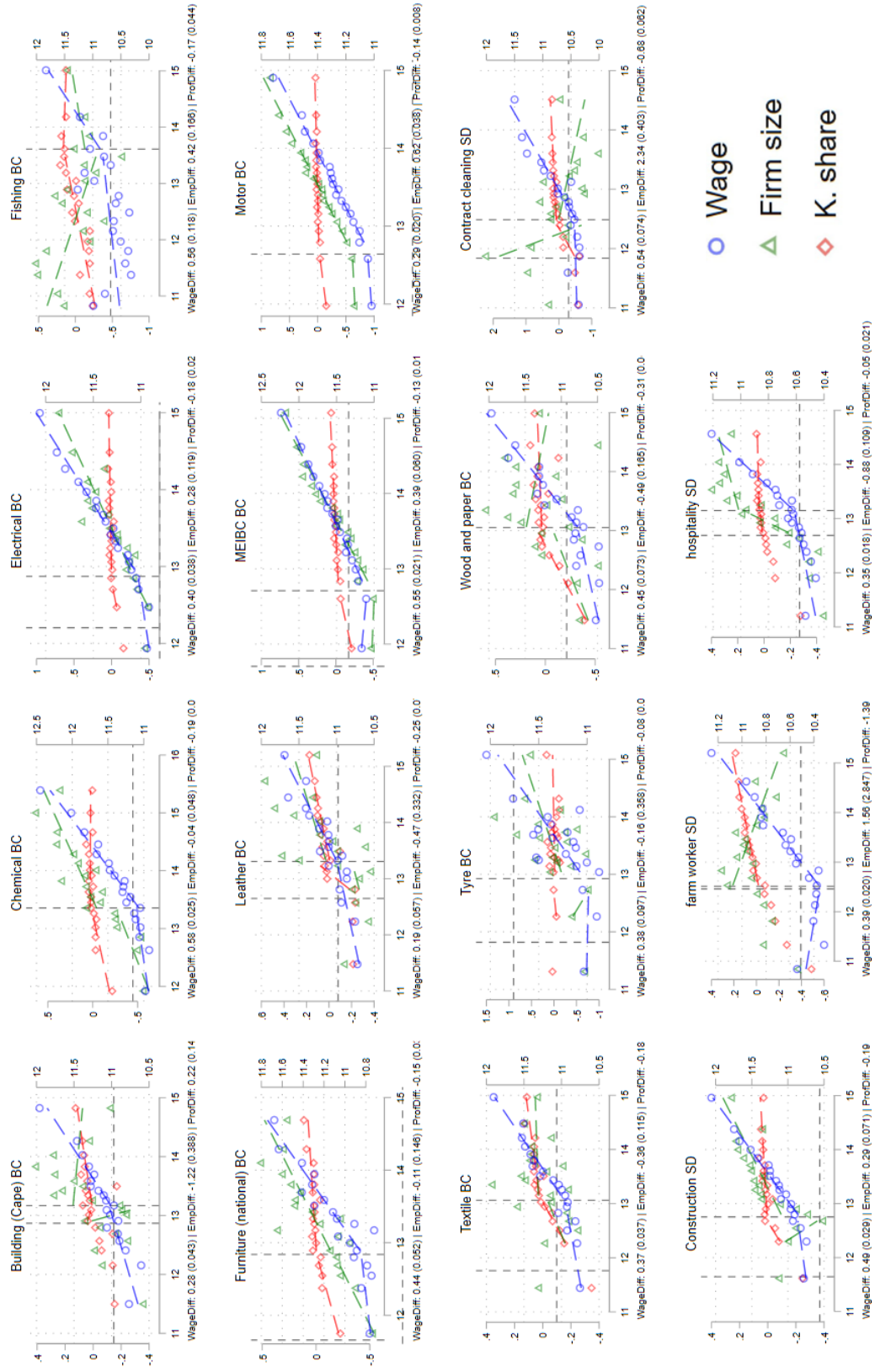
Figure A2: Pooled monopsony model with amenities



Notes: 12.8% of firms supply constrained, 0.1% demand-constrained. Productivity normally distributed. Supply constrained region normalized to end at 0.

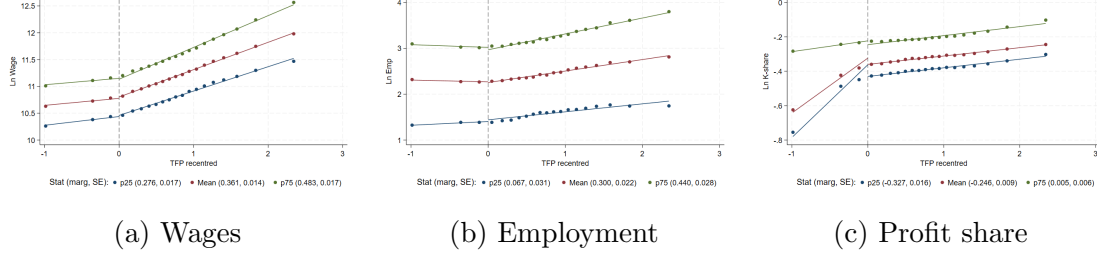
Notes: The plot shows the pooled monopsony simulation as in Figure 3, with the extension that it allows for firm-specific amenities. Amenities are set to positively covary with wages.

Figure A3: Cross-sectional kink design case-studies



Notes: Figure shows firm median wage, employment, and profit share by 20 firm productivity bins (productivity estimated using the ACF method of Section 4; ventiles), for selected BCs and SDs. The algorithm outlined in Section 4 is used to fit underlying firm median wages as a piece-wise continuous linear function of productivity. Analogous linear fits of employment and profit share are then plotted on either side of the identified wage kink. The right-most vertical line is the estimated value of the wage-kink \hat{v}^* . If there is a second left-most line, this identifies the estimated demand-constrained value \hat{v}_1^* . The horizontal line is the BC/SD minimum wage.

Figure A4: Cross-sectional kink design: percentiles of outcome variable



Notes: Figure shows the cross-sectional distribution of firm wages, employment, and profit share (all in logs) against recentered productivity, as in Figure 4, but for different percentiles of the outcome variable. The red markers show the mean value of the outcome for each productivity bin (the same as Figure 4). The green markers show the 75th percentile, and the blue markers the 25th percentile, for each productivity bin.

Table A1: Kink design implied elasticities with respect to value added

	No controls			With controls		
	Constrained	Unconstrained	Difference	Constrained	Unconstrained	Difference
Wage	0.207*** (0.0214)	0.432*** (0.0094)	0.226*** (0.0234)	0.060*** (0.0127)	0.290*** (0.0054)	0.230*** (0.0138)
Employment	0.022 (0.0452)	0.283*** (0.0183)	0.261*** (0.0487)	0.126** (0.0272)	0.415*** (0.0104)	0.289*** (0.0291)
Profit-share	0.505*** (0.0281)	0.035*** (0.0033)	-0.470*** (0.0283)	0.470*** (0.0195)	0.049*** (0.0022)	-0.422*** (0.0196)
F-stat	452.3	2156.9		893.9	3565.5	
N	237763	978665		179035	807201	

Notes: Table shows implied elasticities from cross-sectional kink design (Section 4) for constrained and unconstrained firms, and their differences, with and without controls. The wage (rent-sharing), employment and profit share elasticities are with respect to value-added. Elasticities are estimated by running a cross-sectional regression of the (log) outcome on (log) value-added, where value-added is instrumented with the recentered productivity value. Controls refer to the average AKM worker fixed at the firm, the firm poaching ratio, and industry (2-digit) by region (district council) labour market dummies. 5% of firms above and below the threshold \hat{v}^* are dropped to create a “donut” estimate. Standard errors are shown in parentheses and clustered at labour market by event; differences are calculated using the Delta method. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

