

# Measuring Cross-Country Differences in Misallocation\*

Mitsukuni Nishida<sup>\*</sup>, Amil Petrin<sup>\*\*</sup>, Martin Rotemberg<sup>†</sup>, and T. Kirk White<sup>††</sup>

<sup>\*</sup>Johns Hopkins Carey Business School

<sup>\*\*</sup>University of Minnesota

<sup>†</sup>New York University

<sup>††</sup>Center for Economic Studies, U.S. Census Bureau

October 31, 2016

## Abstract

In this paper, we discuss the role that data processing and collection have for the measurement of misallocation. First, we turn to the raw self-reported data for the US, reflecting what can be found in most developing countries. In the raw data, measured misallocation (following Hsieh and Klenow 2009) is substantially higher than for any other country for which we have census data. For instance, if Indian firms had the same dispersion of distortions as measured in the reported US data, TFP in the Indian manufacturing sector would decrease by around  $\frac{2}{3}$ . Second, we follow a different strategy for editing and imputing missing data than what is used by the US Census Bureau, by using a method that seeks to replicate the true variance in the underlying data generating process known as Classification and Regression Trees (CART). This change raises the potential gains from removing misallocation in the United States manufacturing sector by around 10%.

---

\*Some of the research in this paper was conducted while the fourth author was an employee of Census Bureau. Any opinions and conclusions expressed herein are those of the authors and do not necessarily represent the views of the U.S. Census Bureau. All results have been reviewed to ensure that no confidential information is disclosed.

## I Introduction

The puzzle of large measured cross-country differences in productivity has recently been attributed to differences in within-industry misallocation of factors. However, unlike inputs and outputs, misallocation cannot be observed directly. In order to measure the extent of misallocation, researchers must undertake two steps. First, assumptions must both be made about firm behavior in the absence of distortions, and how to use observed behavior in order to estimate the size and magnitude of existing distortions. Second, the measured distortions are plugged back into the model to calibrate what the aggregate gains would be under different counterfactuals. For instance, the seminal paper of Restuccia and Rogerson (2008) develops a framework under which greater misallocation of resources leads to more dispersion in the distribution of plant-level total factor revenue productivity (TFPR). Hsieh and Klenow (2009) show that under relatively standard assumptions, within-industry variation in the revenue shares of each input is evidence of idiosyncratic firm-input distortions, and provide a simple algorithm for calculating the productivity gains from equalizing those distortions across firms. Taking their model to data, they find that “moving to U.S. efficiency would increase TFP by 30%–50% in China and 40%–60% in India.”

However, as Hsieh and Klenow (2009) note, measurement error may look to the researcher like misallocation of resources. Firms who report inaccurate information may only spuriously appear to be using a socially inefficient quantity of resources. The converse is true as well, since firms may report values which are in line with the model but do not reflect reality on the ground. As a result, the confidence we have in our measures of misallocation - either measurements of “true” values for a particular country, or of cross-country differences - depends on the extent of measurement error. In this paper, we discuss two potential sources of measurement error: firms who potentially misreport

their own characteristics, and subsequent data cleaning which potentially removes actual distortions. We show that both, in particular the former, are dramatically important: if instead of using the US Census Bureau’s cleaned data we use the raw information self-reported by each establishment, we find that moving to the new measured U.S. efficiency would *decrease* measured TFP by around  $\frac{2}{3}$  for both India and China.

We focus on the efforts undertaken by national statistics agencies when firms do not report (or report unlikely) information. Most statistics agencies, as a first pass, ask firms to verify (or send in) the information, but the next steps vary. Unlike its Indian and Chinese counterparts, the U.S. Census Bureau both edits and imputes responses.<sup>1</sup> The exact procedures vary across industries and time (for more information, see White et al. 2015), but broadly take two forms. First, the Census Bureau *edits* information that seems too far away from the other information reported by the firm or plant. If a reported variable fails one or more edit rules, then it is temporarily replaced with a missing value. Second, the Bureau *imputes* missing information, using other information reported by the plant (both in that year and in previous years) and other plants in the same industry.<sup>2</sup>

Furthermore, the Census Bureau also has administrative records from the IRS for payroll and wages (from Form 941, which is explicitly the information asked for in the census form). For those variables, the Census Bureau edits entries using this alternative source of firm-reported data. This information is unavailable for statistics agencies in many developing countries due to their low tax base (Jensen, 2016). For 2002 and 2007, we have access to the original values reported by firms for plants in the Census of Manufactures.<sup>3</sup>

---

<sup>1</sup> We have confirmed this both in the documentation for the data, and in email communications with the relevant national statistics agencies.

