

# Network embeddedness indicates the innovation potential of firms

**Giacomo Vaccario, Luca Verginer, Antonios  
Garas, Mario V. Tomasello and Frank Schweitzer**

Chair of Systems Design, ETH Zurich, Switzerland

[www.sg.ethz.ch](http://www.sg.ethz.ch)

## Abstract

Firms' innovation potential depends on their position in the R&D network. But details on this relation remain unclear because measures to quantify network embeddedness have been controversially discussed. We propose and validate a new measure, *coreness*, obtained from the weighted  $k$ -core decomposition of the R&D network. Using data on R&D alliances, we analyse the change of coreness for 14,000 firms over 25 years and patenting activity. A regression analysis demonstrates that coreness explains firms' R&D output by predicting future patenting.

## 1 Introduction

The ability of firms to innovate considerably depends on their research and development (R&D) collaborations (Schilling and Phelps, 2007). To formalize these collaborations, firms establish R&D alliances which allow them to coordinate their activities, share resources, and exchange knowledge (Nooteboom, 1999). Alliances are publicly announced. Therefore, we know that firms engage in various alliances at the same time and change them over time (Tomasello et al., 2016). Data on 21,500 alliances between 14,000 firms allows us to reconstruct an alliance network, in which nodes represent firms and links their R&D collaborations. This network evolves as new firms enter, incumbents exit, new alliances are formed, and established alliances end. These processes continuously change a firm's position in the R&D network and thus affect its ability to innovate.

This paper focuses on the relation between innovation output, as measured by the number of patents, and the importance of individual firms, as measured by their topological position in the alliance network. It has been established in the literature that the interfirm network to which firms belong is a driver of their innovation output (Freeman, 1991; Powell et al., 1999; Ahuja, 2000; Owen-Smith and Powell, 2004; Schilling and Phelps, 2007; Phelps, 2010). For example, Schilling and Phelps (2007) have argued that global topological properties of alliance networks (i.e., average path length and clustering coefficient) influence the number of patents. Similarly, Owen-Smith and Powell (2004) have shown that firms' innovation output depends on their ability to absorb information flows in these networks.

The relation between innovation ability and network properties points to the central question of how to quantify the position of firms in the R&D network. A common proxy is betweenness centrality (Owen-Smith and Powell, 2004; Gilsing et al., 2008; Paier and Scherngell, 2011). It measures how often a focal node is included in the shortest paths between all existing nodes in a network (Newman, 2010). This quantifies the importance of the focal node in controlling the information flow between other nodes and its access to information. Other centrality measures, e.g., degree, closeness or eigenvector centrality, capture the influence of different topological properties, such as the number of neighbours, the topological distance to all other nodes, or the impact of neighbouring nodes. Several of these measures have been used in various studies to determine the importance of firms and describe their *embeddedness* in the alliance network (Powell et al., 1999; Ahuja, 2000; Schilling and Phelps, 2007).

The concept of embeddedness is pervasive in economics and sociology (Polanyi and MacIver, 1944; Granovetter, 1985). Yet, it is only loosely defined. A central idea is that individuals or firms are *embedded* in social relations, which in turn affect their economic behaviour. Precisely, Granovetter (1985) stresses that “the level of embeddedness of economic behaviour [...] has always been and continues to be more substantial than is allowed for by formalists and economists”. Also, Gulati et al. (2000) write, “the image of atomistic actors competing for profits against each other in an impersonal marketplace is increasingly inadequate in a world in which firms are embedded in networks of social, professional, and exchange relationships with other organizational actors”. Our work builds on these arguments about the importance of embeddedness. We study how firms’ innovation output is affected by their embeddedness in R&D alliances, which we quantify in a novel method called “coreness”.

Importantly, formal agreements such as R&D alliances define not only economic but also social relations, as (Powell et al., 1999) and (Gulati et al., 2000) point out. They provide access to complementary capabilities and the opportunity for learning (Nooteboom, 1999). At the same time, establishing and maintaining social relations is costly (Granovetter, 1985; Uzzi, 1997). For this reason, it is debated whether R&D alliances have a net positive effect on firm innovation. Indeed, when looking at the biotechnology industry, Powell et al. (1999) found that more embedded firms file more patents. However, when considering the subsidiaries of multinational pharmaceutical companies, Al-Laham and Bort (2011) found the opposite effect, i.e., the more a firm is embedded, the fewer patents it files.

We contribute to the ongoing discussion in two ways. Firstly, for our analysis, we introduce a new measure of firms’ embeddedness in the alliance network based on the *weighted k-core* centrality (Garas et al., 2012) This is an extension of the unweighted version introduced by Seidman (1983); Bollobás (1988). The weighted *k-core* centrality has the advantage to control for repeated interactions, while the unweighted centrality treats repeated interactions as a single interaction. For example, if a firm has two alliances with the same firm (i.e., repeated interactions), this would

be treated as only one alliance. Such a simplification would underestimate a firms' involvement in R&D collaborations. The weighted  $k$ -core centrality addresses this problem using the weighted degree, which accounts for repeated interactions.