Table A2: Model-exogenous characteristics of constrained and unconstrained firms

	All firms			Narrow band		
	Const- rained	Uncons- trained	Difference	Const- rained	Uncons- trained	Difference
Monthly minimum wage	4127.00	4501.51	374.284*** (85.0909)	4250.98	4226.41	-28.930 (28.9213)
Metro	0.59	0.68	0.095*** (0.0080)	0.59	0.60	0.014* (0.0076)
Sector proportions						
Primary sector	0.09	0.06	-0.028*** (0.0083)	0.06	0.06	-0.002 (0.0044)
Manufacturing	0.23	0.29	0.061*** (0.0210)	0.26	0.25	-0.009 (0.0067)
Construction	0.14	0.14	0.001 (0.0119)	0.16	0.16	0.008 (0.0065)
Wholesale & Retail	0.42	0.37	-0.047*** (0.0150)	0.42	0.43	0.006 (0.0076)
Infrastructure services	0.02	0.04	0.022*** (0.0066)	0.02	0.02	-0.001 (0.0024)
Bus. & Pers. services	0.10	0.09	-0.010** (0.0048)	0.09	0.08	-0.002 (0.0048)
Province proportions						
Western Cape	0.24	0.21	-0.038*** (0.0124)	0.24	0.23	-0.006 (0.0065)
Eastern Cape	0.07	0.06	-0.015*** (0.0051)	0.07	0.07	0.001 (0.0047)
Northern Cape	0.02	0.02	-0.004** (0.0016)	0.02	0.02	-0.003 (0.0022)
Free State	0.04	0.04	-0.008*** (0.0030)	0.05	0.04	-0.004 (0.0031)
KwaZulu-Natal	0.16	0.15	-0.012* (0.0069)	0.17	0.17	-0.004 (0.0056)
North West	0.03	0.02	-0.008** (0.0033)	0.03	0.03	-0.004 (0.0036)
Gauteng	0.33	0.43	0.100*** (0.0162)	0.33	0.35	0.015** (0.0072)
Mpumalanga	0.06	0.05	-0.009* (0.0051)	0.06	0.06	0.002 (0.0039)
Limpopo	0.03	0.02	-0.007** (0.0028)	0.03	0.03	0.002 (0.0024)

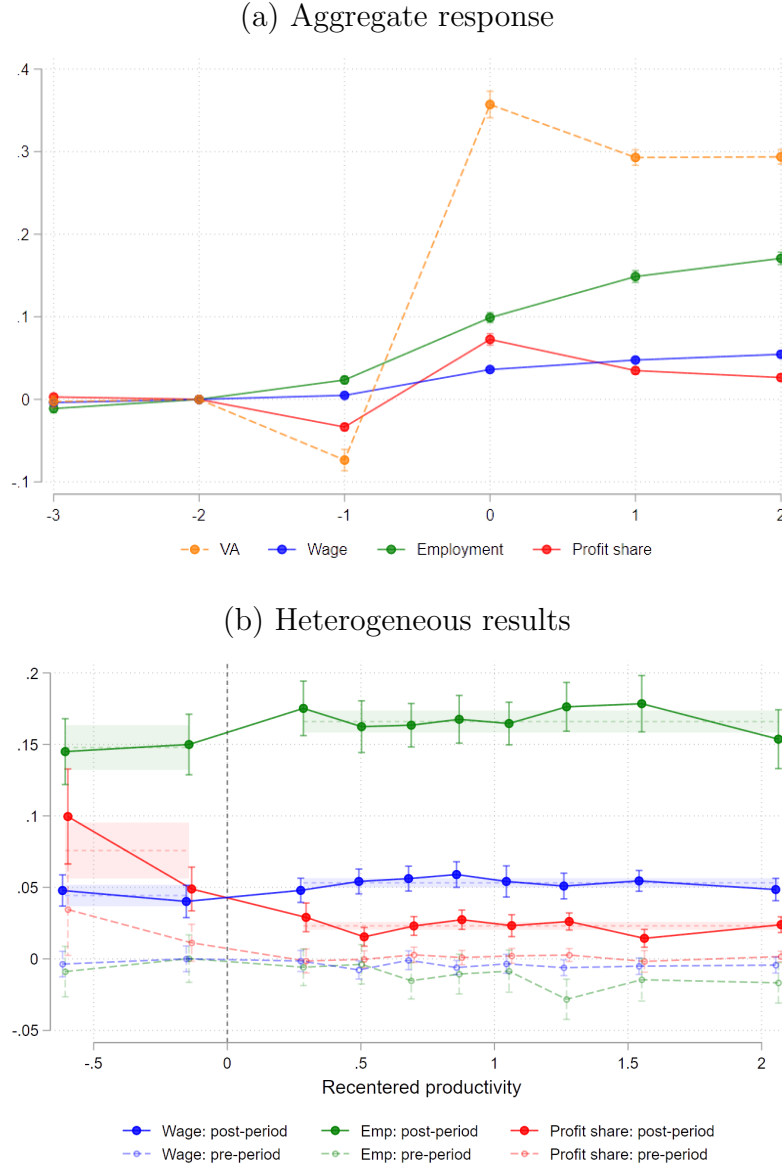
Notes: Table shows descriptive statistics (means) of the pooled cross-sectional sample, by firms' constrained status, for variables exogenous to the model. "All firms" refers to all firms in the cross-section. "Narrow band" refers to a subset of firms which are close to the wage-kink (constrained/unconstrained boundary); specifically the top 10% of the constrained-firm recentered productivity distribution and the same number of lowest-recentered productivity unconstrained firms. For the sectors: "Primary sector" refers to agriculture and mining; "Infrastructure services" refers to electricity, gas, and water supply as well as transport, storage, and communication; while Business and Personal Services refers to financial intermediation, insurance, real Estate, and business services as well as community, social, and personal services. "Metro" is a dummy referring to a firm being located in one of South Africa's eight metropolitan municipalities. The monthly minimum wage refers to the BC/SD minimum in 2018 Rands. Differences between means are shown with standard errors clustered at the labour-market by event level. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

Table A3: Internal IV reduced form results (semi-elasticities)

	Internal IV (large VA change)		
	Constrained	Unconstrained	Difference
Panel (a)			
Wage	0.044*** (0.0038)	0.053*** (0.0016)	0.009** (0.0042)
Employment	0.148*** (0.0079)	0.166*** (0.0039)	0.018** (0.0089)
Profit-share	0.076*** (0.0099)	0.023*** (0.0014)	-0.053*** (0.0100)
N firms	6167	34120	
Obs	81576	589404	
Panel (b)			
Hires	0.134*** (0.0194)	0.198*** (0.0073)	0.064*** (0.0208)
N firms	5986	33577	
Obs	63071	483531	

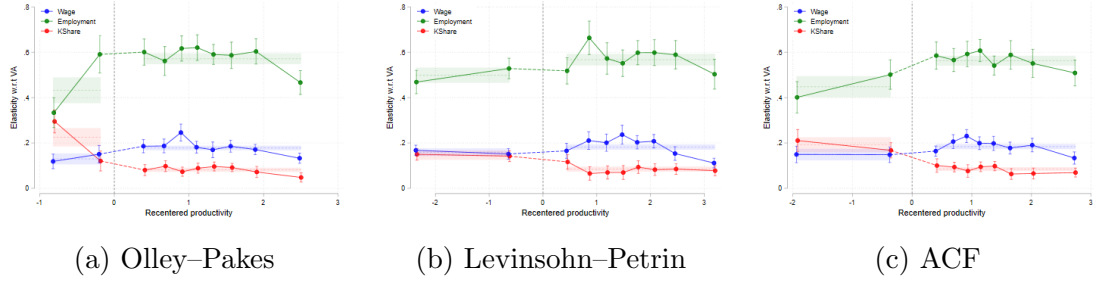
Notes: The table shows the reduced form results for the internal instrument event study, where treatment is above-median increase in firm value-added between periods -1 and 0. Panel (a) shows the results for the main balanced sample with effects on wages (of stayers), employment and profit share (all in logs). Panel (b) shows effects on hires. Estimates are normalised relative to period -2; the effects are average treatment effects across post-periods 1 and 2. See Section 5 for sample restrictions and specification. Estimates are shown separately for constrained and unconstrained firms, as well as the difference between estimates for constrained versus unconstrained. Standard errors are shown in parentheses and clustered at labor market by event; for differences these are calculated using the Delta method. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

Figure A5: Internal IV reduced form results (semi-elasticities)



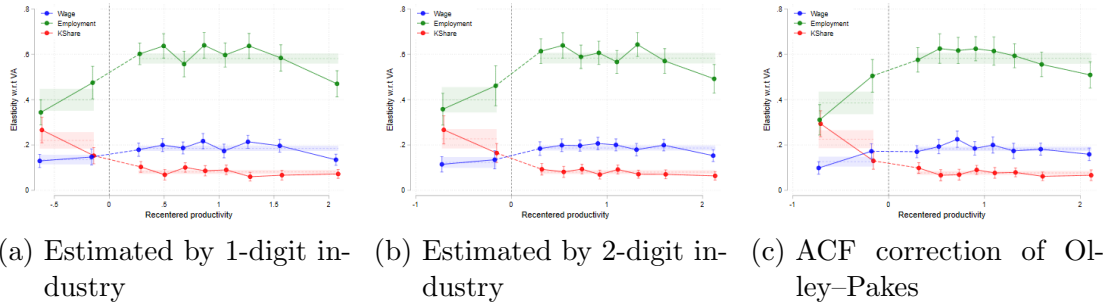
Notes: The figure shows the reduced form results for the internal instrument event study, where treatment is above-median increase in firm value-added between periods -1 and 0. Panel (a) shows the aggregate stacked event study for the full estimation sample. Estimates are normalised relative to period -2. Orange line shows response of log value-added, green log firm employment, blue the log of median wage of firm stayers (incumbents), and red log of profit share of value added. See section 5 for sample restrictions and specification. 95% confidence intervals are shown with vertical bars. Panel (b) shows estimates from event studies as in Panel (a), but estimated by productivity bin. The horizontal axis is firm productivity (estimated using the ACF method) recentered around the estimated productivity wage-kink \hat{v}^* (see Section 4). Recentered productivity equals zero at the wage-kink, shown with a vertical dashed line, and this line divides minimum wage-constrained (to the left) and -unconstrained (to the right) firms. Ten approximately equally-sized productivity bins (deciles) are created. The solid lines and points show the average treatment effect across post-periods 1 and 2. The dashed lines and hollow points show effects estimated for pre-period -3. 95% confidence intervals are shown with vertical bars. The horizontal dashed lines with attendant shaded regions (95% confidence intervals) show applicable post-period treatment effects estimated across the productivity bins, pooled separately below and above the wage-kink value.

