



# Why do industries coagglomerate? How Marshallian externalities differ by industry and have evolved over time

Dario Diodato<sup>a,\*</sup>, Frank Neffke<sup>a</sup>, Neave O'Clery<sup>b</sup>

<sup>a</sup> Center for International Development, Harvard University, 79 JFK street, Cambridge, MA 02138, USA

<sup>b</sup> Mathematical Institute, University of Oxford, Woodstock Road, Oxford, OX2 6GG, England, United Kingdom

## ARTICLE INFO

### Article history:

Received 10 August 2016

Revised 7 February 2018

Available online 29 May 2018

### JEL classification:

J24

O14

R11

### Keywords:

Coagglomeration

Marshallian externalities

Labor pooling

Value chains

Manufacturing

Services

Regional diversification

## ABSTRACT

The fact that firms benefit from close proximity to other firms with which they can exchange inputs, skilled labor or know-how helps explain why many industrial clusters are so successful. Studying the evolution of coagglomeration patterns, we show that the type of agglomeration that benefits firms has drastically changed over the course of a century and differs markedly across industries. Whereas, at the beginning of the twentieth century, industries tended to collocate with their value chain partners, in more recent decades the importance of this channel has declined and collocation seems to be driven more by similarities in industries' skill requirements. By calculating industry-specific Marshallian agglomeration forces, we are able to show that, today, skill-sharing is the most salient motive behind the location choices of services, whereas value chain linkages still explain much of the collocation patterns in manufacturing. Moreover, the estimated degrees to which labor and input-output linkages are reflected in an industry's coagglomeration patterns help improve predictions of city-industry employment growth.

© 2018 Elsevier Inc. All rights reserved.

## 1. Introduction

In spite of congestion, elevated factor costs and the risk that trade secrets leak to competitors, firms in the same industry frequently locate close to one another (Ellison and Glaeser, 1999; Rosenthal and Strange, 2001). As a consequence, in many industries we observe pronounced geographical clusters. The existence of these clusters are often attributed to the presence of three different types of externalities, namely the sharing of inputs, labor and knowledge. However, not all industries might benefit equally from agglomeration. Moreover, given the dramatic decline in transportation and communication costs and the progressive spatial fragmentation of value chains in the previous century, it is likely that the balance of costs and benefits for industries to cluster has changed drastically. However, our understanding of how agglomeration externalities differ across industries and have changed over time is still limited. The aim of this paper is to assess the relative importance of different drivers of agglomeration across industries, and how the main forces behind agglomeration externalities have changed over the course of a century.

Marshall (1920) ascribed “the advantages which people following the same skilled trade get from near neighborhood to one another” (p.225) to three different types of agglomeration externalities: the benefits of a large pool of skilled labor, easy access to local customers or suppliers and local knowledge spillovers.<sup>1</sup> However, in spite of the early recognition and ample subsequent research of this topic, the relative importance of each of these Marshallian externalities has fueled debate for over a century. In part, this is due to the so-called “Marshallian equivalence” (Duranton and Puga, 2004): all three Marshallian agglomeration theories yield the same prediction for the spatial distribution of an industry, namely, that, because they generate benefits for one another, economic establishments engaged in similar activities will tend to agglomerate. This confluence of agglomeration benefits makes it difficult to determine which theory carries most weight as an explanation for the observed tendency of industries to concentrate in space.

<sup>1</sup> Whereas theoretical models tend to divide agglomeration externalities into benefits from sharing, matching and learning in local economies (Duranton and Puga, 2004), most of the empirical literature on the topic categorizes Marshallian externalities as economies in transportation, coordination or communication when acquiring one of three factors: labor, (intermediate) capital goods or knowledge.

\* Corresponding author.

E-mail address: [Dario.Diodato@hks.harvard.edu](mailto:Dario.Diodato@hks.harvard.edu) (D. Diodato).

A major stride forward was achieved by Ellison, Glaeser and Kerr (2010), henceforth EGK, who study not the agglomeration of individual industries but the coagglomeration of pairs of industries. The rationale for this is that industries that are similar along some dimensions may differ along others. For instance, whereas some industries will benefit from being colocated because they employ similar labor, others may colocate because of input-output or technological linkages. By analyzing the relationship between locational similarity (i.e., coagglomeration) and similarities that reflect different Marshallian agglomeration benefits, EGK disentangle the strength of three different types of Marshallian externalities. They find that input-output linkages are the most important explanation for why industries coagglomerate, closely followed by opportunities to share labor. Least empirical support is found for sharing know-how as a rationale for coagglomeration.

The effects EGK report represent averages across all industries, and, as such, may conceal marked differences among industries. For instance, whereas making musical instruments requires specialized workers with years of on-the-job training, workers in food-processing firms are often employed on short-term contracts through temporary work agencies, without much regard for their skills. Similarly, while car manufacturers often closely collaborate with their local suppliers (Morgan and Cooke, 1998), the principal inputs for steel mills, coal and iron ore, are acquired on anonymous exchanges with little need for buyer-supplier interaction. Finally, although knowledge spillovers may be important drivers of the clustering of biotechnology firms (Zucker et al., 1994), they will be less important in industries in which technology progresses less rapidly. Meta-studies reviewing the empirical literature on agglomeration externalities since the foundational papers by Glaeser et al. (1992) and Henderson et al. (1995) confirm the existence of considerable variation in empirical findings (Beaudry and Schiffauerova, 2009; Groot et al., 2016). We expect that these differences are in part driven by variation across industries in how much they rely on specific agglomeration forces, but also by changes in these forces over time.

Building on EGK, we explore whether this hypothesized heterogeneity in agglomeration benefits is expressed in industrial coagglomeration patterns. We start by replicating key parts of the original work by EGK, which was based on US manufacturing industries in the late 1980s and 1990s, using similar data for the 2000s. Mimicking EGK, we do so using both Ordinary Least Squares (OLS) estimation and Instrumental Variables (IV) based strategies that instrument inter-industry linkages among US industries by analogous measures constructed from data external to the US, mainly drawing from data for the Mexican economy. We confirm that labor linkages and input-output linkages are still more or less equally important explanations for coagglomeration, whereas evidence for technological (knowledge) spillovers is relatively weak. However, there are good reasons to expect that coagglomeration patterns in services will be different from those in manufacturing. First, unlike manufactured goods, services are often difficult to trade over large distances. Consequently, service firms will need to colocate with their customers. Second, services tend to be labor intensive, and their quality often crucially depends on face-to-face interactions between a firm's employees and its customers.<sup>2</sup> For both reasons, access to adequate human capital may be particularly important in services. This conjecture is confirmed when we extend the analysis to industries beyond manufacturing: the estimated effects of input-output linkages on coagglomeration are at least as large in services as in manufacturing, while labor-linkages are a much

stronger factor in the coagglomeration of services than of manufacturing industries.<sup>3</sup>

Next, we allow Marshallian externalities to vary fully by industry. Doing so reveals an even wider variation in agglomeration effects. For some manufacturing activities, such as furniture and food production, coagglomeration patterns can be attributed neither to labor nor to input-output linkages. For other manufacturing industries, such as those in electronics or in the pharmaceuticals & medical sector, coagglomeration patterns are primarily driven by input-output linkages. However, the greatest variety in agglomeration effects exists in services, where some industries, such as those in arts and culture, cluster along both externality channels, whereas in other services, like media and knowledge-intensive business services, labor pooling opportunities dominate coagglomeration patterns.

Finally, we turn to changes in coagglomeration forces over time. The US economy and its geography have undergone major changes in the twentieth century. Urbanization rates rose from 40.5% to over 75% (Black and Henderson, 2003), work shifted from over 40% of male employment in agriculture<sup>4</sup> to the current dominance of service jobs, and the US population has become increasingly educated, with school enrolment rates rising from 18% in 1910 to 73% in 1940 (Goldin and Katz, 2009, p. 195). Meanwhile, value chains started fragmenting through outsourcing and offshoring spurred by the drastic reduction in transportation and communication costs due to the roll-out of highway systems, airport and internet infrastructure. Using data on the coagglomeration of industries in US states between 1910 and 2010, we explore how, combined, these processes have altered agglomeration forces. We find that the importance of value chain linkages in industrial coagglomeration patterns has strongly declined, while the relevance of labor pooling opportunities has – if not mildly increased – at least stayed constant. These findings suggest that the aforementioned historical changes indeed led to a shift in agglomeration forces. Arguably, whereas the contraction of space has eroded the importance of geographical proximity among customer and suppliers, the skill accumulation and specialization that came with the expansion of the education system and a deepening of the economy's knowledge base has amplified the relative benefits of sharing a common pool of labor.

By studying the heterogeneity of Marshallian externalities, we contribute first and foremost to the literature on agglomeration externalities in general, and to the fast-expanding strand of research that focuses on coagglomeration to untangle different agglomeration forces in specific (Ellison et al., 2010; Howard et al., 2015; Behrens, 2016; O'Sullivan and Strange, 2017). Our work is closely related to the work of Faggio et al. (2017), who study heterogeneity in Marshallian externalities in the UK. These authors find that agglomeration effects are particularly strong among small firms. Moreover, they report differences in agglomeration forces between high- and low-tech sectors, differences driven by the level of education of an industry's workforce, and – using quantile regressions – differences by industries' general tendency to coagglomerate. Although we also conclude that different industries experience different agglomeration forces, we arrive at this conclusion from a markedly different empirical approach. Whereas Faggio et al. (2017) explicitly look for differences related to the

<sup>3</sup> For reasons explained more fully later on, unlike EGK, we disregard technological linkages and shared natural advantages as explanations for coagglomeration. Whereas our reservations about technological linkages are related to measurement issues – the most readily available source of information on such linkages, patents, are barely used in the majority of services – our reservations about natural advantages (as defined in EGK) are more fundamental: first, many natural advantages in EGK are actually man-made, second, EGK use local prices, which are endogenous to the agglomeration process, to quantify natural advantages.

<sup>4</sup> Authors' own calculations using IPUMS data.

<sup>2</sup> See, for instance, Kolko (1999) on the importance of agglomeration in services.

organization and technology of industries, we take a more agnostic approach and study heterogeneity in agglomeration forces in a nonparametric fashion. One finding that emerges from this is that labor sharing tends to be more important in services than in manufacturing, whereas the opposite holds for value-chain linkages. Given the growing importance of services in the economy and the fact that – using a century worth of data – our historical analysis suggests that labor linkages are gaining in relative importance, labor sharing may soon become the dominant rationale for the cities of the future. However, at the same time, for some services coagglomeration are also strongly driven by value-chain linkages,<sup>5</sup> possibly reflecting the fact that the output of services is to a large extent nontraded and, therefore, associated with very high transportation costs. Moreover, we show that the heterogeneity in coagglomeration translates into differential growth of local industries. With this exploration of how the estimated differences in coagglomeration forces can help gauge how sensitive an industry's local growth rate is to a given type of agglomeration benefit, our paper also contributes to a growing literature on agglomeration externalities and local industry growth in urban economics (Glaeser and Kerr, 2009; Dauth, 2010; Jofre-Monseny et al., 2011, Hanlon and Miscio, 2017), cluster research (Porter, 2003; Delgado et al., 2010), and an emerging literature on related diversification in economic geography (Neffke et al., 2011; Hausmann et al., 2014).

## 2. Methodology

### 2.1. Data

To measure coagglomeration patterns in the current economy, we use datasets that describe employment by region-industry pair in the US and Mexico. For the US, this data is derived from the County Business Patterns (CBP) for the years 2003 and 2008.<sup>6</sup> Employment data for Mexico are taken from the economic censuses in 2003 and 2008.

We analyze this data at three different geographical levels: US counties (of which there are 3190), metropolitan areas (922 including both Metro- and Micropolitan Statistical Areas) and states (51). In Mexico, our geographical units are 2455 municipalities, 58 metropolitan areas and 32 states. Because we regard metropolitan areas (labelled 'cities' hereafter) as the most appropriate spatial unit for defining labor markets and economically integrated regions, we will focus the discussion on the results derived at this level, and report county and state level results as supporting evidence where practical. Furthermore, to define industries, we adjust the North American Industry Classification System (NAICS) to correct for small inconsistencies between US and Mexican data (see Appendix B). From the resulting list of 184 industries, we exclude nontraded industries with spatial distributions that are strongly driven by the distribution of population, such as retail, auto repair, construction work and elementary schools. Note, however, that we do not exclude extractive activities such as mining. Although these activities are obviously restricted in their location choice, it is still informative to see how other industries choose to colocate with them. This leaves us with a set of 120 traded industries.<sup>7</sup> As a robustness check, we repeat the main analysis using the full sample

of 184 industries. The outcomes, which are reported in Appendix F, are in line with those based on the restricted sample.

To analyze historical coagglomeration patterns, we use the US census samples provided by IPUMS USA (Ruggles et al., 2017) for the period 1910–2010. Given the between 1% and 5% sampling rates, estimated industry-region employment data will be too noisy for small spatial units. Therefore, we restrict our analysis in this section to the coagglomeration of industries within US states. The industry classification in IPUMS (IND1990) distinguishes more than 200 activities, from which we select a set of 104 industries following the same logic as for CBP data.<sup>8</sup>

### 2.2. Dependent variable: coagglomeration

Our main object of interest is the degree to which industries coagglomerate. That is, to what extent do pairs of industries employ workers in the same regions? To quantify the tendency of industry  $i$  to coagglomerate with industry  $j$ , EGK propose the following measure:<sup>9</sup>

$$EG_{ij} = \frac{\sum_{r=1}^R (s_{ir} - x_r)(s_{jr} - x_r)}{1 - \sum_{r=1}^R x_r^2}. \quad (1)$$

where  $s_{ir} = \frac{E_{ir}}{\sum_{r'=1}^R E_{ir'}}$  – with  $E_{ir}$  industry  $i$ 's employment in region  $r$  – is the employment share of industry  $i$  in region  $r$ , while  $x_r$  is the mean of these shares in region  $r$  across all industries. Using a model of location choice, Ellison and Glaeser (1997,1999) motivate this index as a measure of the likelihood that establishments in industries  $i$  and  $j$  generate spillovers for one another. The index has the advantage that it should not be affected by the size distribution of establishments in an industry or by the granularity of spatial units. Following this logic, we calculate EG indices for all pairs of industries in the US.

The EG index is similar in spirit to a measure used by Porter (2003), who quantifies the coagglomeration of two industries as the correlation between industries' locational employment vectors:

$$LC_{ij} = \text{corr}(s_{ir}, s_{jr}), \quad (2)$$

Hausmann et al. (2014) show that the locational correlation of industries  $i$  and  $j$ ,  $LC_{ij}$ , can be interpreted as an estimate of the similarity of the industries' technology requirements in a Ricardian trade model. For brevity, we focus our discussion on the analysis that uses the EG index. The results using the LC index are reported in Appendix E, unless they lead to qualitatively different conclusions, in which case we discuss them in the main text.

### 2.3. Independent variables

EGK use four different inter-industry linkages to explain why industries may coagglomerate. The first three represent Marshallian externality channels: value-chain linkages, similarities in labor requirements and technological similarities among industries. The fourth linkage type captures the fact that the location of some industries reflects a need to access natural resources. Although these "natural advantages" do not represent Marshallian externalities, they may result in nonrandom coagglomeration of industries that require similar natural resources.

<sup>5</sup> This evidence is particularly salient when using a locational correlation index instead of the Ellison–Glaeser index of coagglomeration.

<sup>6</sup> The CBP data censor small industry-region cells. For such cells, we assign employment numbers following Holmes and Stevens (2004). For further details, we refer to Appendix G.

<sup>7</sup> Traded industries represent about 38% of employment nationwide. That is, 100% of manufacturing, 100% of extractive activities, 0% of agriculture, utility or construction. Among services, we consider as traded a set of activities that constitute about 30% of employment in the tertiary sector. We list all excluded and included industries in Appendix D.

<sup>8</sup> More details on the construction of variables from IPUMS data can be found in Appendix B, while the list of included IPUMS industries is provided in Appendix D.

<sup>9</sup> This measure is equivalent to the coagglomeration index in Ellison and Glaeser (1999) for the coagglomeration of pairs of industries.

### Input-output links

Value chains allow individual firms to specialize. However, such specialization also creates costs: intermediates need to be shipped between firms, and innovation efforts must be coordinated with suppliers (Richardson, 1972; Abdel-Rahman, 1996; Porter, 1998). Because the costs of transportation and coordination typically rise with distance, colocating different parts of a value chain can be an effective cost-reduction strategy.