<sup>2</sup> Firms who have a variable edited have that variable imputed as if the firm had not reported anything. Note that the imputed data must also pass the editing rules. For most establishments, at least one of the variables needed to calculate TFP is imputed. For payroll and number of employees, the Census Bureau uses administrative records (mainly IRS payroll data) to replace reported data that fails edit rules. The Census Bureau classifies these changes from the reported data as “non-imputes”. However, these non-imputes still change plants’ measured TFP.

<sup>3</sup> It is worth noting that when HK was written, neither imputation flags nor this data were available for the

This allows us to know exactly which entries were imputed or entered in the cleaned Census data.<sup>4</sup> In order to focus our attention on the role of measurement, we follow the Hsieh and Klenow (2009) model exactly in order to remeasure the gains from reallocation.<sup>5</sup>

The HK insight is as follows: with CES demand and a constant returns to scale production function, revenue productivity (TFPR) is equalized across firms in the absence of distortions (regardless of any underlying variation in quantity productivity). That said, in most if not all firm-level or plant-level datasets, there is substantial within-industry variation in TFPR. Hsieh and Klenow (2009) rationalize those differences with idiosyncratic distortions on the firm-specific prices for capital and output. Each firm's distortions can be calibrated using the firm's first-order conditions. They then use the model to generate an elegant expression for the potential gains from reallocation from equalizing TFPR across firms. Other models of misallocation have similar features (Banerjee and Duflo, 2005; Restuccia and Rogerson, 2008; Hopenhayn, 2014).

Taking the model as given, we have two main goals in this paper. Our first goal is to document how much measurement matters for estimates of cross-country differences in misallocation. To this end, we treat the US data as we effectively treat data in other settings: we use the raw (plant-level) data reported by firms, with no additional imputation or editing.<sup>6,7</sup> In both 2002 and 2007, in the raw data the US appears to have substantially

---

Census years used in their study (1977-1997).

<sup>4</sup> Researchers have almost exclusively used the cleaned data for studies on manufacturing in the U.S.

<sup>5</sup> A growing literature has built on Hsieh and Klenow (2009) model by adding in additional features of firm behavior (such as dynamic considerations, as in Asker et al., 2014 and Foster et al., 2016, or entry & exit, as in Jaef (2016), Peters (2013), and Yang (2012). Regardless of what the true model of the world is, our results are unambiguous that measured cross-country differences in productivity are sensitive to features of data collection and cleaning. The Hsieh and Klenow (2009) approach is a convenient vehicle for showing how large the range of measured variation is in the data since it has been applied in so many settings.

<sup>6</sup> If a firm does not report at least one of the variables needed to estimate misallocation, it is not used in the estimation. As a result, we do not need to focus our attention on imputation for this part.

<sup>7</sup> Hsieh and Klenow (2009) implicitly use a similar strategy for the measurement of capital. Instead of using cumulative depreciated investment as their measure of firm capital, which would be the norm in the literature, they use the book value of capital, since that is the only such variable available in their cross-sectional Indian data, and for Census of Manufacturing firms who are not in the Annual Survey of

more misallocation than India and China. We do not take this result literally - we do not think that we have compelling evidence that the US manufacturing sector is characterized by more misallocation than most other countries. Instead, we consider our results a “smoking-gun” that measurement (and data processing in particular) is deeply important to the study of misallocation.

We do this by using some information from the clean data, while using other information from the raw data. First, we replace one variable from the final data for its uncleaned counterpart. Doing this for capital does not substantively affect measured misallocation, but value added and labor do matter. For instance, using cleaned values for capital and payroll, but self-reported value added in 2007 completely eliminates the measured misallocation gap between the US and India, with a similar consequence if instead the only raw data we use is for payroll. Both for 2002 and 2007, using raw data both for value added and payroll raises measured misallocation in the US to multiples of the values for India (for this exercise, which variable we use for capital matters quantitatively little).

