

Immigration and Worker Responses Across Firms: Evidence from Administrative Records in Colombia*

Lukas Delgado-Prieto[†]

Job Market Paper

October, 2023

[\[Click here for latest version\]](#)

Abstract

Many migrants are moving to developing countries where small firms are prevalent in the labor market. The interaction between firms, workers, and labor supply shocks in these contexts is mainly unknown. To address this, I exploit the mass arrival of migrants from Venezuela in Colombia and use administrative records covering the universe of formal workers and firms. As immigrants concentrate in the informal sector, I find a reduction in formal employment for natives earning the minimum wage, explained by their high substitutability with informal workers who become less costly. Across firms, I find negative formal employment and wage effects for natives in small firms. To rationalize this, I construct a model of heterogeneous firms that hire formal and informal labor to show that the response to immigration is more pronounced in small firms as they hire relatively more informal labor. Finally, using causal forests, I show that firm pay premiums explain more the heterogeneity of employment and wage effects than other worker characteristics. Overall, these results suggest that firms play an influential role in determining the impact of immigration on workers' outcomes.

Keywords: Immigration, Formal labor markets, Causal forest.

JEL Codes: F22, O15, O17, R23.

*I am extremely grateful to Jan Stuhler for his guidance and support in this project. I am also thankful to Michael Amior, Jaime Arellano-Bover, Agostina Brinatti, David Card, Juanjo Dolado, Juan Carlos Escanciano, Jesús Fernández-Huertas, Patrick Kline, Luigi Minale, Jesse Rothstein, and Jon Piqueras for their comments and suggestions. In general, I want to thank participants at the EALE Conference, ENTER Jamboree, the 1st HUMANS LACEA Seminar, the Labor Lunch Seminar at UC Berkeley, the IRLE Seminar, the UC3M Applied Reading Group, the 23rd IZA Summer School in Labor Economics, and the Junior Seminar of the Economics of Migration for the helpful discussions. I am indebted to the Ministry of Health and Social Protection of Colombia for providing access to their data. This paper was previously circulated as: "The Role of Workers and Firms in the Impact of Immigration".

[†]Department of Economics, Universidad Carlos III de Madrid, Spain. Email: ludelgad@eco.uc3m.es.

1 Introduction

In the last decade, several countries across the globe have experienced large population exodus.¹ The majority of these migrants are leaving for close destinations that, in most cases, are developing countries where small firms are prevalent (McKenzie, 2017). The labor market in these countries is also characterized by the links between the formal and informal sectors, most often through these small firms (Ulyssea, 2018). As migrants can disproportionately concentrate in small firms (Delgado-Prieto, 2022), it becomes more relevant to study the labor market impacts of immigration exploiting the firm dimension in these contexts. However, the migration literature mostly neglects firms when analyzing immigration effects, partly due to data limitations. This paper shows that firms are an important aspect in order to understand how native workers and the labor market overall adjust to labor supply shocks.

To do so, I study the labor market impacts of one of the most significant episodes of immigration in recent history, the Venezuelan mass migration to Colombia, and use novel administrative data that covers the universe of formal workers and firms in the country.² Exploiting the unequal arrival of migrants across local labor markets with the panel structure of the data, I quantify worker-level impacts across different worker and firm characteristics and construct a proxy for the role of firms in the immigration effects.³ A growing number of studies analyze how firms shape, for instance, wage inequality (Card et al., 2013) or immigrant assimilation (Arellano-Bover and San, 2020), but less is known about how firms determine natives' adjustments to labor supply shocks.

To my knowledge, this is one of the first papers to study the impact of immigration in developing countries equipped with matched employee-employer data.⁴ This helps to uncover important sources of heterogeneity that have not been explored previously, as prior findings primarily focused on the effects across worker characteristics. Although worker and firm characteristics are related (e.g., minimum wage workers often work in smaller firms), most of the heterogeneity in migration

¹For example, Afghanistan, Ukraine, Syria, and Venezuela, among others.

²Throughout this paper, formal workers refer to workers who contribute monthly to the health system in Colombia.

³By following workers over time, I can address more carefully the compositional changes in the employed population after immigrants arrive, which are typically aggregated when constructing regional outcomes with cross-sectional data. For instance, recent papers emphasize that when a specific set of workers move out of employment or to other regions, the wage estimates are not properly identified (Borjas and Edo, 2021; Dustmann et al., 2023).

⁴In developed countries, Bratsberg and Raaum (2012) for Norway, Foged and Peri (2016) for Denmark, Dustmann et al. (2017) for Germany and Orefice and Peri (2020) for France have recently studied immigration shocks exploiting administrative data.

impacts in this setting comes from the firm side. The close interaction between the formal and informal sectors in small and less productive firms provides insights into this heterogeneity. As the labor supply shock is clustered in the informal sector, it reduces informal wages, which in turn decreases labor demand in the formal sector, but mostly in small firms that can substitute formal labor for informal labor more easily (Delgado-Prieto, 2022).⁵ The implication of this is that the wage and employment effects vary massively across the firm size distribution, which motivates the worker-level analysis exploiting the firm dimension done in this paper.

My empirical strategy compares similar workers in areas with different exposure to migration over time. Because migrants endogenously sort into areas that offer the best economic opportunities for them, I use two distinct instruments: past settlements of Venezuelans and distance to the nearest crossing bridge with Venezuela. I exploit these instruments in a differences-in-differences research design (DiD-IV) to find a persistent negative impact on individual formal employment and wages for natives.

The negative employment effect is driven by workers earning the minimum wage before the immigration shock. In this context, the relatively high and binding minimum wage for many formal workers limits the space for downward wage adjustments (around 40% of all formal workers in Colombia earned the minimum wage in 2015) and increases their chances of job displacement. Regarding wages, the negative wage effect mainly affects native workers earning above the minimum wage and working in the smallest firms. The fact that migrants tend to concentrate in small firms, which are more constrained by the minimum wage and employ a higher share of informal workers, helps to explain this finding. Regarding regional mobility, results suggest that formal workers are not systematically moving to other places after immigrants arrive.

To rationalize how firms interact with immigration effects, I construct an imperfectly competitive labor market model based on Card et al. (2018) with heterogeneous firms, but adapting the labor inputs cost similarly to Ulyssea (2018) and allowing for imperfect substitution of labor inputs as in Delgado-Prieto (2022). The model shows that the aggregate substitutability between formal and informal workers must be high for a negative formal employment and wage response. The model also establishes firm-level reactions to immigration depending on their initial share of

⁵This is an effective application of the first Hicks-Marshall rules of Derived Demand: “The demand for anything is likely to be more elastic, the more readily substitutes for the thing can be obtained” (Hicks, 1932).

informal work in production. Hence, the model predicts smaller firms have elasticities of formal employment and formal wages with respect to informal labor that are higher, in absolute terms, relative to larger firms. Notably, these patterns are matched in the empirical findings, that is, workers in the smallest firms are more affected in terms of formal employment and formal wages than workers in larger firms. Consistent with a higher substitutability of formal and informal labor in small firms.

Next, I exploit the canonical [Abowd et al. \(1999\)](#) (AKM hereafter) framework to recover firm fixed effects (FEs), or firm-specific pay premiums, and worker FEs. A significant contribution of this paper is to understand the sources of wage and employment losses stemming from immigration using these constructed variables. I find that workers in middle-paying firms during the pre-shock period suffer the largest wage losses compared to workers in low-paying firms. A binding minimum wage in low-paying firms prevents wage losses, while for the rest of the firms, a potential explanation is the reallocation effects of immigration, that is, workers might be moving from middle- to low-paying firms.⁶ However, I find no differential sorting of native workers after immigrants arrive.

Regarding employment, the finding is the opposite: native workers in low-paying firms are more affected than workers in high-paying firms. I find a similar picture when dividing between low- to high-wage workers. Next, I analyze other firm outcomes changing in response to the immigration shock. Particularly, I show that firms opt-out from the formal sector for new hires, especially those that poach less from other formal firms, and that firm exit is higher in places that receive more migrants.

In the second part of this paper, I estimate the heterogeneity of treatment effects according to a vector of worker and firm characteristics following the recent literature in machine learning ([Athey and Imbens, 2016](#); [Athey et al., 2019](#)). Specifically, I implement different causal forests that quantify a set of reduced-form estimates from random subsamples to determine those variables that explain most of the heterogeneity of immigration effects. From this algorithm, first, I identify the subgroups most affected by immigration, both for employment and wages. Then, based on the frequency that these variables appear in the splits of all the decision trees, I construct a simple measure of variable importance. In this exercise, I consistently find that firm-specific pay premiums

⁶Another explanation, according to the model in [Amior and Stuhler \(2022\)](#), is that when the share of firms in the low-pay sector grows due to immigrants, firms in the high-pay sector increase their monopsony power and reduce workers' wages.

are ranked higher in the algorithm, meaning they are more likely to explain employment and wage losses than worker characteristics (i.e., job tenure, age, sex, and wages) in the pre-shock period. Therefore, firms' role in the impact of immigration appears to be very relevant, which is one of the main findings of this paper. These results suggest direct policy implications. For instance, to mitigate the employment and wage losses in the smallest and least productive firms, initiatives like skill development and training programs for their workforce, in conjunction with tax incentives or accessible financing, can be introduced to enhance their overall productivity.

Literature. This paper contributes to different strands of the labor economics literature. First, I contribute to the literature that analyzes how firms shape native and migrant outcomes. Several papers emphasize that firms affect workers' outcomes through different channels. For instance, [Arellano-Bover and San \(2020\)](#) and [Dostie et al. \(2021\)](#) find that firm-specific pay premiums explain around one-fifth of the immigrant-native earnings gap in Israel and Canada, respectively; [Doran et al. \(2022\)](#) uses the H-1B visa lottery to find that winner firms crowd out their workers for H-1B visa workers; and [Orefice and Peri \(2020\)](#) shows an increase in assortative matching after immigrants arrive in France, with high-quality firms attracting high-quality workers. In this respect, my paper highlights that the effects of immigration are concentrated on natives working in small firms, as in [Amior and Stuhler \(2022\)](#). However, the mechanisms that lead to these implications differ substantially. In my paper, this implication derives from the particular setting and the way small firms can substitute formal for low-priced informal labor after immigrants arrive.⁷

Next, I contribute to the literature that estimates the individual impact of immigration. Two of the first papers that estimated worker-level effects of immigration were [Bratsberg and Raaum \(2012\)](#) and [Foged and Peri \(2016\)](#).⁸ With the universe of formal workers, I integrate into the analysis all the movements of natives between areas, reducing the attenuation of the wage estimates discussed in [Borjas \(2006\)](#), and exclude all inflows into employment that are part of regional-level responses, as documented in [Dustmann et al. \(2023\)](#). Using worker-level responses, I can identify the main drivers behind labor market adjustments to immigration in two novel ways.⁹ First, building a model of

⁷In contrast, the argument in [Amior and Stuhler \(2022\)](#) is that as small firms pay worse wages, they have greater incentives to hire migrants with lower reservation wages, affecting more native workers in small firms.

⁸The first paper uses licensing requirements in the Norwegian construction sector to leverage exogenous immigration shares. These authors find that native wages in this sector decrease as immigrant shares increase, with low-paid native workers leaving this sector more frequently. The second paper exploits a refugee dispersal policy in Denmark to find that low-skilled natives pursue less manual-intensive occupations, upgrading their wages.

⁹Other recent papers that quantify worker-level effects are [Hoen \(2020\)](#) for Norway, [Ortega and Verdugo \(2022\)](#)

immigration with heterogeneous firms (a dimension typically ignored in the migration literature), and second, exploiting a machine learning algorithm that estimates heterogeneous immigration effects in a data-driven way.

I also contribute to the literature on how workers react to other types of labor market shocks (see related papers of [Autor et al. \(2014\)](#) for industry shocks to import competition, [Yagan \(2019\)](#) and [Redondo \(2022\)](#) for local employment shocks, and [Gulyas et al. \(2019\)](#) for mass layoffs). In this respect, I show that firms play a significant role in determining the wage and employment losses to labor supply shocks, both theoretically and empirically. This important result speaks directly to the labor literature that focuses mostly on workers or industry characteristics to understand the sources of adjustments to these shocks. My emphasis on the importance of firm heterogeneity is consistent with [Gulyas et al. \(2019\)](#), which finds that after job displacement, earning losses are higher for workers employed in high-paying firms using causal forests.¹⁰

Lastly, this paper contributes to the literature on the impact of immigration in developing countries (see related papers of [Morales-Zurita et al. \(2020\)](#), [Caruso et al. \(2021\)](#), [Lebow \(2021\)](#), and [Delgado-Prieto \(2022\)](#) for Colombia; [Del Carpio and Wagner \(2015\)](#), [Ceritoglu et al. \(2017\)](#), and [Aksu et al. \(2018\)](#) for Turkey; and [Groeger et al. \(2022\)](#) for Peru). Since having administrative data in developing countries for workers and firms is unusual, most previous studies used cross-sectional surveys to determine the regional impact of immigration. However, [Dustmann et al. \(2023\)](#) discusses how regional and individual effects are conceptually different types of labor market responses to immigration. This helps to explain why my findings relative to other studies that analyze the impact of immigration in the Colombian setting vary. Therefore, with panel administrative data, I quantify for the first time in a developing country individual employment and wage effects of immigration.

The rest of the paper is structured as follows. Section 2 discusses the characteristics of the labor supply shock derived from the Venezuelan crisis and describes the data. Section 3 describes the empirical strategy and the identification assumptions. Section 4 reports results at the worker level by different individual characteristics. Section 5 introduces a model with heterogeneous firms

for France and [Kuosmanen and Meriläinen \(2022\)](#) for Finland.

¹⁰[Yakymovych et al. \(2022\)](#) also uses causal forests, with administrative data from Sweden, to identify sets of workers more vulnerable to job displacement and to uncover the subgroups with the most significant earnings losses after displacement. They find that older, less-educated, and manufacturing workers are the most affected.

and shows results by firm characteristics. Section 6 introduces the machine learning approach and discusses the main findings. Section 7 provides robustness tests. Finally, Section 8 concludes.

2 Institutional Context and Data

2.1 The Venezuelan Mass Migration

Historically, Colombia and Venezuela shared an extensive territorial border characterized by a dynamic relationship with frequent economic and cultural interactions. People often moved back and forth between the two countries, but frequently, Colombians settled in Venezuela. This trend intensified after 1950, fueled by the oil boom in Venezuela and the internal conflict in Colombia. The economic opportunities presented by Venezuela’s oil industry attracted many Colombians to emigrate, seeking better livelihoods and prospects for their families. Recently, the trend reversed with Venezuela’s unprecedented socio-economic and political deterioration that triggered massive outflows of people since 2013, both voluntarily and forcedly. As a result, several countries in Latin America are receiving vast numbers of migrants, especially Colombia, Perú, and Ecuador (UNHCR, 2019). By far, Colombia has been the major receiver country with more than 1 million working-age Venezuelans (4.1% of the working-age population living in Colombia) as of 2019 (DANE, 2019).

The Venezuelan exodus is unmatched in the recent migration history in Latin America. Worldwide, there are only two contexts with similar figures, namely, the Syrian and the Ukrainian exodus. In the first case, Turkey has been the major receiver country of Syrians, with various papers analyzing this labor supply shock (e.g., Del Carpio and Wagner (2015); Ceritoglu et al. (2017); Aksu et al. (2018)). However, the Colombian context is different from the Turkish one. First, Venezuelans speak the same language as Colombians and second, Colombia’s government has implemented an open border policy in which all Venezuelan immigrants can have a work permit. In particular, after 2017, all undocumented Venezuelans in Colombia have access to the Special Permit of Permanence (PEP, by its acronym in Spanish). This allows them to work for a specific period, provides access to basic services, and facilitates their integration into Colombian society.¹¹ Yet, around 90% of Venezuelan immigrants were employed in the informal sector in 2019 –meaning

¹¹To overcome the limitations of the PEP, the government enacted in 2021 a Temporary Protection Statute for Migrants (ETPV, by its acronym in Spanish) that grants up to ten years of regularization for Venezuelan immigrants.

they do not contribute to the social security system— and were concentrated at the bottom of the native wage distribution (Delgado-Prieto, 2022). This fact relates to the occupational downgrading of Venezuelans since they have similar average levels of education compared to their Colombian counterparts and are even more educated in the latest years of arrival.

As described above, the labor supply shock in Colombia occurs mainly in the informal sector, so why focus on the formal sector in this paper? First, Delgado-Prieto (2022) shows that there is a strong negative effect on the wages of the informal sector following the arrival of migrants, and as firms can combine formal and informal labor in production (especially the smallest firms), they will substitute formal for informal employment in response to lower informal wages when both types of labor have a high substitutability. So, formal employment is primarily affected in response to the labor supply shock, even if immigrants mostly work in the informal sector.¹² For that reason, focusing on formal workers’ adjustments across the firm size distribution is of central interest.

2.2 Data

The administrative data source for Colombia is the *Planilla Integrada de Liquidación de Aportes* (PILA), which contains administrative records from the Colombian social security system managed by the Health Ministry (*Ministerio de Salud y Protección Social*). PILA contains information on the universe of formal workers in tax-registered firms. It excludes informal workers and informal firms but includes self-employed formal workers. The PILA is based on the monthly contribution of the worker, according to their reported base income, to the health system in Colombia. Each observation in PILA is a worker-employer match for a given year and month. The dataset contains worker-level information on labor income, sex, age, job type (employee or self-employed), foreign status, municipality, and the firm identifier for each job. I have access to PILA from 2010 to 2019 for the month of August.¹³ In addition, I use the most recent Colombian census (CNPV, by its acronym in Spanish), recollected between January and October of 2018, to construct the immigration shock. The census provides the most reliable source of information on Venezuelan

¹²The asymmetric employment and wage responses across the informal and formal sectors in the face of labor supply shocks have been analyzed in other contexts too (see Corbi et al. (2021) for Brazil and Kleemans and Magruder (2018) for Indonesia).

¹³I choose August to exclude the seasonal characteristics of other months (e.g., December-January or March-April) and because the census recollection ended in October of 2018, omitting arrivals of migrants in November and December of that year.

immigrants.¹⁴

For the analysis, I built a dataset with all the individuals that appeared in PILA between 2012 to 2019 in the rows and their yearly variables on the columns.¹⁵ The total number of workers in this dataset, who appeared in at least one of the eight years, is 18,430,987. Next, I restrict to full-time native workers between 25 and 55 years old in 2015 and assign the immigration shock to all these workers according to their 2015 location, which leaves 7,123,223 workers.¹⁶ Then, I transform the municipality variable into a more standard definition of local labor markets or commuting zones, adjusting the methodology of [Sanchez-Serra \(2016\)](#).¹⁷ This adjusted definition yields 109 functional urban areas (FUAs) after eliminating small or rural municipalities, with a sample of 6,706,035 workers.¹⁸ This is the sample I use for the employment analysis over time (a balanced panel) as the worker can be employed or not in the comparison year. For the wage analysis, I further restrict to workers with 30 days of employment in the month and positive wages; moreover, the worker must be employed in the post-treatment year of comparison. Thus, the sample varies slightly by year (an unbalanced panel). It is worth noting that all the restrictions to the dataset are the typical ones implemented in the literature ([Gulyas et al., 2019](#); [Yagan, 2019](#)).

Descriptive Statistics for Formal Workers. Table 1a, 1c and 1b shows descriptive statistics for natives, foreigners, and Venezuelans with PEP by age, sex, and wages across time.¹⁹ In terms of observable characteristics, Venezuelans with PEP in the formal sector are younger, predominantly

¹⁴The labor force survey (GEIH, by its acronym in Spanish) also measures the number of Venezuelan immigrants in Colombia at a higher frequency but not at the local level I exploit. Furthermore, [Aydemir and Borjas \(2011\)](#) document that surveys can attenuate immigration effects due to measurement error of migrants.

¹⁵The administrative records of PILA are constructed at the level of the contribution, as workers with more than one labor contract must pay contributions for each one of them. To transform to a worker-level dataset, first, I drop the contributions to the health system with type N, which are the ones that present corrections to their base income or changes to their labor status. Second, I aggregate the income for workers with more than one labor contract and leave the characteristics only for the job with the highest base income by the worker.

¹⁶Selecting only the workers observed in the base period rules out inflows of workers in the post-treatment period from the analysis. Also, in 2015, part-time workers in PILA were less than 0.3% of all workers.

¹⁷A shortcoming with the municipality variable in PILA is that certain firms with several establishments across the country report the information for all employed workers in the municipality where the biggest establishment is located, such that the observed employment in smaller cities is understated.

¹⁸The definition of FUAs consists of the 53 most extensive urban areas in the country defined from population grid data, municipal boundaries, inter-municipal commuting flows, plus 56 municipalities with more than 2,500 formal workers according to the restricted sample in 2015. I exclude San Andrés, Cumaribo, Leticia, and San José del Guaviare from this definition as they belong to the islands or Amazonia. In Appendix Table G.1, I show the sample distribution by FUAs; using this definition, only 5.9% of workers are excluded.

¹⁹To identify foreign status in PILA, I exploit the type of document workers have in their health contribution. If workers have a national ID, they are defined as natives, whereas if their document refers to the Special Permit of Permanence (PEP, by its acronym in Spanish), they are defined as Venezuelan migrants. Since the PEP's program started around 2018 to foster the regularization of Venezuelan migrants, it was not possible to identify these migrants before that year. Last, if their document is a foreigner ID or a passport, they are defined as foreigners.

male, and earn lower wages compared to natives and other foreigners (see Table 1a, 1c and 1b). In fact, the group of foreigners earns substantially higher wages than natives. In addition, the share of Venezuelans with the PEP working in the formal sector is small, supporting the fact that the impact of the PEP regularization on the Colombian labor market is limited (Bahar et al., 2021).²⁰ Note that it is not possible to observe informal workers in the administrative data, but they represent around half of all workers employed in Colombia.