Figure A6: Internal IV results: Cobb–Douglas production function variations



Notes: Figure shows robustness versions of the main results in Figure 5. Each panel uses a different method for estimating the production function underlying the recentered productivity term in a Cobb–Douglas framework. Panel (a) uses Olley–Pakes; Panel (b) uses Levinsohn–Petrin; and Panel (c) uses Akerberg–Caves–Frazer (ACF). Specifications are otherwise the same as in Figure 5.

Figure A7: Internal IV results: Alt. production function estimation routines



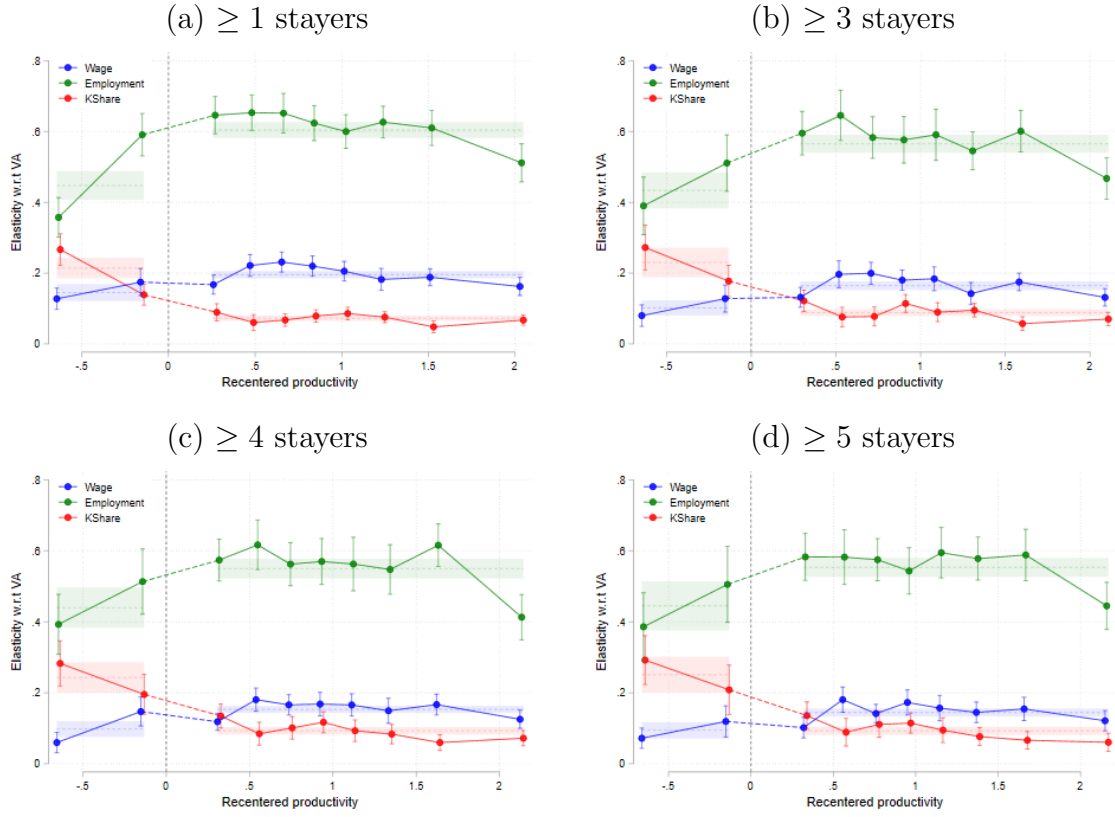
Notes: Figure shows robustness versions of the main results in Figure 5. Panels (a) and (b) show results when production functions are estimated separately by 1-digit and 2-digit industries, respectively. Panel (c) shows results when the ACF correction is applied to Olley–Pakes rather than Levinsohn–Petrin production function estimation. The specification is otherwise the same as in Figure 5, including the translog production function.

Table A4: Rent sharing levels and pass-through estimates

	Internal IV (large VA change)		External IV (trade shock)	
	Constrained	Unconstrained	Constrained	Unconstrained
Rent sharing level	0.796	0.337	0.988	0.327
Rent sharing elasticity	0.132	0.186	0.129	0.188
Pass-through level	0.105	0.063	0.127	0.061

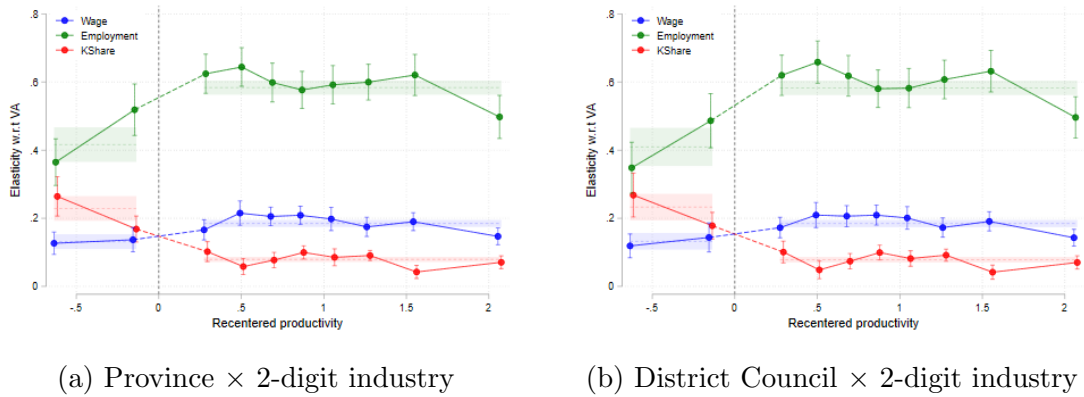
Notes: Table shows the rent sharing levels, rent sharing elasticities, and implied pass-through levels from the main specifications. The rent sharing elasticities come from Table 2. The rent sharing level is calculated as the firm total wage bill divided by firm value added, for the estimation sample upon which the rent sharing elasticity is calculated. The pass-through (the dollar increase in wages for a one dollar increase in value added) is the product of the elasticity and the level. Constrained and unconstrained firms are as defined in Table 2.