One concern with the usefulness of our exercise is that other institutional details unique to the US might be overstate true misallocation. Most saliently, the Census form asks firms the same value of plant-level payroll that the firm reported to the IRS. This might be a complicated request for firms with multiple plants who share one tax ID, since perhaps the firm has no administrative records on plant-specific payroll. However, even within single-plant firms, 13% report different values on the Census form than the IRS payroll data for the same plant, and 7% report values that differ by more than 10%. Furthermore, measured misallocation in the US is even larger if we only use the self-reported labor variable for single-plant firms.

To that end, we explore why measured US misallocation is so sensitive to how the data is cleaned. We focus on two types of tests. First, we show that selection does not play a

---

Manufacturers sample. We are essentially applying this logic to data cleaning.

particularly large role: measured misallocation barely increases when using the cleaned data on only on the sample of plants with fully-reported information. Next, we consider the role of each of the plants' characteristics.<sup>8</sup>

Our second goal is to provide an alternative measure of the gains from reallocation in the United States. For this task, we focus our attention on the data generated from the Census Bureau's mean and regression based imputation strategies which are used to clean several of the components of value-added. We instead impute information following a technique from the epidemiology literature, a non-parametric multiple imputation strategy known as sequential CART (Burgette and Reiter, 2010; White et al., 2015). This method is designed to approximate the conditional distributions of the variables being imputed, which is crucial for measuring dispersion. Using this approach raises the measured gains from reallocation in the United States in 2002 from 44% to 54%, and in 2007 from 52% to 60%. While we are working on generating appropriate country/industry-specific edit and imputation rules, we will never be able to get administrative tax records for payroll in places where it does not exist, so we cannot fully replicate this procedure in other settings.

In the next section, we recap the theory of distortions underlying our analysis. Section 3 discusses the extant data collection and cleaning procedures in the United States. In Section 5, we use the raw data for the US to discuss the role of measurement for understanding cross-country differences in misallocation. In Section 4, we use CART imputation for just the US data. Section 6 concludes.

---

<sup>8</sup> We also consider increasingly aggressive trimming of the data, since it is possible that measured misallocation is driven by the tails. While trimming outliers does lower measured misallocation, we show that the relative effect of the raw data is similar with more outliers removed, and by request are happy to show tables that the relative effect of trimming is similar across countries.

## II A Theory of Misallocation

In this section, we briefly describe the Hsieh and Klenow (2009) approach to measuring misallocation that we follow. First, we start from the firm-side of the problem, showing how firm behavior is affected by idiosyncratic distortions on capital and output. In the model, variation in those distortions is captured by variation in firm-level revenue productivity. We then turn to the aggregate side, and derive how aggregate productivity would be affected in a counterfactual where the variation in revenue productivity were removed.

### II.A Firm-level Distortions

Overall utility  $Y$  is a Cobb-Douglas aggregate over sectoral output  $Y_s$ ,

$$Y = \prod Y_s^{\theta_s},$$

so normalizing the price of the final good to 1, expenditure for each sector is a fixed proportion

$$P_s Y_s = \theta_s Y$$

where  $P_s$  is the price index for sector  $s$ .

Within each sector, output takes a CES form over output of each variety  $Y_{si}$ :

$$Y_s = \left( \sum_{i=1}^M Y_{si}^{\frac{\sigma-1}{\sigma}} \right)^{\frac{\sigma}{\sigma-1}}$$

and each firm produces value added using capital and labor, with Cobb-Douglas production-function elasticities which vary across sectors:

$$Y_{si} = A_{si} K_{si}^{\alpha_s} L_{si}^{1-\alpha_s}.$$

The wage and rental rate are constant in the economy, but firms face idiosyncratic distortions on output and capital. As a result, each firm's profits are:

$$\pi_{si} = (1 - \tau_{Y_{si}}) P_{si} Y_{si} - w L_{si} - (1 + \tau_{K_{si}}) R K_{si}.$$

Marginal revenue productivity for each input is

$$\begin{aligned} MRPL_{si} &= \frac{\sigma - 1}{\sigma} (1 - \alpha_s) \frac{P_{si} Y_{si}}{L_{si}} \\ MRPK_{si} &= \frac{\sigma - 1}{\sigma} (\alpha_s) \frac{P_{si} Y_{si}}{K_{si}}. \end{aligned}$$

and each firm's revenue productivity is

$$TFPR_{si} = \frac{P_{si} Y_{si}}{K_{si}^{\alpha_s} L_{si}^{1-\alpha_s}} = P_{si} A_{si} \propto MRPL_{si}^{1-\alpha} MRPK_{si}^{\alpha}. \quad (1)$$