We measure the strength of input-output relations between a pair of industries using the same indicator as EGK. That is, the input-output proximity of industries  $i$  and  $j$  is defined as the maximum relative importance of  $i$  as a customer or as a supplier of  $j$  and vice versa. Let  $IO_{ij}$  be an input-output matrix, i.e.,  $IO_{ij}$  represents the value of goods and services that industry  $j$  sources from industry  $i$ . We now measure the proximity between  $i$  and  $j$  in terms of input-output linkages as

$$P_{ij}^{IO} = \max \left( \frac{IO_{ij}}{\sum_k IO_{kj}}, \frac{IO_{ji}}{\sum_k IO_{kj}}, \frac{IO_{ij}}{\sum_k IO_{ik}}, \frac{IO_{ji}}{\sum_k IO_{ik}} \right) \quad (3)$$

To estimate current value-chain linkages, we use tables provided by the Bureau of Economic Analysis (BEA) for the year 2002 in the US and by the Mexican statistical office, Instituto Nacional de Estadística y Geografía (INEGI) for the year 2008 in Mexico. For the historical analyses, we use input-output tables provided by BEA for the period 1947–2012. A detailed description of the data processing and subsequent calculations is provided in Appendix B.

### Labor market pooling

A large local pool of specialized labor benefits both firms and workers. First, larger pools of skilled workers (and firms that want to hire them) may result in a better matching of workers to firms (Helsley and Strange, 1990). Second, workers may demand a wage premium as a compensation for moving to regions that offer few alternative employment opportunities in case they lose their jobs. In contrast, having many firms and industries that can absorb one another's redundant workers acts as an implicit insurance scheme and may, therefore, lower wage demands (Marshall, 1920; Duranton and Puga, 2004).

Following EGK, we measure the extent to which two industries can draw from the same pool of workers using industry-occupation employment matrices. In particular, we compute correlation coefficients between the occupational employment of industries  $i$  and  $j$ ,  $E_{io}$  and  $E_{jo}$ , where  $E_{io}$  represents the number of workers in occupation  $o$  employed (nationwide) in industry  $i$ :

$$P_{ij}^L = \text{corr}(E_{io}, E_{jo}). \quad (4)$$

For the US, we use the industry-occupation data for the year 2002 reported in the Occupational Employment Statistics (OES), while for Mexico, we use the Encuesta Nacional de Ocupación y Empleo (ENOE). Historical estimates are derived from the samples of the US census provided by IPUMS USA. For more details, we refer the reader to Appendix B.

### Technological similarity

EGK estimate inter-industry technological similarity using cross-industry patent citations and the technology-flow matrix constructed by Scherer (1984). The latter yields the, conceptually, most attractive metric. To build the technology-flow matrix, science-degree students assessed for a large sample of patents which industries were most likely to be beneficiaries of the inventions described therein. Given that this exercise has not been repeated, we fall back to patent citations. Using the NBER patent citations dataset (Hall et al., 2001), we construct the overrepresentation of industries  $i$  and  $j$  in bilateral citation flows, and measure

the extent to which industries exchange knowledge,  $P_{ij}^T$ , in a similar fashion as  $P_{ij}^{IO}$ , as

$$P_{ij}^T = \max \left( \frac{C_{ij}}{\sum_k C_{kj}}, \frac{C_{ji}}{\sum_k C_{kj}}, \frac{C_{ij}}{\sum_k C_{ik}}, \frac{C_{ji}}{\sum_k C_{ik}} \right) \quad (5)$$

where  $C_{ij}$  is the number of citations from patents associated with industry  $i$  to patents associated with industry  $j$ . We estimate  $P_{ij}^T$  twice: once as a main covariate using patents by US inventors, and once as an instrument using patents by foreign inventors. Details of this procedure are provided in Appendix B.

### Natural advantages

In some industries, firms' location choices are limited by factors that are driven by natural advantages. For instance, shipbuilding requires access to waterways and some tourism activities rely on the beauty of a place's natural environment. To correct for such factors, EGK calculate an index of how similar industries are in terms of their reliance on natural resources. This index is based on work by Ellison and Glaeser (1999) that predicts for each industry what share of its employment will locate in a given region from information on, on the one hand, the industry's intensity of use of a given resource, and, on the other hand, the price of this resource in the region. Next, the predicted shares are used to estimate predicted coagglomeration patterns as a metric of similarity in natural advantages:

$$EG_{ij} = \frac{\sum_{r=1}^R (\hat{s}_{ir} - x_r)(\hat{s}_{jr} - x_r)}{1 - \sum_{r=1}^R x_r^2}. \quad (6)$$

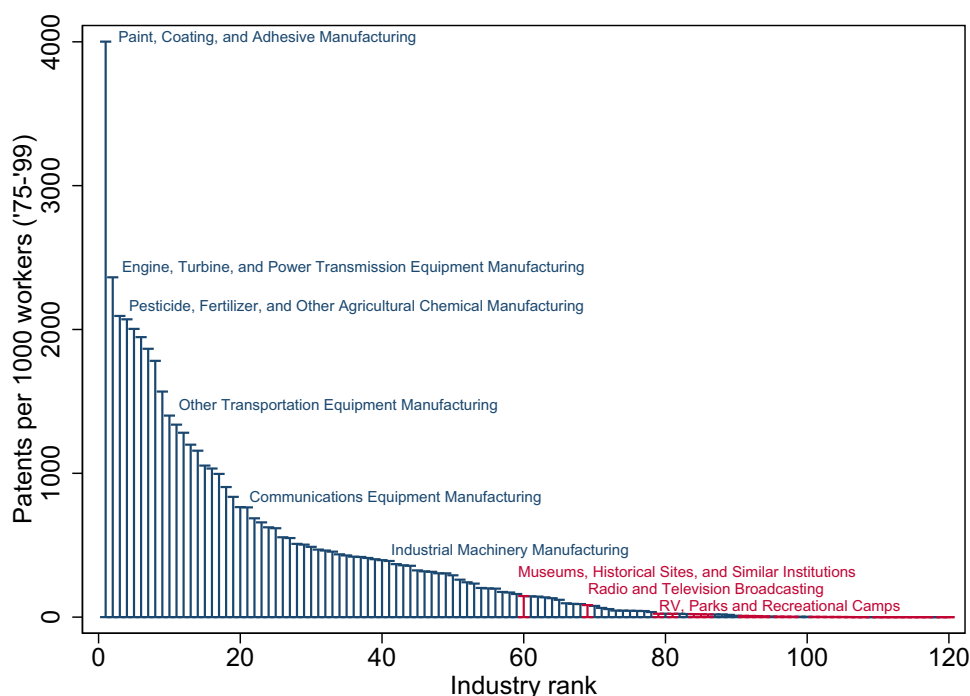
Note that Eq. (6) is the same as Eq. (1) but for the substitution of  $\hat{s}_{jr}$  for  $s_{jr}$ . For reasons explained hereafter, we will not repeat this exercise in full. Instead, we recalculate Eq. (6) using two different versions of  $\hat{s}_{ir}$ . First, to replicate the EGK results, we update the estimated employment shares by industry taken from the supplementary online material of EGK in such a way that they match the US economy's employment composition of 2003 and the NAICS classification system. Next, we create our own estimates to extend the analysis to services, which are disregarded by EGK. These estimates combine information on an industry's use of primary resources (using IO tables) and unskilled labor (using the industry's occupational employment vector) on the one hand, and the regional availability of these primary resources (proxied by the presence of the associated extractive industries) and unskilled labor (taken from BLS tables on urban employment by occupation), on the other hand. Both procedures are explained in more detail in Appendix B.

### Discussion

As pointed out by EGK, a possible cause for concern is that the inter-industry similarity matrices are endogenous. That is, if two industries located close to each other due to an historical accident, this coagglomeration may have prompted them to adjust their production technologies, for instance, in such a way that they could use one another's outputs or skilled workers. Another concern is that similarity matrices are only imperfect proxies for the degree to which industries can exchange products, or use the same workers, technologies or natural advantages. Such measurement error would lead to a downward bias in our estimates. In line with EGK, we therefore also estimate the models in Eqs. (7) and (8) using an IV approach. As instruments, we use analogously constructed variables based on information from outside the US economy.<sup>10</sup> Similar

<sup>10</sup> For value chain and labor similarity linkages, we use data on the Mexican economy, for technological similarity we use patents by inventors who reside outside the US. Natural advantages, should, by definition, be exogenous and therefore don't require instrumenting. Below, however, we will raise some doubts about the exogeneity of the natural advantage variable in EGK.





**Fig. 1.** Patenting rates by industry.

Patents in the period 1975–1999 per 1000 employees by industry. Industries sorted by patenting intensity. Blue bars refer to manufacturing industries, red bars to services. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

to the UK-based instruments in EGK, our instruments are valid, as long as idiosyncratic patterns in the input-output, labor or technological linkages *outside* the US are exogenous to coagglomeration patterns *in* the US.

A second concern is related to the main goal of this paper. To explore the heterogeneity of Marshallian externalities across industries and over time, it is insufficient to generate proxies for each type of inter-industry linkage; we must also ensure that these proxies can be created with some accuracy for all industries. Although value-chain linkages and occupational similarity matrices can be constructed with relative ease for any industry, the same cannot be said for technological linkages because industries differ widely in how much they rely on patents to protect their intellectual property (Cockburn and Griliches, 1988). To illustrate this, Fig. 1 shows the number of patents in the period 1975–1999 per 1000 employees by industry. Clearly, although most manufacturing industries (in blue) produce at least some patents, most service industries do not patent any inventions at all. What is more, the technologies that get patented are arguably unrepresentative of the vast set of technologies firms use in their day-to-day processes.<sup>11</sup> As a result, the accuracy with which inter-industry technological similarity can be estimated from patent data will differ tremendously by industry. Because this heterogeneity in measurement error will translate mechanically into heterogeneity in the estimated effect of inter-industry knowledge spillovers on coagglomeration patterns, it compromises what our paper sets out to do.<sup>12</sup>

<sup>11</sup> EGK are well aware of this: “Our patent citation measure is a proxy for the importance of exchanging technology rather than a proxy for all forms of intellectual spillovers. Since our measures of idea sharing are weaker than our measures of input-output linkages, we anticipate their connection with coagglomeration to be weaker.” (p. 1202)

<sup>12</sup> That is not to say that we believe that the other inter-industry proximity measures are unaffected by heteroscedastic measurement-errors. However, for reasons explained above, we believe that the variation in the precision with which we can characterize labor and input-output linkages is dwarfed by the variation in measurement error for technological linkages. Moreover, although the IV estimations

A third and final concern arises due to the way natural advantages are measured in EGK. Although we agree that natural advantages may have driven some of the historical coagglomeration patterns in the US, the metric EGK use to correct for this has problems of its own. First, it is unclear whether prices should be considered when constructing a measure of natural advantages as in EGK. After all, local prices are endogenously determined by supply and demand in a region. For instance, low prices of a natural resource may make a location attractive to industries that heavily rely on this resource, or reflect that demand for them is low because such industries have shunned the region for other reasons. Second, the majority of production factors EGK list as natural resources are, in fact, man-made. For instance, EGK do not just label cheap electricity and abundant farmland and cattle as natural advantages. They also include low local manufacturing wages, educational attainment, unionization rates and population density. Most of these factors would also be reflected in similarities in industries’ labor requirements. Moreover, a number of other resources, such as coal and timber – although arguably related to natural endowments – are more accurately regarded as outputs of extractive industries. As such, they should be considered as parts of value chains. Upon closer inspection, the only truly exogenous factor in EGK appears to be a coastal dummy.<sup>13</sup> Therefore, the measure of similarities in natural-resource use is – if not endogenous – redundant given that it uses much of the same information that is provided in input-output and labor-use matrices. For these reasons, after having compared our results to the ones in EGK, we will drop both technological and natural-resource-use similarities as explanations for coagglomeration patterns.

will correct for some of the errors-in-variables bias, the scope for this is limited in the case of technological linkages, where the instrument is based on patent data as well.

<sup>13</sup> However, this dummy has no significant impact on industries’ locational choices (see Ellison and Glaeser, 1999).

**Table 1**  
Summary statistics.

Variable	Obs	Mean	Std. dev.	Min	Max
United States					
EG index	7140	0.0001	0.0047	−0.0232	0.0846
LC index	7140	0.7166	0.1209	0.4854	0.9915
Input-output	7140	0.0124	0.0353	0.0000	0.6776
Labor	7140	0.6013	0.0997	0.4956	0.9919
Technology	7140	0.0082	0.0267	0.0000	1.2346
Technology (foreign)	7140	0.0082	0.0266	0.0000	1.1838
Nat. advantage LC (original)	3828	0.9703	0.0357	0.7487	1.0000
Nat. advantage EG (original)	3828	−0.0000	0.0011	−0.0060	0.0084
Nat. advantage LC (new)	7140	0.9784	0.0262	0.7819	1.0000
Nat. advantage EG (new)	7140	0.0000	0.0010	−0.0056	0.0058
Mexico					
EG index	7140	0.0005	0.0477	−0.2579	0.5149
LC index	7140	0.7784	0.1657	0.4399	0.9994
Input-output	7140	0.0143	0.0503	0.0000	1.0000
Labor	7140	0.5565	0.1094	0.4828	1.0000

Summary statistics for coagglomeration (EG and LC metrics) and input-output and labor linkages for  $120 \times (120 - 1)/2 = 7140$  unique industry pairs for the US (top) and Mexico (bottom).

2.4. Descriptive statistics

Excluding the diagonal, there are 7140 unique industry pairs in our sample of 120 industries. Table 1 contains descriptive statistics. It is interesting to note that, whereas means are similar for all variables in the US and Mexico, the dispersion is larger in Mexico, and sometimes substantially so. Although the greater dispersion in the Mexican variables may be structural, it could also mean that the variables constructed from US data are measured more accurately. Table A.6 reports correlation coefficients between the various dependent and independent variables. Interestingly, the correlations of the same variables in a different country are typically higher than of different variables in the same country. For instance, there is a 0.24 correlation between US labor and input-output similarities, whereas these variables exhibit correlations of over 0.5 with their Mexican counterparts. This is reassuring because it suggests that the similarity measures capture distinct, yet general, relations among industries.

2.5. Estimation framework

To infer the strength of agglomeration forces, we follow EGK and analyze the relationship between the coagglomeration patterns and various types of inter-industry linkages using the following econometric model:

$$EG_{ij} = \alpha + \sum_{rel} \beta^{rel} P_{ij}^{rel} + \epsilon_{ij}, \tag{7}$$

where  $EG_{ij}$  is the EG-index of coagglomeration for industries  $i$  and  $j$ , and  $rel \in \{IO, L, T, NA\}$ . Matrices  $P^{IO}$ ,  $P^L$ ,  $P^T$  and  $P^{NA}$  contain the inter-industry value-chain linkages, labor-requirement, technology and natural-advantage similarities defined in the previous sections. While EGK only study coagglomeration among manufacturing industries, we extend the analysis to other sectors of the economy. We do so by expanding both the columns and rows of the EG matrix. Next, we split the rows into a manufacturing and a services section. Within these subsets, we study how industries co-agglomerate with any other industry. That is, we study the coagglomeration patterns of a manufacturing industry (or of a service) with all other industries. Finally, note that each industry can co-agglomerate with 119 other industries. Therefore, there are sufficiently many degrees of freedom to estimate agglomeration effects that vary freely by industry. We do so using the following two-way

**Table 2**  
OLS univariate regressions (EGK replication)..

	EG index (our estimates)			EG index (EGK estimates)		
	(1) State	(2) City	(3) County	(4) State	(5) City	(6) County
Input-output	0.214 (0.035)	0.177 (0.032)	0.144 (0.029)	0.205 (0.037)	0.167 (0.028)	0.130 (0.022)
Observations	3655	3655	3655	7381	7381	7381
R <sup>2</sup>	0.060	0.060	0.047	0.042	0.028	0.017
Labor	0.181 (0.018)	0.164 (0.018)	0.195 (0.016)	0.180 (0.014)	0.106 (0.016)	0.082 (0.013)
Observations	3655	3655	3655	7381	7381	7381
R <sup>2</sup>	0.029	0.035	0.060	0.032	0.011	0.007
Technology	0.108 (0.016)	0.109 (0.013)	0.116 (0.012)	0.081 (0.012)	0.100 (0.016)	0.085 (0.013)
Observations	3655	3655	3655	7381	7381	7381
R <sup>2</sup>	0.015	0.022	0.030	0.007	0.010	0.007
Nat. advantages	0.109 (0.018)	0.085 (0.014)	0.082 (0.009)	0.210 (0.020)	0.188 (0.017)	0.222 (0.014)
Observations	3655	3655	3655	7381	7381	7381
R <sup>2</sup>	0.014	0.012	0.014	0.044	0.036	0.049

Robust standard errors in parentheses. Columns 1–3 replicate the univariate regression results of Ellison et al. (2010) (shown in columns 4–6) using equivalent data for the US from 2003 for manufacturing industries and three levels of spatial aggregation: states, cities and counties. In all cases, the dependent variable is the pairwise EG coagglomeration index (Eq. (7)).

fixed effects model:

$$EG_{ij} = \alpha_i + \delta_j + \sum_{rel} \beta_i^{rel} P_{ij}^{rel} + \epsilon_{ij} \tag{8}$$

This regression yields vectors of industry-specific agglomeration effects,  $\hat{\beta}_i^{rel}$ . The elements of these vectors represent the extent to which a given inter-industry linkage type is expressed in an industry's coagglomeration patterns.