Figure A8: Internal IV results: various “stayer” sample restrictions



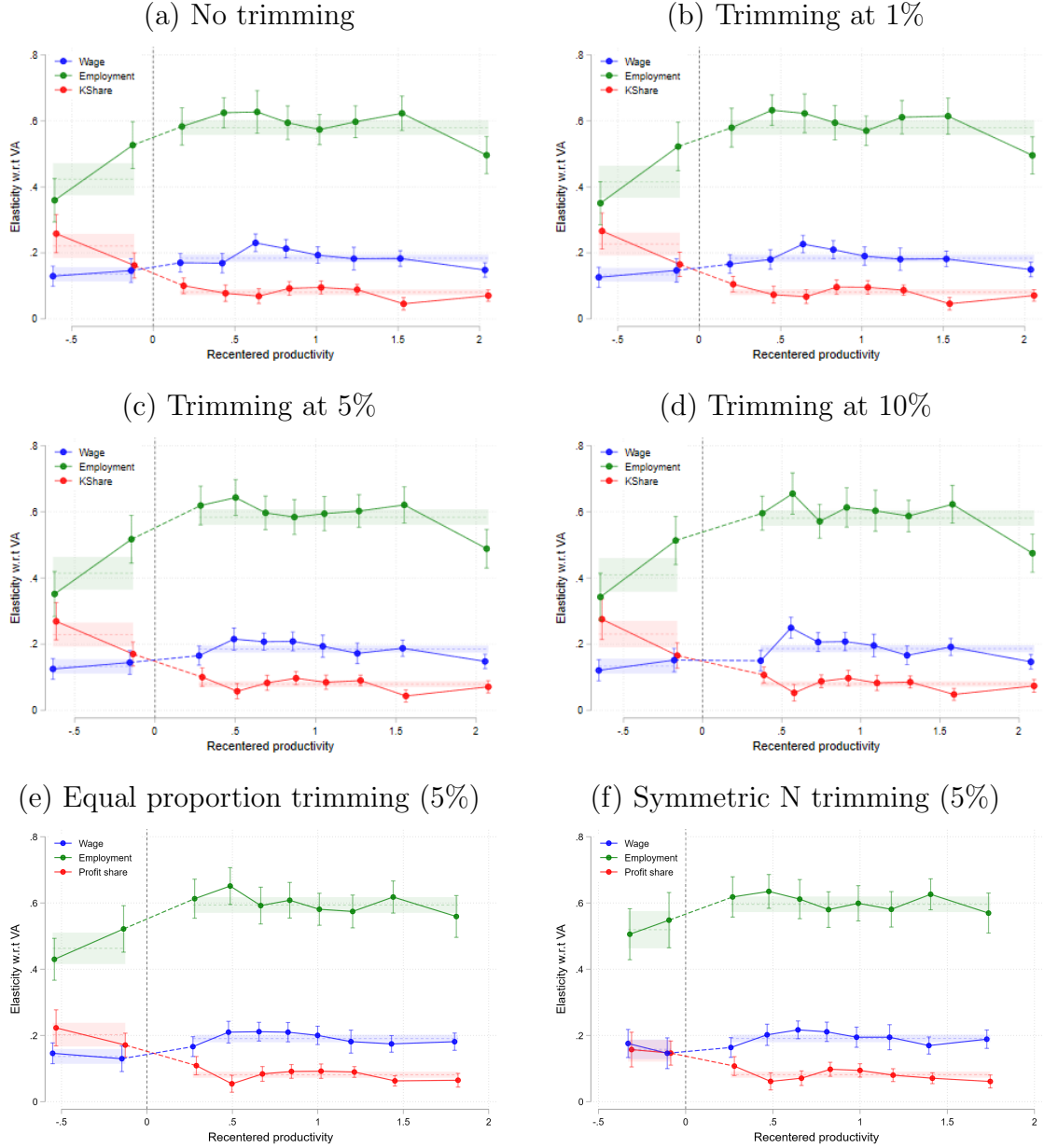
Notes: Figure shows robustness versions of main results in Figure 5. Each panel shows results when the sample is restricted to have at least 1, 3, 4 or 5 stayers over the event-study period. The specification is otherwise the same as in Figure 5.

Figure A9: Internal IV results: Various labour market definitions



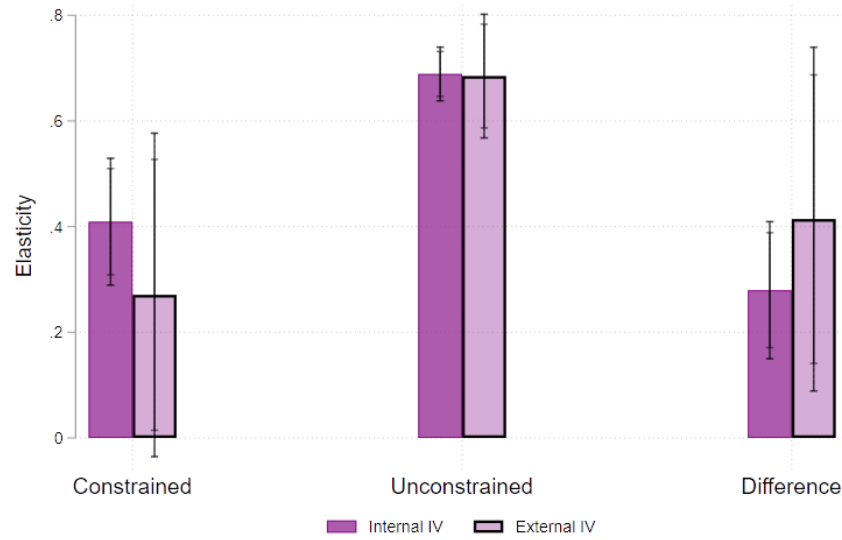
Notes: Figure shows robustness versions of the main results in Figure 5. Each panel shows results when labour market fixed effects are defined by interactions of different geography and industry variables. There are nine provinces and 52 district councils. The one-digit industry level has nine categories, while the two-digit level has 50 categories. The specification is otherwise the same as in Figure 5.

Figure A10: Internal IV: trimming around recentered productivity threshold



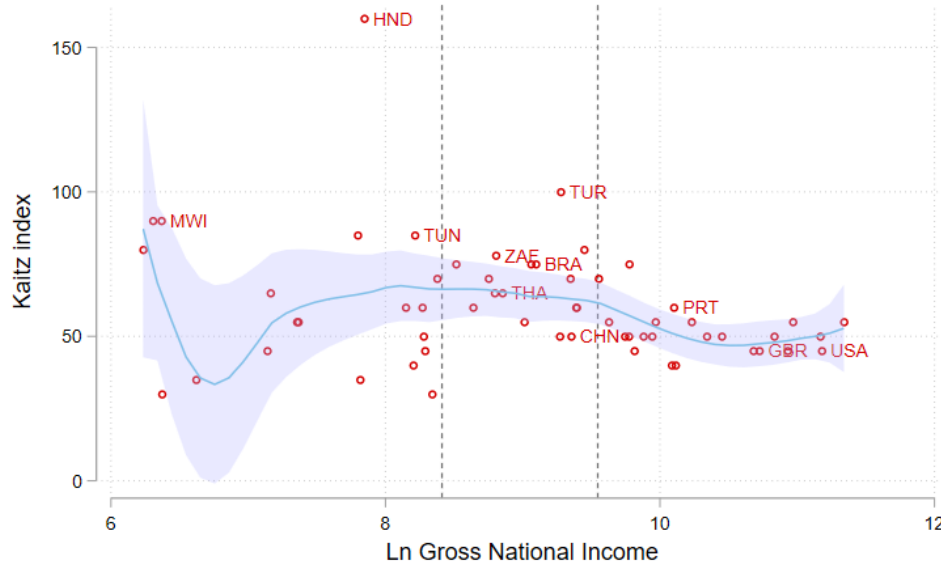
Notes: Figure shows robustness versions of main results in Figure 5. Each of panels (a)-(d) shows results for different choices of dropping firms close to the productivity threshold, where “Trimming at $x\%$ ” means dropping the most productive $x\%$ of constrained firms and least productive $x\%$ of unconstrained firms. Panels (e) and (f) shows results for different approaches to dropping firms in the tails of the recentered productivity distribution. Panel (a) drops the 5% of constrained firms with the lowest recentered productivity and 5% of unconstrained firms with the highest productivity. Panel (b) instead drops the bottom and top 5% of the combined firm distribution, which removes more firms overall—particularly among constrained firms. The specification is otherwise the same as in Figure 5.

Figure A11: Comparison of estimates from internal and external IV shocks: Hires



Notes: The figure plots the Hires estimates from Table 2 Panel (b), where the internal IV refers to above-median increases in firm VA between event periods -1 and 0 while the external IV refers to the shift-share trade shocks. See section 5 for sample and specification details. Vertical bars represent 90% and 95% confidence intervals, where standard errors are clustered at labor market by event and differences these are calculated using the Delta method.

Figure A12: Kaitz Index by cross-country GNI



Note: Kaitz index is minimum / median wage (%). Vertical lines indicate World Bank income groups. Average kaitz in low GNI countries is 63%, mid is 68%, and high is 52%.

Notes: The plot shows the Kaitz index, defined as the ratio of the minimum wage to the median wage, by country. The x-axis is the log gross national income, using the Atlas method. The vertical lines indicate the World Bank classifications of lower (left) and middle (centre) income countries. Data are from ILO and World Development Indicators.

Table A5: External IV shock results using pre- and post-period restrictions

	External IV (trade shock)		
	Constrained	Unconstrained	Difference
Panel (a)			
Rent-sharing	0.079 (0.1161)	0.214*** (0.0449)	0.134 (0.1245)
Employment	0.135 (0.2168)	0.439*** (0.1533)	0.304 (0.2656)
Profit-share	0.348** (0.1455)	0.040 (0.0667)	-0.308* (0.1601)
F-stat	6.0	35.7	
N firms	956	12017	
Obs	6621	112122	
Panel (b)			
Hires	0.296 (0.5010)	1.380** (0.5508)	1.084 (0.7446)
F-stat	5.7	12.5	
N firms	918	11644	
Obs	5573	94257	

Notes: This table presents results analogous to the External IV results of Table 2 (the right-hand super column), but where the trade shock IV specification uses the pre- versus post-period structure of the data. Specifically, the export and import shares are estimated only over the pre-periods used in the internal instrument approach, and GDP shocks and effects on outcomes are only considered in the equivalent post-periods (see Equations 8 and 9). Standard errors are shown in parentheses and clustered at labor market by event; for differences these are calculated using the Delta method. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

Appendix B: Models

B.1 Details of monopsony model

The full model in Manning (2003, pp. 338-345) is different from the simplified model in Section 2.1 mainly in that Manning (2003) incorporates the average market wage as a determinant of aggregate labour supply and a firm-specific supply-shifter b_i (e.g. disamenities), so that the firm-specific labour supply depends on the firm wage premium relative to the market wage and the firm-specific disamenity. The model below is essentially a stripped-down re-presentation of Manning (2003).