## II.A.1 Optimization Behavior

Profit maximization implies that:

$$\begin{aligned} P_{si} &= \frac{\sigma}{\sigma - 1} \left( \frac{R}{\alpha_s} \right)^{\alpha} \left( \frac{w}{1 - \alpha_s} \right)^{1-\alpha} \frac{1}{A_{si}} \frac{(1 + \tau_{K_{si}})^{\alpha_s}}{(1 - \tau_{L_{si}})} \\ A_{si} &\propto \frac{(P_{si} Y_{si})^{\frac{\sigma}{\sigma-1}}}{K_{si}^{\alpha_s} L_{si}^{1-\alpha_s}} \end{aligned} \quad (2)$$

$$\begin{aligned} w L_{si} &= (1 - \tau_{Y_{si}}) \frac{\sigma - 1}{\sigma} (1 - \alpha_s) P_{si} Y_{si} \\ \Rightarrow MRPL_{si} &= \frac{w}{(1 - \tau_{Y_{si}})} \end{aligned} \quad (3)$$

$$\begin{aligned} (1 + \tau_{K_{si}}) R K_{si} &= (1 - \tau_{Y_{si}}) \frac{\sigma - 1}{\sigma} (\alpha_s) P_{si} Y_{si} \\ \Rightarrow MRPK_{si} &= \frac{R (1 + \tau_{K_{si}})}{(1 - \tau_{Y_{si}})} \end{aligned} \quad (4)$$

As a result, combining Equations 1, 3, and 4 gives

$$TFPR_{si} \propto \frac{(1 + \tau_{K_{si}})^\alpha}{(1 - \tau_{Y_{si}})}, \quad (5)$$

so revenue productivity is only a function of the distortions, and not directly a function of firm TFP. As a result, in the absence of distortions, TFPR would be equalized across firms. In the next subsection, we show how variation in TFPR affects aggregate productivity

## II.B Aggregate Distortions

Aggregate productivity in each sector is

$$TFP_s = \frac{Y_s}{K_s^{\alpha_s} L_s^{1-\alpha_s}} = \frac{\overline{TFPR}_s}{P_s}, \quad (6)$$

where, given cost-minimization, the price index for sector  $s$  is:

$$P_s = \left( \sum_{i=1}^M P_{si}^{1-\sigma} \right)^{\frac{1}{1-\sigma}}.$$

From Equation 1, we can rewrite the price index as

$$P_s = \left( \sum_{i=1}^M \left( \frac{A_{si}}{TFPR_{si}} \right)^{\sigma-1} \right)^{\frac{1}{1-\sigma}},$$

and plugging back in to Equation 6 gives the core Hsieh and Klenow (2009) expression for productivity

$$TFP_s = \left( \sum_{i=1}^M \left( A_{si} \frac{\overline{TFPR}_s}{TFPR_{si}} \right)^{\sigma-1} \right)^{\frac{1}{\sigma-1}}, \quad (7)$$

Since we know from Equation 5 that  $TFPR_{si}$  would only be different from  $\overline{TFPR}_s$  in the presence of distortions, the “efficient” counterfactual TFP is  $\bar{A}_s = \left( \sum_{i=1}^M A_{si}^{\sigma-1} \right)^{\frac{1}{1-\sigma}}$ ,

and so (aggregating over all sectors)

$$\frac{Y_s}{Y_{s(\text{efficient})}} = \prod_{s=1}^S \left[ \sum_{i=1}^{M_s} \left( \frac{A_{si} \overline{TFPR}_s}{\overline{A}_s \overline{TFPR}_{si}} \right)^{\sigma-1} \right]^{\frac{\theta_s}{\sigma-1}}. \quad (8)$$

Equation 8 can be calculated from observed data. Instead of measuring how sensitive our calculation of productivity gains are to different underlying assumptions, which has been the primary focus of much of the recent methodological literature on misallocation, we instead calculate Equation 8 using different cuts of the data, which we describe in the next section.

### III Data

We primarily use micro-data from the United States, from the 2002 and 2007 US Census of Manufactures (CM).<sup>9</sup>

The quinquennial survey covers roughly 300,000 manufacturing plants, although information for the smallest plants - roughly a third of the sample - are almost entirely imputed. The standard is to exclude these so-called administrative records plants, which we do for all of our analysis.

