

# Intergenerational Mobility in the Land of Inequality\*

Diogo G. C. Britto<sup>†</sup>    Alexandre Fonseca<sup>‡</sup>    Paolo Pinotti<sup>§</sup>  
Breno Sampaio<sup>¶</sup>    Lucas Warwar<sup>||</sup>

January 8, 2024

## Abstract

We provide the first estimates of intergenerational income mobility using population-wide tax data for a large developing country, namely Brazil. We measure formal income from tax and payroll data, and we train machine learning models on census and survey data to predict informal income. We develop methods to quantify and characterize the estimation bias resulting from income imputation and other sources of measurement error, and show that such bias remains negligible in our context. A 10 percentile increase in parental income rank is associated on average with a 5.5 percentile increase in child income rank, and only 2.5% of children born to parents in the bottom quintile reach the top quintile. Mobility varies widely by gender, race, and geographical areas, and causal place effects explain 57% of variation in mobility across regions.

*JEL codes:* J62, D31, I31, R23.

*Keywords:* Intergenerational mobility, Inequality, Causal place effects

---

\*This paper has benefited from comments by Koray Aktas, Massimo Anelli, Bladimir Carrillo, Christian Dustmann, Eliana La Ferrara, Joana Naritomi, Rodrigo Soares and participants in seminars and conferences at several institutions. Paolo Pinotti gratefully acknowledges ERC funding GA 861881 CLEAN. All data work for this project involving confidential taxpayer information was done at Federal Revenue of Brazil (*Receita Federal do Brasil*, RFB) facilities, on RFB computers, by RFB employees, and at no time was confidential taxpayer data ever outside of the RFB computing environment. All results have been reviewed to ensure that no confidential information is disclosed.

<sup>†</sup>University of Milan-Bicocca, Bocconi University (Baffi Centre), and CEPR, GAPPE/UFPE, IZA, e-mail: diogo.britto@unibocconi.it.

<sup>‡</sup>Federal Revenue of Brazil, GAPPE/UFPE, e-mail: alexandre.fonseca@rfb.gov.br

<sup>§</sup>Bocconi University, CEPR, CESifo, CReAM, and BAFFI-CAREFIN Center, CLEAN Unit for the Economic Analysis of Crime, CEPR, e-mail: paolo.pinotti@unibocconi.it

<sup>¶</sup>Universidade Federal de Pernambuco, GAPPE/UFPE, IZA, e-mail: breno.sampaio@ufpe.br.

<sup>||</sup>Stanford University, e-mail: warwar@stanford.edu.

# 1 Introduction

Intergenerational mobility (IGM) is a long-standing interest in social sciences and the public debate. The extent to which children’s opportunities are determined by parental income is a relevant question from both equity and efficiency perspectives. Moreover, evidence on the actual degree of IGM can shift preferences for redistributive policies, as recently shown by [Alesina et al. \(2018\)](#).

A large literature has used survey data to document that parental income is a relevant predictor of child income in adulthood, and that the strength of this association substantially varies across countries (e.g., see [Narayan et al., 2018](#); [Van der Weide et al., 2021](#)). More recently, a new wave of studies starting with [Chetty et al. \(2014\)](#) have relied on nation-wide tax data to study these questions. The use of large-scale data allowed for innovative analyses and new stylized facts – e.g., showing that mobility sharply varies for children growing up in nearby neighborhoods. Because similar data are typically not available for medium- and low-income economies, such studies have been largely restricted to the context of high-income countries.<sup>1</sup>

In this paper, we study income mobility in Brazil, a large developing country characterized by extreme inequality in socioeconomic conditions. In 2019, the Gini index was as high as 0.53 – the 9<sup>th</sup> highest worldwide – and the top 10% of the population holds 43% of the country’s income ([IBGE, 2019](#)). To conduct our analysis, we combine rich individual-level data from multiple population-wide administrative registries and large-scale household surveys. In addition to estimating IGM at the national level, these data allow us to document in detail how income mobility varies by groups and fine geographical units, and to study the role of causal place effects for upward mobility.

Importantly, we address a major challenge in the estimation of IGM in Brazil that is common to other developing countries. Almost a third of the Brazilian economy is informal ([IBGE, 2019](#)) and, as such, is not reported in administrative registries. We predict informal income by training machine

---

<sup>1</sup>For example, see [Abbas and Sicsic \(2022\)](#); [Acciari et al. \(2021\)](#); [Bratberg et al. \(2017\)](#); [Chetty et al. \(2014\)](#); [Connolly et al. \(2019\)](#); [Deutscher and Mazumder \(2020\)](#); [Heidrich \(2017\)](#); [Helsø \(2021\)](#).

learning (ML) models on rich survey data reporting income from all sources.<sup>2</sup> We use the same method to impute formal non-labor income – notably, dividends and capital gains – for earlier time periods when tax data are not available or for individuals not required to file taxes.<sup>3</sup>

This approach allows us to measure both formal and informal income at the individual level for a large, representative sample of 1.3 million children born between 1988-1990 and their parents. We use these data to estimate several IGM measures. Our main measures are based on the relationship between the percentile income rank of children and their parents', in line with [Chetty et al. \(2014\)](#). The estimated slope coefficient of the rank-rank regression equals 0.55, meaning that a 10 percentile increase in parental income is associated with an average 5.5 percentile increase in child income during adulthood. In terms of absolute mobility, children born to below-median income parents reach on average the 36<sup>th</sup> income percentile in adulthood. A transition matrix between parental and child income quintiles shows that only 2.5% of children born to parents in the bottom quintile surge to the top quintile, and only 4% of those born to parents in the top quintile fall to the bottom quintile. In turn, almost one in two children born in the bottom and top quintiles remain at the same quintiles when adults. We show that our main results are unaffected by several robustness checks that address measurement and estimation issues, notably sample selection, attenuation, and life-cycle bias.

The approach that we develop for predicting income is crucial for obtaining these results, as relying exclusively on payroll and tax data on formal income would result in a much flatter rank-rank regression (slope equal to .35). This large attenuation bias is due to administrative data neglecting informal income for a large share of individuals at the bottom of the distribution and, also, dividends and other types of capital income at the top of the income distribution when tax data are not available.

Importantly, we develop simulation exercises to quantify biases arising from

---

<sup>2</sup>We show that our ML algorithm yields improvements in accuracy relative to saturated OLS regressions. Specifically, we show that saturated OLS predictions perform poorly out-of-sample due to overfitting issues, which are not present for our ML algorithm.

<sup>3</sup>Formal labor income is available, in all these cases, from employee payroll data.

the fact that a significant portion of income in our analysis is imputed. They address different types of potential biases due to income imputation, discussed in earlier literature (e.g., see [Crossley et al., 2022](#); [Inoue and Solon, 2010](#); [Jerrim et al., 2016](#); [Zimmerman, 1992](#)). They also address potential biases related to measurement error in survey-based income measures ([Abowd and Stinson, 2013](#); [Bound et al., 1994](#); [Gottschalk and Huynh, 2010](#); [Kim and Solon, 2005](#)). The central idea is that we can learn about such biases by replacing accurately measured income components (from administrative data sources) with predicted income and studying the impacts of these changes on IGM estimates. These exercises suggest that such biases are quantitatively small in our context and unlikely to significantly affect our estimates. We also develop a formal decomposition of these biases, which is close in spirit to [Gottschalk and Huynh \(2010\)](#) who study intragenerational mobility. We show that different bias components are small in magnitude and partially offset each other, explaining the small overall bias.<sup>4</sup> These exercises make no assumption on the distribution of the measurement error components or on the model determining income, and hold for any model and covariates that one might use to predict child and parental income.

Our main findings are also robust to using two alternative, novel approaches that rank parents and children on socioeconomic status without the need to impute income components unobserved in administrative data. The first method exploits that (i) Brazilian workers move very frequently between formal and informal jobs ([Ulyssea, 2018, 2020](#)), and (ii) more than 80% of individuals in our sample hold at least one formal job in our analysis period. We can thus rank them on the average income earned during periods of formal employment, which is precisely recorded in administrative employment data, as a measure of individual-specific “productivity”. The second approach ranks parents and

---

<sup>4</sup>The key assumption for these exercises is that measurement error on informal income predictions follows a similar structure as those on formal income predictions (based on the same model and survey data). We also replicate the same exercise on a (selected) subset of the survey data where we can link the total income of parents cohabiting with adult children. We show that replacing informal income with predicted informal income has little impact on estimated mobility measures.

children on a “neighborhood-based” income measure, defined by the average (formal) income across 300 thousand census tracts, leveraging data on more than 500 million residential addresses. The rationale for this second approach is that residential choices strongly correlate with socioeconomic status, particularly in highly unequal contexts such as the Brazilian one. Each method has advantages and disadvantages, but the rank-rank curves estimated using both these two approaches are largely consistent with the one obtained using our baseline method. Importantly, these alternative approaches may be viable in other contexts characterized by a paucity of data on informal income.

Our large-scale data allow us to explore how upward mobility varies with individual characteristics and across geographical areas. A girl born to below-median income parents ranks on average 14 percentiles below boys born with the same parental income, and this gap is unaffected when we restrict the comparison to siblings. In turn, whites rank on average 7 percentiles above non-whites with the same parental income, and the gap is larger for below-median income families. While these results are broadly in line with previous evidence on differences in IGM by race in the US ([Davis and Mazumder, 2018](#); [Chetty et al., 2020](#)), they are all the more remarkable in the context of Brazil, where non-whites are not a minority group but instead represent about half of the population. We also document that higher parental income is associated with an improvement in several long-term outcomes – e.g., related to education, mortality, teenage pregnancy, welfare dependency and victimization.

Turning to heterogeneity across local areas, we uncover a mobility divide between the wealthier Center-South regions and the poorer Northern regions. A second key finding is that poor children born in the largest economic centers such as Sao Paulo and Rio de Janeiro do not achieve the best outcomes. Instead, southern regions colonized by European immigrants in the late XIX century and Center-Western regions that recently experienced a “soy boom-driven” economic growth exhibit the highest degrees of upward mobility.

Motivated by these stark regional divides, we estimate causal place effects on absolute mobility leveraging (within siblings) variation in age at move among the children of migrating families ([Chetty and Hendren, 2018a](#)). Movers

converge linearly to the income of permanent residents in the destination area at a rate of 2.4% per year of childhood exposure, meaning that children moving at birth to a place where they are expected to rank 10 percentiles higher will increase their rank by 5.76 percentiles on average due to causal place effects.<sup>5</sup> Hence, these effects explain more than half of the regional differences in absolute mobility across Brazil.

Our paper contributes to a recent body of literature estimating IGM using large-scale tax data. Starting with the seminal paper by [Chetty et al. \(2014\)](#) in the US, this literature has focused exclusively on rich countries, mainly due to data constraints.<sup>6</sup> Our paper is the first one that studies income mobility in a large developing country using population-wide administrative registries, while previous evidence on developing countries largely relied either on survey data ([Dunn, 2007](#); [Ferreira and Veloso, 2003](#); [Leone, 2018](#); [Mahlmeister et al., 2017](#); [Narayan et al., 2018](#)) or on educational mobility as a proxy for income mobility (e.g., see [Alesina et al., 2021](#); [Asher et al., 2021](#); [Saavedra and Andres, 2022b](#)). Finally, [Leites et al. \(2022\)](#) and [Meneses \(2020\)](#) have access to administrative data on income but do not attempt to estimate informal income.

Our second main contribution is methodological, as we devise new approaches to measure income mobility in contexts of high labor informality.<sup>7</sup> Specifically, we show how survey and admin data can be combined for improving income measurement and we develop new methods to quantify and decompose any bias from errors in income imputation, which has been a major concern in the IGM literature since the seminal work by [Solon \(1992\)](#). Additionally, we develop two novel methods for ranking parents and children on economic status without the need to impute informal income. These tools can be adapted to estimate IGM in other countries – including many devel-

---

<sup>5</sup>In this analysis, we measure income at the age of 24. Hence, exposure from birth to the age of 24 implies a  $24 \times 2.4\% = 57.6\%$  convergence.

<sup>6</sup>[Solon \(1999\)](#) and [Black et al. \(2011\)](#) review previous studies relying mainly on household surveys, while [Blanden \(2013\)](#) and [Björklund and Jäntti \(2020\)](#) consider alternative approaches.

<sup>7</sup>On the relationship between informality and economic development, see, e.g., [La Porta and Shleifer \(2014\)](#) and [Ulyssea \(2020\)](#).

oped economies – that are also characterized by a large informal sector.<sup>8</sup> More generally, they may find application in investigations tackling similar measurement challenges. For instance, in research using administrative data to study income dynamics and inequality in contexts where the underground economy is relevant (Engbom et al., 2022; Guvenen et al., 2022), and in studies relying on different forms of income imputation which naturally lead to measurement error (e.g., see Jácome et al., 2021).

Finally, we contribute to the literature studying the impact of places on social mobility (e.g. Chetty and Hendren, 2018a,b; Deutscher, 2020) and other long-term outcomes (e.g. Chetty et al., 2016; Chyn, 2018; Damm and Dustmann, 2014). In line with previous evidence from the US and other rich countries (Chyn and Katz, 2021), we find that causal place effects explain a large share of the total variation in intergenerational mobility. These results add to recent evidence showing that places matter for educational mobility in Africa and Latin America (Alesina et al., 2021; Saavedra and Andres, 2022a).

The remainder of the paper proceeds as follows. Section 2 briefly introduces the Brazilian context, followed by Section 3 describing our mobility measures and Section 4 describing our data, family linkage and income measurement methods. Section 5 tackles measurement error issues, while Section 6 presents our main IGM estimates at the national level and by subgroups. We explore geographic variation in mobility in Section 7 and estimate causal place effects in Section 8. Finally, Section 9 concludes.

## 2 Institutional Background

Brazil is the fifth largest country in the world by area and the sixth by population size, hosting nearly one-third of the population in Latin America, and it has historically been characterized by extreme socioeconomic inequality. In 1990 – roughly the period when our cohorts of children were born – the Gini index was as high as 0.60, placing Brazil as the fifth most unequal country

---

<sup>8</sup>Medina and Schneider (2018) estimate that, during the period 1991-2015, one fourth of Italian GDP is produced in the informal economy, and the size of the informal sector accounts for as much as 15% of GDP in countries like Canada, Denmark, Norway, and Sweden.

in the world, and the first one outside Africa. Although inequality has subsequently followed a mildly decreasing trend, the Gini index remained as high as 0.53 in 2019. According to official estimates, the top 10% of the population holds 43% of the country's income (IBGE, 2019), compared to 31% in the US, 29% in China, and around 25% in European countries.<sup>9</sup>

The country's colonial past, characterized by short-spanned extractive economic cycles and over 350 years of slavery, bestowed strong social disparities. The gap in income per capita between white and non-white households is over 35%. Non-whites represent nearly half of the population but account for 64% of the unemployed, 67% of the incarcerated population, and 75% of the beneficiaries of *Bolsa Família* cash transfers. Socio-economic conditions also vary widely across geographical areas. The country comprises 27 states (and 5,570 municipalities), and GDP per capita is about 40% lower in Northern states relative to the more developed Center-South. The homicide rate ranges from above 50 per 100k inhabitants in poorer states such as Roraima and Ceará to below 12 in the richest states such as São Paulo and Santa Catarina. These facts further motivate an analysis of mobility across subgroups and geographical areas.

Like in most low- and middle-income countries, the labor market is characterized by a large degree of informality. Labor turnover is also very high, with 70% of formal jobs lasting less than a year, and it is common for workers to turnover between the formal and informal sector (Ulyssea, 2018, 2020). In our data, 82.8% of men have held at least one formal job over their lifetime, but about 40% of workers are employed in the informal sector in a given year. Hence, it is crucial to properly measure informal income in our analysis.

The bulk of income taxes in Brazil is collected on formal labor income, although around half of formal workers are fully exempted from filing yearly income taxes because they earn below the first tax bracket (BRL 22,847 in

---

<sup>9</sup>Estimates based on the World Bank's Poverty and Inequality Platform (World Bank, 2021).



2019).<sup>10,11</sup> For the same reason, most informal workers would not pay taxes even if they had an official contract, since the majority of them earn below the first tax bracket. Dividends are fully exempt from income taxes.<sup>12</sup>

Individual income taxes are exclusively levied by the federal government and marginal tax rates range from 7.5% to 27.5%. Tax filings are mandatory for individuals with earnings above the first tax bracket, for all firm owners and for all individuals with any capital gains, any stock market operations, or property wealth above BRL 300,000.<sup>13</sup> Individuals filing taxes must report all (formal) income sources, including tax-exempted ones.

### 3 Mobility measures

Following the recent literature (e.g., [Chetty et al., 2014](#); [Acciari et al., 2021](#)), we focus on the relationship between children and parents' income ranks, as originally proposed by [Dahl and DeLeire \(2008\)](#). Since this relationship tends to be linear, it can be summarized by a few statistical parameters that can be compared across areas and groups. We estimate the linear regression:

$$y_i = \alpha + \beta p_i + \epsilon_i \quad (1)$$

where  $y_i$  and  $p_i$  are, respectively, the income percentile rank of child  $i$  and her parents' at the national level, ordered from 1 to 100. Child ranks are measured relative to their own cohorts, and parents' ranks are measured relative to other parents with children from the same cohorts.

The estimated parameters in equation (1) provide us with two IGM measures. The slope coefficient  $\beta$  measures the (inverse) *relative mobility* of children born to parents who are 1 percentile apart in the parental income distribution. A higher  $\beta$  means a wider gap between the two, thus implying lower IGM. In a perfectly mobile society, the rank-rank slope would equal zero as

---

<sup>10</sup>Throughout the paper, we refer to BRL at 2019 prices. In 2019, the purchasing power parity rate was 2.28 relative to the US dollar.

<sup>11</sup>For instance, in 2015 only slightly more than 27 million tax forms were filled in a universe of over 60 million formal workers.

<sup>12</sup>For simplicity, throughout the paper we refer to all types of withdrawals by firm owners as dividends.

<sup>13</sup>Starting in 2010, a small share of firm owners receiving dividends below 40,000 BRL were no longer required to file taxes.

children’s long-term outcomes would be unrelated to parental income.

The intercept  $\alpha$  equals the expected rank for children at the bottom of the parental income distribution. Combining  $\alpha$  and  $\beta$ , one can recover the expected rank for children born at any point of the income distribution. Following previous literature (e.g., [Chetty et al., 2014](#)), we focus on the expected rank of children born in below-median income families as our main measure of *absolute mobility*, which we also refer to as *upward mobility* throughout the paper. In turn, the latter equals the expected rank for children whose parents are in the 25<sup>th</sup> percentile of the income distribution (i.e.  $\alpha + 25 \times \beta$ ). This measure is particularly useful to characterize geographical variation in mobility patterns, as it compares the outcomes of children born in different regions of the country while holding constant parental income.

In addition, we construct transition matrices from parental income quintiles to child income quintiles. In particular, we focus on the chances of *escaping poverty* – defined as the probability that children born to parents in the bottom quintile do not belong to the same quintile when adults –, and on the probability that children move from the bottom to the top quintile of the income distribution within one generation ([Corak and Heisz, 1999](#)). We also estimate intergenerational income elasticities (IGE), defined by the correlation of children’s and parents’ log incomes. IGE allows for a comparison with earlier survey-based studies in Brazil and other countries (see, e.g., [Dunn, 2007](#); [Lee and Solon, 2009](#); [Black et al., 2011](#)).

Finally, we also document the association between parental income and several children’s long-term outcomes beyond income – namely, education, access to prestigious occupations, victimization, mortality, and teenage pregnancy.

## 4 Data and income measurement

Estimating the mobility measures described in the previous section requires (i) linking one or more cohorts of children to parents at the individual level, and (ii) measuring their individual income. Constructing such data for Brazil faces two main challenges, which are common in the context of developing countries. First, comprehensive registries of family links (of the type available,

e.g., for Scandinavian countries) are not readily available. Second, a large portion of income is earned in the informal economy and, as such, it is not reported in administrative registries. We next describe how we overcome these challenges by combining several sources of individual-level data to recover family relationships, and training supervised ML models on large-scale survey data to impute informal income. In Appendix A.1, we describe all data sources used in the paper and how we link them to our main sample.

#### 4.1 Family links

We aim to link each child’s unique person code (*CPF*) to their parents’. Our starting point is dependent claims in individual tax returns data for the 2006-2020 period, provided by the Brazilian tax authority (*Receita Federal do Brasil*). Parents report children aged 0-24 for the purpose of tax deductions, in which case we can directly link them to each other through the unique person codes available in these data.<sup>14</sup> However, only one-third of Brazilians – mainly in the upper part of the income distribution – file taxes every year (unlike in the context of rich countries, where much larger shares of the population file taxes). Therefore, we rely on additional data sources to link children who are not claimed by their parents in the tax data.

We link unclaimed children to their mothers using the Brazilian person registry (*Cadastro de Pessoas Físicas*), which covers the entire population and is provided by the Brazilian tax authority. All individuals are identified by their person code, full name, and mother’s full name. If the mother can be uniquely identified by her name – as is the case for 52% of Brazilians – we link the child’s person code to her mother’s based on the mother’s name.<sup>15</sup> Since fathers’ names are not available in the person registry, we rely on a welfare registry (*Cadastro Único*) to link children to their fathers. The registry

---

<sup>14</sup>Children aged 22-24 can only be reported if they are enrolled in technical school or higher education.