Specifically, retain Equation 2 for the demand for labour, but now model the share of total employment (N) supplied to firm i (N_i) as a function of its own wage (W_i) relative to an average market-level wage index (W) and the value of its disamenity (B_i): $\frac{N_i}{N} = \left(\frac{W_i}{B_i W}\right)^{1/\varepsilon}$. If one then models the labour supply to the whole market as $N = N_0 W^\phi$ and takes logs (again denoting logs of variables as lower case letters), the labour supply to the individual employer is

$$w_i = (1 - \varepsilon\phi)w + \varepsilon(n_i - n_0) + b_i,$$

or, subsuming n_0 into b_i and defining the coefficient on the average wage as θ ,

$$w_i = \theta w + \varepsilon n_i + b_i. \tag{B1}$$

The marginal cost of labor in the absence of the minimum wage is then

$$\text{mcl}_i = \ln(1 + \varepsilon) + w_i = \ln(1 + \varepsilon) + \varepsilon n_i + \theta w + b_i, \tag{B2}$$

which diverges from the simplified Equation 1 in its two additional terms reflecting the influence of the average wage and the firm-specific disamenity. Equating the expression for the MRPL in Equation 2 to the MCL above (Equation B2), and substituting in Equation B1, the firm's unconstrained wage is given by:

$$w_i^* = \frac{\eta\theta w - \varepsilon \ln(1 + \varepsilon)}{\eta + \varepsilon} + v_i \tag{B3}$$

with

$$v_i = \frac{\varepsilon a_i + \eta b_i}{\eta + \varepsilon}, \tag{B4}$$

while the unconstrained employment level is

$$n_i^* = \frac{-\theta w - \ln(1 + \varepsilon) + a_i - b_i}{\eta + \varepsilon}. \quad (\text{B5})$$

With the introduction of a minimum wage w_m , the discussion of the simplified model in Section 2 explains how the value of a firm's "adjusted productivity" term v_i relative to the thresholds v^* and v_1^* determines which of the qualitative distinct demand-constrained, supply-constrained or unconstrained regions it falls into. Expressions for these threshold values can be derived by noting that v^* is the value of v_i where the unconstrained wage w_i^* is greater than or equal to the minimum wage w_m , so that, from Equation B3,

$$v^* = w_m - \frac{\eta\theta w - \varepsilon \ln(1 + \varepsilon)}{\eta + \varepsilon}. \quad (\text{B6})$$

For firms which have $v_i < v^*$, for some it will be optimal to accept all workers forthcoming at the minimum wage w_m ; these are supply-constrained firms. However for other firms with even lower v_i , it is not profitable to employ all the workers forthcoming at the minimum wage w_m ; these are the demand-constrained firms. To find the threshold value of v_i which delineates these sets of firms, v_1^* , note that these firms set their wage at w_m but choose employment less than the potential supply at that wage so that $\text{mrpl}_i = w_m$. From Equations 2 and B1:

$$v_1^* = w_m - \frac{\theta\eta w}{\eta + \varepsilon}. \quad (\text{B7})$$

In order to find the equilibrium level of employment for supply-constrained firms, one can substitute $w_i = w_m$ into the labour supply, Equation B1, which Manning (2003) shows can be expressed as:

$$n_i^{\text{sc}} = n(w, a_i, b_i) + \frac{1}{\varepsilon}(v^* - v_i), \quad (\text{B8})$$

where $n(w, a_i, b_i)$ is the unconstrained employment level given in Equation B5. For our purposes it is useful to note that a_i does not enter Equation B1, and therefore does not enter the expression for n_i^{sc} , which reflects our main insight that equilibrium employment for supply-constrained firms is unaffected by local shifts in (revenue-) productivity, and in the special case where $b_i = 0$ all supply-constrained firms will have the same employment level, corresponding to the labour supplied at the minimum wage.

To find the equilibrium employment for demand-constrained firms, again use that they will choose employment such that $\text{mrpl}_i = w_m$, and some rearranging leads to

$$n_i^{\text{dc}} = n(w, a_i, b_i) + \frac{\ln(1 + \varepsilon)}{\eta + \varepsilon} - \frac{1}{\eta}(v_1^* - v_i). \quad (\text{B9})$$

B.2 Baseline Diamond-Mortensen-Pissarides (DMP) model

We use a standard presentation of the baseline DMP model (e.g. [Cahuc et al. 2014](#)). A worker is matched with a firm vacancy through the matching function, which depends on labour market tightness $\theta = v/u$ (vacancies v over unemployment u). Workers and firms then split the surplus from the match, determined by worker bargaining power β , firm productivity p , and the cost of posting a vacancy c .

Setting reservation wage and discount rate to zero for simplicity, the standard wage curve is:

$$w = \beta(p + \theta c) \quad (\text{B10})$$

And the job creation equations is as follows, where s is the exogenous job separation rate and η parametrizes the matching function (e.g. $m = u^\eta v^{1-\eta}$):

$$\theta = \left(\frac{p - w}{sc}\right)^{1/\eta} \quad (\text{B11})$$

Finally, the Beveridge curve pins down the ratio of vacancies to unemployment:

$$v = \theta \frac{s}{s + \theta^{1-\eta}} \quad (\text{B12})$$

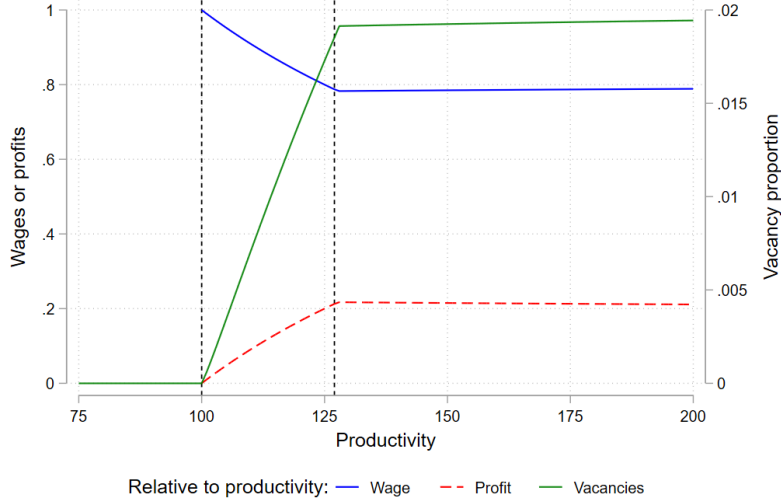
Since we are interested in models of firm heterogeneity, one can think of the relevant matching function as concerning the representative firm of a particular labour market (e.g. industry-region) with its own level of tightness.

To incorporate a minimum wage, we follow a similar procedure to above. [B10](#) becomes a max function between the optimal wage in this equation and the mandated minimum wage. Then while firm productivity is below the minimum wage, there are no vacancies as firms post no matches. While wages are constrained by the minimum wage (i.e. the optimal wage is below the minimum wage), [B11](#) applies with w as the minimum wage. The rest of the model is the same.

In the simulation, we use $\beta = 0.4$, $\eta = .9$, $c = 50$, $s = 0.25$, the minimum

wage is 100, and firm productivity p varies from 75 to 200 in increments of 1. The model is qualitatively similar across a range of these parameter values.

Figure B1: DMP, outcomes relative to productivity



Notes: The plot shows the DMP model with outcomes relative to productivity.