Table 1: **Descriptive statistics for natives and immigrants in the formal sector**

(a) Colombians							
	Age		Male (%)		Real wages (USD)		N
	Mean	Std. dev.	Mean	Std. dev.	Mean	Std. dev.	
2013	37.0	10.8	0.56	0.496	421.4	504.0	7,335,989
2015	37.2	11.1	0.56	0.497	416.8	484.5	8,391,843
2017	37.8	11.4	0.55	0.497	411.4	477.8	8,064,282
2019	38.2	11.7	0.55	0.498	436.1	505.7	8,363,249

(b) Venezuelans with PEP							
	Age		Male (%)		Real wages (USD)		N
	Mean	Std. dev.	Mean	Std. dev.	Mean	Std. dev.	
2018	30.8	7.8	0.68	0.467	243.7	99.0	12,842
2019	31.8	8.1	0.67	0.472	248.7	98.5	42,752

(c) Other foreigners							
	Age		Male (%)		Real wages (USD)		N
	Mean	Std. dev.	Mean	Std. dev.	Mean	Std. dev.	
2013	40.2	10	0.61	0.487	1,310.0	1,535.9	20,978
2015	39.5	10.2	0.64	0.481	1,253.4	1,446.6	27,730
2017	39.3	10.3	0.63	0.483	1,050.8	1,325.0	31,553
2019	39.6	10.4	0.58	0.494	990.5	1,293.0	39,704

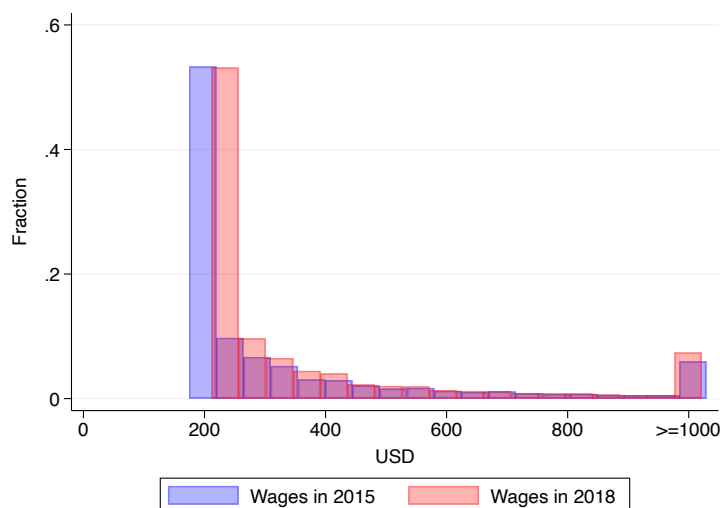
Note: This table reports the descriptive statistics for Colombians, foreigners, and Venezuelans with PEP between 18 and 64 years of age. Only workers with full days of employment recorded in PILA and a positive health contribution are considered for wages and the number of observations. I only observe Venezuelans with PEP since 2018, after the law's enactment. The real wages are deflated using the Consumer Price Index (CPI) from DANE for prices in 2018. Colombian pesos to USD using 2020 exchange rates from the World Bank. For self-employed workers, observed wages in PILA correspond to 40% or more of their actual wages by law, with the minimum wage as a lower bound. Source: PILA, 2013–2019, August.

Figure 1 shows how binding the minimum wage in Colombia is for a large portion of formal

²⁰Note that other Venezuelans can work in the formal sector in the group of foreigners, but it is not possible to correctly identify them.

workers. In 2015, around 40% of all formal workers earned the minimum wage.²¹ Importantly, the national minimum wage in Colombia must increase, by law, more than the inflation in the preceding year. This downward rigidity suggests why, in general, there are no real wage drops (but more layoffs) for minimum wage workers in the face of positive labor supply shocks or negative demand shocks. Last, in the period of analysis (2015-2019), the minimum wage increased in real terms by less than 3% each year, reducing the concern of additional impacts of the minimum wage on employment.

Figure 1: **Histogram of wages by years**



Note: The sample is restricted to native workers between 18 and 64 years with full employment days in the month and positive wages. Wages are in nominal terms. Colombian pesos to USD using 2020 exchange rates from the World Bank. The chosen bin width is 45. Source: PILA, 2015–2018.

Descriptive Statistics for Formal Firms. I aggregate the worker information of the PILA at the firm level to describe patterns in the workforce composition of formal firms.²² Table 2 is split into seven firm size categories to show certain facts. First, regarding sex, male workers are the main group in all formal firm sizes, especially for small-medium firms (between 10 and 999 workers), where more than 60% of formal workers are males. Second, smaller firms have older workers on average (39.6 years), while larger firms have younger workers (36 years). Third, average wages are growing with firm size, from around 272 USD in firms with 1 to 4 workers to 531 USD in

²¹To contribute to the health and pension system in Colombia, the worker must declare a labor income equal or greater to the minimum wage, so many self-employed workers (who decide how much is their observed income in PILA) declare the minimum wage even if they earn more or less.

²²For the firm analysis, I eliminate self-employed workers from the main sample.

firms with more than 1000 workers. Last, Appendix Figure C.1 plots the histogram of firms by size overlay with the total number of employees in each firm size. Although most firms are concentrated in the interval size of 1 to 4 employees, the number of employees is more evenly distributed across different sizes of firms.

Table 2: **Descriptive statistics by firm size**

Firm size (# of workers)	Average				Firms
	Employment	Male (%)	Age	Real wages	
1-4	2	0.56	39.6	271.9	206,456
5-9	7	0.60	37.5	304.2	64,347
10-19	13	0.61	37.2	329.2	42,207
20-49	30	0.63	36.9	360.9	28,625
50-99	69	0.65	36.7	394.8	10,032
100-999	259	0.63	36.9	443.1	10,107
1000 and more	2677	0.58	36.0	530.8	859

Note: This Table reports the descriptive statistics for seven groups of firm size. I deflate real wages using the CPI from DANE for prices in 2018. Then, I transform Colombian pesos to USD using 2020 exchange rates from the World Bank. I only consider employees for constructing firm sizes. Source: PILA, 2015-August.

3 Empirical Strategy

To quantify the evolving impact of immigration on worker outcomes, I estimate the following differences-in-differences specification from $t = \{2012, \dots, 2019\}$ estimating separate yearly regressions of the following form:

$$\Delta Y_{i,l,t} = \delta_t + \theta_t \Delta M_{l,2018} + X_i' \beta_t + \Delta u_{i,l,t}. \quad (1)$$

Here, the outcome $\Delta Y_{i,l,t}$ is the difference in wages, employment, or earnings within workers' pre- and post-treatment years relative to 2015. The immigration shock $\Delta M_{l,2018}$ varies for each local labor market l , and the vector X_i contains individual characteristics in 2015; namely, dummies for six age groups interacted with dummies for sex and self-employment.²³ Broadly, this specification compares individuals with similar observables in the base period but who were working in local labor markets with different exposure to the immigration shock, which I will describe below in detail. Hence, θ_t measures the worker-level impacts of migration, where $\theta_{2015} = 0$ by construction.

²³Education information is not available, and I do not use industry information because the Health Ministry did not verify industry codes in PILA, so their measurement error is relevant.

By taking differences, I can net out any individual constant unobservable that can confound the impact of migration. Lastly, the intercept for each year is δ_t , and I cluster the standard errors in all the specifications at the level of the treatment, which are the FUAs (defined as G and equal to $G = 109$).

The individual outcomes are more precisely defined as follows. First, the employment outcome is $e_{i,l,t} - \sum_{k=2013}^{2015} e_{i,l,k}/3$, where $e_{i,l,t}$ is the indicator of employment in the formal sector for worker i in local labor market l in period t . As in [Yagan \(2019\)](#), I consider the average employment in the pre-shock period to transparently allow for varying labor trajectories of workers in the formal sector. In the event study figures, however, I take the simple difference with the base period ($e_{ilt} - e_{il,2015}$) to avoid pre-treatment coefficients being mechanically around zero. Second, the wage outcome is $\frac{w_{i,l,t} - w_{i,l,2015}}{w_{i,l,2015}}$, so it measures the percentage change in wages $w_{i,l,t}$ for each worker i with respect to 2015, so the worker must be observed both in 2015 and t . Third, the earnings outcome is $\sum_{t=2016}^{t=2018} \frac{Earnings_{it}}{Earnings_{i,2015}}$ and it measures changes in the evolution of earnings normalized by the earnings in the pre-shock period, as in [Autor et al. \(2014\)](#). If the worker is not employed in the formal sector in any given period, the earnings are zero, so this outcome yields a combined effect of the observed changes in employment and wages.²⁴

The immigration shock $\Delta M_{l,2018}$ is defined as follows:

$$\Delta M_{l,2018} = \frac{L_{Ven,l,2018} - L_{Ven,l,2015}}{L_{Total,l,2018}}, \quad (2)$$

where the numerator is the stock of employed migrants from Venezuela (between 18 and 64 years) in local labor market l who arrived in Colombia in the previous 5 years, starting from 2018, minus the stock of employed migrants from Venezuela in l whose year of arrival was 2015 or earlier according to the census. Employed migrants are either Venezuelans or returning Colombians from Venezuela, and the denominator $L_{Total,d,2018}$ is the total employed population in the local labor market. I focus and interpret mainly the coefficient of 2018 in the regressions (i.e., θ_{2018}) to match the year of the census and avoid rescaling the shock as for the coefficients of other periods. Lastly, having this constant in-time immigration shock is useful because it exploits the full count of a census instead of a survey to construct migration shares, and it also allows for placebo tests on pre-trends within

²⁴I define workers with less than 30 days of employment in the social security contribution as missing wages (the wage analysis is focused only on full-time workers).

the same analysis in a transparent manner.

Because migrants self-select into areas where the economic opportunities are better, the immigration rate $\Delta M_{l,2018}$ is likely to be endogenous, and its coefficient is downward biased (see ordinary least squares (OLS) estimates of Figure 3a and 3b). Thus, to consistently estimate the effect of immigration on the outcome variables, I instrument the immigration rate $\Delta M_{l,2018}$ with the distance to the nearest crossing bridge with Venezuela and with past settlements of Venezuelans. The motivation for the IV approach follows.

First, distance is exploited as an instrument since Colombia and Venezuela share 2,220 kilometers of terrestrial borders. Therefore, arrivals to the local labor market l are a function of travel distance between the two countries, as distance acts as a time and economic constraint for Venezuelan immigrants. A threat to this identification strategy is that border departments might be more affected, in terms of economic shocks (such as less trade), than the counterpart far-located states from the Venezuelan crisis (violation of the exclusion restriction).

Appendix Figure C.2a shows suggestive evidence that the trade shock arising from the Venezuelan crisis started years earlier than the immigration shock. Importantly, in the post-treatment period, border department exports to Venezuela are regularly around zero. Another important piece of evidence is that I find insignificant employment and wage effects in the largest firms, presumably more affected by trade shocks and less affected by immigration shocks (as migrants disproportionately concentrate in small firms). In addition, I plot log GDP for border and non-border departments over time to show that it is evolving similarly before the immigration shock, suggesting that any trade impact on economic activity is limited (see Figure C.2b). Last, I exclude border areas from the main sample and find similar point estimates, only not significant for wages. With this suggestive evidence in mind, formally, it is required that distance fulfills the following exogeneity assumption $E[f(dist_l)\Delta u_{lt}] = 0$.

The other instrument constructed uses past settlements of Venezuelans in the spirit of Altonji and Card (1991) and Card (2001), and it is defined as:

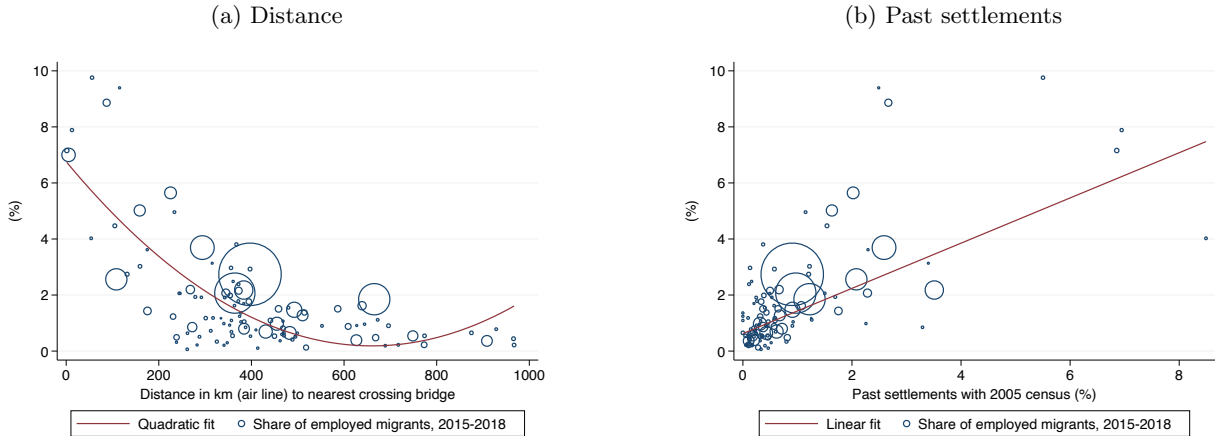
$$z_l = \left(\frac{Ven_{l,2005}}{Ven_{2005}} * M_{2018} \right) / L_{l,2005}, \quad (3)$$

where the first term is the share of Venezuelans in FUA l (according to the 2005 population census),

normalized by the working-age population $L_{l,2005}$ in l at 2005, whereas M_{2018} is the number of migrants in Colombia that arrived between 2018 and 2016 according to the census. I use past settlements as the other instrument because newly arriving immigrants will likely move to areas with previously established Venezuelans. To have a valid instrument, it is required that past settlements are related to new arrivals but not related to time-varying shocks (i.e., $E[z_l \Delta u_{lt}] = 0$).

Figures 2a and 2b show the first stage of the immigration shock $\Delta M_{l,2018}$, for the 109 FUAs defined for this analysis against the instruments. I show the instruments' relevance and functional form in these figures. For the first instrument, a larger distance from a crossing bridge decreases the share of employed migrants in the FUAs until a point where longer distances do not imply lower immigration rates, so the slope of the curve bends downward. For past settlements, there is a positive relationship against the immigration rate that appears to be linear. The immigration shock at the FUA level is quite large, as some areas experience an increase in the share of employed migrants that represent between 7% to 10% of their overall employed population.²⁵

Figure 2: Immigration rates and the two instruments



Note: I weigh dots by formal employment according to the PILA in 2015. In (b), I remove one area to narrow the x-axis values. Functional Urban Areas in Colombia (G=109). Source: CNPV, 2018.

Figures 2a and 2b are constructed at the FUA level. Yet, since this paper aims to estimate the impact of immigration at the individual level, the first stage of the two-stage least squares regression (2SLS) is going to weigh each FUA differently by the number of individual observations

²⁵Delgado-Prieto (2022) uses the department as the area of analysis because of sample limitations of the GEIH survey, but with administrative data there are no sample issues when constructing more detailed areas.

available.²⁶ With this in mind, the first-stage model is:

$$\Delta M_{l,2018} = \delta + f(dist_l) + z_l + v_l \quad (4)$$

Here, $f(dist_l)$ is equal to a linear and quadratic term of distance to the nearest crossing bridge, whereas z_l are the past settlements of Venezuelans. In this equation, the error term is v_l , which captures the endogenous component of $\Delta M_{l,2018}$. I combine the two instruments in the analysis as past settlements or distance capture different exogenous components of migration while increasing the R^2 of the first-stage regression (see Table A.1).²⁷ As a result, I estimate equation (1) throughout the paper using 2SLS with the aforementioned instruments.

4 Worker Responses

This section documents the impact of immigration on formal wages and employment at the worker level, then it focuses on the heterogeneity of the effects across worker characteristics. To start, I show wage and employment event study estimates using OLS and 2SLS. One advantage of the proposed empirical specification is that it is possible to test for differential trends of the outcome before the immigration shock happens. Importantly, there are no significant pre-trends for employment and wages in this setup that can confound the impact of immigration. Figure 3a shows that the OLS coefficients are close to zero, presumably downward biased, as immigrants are expected to arrive in the areas that offer better economic opportunities. The 2SLS helps to reduce this bias, so its coefficients are more negative. Particularly, I find that in 2018 a one percentage point (pp) increase in the share of employed migrants in a given area reduces the probability of employment in the formal sector by 1.1 pp (see Figure 3a).²⁸

To interpret this coefficient, I use the labor force survey to measure the probability of employment in the formal sector, and for workers between 25 and 55 in 2015, it is equal to 0.42. So, a 1.1 pp drop is equivalent to a 2.4% decrease relative to the mean. More broadly, a worker located in

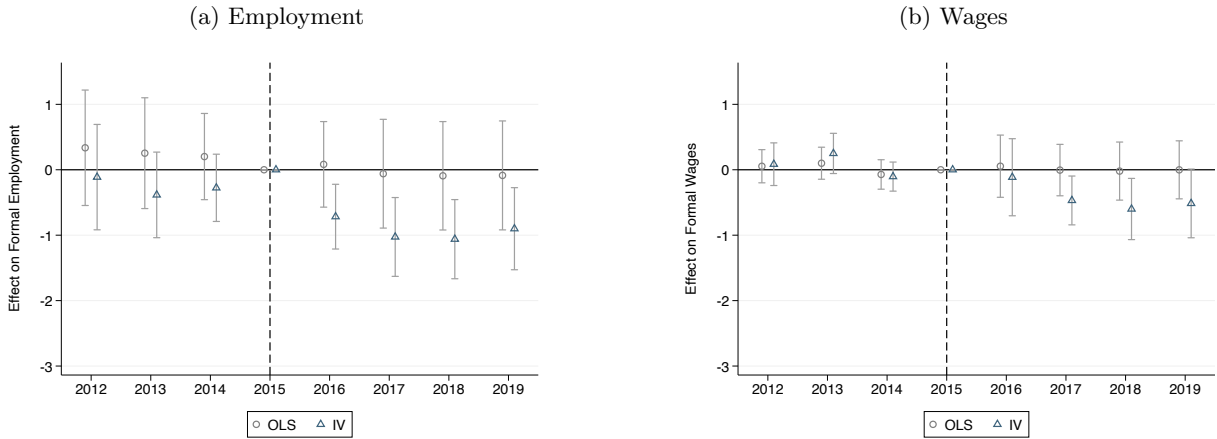
²⁶Hence, the first stage varies slightly depending on the sample used.

²⁷Notably, the main coefficients do not change if I use one instrument instead of both.

²⁸This regression uses as dependent variable $e_{i,l,2018} - e_{i,l,2015}$, which captures the difference in the employment indicator in 2018 relative to the base period. In the heterogeneity analysis, the regression uses as a dependent variable the employment change $e_{i,l,t} - \sum_{k=2013}^{2015} e_{i,l,k}/3$, which yields slightly less negative coefficients as it uses for comparison the average of employment in the pre-shock period.

the LLM in the 75th percentile of exposure relative to one in the 25th percentile of exposure has a relative drop of 3.6% in the probability of formal employment.²⁹ Regarding formal wages, I find a coefficient of -0.6% in 2018 for a one pp increase in the immigration shock (see Figure 3b).³⁰ A worker in the 75th percentile of exposure relative to one in the 25th percentile of exposure has a relative drop of 0.9% in their formal wages. Thus, the impact on wages is minor compared to the one on employment.

Figure 3: **Event study estimates on individual wages and employment**



Note: I estimate equation (1) separately by year. The sample is restricted to natives between 25 and 55 years old. In panel (a), there are 6,706,035 workers, while in panel (b), this varies slightly by year as the worker must be employed in the post-treatment and base year. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). 95% confidence interval. The coefficients for employment (in percentage points) and wages (in percent) are already multiplied by 100. Workers are observed in August of each year. Source: PILA 2012–2019.

For the rest of the paper, I focus on the heterogeneity of wage and employment estimates using workers' and firms' characteristics before the immigration shock, specifically, the characteristics in 2015. In this case, the coefficients for each subgroup come from separate regressions of the main empirical specification (see Equation 1). The first worker characteristic is job type, they can be employees or self-employed.³¹ Self-employment in Colombia represents about half of the employed population, mainly working in the informal sector but with a large share of workers in the formal sector (around 18% of all native formal workers were self-employed in 2015). Figure 4 shows a drop

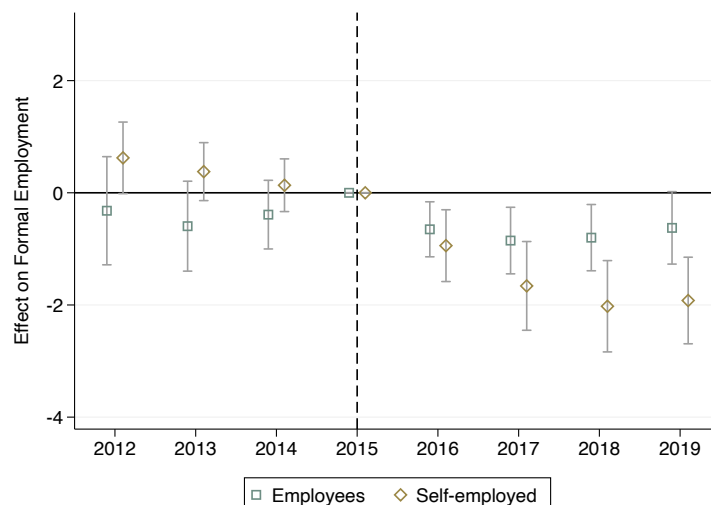
²⁹The 25th and 75th percentile migration rate is 0.6% and 2.1%, respectively. So, $(2.1-0.6)*2.4=3.6$.

³⁰I do not compare formal employment and wage estimates with other countries, as there are no papers, to the best of my knowledge, that estimate worker-level effects in developing countries. Yet, I compare later on these estimates with regional-level estimates from the Colombian setting.

³¹I use only IV hereafter because OLS estimates are inconsistent (see Figures 3a and 3b).

in the probability of being a formal worker for self-employed natives, more negative than the one for employees. Most self-employed in the private sector decide voluntarily whether to contribute or not to the social security system, so opting out from the formal sector is less costly for them than for employees.³²

Figure 4: **Event study estimates on employment by job type**



Note: I estimate equation (1) separately by year and characteristic. The sample is restricted to natives between 25 and 55 years old. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). 95% confidence interval. The coefficients for employment (in percentage points) are already multiplied by 100. Workers are observed in August of each year. Source: PILA 2012–2019.

After showing suggestive evidence that the instruments do not predict native wages or employment trends in the pre-treatment period, I focus, for the rest of the analysis, on the coefficient of 2018 (the year of the immigration shock from the census).³³ I continue the heterogeneity analysis with the standard variables used in the migration literature, but later on, I exploit firm characteristics and develop a more systematic analysis using a machine learning algorithm.

The next results are by age groups and sex, which are also the controls used in the main specification. Figure 5a shows a pronounced decline in the probability of employment in the formal sector as the worker ages. In contrast, the pattern is not equally clear for wages, and I find similar negative estimates in all age groups. I extend the sample of analysis in Appendix Figure B.1 to

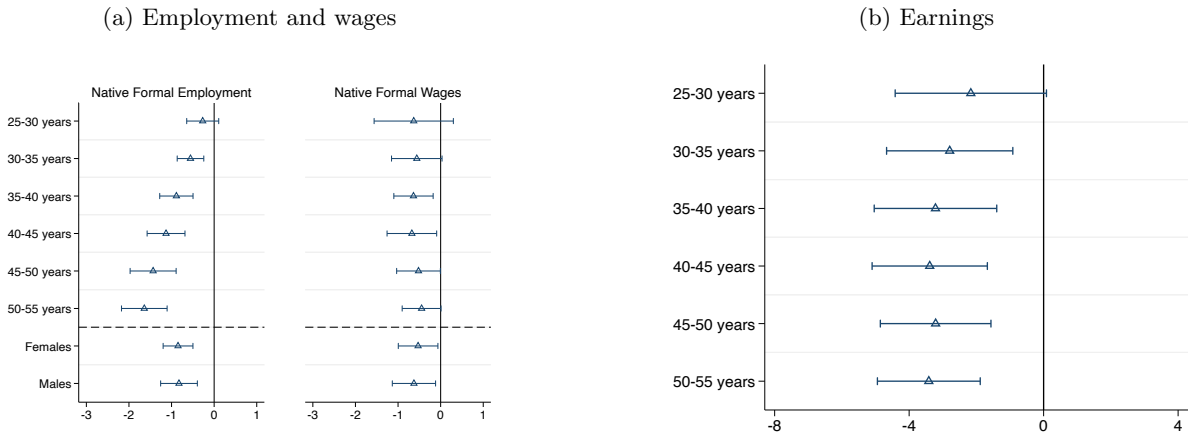
³²The labor income for self-employed is noisy in PILA, but the point-estimates on wages are also more negative than for employees.