<sup>15</sup>The share of individuals with a unique name in the country is large because Brazilians typically carry one or more surnames from both their parents. These individuals are easily identified in the person registry as the latter includes both names and person codes. Britto et al. (2022) show that individuals with unique names do not strongly differ from the overall population along several characteristics.

covers around two-thirds of the Brazilian population and contains the father’s name for all individuals, along with person codes.<sup>16</sup> Since it provides the informative basis for administering social programs such as *Bolsa Família*, the registry mainly covers the low and middle parts of the income distribution. We implement the same procedure as before, linking children to their fathers conditional on the father having a unique name in the country, so that we can precisely identify his person code.

Overall, 49% and 25% of the children of the 1988-1990 cohorts can be linked to their mother and father, respectively.<sup>17</sup> Our main sample is defined by 1.34 million children who can be linked to both parents, accounting for around 15% of the entire 1988-1990 cohorts. In Appendix A.3, we show that our main sample is fairly representative of the population in terms of several individual characteristics, given that our procedure relies on complementary data sources covering different parts of the income distribution.<sup>18</sup>

Importantly, we also show that our main findings are robust to: (i) using the less conservative linkage procedure, which increases sample coverage from 15% to up to 45% of the population (see Appendix A.2); (ii) extending the sample to additional cohorts; and (iii) re-weighting the sample to eliminate any remaining differences in characteristics between our working sample and the general population.

## 4.2 Income

Our mobility measures are based on individuals’ total income, defined as the sum of formal and informal income. Accounting for informal income is crucial given the size of the informal labor market (about 40% of all jobs). For this purpose, we develop a novel approach leveraging rich survey data and ML

---

<sup>16</sup>We combine yearly snapshots of this registry for the 2011-2020 period, with 135.6 million individuals in total.

<sup>17</sup>Using younger cohorts reduces the period in which we can measure their income as adults, whereas using earlier cohorts reduces our sample because tax data on dependent claims starts in 2006. Nevertheless, we show in Appendix C.3.A that our main findings remain similar when using additional cohorts.

<sup>18</sup>Specifically, the tax registry covers the upper part of the income distribution; the person registry covers the entire distribution (for mothers); and the welfare registry covers the lower and middle part of the distribution (for fathers).

methods to estimate income that is unobserved in administrative data.<sup>19</sup>

**FORMAL INCOME.** Tax records cover all sources of formal income earned by an individual in a given year, including both labor and non-labor components. Tax-exempt income (e.g., dividends) must also be reported in tax filings. However, tax data are not always available for two reasons: first, only a third of Brazilians file taxes each year; and, second, tax data are available from 2006 onwards, limiting our ability to measure parental income until children in the main sample – born in 1988-1990 – are aged 16-18.<sup>20</sup>

Whenever tax records are unavailable for a given individual in a given year, we measure formal income as the sum of a labor and non-labor component. The first component – formal labor income – is directly available from administrative employment data covering the population of formal jobs for the 1985-2019 period (*Relação Anual de Informações Sociais*, RAIS).<sup>21</sup> The second component – formal non-labor income – includes dividends, rents, interests, and capital gains, which are not available in administrative registries other than tax data. We thus follow an imputation procedure to predict formal non-labor income leveraging survey data sources. The procedure is the same one used to input informal income, which we describe next.

**INFORMAL INCOME.** While the Brazilian administrative registries allow us to accurately measure formal income, they do not contain – by their very nature – information on informal income. We measure the latter using individual-level data from two large-scale surveys: *Pesquisa Nacional por Amostra de Domicílios* (PNAD), a cross-sectional household survey covering about 400,000 individuals per year for the 1992-1999 and 2001-2019 periods; and Population Census surveys covering 10% of the population in 1991, 2000, and 2010. Both surveys are collected by the Brazilian Institute of Geography and Statistics

---

<sup>19</sup>Our main measure of parental income is the sum of father and mother’s incomes. We proceed by estimating income for all individuals in our data to later aggregate the income of fathers and mothers.

<sup>20</sup>We show that our results remain similar when measuring parental income only for years when tax data are available (Section C.3).

<sup>21</sup>RAIS has been extensively used in previous research on the Brazilian labor market, see e.g. Ferraz et al. (2015) and Gerard and Gonzaga (2021).

(IBGE), which has a long tradition of measuring informal income for estimating GDP and other national aggregates.

We impute informal income based on a rich array of individual characteristics available in both administrative registries and survey data. This is a typical prediction problem, which we address using random forests (RF), a supervised ML algorithm that endogenously splits the space of covariates to generate predictions for a given outcome – see Appendix A.4 for details. The key advantage relative to a fully-saturated OLS is that it avoids excessively splitting the sample.

We grow a separate RF to predict informal income in each year from 1991 to 2019 by training the algorithm developed by [Athey et al. \(2019\)](#) on our survey data. The vector of predictors includes a wide array of individual characteristics: state of residence (27), a dummy identifying metropolitan regions, gender, age, race (white vs. non-white), education dummies (4), and occupation category (dummies for formal worker, formally self-employed, and firm owner, with informal workers being the residual category).

After training the model, we predict informal income for all individuals in our main sample, including formal workers and owners who may earn part of their total income in the informal sector. We repeat the same process for estimating formal non-labor income, which is necessary for measuring total formal income when tax data are not available.

In Appendix A.5, we estimate that our procedure based on the RF model predicts income ranks based on a single year with a fairly high R-squared of .57, which helps mitigating measurement error issues. As a comparison, using a fully-saturated OLS leads significantly smaller R-squared of .29 due to overfitting issues.<sup>22</sup> We also provide evidence that averaging out income over multiples years further increases the precision of our predictions.

**MAIN SAMPLE.** In the main analysis, we measure the average income of children born in 1988-1990 over the period 2015-2019, when they are 25-31 years old, and relate it to the average income of their parents (father plus mother) at

---

<sup>22</sup>R-squared statistics are based on out-of-sample predictions using a random subsample of the survey data – not used for training the prediction models.

the time when children were 3-18 years old.<sup>23</sup> The median parental and child annual income is BRL 47,068 and BRL 19,730, respectively, while the share of total income held by the top decile is around 40% for both populations. Table 1 displays descriptive statistics for the full sample and separately by gender and race. Appendix Figure A.3 shows that the distribution of total income in our sample matches the distribution of total income in the PNAD survey, apart from some (small) differences at the bottom of the father’s income distribution.

**Table 1:** Income Distribution Statistics

	Parents			Children		
	5%	50%	95%	5%	50%	95%
All	16,044	47,068	249,468	8,515	19,730	102,068
Males	9,509	31,997	193,521	10,604	22,226	117,414
Females	5,202	13,046	64,874	7,736	16,005	87,762
White	19,906	53,931	282,546	10,419	22,274	111,560
Non-white	13,797	32,917	187,921	7,388	16,267	81,658

*Notes:* The table reports the average yearly income at the 5<sup>th</sup>, 50<sup>th</sup>, and 95<sup>th</sup> percentiles of both parents and children (in 2019 BRL). The first row refers to the entire sample, while the other rows present separate statistics by gender and race. In the columns for “Parents”, the entries for “Males” and “Females” report individual incomes of fathers and mothers, respectively, while all other entries report household income. The columns for “Children” always report individual income.

## 5 Measurement error

Even though imputation is crucial for properly measuring total income, it carries with it some degree of measurement error that may bias our mobility estimates. Such error is unlikely to be classical for at least three reasons. First, errors are likely correlated across generations. Specifically, our algorithm underestimates the income of individuals with high unobserved ability, who earn

<sup>23</sup>Three years old is the earliest age at which we can measure parental income for our oldest cohort born in 1988, since survey data (in the format that we use) is available since 1991. In turn, we measure child income setting a five-year window as late as possible. Income data from the tax authority cover children in our main cohorts over the period 2015-2019, and their parents over the period 2006-2010. In addition, RAIS data are available until 2019. Appendix A.4 shows that averaging income over several years significantly increases the precision of our predictions, while Appendix C.3 shows that our main results are not affected by life-cycle bias and alternative income definitions.

above the average of their group. Such error will be positively correlated for parents and children if ability is transmitted across generations. Second, measurement error may also be non-classical because of mean-reversion: income errors tend to be negatively correlated with income levels in survey data (Bound et al., 1994; Gottschalk and Huynh, 2010). Third, our imputation procedure is equivalent to an instrumental variable approach discussed in earlier literature (e.g., see Crossley et al., 2022; Inoue and Solon, 2010; Jerrim et al., 2016; Zimmerman, 1992). As such, violations of the exclusion restriction could drive a correlation between measurement error in parental income and unexplained child income ( $\epsilon$ , equation (1)) and generate bias to our IGM estimates.

### 5.1 Bias decomposition

We provide a formal decomposition for the bias that non-classical measurement error may generate to our estimates. When estimating equation (1), we effectively estimate  $y = y^* + \eta$  on  $p = p^* + \mu$ , where  $\{y^*, p^*\}$  are actual child and parental income ranks and  $\{\eta, \mu\}$  are the respective measurement error terms. We make no assumption on the distribution of such errors or on the income generation process. In fact, error components may come from any prediction models using any characteristics which are relevant for income determination. These characteristics could be fix, such as race, or correlated across generations such as ability and education. In Appendix B.1, we show that they lead to the following estimation bias:<sup>24</sup>

$$\hat{\beta} - \beta = -\frac{1}{2}\beta\frac{v(\mu)}{v(p)} + \beta_{\epsilon\mu}\frac{v(\mu)}{v(p)} + \beta_{\eta p^*} + \beta_{\eta\mu}\frac{v(\mu)}{v(p)}, \quad (2)$$

where  $\beta$  is our coefficient of interest (i.e., the regression of  $y^*$  on  $p^*$ ) and  $\beta_{ab}$  denotes the coefficient of a hypothetical OLS regression of  $a$  on  $b$ ;  $v(\cdot)$  denotes the variance operator; and  $\epsilon$  is the error-term in Eq. (1) (i.e., child income that is unexplained by parental income). The decomposition is close in spirit to Gottschalk and Huynh (2010) who study the impacts of measurement error

---

<sup>24</sup>The same appendix provides an intermediate decomposition for the case where estimates are not based on income ranks.



on intragenerational mobility.

The first term in the decomposition is a downward bias, which grows larger in magnitude as our estimates for parental income becomes more imprecise (i.e., larger  $v(\mu)$ ), working in a similar way to attenuation bias caused by classical measurement error. The second term shows that a positive correlation between unexplained child income ( $\epsilon$ ) and measurement error on parental income will lead to an upward bias in the rank-rank slope. This captures biases due to the violation of the exclusion restriction documented in earlier literature.<sup>25</sup> The third term shows that a correlation between measurement error for child ( $v$ ) and parental income ( $p^*$ ) leads to an upward bias in the rank-rank slope. Finally, a positive correlation in measurement error across generations,  $\beta_{\eta\mu}$ , will bias our estimates of the rank-rank slope upward. The biases in the second, third and fourth components are explained by fact that inflating the left- and right-hand-side of the equation (1) at the same time drives a spurious correlation between child and parental income, leading to an upward bias in the estimation of the rank-rank slope.

In Section 6.2, we develop a simulation exercise to learn about the impact of measurement error on our estimates. The key insight is that we can learn about the impact of measurement error by replacing income components which are precisely measured in administrative data with predicted counterparts. This exercise suggests that measurement error leads only to relatively small biases in our context.

## 5.2 *Alternative approaches to rank incomes*

In addition to our main analysis, we rank parents and children on two novel measures of their overall economic conditions that do not require imputating informal income.

**PRODUCTIVITY-BASED MEASURE.** A large share of individuals in Brazil frequently turnover between the formal and informal sector, and about 80% of individuals in our main sample hold at least one formal job throughout their

---

<sup>25</sup>Parental income measurement error directly depends on parental characteristics which may have a direct effect on child income. If that is the case, such characteristics will be part of the unexplained child income component ( $\epsilon$ ), driving this bias.

career. We can thus rank individuals based on the average monthly income earned during employment spells in the formal labor market as a measure of their individual-specific productivity. The underlying assumption is that the average productivity during employment spells in the formal sector – as measured by formal earnings – is a reasonable proxy for individual productivity when employed in the informal sector. Even though this method is unable to cover individuals who have never held formal jobs, it has the key advantage of exclusively relying on high-quality data on formal labor income.

**NEIGHBORHOOD-BASED MEASURE.** Our second approach ranks parents and children based on the average income in the census tract in which they reside. Census tracts are small geographical areas designed to cover homogeneous groups of about 400 families throughout the country. The rationale for this measure is that residential choices are strongly correlated with income, particularly in poorer countries characterized by high inequality and socioeconomic spatial segregation. In addition, neighborhoods have a major impact on living standards and access to opportunities, as determined by access to public goods, job opportunities, and exposure to violence (e.g., see [Bilal and Rossi-Hansberg, 2021](#); [Card et al., 2021](#)). In our data, variation between census tracts explain 29% of total variation in formal labor income. Another important advantage of this measure is that it may better capture the high living standards of individuals benefiting from inherited wealth or living on in-kind and informal family donations. While we acknowledge that this measure may differ in nature from a pure income measure, it tracks a relevant dimension of socioeconomic status and may be useful for validating our main results based on individual income.

To implement this strategy, we geocode unique data from the Brazilian tax authority tracking 500+ million residential addresses for the entire population in the 2000-2020 period, and assign them to a census tract using shape files provided by IBGE.<sup>26</sup> We then measure the average income in each location as the average labor income of residents holding formal jobs.

---

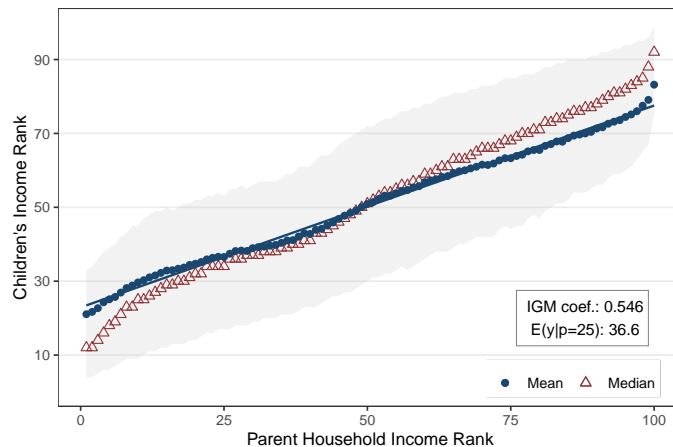
<sup>26</sup>Specifically, we consider the place of residence for children in 2000, when they are aged 10-12, and the place where they live when adults in 2019.

## 6 Income mobility at the national level

### 6.1 IGM estimates

Figure 1 plots the average and median income rank in adulthood for children born to parents in each income percentile, along with the inter-quartile range (i.e., the range between the 25th and 75th percentiles of the children’s income distribution). The ranks are based on our main measure of total income, described in Section 4.2. The rank-rank relationship is approximately linear, with the exception of the very top percentiles of the distribution, which exhibit a steeper slope. Although similar patterns have been documented for Canada, Denmark, Italy, Norway, Sweden, and the United States (Bratberg et al., 2017; Chetty et al., 2014; Corak, 2020), for the case of Brazil the change in slope is more concentrated at the very top of the distribution.

**Figure 1:** Baseline Mobility Curve in Brazil



*Notes:* The figure shows the relationship between parental and child income ranks at the national level, for our main sample (1988-1990 cohorts). For each parental income percentile, it plots the mean (blue dots), median (red triangles) and inter-quartile range (shaded area) of child income rank during 2015-2019, i.e. at the age of 25-31. Parental income is the sum of the father’s and mother’s average income when children are aged 3-18 years old. The figure also displays our absolute ( $\alpha + \beta * 25$ ) and relative mobility ( $\beta$ ) measures based on Equation (1).

The rank-rank slope coefficient in Equation (1) equals 0.546, meaning that a 10 percentile increase in parental income is associated on average with a 5.46 percentile increase in children’s income in adulthood. Based on this estimate, it would require seven generations for a family starting in the 25<sup>th</sup> percentile

to reach the same rank of a family in the 75<sup>th</sup> percentile.<sup>27</sup>

Regarding absolute mobility, a child born to parents in the 25<sup>th</sup> percentile is expected to reach the 36<sup>th</sup> percentile in adulthood. Figure 1 also shows that – even conditional on parental income – there is considerable variation in children’s outcomes. For instance, the inter-quartile range of child ranks for parents at the 25<sup>th</sup> percentile is [17, 53].

Figure 2 shows the transition matrix between quintiles of the parental and child income distributions. The probability of raising from the bottom to the top quintile within one generation is only 2.5%, and the probability of falling from the top to the bottom is only 4%. Indeed, roughly half of the children born to parents in the bottom quintile fail to escape poverty, remaining in the bottom quintile when adult; similarly, half of the children born to parents in the top quintile remain at the top of the income distribution when adult.

**Figure 2:** Transition Probability Matrix by Quintile

5	2.5%	5.1%	15%	29%	48.5%
4	10.2%	16%	23.6%	26.9%	23.4%
3	17.3%	24.5%	23.7%	20.3%	14.3%
2	24%	27.1%	22.8%	16.1%	9.9%
1	46.1%	27.4%	14.9%	7.6%	4%
	1	2	3	4	5
	Parent Household Income Quintile				

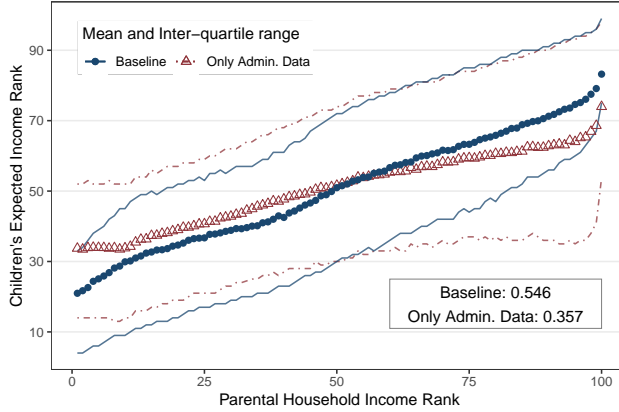
*Notes:* The figure shows the probability that children born to parents in a given quintile of the parental income distribution (horizontal axis) move to a given income quintile in adulthood (vertical axis). Darker red tones indicate higher probabilities.

<sup>27</sup>Assuming that permanent income over generations is an  $AR(1)$  process, the number of generations  $N$  required for families that are  $\Delta$  percentiles apart to converge to the same percentile solves the equation  $\beta^N \Delta = 1$ , where  $\beta$  is the rank-rank slope coefficient (Acciari et al., 2021). This back-of-the-envelope calculation might be a lower bound given that recent empirical estimates find a stronger correlation between the grandparents’ and grandchildren’s incomes than an  $AR(1)$  process would suggest (Lindahl et al., 2015; Braun and Stuhler, 2018).

## 6.2 Measurement error

In Figure 3, we show that the relationship between parent and child rank is severely attenuated when relying exclusively on administrative data, the slope coefficient decreasing from .546 to .357. The difference is particularly marked for informal workers with zero formal income, resulting in a flat relationship over the bottom 10% of the parental income distribution. In addition, the inter-quartile range is substantially larger at the upper side of the parental income distribution, which is likely due to neglecting dividends and other sources of non-labor income for years when tax data are not available. Therefore, imputing all these sources of income unreported in administrative data is crucial for correctly estimating IGM.

**Figure 3:** Mobility Curve: Baseline vs. Administrative Income Data Only



*Notes:* The figure shows our baseline mobility curve displayed in Figure 1 (blue dots) and the mobility curve obtained when solely relying on administrative data sources to measure income (red triangles). For each parental income percentile, we plot the mean child income rank during 2015-2019, i.e. when the cohorts of children in our main sample (1988-1990) were aged 25-31, along with the interquartile range. Parental income is the sum of the father's and mother's average income when the child is aged 3-18 years old. For each curve, the figure also displays the estimated  $\beta$  coefficient in Equation (1).

Nevertheless, as discussed in Section 5, our imputation process leads to some degree of measurement error that may bias our mobility estimates. To gauge the magnitude of such bias, we replace income components that are precisely measured in administrative data with predicted values based on the same ML models and survey data used to impute income in the main analysis. Then, we study how different IGM measures vary with the imputation process and decompose the implied bias in estimated rank-rank slope using equation

(2).

First, we replace formal labor income for parents and children who lie below given income percentiles in our benchmark sample. The goal is emulating the imputation of informal income in our main analysis, which is mainly based on labor income and disproportionately falls upon low-income individuals – see Appendix Figure A.2. Table 2, columns 2-4, presents the results. In all cases, the rank-rank slope remains in the range .549-.570, close to our benchmark estimates (.546). Indeed, all the different bias components tend to be small – below .028 – and the attenuation bias in the first component is more than offset by the other components, resulting in a (small) upward bias. Other mobility measures, reported in the bottom panel of the table, also remain close to our benchmark estimates: absolute mobility lies in the range 36.0-36.5 (vs. 36.6 in the main analysis) and the transition probability from the first to the top income quartile is in the range 1.4%-2.4% (vs. 2.5%).

Second, we address the fact that we impute formal non-labor income for several years when measuring parental income, due to the lack of tax data before 2006. We extend the initial simulation by replacing formal non-labor income with predicted counterparts for measuring parental income in the period 2006-2010 onwards (when tax data is available). The results in Table 2, columns 5-7, reveal similar patterns to the initial exercise. The rank-rank slope remains in the range .539-.551, close to our benchmark, and the same holds true for other mobility estimates. Bias components continue to be small in magnitude and do not change direction.<sup>28</sup>

Appendix Table C.1 provides yet another robustness tackling the issue that we need to impute formal non-labor income for several years when measuring parental income. We show how our results vary when we measure parental income using only years for which tax data is available: we take the 5-year average during the period 2006-2010 for measuring parental income, while child income is measured as in our baseline. We find a rank-rank slope (.537

---

<sup>28</sup>Bias components are only somewhat larger when replacing formal non-labor income for all parents, including those in the top quartile. This can be explained by the fact that survey data is more inaccurate at the top of the income distribution, highlighting the importance of using tax data in the analysis.