Secondly, we empirically test whether firm embeddedness has a positive or a negative effect on firms' innovation. We operationalize innovation output as the number of patents filed by a firm. For our analysis, we combine two databases: (i) the NBER patent database containing information on almost three million patents (Hall et al., 2001), and (ii) the SDC Platinum alliances database listing 21,572 alliances involving 13,936 firms between 1984 and 2009 (Thomson-Reuters, 2013).

From our analysis, we find that firms' embeddedness correlates with their innovation output and significantly affects their future innovation output. Specifically, a regression analysis shows that firms with a higher weighted  $k$ -core centrality in a given year also have a higher innovation output in the following year. The results are robust and significant even after controlling for the number of previously filed patents, industrial sector, and other measures capturing firms position in the alliance network.

The remaining of this paper is divided as follows. We first present the data to infer R&D alliances and measure innovation output via patent data. We then introduce the  $k$ -core centrality to operationalize our notion of embeddedness. We describe the regression analysis used to show the effect of a firm's embeddedness on its patenting activity. In the Results section, we illustrate the evolution of the alliance network and offer results supporting the hypothesis that embeddedness is indeed correlated with patenting activity. Further evidence in support of this result is then provided by showing the results of the regression analysis. Finally, in the discussion, we summarize the findings and discuss how they fit in the extant literature.

## 2 Materials and methods

### 2.1 Data sets

In this paper, we build on two data sets. The first data set, obtained from Thomson Reuters' SDC Platinum alliances database, contains all publicly announced R&D partnerships between firms (Thomson-Reuters, 2013). Because the SDC database does not provide a unique identifier for each firm, we use the firm names reported in the dataset. Therefore, we disambiguate names, i.e., we correct for the cases where two or more entries with different names corresponded to the same firm, by manually controlling for spelling, legal extensions (e.g., LTD, INC), and any other recurrent keywords (e.g., BIO, TECH, PHARMA, LAB). We keep subsidiaries of a firm located in different countries as separated entities.

In total, we use information about 21,572 alliances involving 13,936 firms between 1984 and 2009. We note that these alliances can involve different economic actors, e.g., universities, but we refer to them as firms. We used a firm's 4-digit Standard Industrial Classification (SIC) code to classify the industrial sector. The data allows us to reconstruct a network in which nodes represent firms, and links represent R&D collaborations. On this network, we perform a weighted  $k$ -core decomposition to compute the coreness value of each firm as a proxy of its embeddedness in the R&D network.

The second data set is obtained from the NBER patent database of the National Bureau of Economic Research (Hall et al., 2001). It contains detailed information on almost three million patents granted in the U.S.A. between 1974 and 2006. Every patent is assigned to one or more assignees and is classified according to the International Patent Classification (IPC) system. We use this data to estimate the innovation output of each firm by means of its number of patents in the respective time window.

## 2.2 Reconstructing the R&D network

Because the network of firm collaborations is highly dynamic, there are various ways of studying its evolution over time. Work by Tomasello et al. (2016) focuses on the changing growth pattern in consecutive 5-year time intervals to reveal a remarkable life cycle dynamics of the R&D network over the whole period of 25 years. Here, we focus on how a firm's (current) position, i.e., embeddedness, in the network affects its innovation output.

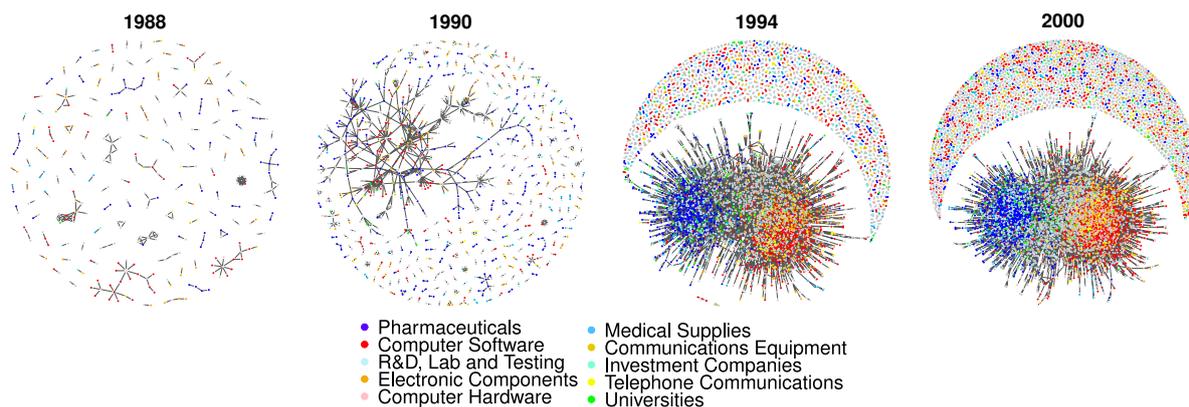


Figure 1: Snapshots of the R&D network showing its evolution and the emergence of a large connected component.