³³Nonetheless, in Appendix Table F.1 I show there are no systematic pre-trends by worker or firm categories on employment or wages.

include labor market entrants (18 to 24 years) and workers close to retirement (56 to 64 years) in the base period. For employment, the highest negative effect is observed in the oldest workers, suggesting that they could be retiring earlier, while for wages, again, there are no stark differences. Last, in terms of sex, the impact on employment and wages is alike; there are no differential effects in this group category.

I then analyze the impact on the earnings outcome, which yields a combined effect of the observed changes in employment and wages. Figure 5b shows that workers above 30 years old present a relatively similar reduction in earnings as their confidence intervals overlap. Indicating that even if older workers are more displaced from the formal sector, younger workers experience a greater reduction in their wages.

Figure 5: **Estimates by age group, 2015–2018**



Note: I estimate equation (1) separately by subgroups. The sample is restricted to natives between 25 and 55 years old. In (a), the dependent variables are employment and wages relative to the base period. In (b) dependent variable is $\sum_{t=2016}^{t=2018} \frac{Earnings_{it}}{Earnings_{i,2015}}$. I use as controls the interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). 95% confidence interval. The coefficients for employment (in percentage points), wages and earnings (in percent) are already multiplied by 100. Workers are observed in August of each year. Source: PILA, 2013–2019 .

To complement the pattern of employment effects by age group, I also calculate the labor supply elasticities, at the extensive margin, for each of these age groups (i.e., $\eta_w^s = \frac{\Delta L}{\Delta w}$). Appendix Table B.1 shows that as native workers age, their labor supply is more elastic. That is, the responsiveness to work from wage changes is greater for older than younger workers. In Germany, Dustmann et al. (2017) estimate the local labor supply elasticity by age groups and document a similar result: it is increasing in workers' age.

To continue understanding which are the most affected types of workers, I now exploit the number of years the worker has been employed in the same firm (i.e., job tenure) up to the base period of 2015.³⁴ Appendix Figure B.2 splits the sample by job tenure of native workers (from zero to more than nine years of tenure). Notably, the shock of employment due to immigration is more severe on workers with fewer years in the same firm, still the coefficient is less negative than the one on older workers.³⁵ This result is partly explained by the fact that the severance payment increases with workers' tenure, so it is more costly for firms to dismiss workers, and partly by the accumulation of firm-specific human capital, as they are less substitutable to migrants with similar characteristics.

The last two figures suggest that older workers and workers with lower tenure have the most significant drop in formal employment from the immigration shock. To better explain how workers react, I combine their age and job tenure. Appendix Table B.2 shows that the age variable is more relevant for employment, as native workers below 35 present an insignificant effect on employment, independent of whether they have low or high job tenure. On the other hand, native workers above 35 present a significant negative effect on employment when they have low and high job tenure, but the effect is much higher for the workers with lower tenure (−1 pp versus −0.3 pp). Regarding wages, there are no clear differential effects across tenure and age.

4.1 Distributional Impacts of Immigration

I then estimate the impact of immigration on workers across the distribution of wages. For this exercise, I divide native workers into seven bins according to their local wage distribution in 2015.³⁶ Figure 6a shows the uneven impacts of immigration: native workers earning the minimum wage suffer the most negative shock on formal employment, while for workers at the rest of the wage distribution, I find insignificant estimates on employment. For these low-wage workers, a one pp increase in the share of employed migrants in a given labor market reduces the probability of employment in the formal sector by 1.5 pp. In contrast, formal workers who earn the minimum wage are the least affected by the immigration shock in terms of wages.

³⁴Self-employed workers are excluded from this analysis as they are not comparable to the average firm.

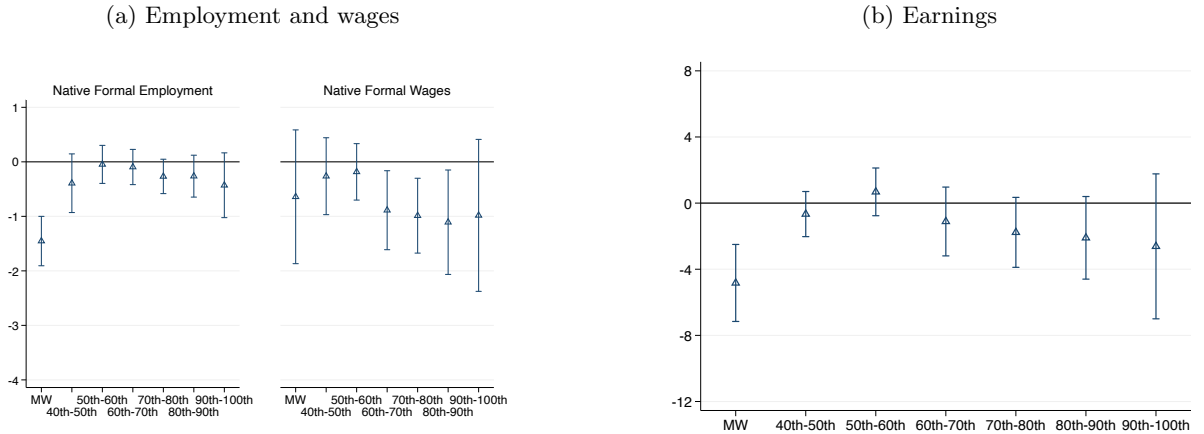
³⁵I construct job tenure from the first year PILA is available (2007).

³⁶As a robustness check, Appendix Table F.2 shows the pre-treatment coefficients by wage categories on employment and wages. Reassuringly, most of these coefficients are insignificant.

Because the minimum wage is relatively high and binding for around 40% of formal workers in the pre-shock period, the chances of job displacement for these workers are higher.³⁷ Furthermore, the existence of a large informal sector explains part of the large coefficient, as minimum-wage workers are the less educated and, thus, the most substitutable with informal workers who become less costly after the arrival of migrants (Delgado-Prieto, 2022).

For the workers between the 60th and 90th percentile of the local wage distribution, I find a drop of around 1% to 1.2%. This does not necessarily mean a decrease in absolute terms of wages. The coefficient measures the average growth of wages of native workers in areas with more exposure to migration compared to areas with less exposure, so the growth of wages in more affected areas is relatively lower. Last, I estimate the earnings outcome to find which set of workers gets more affected overall. Figure 6b shows that the only significant negative impact on earnings is observed on workers who earn the minimum wage before immigrants arrive, reflecting a stronger effect from the employment losses.

Figure 6: **Estimates by individual wage at baseline, 2015–2018**



Note: I estimate equation (1) separately by subgroups. The sample is restricted to natives between 25 and 55 years old. In (a), dependent variables are employment relative to the pre-shock period and wages relative to the base period. In (b) dependent variable is $\sum_{t=2016}^{t=2018} \frac{Earnings_{it}}{Earnings_{i, 2015}}$. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). 95% confidence interval. The coefficients for employment (in percentage points) and wages (in percent) are already multiplied by 100. Workers are observed in August of each year. Source: PILA, 2013–2019 .

³⁷Conditional on being employed in the two periods, around 75% of minimum wage earners still earn the minimum wage after three years.

4.2 Worker-level and Regional-level Effects

Most of the migration literature only considers the regional responses when studying immigration shocks.³⁸ Since regional responses aggregate several margins of adjustment to immigration, they can lead to different findings relative to worker-level responses, as emphasized in [Dustmann et al. \(2023\)](#). For that reason, I adapt to this setup the employment decomposition they introduce to shed light on these different responses. I decompose the changes in regional formal employment into 1) a displacement of incumbent workers –outflows from formal employment–, 2) hiring of new formal workers or inflows from other regions –inflows to formal employment–, and 3) relocation of existing employed formal workers to other regions.

In this analysis, the worker-level estimates of employment capture the outflows or the displacement of incumbent native workers in the formal sector, while the regional-level estimate from cross-sectional data combines the three margins of adjustment. Appendix Figure B.3 shows the decomposition of the regional formal employment response at the FUA-level (–1.3%) along with the three margins: outflows to non-employment or the informal sector (1.1%), inflows from other regions, non-employment or the informal sector (–0.5%) and relocation to other regions (–0.4%).³⁹ In this case, the most important and only significant margin is the outflows from the formal sector, which differs from the findings of [Dustmann et al. \(2023\)](#), where inflows are the most relevant margin. The results can differ due to the existence of an informal sector where firms can hire after they displace formal workers and the less restrictive job protection in Colombia relative to Germany.

Regarding wage estimates, the worker-level response is –0.6%, while the regional-level estimate in [Delgado-Prieto \(2022\)](#) is insignificant and close to zero. These two responses are complementary and answer different policy questions. As stated in [Dustmann et al. \(2023\)](#), the wage estimates of the worker-level regressions capture the change in the price of labor, holding the composition of the population constant, while the regional-level regressions jointly measure the change in the selection and composition of workers and the price of labor. The differential estimate between the two is

³⁸Recent regional-level studies are [Monras \(2020\)](#) in the US and [Muñoz \(2021\)](#) in the EU. The first documents that low-skilled Mexicans who left their country due to the peso crisis had a high transitory impact on local labor markets in the US. The second exploits a trade liberalization in services across Europe to find a negative regional effect on the employment of domestic workers.

³⁹The decomposition is equal to: $\frac{E_{r1}-E_{r0}}{E_{r0}} = -\frac{E_{r,Out}}{E_{r0}} + \frac{E_{r,In}}{E_{r0}} - \frac{E_{r,Move}}{E_{r0}}$. The first term measures the outflow margin, the second term the inflows margin and the third term the relocation margin. The main distinction, in this case, is that the outflows and inflows margins can be decomposed further into non-employment or the informal sector. Unfortunately, there is no panel data for the informal sector to measure these decompositions.

rationalized in this setup as follows. The immigration shock changes the composition of employed natives and positively selects the individuals remaining in the region (see Figure 6a), therefore mechanically increasing regional formal wages. On the other hand, immigration decreases the price of labor in certain mid- and high-wage subgroups (see Figure 6a), reducing regional formal wages. Hence, this suggests why there is an insignificant formal wage effect at the regional level while having a negative wage effect at the worker level, motivating the analysis of immigration not only for the aggregate local labor markets but for individuals within local labor markets.

Another benefit of individual data compared to regional data is the possibility of estimating inter-regional movements of different types of workers to respond to the immigration shock. For instance, Foged and Peri (2016) document that after the arrival of refugees in Denmark, younger workers are much more mobile. Hence, Appendix Table B.3 shows the impact of movements across regions by age groups. Indeed, younger formal workers tend to move more, but the coefficients are insignificant. Overall, the point estimates decrease as the worker ages, but all are insignificant. The mobility margin of adjustment is less important in this setup.

5 Immigration, Workers, and Firms

In this section, I first develop a partial equilibrium model with heterogeneous firms and types of workers to motivate and interpret the empirical findings. Then, I show that immigration effects for natives vary substantially depending on the type of firm they were employed in before the shock.

5.1 Model

The market structure of the model consists of J firms that hire two types of labor inputs. Specifically, firms hire formal workers F paying payroll taxes and informal workers I off the books to avoid paying the payroll taxes, as in Ulyssea (2018). So, each firm $j = \{1, \dots, J\}$ posts a pair of wages (w_{I_j}, w_{F_j}) that all workers i observe and decide to accept.⁴⁰ Importantly, each firm has different work environments, measured by amenities a_{L_j} , workers have idiosyncratic preferences ϵ_{i,L_j} depending on the fixed labor group they belong $L \in \{I, F\}$. This gives a different job valuation at

⁴⁰The transitions of workers between the formal and informal sectors are out of the scope of the model.

each firm.⁴¹ In this case, the indirect utility of worker i employed at firm j is:

$$v_{i,L_j} = \beta_L \ln w_{L_j} + a_{L_j} + \epsilon_{i,L_j}. \quad (5)$$

Under the assumption that ϵ_{i,L_j} follows a type I extreme value distribution for each of the workers' types $L \in \{I, F\}$ and that the number of firms J is sufficiently large, [Card et al. \(2018\)](#) shows that the firm-specific supply functions are expressed as:

$$\ln I_j(w_{I_j}) = \ln(\mathcal{I}\lambda_I) + \beta_I \ln w_{I_j} + a_{I_j}, \quad (6)$$

$$\ln F_j(w_{F_j}) = \ln(\mathcal{F}\lambda_F) + \beta_F \ln w_{F_j} + a_{F_j}. \quad (7)$$

In this case, the total number of informal workers in the market is \mathcal{I} and of formal workers is \mathcal{F} , where λ_I and λ_F are constant parameters across firms. Moreover, $\frac{d \ln L(w_{L_j})}{d \ln w_{L_j}} = \beta_L$ is the elasticity of labor supply to the firm with respect to its wage. Hence, as $\beta_L \rightarrow \infty$, the supply functions become perfectly elastic, and firms have no monopsony power to set wages.⁴²

Regarding firms, there is a productivity shifter T_j , a price of the good P_j , and a production function Q_j for each firm, such that the profit function of firm j is:

$$\max_{I_j, F_j} \pi_j = P_j T_j Q_j - \tau(I_j) w_{I_j}(I_j) I_j - (1 + \tau_F) w_{F_j}(F_j) F_j. \quad (8)$$

Here, $\tau(I_j)$ represents a convex cost that is increasing on the firm's informal labor size. These convex costs are important to match the stylized fact that informal labor decreases with firm size and captures the cost of evasion related to law enforcement exerted by the government. Particularly, I assume that $\tau(I_j) = I_j^\eta$ with $\eta \geq 0$. The τ_F represents the payroll taxes firms must pay for formal workers, and the production function takes the following form: $Q_j = (\alpha_I I_j^\rho + \alpha_F F_j^\rho)^{\frac{1}{\rho}}$. Thus, formal and informal workers are imperfect substitutes, and the aggregate elasticity of substitution common across all firm types is given by $\sigma = \frac{1}{1-\rho}$. To finish the setup, P_j is the inverse demand function

⁴¹For instance, preferences for working in a firm may refer to distance to the workplace or interactions with coworkers ([Card et al., 2018](#)).

⁴²Here, I exclude any market wage offered in an outside competitive sector as the comparative statics focus is on firm-level responses to immigration and not on market-level responses that have been thoroughly analyzed previously.

defined as $P_j = D_j(T_j Q_j)^{-(1-\epsilon)}$, where $\epsilon^D = -1/(1 - \epsilon)$ is the elasticity of product demand and D_j is the firm-specific product demand.⁴³

I then analyze the impact of an immigration shock that shifts the aggregate informal supply outwards ($d\mathcal{I}$).⁴⁴ I study the firms' response across the wage and employment margin, so the wage elasticity for each type of worker in firm j is $\varepsilon_{w_{L_j}, \mathcal{I}}$ and the employment elasticity for each type of worker in firm j is $\varepsilon_{L_j, \mathcal{I}}$. Allowing for firm-level responses to an immigration shock is the main contribution of this framework. Unsurprisingly, in Appendix E I show after some derivations that the elasticity of informal labor is always positive ($\varepsilon_{I_j, \mathcal{I}} > 0$) and the elasticity of informal wages is always negative ($\varepsilon_{w_{I_j}, \mathcal{I}} < 0$) after an aggregate informal supply shock. Independently from whether informal and formal workers are close substitutes or not.

More interestingly, I show how formal wages of firm j change in response to an aggregate informal supply shock:

$$\varepsilon_{w_{F_j}, \mathcal{I}} = \Omega_j s_{I_j} (\epsilon - \rho). \quad (9)$$

Here, $s_{I_j} = \frac{\alpha_I I_j^\rho}{\alpha_I I_j^\rho + \alpha_F F_j^\rho}$ is the relative contribution of informal work to production before immigrants arrive and $\Omega_j = \frac{1}{\xi_{I_j} \xi_{F_j} - (\epsilon - \rho)^2 s_{I_j} \beta_I s_{F_j} \beta_F}$ is a positive parameter.⁴⁵ Firstly, if informal workers are close substitutes to formal workers (such that $\rho > \epsilon$), then the elasticity of formal wages with respect to aggregate informal labor is negative. Importantly, as the contribution of informal labor to production in firm j increases ($s_{I_j} \uparrow$), the elasticity of formal wages is more negative ($\varepsilon_{w_{F_j}, \mathcal{I}} \downarrow$). Note that, for certain low productivity firms, formal wages can be downwardly rigid due to the existence of a minimum wage, so the formal wage margin is muted (i.e., $\varepsilon_{w_{F_j}, \mathcal{I}} = 0$).

In terms of formal employment, the corresponding expression is equal to:

$$\varepsilon_{F_j, \mathcal{I}} = \Omega_j s_{I_j} (\epsilon - \rho) \beta_F. \quad (10)$$

The implications for formal employment in terms of the substitution parameter (i.e., $\rho > \epsilon$) hold similarly as for formal wages, though the response is now adjusted by β_F . Hence, as the

⁴³For simplicity, in this model, I do not distinguish if the produced good is tradable or non-tradable, only that the firm produces a good.

⁴⁴GEIH survey data shows that around 90% of Venezuelan immigrants are employed in the informal sector.

⁴⁵To show that $\Omega_j > 0$, note that this can be simplified as $\Omega_j = (1 + (1 + \eta - \rho)\beta_I)(1 + (1 - \rho)\beta_F) - (\epsilon - \rho)(s_{I_j}\beta_I + s_{F_j}\beta_F + (1 + \eta - \rho)\beta_I s_{F_j}\beta_F + (1 - \rho)\beta_F s_{I_j}\beta_I)$ which is always positive.

relative contribution of informal workers to production increases ($s_{I_j} \uparrow$), the adjustment on formal employment is more negative ($\varepsilon_{F_j, \mathcal{I}} \downarrow$) as long as informal and formal workers are sufficiently close substitutes.

To summarize, the model I propose points to two main conclusions. The first one is that when the substitutability between formal and informal workers is high, there is a negative response in terms of formal wages and employment to an informal supply shock. The second one is the importance of the production structure to determine how responsive the firm is to an informal supply shock. In particular, as a firm's weight on informal labor for production is higher, it will adjust more their formal wages and employment in response to the immigration shock.

In the model, the firm's informal production share is inversely related to the firm's size. This is because, for larger firms, it is marginally more expensive to hire an additional informal worker due to the convex cost of informal labor $\tau(I_J)$. For that reason, in the empirical results, I focus first on worker responses across the firm size distribution to show the patterns predicted from the model are also observed in the data.

5.2 Worker Responses Across Firms

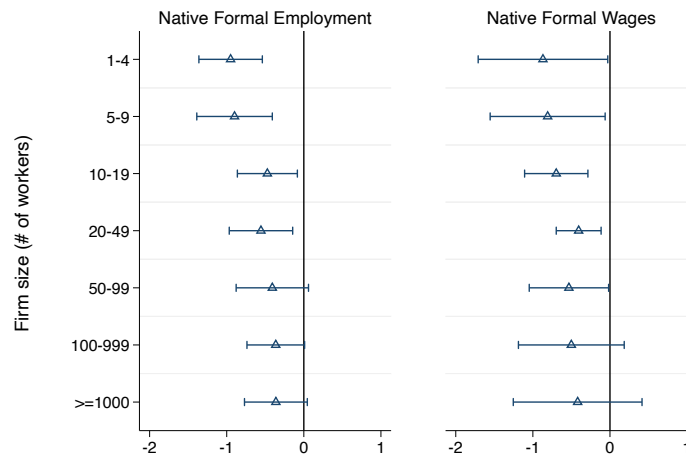
The previous model suggests that workers in certain types of firms should be more affected by an immigration shock, so I turn to the data to test these implications. In this context, using the firm dimension for the heterogeneity analysis is also motivated by three stylized facts. First, Venezuelan immigrants are disproportionally employed in the smallest firms, second, small firms are more likely to pay the minimum wage for their formal workers and, third, small firms employ a higher share of informal workers (Delgado-Prieto, 2022). Hence, the impact of immigration on workers in small firms is more salient, as they can substitute formal for informal labor more easily. In this exercise, I divide workers by firm size categories in 2015 (the year before the immigration shock) and show worker-level employment and wage coefficients for 2018 (the year of the immigration shock from the census).

Figure 7 shows that native workers in firms with less than 50 workers in the pre-shock period suffer the most negative effect on the probability of employment, while workers in bigger firms are less affected. In line with the model's predictions, small firms tend to rely more on informal work for production, and the cost of being caught by authorities is lower in these firms compared to the

largest ones. Thus, when formal and informal workers are close substitutes, it is profitable for the firm to exchange formal for informal labor. [Delgado-Prieto \(2022\)](#) documents, using survey data of the formal and informal sectors, that the share of informal labor increases more in smaller firms after the arrival of migrants, indicating a change in the composition of the firm's workforce.

Regarding wages, workers in the smallest firms (with less than ten workers) present the most negative effect, yet workers in firms with less than 100 workers also present a significant negative effect. A similar pattern is predicted from the model, where smaller firms adjust their wages more to an immigration shock. Lastly, these results are useful to transparently show that trade shocks from the Venezuelan crisis are less of a concern in this setup, as the main effects are observed in the small firms that are directly affected by migration and presumably much less by trade.

Figure 7: **Estimates by firm size, 2015–2018**



Note: I estimate equation (1) separately by subgroups. The sample is restricted to native employees between 25 and 55 years old. Dependent variables are employment relative to the pre-shock period and wages relative to the base period. I use as controls the interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). 95% confidence interval. Workers are observed in August of each year. The coefficients for employment (in percentage points) and for wages (in percent) are already multiplied by 100. Source: PILA, 2013–2019.

Next, I quantify worker-level effects by exploiting other relevant firm characteristics. In this case, I consider the years the firm appears in the administrative records up to the pre-shock period, that is, a proxy of the firm's age. Appendix Figure B.4 shows results for native workers according to the age of their firm in 2015. For employment, workers in younger firms present a more negative impact than workers in older firms, while the pattern is not equally clear for wages. Still, workers

in the youngest firms present the most negative coefficient on wages. The positive correlation between firm size and age helps explain previous negative findings, as smaller firms tend to be younger. However, [Fort et al. \(2013\)](#) document different kinds of responses from young and old firms depending on their size during the business cycle, so I combine these characteristics to measure how worker-level effects vary. Appendix Table [B.4](#) shows that native workers in the youngest firms present a significant negative effect on employment and wages, independent if their firm is small or large, but the coefficient for wages is more negative in younger firms. On the other hand, native workers in older firms present a significant negative effect on employment and wages only in the smallest firms.