Like in most surveys, not all respondents answer all of the questions, and some responses seem inconsistent with with responses to other questions for the same plant. The Census Bureau has created imputation and edit rules for this data, which are described in White et al. (2015). However, until the 2002 census, it was difficult for researchers to identify which, if any, items for a given plant were imputed. We go beyond the imputation flags and use the actual information that firms report to the Census Bureau (the “reported” data). The reported data differs from the final (“cleaned”) data in two re-

<sup>9</sup> We are working on acquiring access to the 2012 vintage (including the original reported data) as well. Furthermore, we report values for cross-country measured misallocation, where for the most part we rely on published sources, but also use microdata from India (in 2009) and Slovenia (in 2004), which we describe in Appendix A.I.

spects.<sup>10</sup> Missing values due to non-response in the reported data are imputed in the cleaned data, using a variety of industry-specific regression-based and other imputation strategies. Actual responses which fail edit rules in the reported data are also imputed in the final data. This editing takes two forms. For most variables, the imputation is done in the same way (and, in fact, at the same time) as the missing variables, using regressions to predict what plant-level behavior would be. For employment (both payroll and number of employees), the Census Bureau has administrative data that it can use, coming primarily from the IRS. The census forms in 2002 and 2007 specifically asked for annual payroll from “full- and part-time employees working at this establishment whose payroll was reported on Internal Revenue Service Form 941, Employer’s Quarterly Federal Tax Return,” and so, for single-plant firms and multi-plant firms that report payroll using a different EIN for each plant, the Census Bureau can use the actual reported information from those forms in order to potentially correct misreporting on the census form. For multi-plant firms that report payroll for multiple plants under the same EIN, when necessary the Census Bureau uses reported data for the same plants in current or prior years to allocate EIN-level payroll to plants.

We use several different versions of the CM data, all of which would plausibly be happily used by researchers in the absence of other alternatives. In addition to using the reported and cleaned data, we also replace the Census Bureau’s imputations using a different strategy known as sequential CART. We do this for two reasons. First, the Census Bureau’s regression methods put the imputed values on regression lines, thus potentially changing the true variability in the data.<sup>11</sup> Second, in some settings the Census Bureau’s imputation methods only use one covariate. This can introduce bias in estimates

---

<sup>10</sup>Another convention for naming the data is “captured” data for the reported information, and “completed” for the cleaned data.

<sup>11</sup>Nevertheless, there is no “spike” of TFPR at the industry mean. This is likely because the imputation is done for sub-components, so when aggregated up to the level of overall capital and value added the relative distortions are decreased but not fully eliminated.

of the relationship between the imputed variable and the variable used to predict it. The CART method is explicitly designed to approximate the conditional distributions of the variables being imputed, and (potentially) uses all available variables. We follow the procedure in White et al. (2015), and describe it here only briefly.

The foundation of regression tree-based-methods such as CART is to partition the covariate space into groups with similar outcomes. Along each branch, the data is split into two subgroups along the value contained in one covariate. The data continue to be split until the partitions contain a minimum number of observations, if those observations have sufficiently similar outcomes. Cross-validation strategies are used to prevent over-fitting.<sup>12</sup> Each of the reported variables used to calculate value added – total value of shipments, total cost of materials, and beginning and end of year finished goods inventories and work-in-progress inventories – are imputed primarily using univariate regressions. For the other variables used to compute TFP – payroll or employment and the book value of assets – it is not clear that the sequential CART imputation method does a better job of approximating the observed distributions than the Census Bureau’s imputations. For this reason, following White et al. (2015), we only replace the imputations for the variables used to calculate value added.

#### **IV Cross-Country Differences in Misallocation**

In this section, we consider how measured misallocation varies across countries, and how data cleaning and management strategies vary across countries along with the potentially fundamental differences in the organization of the manufacturing sector. In order to demonstrate the importance of this issue, we use the raw US data instead of the cleaned version. Much like for Table 3, we do so one variable at a time in order to show how measurement matters along each dimension. We suspect that the Census Bureau’s cleaned

---

<sup>12</sup>Currently we replace each Census Bureau imputation with one CART imputation. In future work we plan to replace each Census Bureau imputation with 100 CART imputations and use the average measure across the 100 CART-cleaned datasets as our measure of misallocation.

data is more accurate than the raw data, but we do not have the luxury to use similar data in other countries.