B.3 Baseline union bargaining model

We focus on the insider-outsider union bargaining model since this seems most relevant to our setting of rent-sharing with productivity. Following [Cahuc et al. \(2014\)](#), a stock of “insider” L_0 workers at a firm are represented by a union, whose objective function is simply the utility of the wage premium (wage above the minimum wage \bar{w} , denoted b). The firm may hire “outsider” workers, and since these workers have no bargaining power the firm pays these workers the minimum wage. Firm profit from total workers L_u is therefore:

$$\pi = R(L_u) - \bar{w}L_u - bL_0 \quad (\text{B13})$$

Given b , firms maximize L_u with respect to the minimum wage \bar{w} with $R'(L_u) = \bar{w}$. Firms and unions maximize the following Nash with respect to b , where bargaining power is given by β , L_u is given as above, and we impose a simple linear utility function for insiders equal to the wage premium:

$$\max_{\{b\}} (R(L_u) - \bar{w}L_u - bL_0)^{1-\beta} (b)^\beta \quad (\text{B14})$$

This gives the simple result that the wage premium is equal to the bargaining

power parameter times by the quasi-rents $(R - \bar{w}L_u)$, shared across insiders:

$$b = \beta(R - \bar{w}L_u)/L_0 \quad (\text{B15})$$

Profits are just the remainder portion of these rents:

$$\pi = R - \bar{w} - bL_0 = (1 - \beta)(R - \bar{w}L_u) \quad (\text{B16})$$

For simplicity, we use a Cobb-Douglas revenue function $R = AL^\alpha$. In the simulation, we use $\beta = 0.4$, $\alpha = .7$, a minimum wage of 1.5, and iterate the productivity shifter A between 1 and 1,000. We also take a simple view of the generating process of insider workers L_0 : we take a maximum to the size of the union (in simulations, 30 workers). One can of course imagine that L_0 grows over time with L_u , perhaps for workers who have been at the firm for a few years, and so one can view the simulated process as a local rent-sharing dynamic conditional on time.³¹ As above, the model is qualitatively similar across a range of these parameter values.

While the above focuses on the insider-outsider model, well-known alternative union bargaining models include the right-to-manage model and weakly efficient bargaining over wages and employment. However, these baseline models do not always generate rent-sharing as they require a non-homogenous production function such as the CES function (see also [Manning 1993](#)). [Cahuc et al. \(2014\)](#) notes regarding the right-to-manage model, and we confirm this through simulations and derivations, “If the revenue function of the firm is homogeneous of degree alpha (0,1), then [...] shocks to productivity or the firm’s selling price do not affect the wage and lead only to employment adjustments.” A similar condition holds for weakly efficient bargaining, for example with the revenue function $R = AL^\alpha$. The intuition for why this fails to deliver rent-sharing is that productivity shifts A do increase L , as well as w if *conditional on* L ; however, marginal productivity of labour decreases with greater L and it turns out the increase in w due to A is exactly offset by the decrease in marginal productivity through L .

³¹If L_0 is equal to $L_{u,t-1}$, then b is negligible as union welfare comes through expansions to the number of workers rather than the premium. This is very similar to the lack of rent-sharing in the right-to-manage model where the union bargains only over the wage which applies to all workers. In fact this is the case in this model when L_0 is less than its maximum, see figure 8.

Appendix C: Data and estimation sample

C.1 Data Access

The data used for this research was accessed from the NT-SDF. Access was provided under a non-disclosure agreement, and our output was checked so that the anonymity of no firm or individual would be compromised. Our results do not represent any official statistics (NT or SARS). Similarly, the views expressed in our research are not necessarily the views of the NT or SARS.

Data used: CIT-IPR5 panel (`citirp5_v5_0`); year-by-year IPR5 job-level data (`v5`), and year-by-year transaction-level Customs data (beta version 5.0). Date of first access for this project: 6 January 2023. Last accessed: 24 October 2025.

C.2 Estimation sample: Cross-sectional kink design

To implement the LMS internal instruments approach of Section 5 we must create a dataset of “stacked” events (or cohorts) from the panel of firms. This division of time periods into pre- and post-periods for specific firms is also useful for our cross-sectional kink design of Section 4 (the results of which are used in the within-firm shock approaches of Section 5), and so we use this basic data structure throughout the empirical analysis.

We start by dividing firms in the 2010-2019 period into four events (or cohorts), with the potential treatment date starting in 2014, 2015, 2016 and 2017 respectively. A firm may be in multiple events/cohorts, if it is observed in the panel for more than one of these treatment date starts, and is observed both prior to and after that treatment date. We call years prior to treatment start the pre-period, and including and after treatment start the post-period.

Separately for each cohort, we estimate productivity for each firm per Section 4 using only pre-period years.

Separately for each cohort (and each BC/SD), we then estimate the wage-kinks per the procedure in Section 4 using only post-period years. From this, we have a recentered productivity measure for each firm in each event. For our aggregate Tables and Figures in Section 4 we pool (or “stack”) across events/cohorts (accounting for the stacked design when clustering our standard errors).

Using this stacked approach for the cross-section kink design as well as the within-firm IV analysis is useful because 1) it makes it straightforward to define event-specific recentered productivity values (and therefore constrained versus unconstrained firms) which we need for our internal IV analysis and 2) it provides a

natural way to estimate productivity in a pre-period, so it is a fixed heterogeneity category for the cross-sectional and shock analyses in a post-period, without simply bifurcating the original panel and losing many observations 3) it allows a firm to have *some* time-varying productivity, so that a firm could be low productivity and constrained in earlier events and higher-productivity and unconstrained in later events. We do check that the results of Section 4 do not depend on this pooling across events.

C.3 Estimation sample: Internal instrument approach

The sample restrictions and key variable restrictions are an important feature of the LMS strategy, and we follow them closely, constructing our sample as follows:

1. Identify “stayers” in the worker-level data who remain employed at the same firm for 8 consecutive years, separately defining stayers for the tax year event periods 2010-2017, 2011-2018, 2012-2019, and 2013-2020. These are the cohorts/events mentioned immediately above, with 2010-2020 covering the usable period of the employment data. Drop stayers’ records in the first and last years of this tenure (when they may have entered or separated), and only keep workers who are full-time employed over this 6 year period at their firm. Count the number of stayers in each firm for each event period and create year-specific firm-level statistics for stayers’ wages (specifically median wage).
2. For each event, only keep firms which have at least 2 stayers. LMS use a 10-stayers minimum as their baseline, but this is overly restrictive when it comes to the South African firm-size distribution and a labour market context defined by high churn. In our baseline specification we use a 2-stayer minimum, to mitigate measurement error in one-stayer firms where the one stayer may be an owner or otherwise unrepresentative of employer/employee dynamics. In Appendix Figure A8 we show that our results are not sensitive to the number of stayers.
3. Over the 6-year period for each event, we treat the first three periods as the pre-period and the latter three as the post. Treatment is defined as an above-median increase in firm value-added between periods -1 and 0 for each event, where the median increase is weighted by firm size. Events are stacked (Cengiz et al. 2019). Period -2 is used as the omitted reference period to allow for some mean reversion dynamics in period -1, as in LMS. For the

same reason, periods 1 and 2 are considered the post periods of interest, rather than period 0 (results are essentially unchanged if we use only period 2 instead). Period -3 is used to assess pre-period parallel trend violations.

When estimating Equation 7 by productivity bin, and for constrained versus unconstrained firms (e.g. Table 2 and Figures 5, 6 and A5) we make a few small data construction/visualization decisions for the internal IV specification:

1. As discussed in Section 5 above, there is measurement and estimation error in the \hat{v}^* threshold, and so we drop firms very close to the threshold before constructing bins. In our baseline specification we drop the 5% of constrained firms with the highest recentered productivity values, and the 5% of unconstrained firms with the lowest such values. Our results are not sensitive to this trimming procedure; see Appendix Figure A10.
2. We divide firms into 10 approximately equally-sized bins to get a sense of the shape of response along the distribution. We require that there be at least 2 bins on either side of \hat{v}^* , as the value of the bins is in seeing the shape of the marginal response against productivity, which means in practice that the bins in the constrained region are smaller than those in the unconstrained region; 15% of firms (or 12% of event-specific observations) fall in the constrained region in this baseline specification.
3. In Figure 5 and Panel (b) of A5, treatment effects are the average response across post-periods 1 and 2; results are essentially unchanged if we only use post-period 2.

C.4 Estimation sample: External instrument approach

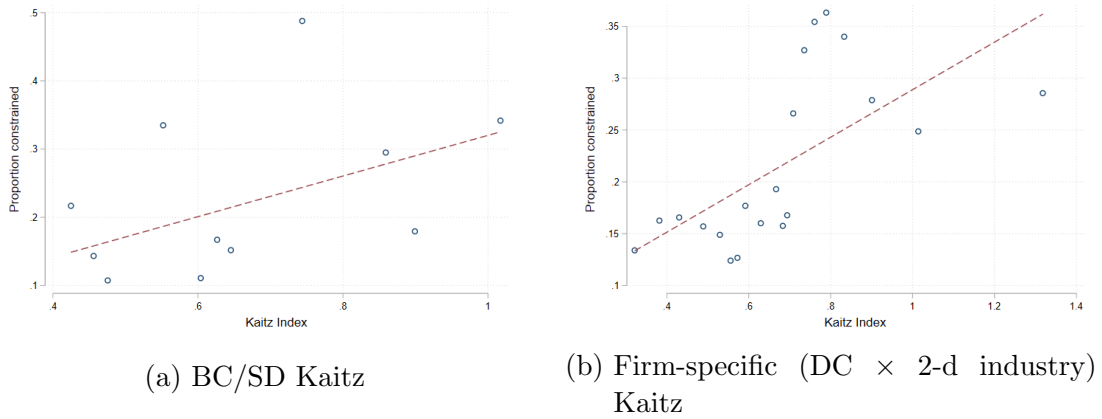
Due to the power issues discussed in Section 5.1.2, for our external instruments approach we do not use the event structure of the dataset for identification in our main results, but simply pool across firms in the internal IV estimation sample. This means we do impose restrictions such as requiring a balanced panel with a two-stayer minimum, but do not separate into pre- and post-periods, except for the supplementary exercise mentioned in footnote 21 and shown in Appendix Table A5.