3. Empirical findings

3.1. Replication of EGK results

We begin our analysis by assuming that agglomeration effects are homogeneous across industries. That is, we assume that Marshallian externalities affect all industries in the same way. Table 2 replicates the results of EGK. The table shows the results of uni-

**Table 3**  
OLS and IV multivariate regressions on extended sample.

	(1) OLS State	(2) OLS City	(3) OLS County	(4) IV State	(5) IV City	(6) IV County
EG index - all traded						
Input-output	0.227 (0.042)	0.149 (0.026)	0.125 (0.028)	0.244 (0.048)	0.128 (0.033)	0.105 (0.036)
Labor	0.144 (0.012)	0.109 (0.010)	0.077 (0.009)	0.173 (0.034)	0.166 (0.023)	0.080 (0.023)
Technology	−0.001 (0.010)	−0.028 (0.008)	−0.026 (0.009)	−0.003 (0.011)	−0.030 (0.008)	−0.024 (0.008)
Nat. advantages	0.239 (0.018)	0.323 (0.023)	0.272 (0.031)	0.223 (0.021)	0.300 (0.023)	0.274 (0.031)
Observations	7140	7140	7140	7140	7140	7140
R <sup>2</sup>	0.136	0.127	0.081	0.135	0.124	0.081

Robust standard errors in parentheses.

**Table 4**  
Growth in local industries – Intensive margin.

City-industry employment growth									
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
$\ln E_{ir03}$	−0.2003 (0.0034)	−0.2037 (0.0034)	−0.2003 (0.0034)	−0.2037 (0.0034)	−0.2007 (0.0034)	−0.2040 (0.0034)	−0.2045 (0.0035)	−0.2046 (0.0035)	−0.2054 (0.0035)
$\ln E_{ir03}^{IO}$	0.0283 (0.0045)		0.0285 (0.0045)		0.0282 (0.0045)		0.0153 (0.0046)	0.0156 (0.0046)	0.0154 (0.0046)
$\ln E_{ir03}^L$		0.0616 (0.0050)		0.0619 (0.0050)		0.0630 (0.0050)	0.0575 (0.0051)	0.0576 (0.0051)	0.0590 (0.0051)
$\ln E_{ir03}^{IO} \hat{\beta}_{i,EG}^{IO}$			0.0053 (0.0022)					0.0045 (0.0022)	
$\ln E_{ir03}^L \hat{\beta}_{i,EG}^L$				0.0045 (0.0020)				0.0043 (0.0020)	
$\ln E_{ir03}^{IO} \hat{\beta}_{i,LC}^{IO}$					0.0090 (0.0020)				0.0098 (0.0020)
$\ln E_{ir03}^L \hat{\beta}_{i,LC}^L$						0.0108 (0.0023)			0.0117 (0.0023)
Obs.	38,124	38,124	38,124	38,124	38,124	38,124	38,124	38,124	38,124
Adj.R <sup>2</sup>	0.1302	0.1332	0.1303	0.1332	0.1305	0.1335	0.1334	0.1335	0.1342

Robust standard errors are reported in parentheses. The dependent variable is the logarithm of employment growth in local industries between 2003 and 2008.  $E_{ir03}$  is the employment in industry  $i$ , city  $r$  and the year 2003.  $E_{ir03}^{IO}$  and  $E_{ir03}^L$  define the employment related to industry  $i$  according to input-output and labor linkages, respectively, in region  $r$  in 2003. The terms  $E_{ir03}^{IO} \hat{\beta}_{i,EG}^{IO}$  and  $E_{ir03}^L \hat{\beta}_{i,EG}^L$  refer to the interaction of related employment and the industry-specific effect-estimates of input-output or labor linkages in the coagglomeration regressions. Industry-specific coagglomeration effects are centered on their means and scaled by their standard deviations. All regressions include industry and region dummies.

variate OLS specifications for manufacturing industries.<sup>14</sup> To facilitate the interpretation of estimated effect sizes, all variables are rescaled such that they are expressed in units of standard deviations. We will apply this rescaling to all subsequent analyses as well. The left part of the table (columns 1–3) shows our own estimates; the original findings by EGK (Table 3, p.1204) are shown on the right (columns 4–6).

In spite of the fact that our sample refers to a different period and more aggregated industry definitions (manufacturing codes are less detailed at the 4-digit NAICS level than at the 3-digit SIC level), most of our estimates are remarkably close to those of EGK. Moreover, we replicate EGK's main conclusions: all three externality channels are important, but the strongest impact on coagglomeration is recorded for value-chain linkages, the weakest for technological linkages. The main difference between our results and the ones in EGK is that we find larger effects of labor pooling at the city and county levels. Furthermore, natural advantages have a weaker influence according to our estimates.

Next, we extend the EGK framework in several directions. First, we expand the sample of industries by including services. This results in a sample of 120 industries in the (adjusted) NAICS classification, roughly two thirds (83) of which belong to the manufacturing and one third to the services sector. Doing so also forces us to abandon the natural-advantages variable taken from EGK. Instead, we use the alternative implementation described in the data section. Furthermore, we switch to a multivariate regression framework and add IV estimations, using the similarity matrices based on non-US data as instruments. Marshallian effects on coagglomeration measured by the EG index are reported in Table 3. For results using the LC index, we refer to Appendix E.

Adding services and moving to a multivariate setting changes some of the point estimates, but does not alter the main conclusions. The most substantive change is that the effect of technological similarities on coagglomeration becomes insignificant, with some of the point estimates even turning negative. This change is in part due to the fact that we now estimate the effects of all externality channels simultaneously. However, also in univariate regressions merely including services in the sample halves the technological similarity effect (see Table F.2 of Appendix F). This reinforces our concerns about the use of this variable outside the manufacturing sector. Furthermore, the effects of the new natural-advantages variable are over three times stronger than those re-

<sup>14</sup> EGK restrict their sample to the manufacturing sector. To enhance comparability, we define manufacturing industries in this table at the 4-digit level of the original US NAICS classification. In the remainder of this section, we switch to the (US-Mexican) harmonized NAICS classification.

**Table 5**  
Trends in Marshallian agglomeration forces .

EG	OLS (1)	OLS (2)	OLS (3)	IV (4)	IV (5)	IV (6)
Labor	0.205 (0.025)		0.188 (0.025)	0.210 (0.026)		0.193 (0.026)
Labor $\times$ year	0.108 (0.034)		0.100 (0.034)	0.116 (0.036)		0.110 (0.036)
Input-output		0.341 (0.032)	0.315 (0.031)		0.384 (0.034)	0.364 (0.033)
Input-output $\times$ year		−0.170 (0.042)	−0.194 (0.041)		−0.211 (0.044)	−0.254 (0.044)
$R^2$	0.05	0.02	0.07	0.05	0.02	0.07
N	31,757	31,757	31,757	31,757	31,757	31,757

Time is expressed in units of 100 years. Regression includes dummies for each decade. Robust standard errors in parentheses.

ported in Table 3.<sup>15,16</sup> Moreover, adding this control variable cuts the estimated effect of labor linkages in half. This is unsurprising in light of our earlier discussion of the fact that what EGK call natural advantages includes characteristics of the local labor force. Given these concerns about the technology and natural-advantages measures, we henceforth drop these variables from the analysis.<sup>17</sup>

Using IV regressions, the labor channel regains some of its prominence in the state- and city-level estimates.<sup>18</sup> Unlike EGK, who base their instruments on data from the UK, our instrumental variables are derived from Mexican data and, therefore, further test the robustness of the original EGK results. Overall, however, the discrepancies between OLS and IV estimates are minor, suggesting that endogeneity plays no major role in these estimations.<sup>19</sup> Balancing the trade-off between bias and efficiency, our preferred estimates are the ones from OLS regressions. At the city level, these estimates imply that a one-standard-deviation increase in similarity in the human capital requirements of two industries increases their coagglomeration by 0.1 standard deviations, whereas a one-standard-deviation increase in the strength of value-chain linkages increases their coagglomeration index by 0.16 standard deviations.

### 3.2. Manufacturing versus services

We start our exploration of sectoral differences in Marshallian externalities by splitting the sample into a manufacturing and a services subsample. Fig. 2 shows the effects of labor market pooling and value chain linkages in each subsample. The different panels show estimated effects, together with their 95% confidence intervals, on coagglomeration patterns for manufacturing industries in blue and services in red. Estimates using the EG index are shown in the top panel, those using the LC index in the bottom panel. On the left, we plot univariate, on the right, multivariate estimates.<sup>20</sup>

The most striking aspect of Fig. 2 is the difference between services and manufacturing. Above the (dashed) 45 degree line, labor-pooling effects exceed value-chain linkages as an explanation for coagglomeration. Below this line, the opposite holds. Whereas

services typically appear above the 45 degree line, manufacturing industries tend to be located on, or below, the line. This means that coagglomeration patterns in services are more driven by labor-linkages than by value chains, whereas the opposite holds for manufacturing industries. Moreover, the point estimates for services always lie above and – in the bottom, LC-based, plots – to the right of the manufacturing estimates. This confirms our earlier contention that the nontraded nature and labor-intensity of services may make them particularly sensitive to Marshallian externalities: both externality channels are as or even more pronounced in the coagglomeration patterns of services than of manufacturing.<sup>21</sup>

### 3.3. Effect heterogeneity

To explore the observed sectoral heterogeneity in greater detail, we allow the effects of labor and value chain linkages to vary freely by industry, as in Eq. (8). To present outcomes in a more compact way, we average the resulting 120 industry-specific labor ( $\beta^L$ ) and input-output ( $\beta^{IO}$ ) effects for 27 broad subsectors and plot these averages in Fig. 3.<sup>22</sup> Estimates for services are colored red, for manufacturing blue. To avoid cluttering the graph, we only show city-level estimates.

The general patterns roughly track those foreshadowed in Fig. 2. Both types of inter-industry similarities tend to exert a greater influence on coagglomeration in the subsectors of services than of manufacturing. However, there is pronounced variation in agglomeration effects between (and, as shown in Appendix D, also within) subsectors, especially in labor-pooling effects. For instance, the labor channel's greatest impact on coagglomeration patterns is found in industries in *arts & culture*, *architecture & engineering*, *media* and *knowledge-intensive business services* (KIBS). In these subsectors, a one standard-deviation rise in the human capital similarity of two industries increases their EG index by up to 1.5 standard deviations.

Although less pronounced than in services, manufacturing industries also exhibit heterogeneous agglomeration effects. For instance, industries in *hardware* and *machinery* manufacturing tend to be driven by value chains and labor pooling relations. In contrast, *pharma & medical* and *electronics* industries coagglomerate with value-chain partners but not with industries that employ similar labor.

Hausman test. This test fails to reject the null hypothesis that regressors are exogenous in all cases, except for the labor pooling channel in city-level estimates.

<sup>20</sup> Tables F.3 and F.4 report the analogous coefficients derived from multivariate IV regressions. However, due to the collinear nature of the inter-industry proximity matrices our instruments are too weak to identify effects in services, where the Kleibergen–Paap statistic drops below 1.

<sup>21</sup> Note that in multivariate regressions of the LC index (Appendix E), the effect of labor-linkages even turns slightly (although not significantly) negative in two of the manufacturing samples.

<sup>22</sup> We report the full list of 120 industry-specific estimates in Appendix D.

<sup>15</sup> To some extent, this increase seems to be driven by some exceptionally large locations. When we drop the twenty largest counties, ten largest cities, or five largest states, the coefficients drop to levels closer to the ones reported in Table 3.

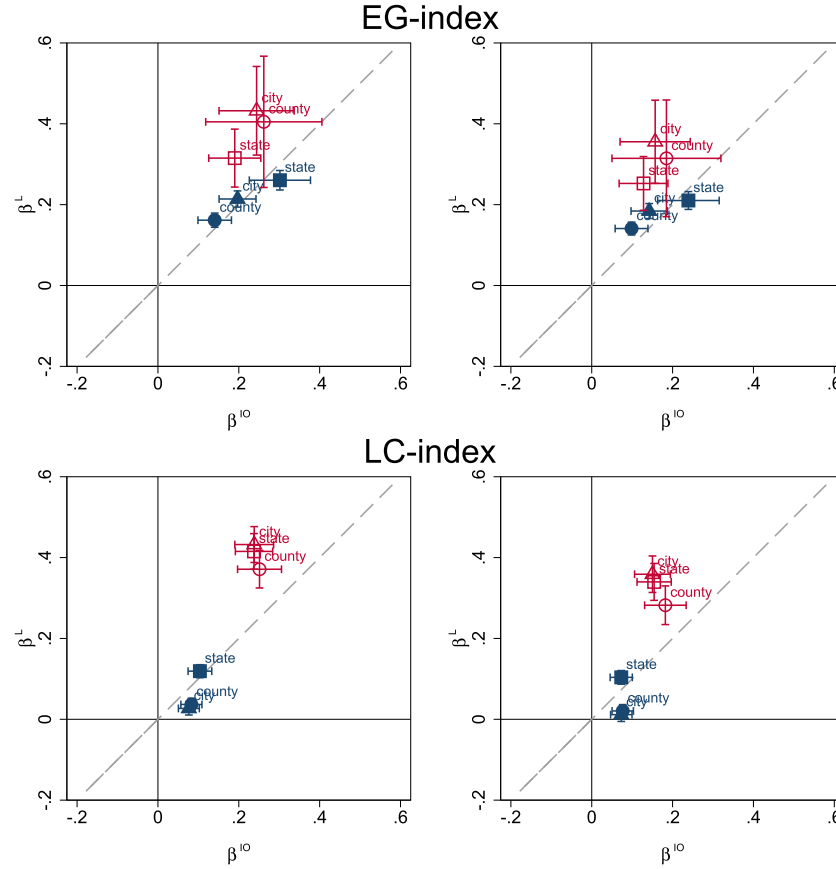
<sup>16</sup> Note that, given that natural advantages should be exogenous, these variables enter the models of columns 4–6 uninstrumented.

<sup>17</sup> We also ran robustness tests, where we retain natural advantages as a control variable to correct for potential omitted variable bias. The only difference we find is that the estimated OLS effect of labor linkages in services is attenuated when using the EG-index as a dependent variable. All other estimates – which are available upon request – are consistent with what we report here.

<sup>18</sup> Note that this also occurs in a univariate analysis (Appendix F). Adding natural-advantages as a control variable thus also lowers estimated labor pooling effects in IV regressions.

<sup>19</sup> We formally test for endogeneity in each Marshallian channel using the difference-in-Sargan statistic, a heteroscedasticity robust variant of the standard





**Fig. 2.** Coagglomeration effects, manufacturing versus services.

Estimated effects on coagglomeration of services are plotted with red hollow markers; of manufacturing industries with blue solid markers. Labor pooling effects are depicted on the vertical axis, value chain effects on the horizontal axis, using the univariate estimations of Eq. (7). State-level estimates are marked by squares, city-level estimates by triangles and municipality-level estimates by circles. IV estimates use labor pooling and input-output linkages based on Mexican data as instruments. The crosshairs represent 95% confidence intervals based on robust standard errors. The dashed gray line is a 45 degree line.

#### 4. Marshallian externalities and regional diversification

Using the EGK framework, we have shown that the strength of labor-market-pooling and value-chain-based agglomeration externalities differs widely across industries. To test the predictive validity of the underlying estimations, we now analyze whether the documented differences across industries manifest themselves in industries' local growth paths. If the measured effect-heterogeneity reflects the importance of agglomeration forces, it should not just affect how industries coagglomerate, but also help predict local growth patterns.

To investigate this, we build on a large and growing literature that focuses on the role of Marshallian externalities in the evolution of regions' industrial structures. According to this literature, inter-industry linkages are important for understanding diversification processes, because regions tend to branch out into new economic activities that build on a region's pre-existing strengths. This claim is supported by the fact that industries grow faster in regions with substantial employment in related industries (Porter, 2003; Greenstone et al., 2008; Delgado et al., 2010; Dauth, 2010; Jofre-Monseny et al., 2011; Neffke et al., 2011; Hausmann et al., 2014).

Most papers in this tradition implicitly assume that all industries benefit equally from inter-industry spillovers. Moreover, they typically do not distinguish among different types of inter-industry linkages. In this section, we use our estimated industry-specific effects on coagglomeration to enhance the econometric models of local-industry growth-rates used in this line of research. To do so,

let the relatedness between industries  $i$  and  $j$  be measured by one of two proximity measures,  $prel_{ij} \in \{p_{ij}^L, p_{ij}^O\}$ . We can now calculate the proximity-weighted employment for a region-industry,  $(i, r)$ , in year  $t$  as:

$$E_{irt}^{rel} = \sum_j \frac{prel_{ij}}{\sum_{k \neq i} prel_{ik}} E_{jrt}.$$

This expression can be interpreted as the amount of related employment that already exists in the local economy, where what we call "related" depends on  $rel$ . For example, in the case of labor,  $E_{irt}^{rel}$  is an index that reflects the size of the local workforce with skills and know-how that are relevant to industry  $i$ .