#### **IV.A Measured Misallocation in the Raw US data**

First, we consider the effects on measured misallocation of replacing cleaned with raw data in the US manufacturing sector. The results are shown in Table 1, with Panel A showing the results for the cleaned labor values, and Panel B for the raw labor values. Starting with the cleaned labor values, it remains clear that post-processing of capital again matters less than for value added. Moving to reported value added increases the gains from removing distortions to around 100% both for 2002 and 2007. If instead we used reported labor, but cleaned value added, the gains also increase to around 100%. If we use the self-reported values for both payroll and capital (which is what implicitly is done in most countries) the gains from removing misallocation rise to around 400%.<sup>13</sup>

#### **IV.B Measured Cross-Country Differences in Misallocation**

We now turn to discussing cross-country differences in measured misallocation. While we have measured misallocation in the United States for a large set of data choices, in order to avoid tedium we only describe cross-country differences in misallocation for two extremes: the 2002 and 2007 average for Census-Cleaned data, and the corresponding average in the Census-Reported data. Our measures of misallocation internationally come from a variety of published sources discussed in Appendix A.I. The results are shown in Table 2. While estimated misallocation in almost every developing country is higher than that for the cleaned US data, the measurement from the raw data the United States yields a larger amount of misallocation than is observed in any other country. Taking the results literally - which we do not - would imply that for, e.g., Argentina going to the US level of

---

<sup>13</sup>Another difference between the all-final and all-reported data samples is the number of observations: we drop plants with item non-response in the latter case. However, constraining the sample to only plants with complete reporting does not affect measured misallocation much for the cleaned data, with the change being consistently under 5%.

misallocation would decrease manufacturing TFP by around  $\frac{2}{3}$ .

In a more speculative approach, we build on Kalemli-Ozcan and Sorensen (2012) and calculate the measured gains within country-years in the World Bank Enterprise Analysis Unit's Enterprise Surveys.<sup>14</sup> The surveys are relatively small, and for the countries for which we have both enterprise surveys and data from national statistics agencies, the former have higher measured gains from reallocation. Nevertheless, Figure 1 shows that the gains from the US are larger than the corresponding gains for around 60% of the enterprise survey.<sup>15</sup>

## V Misallocation in the United States

In this section, we show the sensitivity of measured misallocation in the United States to the imputation strategy used for missing and edited data. In Table 3, we show how measured misallocation changes in the US when using the Census imputation strategy versus CART.<sup>16</sup> Following HK, we drop extreme values of the observed wedges and TFPR (relative to each firms' industry mean in the corresponding year) at the 1% or 2% extremes. We show how sensitive the results are to either form of data cleaning. In 2007, using the standard census-cleaned data we calculate that manufacturing TFP in the United States would be around 44% higher if the distortions were removed. If instead we used CART imputation for each of the components of value added, the calculated gains would instead be around 53%, an increase of around 10%.<sup>17</sup> While the 2% trim decreases measured mis-

---

<sup>14</sup>The raw data is available at <http://www.enterprisesurveys.org/data>. Our version is from August 1 2016, and we use the most recent sample from the 18 countries who a) we don't have access to an actual census for, b) who have at least 250 firms who report sales, labor, materials, and the replacement value of capital, and c) we drop Turkey in 2013 and Nigeria in 2014, since the measured gain from removing distortions are over 58000%, in those countries, which is implausibly high. Further details of data construction are in Appendix A.I

<sup>15</sup>When the World Bank Enterprise Group uses the data to calculate firm TFP, they undertake a careful process to clean the data and remove outliers. We, however, use the raw data.

<sup>16</sup>Note that using CART imputation does not mechanically increase measured misallocation.

<sup>17</sup>Here we are calculating the increase in actual TFP: if the current value of US TFP were denoted by  $x$ , then the counterfactual no-distortions TFP using the Census imputations would be  $1.444x$ , and with CART imputations for value added would be  $1.5309x$ .

allocation by around 5%, the relative changes of CART versus Census imputation strategies are reasonably unchanged. Overall, it is clear that the choice of imputation strategy matters quantitatively.