**Table 2:** Quantifying IGM biases due to formal income imputation

	Replacing income components with predicted counterparts for individuals in different income quartiles						
	Benchmark	Formal labor income (all)			Formal labor income (all) and formal non-labor income (parents only)		
		Q1	Q1-Q3	All	Q1	Q1-Q3	All
		(1)	(2)	(3)	(4)	(5)	(6)
PANEL A. RELATIVE MOBILITY							
Rank-rank slope	0.546	0.549	0.561	0.570	0.549	0.551	0.539
SE	0.001	0.001	0.001	0.001	0.001	0.001	0.001
ME Bias decomposition							
Term 1: $-\frac{1}{2}\beta\frac{v(\mu)}{v(p)}$		-0.002	-0.013	-0.028	-0.002	-0.017	-0.052
Term 2: $\beta_{\epsilon\mu}\frac{v(\mu)}{v(p)}$		0.005	0.013	0.027	0.004	0.007	0.019
Term 3: $\beta_{vp^*}$		0.001	0.014	0.021	0.001	0.014	0.021
Term 4: $\beta_{v\mu}\frac{v(\mu)}{v(p)}$		0.000	0.002	0.005	0.000	0.002	0.006
Total bias		0.004	0.016	0.024	0.003	0.006	-0.006
PANEL B. OTHER IGM MEASURES							
Exp. rank p=25	36.6	36.5	36.2	36.0	36.5	36.4	36.7
Q1Q5	2.5%	2.4%	1.8%	1.4%	2.4%	2.0%	2.1%
Q5Q5	48.5%	48.8%	51.3%	52.6%	48.8%	50.9%	50.4%
IGE	0.500	0.503	0.516	0.534	0.503	0.510	0.584

*Notes:* This table shows how benchmark relative mobility (Panel A) and several IGM measures based on our main sample change after replacing formal income with predicted counterparts for different groups (Panel B). Column 1 reports the benchmark estimates, while columns 2-6 reports IGM estimates after replacing formal income for specific groups of parents and children based on their income quartiles. Q1Q5 (Q5Q1) defines the probability that children born in income quintile 1(5) reach income quintile 5(1) in adulthood. Panel A also provides a decomposition of the total bias resulting from income imputation following the decomposition presented in Section 5.1.

vs. .545) and absolute mobility close to our main estimates (36.8 vs. 36.58).

Next, we repeat our simulation exercise on PNAD data. The key advantage of this exercise is that we can replace exactly the same income components which we impute in our main analysis: informal income and formal non-labor income. We focus on a sample of adult children aged 25-34 who live with their parents, so that we can observe their incomes and estimate a rank-rank regression which we use as the benchmark for this simulation.<sup>29</sup> Given that this is a small and an extremely selected subsample of individuals and that income is observed for a single year only, we are not directly interested in the mobility

<sup>29</sup>We focus on the period 2006-2014, and also restrict the sample to fathers aged 45-64. To avoid overfitting concerns, we train the ML model on the survey data after dropping those households in our estimation sample (children cohabiting with their parents).

estimates.<sup>30</sup> However, studying how estimates change as we impute informal income and formal non-labor income is informative on the size of measurement error biases that income imputation may generate in our main analysis. The results in Appendix Table C.2 show that key mobility estimates remain close to the benchmark after the imputation of these income components. For instance, the rank-rank slope is in the range .514-.52 (vs. .52) and absolute mobility is in the range 35.24-35.37 (vs. 35.23).<sup>31</sup>

These results offer a transparent assessment of the potential consequences of measurement error to our mobility estimates and clarify how different sources of biases interact. Overall, they suggest that the magnitudes of measurement error biases are reasonably small for different mobility measures and unlikely to overturn our key results showing strong persistence in income across generations in Brazil. In Appendix C.2, we show that similar results emerge for group and area-level mobility measures presented in the remainder of the paper.

### 6.3 *Additional robustness exercises*

In Appendix C.3, we show that our main IGM estimates are not significantly affected by other sources of bias in IGM measurement, namely selection, life-cycle, and attenuation bias. In Appendix C.4, we show that focusing on household income to measure child ranks also has little impact on our results.

### 6.4 *Alternative mobility measures*

Figures 4a-4b show how our main results change when ranking parents and children based on the “productivity” and “neighborhood-based” measures described in Section 5.2. The productivity-based curve, which relies only on pre-

---

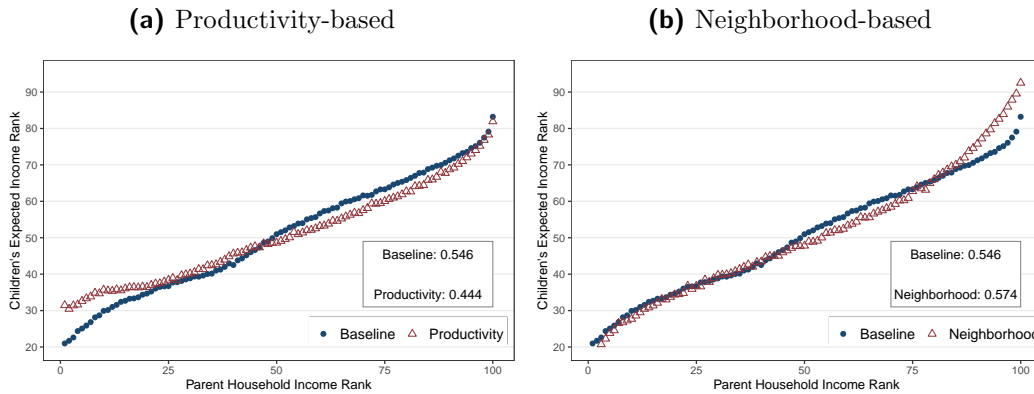
<sup>30</sup>Incidentally, estimates based on this selected sample are similar to our main estimates. However, it would be difficult to draw strong conclusion based on the survey data only given these important limitations. In addition, the small sample size would not allow for analyses across subgroups or small geographical areas as done in the remainder of the paper.

<sup>31</sup>The first and second bias components are larger in this simulation because we measure income based on a single survey year since it is not possible to follow individuals across years in PNAD, so income imputation is less precise than in our main analysis averaging income over many years (see Table ??). In particular,  $v(\mu)$  is larger, increasing the magnitudes of the first two bias components. Thanks to the fact that they compensate each other, the overall bias remains small in magnitude. In turn, IGE estimates are much more unstable, which is an additional reason for focusing on rank based measures.



cisely measured formal labor income, yields a rank-rank slope of .44, somewhat flatter relative to our baseline curve. This is the case because the minimum wage is binding in the bottom part of the distribution and because omitting capital income contributes to flattening the curve in the upper half of the distribution. Although this measure excludes individuals who have never held formal jobs in our sample, it offers some additional support to our main finding of high intergenerational income persistence in Brazil.

**Figure 4:** Alternative Measures



*Notes:* The figure plots mobility curves based on the productivity-based ranking (a) and neighborhood-based ranking (b), along with our baseline mobility curve. The productivity-based measure is based on the average formal labor income for parents and children in periods when they hold formal jobs. The neighborhood-based measure is based on the average formal income in the census tract in which children grew up (parental rank) and where they live as adults (child rank). Section 6.4 provides a detailed description of these measures. For each curve, the figure also displays the estimated  $\beta$  coefficient in Equation (1).

In turn, the neighborhood-based curve has a slope of .57, and it is steeper in the top quartile of the parental income distribution. This is consistent with the intuition that ranking individuals on income may underestimate the high living standards of children raised in affluent families, who may enjoy amenities and transfers beyond the income that they produce.<sup>32</sup> Although it may differ

<sup>32</sup>To the extent that individuals born in a given place may develop preferences for that area, the neighborhood-based measure could underestimate mobility since such preferences will mechanically create persistence in our analysis. Although such selection is endogenous and should be interpreted with some caution, Appendix Figure C.4 shows that dropping children who did not change area flattens the neighborhood-based measure, but the rank-rank slope remains as high as 0.48 and continues to show strong persistence in the top 20% of the distribution.

to some extent in nature to our main income measure, it also supports the notion that there is high intergenerational persistence in socioeconomic status in Brazil.

### 6.5 Cross-country comparisons

Although comparisons of IGM estimates across countries must be interpreted with caution (Bratberg et al., 2017; Heckman and Landersø, 2021), social mobility in Brazil appears to be much lower than in any other country for which similar estimates are available. In particular, the rank-rank slope is estimated at 0.34 in the US (Chetty et al., 2014), and ranges from 0.19 to 0.24 in Australia, Canada, France, Italy, and Scandinavian countries (Abbas and Sicsic, 2022; Acciari et al., 2021; Bratberg et al., 2017; Connolly et al., 2019; Deutscher and Mazumder, 2020; Heidrich, 2017; Helsø, 2021). Figure 5 plots these estimates against the Gini index of income inequality. Interestingly, both intergenerational persistence of income and inequality are much higher in Brazil than the other (richer) countries, and Brazil lies perfectly on the Great Gatsby Curve depicted by other countries.<sup>33</sup>

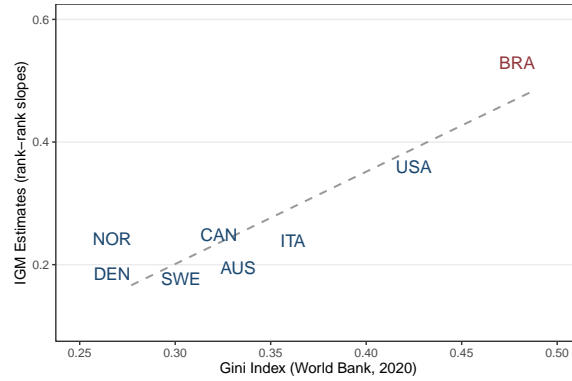
Absolute upward mobility is also lower, as below-median income children reach an income rank around 6 percentiles lower in Brazil than in the US. We reach similar conclusions when comparing the full mobility matrix across income quintiles. For instance, children born in the bottom income quintile in Brazil have only a 2.5% chance of reaching the top quintile, while the same figure is three times larger in the US (7.5%) and 4-6 times larger in Italy (11.2%) and Sweden (15.7%).

The stark contrast between Brazil and developed countries is also evident when we turn to the intergenerational income elasticity (i.e., the log-log relationship between parent and child income) as an alternative measure of income persistence (Figure 6). We estimate an IGE coefficient of .50, significantly larger than the estimates available for high-income countries, e.g., Chetty et

---

<sup>33</sup>See Corak (2013) for a discussion on the factors driving the relationship between inequality and income mobility. In line with previous evidence for the US and Italy (see, respectively, Chetty et al., 2014; Acciari et al., 2021), we also document a within-country Great Gatsby Curve, as mobility is inversely correlated with income inequality across Brazilian areas.

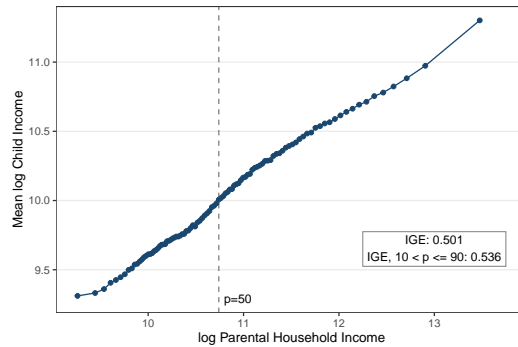
**Figure 5: The Great Gatsby Curve**



*Notes:* The figure plots the relationship between the Gini index (horizontal axis) and relative mobility (vertical axis) using this paper’s estimates for Brazil and available rank-based mobility estimates for developed countries. The latter are obtained from [Deutscher and Mazumder \(2020\)](#) (Australia), [Connolly et al. \(2019\)](#) (Canada), [Helsø \(2021\)](#) (Denmark), [Acciari et al. \(2021\)](#) (Italy), [Bratberg et al. \(2017\)](#) (Norway), [Heidrich \(2017\)](#) (Sweden), and [Chetty et al. \(2014\)](#) (US).

[al. \(2014\)](#) finds .34 for the US and [Acciari et al. \(2021\)](#) .23 for Italy. [Dunn \(2007\)](#) finds an even larger IGE of .69 in Brazil using survey data and instrumenting parental income by education. Such a higher estimate may reflect – among other things – the fact that parental education increases child income through other mechanisms beyond parental income.

**Figure 6: Intergenerational Income Elasticity (IGE)**

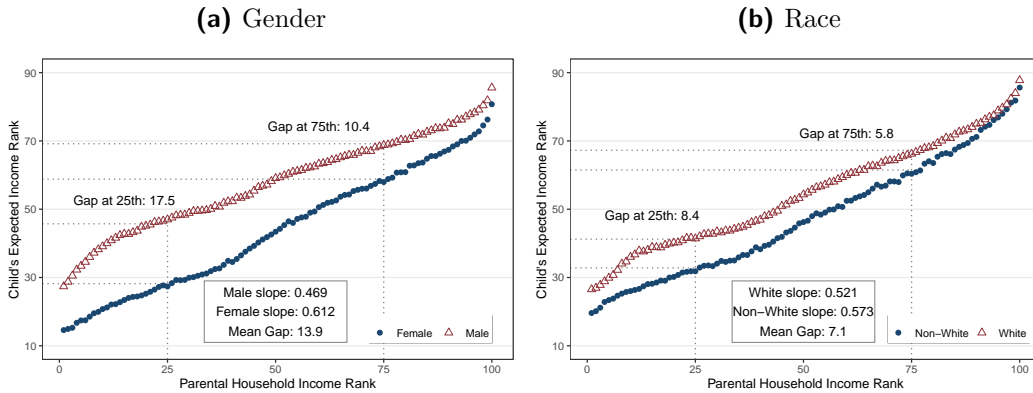


*Notes:* The figure plots the relationship between child and parental log income for our main sample (1988-1990 cohorts). For each level of log parental income (100 bins), it plots the mean log child income during 2015-2019, at the age of 25-31. It also reports the estimated IGE slope across all individuals and when restricting to parents between the 10<sup>th</sup> and 90<sup>th</sup> income percentiles. The vertical dashed line marks the median income in the parental income distribution.

### 6.6 Mobility by gender and race

Opportunities depend not only on parental income but also race and gender, especially in a country characterized by strong segregation such as Brazil. Figures 7a and 7b show the gender- and race-specific mobility curves, respectively. Importantly, the ranks on both axes indicate the positions relative to all individuals within the same cohort (rather than separately by gender and race), so the graphs show between-group differences in child ranks keeping constant parental income.

**Figure 7:** Mobility Curves by Gender and Race

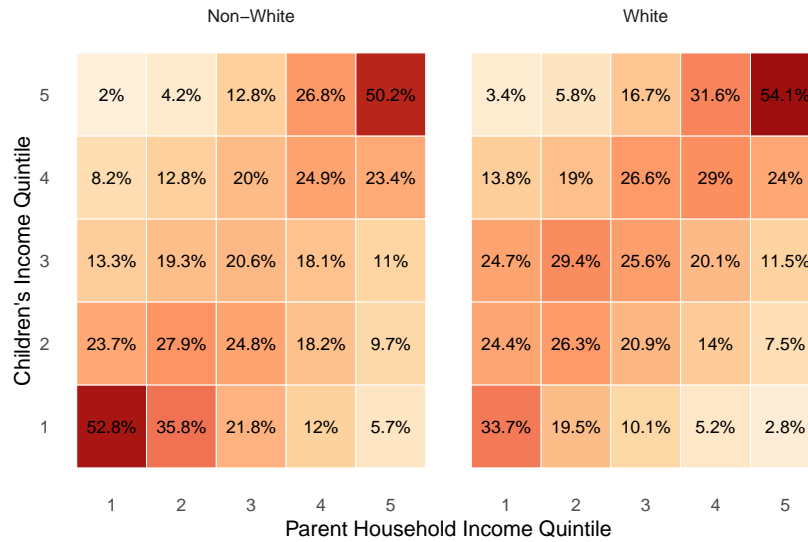


*Notes:* This figure plots separate mobility curves by gender (a) and race (b), for our main sample (1988-1990 cohorts). For each parental income percentile, it plots the mean child income rank in ages 25-31. Parental income is the sum of the father's and mother's average income when the child is aged 3-18 years old. The ranks on both axes indicate the income positions relative to all individuals within the same cohort (rather than separately by gender and race). For each curve, the figure also displays our relative mobility measure based on Equation (1), the between-group gap conditional on having parents at the 25<sup>th</sup> and 75<sup>th</sup> income percentiles, and the average between-group gap across parental income percentiles.

Female children's income is on average 14 percentiles below males with the same parental income. This gap largely reflects gender differences in labor market participation and wages (Appendix Table C.8). The mobility gap is virtually identical when restricting the same comparison to siblings, whereas the gap between siblings unconditional on gender is near zero (Appendix Table C.9). Interestingly, the rank-rank slope is steeper for females than for males (.61 vs. .47). Consequently, the gender gap declines from 17 to 10 percentiles when moving from the 25<sup>th</sup> to the 75<sup>th</sup> percentile of the parental income distribution.

Turning to differences by race, Figure 7b shows that non-white children rank on average 7 percentiles lower than white children with the same parental income. Race-specific transition matrices show that non-whites born in the first income quintile are much more likely to remain at the bottom (52.8% vs. 33.7%) and less likely to climb to the top (2% vs. 3.4%) compared to white children (Figure 8). Although differences are strongly reduced at the top of the distribution, non-whites born in the top quintile are twice as likely to fall to the bottom relative to white children (5.7% vs. 2.8%). The large mobility gap is remarkable given that non-whites – mainly comprising black and mixed-race individuals – are far from a minority in Brazil, representing about half of the population. The gap by race in Brazil is similar to the black-white gap in the US, where the former group is a minority.

**Figure 8:** Racial transition matrix



*Notes:* The figure shows the transition probability matrix by quintiles of the income distribution for the 1988-1990 cohorts separately for each race group. Each cell displays the share of children born in that parental income quintile (horizontal axis) who end up in a given income quintile in adulthood (vertical axis). Income quintiles in both axes indicate the income positions relative to all individuals in their own cohorts (rather than each group). Cells are colored according to the quintile-quintile transition probability, with darker red tones indicating higher likelihoods.

### 6.7 Parental income and children’s long-term outcomes

Next, we show that parental income is associated with improvement in several other long-term outcomes. Figure 9 shows how parental income relates to a wide array of children’s outcomes other than income. Figure 9a plots college attainment over parental income ventiles, showing that it is convex over income: while children in the bottom ventile have almost no chances of completing college, roughly 80% of children in the upper ventile do so. Girls exhibit higher educational attainment than boys over the entire parental income distribution, yet they experience lower income later in life (Figure 7a). Children in higher-income families are disproportionately more likely to hold prestigious occupations – such as doctors and lawyers – and this relationship is highly convex at the top (Figure 9b).

Figures 9c, 9d and 9e show that low parental income is also strongly associated with markers of socioeconomic struggle. Children born to below-median income families are four times more likely to receive conditional cash transfers (*Bolsa Família*), five times more likely to become teenage mothers, and twice as likely to be the victim of a crime leading to hospitalization compared to richer children.<sup>34</sup> Finally, low parental income is associated with early mortality (Figure 9f): children in low-income families are up to three times more likely to die before they turn 30.

These results suggest that income persistence may be explained (or amplified) by gaps in educational achievement and other factors that emerge early in life such as teenage fertility. The fact that all children outcomes are correlated in the expected direction with parental income and that most of these relationships are smooth bolsters our estimates of the rank-rank curve.

## 7 Geographic variation in mobility

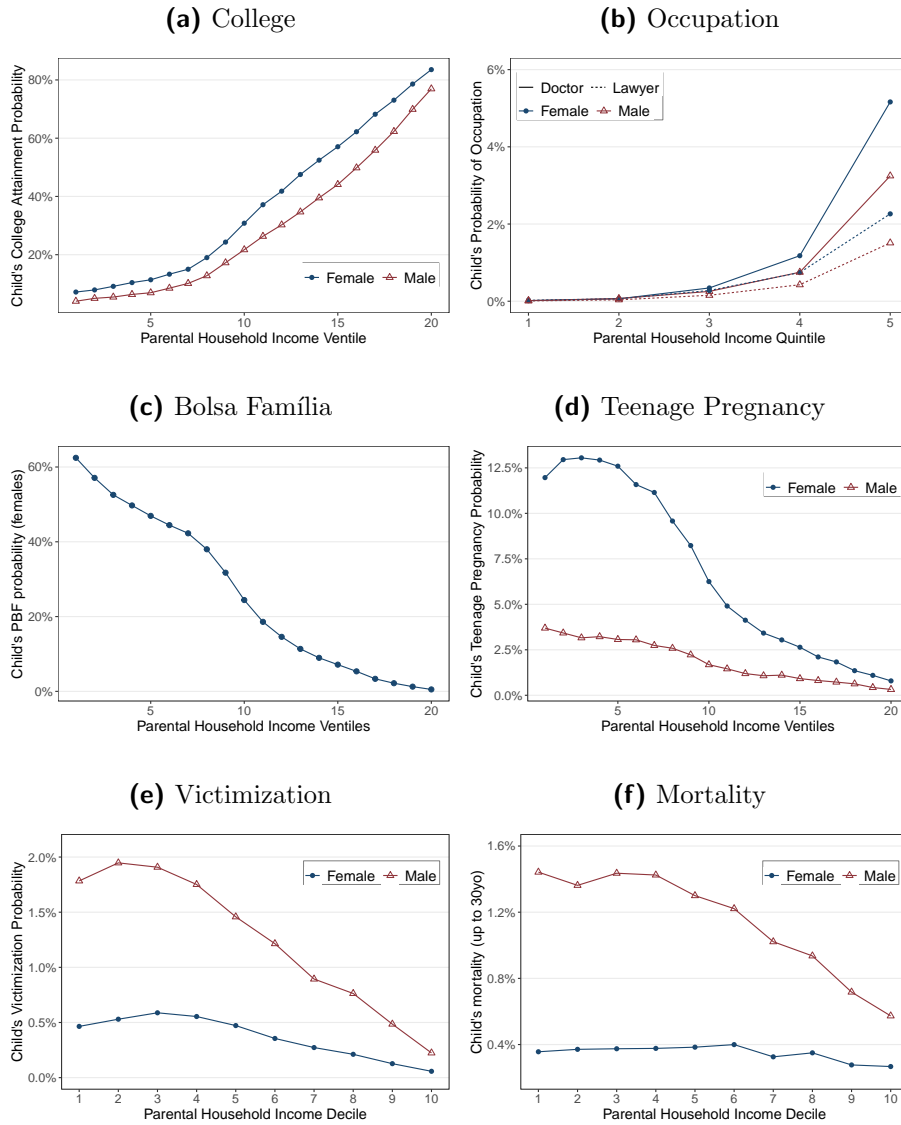
### 7.1 Geographical units and IGM measures

Brazil exhibits extreme variability in local socioeconomic conditions. We investigate social mobility across the 510 “immediate geographic regions” (IGRs),

---

<sup>34</sup>We measure victimization as the probability of hospitalization due to an assault.