Letting  $G_{ir}^{03-08}$  refer to the logarithm of employment growth of industry  $i$  in region  $r$  between 2003 and 2008, i.e.,  $G_{ir} = \ln(E_{ir08}) - \ln(E_{ir03})$  and restricting the sample to local industries with non-zero employment in 2003, we estimate the following:

$$G_{ir}^{03-08} = \delta \ln(E_{ir03}) + \sum_{rel} \gamma^{rel} \ln(E_{ir03}^{rel}) + \iota_i + \rho_r + \epsilon_{ir03} \quad (9)$$

where  $\delta$  captures mean reversion effects, and  $\iota_i$  and  $\rho_r$  are industry and region fixed effects.

The parameter of interest is  $\gamma^{rel}$ . It quantifies to what extent the growth in a local industry can be predicted from the amount of related activity that already exists in the local economy.

Table 4 contains results when using cities as the spatial unit of analysis. In line with prior studies (e.g., Delgado et al., 2010; Hausmann et al., 2014), we find significant and negative mean reversion

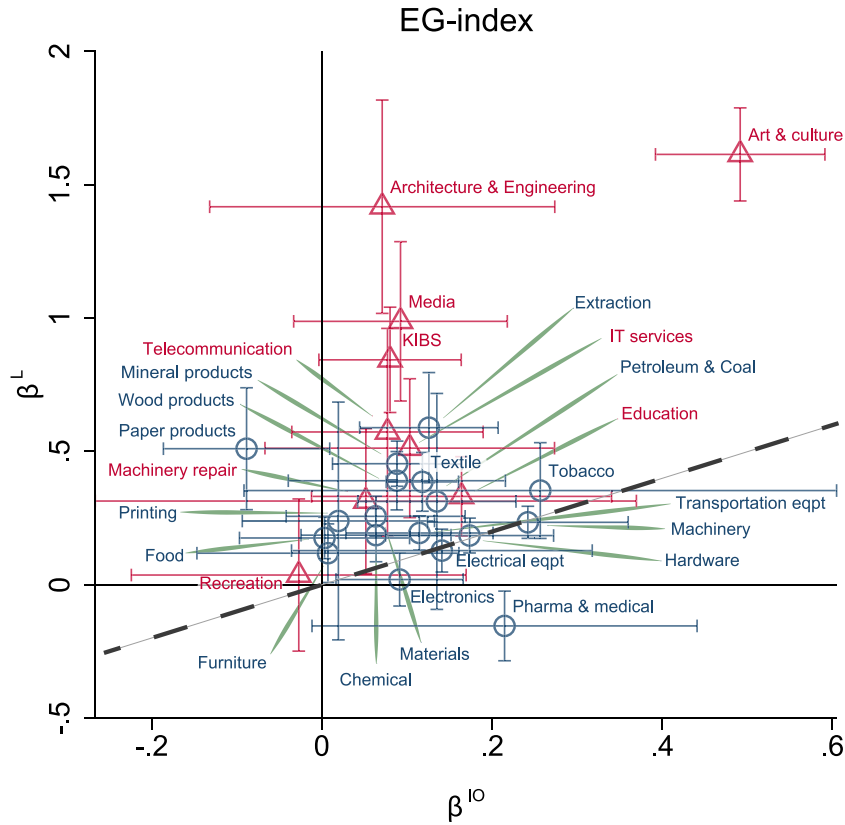


Fig. 3. Coagglomeration effects for 27 subsectors (EG-index).

Average point estimates for industries in 27 subsectors (see Appendix D for results by industry) of labor pooling effects (vertical axis) and value chain effects (horizontal axis) on EG coagglomeration index, using OLS regression. Estimates for services are marked with red triangles, for manufacturing industries with blue circles. The crosshairs represent 95% confidence intervals for these averages, based on robust standard errors.

effects, and positive effects of related employment, regardless of whether we measure relatedness in terms of labor or input-output linkages.

However, if the estimated  $\hat{\beta}_i^{rel}$  parameters of Section 3.3 capture any information on the importance of Marshallian channels, we would expect that industries with a higher  $\hat{\beta}_i^{rel}$  will benefit more from the associated agglomeration externality when it comes to their growth.

To explore this, we augment 9 by interacting the amount of related employment in a city with  $\hat{\beta}_i^{rel}$ . To facilitate interpretation,  $\hat{\beta}_i^{rel}$  values are expressed in units of standard deviations away from their respective means. This yields the following:<sup>23</sup>

$$G_{ir}^{03-08} = \delta \ln(E_{ir03}) + \gamma^{rel} \ln(E_{ir03}^{rel}) + \gamma_{\beta}^{rel} \hat{\beta}_i^{rel} \ln(E_{ir03}^{rel}) + \iota_i + \rho_r + \epsilon_{ir03} \quad (10)$$

where we expect  $\gamma_{\beta}^{rel}$  to be positive. Table 4 shows that this is indeed the case. The interaction effect,  $\gamma_{\beta}^{rel}$ , is always positive and significant.

The marginal effects of related employment are given by

$$\frac{\partial G_{ir}^{03-08}}{\partial E_{ir}^{rel}} = \gamma^{rel} + \gamma_{\beta}^{rel} \hat{\beta}_i^{rel}.$$

Fig. 4 plots these partial derivatives against  $\hat{\beta}_i^{rel}$  on a range that reflects the span of coefficient estimates obtained in Section 3.3. The interaction effects are strongest when the  $\hat{\beta}_i^{rel}$  are based on

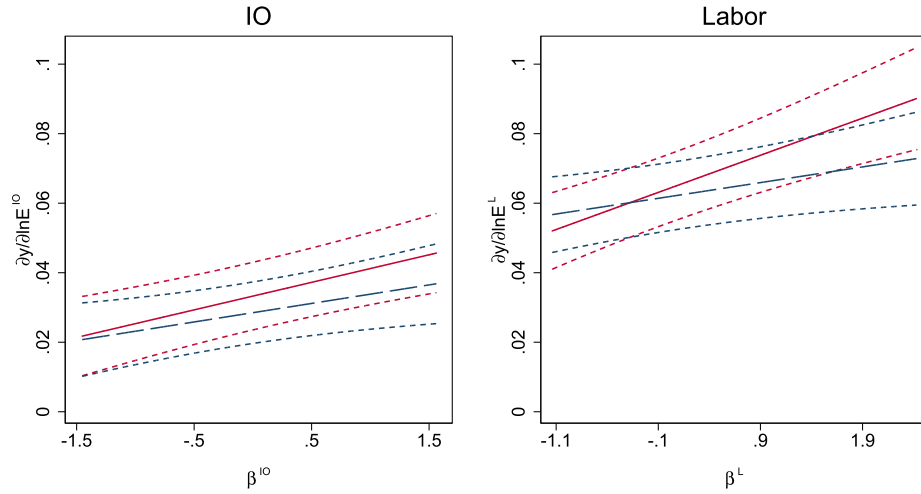
regressions of the LC index. In this case, the elasticities for value-chain related employment run from as low as 2% to as high as 4%. The range of elasticities is slightly wider for employment that is related in terms of human-capital similarities, starting from a low 5% for industries with coagglomeration patterns that are relatively insensitive to human capital similarities and reaching over 8% where such similarities are strong drivers of coagglomeration. The fact that the industry-specific coagglomeration effects are also reflected in industries' local growth patterns validates the effect-heterogeneity reported in Section 3.3.<sup>24</sup> In Appendix C, we re-estimate the growth model in Eq. (9) for growth at the extensive margin, estimating the entry probability for local industries that do not yet exist in 2003. Here, the range of elasticities runs from about -1% and 2% for value-chain effects, whereas labor linkages have a constant effect of about 1%. This suggests that a pool of local suppliers or buyers may also attract new industries to a city, but only those industries whose location patterns are strongly affected by the spatial distribution of value-chain partners.

## 5. Trends in Marshallian forces

So far, we have concentrated on cross-sectional heterogeneity in agglomeration externalities. However, in the past 100 years, the US economy underwent significant transformations, including ur-

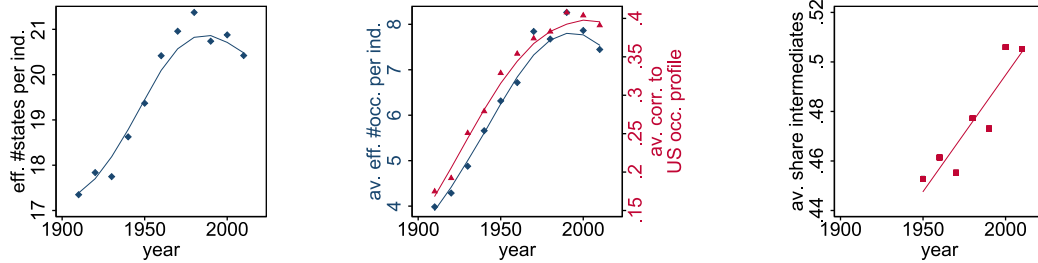
<sup>23</sup> Note that level effects of  $\beta_i^{rel}$  will be absorbed in  $\iota_i$ .

<sup>24</sup> The industry-specific parameters estimated using the LC index outperform the ones using the EG index in these growth projections. This suggests that the LC coagglomeration metric may be a better indicator of inter-industry spillover potentials than the EG index. Although interesting in its own right, we will not pursue this issue further.



**Fig. 4.** Marginal effects of related employment.

Marginal effects of related employment with 95% confidence intervals, computed as the partial derivative of growth with respect to related employment. The left panel depicts IO effects, the right panel labor-pooling effects. The horizontal axis plots  $\beta_{IO}$  and  $\beta_L$ , respectively, on a range that reflects the distribution of coefficient estimates obtained in Section 3.3. Blue dashed lines use estimates from coagglomeration regressions using the EG-index, red solid lines using the LC-index.



**Fig. 5.** Historical trends, descriptive statistics.

Left: geographical concentration. Effective number of states per industry. Center: occupational composition. The left scale (blue diamonds) is the average effective number of occupations per industry, while the right scale (red triangles) is the average correlation to the US occupational profile. Right: Intermediates share. Average share of intermediates for all industries.

banization, mass education, de-industrialization, tremendous reductions in transportation and communication costs, fragmentation and internationalization of supply chains, and a secular shift from agriculture and manufacturing activities to services. Have these phenomena collectively affected the relative strength of different Marshallian externalities? For instance, the increasing spatial fragmentation of supply chains that was facilitated by new transportation and communication technologies may have made spatial proximity to value-chain partners less crucial than it once was. Similarly, if the reliance on a highly skilled and specialized work force has increased over time, is this reflected in an increasing role of labor similarity in determining coagglomeration? To shed light on these issues, we turn to historical data on the US economy.

### 5.1. Descriptive statistics

We start by describing how the agglomeration, occupational structures, and input-output structures of US industries have changed over time. Fig. 5 shows weighted average statistics for industries in the US economy across decades. To ensure that we are picking up changes in industry characteristics, and not in the economy's industrial composition, weights are kept constant at industries' average employment shares across all decades.

The left panel shows the evolution of the geographical concentration of industries. The graph plots the weighted average effective number of US states in which an industry is active, as expressed by the inverse of the Herfindahl–Hirschman Index (HHI) of

an industry's employment shares across states.<sup>25</sup> Until the 1980s, on average, industries have spread their employment more evenly across space. However, since then, the curve has flattened, indicating that this deconcentration process has ceased.

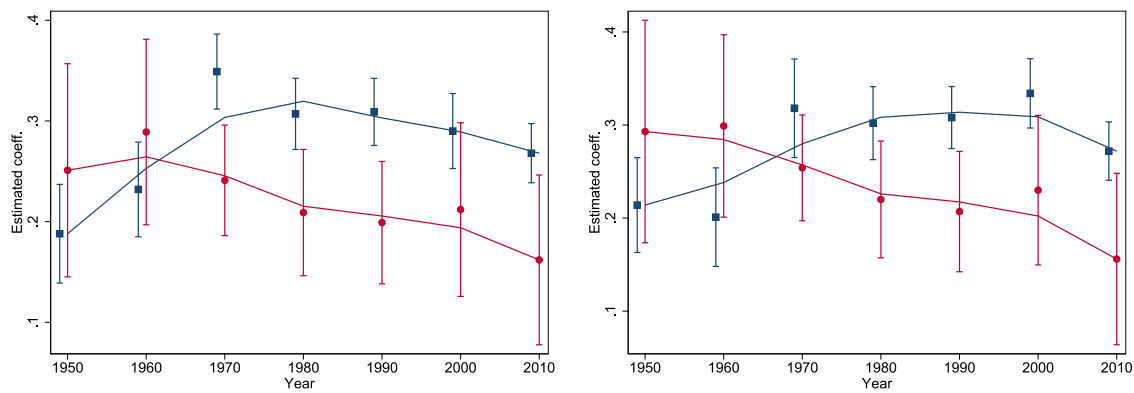
The center graph shows how the occupational structure of industries has changed in the same period. The blue line depicts the effective number of occupations per industry, the red line shows the (weighted) average similarity of industries' occupational structures to the occupational composition of the US economy as a whole.<sup>26</sup> The effective number of occupations per industry rises over the decades – again, until roughly the 1980s – suggesting that the division of labor within industries increases.<sup>27</sup> However, at the same time, this division of labor has become more uniform, with industries' occupational compositions becoming more similar

<sup>25</sup> If an industry's employment is spread out equally across  $N_s$  states, the inverse of the HHI,  $\frac{1}{HHI_{it}} = \frac{1}{\sum_s (\frac{1}{E_{it}})^2}$ , attains a value of  $\frac{1}{\sum_s (\frac{1}{N_s})^2} = \frac{1}{N_s (\frac{1}{N_s})^2} = N_s$ . Therefore, a

given value of  $\frac{1}{HHI_{it}}$  can be interpreted as the “effective” number of states in which an industry is present, i.e., the number of states needed to achieve this value had the industry's employment been distributed equally across them.

<sup>26</sup> The latter is calculated as  $corr(E_{i0t}, \sum_t E_{i0t})$ , where the correlation is taken across occupations.

<sup>27</sup> The flattening of the curve towards the end of the period may reflect real changes in specialization patterns, or simply that the historical occupational classification system becomes less and less adequate to describe our modern economy's job structure.



**Fig. 6.** Trends in Marshallian agglomeration forces - OLS (left) and IV (right).

Estimated impact of input-output (red circles) and labor (blue squares) linkages on EG-coagglomeration in the period 1950–2010, with 95% confidence intervals based on robust standard errors. The lines represent LOWESS smooths.

to the economy's overall employment composition, suggesting that workers can be shared across a wider set of industries.

The final graph shows how value chains have developed over time. It plots the value of intermediates that industries use as a share of total output. Because input-output data are not available before 1947, this graph is limited to the period from 1950 to 2010. In spite of this, in line with the notion that value chains have become progressively fragmented, the trend towards greater outsourcing is clearly visible. Over the course of six decades, the weighted average input share rose by four pp, from 43% to 47%.

## 5.2. Marshallian externalities

Have the changes described in Fig. 5 impacted on agglomeration forces? In particular, we'd expect that increased outsourcing would diminish the importance of geographical proximity to value chain partners. Moreover, the documented rise in the division of labor, together with the fact that workers can be shared more widely across industries, may have increased the benefits of opportunities to share skilled workers.

Fig. 6 shows the evolution of how labor and value-chain proximity have affected coagglomeration patterns in the period 1950–2010. The graph on the left shows parameter estimates using OLS regressions, while the graph on the right is based on IV estimates. Note, however, that we do not have data on the Mexican economy for this period. Therefore, we cannot construct the same instruments as in Section 3. Instead, we instrument the proximities in one decade with the proximities in another decade.<sup>28</sup> This does not mitigate reverse causality problems, which, given our earlier estimates, we believe are relatively unimportant. However, doing so does help address the errors-in-variables bias that may result from using mismeasured proximities.

Both graphs display clear downward trends in the effect of value-chain linkages on coagglomeration. This confirms the conjecture that locating close to value-chain partners has become less relevant as transportation and communication technologies have improved. Instead, the labor-pooling effects have sharply increased in the early decades but then stabilized or even mildly declined towards the end of the period. Overall, however, in line with the notion that workers have become more specialized, labor pooling is more important nowadays than in the 1950s.