5.3 Wage Decomposition

With access to the universe of workers and firms in Colombia, it is possible to construct a measure of the wage premium of firms. I estimate the standard AKM model proposed in [Abowd et al. \(1999\)](#) that decomposes the contribution of firm-specific and worker-specific constant characteristics to log formal wages (lnw_{it}). The AKM model is expressed as:

$$lnw_{it} = \alpha_i + \psi_{j(i,t)} + X'_{it}\beta + \epsilon_{it}. \quad (11)$$

Here, α_i captures the unobserved worker effect, ψ_j captures the unobserved firm effect, and $j(i, t)$ refers to the firm j where worker i is working in t . X_{it} is a vector of controls that are age squared and cubic after being normalized and year FEs. Lastly, ϵ_{it} is the error term. To rule out possible endogenous movement of workers due to the immigration shock, I estimate the model from 2010 to 2015 for August ($T = 6$).

In these types of models, the firm FEs are identified through the movements of workers across firms. These movements are taken as exogenous conditional on worker and firm effects (i.e., $E[\epsilon_{it}|\alpha_i, \psi_{j(i,t)}, X_{it}] = 0$). Still, AKM models can present issues when estimating firm FEs with limited mobility of workers across firms, especially in smaller firms ([Andrews et al., 2008](#); [Bonhomme et al., 2020](#)). Several strategies have been proposed to address this limitation, one of them would be to exclude all the small firms from the estimation.⁴⁶ Since the majority of migrants are

⁴⁶Another solution is to aggregate small firms according to their observable characteristics like industries, but as I observe industries with measurement error, the aggregation could include high-productivity sectors with low-

working in small firms, I prefer to restrict the sample to the largest set of firms connected by the mobility of workers to reduce the concern of limited mobility bias.⁴⁷ With this in mind, I estimate the vector of firm FEs $\hat{\psi}_1, \dots, \hat{\psi}_J$ and worker FEs $\hat{\alpha}_1, \dots, \hat{\alpha}_N$.

To begin, Appendix Table D.2 shows the decomposition of the variance of wages $Var(lnw_{it})$ in the formal sector of Colombia using the leave-out method proposed by Kline et al. (2020). Worker effects explain 50.2% of the variance and firm effects explain 15.7%, in line with the related literature cited in Card et al. (2018). Furthermore, there is a positive sorting of high-wage workers into high-wage firms, which explains an additional 21.6% of the variance. In four European countries and the US, this sorting explains between 10% to 20% of the wage inequality (Bonhomme et al., 2020).⁴⁸

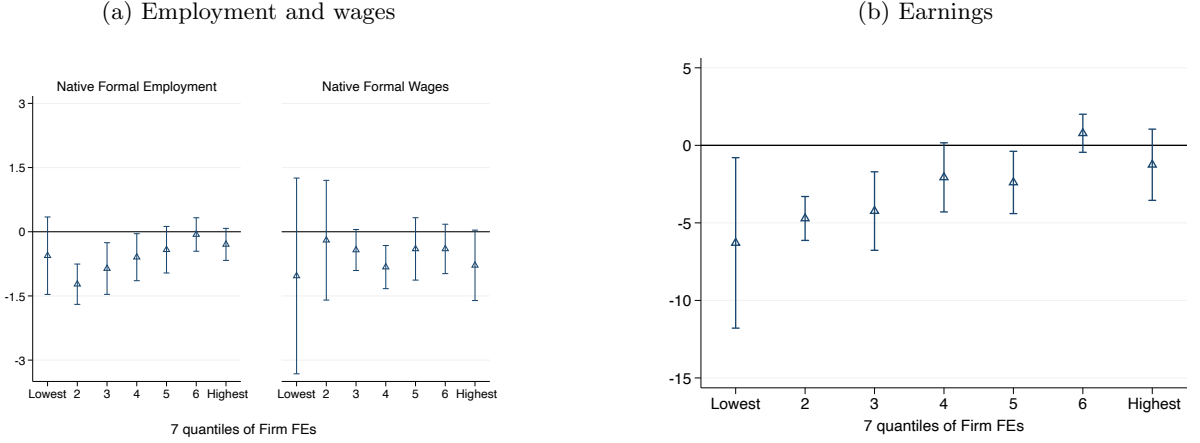
Then, using the estimated $\hat{\psi}_j$, I divide workers by seven quantiles of firm FEs or firm-specific pay premiums, which I now refer to as lowest- or highest-paying firms, to compute the impact of immigration. Figure 8a shows that workers at low-paying firms suffer negative employment losses while having insignificant wage changes. In contrast to workers in middle-paying firms, where the wage and employment response is negative. A possible explanation for this result is that the share of firms in the low-pay sector grows as immigrants work mainly in these firms, and as a response, firms in the high-pay sector extract higher rents from workers and hence reduce their wages, as shown theoretically in Amior and Stuhler (2022). Lastly, to define if wages or employment are decreasing more the earnings, I estimate this outcome across quantiles of firm FEs. Figure 8b shows that workers in the lowest-paying firms present a more pronounced decline in earnings than workers in middle- or high-paying firms.

productivity ones, misleading the estimates.

⁴⁷The leave-out estimation of variance components in Kline et al. (2020) is a different solution to this problem. However, this method yields the corrected moments of interest (i.e., the variance of firm and workers FEs with their corresponding covariance) but does not estimate the corrected vector of $\hat{\psi}_j$ used in this paper.

⁴⁸The four European countries are Austria, Italy, Norway, and Sweden. The method they use for estimating the sorting in 6-year panels is the correlated random effects based on the grouping proposed in Bonhomme et al. (2019).

Figure 8: **Estimates by quantiles of firm FEs, 2015–2018**



Note: I estimate equation (1) separately by subgroups. The sample is restricted to native employees between 25 and 55 years old. In (a), dependent variables are employment relative to the pre-shock period and wages relative to the base period. In (b) dependent variable is $\sum_{t=2016}^{t=2018} \frac{Earnings_{it}}{Earnings_{i,2015}}$. I compute Firm FEs in the first stage using the standard AKM framework, with age squared and its cubic as time-varying controls, for the period 2010–2015. I use as controls in the second stage interactions of sex with six age categories and a dummy for self-employed in the base period. Cluster standard errors ($G=109$). 95% confidence interval. The coefficients for employment (in percentage points) and wages (in percent) are already multiplied by 100. Workers are observed in August of each year. Source: PILA, 2013–2019.

Appendix Figure B.5 shows a similar exercise but dividing by seven quantiles of worker FEs: $\hat{\alpha}_i$. The wage and employment estimates hold similarly as before. High-wage workers present more negative point estimates for wages and the least negative ones for employment. In contrast to low-wage workers, where the wage effect is close to zero, and the employment effects are more negative.

5.4 Heterogeneity by Worker and Firm Characteristics

As shown previously, workers present different employment and wage effects depending on their own characteristics but also on the type of firms they were employed in before immigrants arrived. To illustrate the groups most affected in a more standard way, I restrict the sample to the intersection between subgroups where previous findings indicate a more negative coefficient on workers. In the next section, I present a more systematic analysis of heterogeneity using a machine learning algorithm.

First, Table 3 shows that for minimum wage earners in 2015, immigration reduced the probability of employment in the formal sector by 1.5 pp. For the medium age group, the impact is

less negative (−1.2 pp), while for self-employed workers, the impact is more negative (−2.2 pp). When combining these three characteristics, there are 565,594 workers in the sample, for whom the effect of Venezuelan immigration on the probability of being a formal worker is -2.6 pp, a larger displacement effect.

Table 3: **Most affected native workers in terms of employment, 2015–2018**

	(1)	(2)	(3)	(4)	(5)
Prob. of Employment	-0.841*** (0.192)	-1.453*** (0.231)	-1.188*** (0.227)	-2.194*** (0.327)	-2.647*** (0.388)
Sample restriction					
Minimum wage earners	✗	✓	✗	✗	✓
Medium age (35 years or more)	✗	✗	✓	✗	✓
Self-employed	✗	✗	✗	✓	✓
<i>N</i>	6,706,035	2,205,814	3,915,188	1,103,384	565,594
Clusters	109	109	109	109	109

Standard errors are in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Note: I estimate equation (1) separately by subgroups. The outcome variable is $e_{i,2018} - \sum_{t=2013}^{2015} e_{it}/3$ where e_{it} is the indicator of employment in the formal sector. The sample is restricted to natives between 25 and 55 years old. To understand how large the coefficients are, the size of the formal sector in urban areas, relative to overall employment, was 55.2% in 2015. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). Workers are observed in August of each year. Source: PILA, 2013–2019.

Next, I use the same criteria as in Table 4 to divide the sample by the subgroups with the highest negative coefficient, but for native wages. First, I find that for workers earning more than the minimum wage in 2015, migration reduced average wages by 0.7%. For workers in the smallest firms in 2015, the impact is more negative (−0.8%), while for workers in middle-paying firms in 2015, I find an estimate of -0.8%. When combining these characteristics, there are 53,279 workers in the sample, for whom the effect on wages in 2018 is −1.9% for a one pp increase in the immigration shock. Note that this analysis is subject to arbitrary sample restrictions with a smaller sample size that can lead to differential effects partly due to random variation. Therefore, in the next section, I estimate heterogeneous immigration effects in a data-driven way.

Table 4: **Most affected native workers in terms of wages, 2015–2018**

	(1)	(2)	(3)	(4)	(5)
Wages	-0.600*	-0.711*	-0.827**	-0.804**	-1.908**
	(0.239)	(0.315)	(0.320)	(0.260)	(0.477)
Sample restriction					
Above minimum wage	✗	✓	✗	✗	✓
Small firm (1 and 19 workers)	✗	✗	✓	✗	✓
Middle-paying firm (quantile 4)	✗	✗	✗	✓	✓
<i>N</i>	4,090,973	2,639,040	643,346	195,647	30,772
Clusters	109	109	109	109	109

Standard errors are in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Note: I estimate equation (1) separately by subgroups. The outcome variable is $\frac{w_{i,2018} - w_{i,2015}}{w_{i,2015}}$ where w_{it} are wages in the formal sector. The sample is restricted to natives between 25 and 55 years old. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). Workers are observed in August of each year. Source: PILA, 2015–2018.

5.5 Sorting

Next, I study the reallocation effects of the immigration shock, analyzing changes in the sorting patterns of high- and low-paying workers into high- and low-paying firms.⁴⁹ In this exercise, the outcome is constructed using the values of $\hat{\psi}_j$ from equation (11) and exploiting the movements of workers between firms in the post-treatment period. More concretely, the outcome is the change in the AKM firm FEs in 2018 relative to 2015: $\hat{\psi}_{i,\{j=2018\}} - \hat{\psi}_{i,\{j=2015\}}$. If the worker remains in the same firm during that period, the difference is zero.⁵⁰ Results are shown by seven quantiles of worker FEs to determine if low- or high-wage workers are sorting more into low- or high-paying firms after the immigration event. A positive coefficient means a positive sorting effect from immigration. Figure 9a plots the estimates for these categories, and none of them present significant results.⁵¹ There is no differential sorting due to immigration.⁵² Thus, to explain the negative wage coefficient of workers in high-paying firms, there must be lower wage growth within these firms. In a related exercise, I also measure if workers are moving to larger or smaller firms after the immigration shock,

⁴⁹For France, [Orefice and Peri \(2020\)](#) study the changes in worker-firm sorting after immigrants arrive, they find that high-paying workers are moving more into high-paying firms.

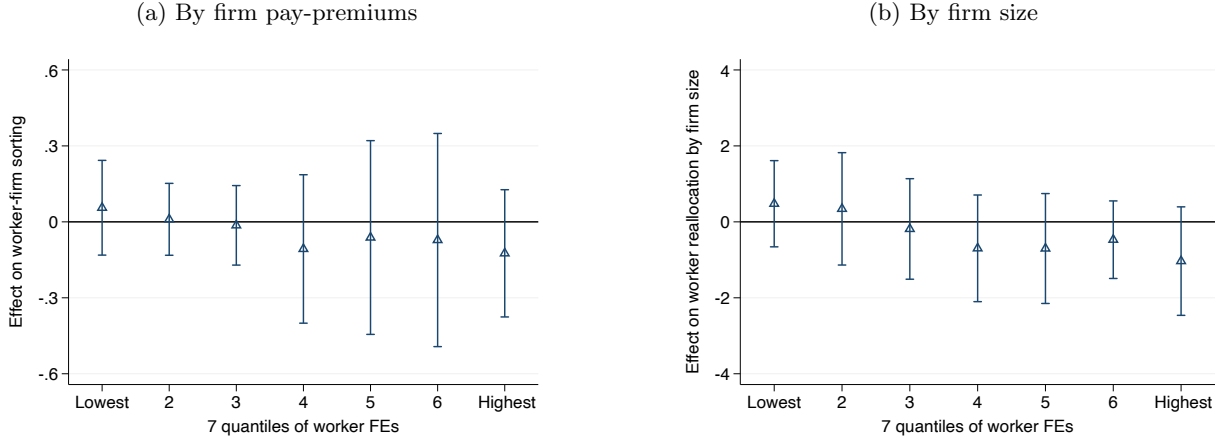
⁵⁰Since the FEs are constructed for the pre-policy period, all workers that belong to firms created after 2015 are not considered in the analysis. Last, the estimated firm FEs are transformed into positive values to construct the outcome.

⁵¹Compared to Germany, the introduction of a nationwide minimum wage led to the reallocation of low-wage workers into higher-paying firms ([Dustmann et al., 2022](#)).

⁵²This is partly attributed to the macroeconomic conditions of the labor market in Colombia during the period studied, as unemployment slightly increased.

and again, there does not seem to be reallocation on this margin (see Figure 9b).

Figure 9: **Reallocation estimates by quantiles of worker FEs, 2015–2018**



Note: The sample is restricted to natives between 25 and 55 years old. The dependent variable in (a) is the change in $\hat{\psi}_{i,\{j=2018\}} - \hat{\psi}_{i,\{j=2015\}}$ and in (b) is the change in the categories of firm size in 2018 relative to 2015, both measured in the base period. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors (G=109). 95% confidence interval. Workers are observed in August of each year. Source: PILA, 2013–2019.

5.6 Hiring Patterns of Formal Firms

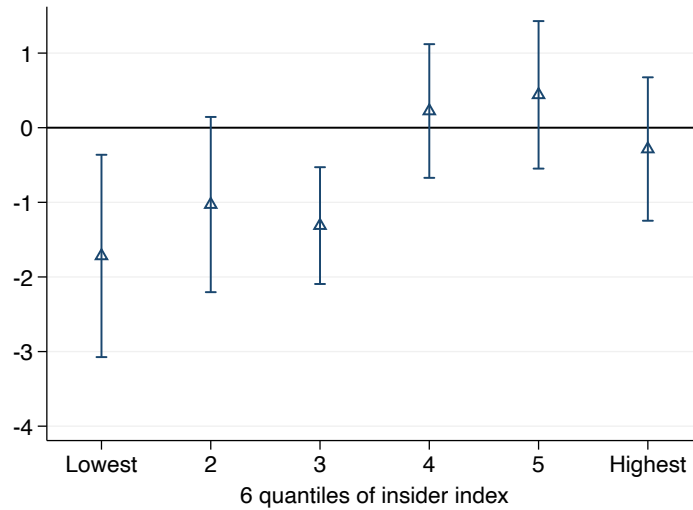
In the absence of informal worker-level data in the administrative records, it is possible to construct a measure of connectedness with the informal sector for formal firms apart from the standard firm size variable. To build this proxy, I develop an insider index similar to the poaching index constructed in [Bagger and Lentz \(2019\)](#). The intuition of the index is that firms are divided by the share of hires that come from *outside* the formal sector, that is, workers who have not been employed in the formal sector before (most likely labor market entrants or workers from the informal sector) and from *inside* the formal sector, that is, workers who were employed in other formal firms at the time of the hire or may be unemployed but have worked previously in the formal sector. In some way, it is a measure of revealed preferences of workers, as firms that hire more from the formal sector are more desirable, while the firms that hire more labor market entrants or from the informal sector act as “gatekeepers” for workers who want to enter the formal sector. The insider index is constructed for every firm j before and after immigrants arrive,

$$\pi_{j,t} = \frac{N_{j,t}^{In}}{N_{j,t}^{In} + N_{j,t}^{Out}}, \quad (12)$$

where $N_{j,t}^{In}$ is the number of firm j in year t hires that have been employed before in the formal sector, and $N_{j,t}^{Out}$ is the number of firm hires that come outside the formal sector.⁵³ Next, I take differences in the insider index between 2018 and 2015 (i.e., $\pi_{j,2018} - \pi_{j,2015}$) at the worker level according to the firm the worker was employed in 2015.

Figure 10 shows results for this outcome by six quantiles of the insider index in the pre-shock period. Interestingly, formal firms that tend to hire workers from the informal sector are having a negative effect on their insider index after immigrants arrive, indicating that these firms are hiring relatively fewer workers that at some point belong to the formal sector (a 1 pp increase in the migration rate reduces the insider index of the lowest type of firms by around 1.7 pp). On the opposite, for firms that have a higher share of hires within the formal sector, the insider index does not change much. This measure is an important way of showing that some firms are opting out or poaching less from the formal sector for new hires, especially the firms that are supposedly more connected to the informal sector, according to the insider index.

Figure 10: **Estimates by quantiles of the insider index, 2015–2018**



Note: Dependent variable is the change in the insider index for workers employed in firm j in the base period between 2015-2018. I use as controls interactions of sex with six age categories in the base period. I cluster standard errors (G=109). 95% confidence interval. Workers are observed in August of each year. Source: PILA, 2013–2019.

⁵³I can record the hires of firms since 2007 and can build the measure up until 2018 for February and August in each year. If the firm did not make any hiring in the year, the index takes a missing value.

5.7 Exit and Entry of Formal Firms

In a related exercise, I test how likely it is that firms disappear entirely from the formal sector after immigrants arrive. Table 5 shows evidence that formal firms present a negative growth in places that receive more immigrants relative to the places that receive fewer immigrants, yet the coefficient is insignificant. If I decompose the growth in the exit and entry margin of formal firms, there is a significantly higher firm exit. A 1 pp increase in the immigration shock increases the firm exit rate by 1.2%. On the other hand, the firm entry rate is close to zero, indicating that the opening rates of formal firms do not change after the labor supply shock.

Table 5: **Decomposition of Firm Growth, 2015–2018**

	(1) Total Firms	(2) Firm Exit	(3) Firm Entry
$\Delta M_{l,2018}$	-1.127 (0.750)	1.190* (0.582)	0.063 (0.935)
N	109	109	109

Standard errors are in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Note: Regressions are estimated at the regional level for 109 FUs weighted by their formal employment in 2015. The outcome variable in (1) is the percent growth in the number of firms, while in (2) and (3), I decompose the percent growth in terms of the exit and entry of firms, respectively. The sample is restricted to firms with at least one native employee. Firms are observed in August of each year. Source: PILA, 2015–2018.

6 Machine-Learning Approach

In this section of the paper, I develop a machine learning algorithm to identify the subgroups most affected by immigration and to determine a proxy for the role of firms in the labor market effects of immigration. In previous sections, I show wage and employment effects for arbitrarily chosen subgroups of the population according to given characteristics. Yet, to determine exactly which variable explains most of the heterogeneity of immigration effects, I turn to a data-driven approach proposed by [Athey and Imbens \(2016\)](#) and generalized in [Athey et al. \(2019\)](#). Recently, it was implemented by [Gulyas et al. \(2019\)](#) and [Yakymovych et al. \(2022\)](#). This framework identifies the subgroups that experience the greatest wage and employment losses by a recursive partitioning method while allowing for nonlinear effects and high-order interactions between firm and worker variables. The generalized random forest (GRF) method in [Athey et al. \(2019\)](#) builds causal

forests, in the spirit of random forests (Breiman, 2001) but splits the data according to a criterion on treatment effect heterogeneity.⁵⁴ The benchmark specification that the algorithm uses is the following:

$$\Delta Y_{i,l,2018} = \tau(x_i)\Delta \hat{M}_{l,2018} + \Delta \epsilon_{i,l,2018} \quad (13)$$

where x_i are the values of the variables in X_i and $\tau(x_i)$ is the treatment effect. Moreover, $\Delta Y_{i,l,2018}$ is the outcome of interest: the difference in individual employment or wages in 2018 relative to the pre-shock period. $\hat{M}_{l,2018}$ is the predicted immigration rate after a regression of the observed one on the instruments. This is done because the algorithm does not allow for more than one instrument. Vectors of worker and firm variables, including the ones constructed from the AKM model, are the partitioning variables f included in the vector X_f . All these features or variables correspond to characteristics in 2015 (before the immigration shock), and they are age, sex, job tenure, wages, firm FEs, worker FEs, and firm size.⁵⁵ Self-employed workers are omitted in this section due to their incomparable information in most of the firm characteristics to employees.

The procedure in Athey and Imbens (2016) and Athey et al. (2019) to build causal trees consists of several steps that are adapted to this setup. Broadly, the algorithm proceeds as follows:

1. Start with 50% of the full sample P .⁵⁶ The remaining out-of-bag (OOB) sample is used for estimation after the algorithm is trained.
2. Take a random subsample, without replacement, of P and choose a variable randomly from X_f and a value, from all possible values, for this selected variable.
3. For every possible value of one variable in X_f , the data is split into two partitions (say P_l and P_r) to run separate regressions of form (13) to estimate treatment effects for each partition. Choose the variable with its cutoff value that maximizes the difference in treatment effects using this formula:

$$(\tau_l - \tau_r)^2. \quad (14)$$

⁵⁴I use the `grf` package in `R` to estimate the causal forests.

⁵⁵The procedure sample varies depending on the features selected but starts with the same sample. For instance, to construct worker effects, the individual must be observed more than once in the sample, so in this case, the sample is smaller.