**Figure 9: Long-Term Outcomes**



*Notes:* This figure plots the relationship between parental income, measured when children are aged 3-18, and several children long-term outcomes in adulthood: college degree attainment (a), the probability of working as a doctor or lawyer (b), the likelihood of receiving *Bolsa Família* transfers when adult (c), teenage pregnancy rates (d), the probability of being hospitalized due to violent assault (e), and mortality rates (f).

which are aggregations of neighboring municipalities sharing the same urban network and a common local hub (similar to the US commuting zones).<sup>35</sup> We assign children to the area where they grew up, which we proxy by their father’s place of residence (or, when the latter is missing, the mother’s) in 2000, i.e. when children in our sample were aged 10-12.<sup>36</sup> Like in the main analysis, we rank parents and children relative to the national income distribution.

## 7.2 Regional mobility patterns

The rank-rank relationship between parental and child income remains linear within regions – see, e.g., the plots for Belo Horizonte and Fortaleza, two of the largest metropolitan areas in the country, in Appendix Figure D.1. Therefore, we can compare mobility between regions using the measures of relative and absolute mobility introduced in Section 3, which rely on such linearity.

Figure 10 visualizes spatial variation in absolute mobility across IGRs. The map highlights three striking patterns. The first pattern is that absolute mobility strongly varies across regions, with the expected rank of below-median income children ranging between the 10<sup>th</sup> and the 51<sup>st</sup> percentile. More developed areas in the Center-South display higher upward mobility relative to the less affluent North and Northeast regions. A natural concern is that this map reflects different costs of living across regions. In Appendix Figure D.2, we show that adjusting for prices does not alter the main patterns in the map.<sup>37</sup>

The second striking pattern is that several regions in the countryside display higher absolute mobility than large and rich metropolitan areas such as São Paulo and Rio de Janeiro. By contrast, children in high-income families in these two areas achieve excellent outcomes – see Table D.1 reporting mobility estimates for the 50 largest metropolitan areas of the country.

The third pattern is that the top 5% areas in terms of absolute mobility are all concentrated in a large mobility hotspot crossing three southern states:

---

<sup>35</sup>IGRs replaced the *microrregião* used in earlier studies on Brazil.

<sup>36</sup>In our main sample, the father and mother live in the same IGR in 83% of cases.

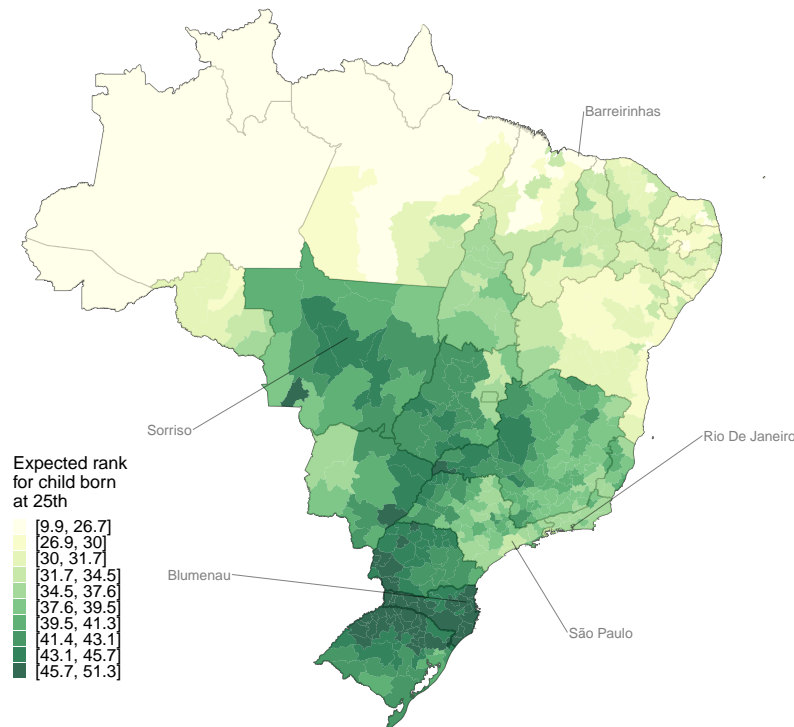
<sup>37</sup>The correlation between baseline and price-adjusted estimates of both absolute and relative mobility across regions is above .9. This high correlation is explained by the fact that, although prices significantly vary across regions, most children live in the same area where they grew up (or in areas with similar price levels).



Paraná, Santa Catarina, and Rio Grande do Sul. This region has historically been characterized by the presence of agricultural communities established by European settlers maintaining a strong cultural heritage. In such regions, below-median income children reach on average the 47<sup>th</sup> percentile in adulthood and about 80% of children born in the bottom quintile escape poverty, transiting to higher income quintiles (see Appendix Table D.1).

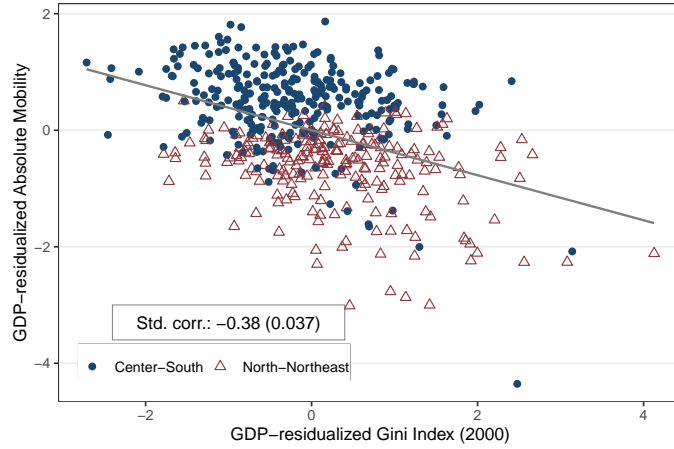
Figure 11 documents a Great Gatsby curve within Brazil, even after con-

**Figure 10:** Absolute Mobility Map: Predicted Rank for a Below-Median Income Child



*Notes:* The figure visualizes spatial variation absolute mobility (in deciles) across Brazil's 510 immediate geographical regions (IGRs) for our main sample (1988-1990). Parent and child incomes are ranked in the national income distribution and measured when children are aged 3-18 and 25-31, respectively. Absolute mobility indicates the expected rank for children in below-median income families, based on Equation (1). Darker green tones indicate higher absolute mobility. Children are assigned to IGRs according to the location of their fathers in 2000.

**Figure 11:** The Great Gatsby Curve across Brazilian regions



*Notes:* The figure plots the relationship between income inequality measured by the Gini index in 2000 (horizontal axis) and absolute mobility (vertical axis) across Brazilian regions. Both variables are residualized with respect to GDP per capita in 2002. The series in blue (dots) displays regions in the Center-South of the country and the series in red (triangles) displays regions in the North-Northeast. The figure also reports the correlation coefficient between the two (residualized) variables.

trolling for variation in GDP per capita. Appendix Section D.4 presents an analysis of the factors that better explain the substantial regional variation in mobility. Although entirely correlational, this analysis may inform future work aimed at understanding the causal determinants of upward mobility. Interestingly, we find that factors related to the quality of education provision yield by far the highest explanatory power on absolute mobility across IGRs, followed by indicators related to family structure, demographics (including the racial composition), household characteristics, and the local infrastructure. Although there is some overlap with the main mobility predictors found by Chetty et al. (2014) and Acciari et al. (2021) for the US and Italy, in Brazil the quality of education stands out as the strongest factor.

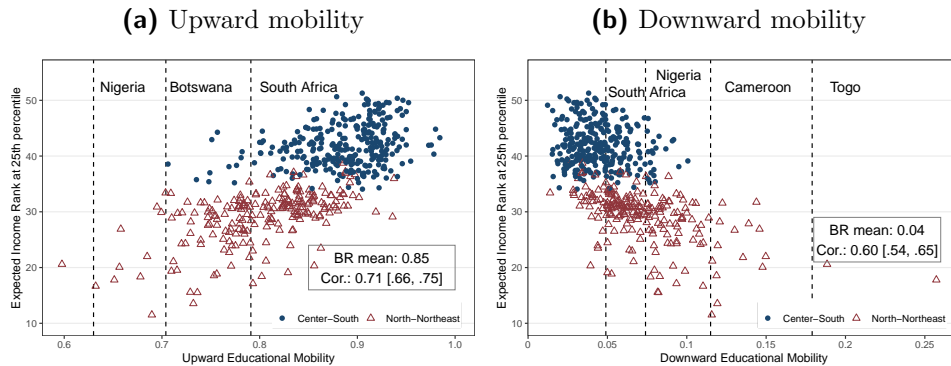
### 7.3 Income mobility and educational mobility

Figure 12 shows the relationship between income mobility and educational mobility across regions. Following Alesina et al. (2021), we compute upward mobility as the likelihood that a child born to parents who did not complete primary school manages to do so; similarly, we compute downward mobility as the likelihood that a child born to parents who completed primary school

fails to achieve the same level of education. These measures of mobility have the advantage of being available for a much larger number of countries than income-based measures, including many developing economies. [Alesina et al. \(2021\)](#) estimate educational mobility across 2,800 regions in 27 African countries; the two graphs in [Figure 12](#) also plot their estimates for some of these countries.

Three striking patterns emerge. First, the stark contrast between the North-Northeast and the Center-South of Brazil emerges for both upward and downward mobility. Second, educational mobility varies widely across Brazilian regions and, overall, it is comparable to that observed in some of the most mobile African countries – Nigeria, Botswana, and South Africa.<sup>38</sup> Finally, although income and educational mobility are strongly correlated with each other, there is a large amount of variation in income mobility for given levels of educational mobility, which further motivates the use of income-based measures.

**Figure 12:** Educational mobility across Brazilian regions



*Notes:* The figure plots estimates of upward (a) and downward (b) educational mobility across Brazilian regions (horizontal axis) versus baseline regional absolute mobility measures estimated in [Section 7](#). Blue dots (red triangles) indicate regions in the Center-South (North-Northeast) region of Brazil. Vertical lines mark estimates of educational mobility for selected African countries from [Alesina et al. \(2021\)](#). Upward (downward) mobility is the likelihood of a child born to parents who did not (did) complete primary school succeeding (failing) to do so. The figure also reports the average upward (downward) educational mobility in Brazil and the cross-regional correlation between educational mobility and income mobility.

<sup>38</sup>We achieve a similar conclusion when comparing [Saavedra and Andres \(2022b\)](#) estimates for Latin American countries.

## 8 Causal place effects

Motivated by the stark regional disparities in IGM documented in the previous section, we next estimate the causal effect of the place where children grew up on their perspectives of upward mobility. To disentangle such effect from sorting, we compare migrant children (or siblings) who moved to new areas at different ages (Chetty and Hendren, 2018a).

### 8.1 Data and research design

For this analysis, we use a sample that covers all children born during the 1983-1992 period that can be linked to their fathers. We distinguish between permanent residents and movers based on parents' residency in the 1992-2019 period. Like in Section 7, the geographical unit of analysis is the IGR. We track moves using formal employment data, because address coverage prior to 2000 is low in the person registry (see Appendix E.1 for details).

Our empirical strategy and specifications closely follow Chetty and Hendren (2018a) (see also Deutscher, 2020, for an application to Australian data). We first characterize the predicted outcomes of permanent residents using rank-rank regressions for each cohort and region (see Appendix E.2 for additional details). We then use these estimates to compute the predicted rank difference for each mover based on the origin and destination region, the child's cohort, and parental income rank. Finally, we estimate causal place effects by relating movers' income rank at the age of 24 to their predicted difference in ranks across children moving at different ages.<sup>39</sup> Intuitively, to the extent that location exert causal effects, movers' outcomes should display greater convergence to that of permanent residents the earlier they move (and the longer they are exposed) to the destination place. Specifically, our main analysis is based on the following equation:

$$y_i = \alpha_{ocpa} + \sum_{a=1}^{33} b_a I_a(a_i = a) \Delta_{odpc} + \sum_{c=1983}^{1991} \kappa_c I_c(c_i = c) \Delta_{odpc} + \epsilon_i, \quad (3)$$

---

<sup>39</sup>Like Chetty and Hendren (2018a), we focus on income at an earlier age relative to our main analysis (Section 6), so that we can measure income for (older) cohorts who move at older ages.

where  $y_i$  is the child’s income rank at the age of 24;  $\alpha_{ocpa}$  is a fixed effect by origin  $o$ , cohort  $c$ , parental income decile  $p$ , and age at move  $a$ ;  $I_a$  and  $I_c$  are indicators for each age at move  $a$  and cohort  $c$ ; and  $\Delta_{odpc}$  is the difference in permanent residents’ predicted outcomes between origin  $o$  and destination  $d$  for parental income decile  $p$  and cohort  $c$ . The coefficients of interest  $b_a$  for  $a \leq 24$  give the expected increase in rank associated with moving at age  $a$  to a destination with a 1 percentile higher predicted rank. Since we measure income at 24, moving at an older age cannot possibly have a causal effect on income, so  $b_a$  for  $a > 24$  captures solely selection effects. The coefficients  $\kappa_c$  control for our varying ability to track moves across cohorts, ensuring that we only use within-cohort variation in the age at move.<sup>40</sup> One advantage relative to [Chetty and Hendren \(2018a\)](#) is that we track moves from age 1 (instead of 11). This implies that we can flexibly study how convergence varies by age from early life, without the need to rely on linear extrapolations.

Our econometric specification effectively compares the extent of convergence to destination outcomes by children who move at different ages. The key identifying assumption is that selection effects driving some children to move to better (or worse) areas are orthogonal to the child’s age when they move. Importantly, fixed effects  $\alpha_{ocpa}$  ensure that the estimation of convergence coefficients for moves at each age exclusively relies on variation between children who have the same parental income background, and belong to the same cohort and place of origin. We provide several pieces of evidence that strongly support our main identification assumption in the next Section [8.2](#).

## 8.2 Results

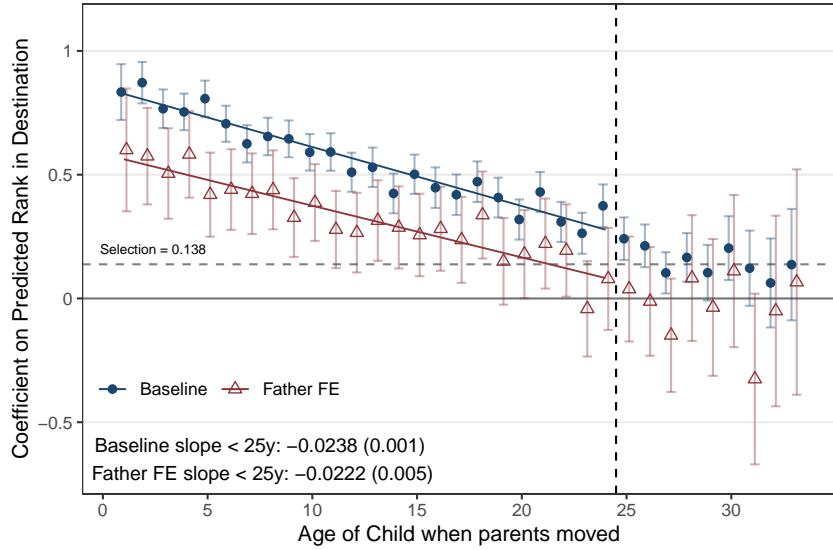
Causal place effects are summarized in [Figure 13](#). The series in blue displays the estimated coefficients  $b_a$  on the predicted difference in outcomes for each child given her origin, destination, cohort, and age at move. Supporting the intuition behind our identification strategy, the extent of convergence decreases with age, following a roughly linear pattern: children who move earlier to

---

<sup>40</sup>For instance, since we track parents’ locations starting from 1993, we observe children born in 1983 moving from 10 years old onwards, while children born in 1992 since the age of 1. To avoid collinearity, we omit the indicator for the 1992 cohort.

better places benefit more. The positive coefficients from the age of 25 purely reflect positive selection into migration, as child income measured at the age of 24 – by construction – cannot be affected by future moves. The flat pattern from the age of 24 supports our identifying assumption that parental selection into migration to a given destination does not vary with age.

**Figure 13:** Exposure Effect Estimates for Children’s Income Rank in Adulthood



*Notes:* This figure plots the estimated  $b_a$  coefficients in equation (3) (blue dots) and in an alternative specification including family fixed effects in Appendix equation (E.6) (red triangles). The sample includes all father-linked children from the 1983-1992 cohorts whose father moved once between 1993-2019, and the dependent variable – child income – is measured at the age of 24 (dashed vertical line). Vertical bars report 95% confidence intervals. Each coefficient  $b_a$  indicates the degree of convergence to the outcomes of permanent residents in the destination relative to those at the origin. Coefficients  $b_a$  for moves until the age of 24 estimate causal place effects, while  $b_a$  coefficients for older ages estimate selection effects as moves when aged 25 or older cannot explain income at age 24. The slope of the blue and red lines, as estimated by linear approximations for  $b_a$  in equations (3) and (E.6), summarize convergence per year of exposure.

Since exposure effects decline linearly with age at move, we substitute the coefficients  $b_a$  with a linear counterpart to estimate an average convergence by year of exposure.<sup>41</sup> Each additional year of exposure to the destination area

<sup>41</sup>Specifically, we substitute the non-parametric term  $\sum_{a=1}^{33} b_a I_a(a_i = a) \Delta_{odpc}$  in eq. 3 with a linear counterpart  $I(a_i \leq 24)(b_0 + (24 - a_i)\gamma) \Delta_{odpc} + I(a_i > 24)(\delta + (24 - a_i)\delta') \Delta_{odpc}$ . The linear term is split at the age of 24, since moves above this age capture selection effects.  $\gamma$  is the main coefficient of interest identifying convergence by year of exposure. We repeat the same procedure to summarize convergence for an alternative specification used in the robustness analyses.

increases convergence in children’s outcomes by .024 (baseline slope), meaning that children moving at birth to a place where they are expected to rank 10 percentiles higher will increase their rank in the national income distribution by  $.024 \times 24 \times 10 = 5.76$  percentiles on average due to causal place effects.<sup>42</sup>

These estimates suggest that about 57% of the substantial mobility gap across Brazilian regions (Section 7) is due to causal place effects. Hence, some areas in Brazil offer significantly better opportunities for low-income children, notwithstanding the high levels of income persistence at the national level.

Importantly, Figure 13 shows that our results are virtually unaffected by the addition of family fixed effects, thus relying exclusively on within-family variation (Appendix E.3 provides the details on this specification). The latter rules out the possibility that our findings are driven by family selection over child age at move. Our results are also robust to overidentification tests – see Appendix E.4. They follow from the intuition that children’s outcomes should converge to the average outcomes of their own group, whereas predicted outcomes of other groups are used as placebos.

## 9 Conclusion

In this paper, we provide the first estimates of income intergenerational mobility using large-scale tax data for a developing country, while addressing in details income measurement issues related to the informal economy. We find that income mobility in Brazil is much lower than comparable estimates available only for developed countries. Moreover, we uncover wide disparities across areas, genders, and racial groups, depicting a “*land of inequality*” in which children’s opportunities are deeply dependent on their parents’ socioeconomic status.

Importantly, we develop new methods for imputing unobserved income and studying the consequences of measurement error for IGM measures, in addition to providing alternative measures for ranking individuals on socioeconomic status. The same methods could be applied to estimate income and

---

<sup>42</sup>Focusing on educational mobility and moves between 1-11 years old, [Alesina et al. \(2021\)](#) find yearly exposure effects of 0.03 in Africa while [Saavedra and Andres \(2022a\)](#) report 0.035 for Latin America.

social mobility in other contexts characterized by a large unofficial sector. This is typically the case in low- and middle-income countries, but it is also relevant in several high income contexts (e.g., see [Medina and Schneider, 2018](#)). More generally, these methods may find application in any study where the underground economy is a challenge for income measurement.

This work is also relevant in public and policy debates. Even though Brazil has long been perceived as a place of high inequality and low mobility, hard evidence on IGM may contribute to shifting people’s perceptions and potentially their preferences for distributive policies ([Alesina et al., 2018](#)). Moreover, revealing dramatic penalties for long-neglected groups and places – in particular, non-whites and the North-Northeast of the country – can encourage public policies targeted at increasing access to opportunities. In particular, our results on causal place effects and drivers of mobility across regions can motivate placed-based policies aimed at improving the quality of public education provision in the poorest areas of Brazil.