To test the existence of these trends more formally, we pool all data across decades and interact the labor and value-chain

proximities with a linear time-trend. We express time in units of 100 years and – for reasons of consistency with the subsequent section – set the year 1910 to zero. Consequently, the level effects ( $t = 0$ ) reflect the (extrapolated) situation in 1910, whereas the sum of level and interaction effect (i.e.,  $t = 1$ ) refers to the decade of 2010. As shown in Table 5, both trends are statistically significant. The estimated parameters imply that, in 1950, a one-standard-deviation increase in labor similarities increased industries coagglomeration by 0.259 ( $0.205 + 0.5 \times 0.108$ ) standard deviations. This effect rises to 0.324 ( $0.205 + 1.1 \times 0.108$ ) in 2010. In contrast, the effect of a one-standard-deviation stronger input-output linkages drops from a 0.256 to a 0.154 standard-deviations increase in the (EG-)coagglomeration of two industries.<sup>29</sup>

In principle, the IPUMS data allow us to measure coagglomeration and labor-similarity matrices for industries from 1910 onwards, albeit using a different industry classification system.<sup>30</sup> The only constraint is that we have no input-output tables for years before 1947. However, it turns out that the results of Table 5 are not driven by changes in the value-chain matrices. In fact, as shown in Appendix F, Table F.5, fixing the input-output matrices to the 1950 version, we find the same results as when using contemporaneous matrices. That is, input-output linkages have been constant enough to reliably estimate their impact on coagglomeration using a time-invariant matrix. We exploit this fact to extend our historical analysis backwards to 1910.

The results of this exercise are depicted in Fig. 7.<sup>31</sup> They show that the decline of input-output linkages as a driver of agglomeration set in long before 1950, whereas, viewed over the entire period, there is no significant increase (nor decrease) in the effect of labor linkages.<sup>32</sup> Overall, these analyses confirm that Marshallian agglomeration forces have changed significantly over the

<sup>29</sup> In Appendix F, Table F.5 we add 2-way fixed effects to the estimation to control for any industry idiosyncrasies. This leads to even more pronounced time trends. Moreover, we find similar trends if we exclude services, as shown in Fig. F.1 which uses a sample of only manufacturing industries. Consequently, the changes in labor-pooling and input-output effects are not simply the result of the US economy shifting towards more and more services.

<sup>30</sup> Before 1910, the US census only recorded an individual's occupation, not her industry. The industry codes recorded in IPUMS in that period are imputed from these occupation codes, making them unusable for our purposes.

<sup>31</sup> Table F.5 of Appendix F shows the results of a regression that interacts externality channels with a time trend, analogous to Table 5.

<sup>32</sup> Closer inspection reveals that the mild positive trend between 1950 and 2010 that was reported in Table 5 is balanced by a negative trend between 1910 and 1950. This U-shaped pattern is weakly corroborated when we include an interaction of labor-linkages with a squared time-trend, which locates the minimum of the labor-pooling effect in 1967. However, the interacted coefficients are only jointly significant with the trend variables themselves.

<sup>28</sup> For all decades but the first, we use as instruments the proximities measured in the preceding decade. In the first decade, 1950, we use the matrices for 1960 as instruments.



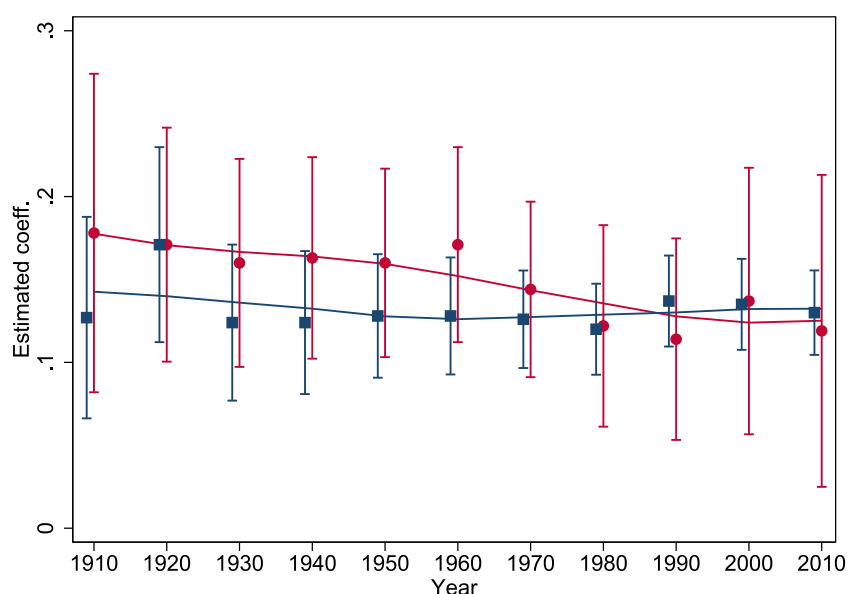


Fig. 7. Trends in Marshallian agglomeration forces, long-run perspective (OLS).

Estimated impact of input-output (red circles) and labor (blue squares) linkages on EG-coagglomeration in the period 1910–2010, with 95% confidence intervals based on robust standard errors. The lines represent LOWESS smooths.

course of a century. Moreover, if these estimated trends were to continue, value chains will lose their potency as an agglomerating force, leaving the skills of the workforce as the principal remaining driver of agglomeration in the cities of the future.

## 6. Conclusion

Extending the work by Ellison et al. (2010), we find substantial variation in why industries coagglomerate. Whereas in some industries, firms tend to locate close to their value-chain partners, in others, they tend to locate near industries that share their labor requirements. In fact, the largest labor pooling effects are found in the coagglomeration patterns of services. Extreme examples are industries in *architecture & engineering, media and knowledge-intensive business services*. Value-chain effects are generally weaker and comparatively strong in manufacturing. However, overall we find that human capital similarities and, although less so, input-output linkages, have a stronger effect on the coagglomeration of services than of manufacturing industries.

This heterogeneity not only affects coagglomeration, but also helps improve predictions of industries' local growth rates: whereas there are ample studies that show that local industries tend to grow faster in regions with substantial employment in related industries, we are able to show that this relatedness matters more whenever the corresponding inter-industry proximity is more strongly expressed in an industry's coagglomeration patterns.

Finally, we find that the relative importance of different Marshallian agglomeration channels has changed over time. Using coagglomeration patterns in the period 1910–2010, we find that value-chain linkages were dominant in determining which industries coagglomerated at the beginning of the twentieth century. However, their influence has weakened significantly since. At the same time, labor pooling effects have, on the whole, maintained their importance in coagglomeration patterns. Overall, we, therefore, conclude that whereas value chains were the main organizing principle of the spatial configuration of industries at the start of the twentieth century, their importance has dropped to a point where they have been overtaken by local pools of specialized labor as the main driver of industrial coagglomeration.

The fact that Marshallian forces operate more strongly on the locational patterns of services than of manufacturing has important implications for the future of cities. In spite of tremendous improvements in transportation and communication technologies – which should have allowed economic activities to become more dispersed – the increasing dominance of services in modern economies means that agglomeration externalities are likely to become more, not less, important. What is more, not only do services benefit particularly strongly from labor pooling externalities, labor pooling has also overtaken access to value chain partners as the dominant reason for why industries coagglomerate. As a result, our findings suggest that large cities increasingly derive their strength from the ease with which skills circulate in their economies.

## Acknowledgments

Dario Diodato acknowledges the financial support of NWO, Innovation Research Incentive Scheme (VIDI). Project number: 45211013. Frank Neffke is grateful to the MasterCard Center for Inclusive Growth & Financial Inclusion for their financial support. The authors would like to thank Ron Boschma, Ricardo Hausmann, David Rigby, Nicoletta Corrocher and the attendees at the Global Conference on Economic Geography 2015 in Oxford, the Geography of Innovation 2016 conference in Toulouse and the CID Growth Lab seminar for their useful comments and suggestions.

## Appendix A. Descriptives

Table A.1 and shows the top 10 most coagglomerated industry pairs. When calculating coagglomeration scores, we exclude the ten largest cities to avoid that this top 10 gets cluttered by patterns generated by industries that locate in extremely large cities. Although all our estimations use coagglomeration scores calculated with the full set of cities, reported results (available on request (see supplementary material online)) are robust to excluding the largest cities.

**Table A.1**

Top-10 industry pairs, EG coagglomeration index.

Top-10 co-location (EG index)		
Industry i	Industry j	Value
Software publishers	Aerospace product manufacturing	0.0276
Textile furnishings mills	Fiber, yarn, and thread mills	0.0233
Motor vehicle manufacturing	Metalworking machinery manufacturing	0.0231
Tobacco manufacturing	Fiber, yarn, and thread mills	0.0225
Motor vehicle parts manufacturing	Motor vehicle manufacturing	0.0212
Seafood preparation and packaging	Aerospace product manufacturing	0.0178
Motor vehicle manufacturing	Leather and hide tanning	0.0170
Software publishers	Communications equipment manufacturing	0.0167
Software publishers	Computer equipment manufacturing	0.0165
Software publishers	Seafood preparation and packaging	0.0163

Top-10 most coagglomerated industry pairs in terms of EG index (Eq. (1)), using city-industry employment data for the US (County Business Patterns, 2003). The industry classification is based on NAICS 4-digit industries and includes 120 distinct industry codes containing both services and manufacturing. The top-10 cities have been excluded from the computation to reduce noise.

**Table A.2**

Top-10 industry pairs, occupational similarity.

Top-10 labor links		
Industry i	Industry j	Value
Other textile product mills	Other apparel manufacturing	0.9919
Other apparel manufacturing	Cut and sew apparel manufacturing	0.9885
Other textile product mills	Cut and sew apparel manufacturing	0.9805
Other transportation equipment	Agri/construction/mining machinery	0.9802
Motor vehicle parts manufacturing	Hardware manufacturing	0.9766
Other transportation equipment	Motor vehicle body manufacturing	0.9760
Motor vehicle body manufacturing	Heating/cooling equipment	0.9756
Other machinery manufacturing	Agri/construction/mining machinery	0.9742
Household appliance manufacturing	Hardware manufacturing	0.9737
Heating/cooling equipment	Hardware manufacturing	0.9715

Top-10 industry pairs in terms of occupational similarity, computed as the correlation of industries' occupational employment shares (Eq. (4)), using industry-occupation data for the year 2002 as reported in the US Occupational Employment Statistics (OES).

**Table A.3**

Top-10 industry pairs, input-output links.

Top-10 input-output links		
Industry i	Industry j	Value
Petroleum and coal manufacturing	Oil and gas extraction	0.6776
Other apparel manufacturing	Cut and sew apparel manufacturing	0.6001
Leather and hide tanning	Animal slaughtering and processing	0.5811
Motor vehicle parts manufacturing	Motor vehicle manufacturing	0.5720
Motor vehicle manufacturing	Motor vehicle body manufacturing	0.5351
Pulp, paper, and paperboard mills	Paper product manufacturing	0.4617
Support activities for mining	Oil and gas extraction	0.4575
Motor vehicle manufacturing	Audio-video equipment manufacturing	0.4343
Spectator sports	Radio and television broadcasting	0.4269
Motor vehicle parts manufacturing	Leather and hide tanning	0.4227

Top-10 industry pairs in terms of input-output linkages, defined as the maximum relative importance of one industry as a customer or as a supplier of the other and vice versa (Eq. (3)), based on make-and-use tables provided by the US Bureau of Economic Analysis (BEA) for the year 2002.

**Table A.4**  
Top-10 industry pairs, technology links.

Top-10 technology links		
Industry i	Industry j	Value
Finance and insurance	Computer equipment manufacturing	1.2346
Finance and insurance	Communications equipment manufacturing	0.8123
Other telecommunications	Communications equipment manufacturing	0.4905
Other telecommunications	Computer equipment manufacturing	0.4622
Radio and television broadcasting	Communications equipment manufacturing	0.3605
Service industry machinery manufacturing	Other telecommunications	0.3207
Management consulting services	Computer equipment manufacturing	0.2913
Manufacturing of magnetic and optical media	Finance and insurance	0.2892
Management consulting services	Communications equipment manufacturing	0.2742
Computer equipment manufacturing	Communications equipment manufacturing	0.2595

Top-10 industry pairs in terms technological linkages, defined as the maximum overrepresentation of citation flows between industries' patents (Eq. (5)), using USPTO data for the years 1975–1999.

**Table A.5**  
Top-10 industry pairs, natural advantage similarities.

Top-10 natural advantage (EG)		
Industry i	Industry j	Value
Lime and gypsum product manufacturing	Iron and steel mills	0.0084
Wood product manufacturing	Sawmills and wood preservation	0.0078
Pulp, paper, and paperboard mills	Iron and steel mills	0.0070
Sawmills and wood preservation	Other wood product manufacturing	0.0067
Steel product manufacturing	Iron and steel mills	0.0067
Iron and steel mills	Fiber, yarn, and thread mills	0.0062
Iron and steel mills	Alumina and aluminum production	0.0061
Sawmills and wood preservation	Iron and steel mills	0.0059
Petroleum and coal manufacturing	Iron and steel mills	0.0054
Wood product manufacturing	Iron and steel mills	0.0053

Top-10 industry pairs in terms of shared use of natural advantages (Eq. (B.3)).

**Table A.6**  
Correlation coefficients.

	EG (US)	LC (US)	IO (US)	Labor (US)	Tech (US)	N-EG (US)	N-LC (US)	EG (MX)	LC (MX)	IO (MX)	Labor (MX)	Tech (IV)
EG (US)	1.00											
LC (US)	0.27	1.00										
IO (US)	0.20	0.15	1.00									
Labor (US)	0.23	0.04	0.24	1.00								
Tech (US)	0.04	0.05	0.13	0.15	1.00							
NA EG (new)	0.31	0.03	0.15	0.37	0.11	1.00						
NA LC (new)	0.14	0.21	0.06	0.25	0.09	0.45	1.00					
EG (MX)	0.19	0.04	0.10	0.16	−0.02	0.16	0.07	1.00				
LC (MX)	0.04	0.49	0.06	0.05	−0.02	−0.00	0.23	0.20	1.00			
IO (MX)	0.13	0.10	0.53	0.21	0.11	0.15	0.03	0.14	0.05	1.00		
Labor (MX)	0.16	−0.02	0.29	0.54	0.13	0.17	0.07	0.19	−0.02	0.28	1.00	
Tech (IV)	0.04	0.05	0.14	0.15	0.99	0.11	0.09	−0.02	−0.02	0.11	0.14	1.00

Correlation table of pairwise inter-industry linkages. Variables are created using US (US) or Mexican (MX) data sources. NA EG (new) and NA LC (new) refer to the EG- and LC-based measures of shared natural advantages using the authors' own calculations (see Section Appendix B). Tech (US) refers to technological similarity between industries measured by patent citations among patents by US inventors. Tech (IV) is an analogous measure using patents by non-US inventors.

## Appendix B. Variables construction

### Industry classification

In Section 3, US industries are classified according to the North American Industry Classification System (NAICS). Whereas data for 2003 use the NAICS2002 classification, 2008 data – which are used in Section 4 to calculate 5-year growth rates – are recorded using NAICS2007 classes. Moreover, although the Mexican classification systems are based on NAICS as well, they are not fully harmonized with those in the US. Therefore, we aggregate industries where necessary and create a new composite industry classification, reducing the original 317 4-digit NAICS codes to 215 harmonized industry classes. After dropping industries with missing data, we are left with 184 industries. In our main analysis, we restrict this number further to 120 traded industries (see Appendix D)

To construct the coagglomeration measure for the historical analysis, we use industry-state employment matrices by aggregating US census samples provided by IPUMS (Ruggles et al., 2017). IPUMS provides harmonized industry codes across decades in two classification systems: IND1950, available for all decades and IND1990, which is available between 1950 and 2010. When studying the period 1950–2010, we use the IND1990 classification, with some small corrections proposed by Autor and Dorn (2013). When extending the analysis back to 1910, we use the somewhat coarser IND1950 classification.

### Input-output

We construct inter-industry value-chain matrices using three different data sources. For the US, input-output linkages are based on make and use tables provided by the Bureau of Economic Anal-

ysis (BEA) for the year 2002. To merge these data to our coagglomeration data, we create a concordance between the 337 IO codes used by the BEA and our (harmonized) NAICS codes. Whenever one IO code corresponds to multiple NAICS codes, we distribute the flow of intermediates to or from the IO industry across the corresponding NAICS codes in proportion to each NAICS code's employment as reported in the CBP. Next, we use the resulting aggregated supply (make) matrix,  $S$ , and demand (use) matrix,  $U$ , to construct an IO-matrix using  $IO = SD_S^{-1}U$ , where  $D_S$  is a matrix with the column sums of  $S$  on its diagonal and all off-diagonal elements equal to zero.

Similar calculations are carried out for Mexico, using input-output data for the year 2008 provided by the Mexican statistical office, Instituto Nacional de Estadística y Geografía (INEGI). The ten strongest input-output linkages (based on US data) are shown in Table A.3.