Appendix D: Minimum wage bindingness

D.1 Kaitz Index heterogeneity

For each worker, we give the ratio between the minimum wage from their BC/SD (Bargaining Council or Sectoral Determination) and one of two measures of local wages: either the median annualized wage in the BC/SD (which we refer to as the BC/SD Kaitz), or the median annualized wage in the detailed geography (DC) by 2-digit industry intersection (referred to as the firm-specific Kaitz).³²

Figure D1: Probability of being in constrained region, by Kaitz Index



Notes: Figure shows bin-scatter plots of firm Kaitz Index measures (x-axis) against firm constrained status (y-axis; a dummy variable equal to 1 for constrained firms).

Firms with higher Kaitz Indices are more likely to be found in the constrained region, which is consistent with the estimated wage kink identifying the threshold below which firms are bound by the minimum wage. Figure D1 shows clear positive relationships between both Kaitz measures and the probability of a firm being found in the constrained region. Table D1 shows these relationships are statistically significant, for both the underlying continuous Kaitz Index and a binary transformation which equals 1 if the Kaitz Index is above its median value.

D.2 Empirical strategy using the minimum wage bite

One could imagine doing the empirical tests of Sections 4 and 5 with the minimum wage bite as the “treatment” or x-variable instead of firm productivity. However, this would not be appropriate. Firstly, in the model, firm productivity determines

³²The Kaitz index of about 0.68 is higher than in Figure A12, about 0.78 for South Africa using ILO and World Development data, because the median wage in our data is higher given we only observe formal firms.

Table D1: Constrained probability by Kaitz Index

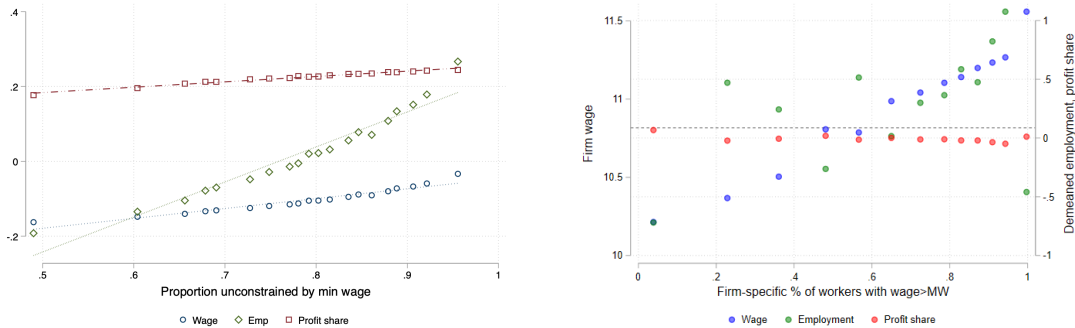
	BC/SD		Firm-specific	
Kaitz coefficient	.164 (.0696)	.298 (.1534)	.123 (.0419)	.229 (.0815)
N	324608	324608	324608	324608
Binary	Y		Y	
Continuous		Y		Y

Notes: Table shows results of regressions of each firm's constrained status (a dummy variable equal to 1 for constrained firms) on its Kaitz Index, for both the underlying continuous Kaitz Index and a binary transformation. Standard errors are shown in parentheses.

the optimal wage and therefore employment and profit; a wage set at the minimum wage is an *outcome* of this process. Secondly, the corresponding empirical test is unclear. In the simulated model, if we plot the wage on the x-axis, all minimum wage bound firms will simply be at one dot corresponding to the minimum wage with the average corresponding level of employment and profits. This would hide precisely the variation in employment and profit that we are studying.

We could instead plot minimum wage bite by groups of firms; figure D2 panel (a) illustrates the results using the simulated model. The figure shows upwards sloping lines for wages and employment, with a slight positive slope on profit share too. Panel (b) shows the results from our observed data, with similar though noisy patterns. However, correlated variables could explain such positive slopes, such as worker quality. Further, recall that we only observe an approximate minimum wage per firm, meaning the minimum wage bite is measured with substantial error. It is therefore our view that this is not a good empirical strategy. Our kink tests are both clearer and more demanding, and so are much preferred.

Figure D2: Results using the minimum wage bite



(a) Simulation from theoretical model

(b) Results from empirical data

Notes: The x-axis is the proportion of firms paying at or below the minimum wage. Panel (a) uses the simulated model, and Panel (b) uses the empirical data.

Appendix E: Structural model

E.1 Estimation details

We use the same model as in Appendix section B.1, though abstracting from firm amenities b_i . The v^* threshold is set at the empirically detected kink-point, as explained in section 4.2, with productivity a_i normalized to zero at the kink-point.

The marginal revenue product of labour is set to equal the marginal cost of labour, incorporating a worker-firm productivity covariance ρ following the discussion in section 6.1. Table E1 shows the wage and employment model equations, for the unconstrained and constrained regions. The firm rent-sharing elasticity is $\partial w_i / \partial a_i = \varepsilon / (\varepsilon + \eta)$, with ρ representing the increase in wages associated with worker productivity. The employment intercept includes market-level factors \bar{n} influencing average employment (as in B.1).

Table E1: Structural model equations by firm outcome and region

Firm outcome	Unconstrained: $a_i \geq 0$	Constrained: $a_i < 0$
Wage	$w_i = w_{\min} + \left(\frac{\varepsilon}{\varepsilon + \eta} + \rho \right) a_i$	$w_i = w_{\min}$
Employment	$n_i = \bar{n} + \frac{\ln(1 + \varepsilon) + w_{\min}}{\eta} + \frac{1}{\varepsilon + \eta} a_i$	$n_i = \bar{n} + \frac{\ln(1 + \varepsilon) + w_{\min}}{\eta}$

The profit share equation requires additional assumptions. Value added is specified as a simple Cobb-Douglas function with only labour n_i and firm productivity shifter a_i . Then the profit share is one minus the wagebill share of value added, with a constant \bar{p} reflecting market factors, a max-function reflecting a positive minimum, and p_{slope} allowing for attenuation between the model-implied and observed profit shares as discussed in the text, for example due to rental payments to capital or measurement error.

$$va_i = \exp(a_i + (1 - \eta) \cdot n_i), \quad (C1)$$

$$profshare_i = \bar{p} + p_{slope} \cdot \log \left(\max \left(1 - \frac{\exp(n_i + w_i)}{va_i}, 0.01 \right) \right). \quad (C2)$$