## Bibliography

- Abbas, Hicham and Michael Sicsic**, “Who Moves Up the Income Ladder Relative to their Parents? An Analysis of Intergenerational Income Mobility in France,” 2022.
- Abowd, John M and Martha H Stinson**, “Estimating measurement error in annual job earnings: A comparison of survey and administrative data,” *Review of Economics and Statistics*, 2013, 95 (5), 1451–1467.
- Acciari, Paolo, Alberto Polo, and Gianluca Violante**, “And Yet, it Moves’: Intergenerational Mobility in Italy,” *American Economic Journal: Applied Economics*, 2021.
- Alesina, Alberto, Sebastian Hohmann, Stelios Michalopoulos, and Elias Papaioannou**, “Intergenerational Mobility in Africa,” *Econometrica*, 2021, 89 (1), 1–35.
- , **Stefanie Stantcheva, and Edoardo Teso**, “Intergenerational Mobility and Preferences for Redistribution,” *American Economic Review*, 2018, 108 (2), 521–554.
- Asher, Sam, Paul Novosad, and Charlie Rafkin**, “Intergenerational Mobility in India: Estimates From New Methods and Administrative Data,” 2021.
- Athey, Susan, Julie Tibshirani, and Stefan Wager**, “Generalized random forests,” *The Annals of Statistics*, 2019, 47 (2), 1148 – 1178.
- Bilal, Adrien and Esteban Rossi-Hansberg**, “Location as an Asset,” *Econometrica*, 2021, 89 (5), 2459–2495.



- Björklund, Anders and Markus Jäntti**, “Intergenerational mobility, intergenerational effects, sibling correlations, and equality of opportunity: a comparison of four approaches,” *Research in Social Stratification and Mobility*, 2020, *70*, 100455.
- Black, Sandra E, Paul J Devereux et al.**, “Recent Developments in Intergenerational Mobility,” *Handbook of Labor Economics*, 2011, *4*, 1487–1541.
- Blanden, Jo**, “Cross-country rankings in intergenerational mobility: a comparison of approaches from economics and sociology,” *Journal of Economic Surveys*, 2013, *27* (1), 38–73.
- Bound, John, Charles Brown, Greg J Duncan, and Willard L Rodgers**, “Evidence on the validity of cross-sectional and longitudinal labor market data,” *Journal of Labor Economics*, 1994, *12* (3), 345–368.
- Bratberg, Espen, Jonathan Davis, Bhashkar Mazumder, Martin Nybom, Daniel D. Schnitzlein, and Kjell Vaage**, “A Comparison of Intergenerational Mobility Curves in Germany, Norway, Sweden, and the US,” *The Scandinavian Journal of Economics*, 2017, *119* (1), 72–101.
- Braun, Sebastian Till and Jan Stuhler**, “The transmission of inequality across multiple generations: testing recent theories with evidence from Germany,” *The Economic Journal*, 2018, *128* (609), 576–611.
- Britto, Diogo GC, Paolo Pinotti, and Breno Sampaio**, “The effect of job loss and unemployment insurance on crime in Brazil,” *Econometrica*, 2022, *90* (4), 1393–1423.
- Card, David, Jesse Rothstein, and Moises Yi**, “Location, Location, Location,” Working Papers, U.S. Census Bureau, Center for Economic Studies 2021.
- Chetty, Raj and Nathaniel Hendren**, “The Impacts of Neighborhoods on Intergenerational Mobility I: Childhood Exposure Effects,” *The Quarterly Journal of Economics*, 2018, *133* (3), 1107–1162.
- and –, “The Impacts of Neighborhoods on Intergenerational Mobility II: County-Level Estimates,” *The Quarterly Journal of Economics*, 2018, *133* (3), 1163–1228.
- , –, and **Lawrence F Katz**, “The effects of exposure to better neighborhoods on children: New evidence from the moving to opportunity experiment,” *American Economic Review*, 2016, *106* (4), 855–902.
- , –, **Maggie R Jones, and Sonya R Porter**, “Race and Economic Opportunity in the United States: an Intergenerational Perspective,” *The Quarterly Journal of Economics*, 2020, *135* (2), 711–783.
- , –, **Patrick Kline, and Emmanuel Saez**, “Where is the land of Opportunity? The Geography of Intergenerational Mobility in the United States,” *The Quarterly Journal of Economics*, 2014, *129* (4), 1553–1623.
- Chyn, Eric**, “Moved to opportunity: The long-run effects of public housing demolition on children,” *American Economic Review*, 2018, *108* (10), 3028–56.
- and **Lawrence F Katz**, “Neighborhoods matter: Assessing the evidence for place effects,” *Journal of Economic Perspectives*, 2021, *35* (4), 197–222.
- Cohen, Jacob**, *Statistical power analysis for the behavioral sciences*, Routledge, 2013.

- Connolly, Marie, Miles Corak, and Catherine Haeck**, “Intergenerational Mobility Between and Within Canada and the United States,” *Journal of Labor Economics*, 2019, *37*, S595–S641.
- Corak, Miles**, “Income Inequality, Equality of Opportunity, and Intergenerational Mobility,” *Journal of Economic Perspectives*, 2013, *27* (3), 79–102.
- , “The Canadian geography of intergenerational income mobility,” *The Economic Journal*, 2020, *130* (631), 2134–2174.
- and **Andrew Heisz**, “The Intergenerational Earnings and Income Mobility of Canadian Men: Evidence from Longitudinal Income Tax Data,” *The Journal of Human Resources*, 1999, *34* (3), 504.
- Crossley, Thomas F, Peter Levell, and Stavros Poupakis**, “Regression with an imputed dependent variable,” *Journal of Applied Econometrics*, 2022.
- Dahl, Molly and Thomas DeLeire**, “The Association between Children’s Earnings and Fathers’ Lifetime Earnings: Estimates Using Administrative Data,” *Institute for Research on Poverty, University of Wisconsin*, 2008, *1342*.
- Damm, Anna Piil and Christian Dustmann**, “Does growing up in a high crime neighborhood affect youth criminal behavior?,” *American Economic Review*, 2014, *104* (6), 1806–32.
- Davis, J. M. and B. Mazumder**, “Racial and Ethnic Differences in the Geography of Intergenerational Mobility,” *Unpublished Working Paper*, 2018.
- der Weide, Roy Van, Christoph Lakner, Daniel Gerszon Mahler, Ambar Narayan, and Rakesh Ramasubbaiah**, “Intergenerational mobility around the world,” *Available at SSRN 3981372*, 2021.
- Deutscher, Nathan**, “Place, Peers, and the Teenage Years: Long-Run Neighborhood Effects in Australia,” 2020, *12* (2), 220–249.
- and **Bhashkar Mazumder**, “Intergenerational mobility across Australia and the stability of regional estimates,” *Labour Economics*, 2020, *66*.
- Dunn, Christopher**, “The Intergenerational Transmission of Lifetime Earnings: Evidence from Brazil,” *The B.E. Journal of Economic Analysis & Policy*, 2007, *7* (2), 1–42.
- Engbom, Niklas, Gustavo Gonzaga, Christian Moser, and Roberta Olivieri**, “Earnings inequality and dynamics in the presence of informality: The case of Brazil,” *Quantitative Economics*, 2022, *13* (4), 1405–1446.
- Ferraz, Claudio, Frederico Finan, and Dimitri Szerman**, “Procuring firm growth: the effects of government purchases on firm dynamics,” Technical Report, National Bureau of Economic Research 2015.
- Ferreira, Sergio and Fernando Veloso**, “Mobilidade Intergeracional de Educação no Brasil,” *Pesquisa e Planejamento Econômico*, 2003, *33* (3).
- Gerard, François and Gustavo Gonzaga**, “Informal Labor and the Efficiency Cost of Social Programs: Evidence from Unemployment Insurance in Brazil,” *American Economic Journal: Economic Policy*, 2021.
- Gottschalk, Peter and Minh Huynh**, “Are earnings inequality and mobility overstated? The impact of nonclassical measurement error,” *The Review of Economics and Statistics*, 2010, *92* (2), 302–315.

- Guvenen, Fatih, Luigi Pistaferri, and Giovanni L Violante**, “Global trends in income inequality and income dynamics: New insights from GRID,” *Quantitative Economics*, 2022, 13 (4), 1321–1360.
- Haider, Steven and Gary Solon**, “Life-Cycle Variation in the Association between Current and Lifetime Earnings,” *American Economic Review*, 2006, 96 (4), 1308–1320.
- Hainmueller, Jens**, “Entropy Balancing for Causal Effects: A Multivariate Reweighting Method to Produce Balanced Samples in Observational Studies,” *Political Analysis*, 2012, 20 (1), 25–46.
- Heckman, James J. and Rasmus Landersø**, “Lessons from Denmark about Inequality and Social Mobility,” 2021.
- Heidrich, Stefanie**, “Intergenerational mobility in Sweden: a regional perspective,” *Journal of Population Economics*, 2017, 30 (4), 1241–1280.
- Helsø, Anne-Line**, “Intergenerational income mobility in Denmark and the United States,” *The Scandinavian Journal of Economics*, 2021, 123 (2), 508–531.
- IBGE**, “Pesquisa Nacional por Amostra de Domicílios (PNAD),” 2019.
- Inoue, Atsushi and Gary Solon**, “Two-sample instrumental variables estimators,” *The Review of Economics and Statistics*, 2010, 92 (3), 557–561.
- Jácome, Elisa, Ilyana Kuziemko, and Suresh Naidu**, “Mobility for all: Representative intergenerational mobility estimates over the 20th century,” Technical Report, National Bureau of Economic Research 2021.
- Jerrim, John, Álvaro Choi, and Rosa Simancas**, “Two-Sample Two-Stage Least Squares (TSTSLS) estimates of earnings mobility: how consistent are they?,” in “Survey Research Methods,” Vol. 10 2016, pp. 85–101.
- Johannemann, Jonathan, Vitor Hadad, Susan Athey, and Stefan Wager**, “Sufficient Representations for Categorical Variables,” *Unpublished Working Paper.*, 2019.
- Kim, Bonggeun and Gary Solon**, “Implications of mean-reverting measurement error for longitudinal studies of wages and employment,” *Review of Economics and statistics*, 2005, 87 (1), 193–196.
- Lee, Chul-In and Gary Solon**, “Trends in intergenerational income mobility,” *The Review of Economics and Statistics*, 2009, 91 (4), 766–772.
- Leites, Martín, Xavier Ramos, Cecilia Rodríguez, and Joan Vilá**, “Intergenerational mobility and top income persistence for a developing country: estimates using administrative data from Uruguay,” 2022.
- Leone, Tharcisio**, “The geography of intergenerational mobility: Evidence of educational persistence and the “Great Gatsby Curve” in Brazil,” *Review of Development Economics*, 2018.
- Lindahl, Mikael, Mårten Palme, Sofia Sandgren Massih, and Anna Sjögren**, “Long-term intergenerational persistence of human capital an empirical analysis of four generations,” *Journal of Human Resources*, 2015, 50 (1), 1–33.
- Mahlmeister, Rodrigo, Sergio Ferreira, Fernando Veloso, Naercio Menezes Filho, and Komatsu Bruno**, “Revisitando a Mobilidade Intergeracional de Educação no Brasil,” *Inspere Policy Paper*, 2017, (26).

- Medina, Leandro and Mr Friedrich Schneider**, *Shadow economies around the world: what did we learn over the last 20 years?*, International Monetary Fund, 2018.
- Mello, Ursula, Martin Nybom, and Jan Stuhler**, “A Lifecycle Estimator of Intergenerational Earnings Mobility,” *Unpublished Working Paper*, 2021.
- Meneses, Francisco**, “Intergenerational Mobility in Chile: A year-to-year analysis of a national cohort of students (RR),” *Unpublished Working Paper*, 2020.
- Narayan, Ambar, Roy Van der Weide, Alexandru Cojocaru, Christoph Lakner, Silvia Redaelli, Daniel Gerszon Mahler, Rakesh Gupta N Ramasubbaiah, and Stefan Thewissen**, *Fair progress?: Economic mobility across generations around the world*, World Bank Publications, 2018.
- Porta, Rafael La and Andrei Shleifer**, “Informality and development,” *Journal of Economic Perspectives*, 2014, 28 (3), 109–26.
- Saavedra, Munoz and Ercio Andres**, “Does It Matter Where You Grow Up? Childhood Exposure Effects in Latin America and the Caribbean,” Working Paper, World Bank, Washington, DC May 2022. Accepted: 2022-05-13T15:39:38Z.
- and –, “The Geography of Intergenerational Mobility in Latin America and the Caribbean,” Working Paper, World Bank, Washington, DC May 2022. Accepted: 2022-05-13T15:36:20Z.
- Solon, Gary**, “Intergenerational Income Mobility in the United States,” *The American Economic Review*, 1992, 82 (3), 393–408. Publisher: American Economic Association.
- , “Intergenerational mobility in the labor market,” in “Handbook of labor economics,” Vol. 3, Elsevier, 1999, pp. 1761–1800.
- Ulyssea, Gabriel**, “Firms, Informality, and Development: Theory and Evidence from Brazil,” *American Economic Review*, August 2018, 108 (8), 2015–47.
- , “Informality: Causes and Consequences for Development,” *Annual Review of Economics*, 2020, 12 (1), 525–546.
- World Bank**, “Poverty and Inequality Platform,” 2021.
- Zimmerman, David J**, “Regression toward mediocrity in economic stature,” *The American Economic Review*, 1992, pp. 409–429.

## A Appendix to Section 4

### A.1 Description of data sources

- Person Registry: The *Cadastro de Pessoa Física* is the administrative population registry maintained by *Receita Federal*, the Brazilian tax authority. It contains all individuals who have ever held a Brazilian person code (CPF) – 255 million people in total. The CPF is similar to the social security number in the United States. Every individual in the country is identified by this unique and non-exchangeable code. Besides the person code, each observation has the person’s full name, date of birth, gender, and the full name of the mother. If the person is dead, it contains the death year, which we use to create mortality outcomes.
- Address Registry: The tax authority provided us with a dataset containing the history of individuals’ place of residence. The tax authority updates these addresses from several administrative sources, such as electoral registries and tax declarations, and when individuals autonomously update their information in the person registry. Each observation is identified by the individual’s person code, the year when the address was updated and the full residential address (street name, number, apartment/house/unit, neighborhood and postal code). Overall, there are more than 500 million addresses, which we geocoded to longitude and latitude coordinates.
- Tax Returns: The tax authority also provided us with all personal income tax returns filed during the period 2006-2020. Each observation is identified by the returnee person code, all dependents’ tax codes, and reported income divided into three categories: taxable income (mainly labor earnings and rents), tax-exempted income (mainly dividends, donations, and bequests), and income subjected to withheld or definitive taxation (mainly investment earnings and capital gains from real estate transactions). These data cover the period 2015-2019 for children in our main sample (cohorts 1988-1990) and the period 2006-2010 for their parents.

- Firm Ownership: The *Cadastro Nacional de Pessoa Jurídica* (CNPJ) is maintained by the tax authority and contains the universe of (formal) firms in Brazil, which are identified by a unique code (*CNPJ*), dates of opening and (eventual) closing, tax regime, city of registry, and a list of all shareholders identified by their person codes.
- Formal Employment: The *Relação Anual de Informações Sociais* (RAIS) is a linked employer-employee administrative dataset covering the universe of firms and workers in the formal labor market, provided by the Ministry of Labor. We use all years of RAIS available, from 1985 to 2019. Employment spells are identified by the worker’s person code and the firm’s unique identifier (*CNPJ*),<sup>43</sup> workers’ full name, gender, race, date of birth, and education; and complete information on the work contract such as dates of start and (eventual) termination, hours, wages, occupation.
- Welfare Registry: The *Cadastro Único* (CadÚnico) is an administrative registry maintained and constantly updated by the Ministry of Social Development to track the socioeconomic conditions of families with per capita income below half minimum wage or with total income below three minimum wages. It also includes all individuals of every family that has ever been a beneficiary of a federal social welfare program. We construct a yearly panel of CadÚnico from 2011 to 2020 with the individual’s full name, gender, birth year, race, education, and mother’s and father’s full names for more than 135 million individuals identified by their person codes. Each household is also identified by a unique identifier, allowing for the recovery of family structures.
- Hospitalization Records: Individual-level data on admissions to public hospitals SIH-SUS (*Sistema de Internações Hospitalares*) for the period 2002-2019. It includes information on individual characteristics such as age, sex, municipality and zip code of residence, and descriptive infor-

---

<sup>43</sup>From 1985 to 2001, workers are identified by a different (unique) code, the PIS. We retrieve PIS-CPF pairs for all workers matching individuals across RAIS waves by their full name and date of birth.

mation on the hospital admission, including the ICD-10 diagnostic, and date of admission. We use ICD-10 codes on hospitalization due to assaults to generate a measure of crime victimization. To merge these records to other datasets, we focus on individuals who can be uniquely identified by their postal code, gender and birth date – all of which can be observed for the entire population in the person registry, maintained by the Brazilian Tax Authority.

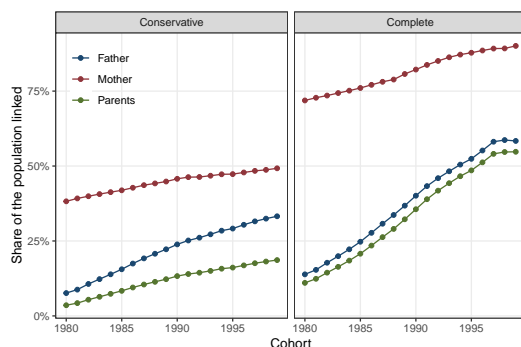
## A.2 Family links

Our main analysis is based on a conservative family linkage procedure focused on avoiding any erroneous links. For robustness purposes, we assemble an expanded sample based on additional, less conservative linkages between parents and children. Namely, we proceed by rounds and expand our conservative family linkages by additionally matching children to mothers using the following information: (i) mother’s name in the person registry, conditional that mother has a unique name within the postal code where the child lives; (ii) mother’s name and state (uf) in the person registry, conditional that mother has a unique name in the state of residence; (iii) mother’s name in the welfare registry, conditional that mother has a unique name in the country; (iv) mother’s name in the welfare registry, conditional that mother has a unique name in the state of residence; (v) mother’s name in the 2014 School Census, covering all enrolled students in Brazilian schools, conditional that the mother has a unique name in the country; (vi) mother’s name in the 2014 School Census, conditional that the mother has a unique name in the the state of residence; (vii) household composition in welfare registry. We follow the same procedure for fathers, with the exception that rounds i-iii are not available because fathers’ names are not available in the person registry. Although these linkages are somewhat less conservative relative to our baseline, conservative linkage, they remain highly accurate as they are based on high quality data on names and addresses.

Figure A.1 plots the share of children from each cohort linked to their parents when using either procedure. Many more children can be linked to

their parents – particularly mothers, since mothers’ names are available for the entire population whereas fathers’ names are available for roughly two-thirds of the population in the welfare registry. In addition, the share of successful links is increasing over time because younger cohorts can be claimed throughout more childhood years in the tax data, which start in 2006.

**Figure A.1:** Number of parent-child links relative to the population by cohort



*Notes:* The figure plots the share of the population that can be linked to their parents by cohort following our baseline, more conservative method (left graph) and the alternative, less conservative method (right graph). The first method links children to parents using unique person codes in dependent claims tax data and using names for uniquely named parents in population and welfare registries. The second method allows for additional links using individual names and addresses.

### A.3 Sample selection

Columns 1-2 in Table A.1 provide descriptive statistics for the population and our main sample. The standardized differences in column 3 are below the critical value .2 for all but three variables that slightly exceed the cutoff (race, college education and living in the North-East), indicating only small differences in the underlying distributions (Cohen, 2013). Nevertheless, in light of these small differences, we will show as a robustness test that our main findings are unaffected: (i) when substantially enlarging the sample by using the less conservative procedure to link families (see Section A.2 above), and (ii) re-weighting the sample to perfectly match the first and second moments of several characteristics in the population using the entropy algorithm by Hainmueller (2012) (Table A.1, columns 4-5) – see Section C.3.A.



**Table A.1:** Descriptive statistics of main sample

	Population	Sample	Std. Diff.	Weighted	Std. Diff., W
Female	0.494	0.512	0.037	0.494	0.000
Non-White	0.532	0.417	0.232	0.532	0.000
Primary	0.060	0.033	0.130	0.060	0.000
Elementary	0.350	0.270	0.172	0.350	0.000
High School	0.501	0.503	0.004	0.501	0.000
College	0.089	0.194	0.305	0.089	0.000
Welfare	0.624	0.581	0.086	0.622	0.003
Formal Job	0.858	0.899	0.124	0.856	0.006
Cohort 1988	0.344	0.319	0.054	0.344	0.000
Cohort 1989	0.337	0.336	0.001	0.337	0.000
Cohort 1990	0.319	0.345	0.055	0.319	0.000
North	0.089	0.077	0.044	0.089	0.000
Northeast	0.277	0.192	0.201	0.277	0.000
Southeast	0.414	0.453	0.080	0.414	0.000
South	0.143	0.195	0.140	0.144	0.005
Center-West	0.077	0.082	0.018	0.076	0.006
State capital	0.251	0.290	0.088	0.251	0.000

*Notes:* The table compares the average characteristics of our main sample of children born in 1988-1990 (omitting missing values) with the average characteristics of the same cohorts in the general population. The means for each variable are presented (columns 1-2), along with the standardized difference (column 3), the mean in the main sample after re-weighting observations to match the first and second moments of population characteristics (Hainmueller, 2012) (column 4), and the standardized difference between the samples in columns 1 and 4. All variables are recorded as dummy indicators.