## VI Discussion

In this paper, we use previously unexplored versions of the United States Census of Manufacturers for 2002 and 2007 in order to investigate the role that measurement plays for estimating misallocation. We have two complimentary goals. The first is to measure misallocation for the United States using a different imputation strategy for missing data. Estimated misallocation using the CART-imputed data is around 10% higher than in the standard mean-imputed data. Our second goal is show measured cross-country differences in misallocation when using uncleaned data for the United States. For this, we instead use the data that is reported directly to the Census Bureau by the establishments, thereby avoiding all of the careful work done by Census Bureau statisticians in order to reduce the extent of misreporting. Here, the result is striking: measured misallocation in the United States is substantially higher than for any other country on earth for whom we have data from an official statistics agency. We do not take this result literally: there are many reasons to believe that comparing the raw US data to its counterparts in other countries is not like-for-like.<sup>18</sup> Our point is somewhat more cautionary: we demonstrate that there is a large scope for different measurement choices to affect the estimation of misallocation in manufacturing. Without much stronger assumptions, or collaborative cross-country efforts in order to ensure comparable data, it is difficult to know what cross-country differences in misallocation are, or if they exist at all.

---

<sup>18</sup>In no country do we know if measured misallocation in the raw data is larger or smaller than it is in reality, nor do we have a way of comparing the relative precision of self-reported information across countries. We do know that the vast majority of Indian firms are unable to fill out their survey forms on a computer.

## References

- Asker, J., A. Collard-Wexler, and J. De Loecker (2014). Dynamic Inputs and Resource (Mis)Allocation. *Journal of Political Economy* 122(5), 1013–1063.
- Banerjee, A. V. and E. Duflo (2005). Growth Theory through the Lens of Development Economics. *Handbook of Development Economics* 1(05), 473–552.
- Bartelsman, E. J. and W. Gray (1996). The NBER Manufacturing Productivity Database.
- Burgette, L. F. and J. P. Reiter (2010). Multiple Imputation Via Sequential Regression Trees. *American Journal of Epidemiology* 172(9), 1070 – 1076.
- Busso, M., L. Madrigal, and C. Pages (2013). Productivity and Resource Misallocation in Latin America. *B.E. Journal of Macroeconomics* 13(1), 903–932.
- Foster, L., C. Grim, J. Haltiwanger, and Z. Wolf (2016). Firm-Level Dispersion in Productivity: Is the Devil in the Details ? *American Economic Review: Papers & Proceedings* 106(5).
- Hopenhayn, H. A. (2014). On the Measure of Distortions. *Working Paper*.
- Hsieh, C.-T. and P. J. Klenow (2009). Misallocation and Manufacturing Tfp in China and India. *Quarterly Journal of Economics* 124(4), 1–55.
- Jaef, R. N. F. (2016). Entry Exit and Misallocation Frictions. *Working Paper*.
- Jensen, A. (2016). Employment Structure and the Rise of the Modern Tax System. *Working Paper*.
- Kalemli-Ozcan, S. and B. Sorensen (2012). Misallocation, Property Rights, and Access to Finance: Evidence from Within and Across Africa. *Working Paper*.
- Nishida, M., A. Petrin, M. Rotemberg, and T. K. White (2015). Are We Undercounting Reallocation's Contribution to Growth? *Working Paper*.
- Peters, M. (2013). Heterogeneous Mark-Ups , Growth and Endogenous Misallocation. *Working Paper*.

Randy Becker , Wayne Gray, J. M. (2016). NBER-CES Manufacturing Industry Database: Technical Notes. *National Bureau of Economic Research Technical Working Paper Series*.

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

White, T. K., J. P. Reiter, and A. Petrin (2015). Plant-level Productivity and Imputation of Missing Data in U.S. Census Manufacturing. *mimeo, Center for Economic Studies*.

Yang, M.-J. (2012). Micro-level Misallocation and Selection: Estimation and Aggregate Implications. *Working Paper*.

## **A Appendix**

### **A.I Cross-Country Estimates of Misallocation**

We use a variety of published sources to report the measured gains from reallocation in Table 2. With the exceptions of Chile and Columbia, our estimated potential gains from reallocation for South America come directly from Busso et al. (2013), and in Indonesia from Yang (2012). Both of those sources use information from national firm censuses. In addition to those values, we report information from Chile, Columbia, India and Slovenia, using the same micro-data described in Nishida et al. (2015).

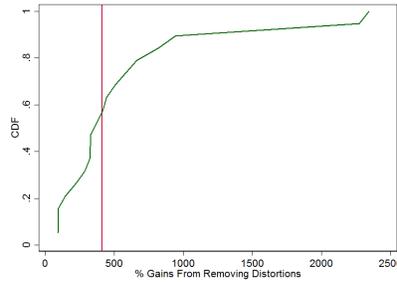
The Chilean and Colombian data are annual and we use 1995 and 1991, respectively. The Chilean data, provided by Chile’s Instituto Nacional de Estadística (INE), cover all manufacturing plants with at least 10 employees. The Colombian data from the Annual Manufacturing Survey, provided by Colombia’s Departamento Administrativo Nacional de Estadística (DANE), cover all plants with at least 10 employees.