For the analyses in Section 5, we use BEA historical IO tables from 1947 to 2012. Throughout this period, BEA uses a consistent industry classification that is closely linked to NAICS. Data are divided into three separated waves (1947–1962, 1963–1996, and 1997–2012). The level of detail of the classification system differs by wave. Therefore, we aggregate all data to the industrial classification used in the wave of 1947–1962. We convert this classification to the IPUMS IND1990 and IND1950 classifications.<sup>33</sup> Finally, because the census data in IPUMS are collected once a decade, we compute 5-year averages of these input-output tables, centered on the years 1950, 1960, 1970, 1980, 1990, 2000 and 2010.

#### Labor linkages

The labor similarity variable used in Section 3 is derived from the US Occupational Employment Statistics (OES) of the year 2002 as published by the Bureau of Labor Statistics. To instrument this variable, we aggregate the Mexican Encuesta Nacional de Ocupación y Empleo (ENOE) for the year 2005 at the industry-occupation level. The fact that the US and Mexico use different occupational classification systems does not complicate the calculation of occupational similarity matrices, because industries are recorded in NAICS codes in both data sources. However, because the ENOE uses a mix of 3- and 4-digit industry classes (consisting of 68 3-digit and 113 4-digit codes), we have to split some of the Mexican 3-digit industries into 4-digit codes. As a result, some industry pairs are mechanically attributed the same  $P^L$  values, and somewhat under 3% of off-diagonal elements are equal to 1.

For the historical analysis in Section 5, we aggregate the IPUMS census samples to the industry-occupation level. For the period 1950–2010, we use the harmonized OCC1990 codes, whereas in the 1910–2010 period, we use the OCC1950 classification. With up to over 380 different codes, the OCC1990 classification is very detailed. Due to the limited size of the IPUMS samples, this risks imprecise estimates of industries' occupational structures, especially if industries are small. Therefore, we aggregate these codes into 82 occupational segments. In contrast, the OCC1950 industries are grouped by only a dozen occupational segments. Moreover, because, with about 200 different occupation codes in a typical decade, the OCC1950 classification itself is not as detailed as the OCC1990 classification, we do not aggregate occupations further in the 1910–2010 period.

#### Technological similarity

The technological similarity matrices are based on patent citations extracted from the NBER database (Hall et al., 2001). This data set contains 16 million patents filed between 1975 and 1999.

We divide these patents into two samples, depending on whether they were filed by inventors residing inside or outside the US, dropping patents that list both, US and non-US inventors. After grouping the patents by their main patent class, we compute the total number of citations among these patent classes. The result is a  $428 \times 428$  matrix of cross-patent-class citations. Using a concordance developed by Goldschlag et al. (2016), we map patent classes to the NAICS industry classification. Next, following EGK, we estimate the degree of industries' technological similarity as described in Eq. (5). The instrument for this metric is created by repeating this process with the sample of patents filed by non-US inventors.

#### Natural advantages

The estimate of natural advantages in EGK is based on the following nonlinear estimation:

$$s_{ir} = \frac{E_r^{\alpha_0} E_{r(m)}^{\alpha_1} \exp(\mathbf{Xb})}{\sum_r E_r^{\alpha_0} E_{r(m)}^{\alpha_1} \exp(\mathbf{Xb})} \quad (\text{B.1})$$

where  $s_{ir}$  is the employment share of industry  $i$  in region  $r$ , while  $E_r$  and  $E_{r(m)}$  are the total employment and manufacturing employment in region  $r$ .<sup>34</sup> The term  $\mathbf{Xb}$  is a linear combination of natural advantage factors  $k \in K$ , such that  $\mathbf{Xb} = \sum_k \beta_k y_{kr} z_{ki}$ . In this expression,  $y_{kr}$  measures the abundance of  $k$  in region  $r$ , while  $z_{ki}$  quantifies how strongly industry  $i$  relies on factor  $k$ .

The estimated employment shares,  $\hat{s}_{ir}$ , are subsequently substituted for their actual counterparts in the EG-coagglomeration index:

$$EG_{ij} = \frac{\sum_{r=1}^R (\hat{s}_{ir} - x_r)(\hat{s}_{jr} - x_r)}{1 - \sum_{r=1}^R x_r^2}. \quad (\text{B.2})$$

This expression can be seen as the *expected* coagglomeration of two industries based on their dependence on the natural advantages in  $K$ .

To replicate the natural-advantage estimates in EGK, we plug two different estimates of  $s_{ir}$  into (B.2). First, we retrieve EGK's own nonlinear estimates of  $\hat{s}_{ir}$  from their supplementary material. Next, we convert these shares to the NAICS classification using the following formula

$$\tilde{s}_{ir} = \frac{\sum_i (E_r(m) \hat{s}_{ir} \times T_{i,i'})}{\sum_{i'} (E_r(m) \hat{s}_{ir} \times T_{i,i'})} \quad (\text{B.3})$$

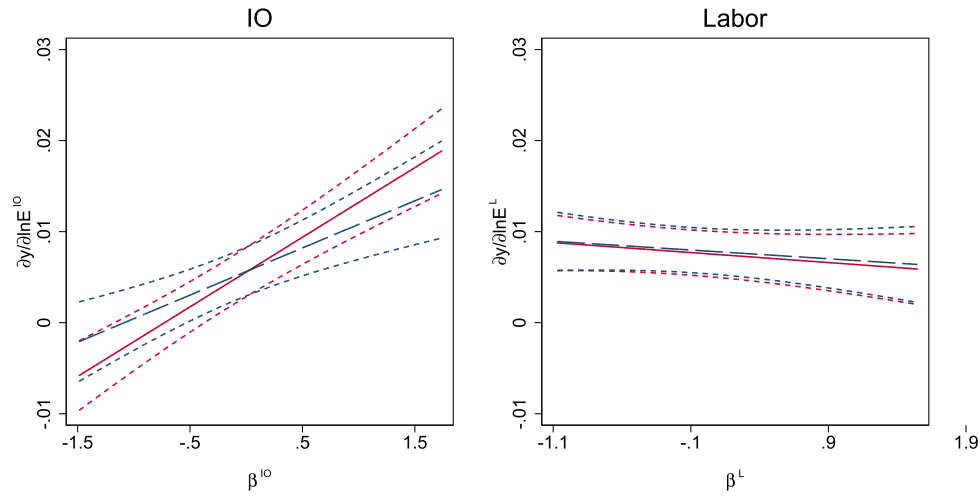
where  $E_r(m)$  is the manufacturing employment recorded in the CBP and  $T_{i,i'}$  is a correspondence matrix that converts employment from SIC industry  $i$  to NAICS industry  $i'$ . That is, we calculate the employment in region  $r$  and industry  $i$  that EGK would have predicted in our data by computing  $E_r(m) \hat{s}_{ir}$ . This predicted employment,  $\hat{E}_{ir}$ , is still recorded in the SIC classification. We, thus, convert it to NAICS using the formula  $\sum_i (\hat{E}_{ir} \times T_{i,i'})$ . Next, we express this updated employment as a share of an industry's total US employment. This share,  $\tilde{s}_{ir}$ , is then used in Eqs. (1) and (2) to compute the predicted EG and LC coagglomeration indices which feature as measures of similarity in natural advantages.

Since this procedure relies on the original estimates in EGK, services are excluded. To arrive at natural-advantage estimates for services, we would need the price information used in EGK, which we do not have. Instead, we follow a procedure that is close in spirit to the calculations in EGK, albeit with some differences. First, we use a different list of production factors. In particular, we use a combination of extractive industries, whose location is pinned down by

<sup>33</sup> The concordance was constructed by hand and can be found in the code provided in the supplementary materials.

<sup>34</sup> Ellison and Glaeser (1999) use population, but we find total employment to be more appropriate.





**Fig. C.1.** Marginal effects of related employment – Extensive margin.

Marginal effects of related employment with 95% confidence intervals, computed as the partial derivative of entry probability from Linear Probability Models (LPMs) with respect to related employment. The left panel depicts IO effects, the right panel labor-pooling effects. The horizontal axis plots  $\beta_{IO}$  and  $\beta_L$ , respectively, on a range that reflects the distribution of coefficient estimates obtained in Section 3.3. Blue dashed lines were constructed with estimates from agglomeration regressions using the EG-index, red solid lines using the LC-index.

**Table C.1**  
Growth in local industries.

	City-industry appearances								
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
$\ln E_{ir03}^{IO}$	0.0054 (0.0014)		0.0056 (0.0014)		0.0061 (0.0014)		0.0039 (0.0014)	0.0041 (0.0014)	0.0047 (0.0014)
$\ln E_{ir03}^L$		0.0080 (0.0013)		0.0079 (0.0013)		0.0078 (0.0013)	0.0072 (0.0013)	0.0072 (0.0013)	0.0069 (0.0013)
$\ln E_{ir03}^{IO} \hat{\rho}_{i,EG}^{IO}$			0.0053 (0.0013)					0.0053 (0.0013)	
$\ln E_{ir03}^L \hat{\rho}_{i,EG}^L$				-0.0010 (0.0011)				-0.0004 (0.0011)	
$\ln E_{ir03}^{IO} \hat{\rho}_{i,LC}^{IO}$					0.0078 (0.0011)				0.0077 (0.0011)
$\ln E_{ir03}^L \hat{\rho}_{i,LC}^L$						-0.0011 (0.0009)			-0.0009 (0.0009)
Obs.	43,082	43,082	43,082	43,082	43,082	43,082	43,082	43,082	43,082
Adj. $R^2$	0.1162	0.1166	0.1166	0.1166	0.1175	0.1166	0.1167	0.1171	0.1180

Idem Table 4, but the dependent variable is an indicator variable that takes value one if a local industry came into existence between 2003 and 2008. All independent variables are log-transformed.

a natural resource,<sup>35</sup> and unskilled labor categories.<sup>36</sup> Second, we change the way we measure  $y_{kr}$  and  $z_{ki}$ . For  $y_{kr}$ , which proxies the availability of a “natural” resource in a location, we either use the employment share of extractive industry  $k$  in region  $r$ , or, to proxy access to unskilled labor, the employment share of occupation  $k$  in region  $r$ . To proxy the intensity with which industry  $i$  relies on natural advantage  $k$ ,  $z_{ki}$ , we use  $i$ 's input share sourced from extractive industry  $k$  or the share of occupation  $k$  in the industry's overall employment.

### Appendix C. Growth at the extensive margin

See Table C.1 and Fig C.1.

<sup>35</sup> These five NAICS industries are 1141 (Fishing), 2111 (Oil and Gas Extraction), 2121 (Coal Mining), 2122 (Metal Ore Mining), 2123 (Nonmetallic Mineral Mining and Quarrying).

<sup>36</sup> We select six unskilled occupations (which we aggregate into one ‘unskilled’ factor): 35 (Food Preparation and Serving Related Occupations), 37 (Building and Grounds Cleaning and Maintenance Occupations), 47 (Construction and Extraction Occupations), 49 (Installation, Maintenance, and Repair Occupations), 51 (Production Occupations), 53 (Transportation and Material Moving Occupations).

### Appendix D. Industry lists

#### Industries used from IPUMS classifications

IND1990={40 41 42 50 60 100 101 102 110 111 112 120 121 122 130 132 140 141 142 150 151 152 160 161 162 171 172 180 181 182 190 191 192 200 201 210 211 212 220 221 222 230 231 232 241 242 250 251 252 261 262 270 271 272 280 281 282 290 291 292 300 301 310 311 312 320 321 322 331 332 340 341 342 350 351 352 360 361 362 370 371 372 380 381 390 391 392 700 701 702 710 711 712 721 722 731 732 740 741 762 770 800 802 810 841 850 872 882 890 891 892 893}.

IND1950={206 216 226 236 239 306 307 308 309 316 317 318 319 326 336 337 338 346 347 348 356 357 358 367 376 377 378 379 386 387 388 399 406 407 408 409 416 417 418 419 426 429 436 437 438 439 446 448 449 456 457 458 459 466 467 468 469 476 477 478 487 488 489 499 716 726 736 746 756 806 807 808 856 857 858 859 879}.

**Table D.1**  
Industry list with IO and labor coefficients (EG).

	Name	Group	Data avail.	Traded	IO coef.	Labor coef.
1	Crop production	.	0	0	.	.
2	Animal production	.	0	0	.	.
3	Aviculture	.	0	0	.	.
4	Mixed agriculture	.	0	0	.	.
5	Silviculture	.	0	0	.	.
6	Forestry	.	0	0	.	.
7	Logging	.	0	0	.	.
8	Fishing	.	0	0	.	.
9	Hunting and trapping	.	0	0	.	.
10	Support activities for crop production	.	0	0	.	.
11	Support activities for animal production	.	0	0	.	.
12	Support activities for forestry	.	0	0	.	.
13	Oil and gas extraction	1	1	1	0.27	−0.26
14	Coal mining	1	1	1	0.06	0.97
15	Metal ore mining	1	1	1	−0.02	0.82
16	Nonmetallic mineral mining and quarrying	1	1	1	−0.04	0.47
17	Support activities for mining	1	1	1	0.36	0.95
18	Utilities	.	1	0	.	.
19	Utility system construction	.	1	0	.	.
20	Land subdivision	.	1	0	.	.
21	Highway, street, and bridge construction	.	1	0	.	.
22	Heavy and civil engineering construction	.	1	0	.	.
23	Building exterior contractors	.	1	0	.	.
24	Building equipment contractors	.	1	0	.	.
25	Building finishing contractors	.	1	0	.	.
26	Other specialty trade contractors	.	1	0	.	.
27	Animal food manufacturing	2	1	1	0.01	0.24
28	Grain and oilseed milling	2	1	1	−0.09	0.35
29	Sugar product manufacturing	2	1	1	0.01	0.19
30	Fruit and vegetable preserving	2	1	1	−0.00	0.23
31	Dairy product manufacturing	2	1	1	−0.16	0.16
32	Animal slaughtering and processing	2	1	1	0.02	0.54
33	Seafood preparation and packaging	2	1	1	0.14	−0.06
34	Bakeries and tortilla manufacturing	2	1	1	0.08	−0.10
35	Other food manufacturing	2	1	1	0.01	0.07
36	Beverage manufacturing	2	1	1	0.02	0.13
37	Tobacco manufacturing	3	1	1	0.26	0.35
38	Fiber, yarn, and thread mills	4	1	1	−0.11	2.10
39	Textile and fabric mills	4	1	1	0.28	−0.18
40	Textile furnishings mills	4	1	1	0.16	0.21
41	Other textile product mills	4	1	1	−0.01	0.12
42	Apparel knitting mills	4	1	1	−0.04	0.69
43	Cut and sew apparel manufacturing	4	1	1	0.31	0.21
44	Other apparel manufacturing	4	1	1	0.37	−0.09
45	Leather and hide tanning	4	1	1	0.10	0.38
46	Footwear manufacturing	4	1	1	0.13	0.31
47	Other leather manufacturing	4	1	1	0.00	0.11
48	Sawmills and wood preservation	5	1	1	0.05	0.64
49	Wood product manufacturing	5	1	1	0.11	0.31
50	Other wood product manufacturing	5	1	1	0.10	0.22
51	Pulp, paper, and paperboard mills	6	1	1	−0.13	0.74
52	Paper product manufacturing	6	1	1	−0.05	0.28
53	Printing and related support activities	7	1	1	0.02	0.24
54	Petroleum and coal manufacturing	8	1	1	0.14	0.31
55	Basic chemical manufacturing	9	1	1	0.04	0.36
56	Synthetic fibers manufacturing	9	1	1	0.04	0.43
57	Agricultural chemical manufacturing	9	1	1	0.01	0.35
58	Pharmaceutical and medicine manufacturing	10	1	1	0.30	−0.33
59	Paint, coating, adhesive manufacturing	9	1	1	0.09	0.09
60	Soap, cleaning compound manufacturing	9	1	1	0.17	−0.27
61	Other chemical product manufacturing	9	1	1	0.03	0.16
62	Plastics product manufacturing	11	1	1	−0.02	0.20
63	Rubber product manufacturing	11	1	1	0.09	0.32
64	Clay product and refractory manufacturing	11	1	1	0.11	0.40
65	Glass and glass product manufacturing	11	1	1	0.06	0.26
66	Lime and gypsum product manufacturing	11	1	1	0.07	0.11
67	Other mineral product manufacturing	12	1	1	−0.02	0.46
68	Iron and steel mills	12	1	1	0.13	0.81
69	Steel product manufacturing	12	1	1	0.18	0.35
70	Alumina and aluminum production	12	1	1	0.04	0.36
71	Nonferrous metal production	12	1	1	0.07	0.33
72	Foundries	12	1	1	0.04	0.55

(continued on next page)

Table D.1 (continued)