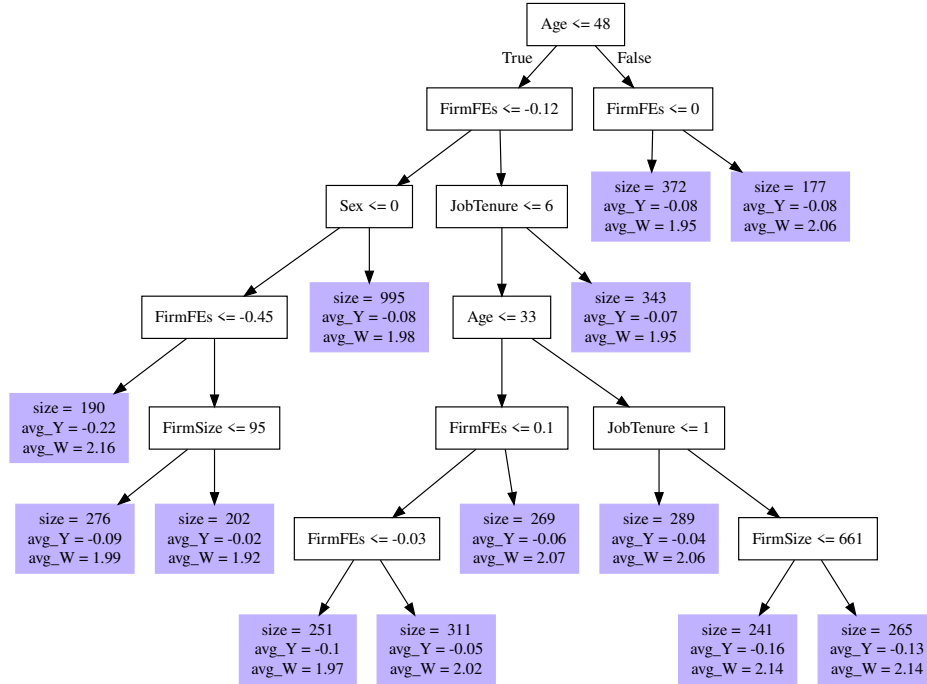
⁵⁶The 50% threshold is selected due to computational burden. This subsample is further cut by 50% to do sub-sample splitting to create similarity matrices.

⁵⁷There are penalties in the algorithm for the imbalance of the splits. For instance, the squared difference criterion can include an additional term $\frac{n_l n_r}{N^2}$ to adjust for more balanced splits (n_l and n_r refer to the sample size of each partition, and total subsample refers to N).

4. Observations with a value below or equal to the cutoff value are placed into a new left node, and observations with a value above are placed into a new right node of the decision tree.
5. Recursively forms the resulting nodes with this algorithm until the nodes reach a minimum node size, the difference in sample size between the two partitions is large, or when the split would only yield a difference in treatment effects relatively small.

As an illustration of a decision tree in the causal forest algorithm, I use a 1% random sample of the main data. Figure 11 shows how observations of certain characteristics are placed to the right and the left of the tree. For the main algorithm, I estimate the causal forest using 2,000 decision trees with a minimum node size of 300, while clustering observations in FUAAs.⁵⁸ Having many trees with a minimum node size reduces overfitting concerns.

Figure 11: **Illustration of decision tree**



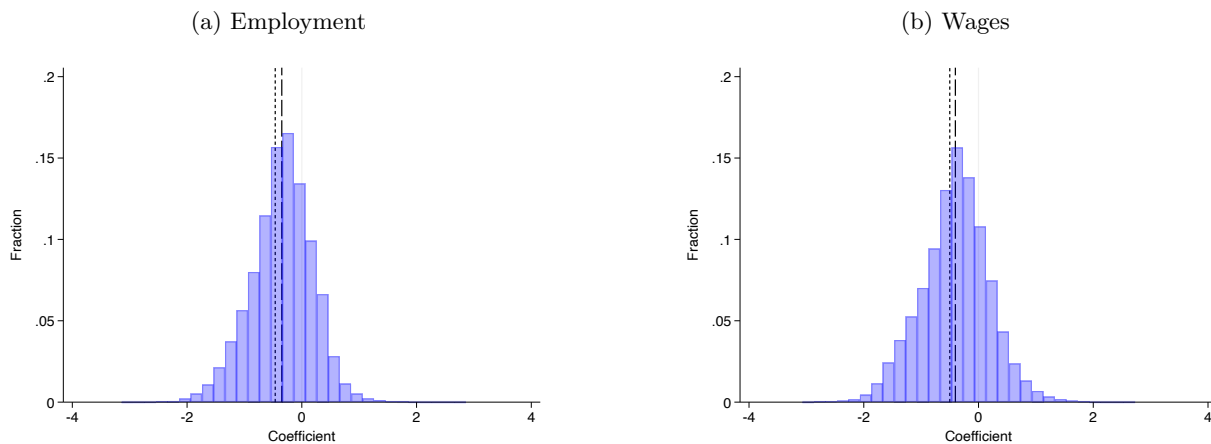
Note: Dependent variable Y is employment changes in 2018 relative to the pre-shock period, and the predicted immigration shock in 2018 is W . This decision tree uses a 1% random sample of the data.

The first output of this procedure is shown in Figures 12a and 12b. According to the trained causal forest, these histograms plot the predicted individual treatment effects for both outcomes, wages, and employment. These treatment effects come from the OOB sample not used in the

⁵⁸I set the tunable parameters from the algorithm to default values, including the honest splitting, while the selected minimum node size is fairly small for precision. In a further cross-validation exercise, results hold when I substantially increase the minimum node size.

algorithm. To estimate the individual treatment effects, each OOB observation is first assigned (according to their characteristics) into a final node of each tree in this forest. Then, for all trained trees, it counts the times these observations fall in the same terminal node as the training sample to calculate the similarity weights. Using these weights, it gets the weighted mean of τ across trees to calculate the individual treatment effect $\tau(x_i)$. In the histograms, the average individual treatment effect is the long dashed line, and the average treatment effect from the standard regression of Equation (1) is the short dashed line. For both outcomes, the average coefficient from the causal forest is similar, reflecting the accuracy of the average prediction.

Figure 12: **Histogram of treatment effects for formal employment and formal wages in the causal forest, 2015–2018**



Note: The short dashed line refers to the coefficient from the benchmark specification, and the long dashed line refers to the average predicted treatment effects that are estimated with the trained causal forest using the OOB sample. The number of trees is 2,000. The sample is restricted to natives between 25 and 55 years old. I use clusters at the FUA level for the causal forest. The causal forest uses 50% of the main sample due to computational burden. The minimum node size is 300.

Next, I exploit the construction of treatment effects from the algorithm to describe which subgroups are most affected by immigration. In this exercise, I divide native workers into quintiles of treatment effects of employment and wages (quintile 1 refers to the most negative effect and quintile 5 to the most positive one). The aim of this exercise is to compare characteristics between groups, not to make inference from the estimated individual treatment effects.

Tables 6a and 6b show worker and firm characteristics in the pre-shock period. First, native workers with the most negative employment effects are the oldest, with the lowest tenure, and earning the lowest initial wages. Besides, these workers are employed in the smallest firms with the

lowest pay premiums. Conversely, workers that suffer the most negative wage effect are relatively younger, with few years of tenure, and earn the highest initial wages. In terms of firm characteristics, these workers are employed in the smallest firms, and in terms of pay premiums, they are in the middle-high part.⁵⁹ From a policy perspective, the distribution of individual treatment effects is useful for targeted measures that aim to decrease the losses from immigration in the most affected subgroups.

Table 6: **Descriptive statistics for native workers by quintiles of treatment effects**

(a) Formal employment					
	Q1	Q2	Q3	Q4	Q5
Male (%)	0.7	0.6	0.5	0.5	0.5
Age of worker	42.8	40.3	38.5	35.1	31.1
Job tenure (1-9 years)	2.3	3.6	4.4	4.1	2.8
Monthly wages (USD)	324.8	462.6	521.8	478.4	336.2
Median firm size	79	105	276	510	1109
Quantiles of firm FEs (1-7)	3.8	5.3	6.0	6.3	6.5

(b) Formal wages					
	Q1	Q2	Q3	Q4	Q5
Male (%)	0.6	0.6	0.6	0.6	0.5
Age of worker	36.6	38.5	38.8	38.1	37.5
Job tenure (1-9 years)	3.2	3.9	4.0	3.8	3.5
Monthly wages (USD)	559.5	466.2	419.3	379.0	393.7
Median firm size	86	189	242	309	892
Quantiles of firm FEs (1-7)	5.7	5.8	5.6	5.5	5.5

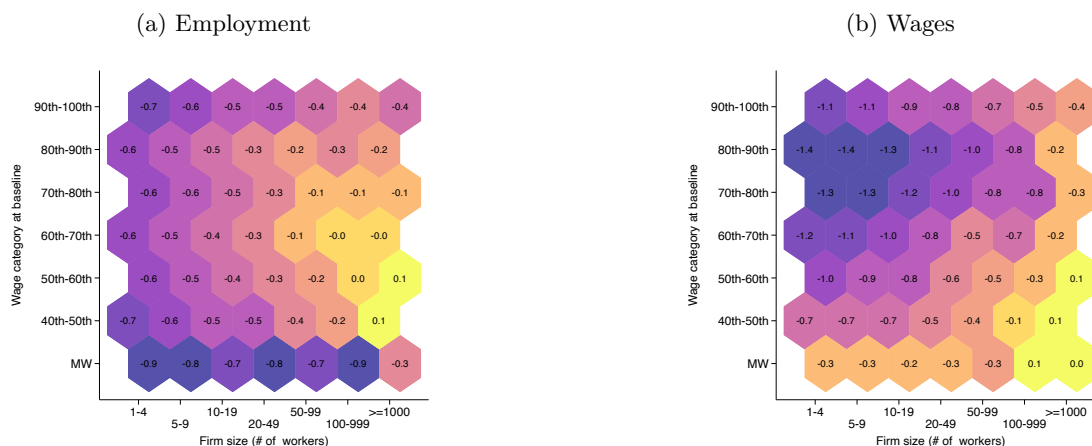
Note: These tables report the average or median statistics for quintiles of treatment effects (Q1 is the most affected and Q5 is the least affected) in terms of employment and wages, according to the predictions of the trained causal forest using the OOB sample. The wages are transformed from Colombian pesos to USD using 2020 exchange rates from the World Bank. Source: PILA, August 2015.

Following up, and to have a better illustration of the subgroups most affected by immigration, I construct heat plots. Figures 13a and 13b show the average of individual treatment effects by individual wages at baseline for different firm sizes. The idea here is to do a relative comparison of effects across these two dimensions. Interestingly, most negative employment effects are concentrated on the intersection of minimum wage earners employed in small and medium firms. Opposite from the most negative wage effects, which are concentrated in the upper part of the wage

⁵⁹In Appendix Figures C.3a and C.3a, I test if the quintiles of treatment effects from the causal forest yield the same order when using the main empirical specification. Importantly, the estimates follow the same order as the quintiles for wages and employment.

distribution, but again in small firms (negative wage effects are smoothly disappearing as the firm is larger).

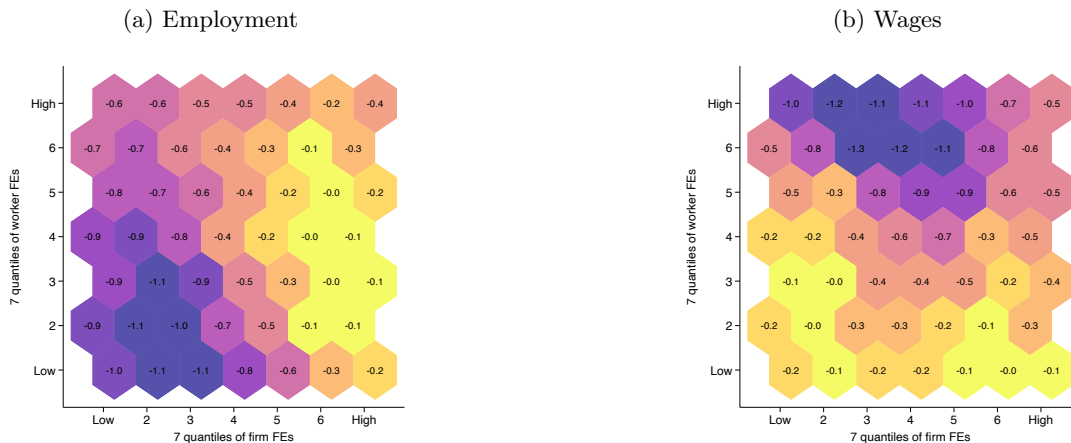
Figure 13: **Heat plot of treatment effects by wage category at baseline and firm size, 2015–2018**



Note: Each hexagon is the average of individual treatment effect in the subgroup according to the trained causal forest using the OOB sample. The sample is restricted to natives between 25 and 55 years old. I use clusters at the FUA level for the causal forest. The causal forest uses 50% of the main sample due to computational burden.

Next, Figures 14a and 14b show average treatment effects by quantiles of firm FEs intersected with quantiles of worker FEs. Interestingly, most negative employment effects are concentrated on the low-wage workers in the lowest-paying firms. Opposite from the most negative wage effects, which tend to be concentrated in high-wage workers in middle-paying firms.

Figure 14: **Heat plot of treatment effects by quantiles of workers and firm FEs, 2015–2018**

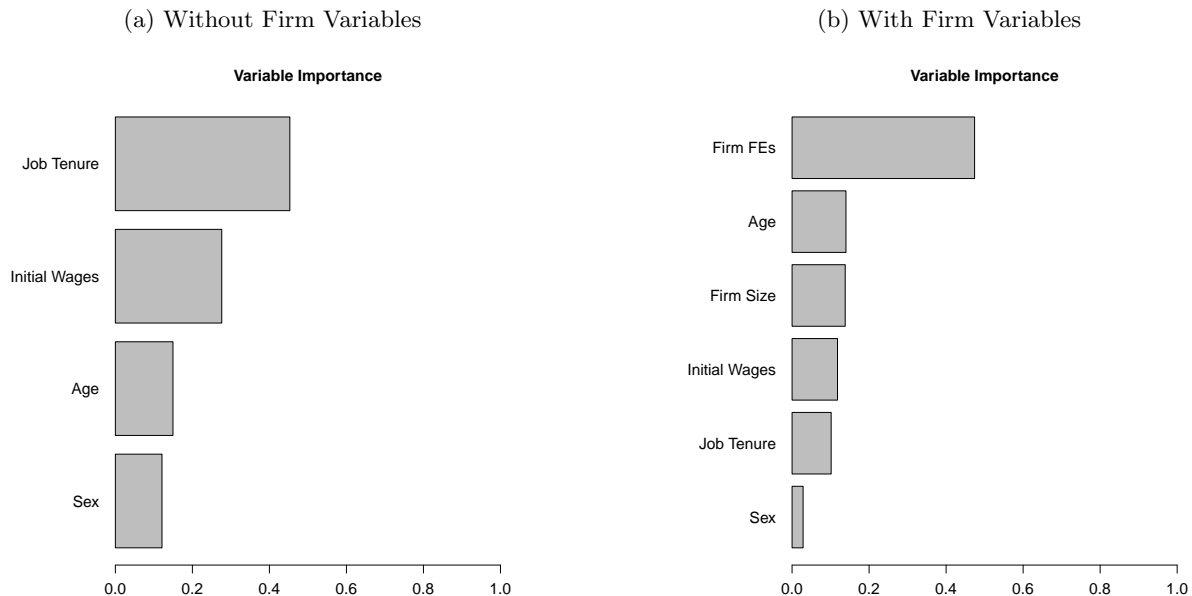


Note: Each hexagon is the average of individual treatment effect in the subgroup according to the trained causal forest using the OOB sample. The sample is restricted to natives between 25 and 55 years old. I use clusters at the FUA level for the causal forest. The causal forest uses 50% of the main sample due to computational burden.

A complementary way of summarizing these findings is with the variable importance measure. In this case, the variables that appear more frequently as splits in the forest are categorized as more important to explain treatment effect heterogeneity. This naive measure yields a ranking that serves as a proxy to classify the sources of heterogeneity. To start, I perform the algorithm excluding and including firms' variables to show how different is the importance measure.⁶⁰ Hence, when excluding firms' variables, I find that job tenure, followed by initial wages and age, are more important to determine the heterogeneity on employment impacts of migration (see Figure 15a). However, when including firms' variables, the most important variable becomes firm-specific pay premiums or firm FEs, followed by age and firm size. Note that this measure does not indicate the sign or magnitude of the effect of each variable on employment, only that it helps to explain most of the heterogeneity in treatment effects. Thus, the relevance of firms for the heterogeneity of immigration effects on natives is notable. Like the findings of [Arellano-Bover and San \(2020\)](#), that shows the important role of firms in the assimilation of immigrants in the labor market.

⁶⁰For employment, I use the individual change in employment between 2018 and the average pre-shock period employment as the outcome.

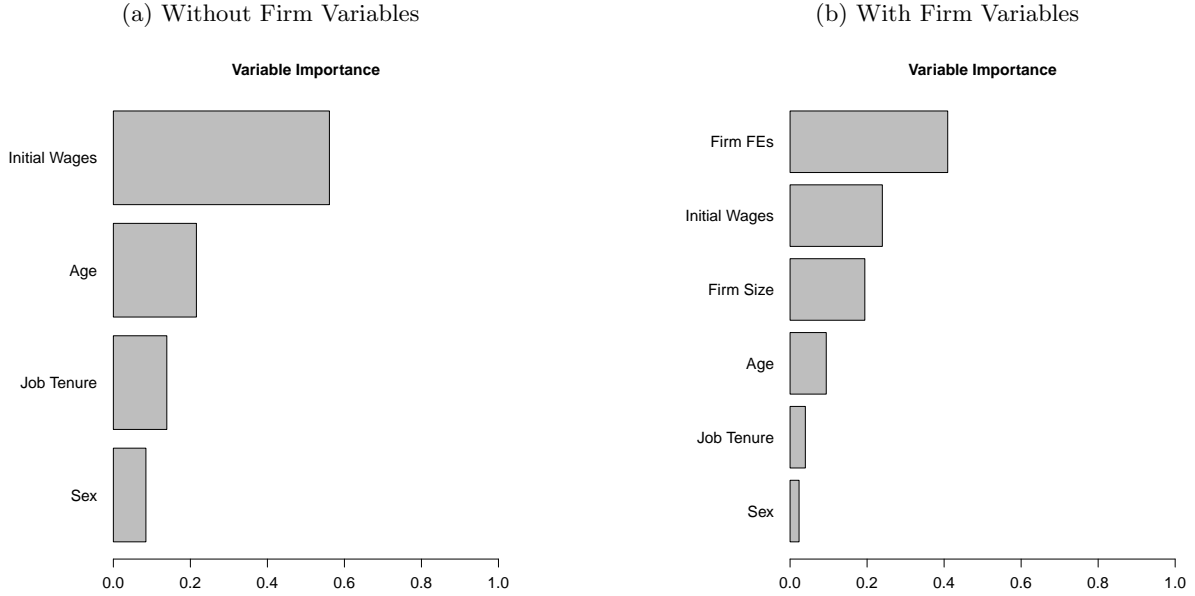
Figure 15: **Variable importance for formal employment in the causal forest, 2015–2018**



Note: Variable importance is a weighted sum of how many times the feature f appears in the split of each leaf of every tree in the forest. The number of trees is 2,000. The sample is restricted to natives between 25 and 55 years old. The importance measure sums up to 1. I use clusters of FUA for the causal forest estimation. The minimum node size is 300.

Following up, I use the individual wage growth between 2018 and 2015 to perform the same exercise. Without firms' variables, the most important variables are initial wages followed by age and job tenure (see Figure 16a). However, when including firms' variables in the causal forest, firm-specific pay premiums followed by initial wages and firm size are the most important (see Figure 16b). Note that the variables of firm size and pay premiums are positively correlated, but the correlation is not so strong (0.19). Conversely, the variables that explain the least are job tenure and sex. To summarize, the most important variable to explain wage and employment changes relates to the firm-specific pay premiums or firm FEs more than any worker characteristics. In the causal forest of wages, firm effects appear in 37% of all splits; for employment, firm effects appear in 30% of the splits.

Figure 16: **Variable importance for formal wages in the causal forest, 2015–2018**



Note: Variable importance is a weighted sum of how many times the feature f appears in the split of each leaf of every tree in the forest. The number of trees is 2,000. The sample is restricted to natives between 25 and 55 years old. The importance measure sums up to 1. I use clusters of FUA for the causal forest estimation. The minimum node size is 300.

Lastly, as initial wages are a function of the unobserved firm and worker FEs, the next result I show includes in the algorithm the constructed worker FEs $\hat{\alpha}_i$ instead of initial wages. This reduces the sample as every worker must be observed more than once. After adding worker FEs, again, the firm FEs are the most important variable to explain the heterogeneity of treatment effects for employment and wages (see Appendix Figures B.6a and B.6b). In conclusion, firms' role in the impact of immigration is still very relevant even after conditioning with the constant quality of workers.

7 Robustness Checks

To start, the exclusion restriction of the distance instrument can fail due to border areas as they are more prone to be affected by other time-varying shocks arising from the Venezuelan crisis. Therefore, I remove all the border areas from the estimation sample to find similar point estimates but not significant for wages (see Appendix Table C.1, row 2). Next, another concern is the relevance of Bogotá as the capital of Colombia (the proportion of observations from the capital

is 32.7% of the whole sample). Hence, I also remove it from the estimation sample and find that coefficients are less negative, especially for wages, but both are significant (see Appendix Table C.1, row 3).

Next, I add further controls to the regression to compare more accurately workers across local labor markets. The additional controls are seven groups of wage categories, according to the local wage distribution in 2015, and job tenure. Reassuringly, results are similar for wages but much less negative for employment, mainly because self-employed workers are omitted from the analysis as there is no comparable measure of job tenure. The next robustness test is the adjustment of nominal wages to real terms using the national CPI. In this case, the results of wages are slightly less negative. Last, to omit outliers driving the wage results, I top code wages after the 99% percentile of the wage distribution to find that estimates are unaltered.

Third, I perform robustness checks for the machine learning algorithm. The first one is that firm pay-premiums might be correlated with the type of industry the firm belongs to, reflecting that some industries generally have higher or lower wage premia (Card et al., 2022). For that reason, I include in the algorithm the industry of the firm, along with the firm FEs, to find that for wages and employment, the most important variable is still the firm FEs (see Appendix Figures C.5a and C.5a). The second one deals with the fact that the frequency of splits in the first nodes of the trees is weighted the same as the frequency of splits in the last nodes of the tree, where the sample size is much smaller. This critique is alleviated by using a decay exponent in the variable importance that puts more weight on the splits selected first.⁶¹ After computing the variable importance, the order is fairly similar for wages and for employment firm size is now the second most important variable, preceded by firm FEs. Interestingly, both variables capture the role of firms (see Appendix Figures C.4a and C.4a).