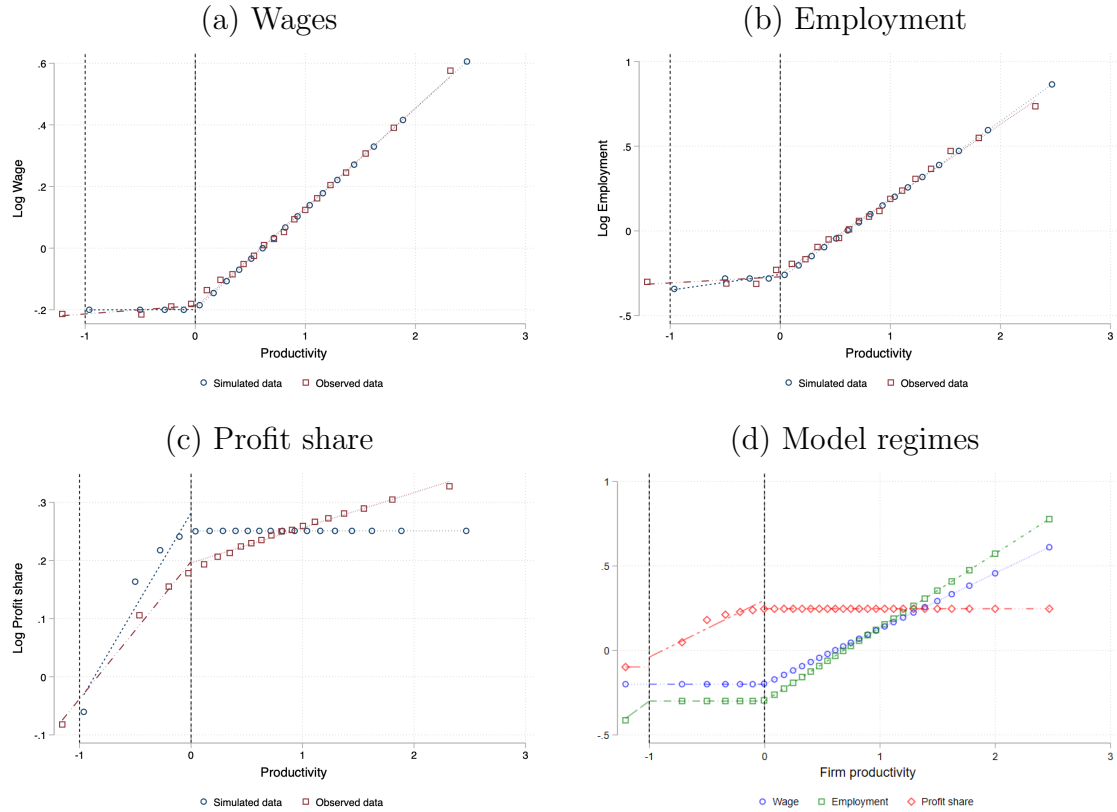
In terms of the actual estimation, see Table E2. We use the same datapoints presented in figure 4 panel (b), i.e. the residualized firm wage, employment and profit share by estimated firm productivity, averaged over 20 quantiles. The grid of parameter values for the structural parameters is scanned, in each case computing the model-implied outcomes given the inputs of observed productivity a_i and

Table E2: Parameter estimates from structural model

Parameter	Symbol	Source of identification	Value
<i>Panel A: External to model (estimated elsewhere in the paper)</i>			
Firm labour supply elasticity	ε	Reduced form IV estimate (see table 2)	3.4
Minimum wage	w_{min}	Cross-sectional kink design (see figure 1)	-0.2
<i>Panel B: Structural parameters estimated via Maximum Likelihood</i>			
Inv. elasticity of lab. demand	η	Slope of employment curve	1.86
Worker-firm productivity covar.	ρ	Ratio of employment to wage slopes	0.19
<i>Panel C: Nuisance parameters estimated via Maximum Likelihood</i>			
Intercept for employment	\bar{n}	Intercept of employment curve	-.25
Intercept for profit-share	\bar{p}	Right intercept of profit-share curve	0.31
Slope for profit-share	p_{slope}	Left slope of profit-share curve	.089

Notes: See text for model details. Figure E1 shows the fit corresponding to the parameter values, and Table E1 gives the model equations.

Figure E1: Structural estimation: Fit of outcomes in model vs. data



Notes: Panels (a) to (c) show the simulated data (blue) based on the optimal parameters against the observed data (red). Panel (d) shows the fit of all three together, replicating the main paper figure 4 in the simulation. The right vertical line represents the threshold for constrained firms, with the left vertical line demarcating the demand-constrained firms.

$\varepsilon = 3.4$, and the values selected that minimize the sum of the squared distances between the model and the observed quantiles. The demand constrained region is then set at the left tail of the observed distribution, on the assumption given the observed patterns that this region is small in the data. The optimal values are $\rho = 0.2$ and $\eta = 2$; figure E1 shows the fit.

E.2 Mixture model

This subsection presents further details on the mixture models discussed in section 6.2. Table E3 gives the estimated equations. Beginning with the case of noncompliance, assume that there is a share α of firms that are compliant. Then, using the within-firm shock results from table 2 (below Table E3 column 3), the first-differenced model equations from Table E1 imply the equations in Table E3 column 4, noting that $1 - \alpha$ firms will not be constrained at low productivities given they are noncompliant.

The optimal parameters are as in the text, with the firm labour supply elasticity $1/\varepsilon = 3.3$, and the inverse elasticity of labour demand $\eta = 1.33$. The case of firm productivity mis-classification can be modelled similarly, with $1 - \alpha$ firms in the constrained regions that are in reality unconstrained.

Table E3: Mixture model equations for estimating structural parameters

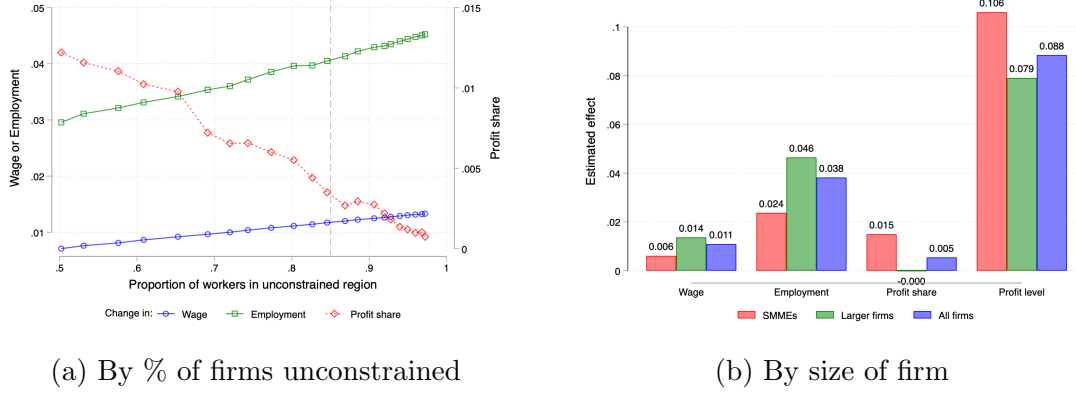
Outcome	Region	Estimate	Mixture model	
			Non-compliance	Union
Wage	Unconstrained	0.19	$\frac{\varepsilon}{\varepsilon + \eta}$	$\alpha \frac{\varepsilon}{\varepsilon + \eta} + (1 - \alpha)\beta$
Wage	Constrained	0.13	$(1 - \alpha) \frac{\varepsilon}{\varepsilon + \eta}$	$(1 - \alpha)\beta$
Employment	Unconstrained	0.61	$\frac{1}{\varepsilon + \eta}$	$\frac{\alpha}{\varepsilon + \eta} + \frac{1 - \alpha}{\eta}$
Employment	Constrained	0.45	$(1 - \alpha) \frac{1}{\varepsilon + \eta}$	$\frac{1 - \alpha}{\eta}$

In section 7, we discuss a model of the labour market where there is a mixture of monopsonistic and non-monopsonistic firms. Table E3 column 5 provides the equations for a simple example using the union model (Appendix section B.3), keeping α as the share of monopsonistic firms, the same value-added function above, and bargaining power β . The optimal parameters are a monopsony share $\alpha = .31$, union power $\beta = 0.19$, firm labour supply elasticity $1/\varepsilon = 2.6$, and inverse elasticity of labour demand $\eta = 1.5$.

E.3 Policy simulations

Productivity shocks. Figure E2 simulates a partial equilibrium increase in firm productivity, which may exemplify policymakers’ industrial policy objectives. Panel (a) considers how effects vary when different proportions of workers are in the unconstrained region. Panel (b) shows outcomes separately for smaller (SMMEs) versus larger firms.

Figure E2: Simulations based on structural estimation: Productivity shocks



Notes: See Section 6. Both panels show the effects of increasing firm productivity by 10%. Panel (a) varies the baseline simulation by minimum wage levels, generating different proportions of workers in unconstrained regions. Panel (b) uses the estimated model as the baseline, with outcomes shown separately for smaller (SMMEs) versus larger firms.

Employment Tax Incentive. Our ETI policy simulation gives a subsidy of 25% of the worker’s wage to the firm for every worker paid less than ZAR5,000. The wage, employment and profit share equations E1 are therefore modified to reflect this lower marginal cost of labour per worker to the firm, i.e. $n_i = 1/(\varepsilon + \eta) * (a_i - \ln(1 + \varepsilon) - \ln(.75))$ and $w = \varepsilon \cdot emp_i$ as before. The effective minimum wage cost becomes $w_{\min} + \ln(.75)$.

South Africa’s actual ETI subsidy phases out from 50% of the wage at ZAR 2,000 to 0% at ZAR 6,000; our simplified cutoff of ZAR 5,000 with a uniform 25% subsidy is an approximation. Additionally, only hires are eligible and the subsidy lasts for 2 years, meaning our simulation may better reflect hiring decisions in a dynamic setting. Finally, other explanations for the limited employment effects found in existing studies include the risks firms face if subsidised workers are difficult to dismiss after the subsidy ends, and principal-agent frictions in local managerial decisions (Budlender and Ebrahim, 2021; Ranchhod and Finn, 2016), or inelastic labour demand — though our reduced-form estimates (Table 2) indicate that firm employment is somewhat responsive to shocks.