	Name	Group	Data avail.	Traded	IO coef.	Labor coef.
73	Forging and stamping	12	1	1	0.19	0.32
74	Cutlery and handtool manufacturing	13	1	1	0.22	0.17
75	Structural metals manufacturing	13	1	1	0.02	0.10
76	Boiler, tank, container manufacturing	13	1	1	0.09	0.28
77	Hardware manufacturing	13	1	1	0.21	0.19
78	Spring and wire product manufacturing	13	1	1	0.19	0.21
79	Screw, nut, and bolt manufacturing	13	1	1	0.20	0.21
80	Coating, engraving, heat treating	13	1	1	0.28	0.22
81	Other metal product manufacturing	13	1	1	0.17	0.10
82	Agri/construction/mining machinery	14	1	1	0.46	0.18
83	Industrial machinery manufacturing	14	1	1	0.09	0.17
84	Service industry machinery manufacturing	14	1	1	0.18	0.08
85	Heating/cooling equipment	14	1	1	0.21	0.14
86	Metalworking machinery manufacturing	14	1	1	0.31	0.46
87	Engine, turbine, transmission manufacturing	14	1	1	0.16	0.46
88	Other machinery manufacturing	14	1	1	0.28	0.15
89	Computer equipment manufacturing	15	1	1	0.15	0.06
90	Audio-video equipment manufacturing	15	1	1	−0.04	−0.10
91	Semiconductor manufacturing	15	1	1	0.04	0.13
92	Communications equipment manufacturing	15	1	1	0.23	0.01
93	Manufacturing of magnetic and optical media	15	1	1	0.07	0.00
94	Electric lighting equipment manufacturing	16	1	1	0.10	0.01
95	Household appliance manufacturing	16	1	1	0.13	0.23
96	Electrical Equipment Manufacturing	16	1	1	0.20	0.15
97	Other electrical equipment manufacturing	17	1	1	0.04	0.12
98	Motor vehicle manufacturing	17	1	1	0.11	0.13
99	Motor vehicle body manufacturing	17	1	1	0.00	0.25
100	Motor vehicle parts manufacturing	17	1	1	0.17	0.22
101	Aerospace product manufacturing	17	1	1	0.24	−0.08
102	Railroad rolling stock manufacturing	17	1	1	0.24	0.47
103	Ship and boat building	17	1	1	0.01	0.27
104	Other transportation equipment	17	1	1	0.12	0.17
105	Household furniture manufacturing	18	1	1	−0.02	0.29
106	Office furniture manufacturing	18	1	1	0.05	0.10
107	Other furniture related manufacturing	18	1	1	−0.01	−0.04
108	Medical supplies manufacturing	10	1	1	0.13	0.02
109	Other miscellaneous manufacturing	99	1	1	0.09	−0.04
110	Wholesale trade	.	1	0	.	.
111	Retail trade	.	0	0	.	.
112	Scheduled air transportation	.	1	0	.	.
113	Nonscheduled air transportation	.	1	0	.	.
114	Rail transportation	.	0	0	.	.
115	Water transportation	.	1	0	.	.
116	Inland water transportation	.	1	0	.	.
117	General freight trucking	.	1	0	.	.
118	Specialized freight trucking	.	0	0	.	.
119	Urban transit systems	.	1	0	.	.
120	Interurban bus transportation	.	1	0	.	.
121	School and employee bus	.	1	0	.	.
122	Charter bus industry	.	1	0	.	.
123	Taxi and limousine service	.	1	0	.	.
124	Pipeline transportation of crude oil	.	0	0	.	.
125	Pipeline transportation of gas	.	1	0	.	.
126	Other pipeline transportation	.	1	0	.	.
127	Sightseeing transportation, land	.	1	0	.	.
128	Sightseeing transportation, water	.	1	0	.	.
129	Sightseeing transportation, other	.	1	0	.	.
130	Support for air transportation	.	1	0	.	.
131	Support for rail transportation	.	1	0	.	.
132	Support for water transportation	.	1	0	.	.
133	Support for road transportation	.	1	0	.	.
134	Freight transportation arrangement	.	1	0	.	.
135	Support activities for transportation	.	1	0	.	.
136	Postal services	.	0	0	.	.
137	Couriers and express delivery services	.	1	0	.	.
138	Local messengers and local delivery	.	0	0	.	.
139	Warehousing and storage	.	1	0	.	.
140	Publishers	19	1	1	0.18	0.81
141	Software publishers	21	1	1	0.17	0.61
142	Sound recording industries	25	1	1	1.91	1.71
143	Radio and television broadcasting	19	1	1	0.01	1.16
144	Satellite telecommunications	20	1	1	0.11	0.47
145	Data processing services	21	1	1	0.04	0.42

(continued on next page)

Table D.1 (continued)

	Name	Group	Data avail.	Traded	IO coef.	Labor coef.
146	Other telecommunications	20	1	1	0.05	0.68
147	Real estate and construction	.	1	0	.	.
148	Offices of real estate agents	.	1	0	.	.
149	Activities related to real estate	.	1	0	.	.
150	Automotive equipment rental	.	1	0	.	.
151	Consumer goods rental	.	1	0	.	.
152	General rental centers	.	0	0	.	.
153	Machinery and equipment rental	.	0	0	.	.
154	Lessors of nonfinancial intangible assets	23	1	1	−0.20	1.11
155	Legal services	23	1	1	0.13	2.58
156	Accounting services	23	1	1	0.44	0.68
157	Architectural and engineering services	22	1	1	0.09	0.59
158	Specialized design services	22	1	1	0.05	2.25
159	Computer systems design	.	0	0	.	.
160	Scientific and R&D services	.	0	0	.	.
161	Advertising and related services	.	0	0	.	.
162	Professional and scientific services	.	0	0	.	.
163	Management of companies and enterprises	23	1	1	0.07	0.46
164	Office administrative services	23	1	1	−0.04	0.54
165	Facilities support services	23	1	1	0.07	0.71
166	Management consulting services	23	1	1	0.12	0.13
167	Business support services	24	1	1	0.04	0.21
168	Travel arrangement services	.	1	0	.	.
169	Investigation and security services	.	1	0	.	.
170	Services to buildings and dwellings	.	0	0	.	.
171	Other support services	.	0	0	.	.
172	Waste treatment and disposal	.	1	0	.	.
173	Elementary and secondary schools	.	1	0	.	.
174	Junior colleges	24	1	1	0.32	0.57
175	Colleges, universities, professional schools	24	1	1	0.07	0.38
176	Business schools and computer training	24	1	1	0.27	0.39
177	Technical and trade schools	24	1	1	0.13	0.10
178	Other schools and instruction	.	1	0	.	.
179	Educational support services	.	1	0	.	.
180	Offices of physicians	.	1	0	.	.
181	Offices of dentists	.	1	0	.	.
182	Offices of other health practitioners	.	0	0	.	.
183	Medical and diagnostic laboratories	.	0	0	.	.
184	Home health care services	.	0	0	.	.
185	Other ambulatory health care services	.	1	0	.	.
186	General medical and surgical hospitals	.	1	0	.	.
187	Psychiatric and substance abuse hospitals	.	1	0	.	.
188	Specialty hospitals	.	1	0	.	.
189	Nursing care facilities	.	1	0	.	.
190	Mental health and substance abuse facilities	.	1	0	.	.
191	Individual and family services	.	1	0	.	.
192	Community food and housing	.	1	0	.	.
193	Vocational rehabilitation services	.	1	0	.	.
194	Child day care services	.	1	0	.	.
195	Performing arts companies	25	1	1	−0.01	3.41
196	Spectator sports	25	1	1	0.03	0.17
197	Promoters of performing arts	25	1	1	0.23	0.42
198	Agents and managers	25	1	1	0.85	2.96
199	Independent artists	25	1	1	0.26	2.19
200	Museums and historical sites	25	1	1	0.18	0.44
201	Amusement parks and recreation industry	26	1	1	0.16	0.15
202	Traveler accommodation	26	1	1	−0.00	0.13
203	RV parks and recreational camps	26	1	1	−0.24	−0.17
204	Residential care facilities	.	1	0	.	.
205	Restaurants	.	1	0	.	.
206	Special food services	.	1	0	.	.
207	Drinking places (alcoholic beverages)	.	0	0	.	.
208	Automotive repair and maintenance	.	1	0	.	.
209	Machinery repair and maintenance	27	1	1	0.05	0.31
210	Household goods repair and maintenance	.	1	0	.	.
211	Other personal services	.	1	0	.	.
212	Associations and organizations	.	1	0	.	.
213	Household services	.	0	0	.	.
214	Other public services	.	1	0	.	.
215	Finance and insurance	23	1	1	0.07	0.53



**Table D.2**

Average point-estimates by subsector.

Group	Name	IO coef.	labor coef.
1	Extraction	0.13	0.59
2	Food manufacturing	0.00	0.18
3	Tobacco manufacturing	0.26	0.35
4	Textile manufacturing	0.12	0.39
5	Wood products	0.09	0.39
6	Paper products	−0.09	0.51
7	Printing	0.02	0.24
8	Petroleum and coal manufacturing	0.14	0.31
9	Chemical manufacturing	0.06	0.19
10	Pharmaceutical and medical supply	0.21	−0.15
11	Materials	0.06	0.26
12	Mineral products manufacturing	0.09	0.45
13	Hardware manufacturing	0.17	0.18
14	Machinery	0.24	0.23
15	Electronics	0.09	0.02
16	Electrical equipment manufacturing	0.14	0.13
17	Transportation equipment	0.11	0.19
18	Furniture	0.01	0.12
19	Media	0.09	0.99
20	Telecommunication	0.08	0.57
21	IT services	0.10	0.51
22	Architecture and engineering	0.07	1.42
23	Professional KIBS	0.08	0.84
24	Educational	0.16	0.33
25	Art and culture	0.49	1.61
26	Recreation	−0.03	0.04
27	Machinery repair	0.05	0.31
99	Other	0.09	−0.04

**Table E.2**

OLS and IV univariate regressions, extended sample.

	(1) OLS State	(2) OLS City	(3) OLS County	(4) IV State	(5) IV City	(6) IV County
LC index - All traded						
Input-output	0.132 (0.015)	0.119 (0.015)	0.122 (0.015)	0.162 (0.027)	0.148 (0.028)	0.139 (0.027)
Observations	7140	7140	7140	7140	7140	7140
R <sup>2</sup>	0.024	0.021	0.026	0.023	0.020	0.026
Labor	0.120 (0.008)	0.033 (0.008)	0.044 (0.008)	0.061 (0.017)	−0.025 (0.016)	0.005 (0.015)
Observations	7140	7140	7140	7140	7140	7140
R <sup>2</sup>	0.024	0.002	0.004	0.018	−0.004	0.001
Technology	0.087 (0.023)	0.046 (0.018)	0.061 (0.020)	0.085 (0.022)	0.041 (0.016)	0.056 (0.018)
Observations	7140	7140	7140	7140	7140	7140
R <sup>2</sup>	0.010	0.003	0.006	0.010	0.003	0.006
Nat. advantages	0.677 (0.039)	0.487 (0.033)	0.477 (0.031)	2.414 (0.207)	2.380 (0.219)	3.058 (0.237)
Observations	7140	7140	7140	3828	3828	3828
R <sup>2</sup>	0.077	0.043	0.048	−0.575	−0.756	−1.342

Robust standard errors in parentheses.

**Table E.3**

OLS and IV multivariate regressions, extended sample, controlling for all four channels.

	(1) OLS State	(2) OLS City	(3) OLS County	(4) IV State	(5) IV City	(6) IV County
LC index - All traded						
Input-output	0.103 (0.014)	0.115 (0.015)	0.115 (0.014)	0.167 (0.033)	0.201 (0.037)	0.171 (0.034)
Labor	0.042 (0.009)	−0.034 (0.009)	−0.023 (0.009)	−0.059 (0.024)	−0.151 (0.024)	−0.107 (0.023)
Technology	0.046 (0.013)	0.019 (0.013)	0.034 (0.014)	0.049 (0.015)	0.019 (0.015)	0.033 (0.015)
Nat. advantages	0.614 (0.040)	0.490 (0.034)	0.468 (0.032)	0.684 (0.043)	0.570 (0.039)	0.527 (0.036)
Observations	7140	7140	7140	7140	7140	7140
R <sup>2</sup>	0.103	0.063	0.072	0.085	0.036	0.057

Robust standard errors in parentheses.

**Appendix E. Analysis using location correlation**

Table E.1 and shows the top 10 most coagglomerated industry pairs. When calculating coagglomeration scores, we exclude the ten largest cities to avoid that this top 10 gets cluttered by patterns generated by industries that locate in extremely large cities. Although all our estimations use coagglomeration scores calculated with the full set of cities, reported results (available on request) are robust to excluding the largest cities.

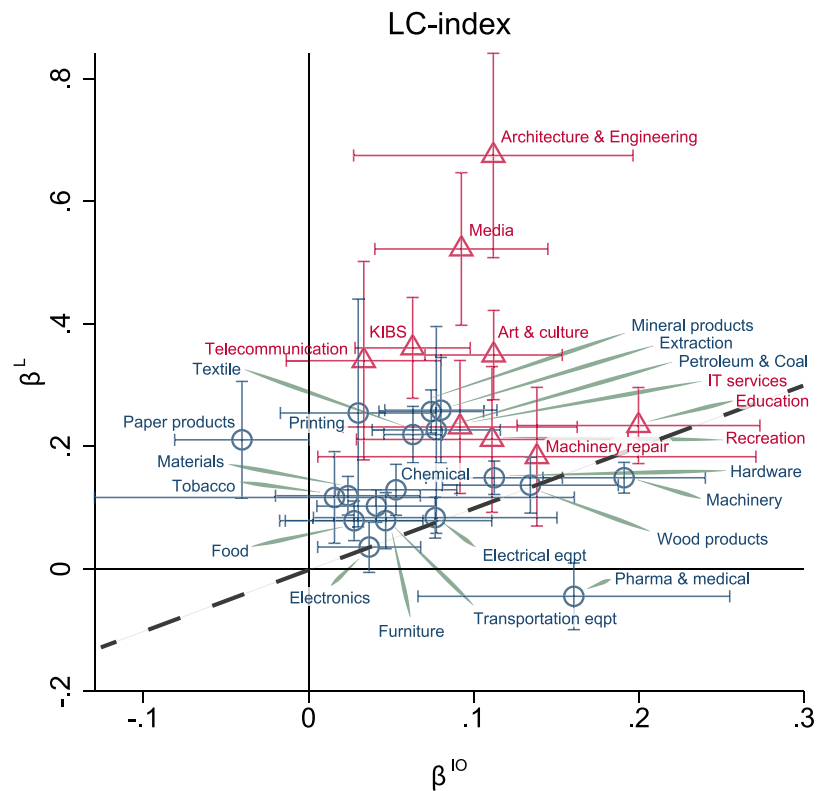
**Table E.1**

Top-10 industry pairs by coagglomeration (Locational Correlation index).

Top-10 co-location (LC index)		
Industry i	Industry j	Value
Legal services	Accounting services	0.9712
Management consulting services	Accounting services	0.9684
Legal services	Finance and insurance	0.9681
Screw, nut, and bolt manufacturing	Coating, engraving, heat treating	0.9636
Specialized design services	Accounting services	0.9614
Specialized design services	Legal services	0.9611
Publishers	Legal services	0.9601
Management consulting services	Architectural and engineering services	0.9565
Finance and insurance	Accounting services	0.9565
Architectural and engineering services	Accounting services	0.9552

Idem Table A.1, but here, we show the top-10 industry pairs using the LC-based coagglomeration metric (see Eq. (2)) using city-industry employment data for the US (County Business Patterns, 2003). The top-10 cities have been excluded from the computation to reduce noise.

The correlation coefficient between the vectors of estimated labor-similarity coefficients using the LC index, respectively EG index, is 0.85. For the estimated input-output proximity effects, the correlation is 0.44.

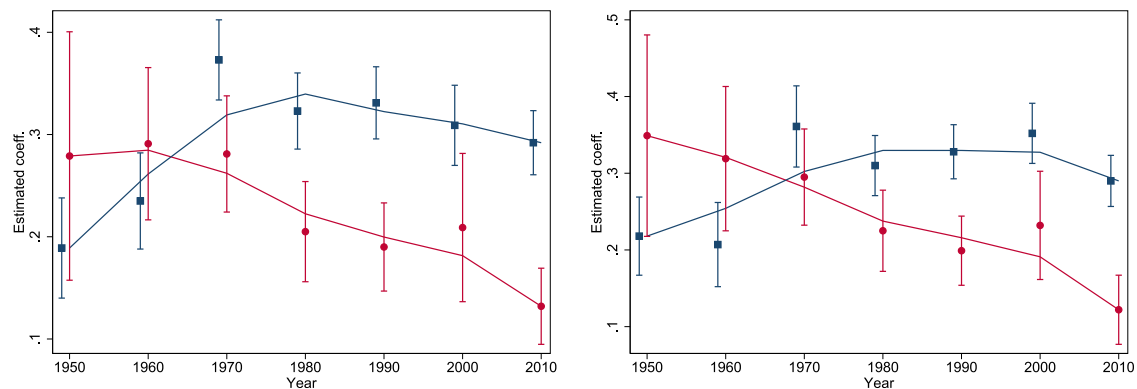


**Fig. E.1.** Coagglomeration effects for 27 subsectors (LC-index).

Replication of Fig. 3 using the LC index. The figure depicts labor pooling effects (vertical axis) and value chain effects (horizontal axis). Estimates for coagglomeration patterns of services are marked with red triangles, for manufacturing with blue circles.