Next is that, in the causal forests, the number of possible values a variable takes might alter the variable importance weighted sum (Strobl et al., 2007). For instance, when variables have a small set of values, they might mechanically appear in fewer nodes further in the tree. For that reason, Appendix Figures C.6a and C.6b show that when transforming all the continuous variables into seven or six categories, as the ones in previous results, the order of importance is similar for

⁶¹The decay exponent is -2, meaning that split frequencies in node k are weighted $1/2$ compared to those in node $k - 1$.

employment but for wages changes slightly, with firm FEs being second. Still, even if the main results hold, one of the benefits of the algorithm comes from exploiting all the possible values a variable takes for allowing non-linear and interaction effects, not from arbitrarily aggregating into categories. Another critique is that the tree is built to maximize the squared difference in treatment effects without analyzing whether pre-trends are significant for all these subgroups, it just assumes strict exogeneity of the instrument. Probably when the treatment effects are higher, there could be differing pre-trends. However, as the algorithm is constructed, it does not allow correcting or checking for pre-trends in every possible subgroup, so what I do is check for pre-trends in several subgroups in Appendix F to show insignificant estimates in the majority of categories. Finally, there is a recent statistical literature that proposes hypothesis testing of variable importance measures of random forests (see, for instance, [Hapfelmeier et al. \(2023\)](#)). The main idea is to perform sequential permutation tests to get the p-value of each variable used in the algorithm. The main issue is that it has not been developed for causal forests, and even if available, it is inefficient in this high-dimensional setup and computationally infeasible.

8 Conclusion

This is the first paper that exploits the labor supply shock of immigrants from Venezuela equipped with data covering the universe of formal workers and firms in Colombia. This is an advantage in several dimensions. First, with administrative panel data, I can follow workers over time and address compositional changes that arise in the standard regional-level analysis of immigration using survey data. Second, with the matched employee-employer dimension of the data, I can uncover heterogeneity across firm and worker characteristics that help to understand additional mechanisms after labor supply shocks in developing countries. Third, with the full count of formal firms and a machine learning method, I can construct a proxy measure for the role of firms in the impact of immigration on workers.

Overall, the findings suggest that after immigrants arrive, there is a negative impact on individual employment in the formal sector. However, this coefficient masks many heterogeneous responses. Specifically, minimum-wage workers are crowded out from the formal sector, while workers above in the wage distribution are not displaced, but instead, they have negative wage growth. Regarding

firm characteristics, the negative effect on employment and wages is concentrated in small firms, which aligns with the predictions from the model that predicts that firms that use more informal work in production are more affected regarding formal employment and wages. Next, I find that workers in middle-paying firms present a negative wage impact, while workers in low-paying firms do not experience a reduction in wages but more on employment, partly due to the existence of a relatively high minimum wage that prevents further wage cuts.

To uncover the mechanisms behind these impacts, I then use causal forests to classify which variable is most important to explain the heterogeneity in treatment effects. Throughout this analysis, firm-specific pay premiums appear prominently as the most important variable to explain heterogeneity in wage and employment effects, followed by firm size in most cases. Thus, using only workers' characteristics when analyzing the labor market impacts of immigration might lead to an incomplete view of the sources of adjustments to immigration. Suggesting that after immigrants arrive, the focus should not be, at least for Colombia, in *who* the worker is but more on which *type of firm* the worker is employed.

References

- Abowd, J. M., Kramarz, F., and Margolis, D. N. (1999). High wage workers and high wage firms. *Econometrica*, 67(2):251–333.
- Aksu, E., Erzan, R., and Kırdar, M. G. (2018). The impact of mass migration of Syrians on the Turkish labor market. Technical report, Working Paper.
- Altonji, J. G. and Card, D. (1991). The effects of immigration on the labor market outcomes of less-skilled natives. In *Immigration, trade, and the labor market*, pages 201–234. University of Chicago Press.
- Amior, M. and Stuhler, J. (2022). Immigration and monopsony: Evidence across the distribution of firms. Technical report, Working paper.
- Andrews, M. J., Gill, L., Schank, T., and Upward, R. (2008). High wage workers and low wage firms: negative assortative matching or limited mobility bias? *Journal of the Royal Statistical Society Series A: Statistics in Society*, 171(3):673–697.
- Arellano-Bover, J. and San, S. (2020). The role of firms in the assimilation of immigrants. *Available at SSRN 3594778*.
- Athey, S. and Imbens, G. (2016). Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113(27):7353–7360.
- Athey, S., Tibshirani, J., and Wager, S. (2019). Generalized random forests. *The Annals of Statistics*, 47(2):1148–1178.
- Autor, D. H., Dorn, D., Hanson, G. H., and Song, J. (2014). Trade adjustment: Worker-level evidence. *The Quarterly Journal of Economics*, 129(4):1799–1860.
- Aydemir, A. and Borjas, G. J. (2011). Attenuation bias in measuring the wage impact of immigration. *Journal of Labor Economics*, 29(1):69–112.
- Bagger, J. and Lentz, R. (2019). An empirical model of wage dispersion with sorting. *The Review of Economic Studies*, 86(1):153–190.
- Bahar, D., Ibáñez, A. M., and Rozo, S. V. (2021). Give me your tired and your poor: Impact of a large-scale amnesty program for undocumented refugees. *Journal of Development Economics*, 151:102652.
- Bonhomme, S., Holzheu, K., Lamadon, T., Manresa, E., Mogstad, M., and Setzler, B. (2020). How much should we trust estimates of firm effects and worker sorting? Technical report, National Bureau of Economic Research.
- Bonhomme, S., Lamadon, T., and Manresa, E. (2019). A distributional framework for matched employer employee data. *Econometrica*, 87(3):699–739.
- Borjas, G. J. (2006). Native internal migration and the labor market impact of immigration. *Journal of Human resources*, 41(2):221–258.
- Borjas, G. J. and Edo, A. (2021). Gender, selection into employment, and the wage impact of immigration. Technical report, National Bureau of Economic Research.
- Bratsberg, B. and Raaum, O. (2012). Immigration and wages: Evidence from construction. *The economic journal*, 122(565):1177–1205.
- Breiman, L. (2001). Random forests. *Machine learning*, 45(1):5–32.
- Card, D. (2001). Immigrant inflows, native outflows, and the local labor market impacts of higher immigration. *Journal of Labor Economics*, 19(1):22–64.
- 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., Heining, J., and Kline, P. (2013). Workplace heterogeneity and the rise of west german wage inequality. *The Quarterly journal of economics*, 128(3):967–1015.
- Card, D., Rothstein, J., and Yi, M. (2022). Industry wage differentials: A firm-based approach.

Unpublished draft, University of California, Berkeley.

- Caruso, G., Canon, C. G., and Mueller, V. (2021). Spillover effects of the venezuelan crisis: migration impacts in colombia. *Oxford Economic Papers*, 73(2):771–795.
- Ceritoglu, E., Yunculer, H. B. G., Torun, H., and Tumen, S. (2017). The impact of syrian refugees on natives? labor market outcomes in turkey: evidence from a quasi-experimental design. *IZA Journal of Labor Policy*, 6(1):1–28.
- Corbi, R., Ferraz, T., and Narita, R. (2021). Internal migration and labor market adjustments in the presence of nonwage compensation.
- DANE (2019). Gran Encuesta Integrada de Hogares (GEIH)-Módulo de Migración. https://microdatos.dane.gov.co/index.php/catalog/641/get_microdata.
- Del Carpio, X. V. and Wagner, M. C. (2015). The impact of Syrian refugees on the Turkish labor market. *World Bank policy research working paper*, (7402).
- Delgado-Prieto, L. (2022). Immigration, wages, and employment under informal labor markets.
- Doran, K., Gelber, A., and Isen, A. (2022). The effects of high-skilled immigration policy on firms: Evidence from visa lotteries. *Journal of Political Economy*, 130(10):2501–2533.
- Dostie, B., Li, J., Card, D., and Parent, D. (2021). Employer policies and the immigrant–native earnings gap. *Journal of Econometrics*.
- Dustmann, C., Lindner, A., Schönberg, U., Umkehrer, M., and Vom Berge, P. (2022). Reallocation effects of the minimum wage. *The Quarterly Journal of Economics*, 137(1):267–328.
- Dustmann, C., Otten, S., Schönberg, U., and Stuhler, J. (2023). The Effects of Immigration on Places and Individuals - Identification and Interpretation. Technical report.
- Dustmann, C., Schönberg, U., and Stuhler, J. (2017). Labor supply shocks, native wages, and the adjustment of local employment. *The Quarterly Journal of Economics*, 132(1):435–483.
- Foged, M. and Peri, G. (2016). Immigrants’ effect on native workers: New analysis on longitudinal data. *American Economic Journal: Applied Economics*, 8(2):1–34.
- Fort, T. C., Haltiwanger, J., Jarmin, R. S., and Miranda, J. (2013). How firms respond to business cycles: The role of firm age and firm size. *IMF Economic Review*, 61(3):520–559.
- Groeger, A., León-Ciliotta, G., and Stillman, S. (2022). Immigration, labor markets and discrimination.
- Gulyas, A., Pytko, K., et al. (2019). Understanding the sources of earnings losses after job displacement: A machine-learning approach. Technical report, University of Bonn and University of Mannheim, Germany.
- Hapfelmeier, A., Hornung, R., and Haller, B. (2023). Efficient permutation testing of variable importance measures by the example of random forests. *Computational Statistics & Data Analysis*, page 107689.
- Hicks, J. (1932). The theory of wages, london: Macmillan.
- Hoen, M. F. (2020). Immigration and the tower of babel: Using language barriers to identify individual labor market effects of immigration. *Labour Economics*, 65:101834.
- Kleemans, M. and Magruder, J. (2018). Labour market responses to immigration: Evidence from internal migration driven by weather shocks. *The Economic Journal*, 128(613):2032–2065.
- Kline, P., Saggio, R., and Sølvssten, M. (2020). Leave-out estimation of variance components. *Econometrica*, 88(5):1859–1898.
- Kuosmanen, I. and Meriläinen, J. (2022). Labor market effects of open borders: Evidence from the finnish construction sector after eu enlargement. *Journal of Human Resources*.
- Lebow, J. (2021). Immigration and occupational downgrading in colombia. *Unpublished Manuscript*.
- McKenzie, D. (2017). Identifying and spurring high-growth entrepreneurship: Experimental evidence from a business plan competition. *American Economic Review*, 107(8):2278–2307.
- Monras, J. (2020). Immigration and wage dynamics: Evidence from the mexican peso crisis. *Journal*

- of *Political Economy*, 128(8):3017–3089.
- Morales-Zurita, L. F., Bonilla-Mejía, L., Hermida, D., Flórez, L. A., Bonilla-Mejía, L., Morales, L. F., Hermida-Giraldo, D., and Flórez, L. A. (2020). The labor market of immigrants and non-immigrants evidence from the venezuelan refugee crisis. *Borradores de Economía; No. 1119*.
- Muñoz, M. (2021). Trading non-tradables: The implications of europe’s job posting policy.
- Orefice, G. and Peri, G. (2020). Immigration and worker-firm matching. Technical report, National Bureau of Economic Research.
- Ortega, J. and Verdugo, G. (2022). Who stays and who leaves? immigration and the selection of natives across locations. *Journal of Economic Geography*, 22(2):221–260.
- Redondo, H. (2022). From bricklayers to waiters: Reallocation in a deep recession.
- Sanchez-Serra, D. (2016). Functional urban areas in colombia.
- Strobl, C., Boulesteix, A.-L., Zeileis, A., and Hothorn, T. (2007). Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC bioinformatics*, 8(1):1–21.
- Ulyssea, G. (2018). Firms, informality, and development: Theory and evidence from brazil. *American Economic Review*, 108(8):2015–47.
- UNHCR (2019). Venezuela situation. Technical report, Office of the United Nations High Commissioner for Refugees. Available on: <https://www.unhcr.org/venezuela-emergency.html> (last accessed 13th of June 2020).
- Yagan, D. (2019). Employment hysteresis from the great recession. *Journal of Political Economy*, 127(5):2505–2558.
- Yakymovych, Y., Nordström Skans, O., Vikström, J., Simon, L., and Athey, S. (2022). Worker attributes, aggregate conditions and the impact of adverse labor market shocks.

Online Appendix

A First stage of the instruments

Table A.1: **First stage: The inflow of Venezuelan immigrants and the two instruments**

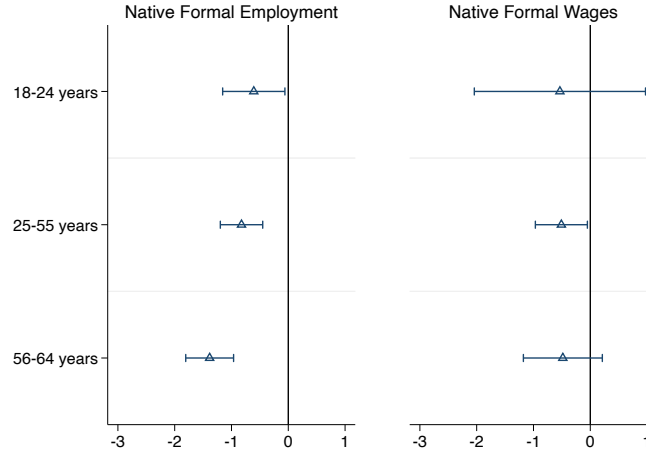
	(1)	(2)	(3)
	$\Delta M_{l,2018}$	$\Delta M_{l,2018}$	$\Delta M_{l,2018}$
Distance (/100)	-1.992*** (0.272)		-1.455*** (0.350)
Distance (/100) squared	0.151*** (0.024)		0.107*** (0.029)
Past settlements		0.703*** (0.160)	0.280* (0.130)
Constant	6.762*** (0.715)	1.040*** (0.149)	5.184*** (1.000)
R^2	0.583	0.450	0.618
F	34.53	19.37	23.68
N	109	109	109

Standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Note: This table reports the coefficient of the first-stage of the share of employed migrants $\Delta M_{l,2018} * 100$ with distance and distance squared to the nearest crossing bridge and past settlements as explanatory variables.

B Additional Results

Figure B.1: Estimates by extended age categories, 2015–2018



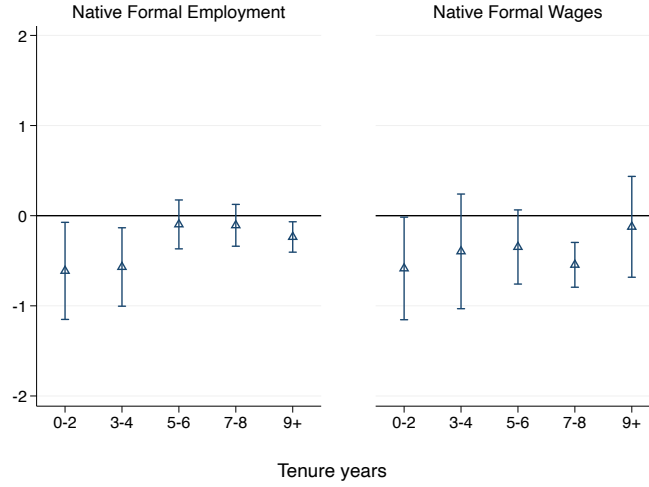
Note: I estimate equation (1) separately by subgroups. The sample is restricted to natives between 18 and 64 years old. Dependent variables are employment relative to the pre-shock period and wages relative to the base period. I use as controls sex with a dummy for self-employed in the base period. I cluster standard errors ($G=109$). 95% confidence interval. The coefficients for employment (in percentage points) and wages (in percent) are already multiplied by 100. Workers are observed in August of each year. Source: PILA, 2013–2019.

Table B.1: Labor supply elasticities by age group

Age group	25-30	30-35	35-40	40-45	45-50	50-55
η_w^s	0.42	0.99	1.39	1.67	2.76	3.67

Note: The elasticity of labor supply is given by the reduced-form results from changes in native employment over changes in native wages in 2018.

Figure B.2: Estimates by job tenure, 2015–2018



Note: I estimate equation (1) separately by subgroups. The sample is restricted to native employees between 25 and 55 years old. Dependent variables are employment and wages relative to the base period. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). 95% confidence interval. The coefficients for employment (in percentage points) and wages (in percent) are already multiplied by 100. Workers are observed in August of each year. Source: PILA, 2013–2019.

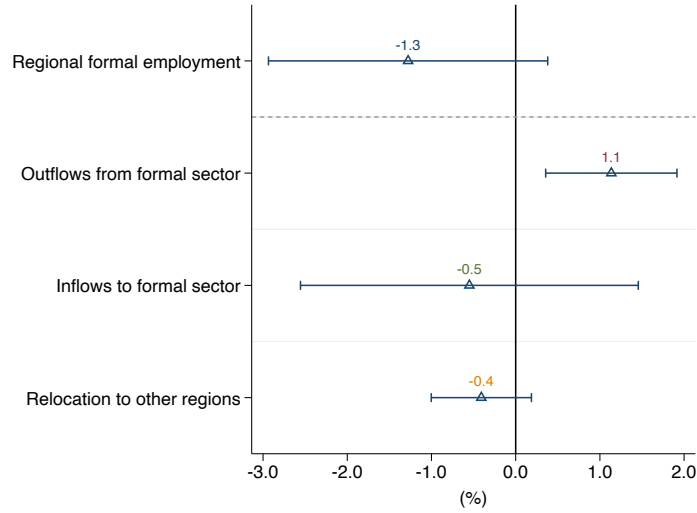
Table B.2: Employment and wage estimates by age and job tenure, 2015–2018

Worker's age	Below 35 years		Above 35 years	
Job tenure	0 to 4 years	5 to 9+ years	0 to 4 years	5 to 9+ years
Prob. of Employment	-0.138 (0.195)	0.209 (0.226)	-1.009** (0.315)	-0.302*** (0.086)
<i>N</i>	2,099,147	344,156	2,075,913	1,083,435
Wages	-0.479 (0.344)	-0.664* (0.279)	-0.556 (0.354)	-0.194 (0.182)
<i>N</i>	1,094,691	240,058	1,170,322	785,839
Clusters	109	109	109	109

Standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Note: I estimate equation (1) separately by subgroups. The sample is restricted to native employees between 25 and 55 years old. Dependent variables are employment relative to the pre-shock period and wages relative to the base period. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). Workers are observed in August of each year. Source: PILA, 2013–2018.

Figure B.3: **Decomposition of formal employment, 2015–2018**



Note: Regressions are estimated at the regional level for 109 FUA's weighted by their formal employment in 2015. 95% confidence interval. The sample is not restricted by age groups. Regional formal employment is decomposed into outflows from formal employment in that region, inflows from non-employment or the informal sector, employed people in other regions, and relocation of formal workers to other regions. Source: PILA, 2015–2018.

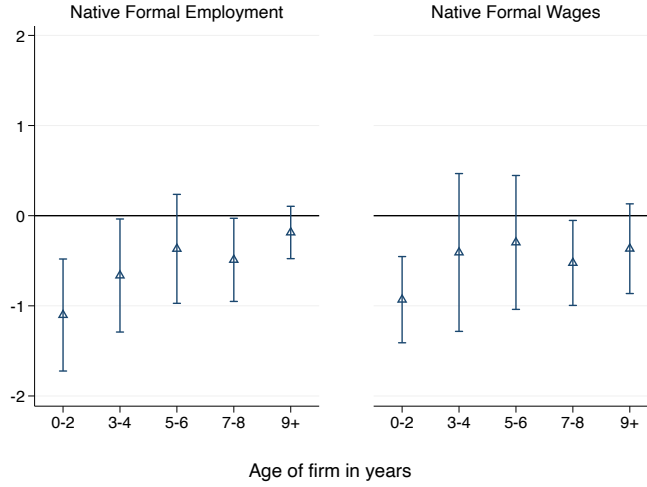
Table B.3: **IV estimates on regional changes of formal workers by age group, 2015–2018**

Age group	25-30	30-35	35-40	40-45	45-50	50-55
Prob. of changing region	0.200	0.088	-0.035	-0.156	-0.211	-0.254
	(0.400)	(0.404)	(0.354)	(0.307)	(0.266)	(0.209)
<i>N</i>	1,255,301	1,041,726	873,437	732,208	674,945	561,949
Clusters	109	109	109	109	109	109

Standard errors are in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Note: The outcome variable is an indicator that takes value one for workers that changed region in 2018 relative to 2015, and zero otherwise. The sample is restricted to natives between 25 and 55 years old. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). 95% confidence interval. The PILA had a measurement error with the regional variable in 2018, so the worker's location in February 2020 (when the health ministry started to verify this information) is used for the workers who present this error. Workers are observed in August of each year. Source: PILA, 2015–2018.

Figure B.4: **Estimates by age of firm, 2015–2018**



Note: I estimate equation (1) separately by subgroups. The sample is restricted to native employees between 25 and 55 years old. Dependent variables are employment relative to the pre-shock period and wages relative to the base period. The firm's age is the number of years the firm appears discontinuously in PILA. I use as controls the interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors (G=109). 95% confidence interval. The coefficients for employment (in percentage points) and wages (in percent) are already multiplied by 100. Workers are observed in August of each year. Source: PILA, 2013–2019.

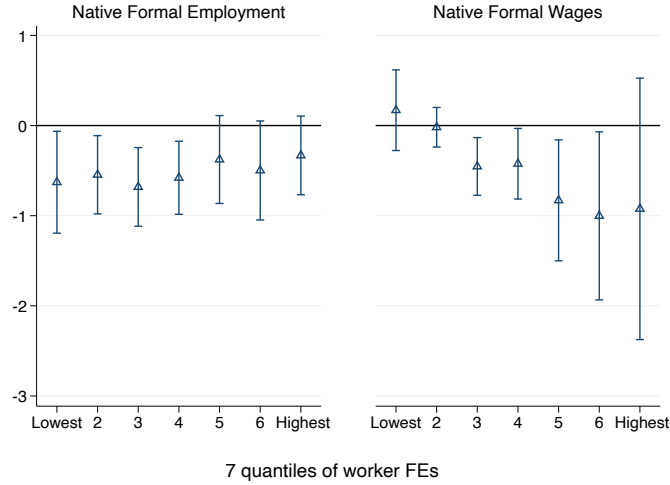
Table B.4: **Employment and wage estimates by firm size and age of firm, 2015–2018**

Firm's size	1 to 19 workers		Above 19 workers	
Age of firm	0 to 4 years	5 to 9+ years	0 to 4 years	5 to 9+ years
Prob. of Employment	-0.762** (0.279)	-0.757*** (0.156)	-1.015** (0.347)	-0.176 (0.170)
N	479,715	498,842	923,272	3,700,822
Wages	-1.021* (0.432)	-0.554 (0.305)	-0.603* (0.304)	-0.395 (0.286)
N	274,728	352,015	444,586	2,219,581
Clusters	109	109	109	109

Standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

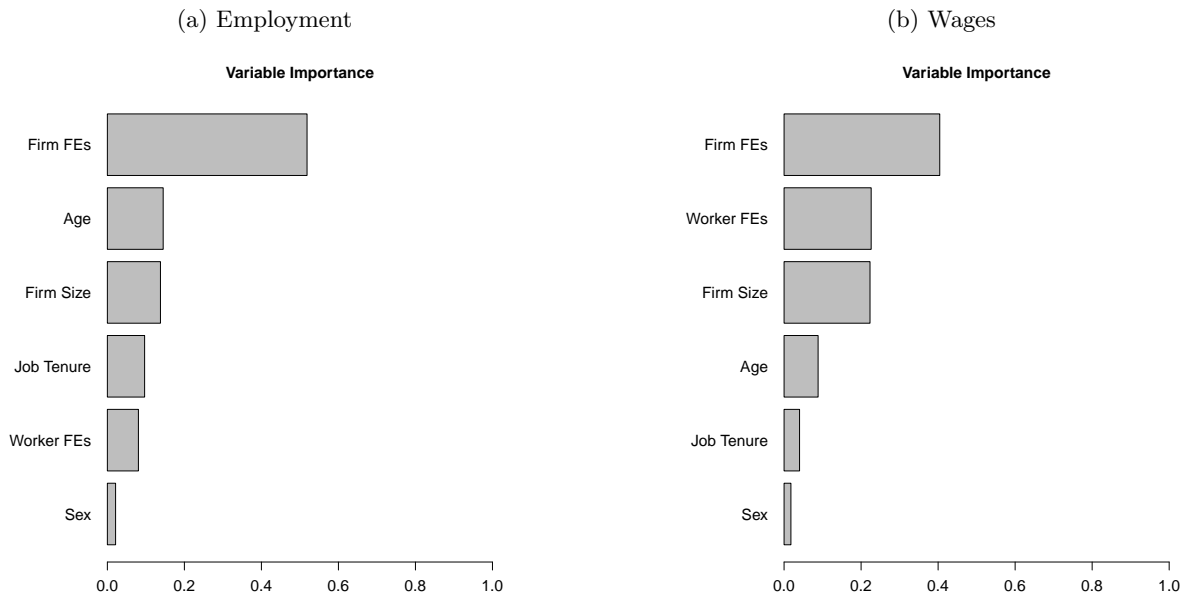
Note: I estimate equation (1) separately by subgroups. The sample is restricted to native employees between 25 and 55 years old. Dependent variables are employment relative to the pre-shock period and wages relative to the base period. The age of the firm is the number of years the firm appears discontinuously in PILA. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors (G=109). Workers are observed in August of each year. Source: PILA, 2013–2018.