#### A.4 Imputation method based on random forests

We use PNAD and population censuses to assemble a repeated cross-section from 1991 to 2019 of all adults aged 18-65 in any occupation – formal, informal, firm owner, or self-employed. We leverage the high-quality information contained in these surveys to train a generalized RF model (Athey et al., 2019) to predict informal income each year. We repeat the same process for estimating formal non-labor income (which is used when tax data are not available).

An RF is a collection of trees, each one endogenously splitting the covariate space to predict our outcomes of interest. To generate each tree, the algorithm starts by sampling without replacement from the survey dataset defining a root node. The root node is split over the space of covariates into child nodes as follows. A random subset of covariates are selected as candidates to split on, and the algorithm selects the split that maximizes heterogeneity in the prediction outcome. The sample splits take place recursively until a stop criterion met, so that overfitting is avoided.

For each outcome (informal income and formal non-labor income) and each year of data, we grow a RF with 5,000 trees, feeding the algorithm with a random subset of the original survey data. We grow *honest* forests, meaning that

separate sets of data, selected at random, are used for splitting the node and estimating prediction improvements (Athey et al., 2019). We tune all parameters of the model via cross-validation – namely, the number of variables selected at each split, the minimum node size, the penalization for imbalanced splits, and the subsample fraction for honest splits. The model covariates are state of residence (27), state capital dummy, gender, age, four education dummies, a white/non-white dummy, and worker category. Following the literature on the representation of categorical variables in ML models (Johannemann et al., 2019), we encode state of residence dummies as a set of real-valued covariates related to demographic and socioeconomic indicators.<sup>44</sup> We build the RF model to perform out-of-sample predictions of unobserved income components: namely, informal and formal non-labor income.

### A.5 Imputation method accuracy

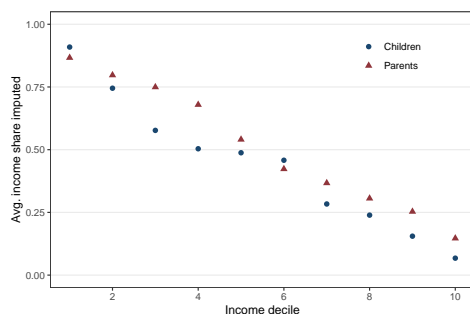
We use our RF prediction model to impute informal income and, when not available in administrative records, formal non-labor income. Figure A.2 plots the share of imputed income by income deciles for parents and children in our main sample. As expected, this share is highest at the bottom and decreases over the income distribution due to informal income being more prevalent among low-income individuals.

To assess the accuracy of our measure for income ranks, we use a random test sample from our survey for the period 1991-2019 – which has not been used to train the model so to avoid overfitting issues.<sup>45</sup> For each year in this data, we rank individuals into income percentiles based on their total income. Next we rank them again after replacing their informal and formal non-labor income with the model predictions – emulating the procedure for estimating income

---

<sup>44</sup>We define formally self-employed as owners of formal firms with zero employees, and firm owners as those owning firms with at least one formal spell in the year. We consider all working-age adults who are not formal workers or firm owners as informal workers (i.e., the residual occupation category). For census years, we use more detailed information on the municipality of residence instead of the state of residence, since the former is available in the data.

<sup>45</sup>Since we use all survey observations to train the model for our main analysis, we re-estimate again the prediction models on a 50% random sample of the survey data, and use the remaining observations to generate out-of-sample accuracy statistics.

**Figure A.2:** Average share of imputed income by income decile

*Notes:* The figure plots the average share of income imputed by income decile for children (blue dots) and parents (red triangles) in our main sample. Imputed income comprises informal income and, when not available in administrative data sources, formal non-labor income.

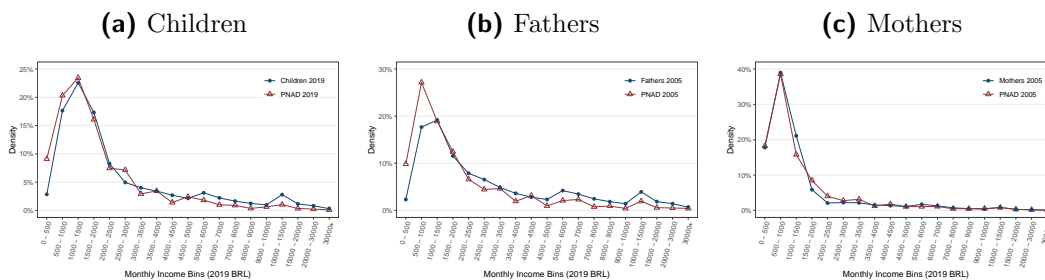
ranks in our main analysis – and then compute accuracy statistics. These estimates yield a R-squared of .57 and a root mean squared error (RMSE) of 19 ranks. For comparison, we repeat the procedure and compute the same statistics using a fully-saturated OLS prediction model – i.e., interacting fixed-effects for all covariates –, which yields a .29 R-squared and RMSE of 24 ranks. The OLS model is less accurate because of overfitting issues. Repeating this exercise on the sample used to train the model results in a much higher R-square for the OLS model (.75 in-sample vs. .29 out-of-sample), while it makes little difference for our RF model (.59 in-sample vs. .57 out-of-sample). The fairly high R-squared of our model helps reducing potential biases in mobility estimates, as pointed by earlier literature studying the consequences of income imputation for mobility estimates (e.g., see [Solon, 1992](#); [Inoue and Solon, 2010](#); [Jerrim et al., 2016](#)).

Moreover, this exercise estimates accuracy statistics on income rank predictions for a single year, while our main analysis averages out income from several years before ranking individuals. We argue that such procedure further increases the accuracy of our income measures, as it reduces the influence of transitory measurement error components. We cannot demonstrate this point based on the (cross-sectional) survey data because it does not follow the same individuals over time. Hence, we provide evidence on this by predicting formal labor income for individuals in administrative employment data (RAIS) using

our RF model. Then, we compare the R-square and RMSE for income rank predictions for parents and children in our main sample averaged over the periods 1995-1999 and 2015-2019, respectively. This exercise shows that averaging income over these periods for parents and children increases the R-squared by 27.5% and 12.2% and reduces the RMSE by 18.6% and 6.7%, respectively. Hence, averaging predictions over multiple years leads to significantly lower measurement error.

## A.6 Income distribution in our main sample and survey data sources

**Figure A.3:** Income Distribution



*Notes:* The graphs compare the income distribution in our main sample (blue dots) and in the PNAD survey (red triangles), separately for children, fathers, and mothers. The PNAD sample includes “children” born in 1988-1990 that were interviewed in the wave 2019, and “parents” of children born in 1988-1990 that were interviewed in the wave 2005. PNAD sample weights are used to compute the income distribution in PNAD data.

## B Appendix to Section 5

### B.1 Decomposition of biases due to measurement error

We formally study how measurement error in child and parental income may bias our estimates of the rank-rank slope based on equation (1). Since  $\alpha = 50(1 - \beta)$ , it is straightforward to extend the decomposition to other IGM measures based on the rank-rank regression. Our measures of child ( $y$ ) and parental ( $p$ ) income ranks in equation (1) are measured with error as  $y = y^* + \eta$  and  $p = p^* + \mu$ . Actual child and parental income are defined by  $y^*$  and  $p^*$ , and measurement errors are defined by  $\eta$  and  $\mu$ , respectively. We make no assumption on the distribution of the errors. Since we observe income with

measurement error, in practice we regress  $y^* + \eta$  on  $p^* + \mu$ . This yields a OLS estimate for  $\beta$  with the following bias:

$$\begin{aligned}\hat{\beta} - \beta &= \frac{c(y^* + \eta, p^* + \mu)}{v(p^* + \mu)} - \beta \\ &= \frac{c(y^*, p^*) + c(y^*, \mu) + c(\eta, p^*) + c(\eta, \mu)}{v(p)} - \beta \\ &= \beta \left( \frac{v(p^*)}{v(p)} - 1 \right) + \beta_{y^*\mu} \frac{v(\mu)}{v(p)} + \beta_{\eta p^*} \frac{v(p^*)}{v(p)} + \beta_{\eta\mu} \frac{v(\mu)}{v(p)}\end{aligned}\quad (\text{B.4})$$

where  $\beta$  is our coefficient of interest (i.e., the regression of  $y^*$  on  $p^*$ ) and  $\beta_{ab}$  denotes the coefficient of a hypothetical OLS regression of  $a$  on  $b$ , and  $v(\cdot)$  and  $c(\cdot)$  denote the variance and covariance operators, respectively.

This decomposition can be applied for different mobility measures based on income measured in any form: ranks, levels or logs. For the case of ranks, the variance of income ranks measured with error equals, by construction, the variance of actual income ranks. This is because the rank measure always takes the same values: they range from 1 to 100, grouping the population into their income percentile. Hence,  $v(p) = v(p^*)$  which implies that:

$$\begin{aligned}v(p) &= v(p^* + \mu) = v(p^*) + v(\mu^*) + 2c(p, \mu) \iff \\ -v(\mu) &= +2c(p^*, \mu) \iff \\ \beta_{p^*\mu} &= -1/2\end{aligned}$$

Using this result and  $v(p^*) = v(p)$  in eq. (B.4), we immediately have our final bias-decomposition formula:

$$\hat{\beta} - \beta = -\frac{1}{2}\beta \frac{v(\mu)}{v(p)} + \beta_{\epsilon\mu} \frac{v(\mu)}{v(p)} + \beta_{\eta p^*} + \beta_{\eta\mu} \frac{v(\mu)}{v(p)}$$

## C Appendix to Section 6

### C.1 Estimates fully based on years when tax data are available

**Table C.1:** IGM estimates fully based on years when tax data are available

	Rank-rank slope (1)	Exp. rank p=25 (2)	Q1Q5 (3)	Q5Q5 (4)	IGE (5)
Estimate	0.537	36.8	2.0%	49.5%	0.557

*Notes:* The table shows mobility estimates obtained using only years when tax data are available, using our main sample (1988-1990 cohorts). Parental income is measured in the period 2006-2010, when tax data are available and children are in age ranges 16-20 (1990 cohort), 17-21 (1989) and 18-22 (1988). Children's income is measured as in our baseline at ages 25-29. Q1Q5 (Q5Q1) defines the probability that children born in income quintile 1(5) reach income quintile 5(1) in adulthood.

### C.2 Quantifying measurement error biases, additional IGM measures

**Table C.2:** Quantifying IGM estimation biases due to measurement error: the impacts of informal and formal non-labor income imputation

	Replacing income components with predicted counterparts			
	Benchmark (based on survey data) (1)	Informal and formal non-labor (2)	Informal (3)	Formal non-labor (4)
PANEL A. RELATIVE MOBILITY				
Rank-rank slope (RRS)	0.520	0.516	0.514	0.520
SE	0.004	0.004	0.004	0.004
ME bias decomposition				
Term 1: $-\frac{1}{2}\beta\frac{v(\mu)}{v(p)}$		-0.150	-0.136	-0.013
Term 2: $\beta\mu\frac{v(\mu)}{v(p)}$		0.134	0.123	0.011
Term 3: $\beta_{np}$		-0.016	-0.017	0.001
Term 4: $\beta\frac{v(\mu)}{v(p)}$		0.028	0.024	0.001
Total bias		-0.005	-0.007	-0.001
PANEL B. OTHER IGM MEASURES				
Exp. rank p=25	35.23	35.32	35.37	35.24
Q1Q5	5.5%	3.5%	3.9%	5.0%
Q5Q5	47.1%	48.1%	48.1%	46.2%
IGE	0.100	0.393	0.300	0.102
Observations	45,718	45,718	45,718	45,718

*Notes:* This table shows how benchmark relative mobility estimates (Panel A) and several IGM measures (Panel B) change after replacing income components with predicted counterparts for different groups in a sample of cohabiting parents and working children aged 25-34 in PNAD survey data. Column 1 reports the benchmark estimates, while columns 2-5 reports IGM estimates after replacing income components. Q1Q5 (Q5Q1) defines the probability that children born in income quintile 1(5) reach income quintile 5(1) in adulthood. Panel A also provides a decomposition of the total bias resulting from income imputation following the decomposition presented in Section 5.1.

**Table C.3:** Quantifying IGM estimation biases due to measurement error: the impacts of formal income imputation, additional IGM measures

	Replacing income components with predicted counterparts for individuals in different income quartiles						
	Benchmark (1)	Formal labor income (all)			Formal labor income (all) and formal non-labor income (parents only)		
		Q1	Q1-Q3	All	Q1	Q1-Q3	All
		(2)	(3)	(4)	(5)	(6)	(7)
ADDITIONAL IGM MEASURES							
RRS - Female	0.611	0.615	0.628	0.632	0.614	0.616	0.597
RRS - Male	0.469	0.472	0.484	0.496	0.472	0.476	0.472
Gender gap, p=25	-17.5	-17.9	-18.8	-19.0	-17.9	-18.7	-18.8
Gender gap, p=75	-10.4	-10.8	-11.6	-12.2	-10.8	-11.7	-12.6
RRS - Non-White	0.571	0.580	0.589	0.585	0.580	0.577	0.550
RRS - White	0.520	0.528	0.544	0.545	0.528	0.533	0.509
Race gap, p=25	-8.48	-8.84	-9.24	-9.89	-8.85	-9.50	-10.56
Race gap, p=75	-5.92	-6.25	-6.96	-7.88	-6.25	-7.29	-8.48
Cor. Regional Ab. p=25	1.000	0.998	0.994	0.989	0.998	0.993	0.986

*Notes:* This table shows how additional IGM measures change after replacing formal income with predicted counterparts for different groups in our main sample (Panel B). Column 1 reports our main estimates, while columns 2-6 reports IGM estimates after replacing formal income for specific groups of parents and children based on their income quartiles. The last row shows the correlation of relative and absolute mobility across Brazil’s 510 IGRs in each simulation (columns 2-6) with our benchmark estimates (column 1).

**Table C.4:** Quantifying IGM estimation biases due to measurement error: the impacts of informal income and formal non-labor income imputation, additional IGM measures

	Benchmark (based on survey data) (1)	Predicting (replacing) income components		
		Informal and formal non-labor (2)	Informal (3)	Formal non-labor (4)
		PANEL A. ADDITIONAL MEASURES		
RRS - Female	0.511	0.536	0.533	0.511
RRS - Male	0.534	0.513	0.511	0.533
Gender gap, p=25	-3.58	-5.57	-5.38	-3.74
Gender gap, p=75	-4.69	-4.39	-4.28	-4.83
RRS - Non-White	0.476	0.430	0.439	0.469
RRS - White	0.481	0.501	0.490	0.485
Race gap, p=25	-8.86	-7.08	-7.56	-8.61
Race gap, p=75	-9.12	-10.63	-10.12	-9.38

*Notes:* This table shows how benchmark relative mobility estimates (Panel A) and several IGM measures (Panel B) change after replacing income components with predicted counterparts for different groups in a sample of cohabiting parents and working children aged 25-34 in PNAD survey data. Column 1 reports the benchmark estimates, while columns 2-5 reports IGM estimates after replacing income components. Panel A also provides a decomposition of the total bias resulting from income imputation following the decomposition presented in Appendix B.1.

### C.3 Additional robustness exercises

We now present a series of additional robustness exercises that further support our main results. In some of these analyses, we use additional birth cohorts

born in the 1983-1990 period, additional parent-child links, and also vary the period when income is measured. Since we were granted access to tax data on our cohort of parents only for the period 2006-2010 and on our cohort of children only for the period 2015-2019, we rely on the procedure laid out in Section 4.2 to measure formal income when running these robustness tests. The procedure sums up formal labor income based on formal employment data and predicted formal non-labor income based on our RF prediction model. Therefore, we first show our benchmark mobility estimates for running these robustness tests in column 1 of Table C.5, Panel A. The rank-rank slope is .453, somewhat smaller relative to our main estimates (.546).

### C.3.A Sample selection

**LARGER SAMPLES.** Our baseline sample comprises 1.3 million children born during the period 1988-1990 that we can link to both parents, comprising 15% of all children in such cohorts. We show that our main results are robust to expanding the sample along three dimensions. First, we include all children that can be linked to their father (regardless of whether they are linked to their mother), which increases the sample size by 1 million, and run the analysis solely based on the father’s income. The results in columns 2-3 of Table C.5, Panel A, show that the father-child rank correlations in the baseline and enlarged samples are nearly identical (.44 and .45), and they are also identical to the baseline rank-rank slope estimated without tax data, reported in column 1.

Second, we expand the sample to include all cohorts born in 1983-1992, for a total of 6.9 million children. In this case as well, the estimated rank-rank coefficient remains identical.

Finally, Panel B of Table C.5 replicates the analysis using the less conservative linking procedure described in Appendix Section A.2. This increases the data coverage for our main cohorts (1988-1990) from 15% to 45%. Once again, all estimated coefficients are virtually unaffected (0.44 - 0.47).

**REWEIGHTING.** To address any residual concern about the representativeness of our sample, we re-weight the data to match a rich set of characteristics (up to



**Table C.5:** Robustness to larger samples

	1988-1990 cohorts			1983-1992 cohorts		
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A. Baseline links</i>						
Parent Rank	0.459***			0.480***		
Father Rank		0.445***	0.455***		0.467***	0.472***
Observations	1,304,586	1,304,586	2,361,010	3,797,639	3,797,639	6,949,075
Only father links			Yes			Yes
<i>Panel B. Complete links</i>						
Parent Rank	0.459***			0.468***		
Father Rank		0.440***	0.448***		0.448***	0.456***
Observations	3,416,131	3,416,131	3,901,433	9,976,431	9,976,431	11,478,370
Only father links			Yes			Yes

*Notes:* The table reports the estimated slope of the rank-rank regressions in equation (1), i.e. our (inverse) measure of income mobility, in different samples. In Panel A, we link parents to children using our baseline, conservative procedure; in Panel B, we expand the sample using the less conservative procedure described in Section A.2. Columns 1-3 cover the 1988-1990 cohorts – as in our main sample – while columns 4-6 cover the 1983-1990 cohorts. Columns 1-2 and 4-5 are based on children who can be linked to both parents, while columns 3 and 6 are based on children who can be linked at least to their fathers. The dependent and explanatory variables are always the child and parental income percentile rank. For consistency, in all specifications we measure income without using tax data, which are only available for the 1988-1990 cohorts. Income for the 1988-1990 cohorts is measured between 2015-2019, at ages 25-31 (as of baseline), while income for other cohorts is measured when they are 25-29 years old (\*p<0.1; \*\*p<0.05; \*\*\*p<0.01).

their 2<sup>nd</sup> moment) in the general population, using the algorithm proposed by Hainmueller (2012). Specifically, we balance our baseline sample with respect to gender, race, month and year of birth, state of residence (27), state capital dummy, education (4), and indicators for being in social welfare registries, formal labor market participation, and having a unique name in the country.<sup>46</sup> In Table C.6, we report estimates of the rank-rank coefficient on the raw and re-weighted data (row 1 and 2, respectively) for the full sample (column 1), males (2), females (3), whites (4), and non-whites (5). Reweighting leads only to small changes in the estimated rank-rank slope in our full sample (from .546 to .566) or subgroups of the population.

<sup>46</sup>Table A.1 in Section A.3 displays descriptive statistics of the sample before and after reweighting together with standardized differences with respect to the population.

**Table C.6:** Reweighting Procedure

	Full Sample	Males	Females	Whites	Non-whites
	(1)	(2)	(3)	(4)	(5)
Baseline	0.546***	0.469***	0.613***	0.521***	0.573***
Weighted	0.566***	0.504***	0.624***	0.546***	0.574***
Observations	1,304,586	633,489	671,097	672,015	546,773

*Notes:* The table reports the estimated slope of the rank-rank regressions in equation (1), i.e. our (inverse) measure of income mobility, in the raw data (first row) and in the re-weighted data (second row). The re-weighted data match the first and second moments of the population of children in the same birth cohorts along several characteristics, using the entropy algorithm by [Hainmueller \(2012\)](#) (see Table A.1). All samples cover the 1988-1990 cohorts and the dependent and explanatory variables are the child and parental income percentile rank (\*p<0.1; \*\*p<0.05; \*\*\*p<0.01).

### C.3.B Timing of income measurement

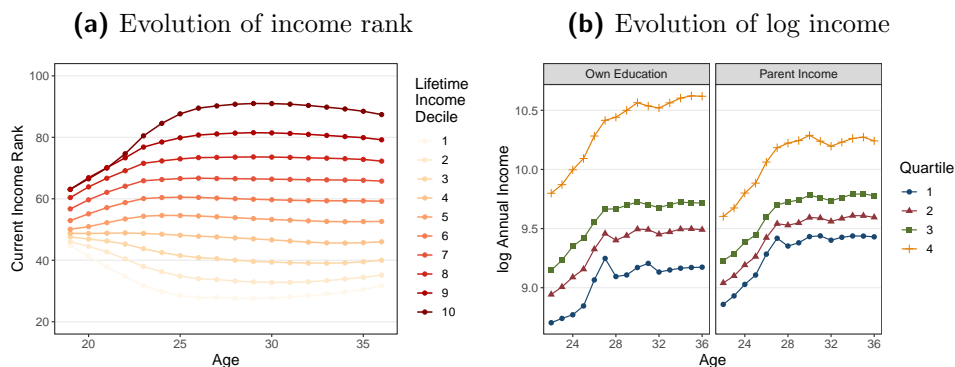
The timing when we measure income may lead to two main types of estimation bias, which we describe next. We provide several tests showing that our main results are not substantially affected by the timing of income measurement. We estimate income mobility without relying on tax data, as the latter are not available for some of the cohorts and years required for such tests.