## Appendix F. Additional robustness analysis

See Tables F.1, F.2, F.3, F.4, F.5 and Fig. F.1.



**Fig. F.1.** Trends in Marshallian agglomeration forces for manufacturing - OLS (left) and IV (right).

Estimated impact of input-output (red circles) and labor (blue squares) linkages on EG-coagglomeration in the period 1950–2010 for manufacturing, with 95% confidence intervals based on robust standard errors. The lines represent LOWESS smooths.

**Table F.1**  
OLS and IV univariate regressions by sector (EG-index).

	EG index					
	(1) OLS State	(2) OLS City	(3) OLS County	(4) IV State	(5) IV City	(6) IV County
All industries						
Input-output	0.174 (0.023)	0.138 (0.015)	0.112 (0.015)	0.237 (0.027)	0.172 (0.020)	0.132 (0.019)
Observations	16836	16836	16836	16836	16836	16836
R <sup>2</sup>	0.033	0.022	0.015	0.028	0.020	0.015
Labor	0.201 (0.009)	0.175 (0.007)	0.129 (0.007)	0.322 (0.023)	0.282 (0.016)	0.191 (0.015)
Observations	16836	16836	16836	16836	16836	16836
R <sup>2</sup>	0.042	0.034	0.019	0.027	0.021	0.015
Manufacturing						
Input-output	0.239 (0.028)	0.161 (0.018)	0.115 (0.016)	0.293 (0.036)	0.199 (0.024)	0.135 (0.019)
Observations	11786	11786	11786	11786	11786	11786
R <sup>2</sup>	0.048	0.025	0.017	0.045	0.024	0.017
Labor	0.229 (0.010)	0.191 (0.007)	0.121 (0.006)	0.302 (0.025)	0.254 (0.015)	0.137 (0.011)
Observations	11786	11786	11786	11786	11786	11786
R <sup>2</sup>	0.061	0.049	0.027	0.055	0.044	0.026
Services						
Input-output	0.127 (0.019)	0.171 (0.028)	0.175 (0.041)	0.209 (0.038)	0.290 (0.053)	0.272 (0.078)
Observations	5360	5360	5360	5360	5360	5360
R <sup>2</sup>	0.015	0.016	0.013	0.009	0.008	0.009
Labor	0.213 (0.024)	0.275 (0.037)	0.294 (0.053)	0.421 (0.057)	0.507 (0.075)	0.515 (0.100)
Observations	5360	5360	5360	5360	5360	5360
R <sup>2</sup>	0.019	0.018	0.016	0.001	0.005	0.007

This table contains robustness checks, using all 184 industries as the “destination” industry. Robust standard errors in parentheses.

**Table F.2**  
OLS and IV univariate regressions, extended sample.

	(1) OLS State	(2) OLS City	(3) OLS County	(4) IV State	(5) IV City	(6) IV County
EG index - All traded						
Input-output	0.295 (0.041)	0.215 (0.025)	0.176 (0.028)	0.372 (0.045)	0.267 (0.032)	0.201 (0.032)
Observations	7140	7140	7140	7140	7140	7140
R <sup>2</sup>	0.069	0.038	0.025	0.064	0.036	0.024
Labor	0.261 (0.012)	0.229 (0.011)	0.178 (0.011)	0.347 (0.031)	0.296 (0.020)	0.194 (0.018)
Observations	7140	7140	7140	7140	7140	7140
R <sup>2</sup>	0.065	0.052	0.030	0.058	0.048	0.030
Technology	0.077 (0.029)	0.041 (0.021)	0.030 (0.015)	0.083 (0.030)	0.046 (0.022)	0.031 (0.015)
Observations	7140	7140	7140	7140	7140	7140
R <sup>2</sup>	0.005	0.001	0.001	0.005	0.001	0.001
Nat. advantages	0.348 (0.017)	0.398 (0.022)	0.328 (0.030)	1.053 (0.153)	0.833 (0.120)	0.584 (0.082)
Observations	7140	7140	7140	3828	3828	3828
R <sup>2</sup>	0.069	0.094	0.061	−0.030	0.006	0.003

Robust standard errors in parentheses.

**Table F.3**  
OLS and IV multivariate regressions, manufacturing.

	(1) OLS State	(2) OLS City	(3) OLS County	(4) IV State	(5) IV City	(6) IV County
EG index - Manufacturing						
Input-output	0.239 (0.039)	0.142 (0.023)	0.098 (0.021)	0.315 (0.056)	0.171 (0.034)	0.135 (0.028)
Labor	0.210 (0.011)	0.184 (0.009)	0.141 (0.008)	0.188 (0.032)	0.194 (0.020)	0.103 (0.015)
Observations	6474	6474	6474	6474	6474	6474
R <sup>2</sup>	0.110	0.073	0.053	0.106	0.072	0.050

Robust standard errors in parentheses.

**Table F.4**  
OLS and IV multivariate regressions, services.

	(1) OLS State	(2) OLS City	(3) OLS County	(4) IV State	(5) IV City	(6) IV County
EG index - Services						
Input-output	0.128 (0.031)	0.157 (0.044)	0.185 (0.069)	−0.050 (0.186)	0.296 (0.258)	0.221 (0.420)
Labor	0.252 (0.034)	0.355 (0.052)	0.315 (0.074)	0.554 (0.248)	0.323 (0.338)	0.388 (0.526)
Observations	3312	3312	3312	3312	3312	3312
R <sup>2</sup>	0.043	0.042	0.030	0.009	0.036	0.029

Robust standard errors in parentheses.

**Table F.5**  
Historical analysis of labor and value-chain linkages in coagglomeration patterns.

Historical regressions										
	[1] standard 1950–2010		[2] 2-way FE 1950–2010		[3] IO fixed year 1950–2010		[4] standard 1910–2010		[5] 2-way FE 1910–2010	
	OLS	IV	OLS	IV	OLS	IV	OLS	IV	OLS	IV
Labor	0.188 (0.025)	0.193 (0.026)	0.193 (0.032)	0.225 (0.034)	0.188 (0.031)	0.183 (0.046)	0.115 (0.016)	0.108 (0.025)	0.163 (0.021)	0.171 (0.049)
Labor $\times$ year	0.100 (0.034)	0.110 (0.036)	0.178 (0.040)	0.193 (0.041)	0.101 (0.041)	0.101 (0.057)	−0.001 (0.023)	−0.003 (0.038)	0.005 (0.023)	−0.004 (0.036)
Input-output	0.315 (0.031)	0.364 (0.033)	0.362 (0.068)	0.415 (0.073)	0.332 (0.063)	0.424 (0.122)	0.209 (0.027)	0.274 (0.036)	0.189 (0.028)	0.228 (0.039)
Input-output $\times$ year	−0.194 (0.041)	−0.254 (0.044)	−0.272 (0.095)	−0.331 (0.101)	−0.217 (0.084)	−0.268 (0.156)	−0.097 (0.047)	−0.117 (0.065)	−0.091 (0.050)	−0.121 (0.071)
$R^2$	0.067	0.067	0.185	0.185	0.067	0.068	0.029	0.030	0.119	0.122
N	31757	31757	31757	31757	31757	26401	30356	29403	30356	29403

Year is divided by 100. Time dummies for each decade are included, but not reported in the table. Robust standard errors in parentheses. The first block (column 1 and 2) replicates columns 3 and 6 of Table 5. The second block controls for (2-way) industry fixed effects. In the third block, the input-output table used to construct value-chain proximity does not change over time but is fixed to the 1950 table (IV: 2010 table). Block 4 and 5 extend the estimates to the period 1910–2010.



## Appendix G. Correction to CBP censoring

To avoid disclosure of sensitive information, the Census Bureau

**Table G.1**

Mean employment by size category in 2000 CBP, as in Holmes and Stevens (2004).

Employment range	Average employment
1–4	1.7
5–9	6.6
10–19	13.5
20–49	30.2
50–99	68.8
100–249	150.1
250–499	340.7
500–999	681.3
1000–1499	1208.8
1500–2499	1892.9
2500–4999	3374.7
5000 or more	9592.0

**Table G.2**

Mean employment by size class in 2003 and 2008 CBP, as used in this paper.

Class	Employment range	Average employment
A	0–19	5
B	20–99	40
C	100–249	150
E	250–499	341
F	500–999	681
G	1000–2499	1488
H	2500–4999	3375
I	5000–9999	5980
J	10,000–24,999	13,955
K	25,000–49,999	29,904
L	50,000–99,999	59,804
M	100,000 or More	125,000

withholds precise data in the CBP, whenever the information in a cell (by geography and by industry) reveals information about a firm. Instead, in such cases, the CBP report employment in terms of size categories. To translate these categories into single employment values, we follow Holmes and Stevens (2004), who provide the following estimate<sup>37</sup> of the mean employment by CBP class in the year 2000.

The CBP categories have changed slightly since Holmes and Stevens (2004). We adapt their estimates to the new categories as follows: categories C, E, F and H directly follow the estimates of Holmes and Stevens. Categories A, B and G use weighted averages from the subcategories in Holmes and Stevens (2004). Categories I, J, K and L use the average ratio between the estimated employment of classes A to G and the mid-point of employment classes A to G (0.7974). This ratio is then multiplied by the midpoint of I, J, K and L to obtain an estimate of the average employment. For category M, we use an average employment of 125,000, which is obtained by multiplying 100,000 by the average of the ratio between the minimum and estimated mean of classes I through L. The resulting scheme is shown below.

<sup>37</sup> Estimates for categories below 1000 employees are directly computed from national totals, dividing total employment by the number of establishments (within each category). For categories with over 1000 employees, the average employment in each category is imputed estimating the parameters of a log-normal distribution of employment across establishments. See Appendix A.2 in Holmes and Stevens (2004).

## Supplementary material

Supplementary material associated with this article can be found, in the online version, at doi:10.1016/j.jue.2018.05.002.

## References

- Abdel-Rahman, H.M., 1996. When do cities specialize in production? *Reg. Sci. Urban Econ.* 26 (1), 1–22.
- Autor, D., Dorn, D., 2013. The growth of low-skill service jobs and the polarization of the US labor market. *Am. Econ. Rev.* 103 (5), 1553–1597.
- Beaudry, C., Schiffauerova, A., 2009. Who's right, Marshall or Jacobs? The localization versus urbanization debate. *Res. Policy* 38 (2), 318–337.
- Behrens, K., 2016. Agglomeration and clusters: tools and insights from coagglomeration patterns. *Can. J. Econ. Revue canadienne d'économie* 49 (4), 1293–1339.
- Black, D., Henderson, V., 2003. Urban evolution in the USA. *J. Econ. Geogr.* 3 (4), 343–372.
- Cockburn, I., Griliches, Z., 1988. The estimation and measurement of spillover effects of R&D investment-industry effects and appropriability measures in the stock market's valuation of R&D and patents. *Am. Econ. Rev.* 78 (2), 419–423.
- Dauth, W., 2010. Agglomeration and Regional Employment Growth. IAB discussion paper 07/2010. Institut für Arbeitsmarkt- und Berufsforschung.
- Delgado, M., Porter, M., Stern, S., 2010. Clusters and entrepreneurship. *J. Econ. Geogr.* 10 (4), 495–518.
- Duranton, G., Puga, D., 2004. Micro-foundations of urban agglomeration economies. In: Henderson, J., Thisse, J. (Eds.), *Handbook of Regional and Urban Economics*, 4. Elsevier, pp. 2063–2117.
- Ellison, G., Glaeser, E.L., 1997. Geographic concentration in US manufacturing industries: a dartboard approach. *J. Polit. Econ.* 105 (5), 889–927.
- Ellison, G., Glaeser, E.L., 1999. The geographic concentration of industry: does natural advantage explain agglomeration? *Am. Econ. Rev.* 89 (2), 311–316.
- Ellison, G., Glaeser, E.L., Kerr, W.R., 2010. What causes industry agglomeration? evidence from coagglomeration patterns. *Am. Econ. Rev.* 100 (3), 1195–1213.
- Faggio, G., Silva, O., Strange, W.C., 2017. Heterogeneous agglomeration. *Rev. Econ. Stat.* 99 (1), 80–94.
- Glaeser, E., Kallal, H.D., Scheinkman, J., Shleifer, A., 1992. Growth in cities. *J. Polit. Econ.* 100 (6), 1126–1152.
- Glaeser, E.L., Kerr, W.R., 2009. Local industrial conditions and entrepreneurship: how much of the spatial distribution can we explain? *J. Econ. Manag. Strat.* 18 (3), 623–663.
- Goldin, C., Katz, L., 2009. *The Race Between Education and Technology*. Harvard University Press.
- Goldschlag, N., Lybbert, T.J., Zolas, N.J., 2016. An Algorithmic Links with Probabilities Crosswalk for USPC and CPC Patent Classifications with an Application Towards Industrial Technology Composition. CES Working Papers 16-15. Center for Economic Studies.
- Greenstone, M., Hornbeck, R., Moretti, E., 2008. Identifying Agglomeration Spillovers: Evidence from Million Dollar Plants. NBER Working Papers 13833. National Bureau of Economic Research.
- Groot, H.L., Poot, J., Smit, M.J., 2016. Which agglomeration externalities matter most and why? *J. Econ. Surv.* 30 (4), 756–782.
- Hall, B.H., Jaffe, A.B., Trajtenberg, M., 2001. The NBER Patent Citation Data File: Lessons, Insights and Methodological Tools. NBER Working Papers 8498. National Bureau of Economic Research.
- Hanlon, W.W., Miscio, A., 2017. Agglomeration: a long-run panel data approach. *J. Urban. Econ.* 99, 1–14.
- Hausmann, R., Hidalgo, C., Stock, D.P., Yildirim, M.A., 2014. Implied Comparative Advantage. CID Working Paper 276. Center for International Development.
- Helsley, R.W., Strange, W.C., 1990. Matching and agglomeration economies in a system of cities. *Reg. Sci. Urban Econ.* 20 (2), 189–212.
- Henderson, V., Kuncoro, A., Turner, M., 1995. Industrial development in cities. *J. Polit. Econ.* 103 (5), 1067–1090.
- Holmes, T.J., Stevens, J.J., 2004. Spatial distribution of economic activities in North America. In: Henderson, J., Thisse, J. (Eds.), *Handbook of Regional and Urban Economics*, 4. Elsevier, pp. 2797–2843. chapter 63.
- Howard, E., Newman, C., Tarp, F., 2015. Measuring industry coagglomeration and identifying the driving forces. *J. Econ. Geogr.* 16 (5), 1055–1078.
- Jofre-Monseny, J., Marín-López, R., Viladecans-Marsal, E., 2011. The mechanisms of agglomeration: evidence from the effect of inter-industry relations on the location of new firms. *J. Urban Econ.* 70 (2), 61–74.
- Kolko, J., 1999. Can I Get Some Service Here? Information Technology, Service Industries, and the Future of Cities. SSRN Working Paper. SSRN.
- Marshall, A., 1920. *Principles of Economics*, eighth ed. MacMillan, London. 1982 reprint.
- Morgan, K., Cooke, P., 1998. *The Associational Economy: Firms, Regions, and Innovation*. University of Illinois at Urbana-Champaign's Academy for Entrepreneurial Leadership Historical Research Reference in Entrepreneurship.
- Neffke, F., Henning, M., Boschma, R., 2011. How do regions diversify over time? industry relatedness and the development of new growth paths in regions. *Econ. Geogr.* 87 (3), 237–265.
- O'Sullivan, A., Strange, W.C., 2017. The emergence of coagglomeration. *J. Econ. Geogr.* 18 (2), 293–317.
- Porter, M., 1998. Cluster and the new economics of competition. *Harv. Bus. Rev.* 77–90.

- Porter, M., 2003. The economic performance of regions. *Reg. Stud.* 37 (6–7), 545–546.
- Richardson, G.B., 1972. The organisation of industry. *Econ. J.* 82 (327), 883–896.
- Rosenthal, S.S., Strange, W.C., 2001. The determinants of agglomeration. *J. Urban Econ.* 50 (2), 191–229.
- Ruggles, S., Katie Genadek, R.G., Grover, J., Sobek, M., 2017. Integrated Public Use Microdata Series: Version 7.0. Minneapolis: University of Minnesota. <https://doi.org/10.18128/D010.V7.0>.
- Scherer, F.M., 1984. Using linked patent and r&d data to measure interindustry technology flows. In: Griliches, Z. (Ed.), *R&D, Patents, and Productivity*. University of Chicago Press, pp. 417–461.
- Zucker, L.G., Darby, M.R., Brewer, M.B., 1994. Intellectual Capital and the Birth of US Biotechnology Enterprises. NBER Working Papers 4653. National Bureau of Economic Research.