Figure B.5: Estimates by quantiles of worker FEs, 2015–2018



Note: I estimate equation (1) separately by subgroups. The sample is restricted to native employees between 25 and 55 years old who appear more than once in PILA. Dependent variables are employment relative to the pre-shock period and wages relative to the base period. I compute Worker FEs in the first stage using the standard AKM framework, with age squared and its cubic as time-varying controls, for the period 2010–2015. I use as controls in the second stage are interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). The coefficients for employment (in percentage points) and for wages (in percent) are already multiplied by 100. Workers are observed in August of each year. 95% confidence interval. Source: PILA, 2013–2019.

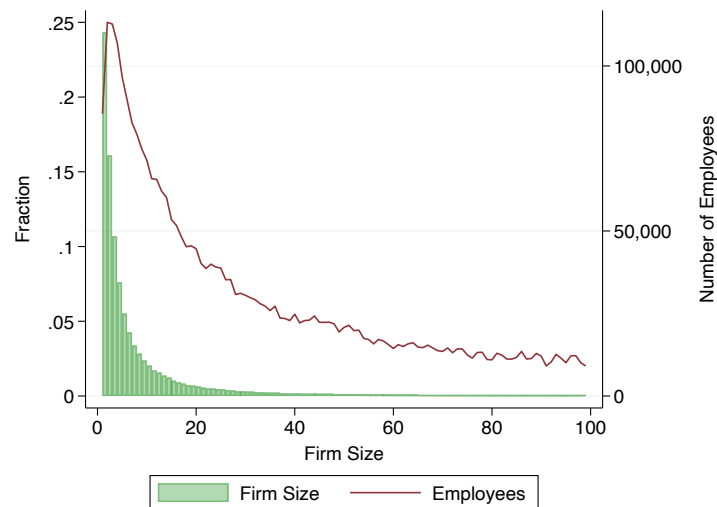
Figure B.6: Variable importance for formal employment and formal wages in causal forest with worker and firm FEs, 2015–2018



Note: Variable importance is a weighted sum of how many times the feature f appears in the split of each leaf of every tree in the forest. The number of trees is 2,000. The sample is restricted to natives between 25 and 55 years old. The importance measure sums up to 1. I use clusters for the causal forest estimation. The minimum node size is 300.

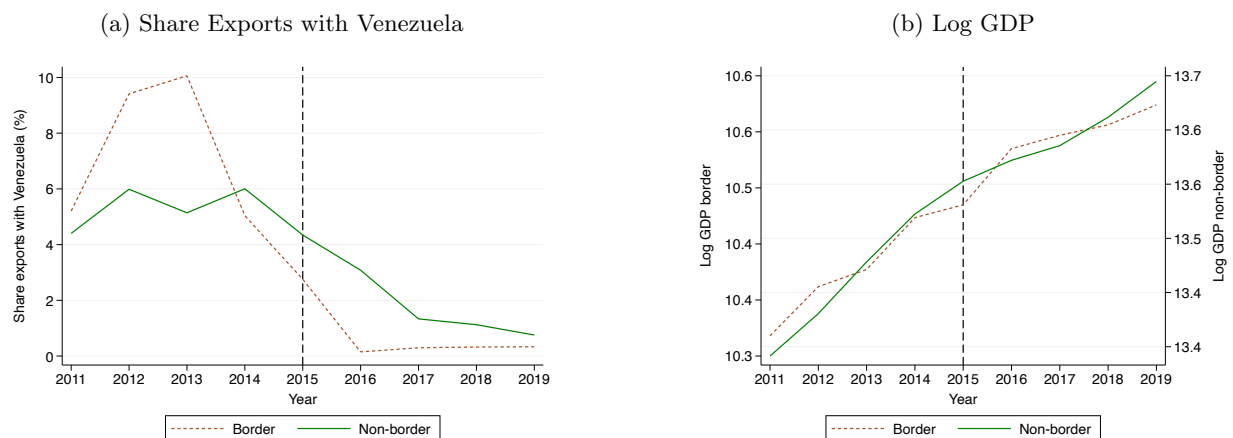
C Robustness Checks

Figure C.1: Firm size distribution and total employees



Note: The upper bound of firm size is restricted to 100 workers for the figure. The chosen bin width is 1. Only workers who contribute as employees are taken into account. Source: PILA, August 2015.

Figure C.2: Evolution of trade and GDP for border and non-border departments



Note: Border departments are *Norte de Santander*, *La Guajira*, and *César*. Non-border departments are the rest. Source: Panel (a) Exportaciones-DANE, 2013–2019. Panel (b) DANE-Cuentas Nacionales, 2011–2019.

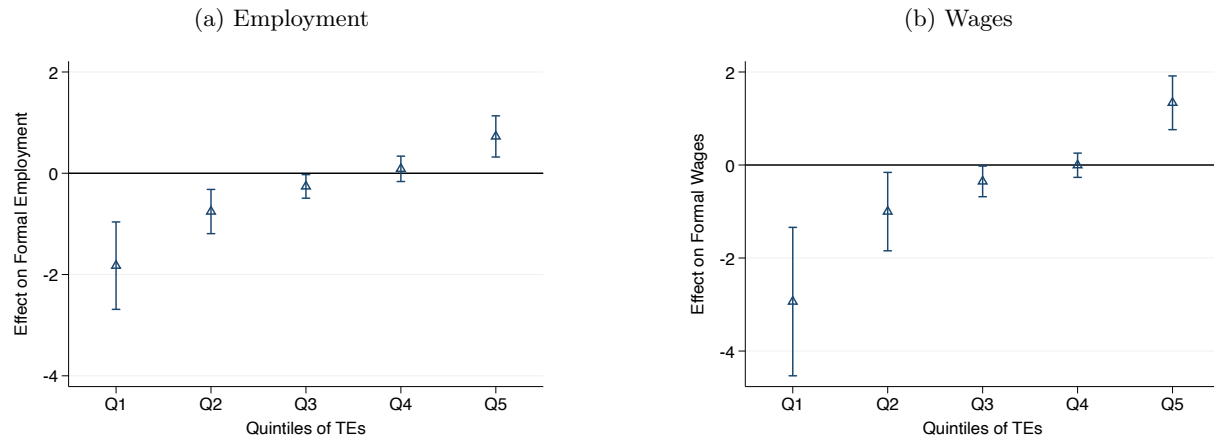
Table C.1: **Robustness checks for formal wages and formal employment, 2015–2018**

	Employment	Wages
Baseline	-0.841*** (0.192)	-0.600* (0.239)
<i>N</i>	6,706,035	4,090,973
Removing border areas*	-1.019* (0.414)	-0.768 (0.559)
<i>N</i>	6,577,923	4,015,648
Removing Bogotá	-0.777*** (0.180)	-0.470** (0.173)
<i>N</i>	4,338,192	2,619,237
Further controls★	-0.484* (0.191)	-0.556 (0.286)
<i>N</i>	4,884,993	3,217,398
Real wages		-0.520* (0.207)
<i>N</i>		4,090,973
Top code local wages above 99%		-0.605* (0.241)
<i>N</i>		4,090,973
Clusters	109	109

Standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

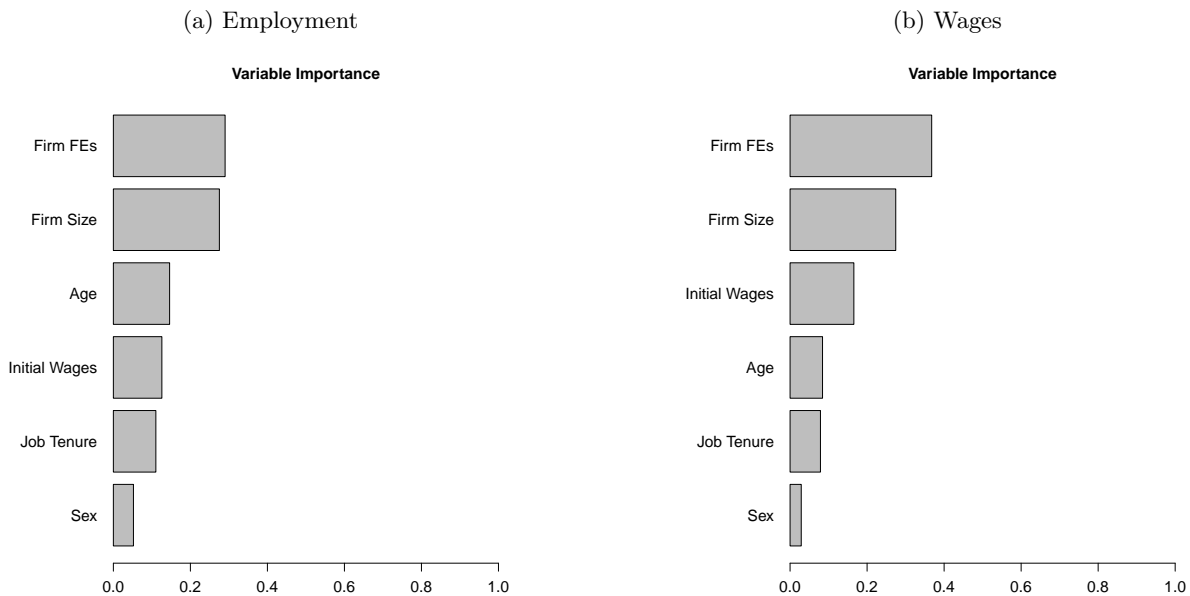
Note: This Table reports the coefficients of the second-stage regression of the instruments with the immigration rate $\Delta M_{i,2018}$. The outcome is the difference with the base period. Controls used are interactions of sex with six age categories and a dummy for self-employed in the base period. *The border areas are Cucutá, Maicao and Arauca. ★ Further controls refer to FEs of seven wage quantiles and job tenure, omitting self-employed workers. The sample is restricted to natives between 25 and 55 years old. I cluster standard errors (G=109). Workers are observed in August of each year. Source: PILA, 2015–2018.

Figure C.3: **Quintiles of treatment effects for formal employment and formal wages in the causal forest, 2015–2018**



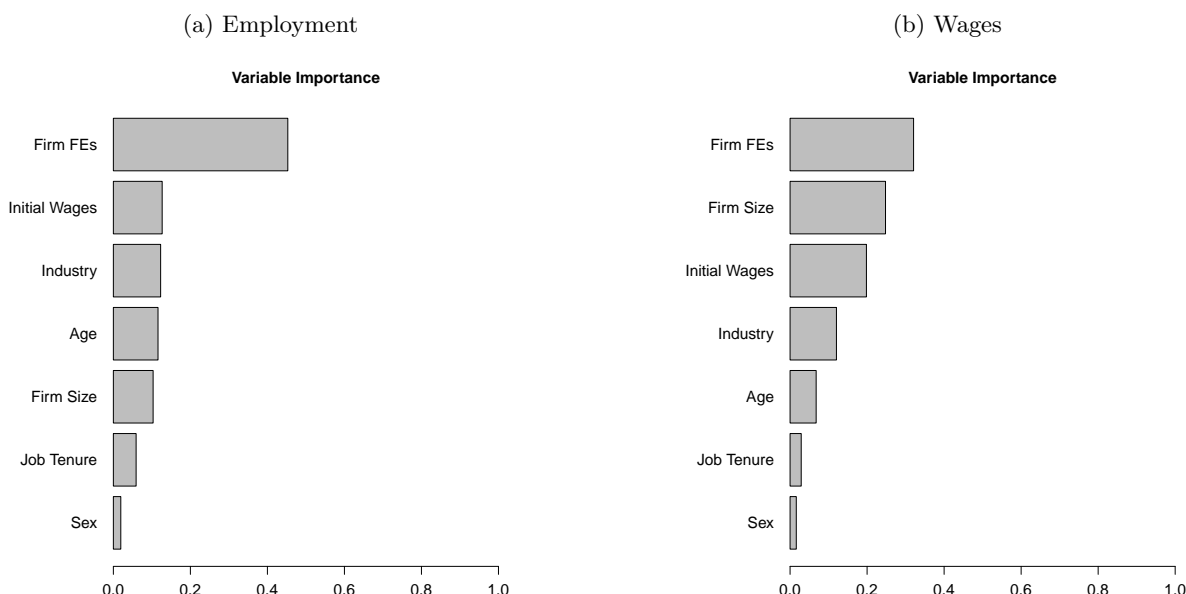
Note: The quintiles of treatment effects are constructed using the individual treatment effects from the trained causal forest. The coefficients come from separate regressions of equation (1). The sample is restricted to natives between 25 and 55 years old. I use clusters at the FUA level for the causal forest. I cluster standard errors ($G=109$). 95% confidence interval. The causal forest uses 50% of the main sample due to computational burden. The coefficients for employment (in percentage points) and wages (in percent) are already multiplied by 100.

Figure C.4: **Variable importance for formal employment and formal wages in the causal forest with decay exponent, 2015–2018**



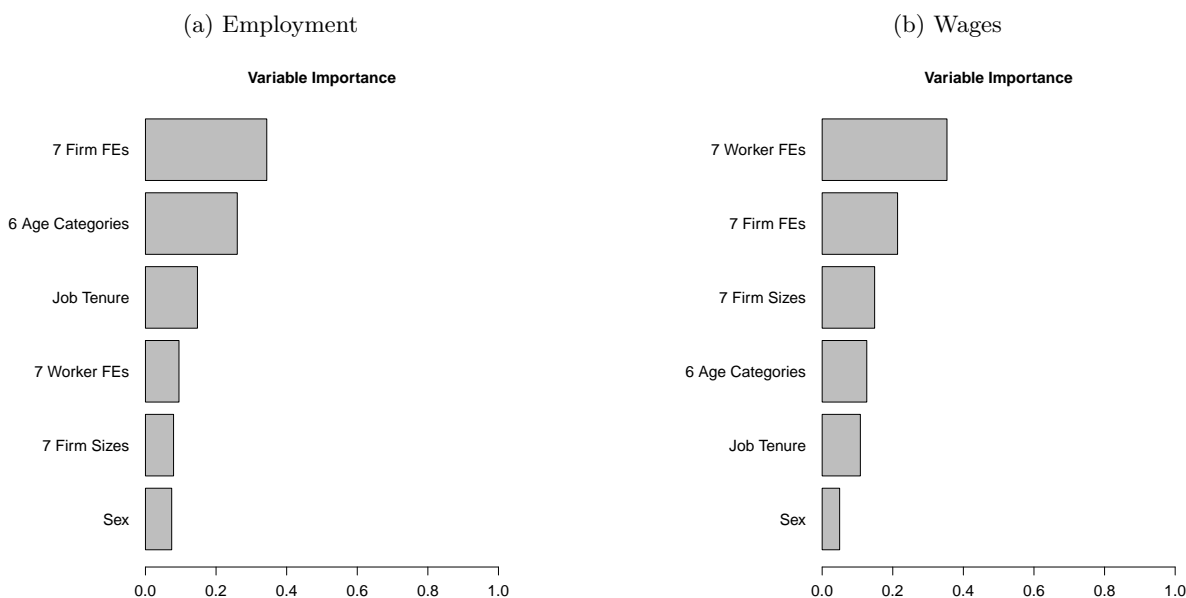
Note: Variable importance is a weighted sum of how many times the feature f appears in the split of each leaf of every tree in the forest. The number of trees is 2,000. The sample is restricted to natives between 25 and 55 years old. The importance measure sums up to 1. The decay exponent is -2. I use clusters for the causal forest estimation. The minimum node size is 300.

Figure C.5: **Variable importance for formal employment and formal wages in causal forest with industry, 2015–2018**



Note: Variable importance is a weighted sum of how many times the feature f appears in the split of each leaf of every tree in the forest. Industry information is aggregated in 16 industries. The number of trees is 2,000. The sample is restricted to natives between 25 and 55 years old. The importance measure sums up to 1. The decay exponent is -2. I use clusters for the causal forest estimation. The minimum node size is 300.

Figure C.6: **Variable importance for formal employment and formal wages in causal forest for categories, 2015–2018**



Note: Variable importance is a weighted sum of how many times the feature f appears in the split of each leaf of every tree in the forest. The number of trees is 2,000. The sample is restricted to natives between 25 and 55 years old. The importance measure sums up to 1. I use clusters for the causal forest estimation. The minimum node size is 300.

D Construction of AKM sample

To construct the sample for the AKM estimation, I restrict it to six years before the immigration shock to capture more movements of workers between firms. This sample uses the years 2010 to 2015 for August. The total sample consists of 32,195,048 worker-year observations after eliminating workers with non-positive wages, with less than 30 employment days per month, restricting to employees between 20 and 60 years, and leaving the highest wage job for workers with more than one contribution to the social security system.^{D.1} Also, I eliminate 3,931,843 additional workers because they do not belong to the largest connected set of firms and workers or appear only once in the estimation sample. Then, the nominal wages are transformed to real terms using the monthly CPI from DANE (with the base year 2018) and use logarithms of the final expression ($\ln w_{it}$). Table D.1 shows descriptive statistics by the seven quantiles of firm FEs and Table D.2 shows the decomposition of the variance of wages $Var(\ln w_{it})$.

Table D.1: **Descriptive statistics by firm FEs**

7 quantiles of $\hat{\psi}_j$	Average				N
	Employment	Male (%)	Age	Real wages (USD)	
1	8	0.6	37.7	239.2	40,201
2	18	0.7	37.1	224.0	41,628
3	14	0.6	37.2	232.5	37,703
4	13	0.6	37.5	248.3	36,223
5	18	0.5	38.0	276.4	36,599
6	40	0.5	38.3	342.0	38,524
7	81	0.5	38.4	616.1	42,455

Note: This table reports the descriptive statistics for different firm sizes recorded in PILA. Real wages are deflated using the CPI from DANE for prices in 2018. Colombian pesos to USD using 2020 exchange rates from the World Bank. Only workers who contribute as employees are taken into account. Source: PILA, August 2015.

Table D.2: Variance decomposition of $\ln w_{it}$

Share of variance explained by:	
$Var(\alpha_i)$	50.2%
$Var(\psi_{j(i)})$	15.7%
$2Cov(\alpha_i, \psi_{j(i)})$	21.6%
$Corr(\alpha_i, \psi_{j(i)})$.38

Note: This Table reports the variance decomposition of wages in the formal sector in Colombia using the largest connected set of workers and firms with the leave-out method proposed in [Kline et al. \(2020\)](#) with year FEs as the control variable. Source: PILA, August 2010–August 2015.

^{D.1} Around 5% of workers in PILA have more than one contribution.

E Derivations of Model in 5.1

In this Appendix section, I explain the derivations of the equations in subsection 5.1. First, to derive the firm-specific optimal wages, I maximize the profit equation (8) for each type of worker:^{E.1}

$$\frac{d\pi_j}{dI_j} = 0 \Leftrightarrow w_{I_j} = \left(\frac{\beta_I(1+\eta)}{1+\beta_I(1+\eta)} \right) D_j T_j^\epsilon \epsilon \alpha_I I_j^{\rho-1-\eta} (1+\eta)^{-1} (\alpha_I I_j^\rho + \alpha_F F_j^\rho)^{\frac{\epsilon-\rho}{\rho}}, \quad (\text{E.1})$$

$$\frac{d\pi_j}{dF_j} = 0 \Leftrightarrow w_{F_j} = \left(\frac{\beta_F}{1+\beta_F} \right) D_j T_j^\epsilon \epsilon \alpha_F F_j^{\rho-1} (1+\tau_F)^{-1} (\alpha_I I_j^\rho + \alpha_F F_j^\rho)^{\frac{\epsilon-\rho}{\rho}}. \quad (\text{E.2})$$

Here, workers' wages not only depend on their marginal productivity but also on the labor supply elasticities to the firm.^{E.2} For clarity, I take logarithms of the wages in equation (E.1) and (E.2):

$$\ln w_{I_j} = \ln \left(\frac{\beta_I(1+\eta)}{1+\beta_I(1+\eta)} \right) + \ln(D_j T_j^\epsilon \epsilon \alpha_I) + (\rho-1-\eta) \ln I_j - \ln(1+\eta) + \left(\frac{\epsilon-\rho}{\rho} \right) \ln(\alpha_I I_j^\rho + \alpha_F F_j^\rho), \quad (\text{E.3})$$

$$\ln w_{F_j} = \ln \left(\frac{\beta_F}{1+\beta_F} \right) + \ln(D_j T_j^\epsilon \epsilon \alpha_F) + (\rho-1) \ln F_j - \ln(1+\tau_F) + \left(\frac{\epsilon-\rho}{\rho} \right) \ln(\alpha_I I_j^\rho + \alpha_F F_j^\rho). \quad (\text{E.4})$$

In general, if I introduce a minimum wage for formal workers ($w_{F_{Min}}$) in this model such that $w_{F_{Min}} \leq w_{F_j}$, then formal workers must be paid the minimum wage and firms' optimal choices would be distorted. This is more likely to happen in low-productivity firms. Broadly, this model predicts firms with higher productivity (T_j) or demand (D_j) will pay higher wages, holding constant amenities. I then study how firm-level wages respond to an immigration shock that shifts the aggregate informal labor supply outwards ($d\mathcal{I}$)^{E.3}:

$$\frac{d \ln w_{I_j}}{d\mathcal{I}} \cdot \mathcal{I} = (\rho-1-\eta) \frac{d \ln I_j}{d \ln \mathcal{I}} + \left(\frac{\epsilon-\rho}{\rho} \right) \frac{(\alpha_I \rho I_j^{\rho-1} \frac{dI_j}{d\mathcal{I}} + \alpha_F \rho F_j^{\rho-1} \frac{dF_j}{d\mathcal{I}})}{\alpha_I I_j^\rho + \alpha_F F_j^\rho} * \mathcal{I}, \quad (\text{E.5})$$

$$\frac{d \ln w_{F_j}}{d\mathcal{I}} \cdot \mathcal{I} = (\rho-1) \frac{d \ln F_j}{d \ln \mathcal{I}} + \left(\frac{\epsilon-\rho}{\rho} \right) \frac{(\alpha_I \rho I_j^{\rho-1} \frac{dI_j}{d\mathcal{I}} + \alpha_F \rho F_j^{\rho-1} \frac{dF_j}{d\mathcal{I}})}{\alpha_I I_j^\rho + \alpha_F F_j^\rho} * \mathcal{I}. \quad (\text{E.6})$$

Simplifying the last expressions and defining the derivatives as the elasticities, I find that:

$$\varepsilon_{w_{I_j}, \mathcal{I}} = -(1+\eta-\rho) \varepsilon_{I_j, \mathcal{I}} + (\epsilon-\rho) (s_{I_j} \varepsilon_{I_j, \mathcal{I}} + s_{F_j} \varepsilon_{F_j, \mathcal{I}}), \quad (\text{E.7})$$

^{E.1}In the derivations, I multiply by $\frac{w(L_j)}{w(L_j)}$ in the last term of FOCs to find the equations on the text.