In India, we use the Annual Survey of Industries (the ASI). Factories with over 100 workers are surveyed every year, while smaller establishments are surveyed every few years (the ASI is designed to be representative at the State by Industry level, so firms without local competitors are more likely to be surveyed). Hsieh and Klenow (2009) use

the same dataset, and we follow standard practice in generating measures of value added, capital, and payroll. Industries are grouped using India's NIC (National Industrial Classification) codes, and we report the value of reallocation for 2009. For Slovenia we rely on annual accounting data provided by the Slovenian Statistical Office which covers all manufacturing firms. The Slovenian data, unlike its counterparts in most other countries, is at the firm (not establishment) level, and we report the estimates for 2004. For the US, Slovenia, and India, we use cost-shares from the NBER-CES Manufacturing Industry Database as our measures of industry production elasticities, and multiply the book value of capital by 10% in order to impute the cost of capital.

The countries we use in the World Bank Enterprise Surveys (used for Figure 1) are Bangladesh, Colombia, Egypt, Iraq, Jordan, Kenya, Malaysia, Mozambique, Peru, Philippines, Russia, South Africa, Sri Lanka, Sweden, Thailand, Tunisia, Vietnam, and Zimbabwe. We use "Total annual cost of labor" as the measure of labor, the sum of the cost of materials, electricity, communications services, fuel, transport, and water for materials, the sum of the (self-reported) replacement costs for machinery and land for capital, and total sales for gross output. For each manufacturing sector we assume that the capital elasticity is  $\frac{1}{3}$ . Across all surveys, we drop firms who report non-positive values for any of those four variables.

Figure 1: Potential Gains from Reallocation in the World Enterprise Surveys



*Notes:* This figure plots the gains from removing distortions for 19 countries in the World Bank Enterprise Surveys, described in Appendix A.I. The vertical line corresponds to the gains in the United States in the raw data for 2007.

Table 1: Potential Gains from Reallocation in the United States, Raw Data

	Capital Imputation Strategy, 2002				Capital Imputation Strategy, 2007			
	1% trimming		2% trimming		1% trimming		2% trimming	
	Census	Raw	Census	Raw	Census	Raw	Census	Raw
<b>Panel A: Cleaned Labor</b>								
Census Value Added	44.4%	46.1%	34.17%	36.49%	52.08%	60.4%	38.92%	43.88%
Raw Value Added	69.46%	65.59%	50.63%	47.78%	95.48%	95.76%	63.15%	66.55%
<b>Panel B: Raw Labor</b>								
Census Value Added	76.67%	72.16%	64.12%	62.45%	99.02%	110.63%	87.73%	96.51%
Raw Value Added	454.41%	333.24%	337.69%	201.61%	413.1%	372.5%	318.87%	264.06%

Each cell calculates Equation 8 using the the corresponding values for capital value added, and labor.

Table 2: Measured Cross-Country Differences in Misallocation

Country	Gains in Most Recent Year	Gains Relative to:	
		Final US	Raw US
India	100%	35%	-56%
Mexico	95%	32%	-57%
China	87%	26%	-59%
Chile	77%	19%	-61%
Indonesia	68%	13%	-63%
Venezuela	65%	11%	-64%
Bolivia	61%	8%	-65%
Uruguay	60%	8%	-65%
Argentina	60%	8%	-65%
Ecuador	58%	6%	-65%
Slovenia	57%	6%	-65%
El Salvador	57%	6%	-65%
Colombia	49%	1%	-67%
Brazil	41%	-5%	-69%

Each cell shows measured misallocation and the “gains” from moving to measured-US levels. Data sources are discussed in Appendix A.I

Table 3: Potential Gains from Reallocation in the United States, CART Imputation

	2002		2007	
	1% trimming	2% trimming	1% trimming	2% trimming
Census Value Added	44.4%	34.17%	52.08%	38.92%
CART Value Added	53.09%	40.69%	59.97%	45.72%

Each cell calculates Equation 8 using the Census cleaned entries for payroll and capital, and the corresponding values for value added.