**ATTENUATION BIAS.** Measuring income for short time spans may attenuate estimates of income mobility due to temporary income shocks. This is not a main concern when we measure parental income since we virtually cover children’s entire childhood (from age 3 to 18). However, this could be a more relevant issue for child income in our main analysis, which uses a five-year window (age ranges 25-29, 26-30, 27-31 for the cohorts born in 1988, 1989 and 1990, respectively). In light of this, we show in Figure C.1a how the rank-rank slope changes as we vary the number of years used to measure parental income. The estimates are remarkably stable regardless of how many years are used in the analysis. Moreover, in Figure C.1b we show that estimates also remain largely stable when using 1 to 5 years to measure children’s income. The estimates vary by less than 5 percentage points in both exercises, relative to the 0.453 rank-rank slope benchmark (without tax data). Overall, these results support the idea that the five-year window used to measure child income is sufficient to prevent meaningful attenuation bias in our main analysis.

**LIFE-CYCLE BIAS.** Measuring income too early may not adequately capture

permanent income, possibly leading to life-cycle bias (Haider and Solon, 2006; Mello et al., 2021). Again, this could be relevant when measuring children’s income with a five-year window at the age range 25-31 in our main analysis. We use the fact that we can track parental income for a long period of time to study how our estimates change when measuring parental income at different ages. In particular, we focus on parents in our main sample born in 1960-1965. Figure C.1c shows that our estimates do not vary much when using a 3-year window to measure father’s income centered from age 31 to 45. Next, in Figure C.1d, we show how the rank-rank slope changes as we center a three-year window around different ages for measuring children’s income. We focus on the 1988 cohort, for which we can track income up to age 31. Again, the estimated rank-rank slope remains fairly constant within cohorts when income is measured at varying ages.

This result can be explained by the fact that there are little positional changes in annual income in Brazil from early ages, especially from the age of 24 (Figure C.2a); this is due, in turn, to the fact that most Brazilians enter the labor market relatively early, given that college enrollment is low. Positional changes below the age of 24 are concentrated at the very top of the income distribution, and they are driven by a large share of high-income children who attend college and delay entry into the labor market (see also Figure 9a). Figure C.2b provides additional evidence on these aspects by showing a near parallel evolution of income by quartiles of completed education and parental income. These patterns are in contrast with the case of developed countries – as documented by Mello et al. (2021) for the US and Sweden – where a much larger share of individual attend college. In fact, Guvenen et al. (2022) document that Brazil displays the highest intragenerational persistence in income in a group of 13 middle- and high-income countries – for instance, 16% and 30% larger than the U.S. and Sweden respectively.

**Figure C.2:** Evolution of children’s position at the income distribution

*Notes:* The figure plots the evolution of children’s income distribution over time. Panel (a) shows the mean income percentile rank (on the vertical axis) when aged 18 to 36 (horizontal axis) for individuals in each decile of total lifetime income distribution. In turn, panel (b) shows the evolution of log incomes of the 1983 cohort when 22 to 36 years old, by quartiles of children’s educational level (left) and parental income (right).

### C.3.C Alternative income and occupation definitions

We now show that two potentially relevant choices that we take to define total annual income have virtually no impact on our IGM estimates. Table C.7 presents our baseline mobility estimate (column 1), along with alternative estimates that we describe next (columns 2-4). First, we rely on survey questions on “normal” monthly income to predict annual informal income and formal non-labor income (used when tax data are not available). In our main analysis, we extrapolated such income to the entire year, multiplying it by 12. We show that results do not change if we adopt a similar procedure for measuring formal labor income (derived from administrative employment data). Namely, we take the average monthly formal income while formally employed and multiply by 12 each year, instead of considering the sum over the year (column 2). Alternatively, we move back to our baseline but multiply predicted monthly informal income and formal non-labor income by the number of months that individuals spend out of formal employment in the year (rather than by 12) (column 3). Second, when predicting unobserved income in the main analysis, we label as informal workers in the administrative data those who do not hold any formal job in the entire year and who are not firm owners. We vary this assumption by defining informal workers as those who work formally for

less than three consecutive months in the year and who are not firm owners (column 4).

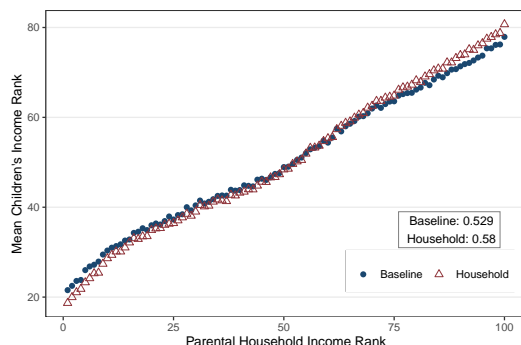
**Table C.7:** Alternative Income Measures

	(1)	(2)	(3)	(4)
Baseline	0.546*** (0.001)	0.545*** (0.001)	0.547*** (0.001)	0.547*** (0.001)
Income definition	Baseline	Alternative 1	Alternative 2	Baseline
Informal workers definition	Baseline	Baseline	Baseline	Alternative
Observations	1,304,586	1,304,586	1,304,586	1,304,586

*Notes:* The table reports relative mobility estimates based on the slope of rank-rank regressions – as in eq. (1) – using different income and occupation definitions. It presents estimates when using our baseline income and occupation definitions (column 1); when measuring formal labor income by multiplying its monthly average in each year by 12 (column 2), when measuring predicted informal and formal non-labor income by multiplying the predicted monthly quantities by the number of months out of formal employment in the year (column 3); when defining informal workers as those who are formally employed for less than three months in the year and who are not firm owners. All samples cover the 1988-1990 cohorts and the dependent and independent variables are the child and parental income percentile rank (\* $p < 0.1$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ ).

#### C.4 Individual vs. household income for child ranks

We assess whether using household instead of individual income for children affects our results. We focus on the 38% of children who can be linked to their spouses in tax declarations and welfare registries (*CadÚnico*). We compute our baseline income measures for spouses in the same calendar years that their partner’s income is measured starting from the first year when we observe both together. Household income is defined as the sum of both partners’ individual income. Figure C.3 plots the mobility curves using individual (baseline) and household income for children in the married sample.

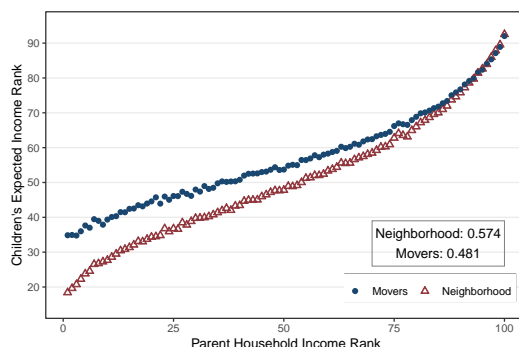
**Figure C.3:** Mobility curve using household income for children

*Notes:* The figure plots the mobility curves using individual (blue dots) and household (red triangles) income for children, while parental income is the sum of father and mother's income, as in our baseline. The sample comprises the individuals in the baseline sample for which we are able to recover partners from tax declarations or *CadÚnico*. Household income is the sum of both partners' individual annual income starting from the year we observe them as a couple. For each curve, the figure also reports our relative mobility measure based on Equation 1.

Both curves perfectly overlap and are similar to the baseline mobility curve based on our main sample (Figure 1). Overall, the exercise indicates that taking household or individual income has little impacts on our IGM estimates, which only become slightly larger, increasing from .529 to .58. They are also similar to our main estimates based on children's individual income, using the full sample rather than the married sample (.546).

### C.5 Neighborhood-based measure for movers

**Figure C.4:** Neighborhood-based measure: children changing address



*Notes:* This figure plots mobility curves based on the neighborhood-based measure, for our main sample (red triangles) and for a sample of movers who live in a different zip code than the place where they grew up (blue dots), both covering the 1988-1990 cohorts. For each parental rank in each curve, it plots the mean child rank. The neighborhood-based measure is given by the average formal income in the census tract where children grew up (parental rank) and where they live as adults (child rank). See Section 6.4 for a detailed description of these measures. For each curve, the figure also displays our relative mobility measure based on Equation (1).

### C.6 Additional results

**Table C.8:** Labor market differences by gender and race

Parent Quintile	Gender			Race		
	Rank Gap	LFP Gap (pp.)	Wage Ratio	Rank Gap	LFP Gap (pp.)	Wage Ratio
1	17.7	17.6	0.84	10.4	10.0	0.95
2	18.8	13.2	0.84	9.2	8.4	0.95
3	15.2	8.2	0.84	8.6	7.2	0.96
4	10.8	4.6	0.87	6.4	6.3	1.01
5	7.1	1.9	0.89	3.4	6.5	1.12

*Notes:* The table reports average gaps in child income ranks and labor market outcomes over gender and race, for each parental income quintile. Income rank gaps are calculated as the difference between average adult ranks for males (whites) and females (non-whites). The labor force participation (LFP) gap is the difference in average participation rate in the formal labor market between the two groups, in percentage points. Finally, the wage ratio is the ratio of the formal average wages of females (non-whites) to males (whites).

**Table C.9:** Siblings comparisons by parental income quintile

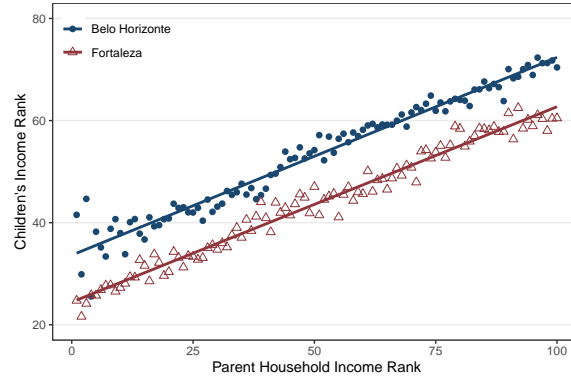
Parental Quintile	Siblings Gap	Brother-Sister Gap
1	0.11	16.23
2	0.22	19.00
3	0.19	16.67
4	1.39	11.79
5	2.36	7.29

*Notes:* The table reports average gaps in income ranks between siblings for each parental income quintile. Siblings gaps are calculated as the difference between adult income rank of the older and younger siblings, regardless of gender. Brother-sister gaps are calculated as the difference between the male and female siblings, regardless of birth order. Both are calculated for individuals in our baseline sample of the 1988-1990 cohorts.

## D Appendix to Section 7

### D.1 Individual mobility curves

**Figure D.1:** Individual mobility curves in Fortaleza and Belo Horizonte



*Notes:* The figure plots separate mobility curves for Fortaleza (red triangles) and Belo Horizonte (blue dots). Both curves are non-parametric binscatters constructed by plotting mean child income rank for children born in each parental percentile income rank in both regions. The figure is based on our main 1988-1990 cohorts sample. Income is defined as our baseline measure and both children and parents continue to be ranked according to the respective national income distribution. Children are assigned to regions according to the location of their parents in 2000, regardless of where they live in adulthood.

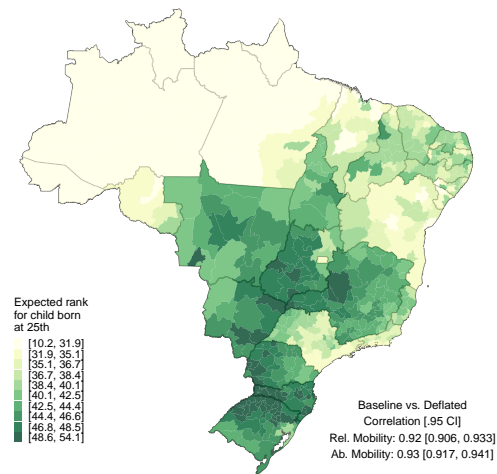
### D.2 Robustness of subnational estimates, price differences

We construct a regional price index to generate mobility maps that account for price differences across Brazil. To create the index, we use the POF (*Pesquisa de Orçamentos Familiares*), a household budget survey conducted by IBGE. POF gathers rich demographic and expenditure data at a fine geographic level. We use the 2003 and 2019 editions to calculate – across Brazilian areas – the average price of the reference basket used to compute the main Brazilian consumer inflation index (IPCA). Specifically, we compute prices at the state level, distinguishing between the (state) capital, the metropolitan areas around the capital, and the countryside. Next, we rescale parents' and children's income by the index computed for 2003 and 2019, respectively, according to their location. Finally, we re-estimate our regional mobility measures based on the price-adjusted income and correlate them with our main estimates. Figure D.2 shows that such adjustment has little impact on regional mobility patterns. It shows that absolute mobility remains similar across space, and that both



absolute and relative mobility are strongly correlated with our original, region-specific, mobility measures (see Figure 10). The correlations for our relative and absolute mobility measures are as high as .92 (.906,.933 - 95% C.I.) and .93 (.917,.941 - 95% C.I.), respectively.

**Figure D.2:** Price-Adjusted Absolute Mobility Map



*Notes:* The figure displays price-adjusted absolute mobility – scaled by deciles – in Brazil’s 510 immediate geographical regions (IGRs) for our main sample (1988-1990). Parent and child incomes are deflated by regional price indexes constructed with POF survey data and ranked in the national income distribution (measured when children are aged 3-18 and 25-31, respectively). Absolute mobility indicates the expected rank for children in below-median income families, based on Equation (1). Darker green tones indicate higher absolute mobility. Children are assigned to IGRs according to the location of their fathers in 2000.

### D.3 Mobility estimates for the 50 largest metropolitan areas

**Table D.1:** Summary of mobility estimates for the 50 largest metropolitan areas

Region	2021 pop. (thousands)	Slope	$E[y p = 25]$	$E[y p = 75]$	Q1Q1	Q1Q5
São Paulo, SP	22,049	0.66	34.3	67.4	44.4	1.7
Rio de Janeiro, RJ	12,901	0.55	35.6	63.0	30.1	1.6
Belo Horizonte, MG	5,348	0.50	39.3	64.4	23.2	2.2
Fortaleza, CE	4,179	0.49	33.1	57.8	45.8	2.0
Recife, PE	4,108	0.54	31.8	58.6	48.5	1.7
Salvador, BA	4,065	0.54	29.1	56.3	52.6	1.5
Curitiba, PR	3,732	0.53	41.7	68.3	17.9	2.4
Porto Alegre, RS	3,267	0.50	39.5	64.5	21.8	2.4
Campinas, SP	3,201	0.53	40.4	66.9	33.1	2.3
Distrito Federal, DF	3,094	0.54	36.7	63.6	41.5	2.4
Belém, PA	2,773	0.58	25.1	54.3	63.0	0.8
Goiânia, GO	2,628	0.49	42.4	66.8	24.4	3.8
Manaus, AM	2,605	0.57	24.5	53.0	58.1	1.4
Vitória, ES	2,100	0.48	40.1	64.2	25.4	3.1
Santos, SP	1,927	0.58	34.1	62.9	44.3	2.9
Sorocaba, SP	1,840	0.53	39.3	65.6	32.3	2.8
Natal, RN	1,734	0.49	33.9	58.6	44.9	1.7
São Luís, MA	1,657	0.46	32.9	55.8	49.6	2.1
Ribeirão Preto, SP	1,534	0.60	36.3	66.3	42.8	2.7
João Pessoa, PB	1,430	0.48	33.9	57.9	46.8	2.6
Maceió, AL	1,316	0.52	30.9	56.8	52.7	0.7
Feira de Santana, BA	1,242	0.44	31.1	53.0	57.9	2.4
Aracaju, SE	1,233	0.51	31.9	57.3	49.1	2.1
Florianópolis, SC	1,181	0.42	46.5	67.5	11.2	4.0
Campo Grande, MS	1,131	0.49	41.1	65.8	31.5	3.4
São José dos Campos, SP	1,125	0.54	36.3	63.2	37.3	2.2
Teresina, PI	1,116	0.48	37.1	61.1	43.8	2.8
Londrina, PR	1,114	0.44	43.9	65.7	25.7	2.8
Cuiabá, MT	1,105	0.48	40.3	64.2	24.5	2.3
Joinville, SC	1,044	0.46	46.6	69.7	15.7	2.0
Jundiaí, SP	973	0.53	41.2	67.7	37.5	9.1
Uberlândia, MG	959	0.43	43.3	64.7	30.5	4.4
São José do Rio Preto, SP	934	0.53	41.5	67.8	30.9	4.0
Novo Hamburgo - São Leopoldo, RS	908	0.48	42.0	65.9	17.8	2.3
Pelotas, RS	845	0.39	42.1	61.8	28.1	3.3
Caxias do Sul, RS	841	0.42	46.6	67.7	23.3	3.4
Maringá, PR	801	0.43	45.3	66.8	25.7	1.4
Montes Claros, MG	770	0.44	40.1	62.0	33.4	2.7
Juiz de Fora, MG	753	0.44	38.7	60.6	31.7	2.3
Macapá, AP	675	0.50	23.9	49.0	59.1	0.5
Bauru, SP	668	0.54	38.4	65.4	34.8	0.6
Volta Redonda - Barra Mansa, RJ	668	0.51	36.5	61.8	36.4	1.6
Porto Velho, RO	667	0.49	30.0	54.6	42.9	2.4
Campos dos Goytacazes, RJ	661	0.48	36.5	60.5	32.0	2.9
Ipatinga, MG	651	0.41	39.1	59.5	34.0	3.5
Ponta Grossa, PR	648	0.51	41.6	67.1	31.1	2.2
Taubaté - Pindamonhangaba, SP	637	0.48	37.1	61.3	26.1	0.8
Araraquara, SP	631	0.55	38.6	66.0	28.7	2.3
Piracicaba, SP	617	0.56	39.1	66.9	31.4	3.9
Santa Maria, RS	485	0.41	45.2	65.7	28.0	4.5

*Notes:* The table summarizes mobility estimates in the 50 largest metropolitan areas (IGRs) of Brazil, according to IBGE's population count in 2021. Mobility estimates are the rank-rank slope (relative), the expected income rank of below- and above-median income children (absolute), the bottom-bottom persistence probability (Q1Q1), and the bottom-to-top quintile transition probability (Q1Q5). Mobility measures are based on our baseline sample of the 1988-1990 cohorts. Children are assigned to IGRs based on the location of their parents in 2000.

### D.4 Mobility correlates

We explore the correlates of social mobility by estimating univariate regressions of absolute mobility on a wide range of local indicators covering thirteen broad categories: demographics, economic structure, education, family structure, health, household, income, inequality, local infrastructure, labor market,

municipal budget, public safety, and social capital. Table D.2 provides a detailed description and data sources for the variables in each category. Figure D.3 plots the results of these regressions when normalizing both the dependent and explanatory variables so that coefficients can be interpreted as correlations and more easily compared with each other.<sup>47</sup> Overall, coefficients have the expected sign, and nearly all of them are statistically significant. Several variables related to education quality show up among the top mobility predictors – in particular, literacy rates and students’ performance in standardized test scores. In line with the analysis by race in Section 6.6, the racial composition is also a strong mobility predictor: the share of white population displays the second highest correlation, while the share of black and mixed-race individuals yield negative coefficients. Other variables related to the number of formal firms per capita, the number of bank agencies, and labor market participation by men are also among the strongest mobility predictors. In turn, markers of socioeconomic struggle such as large or high density households, and the share of individuals without earnings are among the top predictors of low mobility, followed by the GDP share of the public sector.

One difficulty when interpreting these results is the strong correlation between the indicators considered. Thus, we reduce the dimensionality of the problem in two steps. First, we create a single index for each category based on the principal components of the initial variables, similarly to Acciari et al. (2021).<sup>48</sup> Figure D.4 report the results of multivariate regressions of absolute mobility on such indexes. Education quality yields the largest correlation with absolute mobility by far, with a positive sign (blue coefficients). Other categories showing strong correlation with mobility are the indexes related to the family structure, demographics (including the racial composition), household characteristics and local infrastructure. Once we control for region fixed effects (5 categories), the education index continues to stand out relative to other fac-

---

<sup>47</sup>Specifically, we recenter them around the mean and rescale them so that their standard deviation is equal to one.

<sup>48</sup>Specifically, for each group of variables, we keep the number of principal components needed to explain 90% of the variation in them. Subsequently, we compute the index as an average weighted by the amount of variation each component absorbs.

tors, being by far the strongest mobility predictor within regions (light blue coefficients).