^{E.2}If $\beta_L = 9$ then workers are paid 90% of their marginal productivity to the firm.

^{E.3}Assuming that the supply shock does not affect the firm-specific demand and the firm-specific amenities for each group of workers. Besides, the number of firms is sufficiently large such that there are no strategic interactions between firms.

$$\varepsilon_{w_{F_j}, \mathcal{I}} = -(1 - \rho)\varepsilon_{F_j, \mathcal{I}} + (\epsilon - \rho)(s_{I_j}\varepsilon_{I_j, \mathcal{I}} + s_{F_j}\varepsilon_{F_j, \mathcal{I}}). \quad (\text{E.8})$$

In these expressions, $s_{L_j} = \frac{\alpha_L L_j^\rho}{\alpha_I I_j^\rho + \alpha_F F_j^\rho}$ is the relative contribution of type of worker $L \in \{I, F\}$ to production. To solve the model, I derive the changes from the immigration shock using the firm-specific supply functions (6) and (7):

$$\varepsilon_{I_j, \mathcal{I}} = 1 + \beta_I \varepsilon_{w_{I_j}, \mathcal{I}}, \quad (\text{E.9})$$

$$\varepsilon_{F_j, \mathcal{I}} = \beta_F \varepsilon_{w_{F_j}, \mathcal{I}}. \quad (\text{E.10})$$

This yields a direct relationship between wages and employment as a function of the elasticities of supply to the firm.^{E.4} Then, I replace equations (E.9) and (E.10) into (E.7) and into (E.8):

$$\varepsilon_{w_{I_j}, \mathcal{I}} = -(1 + \eta - \rho)(1 + \beta_I \varepsilon_{w_{I_j}, \mathcal{I}}) + (\epsilon - \rho)(s_{I_j}(1 + \beta_I \varepsilon_{w_{I_j}, \mathcal{I}}) + s_{F_j} \beta_F \varepsilon_{w_{F_j}, \mathcal{I}}), \quad (\text{E.11})$$

$$\varepsilon_{w_{F_j}, \mathcal{I}} = -(1 - \rho)\beta_F \varepsilon_{w_{F_j}, \mathcal{I}} + (\epsilon - \rho)(s_{I_j}(1 + \beta_I \varepsilon_{w_{I_j}, \mathcal{I}}) + s_{F_j} \beta_F \varepsilon_{w_{F_j}, \mathcal{I}}). \quad (\text{E.12})$$

Rearranging these expressions, I find that:

$$\varepsilon_{w_{I_j}, \mathcal{I}} = \left(\frac{1}{\xi_{I_j}} \right) (-(1 + \eta - \rho) + (\epsilon - \rho)(s_{I_j} + s_{F_j} \beta_F \varepsilon_{w_{F_j}, \mathcal{I}})), \quad (\text{E.13})$$

$$\varepsilon_{w_{F_j}, \mathcal{I}} = \left(\frac{1}{\xi_{F_j}} \right) (\epsilon - \rho)s_{I_j}(1 + \beta_I \varepsilon_{w_{I_j}, \mathcal{I}}). \quad (\text{E.14})$$

Here, I define $\xi_{I_j} = 1 + (1 + \eta - \rho)\beta_I - (\epsilon - \rho)s_{I_j}\beta_I$ and $\xi_{F_j} = 1 + (1 - \rho)\beta_F - (\epsilon - \rho)s_{F_j}\beta_F$. Then, replacing equation (E.13) into (E.14) yields:

$$\varepsilon_{w_{F_j}, \mathcal{I}} = \Omega_j s_{I_j} \beta_I (\epsilon - \rho) \left(\frac{\xi_{I_j}}{\beta_I} - (1 + \eta - \rho) + (\epsilon - \rho)s_{I_j} \right). \quad (\text{E.15})$$

Here, I define $\Omega_j = \frac{1}{\xi_{I_j} \xi_{F_j} - (\epsilon - \rho)^2 s_{I_j} \beta_I s_{F_j} \beta_F}$. Last, I replace ξ_{I_j} inside of (E.15) to find the equation (9) in the main text. Next, I plug equation (9) inside equation (E.13) to find that:

$$\varepsilon_{w_{I_j}, \mathcal{I}} = \left(\frac{1}{\xi_{I_j}} \right) (-(1 + \eta - \rho) + (\epsilon - \rho)s_{I_j}(1 + s_{F_j} \Omega_j (\epsilon - \rho)\beta_F)). \quad (\text{E.16})$$

In this case, the elasticity is going to be negative $\varepsilon_{w_{I_j}, \mathcal{I}} < 0$.^{E.5} Finally, after finding that

^{E.4} Here, the total number of formal workers \mathcal{F} in the market is held constant. Besides, in this partial equilibrium framework, the response of one firm does not have spillover effects on other firms.

^{E.5} To find that $\varepsilon_{w_{I_j}, \mathcal{I}} < 0$ it is sufficient that $1 \geq s_{I_j}(1 + s_{F_j} \Omega_j (\epsilon - \rho)\beta_F)$, which always happens when $\rho > \epsilon$. On the other hand, if $\rho < \epsilon$ then $\varepsilon_{w_{I_j}, \mathcal{I}} < 0$ is also negative as $1 + \eta - \rho > \epsilon - \rho$.

informal wages always decrease with the informal supply shock, the last adjustment to analyze is what happens to informal employment within the firm. For that, I plug equation (E.16) into equation (E.9):

$$\varepsilon_{I_j, \mathcal{I}} = 1 + \left(\frac{\beta_I}{\xi_{I_j}} \right) (-(1 + \eta - \rho) + (\epsilon - \rho) s_{I_j} (1 + s_{F_j} \Omega_j (\epsilon - \rho) \beta_F)). \quad (\text{E.17})$$

After simplifying the previous expression, I find that:

$$\varepsilon_{I_j, \mathcal{I}} = \frac{1}{\xi_{I_j}} (1 + (\epsilon - \rho)^2 s_{I_j} \beta_I s_{F_j} \beta_F \Omega_j). \quad (\text{E.18})$$

Thus, in this case, a positive aggregate informal supply shock always increases informal labor within the firm ($\varepsilon_{I_j, \mathcal{I}} > 0$), independent of whether formal and informal workers' being close substitutes or not.

F Additional tests for Pre-Trends

This subsection of the Appendix tests for differential trends in the outcomes according to different workers' and firms' characteristics.

Table F.1: Event study estimates on pre-treatment periods of Figure 5a

	Employment			Wages		
	2012	2013	2014	2012	2013	2014
25 to 30 years	-0.312 (0.702)	-0.458 (0.655)	-0.412 (0.428)	0.284 (0.461)	0.370 (0.234)	-0.019 (0.124)
30 to 35 years	-0.249 (0.499)	-0.512 (0.415)	-0.305 (0.322)	0.104 (0.166)	0.226 (0.235)	-0.011 (0.272)
35 to 40 years	-0.212 (0.364)	-0.381 (0.327)	-0.060 (0.196)	0.032 (0.296)	0.166 (0.253)	-0.011 (0.273)
40 to 45 years	0.043 (0.367)	-0.155 (0.318)	-0.328 (0.240)	-0.012 (0.297)	-0.068 (0.275)	-0.501** (0.160)
45 to 50 years	0.191 (0.303)	-0.092 (0.266)	-0.101 (0.215)	0.110 (0.290)	0.566 (0.356)	-0.032 (0.187)
50 to 55 years	0.121 (0.335)	0.005 (0.262)	0.103 (0.191)	0.413 (0.293)	0.369 (0.313)	-0.209 (0.201)
Males	-0.449 (0.512)	-0.715 (0.481)	-0.473 (0.322)	0.272 (0.235)	0.317 (0.285)	-0.018 (0.170)
Females	0.297 (0.376)	0.193 (0.310)	0.094 (0.223)	-0.011 (0.209)	0.199 (0.178)	-0.225* (0.095)

Standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Note: I reduce the sample to a 10% random subsample of the entire dataset due to computational burden. The sample is restricted to natives between 25 and 55 years old. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors (G=109). I observe workers in August of each year. Source: PILA, 2012–2015.

Table F.2: Event study estimates on pre-treatment periods of Figure 6a

	Employment			Wages		
	2012	2013	2014	2012	2013	2014
Minimum wage	0.313 (0.269)	0.296 (0.213)	0.217 (0.181)	-0.538 (0.825)	-0.409 (0.700)	-0.526 (0.435)
40th–50th	-0.295 (0.502)	-0.886 (0.471)	-0.190 (0.403)	0.262 (0.388)	0.190 (0.449)	-0.107 (0.235)
50th–60th	-0.235 (0.469)	-0.766* (0.382)	-0.284 (0.240)	0.058 (0.335)	0.234 (0.298)	-0.141 (0.170)
60th–70th	-0.244 (0.319)	-0.100 (0.321)	-0.136 (0.226)	-0.360 (0.264)	0.401* (0.186)	0.150 (0.120)
70th–80th	-0.243 (0.301)	-0.475 (0.281)	-0.553** (0.211)	0.918 (0.481)	0.730* (0.341)	0.033 (0.259)
80th–90th	-0.130 (0.330)	-0.432 (0.241)	-0.385* (0.167)	0.435 (0.596)	0.367 (0.485)	-0.026 (0.268)
90th–100th	0.330 (0.483)	-0.220 (0.173)	-0.146 (0.136)	-0.039 (0.288)	0.132 (0.269)	-0.477 (0.297)

Standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Note: I reduce the sample to a 10% random subsample of the entire dataset due to computational burden. The sample is restricted to natives between 25 and 55 years old. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). I observe workers in August of each year. Source: PILA, 2012–2015.

Table F.3: Event study estimates on pre-treatment periods of Figure 7

	Employment			Wages		
	2012	2013	2014	2012	2013	2014
1-4 workers	-0.070 (0.353)	-0.022 (0.456)	-0.193 (0.368)	0.263 (0.534)	0.476 (0.389)	0.264 (0.262)
5-9 workers	-0.044 (0.481)	0.131 (0.433)	-0.484 (0.360)	-0.041 (0.300)	0.176 (0.570)	-0.088 (0.222)
10-19 workers	-0.314 (0.736)	-0.352 (0.493)	-0.446 (0.322)	0.646 (0.500)	1.156** (0.356)	0.240 (0.188)
20-49 workers	-0.525 (0.607)	-0.573 (0.622)	-0.397 (0.384)	0.511* (0.220)	0.638** (0.213)	0.398* (0.177)
50-99 workers	-0.178 (0.656)	-0.565 (0.543)	-0.497 (0.435)	0.708** (0.240)	0.877*** (0.193)	0.199 (0.186)
100 to 999 workers	-0.168 (0.695)	-0.499 (0.608)	-0.211 (0.413)	0.648 (0.620)	0.583 (0.443)	-0.137 (0.223)
More than 1000 workers	-0.239 (0.478)	-0.462 (0.474)	-0.168 (0.362)	0.202 (0.390)	0.465 (0.344)	0.134 (0.212)

Standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Note: I reduce the sample to a 10% random subsample of the entire dataset due to computational burden. The sample is restricted to native employees between 25 and 55 years old. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). I observe workers in August of each year. Source: PILA, 2012–2015.

Table F.4: Event study estimates on pre-treatment periods of Figure 8a

	Employment			Wages		
	2012	2013	2014	2012	2013	2014
Lowest quantile	-0.170 (0.469)	-0.169 (0.483)	-0.217 (0.426)	0.121 (1.819)	0.352 (1.369)	-0.867 (0.750)
2nd quantile	0.040 (0.433)	0.273 (0.401)	0.435 (0.340)	0.390 (0.355)	0.370 (0.315)	0.268 (0.150)
3rd quantile	-0.695 (0.548)	-0.900 (0.537)	-1.053* (0.457)	0.574** (0.201)	0.369 (0.230)	0.049 (0.120)
4th quantile	0.196 (0.482)	0.455 (0.438)	-0.084 (0.272)	-0.117 (0.235)	0.135 (0.277)	-0.248 (0.140)
5th quantile	-0.470 (0.462)	-0.772 (0.451)	-0.466 (0.287)	0.480 (0.447)	0.645 (0.407)	0.094 (0.188)
6th quantile	0.102 (0.416)	-0.385 (0.430)	0.119 (0.292)	0.143 (0.221)	0.383 (0.249)	0.244 (0.204)
Highest quantile	-0.342 (0.378)	-0.650 (0.362)	-0.456 (0.263)	0.394 (0.446)	0.539 (0.341)	-0.063 (0.157)

Standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Note: I reduce the sample to a 10% random subsample of the entire dataset due to computational burden. The sample is restricted to native employees between 25 and 55 years old. I use as controls interactions of sex with six age categories and a dummy for self-employed in the base period. I cluster standard errors ($G=109$). I observe workers in August of each year. Source: PILA, 2012–2015.

G Information of FUAs

Table G.1: Number of observations by FUA I

	Observations	Percent			
1. Bogotá	2,327,306	(32.7)	28. Apartadó	26,268	(0.4)
2. Medellín	983,096	(13.8)	29. Giradot	14,920	(0.2)
3. Cali	593,447	(8.3)	30. Cartago	17,006	(0.2)
4. Barranquilla	341,211	(4.8)	31. Maicao	6,263	(0.1)
5. Cartagena	205,150	(2.9)	32. Magangué	5,327	(0.1)
6. Bucaramanga	273,090	(3.8)	33. Sogamoso	18,220	(0.3)
7. Cúcuta	110,123	(1.5)	34. Buga	21,072	(0.3)
8. Pereira	140,791	(2.0)	35. Ipiales	8,754	(0.1)
9. Ibagué	100,823	(1.4)	36. Quibdó	15,687	(0.2)
10. Manizales	103,401	(1.5)	37. Fusagasugá	12,899	(0.2)
11. Santa Marta	84,705	(1.2)	38. Facatativá	18,796	(0.3)
12. Pasto	70,170	(1.0)	39. Duitama	18,427	(0.3)
13. Armenia	71,314	(1.0)	40. Yopal	43,279	(0.6)
14. Villavicencio	106,493	(1.5)	41. Ciénaga	4,701	(0.1)
15. Montería	71,007	(1.0)	42. Zipaquirá	12,908	(0.2)
16. Valledupar	76,072	(1.0)	43. Rionegro	29,601	(0.4)
17. Buenaventura	24,514	(0.3)	44. Ocaña	8,966	(0.1)
18. Neiva	71,376	(1.0)	45. La Dorada	8,563	(0.1)
19. Palmira	41,687	(0.6)	46. Caucasia	7,372	(0.1)
20. Popayán	62,422	(0.9)	47. Sabanalarga	2,434	(0.03)
21. Sincelejo	39,859	(0.6)	48. Aguachica	9,748	(0.1)
22. Barrancabermeja	35,095	(0.5)	49. Espinal	6,439	(0.1)
23. Tuluá	25,123	(0.3)	50. Arauca	11,726	(0.2)
24. Tunja	52,987	(0.7)	51. Santa Rosa de Cabal	4,887	(0.1)
25. Riohacha	31,134	(0.4)	52. El Carmen de Bolívar	1,411	(0.02)
26. San Andres de Tumaco	7,960	(0.1)	53. Fundación	3,881	(0.1)
27. Florencia	19,704	(0.3)	Continues in Table G.2		
			No FUA assigned	417,188	(5.9)
			Total	7,123,223	(100)

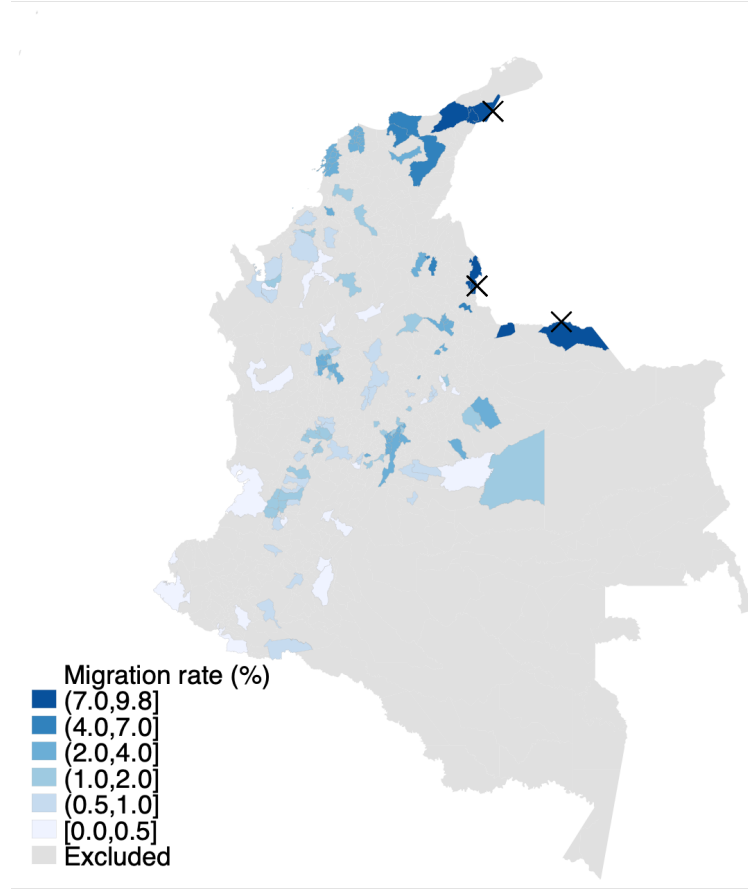
Note: This Table reports the number of workers from PILA by FUAs 1 to 53. The name represents the main city of FUA, but often they aggregate multiple municipalities according to [Sanchez-Serra \(2016\)](#). The sample is restricted to natives between 25 and 55 years old. Workers are observed in August of each year. Source: PILA, 2015.

Table G.2: Number of observations by FUA II

	Observations	Percent			
54. Acacías	12,472	(0.2)	81. Segovia	4,016	(0.1)
55. Madrid	8,922	(0.1)	82. Puerto Berrío	3,989	(0.1)
56. La Ceja	8,662	(0.1)	83. Lorica	3,875	(0.1)
57. Santander de Quilichao	8,505	(0.1)	84. Sopó	3,832	(0.1)
58. San Gil	8,268	(0.1)	85. Aguazul	3,627	(0.1)
59. Mocoa	7,974	(0.1)	86. Santa Fé de Antioquia	3,589	(0.1)
60. Pitalito	7,852	(0.1)	87. Cereté	3,526	(0.0)
61. Albania	7,020	(0.1)	88. Puerto López	3,412	(0.0)
62. Tocancipá	7,007	(0.1)	89. Pradera	3,388	(0.0)
63. Los Patios	6,137	(0.1)	90. La Cruz	3,387	(0.0)
64. Montelíbano	6,083	(0.1)	91. La Virginia	3,375	(0.0)
65. Turbo	5,830	(0.1)	92. San Pedro de los Milagros	3,170	(0.0)
66. Granada	5,298	(0.1)	93. Tenjo	3,166	(0.0)
67. El Carmen de Viboral	5,047	(0.1)	94. Villanueva	3,136	(0.0)
68. Chinchiná	4,903	(0.1)	95. Sahagún	3,126	(0.0)
69. Puerto Boyacá	4,761	(0.1)	96. Melgar	3,099	(0.0)
70. Guarne	4,697	(0.1)	97. Barbosa, Santander	3,042	(0.0)
71. Zarzal	4,584	(0.1)	98. Socorro	3,026	(0.0)
72. Puerto Asís	4,568	(0.1)	99. Carepa	2,999	(0.0)
73. Chiquinquirá	4,526	(0.1)	100. Planeta Rica	2,893	(0.0)
74. Villa de San Diego de Ubaté	4,522	(0.1)	101. Chigorodó	2,880	(0.0)
75. Garzón	4,454	(0.1)	102. Yarumal	2,874	(0.0)
76. Santa Rosa de Osos	4,406	(0.1)	103. Paipa	2,873	(0.0)
77. Puerto Gaitán	4,380	(0.1)	104. Samacá	2,782	(0.0)
78. Pamplona	4,348	(0.1)	105. Barbosa, Antioquia	2,781	(0.0)
79. Puerto Tejada	4,279	(0.1)	106. Saravena	2,730	(0.0)
80. Caloto	4,136	(0.1)	107. El Cerrito	2,597	(0.0)
			108. Amagá	2,534	(0.0)
			109. Villeta	2,518	(0.0)

Note: This Table reports the number of workers from PILA by FUAs 54 to 109. The name represents the main municipality. The sample is restricted to natives between 25 and 55 years old. Workers are observed in August of each year. Source: PILA, 2015.

Figure G.1: Map of FUAs with the immigration shock $\Delta M_{l,2018}$



Note: The X represents the main three crossing bridges with Venezuela. The distance instrument is according to the nearest crossing bridge. Source: CNPV, 2018.

Definition of Variables

Formal wages. I use the nominal contribution to the health system of each worker in August. I only consider positive contributions, as zero indicates workers on leave for several reasons unrelated to wages or jobs. I focus on workers who reported 30 days of employment.

Natives with formal employment. I count all individuals who appear in PILA with a national identity card as natives. I take all the natives in the sample with a non-negative wage as employed.

Firms. I only leave workers classified as employees for the firm-level data and then aggregate by the firm identifier.