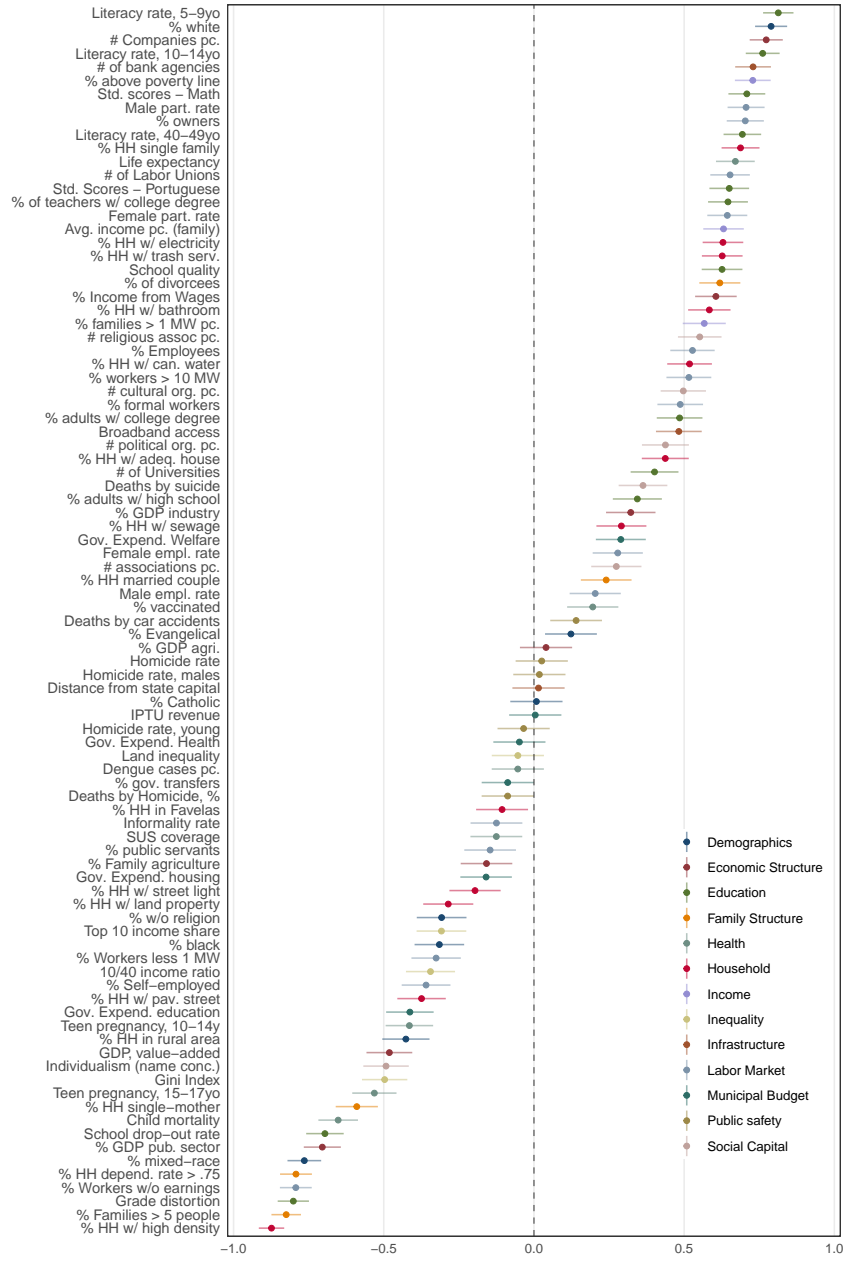
Nonetheless, the regressors in Figure D.4 still have a high degree of multicollinearity, which could harm the interpretation of the results. Hence, we employ a standard LASSO regularization procedure to select robust predictors of mobility. We validate the choice of parameters for the LASSO regression via ten-fold cross-validation. Figure D.5 summarizes the results of such an exercise, displaying the regressors' coefficients (y-axis) against increasing values of the regularization parameter (x-axis). As we increase the penalization for the number of regressors, coefficients are shrunk toward zero and variables leave the model. Again, the education index dominates the other variables.

**Table D.2:** List of municipal socioeconomic indicators and data sources

Group	Indicator	Year	Source
Demographics	% evangelical, % catholic, % w/o religion, % black, % mixed-race, % white, % HH in rural area	2000	IBGE
Economic Structure	% GDP agriculture, % GDP pub. sector, % GDP industry, (value-added), companies pc, % income from wages, % family agriculture	2000	IBGE
Education	Test scores (portuguese and math), % teachers w/ college degree, school quality, drop-out rate, grade distortion, literacy rate (5-9yo, 10-14yo, and 40-49yo), % of adults w/ high school	2000/2005	IBGE/Inep
Family Structure	% families > 5 people, % HH single mother, % HH married couple, % HH dependency rate > .75, % of divorcees	2000	IBGE
Health	Dengue cases pc, teen pregnancy (10-14yo and 15-17yo), child mortality, life expectancy, % of vaccinated, SUS coverage	2000	DataSUS/IBGE
Household	% HH in favelas, % HH single family, % HH w/ land property, % HH w/ trash service, % HH w/ paved street, % HH w/ bathroom, % HH w/ piped water, % HH w/ electricity, % HH w/ street light, % HH w/ adequate housing, % HH w/ sewage, % HH w/ high people/room density	2000	IBGE
Income	Average family income pc, % families above the poverty line, % families earning more than 1 MW	2000	IBGE
Inequality	Gini Index, top 10/bottom 40 income ratio, top 10 income share, land inequality	2000	IBGE
Infrastructure	Broadband access, distance from state capital, of bank agencies	2000/2007	ANATEL/IBGE
Labor Market	% of workers > 1 MW, % of workers w/o earnings, % of workers > 10 MW, female/male participation rate, female/male employment rate, % public servants, % firm owners, % self-employed, % formal workers, % employees, informality rate, of labor unions	2000	IBGE
Municipal Budget	Government spending (health, welfare, education, housing), % of federal government transfers, IPTU revenue (property-tax)	2000	FINBRA
Public Safety	Homicide rate (total, males, young), % of deaths by homicide, % of deaths by car accident	2007	Ipeadata
Social Capital	Religious associations pc, cultural organizations pc, political organizations pc, civil associations pc, % of deaths by suicide	2000	IBGE

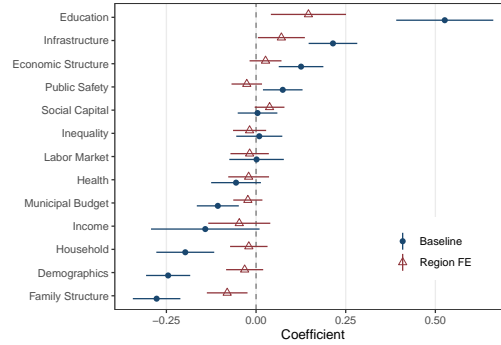
*Notes:* The table list all indicators used in the mobility correlates analysis, along with their source, year, and category group (used in the principal components analysis). All of them are obtained at the municipal level and then aggregated to immediate region level by population-weighted averages.

Figure D.3: Mobility Correlates



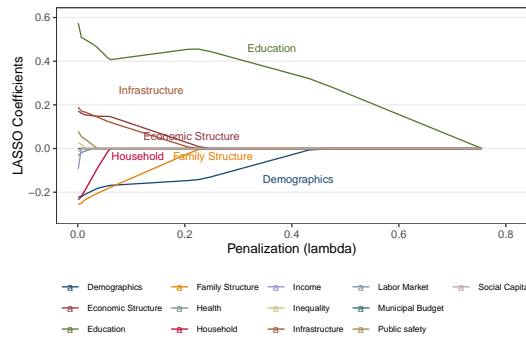
Notes: The figure summarizes a series of cross-regional univariate regressions of absolute mobility on a series of demographic, political, and socioeconomic indicators. The horizontal axis marks coefficients and .95 confidence intervals for each indicator, which are labelled in the vertical axis. Both dependent and independent variables are normalize so coefficients can be interpreted as straightforward correlations. Indicators are colored according to broad categories.

**Figure D.4:** Mobility Correlates: Principal Components



*Notes:* The figure reports the results of two multivariate regressions between absolute mobility and principal components of regional characteristics. The horizontal axis marks coefficients and .95 confidence intervals for each indicator, which are labelled in the vertical axis. The blue dots includes only the principal components and the red triangles include region fixed-effects (North, Northeast, Center-West, Southeast, and South).

**Figure D.5:** Mobility Correlates: LASSO Regularization



*Notes:* The figure plots the results of LASSO regularization for the correlation between absolute mobility and indexes constructed from socioeconomic indicators. The horizontal axis plots the regularization parameter lambda, while coefficients of each index are represented in the vertical axis. As lambda grows, the penalization for the number of regressors grow and coefficients are shrunk towards zero.

## E Appendix to Section 8

### E.1 Sample construction

For the analysis on causal place effects, we focus on children born in the 1983-1992 period who we can link to their fathers (following the baseline, conservative method described in Section 4). For every father in this sample, we retrieve all regions in which they worked during the 1992-2019 period using (formal) employment data (RAIS). We focus on the latter to track moves since address updates in the Brazilian person registry are largely incomplete before

2000. Although this choice implies that our sample is more representative of formal workers, doing so allows us to track migration more precisely, which is crucial for this analysis.

We define a mover as someone leaving a job in region  $a$  and taking a job in a new region  $b$  for at least two years. To increase precision, we focus on fathers showing up in employment data for at least five years in the 1993-2019 period.<sup>49</sup> Fathers who never move are defined as permanent residents of their regions. In turn, movers are those who move at least once. To simplify the analysis, we focus on families moving only once. Our final sample comprises 3,172,145 children and 2,260,645 fathers, with around 18% of them being movers.

## E.2 Defining the predicted outcomes of permanent residents

We closely follow the research design and specifications in [Chetty and Hendren \(2018a\)](#) and [Deutscher \(2020\)](#). First, we characterize outcomes of permanent residents of each region  $m$  and cohort  $c$  by running several rank-rank regressions of the type, for each region and cohort:

$$y_{imc} = \alpha_{mc} + \beta_{mc}p_{imc} + \epsilon_{imc} \quad (\text{E.5})$$

where  $y_{imc}$  denotes the income percentile rank at the age of 24 of a child from cohort  $c$  who spent her entire childhood in region  $m$ . We focus on income at the age of 24 that we can measure for all cohorts (1983-1992) in our sample. To ensure precision, we keep only region-cohort pairs for which we have at least 400 observations. We then calculate the predicted income rank of residents for every parental income rank  $p$ , region  $m$ , and cohort  $c$ :  $\hat{y}_{pmc} = \hat{\alpha}_{mc} + \hat{\beta}_{mc} \times p$ .

## E.3 Parametric specification and family fixed effects

Our baseline specification (3) includes nearly 180 thousands of fixed effects  $\alpha_{ocpa}$ . While they ensure that we exclusively compare very similar children (with the same origin, cohort, parental income decile and age at move) to estimate place effects, they also strongly restrict the variation used in the analysis. Consequently, they leave little space for adding family fixed effects, which also

---

<sup>49</sup>Our results remain similar when varying this threshold to ten or fifteen years. We use the five-year cut-off to enlarge the final sample and enhance precision.



strongly restricts the variation in the analysis. Hence, we follow [Chetty and Hendren \(2018b\)](#) and rely on a less restrictive, parametric specification to assess the robustness of our findings to family fixed effects:

$$y_i = \sum_{a=1}^{33} b_a I_a(a_i = a) \Delta_{odpc} + \sum_{c=1983}^{1991} \kappa_c I_c(c_i = c) \Delta_{odpc} + \sum_{a=1983}^{1992} I_c(c_i = c) (\eta_c^1 + \eta_c^2 \hat{y}_{poc}) + \sum_{a=1}^{33} I_a(a_i = a) (\zeta_a^1 + \zeta_a^2 p_i) + \lambda_f + \epsilon_i \quad (\text{E.6})$$

Rather than controlling for fixed effects  $\alpha_{ocpa}$ , this specification linearly controls for the quality of origin – which is allowed to vary by parental income and cohort – and age at move by parental income, accounting for the disruption effects of moving at different ages. Specifically, the first term in the second line is defined by cohort fixed effects  $\eta_c^1$  and an interaction between the cohort dummies  $\eta_c^2$  and the quality of origin  $\hat{y}_{poc}$ , modeled as the predicted income of permanent residents at origin  $o$ . In turn, the second term is defined by age at move fixed effects  $\zeta_a^1$  and age at move dummies  $\zeta_a^2$  interacted with parental percentile rank,  $p_i$ . Finally, the specification controls for family fixed effects  $\lambda_f$ , ensuring that causal place effects solely rely on variation across siblings.

#### E.4 Overidentification tests

The next results confirm that children’s outcomes converge to those of permanent residents with precisely the same age, gender, and race, while coefficients on other groups are generally an order of magnitude smaller, close to zero, and statistically insignificant. In addition, since children’s outcomes in different areas not only differ at the mean but over the entire distribution, we show that movers’ outcomes track different moments of the distribution of permanent residents’ outcomes beyond the mean. For instance, two areas may have the same mean child rank for low-income children but different probabilities that children end up in the top decile of the income distribution. These tests address additional concerns such as the possibility that moves to better places are driven by different shocks producing positive effects on children that decrease with age, e.g., positive income or wealth shocks. Specifically,

they indicate that following these shocks, parents would need highly accurate knowledge to select better places for our results to be driven by selection. Accordingly, for them to drive our main findings, parents would need to select places that offer better opportunities for children from the same cohort, gender and race. Finally, the potential shocks driving such selection process would need to replicate not only the mean outcomes but also the distribution of outcomes for children in the destination.

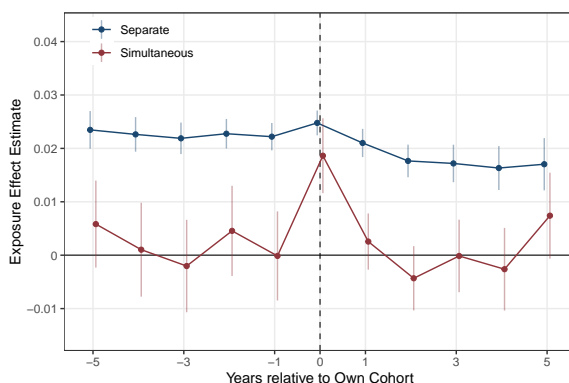
We start by showing that place effects are cohort-specific. The blue line in Figure E.1 displays the estimates for the exposure rate obtained from eleven separate regressions. In each of them, we replace the main independent variable – the predicted difference in outcomes for children of the *same cohort* – by the predicted difference for *other* cohorts, born from five years before to five years after. The coefficients obtained using adjacent cohorts are quite similar to the baseline, as regions with better opportunities for a given cohort usually are also good for other cohorts. In turn, the red line in Figure E.1 plots estimates when all cohort-specific predictions are simultaneously included in a single regression. Conditional on the predicted outcomes of their *own* cohort, all other cohorts’ predictions are statistically insignificant, while the true cohort coefficient approaches the baseline estimate. Hence, children’s outcomes converge to the outcomes of permanent residents of their *own* cohort and other cohorts’ outcomes have little explanatory power. Thus, any omitted variable possibly driving our results would have to precisely emulate cohort-specific place effects.

In Table E.1, we conduct an analogous exercise for gender and race. For this purpose, we construct gender-specific predictions and estimate exposure rates in three different ways: using the predicted outcomes of the child’s own gender (column 1), the opposite gender’s prediction (column 2), and both together (column 3). Like the blue line in the previous exercise, both regressors yield statistically significant estimates, since there is a considerable correlation between outcomes for boys and girls within regions. Nonetheless, the one based on the child’s own gender has higher explanatory power. Column 3 replicates the red line in Figure E.1: controlling for the own gender’s prediction, the

coefficient on the opposite gender (placebo) is negligible. Columns 4-6 conduct the same exercise for race, showing similar results and supporting our main analysis.

All estimates up to now have been based on the predicted differences in *mean* outcomes across locations. Now we show that place effects also replicate permanent residents’ outcomes along the income distribution. We construct permanent residents’ predictions for the probability of being in the top and bottom deciles of the national income distribution in adulthood. Subsequently, we estimate exposure rates contrasting the distributional predictions with the mean prediction (placebo). Columns 1-2 in Table E.2 show that the top ten probability is better explained by the distributional prediction than by the mean prediction when running separate regressions. In column 3, a simultaneous regression yields a significant coefficient for the distributional prediction, while the coefficient on the placebo is zero. Columns 4-6 replicate 1-3 but for the bottom ten probability, with equivalent results. Thus, the distribution of children’s incomes converges to the distribution of incomes in the destination in proportion to exposure time.

**Figure E.1:** Placebo tests: Cohort-specific convergence



*Notes:* This figure presents estimates of the annual childhood exposure effect on children’s income ranks in adulthood using permanent resident predictions for the child’s own birth cohort and surrounding “placebo” birth cohorts. The series in blue plots estimates of the exposure effect  $\gamma_t$  from nine separate regressions, using permanent resident predictions from cohort  $c + t$  (where  $t$  ranges between -5 and 5) as the key independent variables and the outcomes of children in birth cohort  $c$  as the dependent variable. The series in red plots estimates from a single multivariate regression that simultaneously includes all nine permanent resident predictions  $t = 5, \dots, -5$ .

**Table E.1:** Placebo test: Gender- and race-specific convergence

	Exposure effect $\gamma$					
	Gender			Race		
	(1)	(2)	(3)	(4)	(5)	(6)
Own Group	0.025*** (0.001)		0.025*** (0.002)	0.022*** (0.001)		0.022*** (0.002)
Opposite Group		0.020*** (0.001)	-0.000 (0.001)		0.017*** (0.001)	0.000 (0.002)
Observations	285,912	285,912	285,912	267,614	267,614	267,614

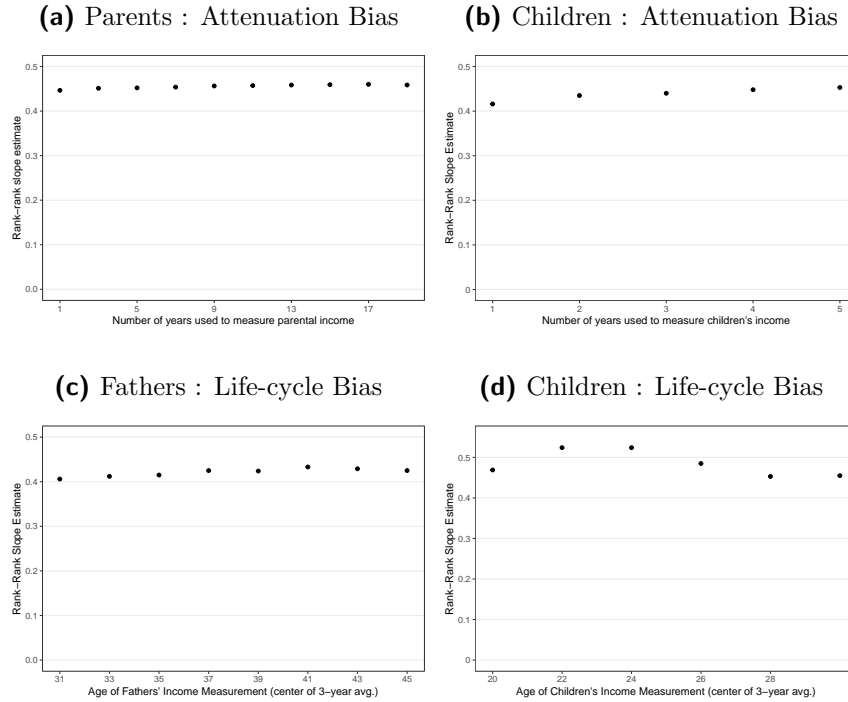
*Notes:* The table reports estimates of annual childhood exposure effects  $\gamma$  using gender- (Columns 1-3) and race-specific (Columns 4-6) permanent resident predictions. In all columns, the dependent variable is the child's family income rank at the age of 24. In both panels, column 1 (4) replaces the predicted outcomes based on all permanent residents in the origin and destination with predictions based on the outcomes of children who have the same gender (race) as the child who moves. Column 2 (5) replicates column 1, replacing the own-gender (race) predicted outcomes with the predicted outcomes of the opposite gender (race). Column 3 (6) combines the variables in columns 1 and 2, including both the own-gender (race) and placebo other-gender (race) predictions (\*p<0.1; \*\*p<0.05; \*\*\*p<0.01).

**Table E.2:** Placebo test: Distributional convergence

	Exposure effect $\gamma$					
	Upper Tail			Lower Tail		
	(1)	(2)	(3)	(4)	(5)	(6)
Distributional prediction	0.031*** (0.002)		0.034*** (0.002)	0.021*** (0.003)		0.022*** (0.003)
Mean Rank prediction		0.000*** (0.000)	0.000 (0.000)		0.000*** (0.000)	0.000** (0.000)
Observations	285,912	285,912	285,912	285,912	285,912	285,912

*Notes:* This table reports estimates of annual childhood exposure effects  $\gamma$  for upper- and lower-tail outcomes: being in the top or bottom 10% of the cohort-specific income distribution at the age of 24. Column 1 reports estimates from a regression of an indicator for being in the top 10% on the difference between permanent residents' predicted probabilities of being in the upper tail in the destination vs. the origin. Column 2 replicates column 1 but uses the difference between permanent residents' predicted *mean* ranks on the right-hand side of the regression. Column 3 includes both the (distributional) and the mean rank prediction. Columns 4-6 replicate columns 1-3 using an indicator for being at the bottom 10% at the age of 24 as the outcome. In all columns, the sample comprises all children in the primary analysis sample of one-time movers (\*p<0.1; \*\*p<0.05; \*\*\*p<0.01).

**Figure C.1:** Sensitivity of child and parental income to timing



*Notes:* This figure plots robustness exercises for attenuation bias (a, b) and life-cycle bias (c, d). In Panels (a) and (c) child income is held constant and measured as in our baseline estimates, and we vary how parental income is measured. In Panels B and D, we measure parental income as in our baseline and vary how child income is measured. Panel A displays estimates of the rank-rank slope from separated rank-rank regressions in which we vary the number of years used to compute parental income, from 1 to 17 years, and centered at the age of 11. Panel (b) displays an analogous exercise in which we measure children's income using from 1 to 5 years, centered at age 27. In Panel (c), we run rank-rank estimates using father's income (rather than parental income) and vary the age when father's income is measured using a three-year window from ages 31 to 45. Finally, in Panel (d) we vary the age when we center the three-year window to measure children's income, from 20 to 30 years old. In Panels (a) and (b) we use our full baseline sample of the 1988-1990 cohorts. In Panel (c), we restrict the sample to children whose fathers are born between 1960-65 and focus on fathers' rather than parental income – to precisely gauge the sensitivity concerning different age windows. In Panel (d), the working sample is the 1988 cohort since income data is only available until 2019.