We start by reconstructing the cumulative R&D network: we use as time resolution one year and add a new link to the cumulative network every time an alliance of two firms is listed in the dataset in this time window. When an alliance involves more than two firms (consortium), all

the firms involved are connected in pairs, resulting in a fully connected clique. The weight  $w^{ij}(t)$  of a link indicates the total number of alliances between firms  $i$  and  $j$  up to time  $t$ . If, during the same time interval, two firms  $i$  and  $j$  have more than one collaboration on different projects, such multiple links are also considered in the weight.

### 2.3 Quantifying network embeddedness

From the reconstructed cumulative R&D network, we can compute the embeddedness of nodes in this network. The measure used to proxy embeddedness is based on *coreness*  $C_C^i$ . For a given network, a coreness value can be assigned to every node using the  $k$ -core decomposition. This procedure recursively removes all nodes with a degree less than  $d$  until only nodes with a degree equal to or larger than  $d$  remain. The procedure starts with  $d = 1$ , i.e., it removes all nodes that have only one neighbour in the networks. The removal may leave these neighbouring nodes with a single neighbour. In the second step of the procedure, such nodes are also removed. Their removal again may leave other nodes with one remaining neighbour. Thus, in the third step, they are also removed and so forth, unless no nodes can be removed. Then, all nodes that have been removed during this procedure are assigned a shell number  $k_s$  equal to  $d$ .

Nodes with a small values of  $k_s$  are removed very early because they are topologically weakly embedded in the network. Conversely, nodes with the largest value of  $k_s = k_s^{\max}$  form the core of the network. The difference between the  $k_s^i$  value of a node  $i$  and the value of the core is called *coreness*, or distance from the core,  $C_C^i = k_s^{\max} - k_s^i$ . Nodes close to the core, i.e. with *small* coreness values, are well embedded in the network, while nodes with *large* coreness values are not.

Since the method described uses only information about the node degree and ignores the link weights, it is called unweighted  $k$ -core decomposition. This method has been successfully applied to characterize various real-world networks (Carmi et al., 2007; Garas et al., 2010). Moreover, Kitsak et al. (2010) showed that the coreness value of a node is a more accurate predictor of its spreading potential than, for example, its degree.

This paper uses an extension of the unweighted  $k$ -core decomposition, which also considers the link weights. The weighted  $k$ -core decomposition (Garas et al., 2012) uses the same procedure to remove nodes as the unweighted version, but a refined measure for the node degree, called weighted degree,  $d'$ . The weighted degree,  $d'$ , depends on two free parameters  $\alpha$  and  $\beta$  balancing the influence of the weights  $w_{ij}$  which indicate multiple alliances between the same firms. Similar to Garas et al. (2012), we consider the case when  $\alpha = \beta = 1$ . With this choice, we assign to the

weight and the degree the same importance, and the equation for the weighted degree of node  $i$  becomes a geometric mean:

$$d'_i = \left( d_i^\alpha \left( \sum_j^{d_i} w_{ij} \right)^\beta \right)^{\frac{1}{\alpha+\beta}} \xrightarrow{\alpha=\beta=1} d'_i = \sqrt{d_i \left( \sum_j^{d_i} w_{ij} \right)} \quad (1)$$

$d^i$  is the degree of node  $i$  and  $w^{ij}$  is the weight of the link between nodes  $i$  and  $j$ . The summation goes over all neighbours of  $i$ .

Note that the coreness value in a dynamic network, such as the alliance network, changes over time. Specifically, there are two processes affecting the coreness of a firm  $i$ : (i) formation of alliances involving  $i$  (direct) and (ii) formation of alliances which not involving  $i$  (indirect).

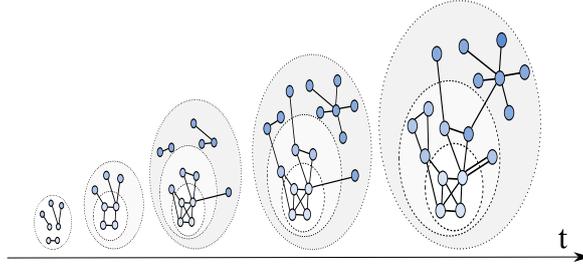


Figure 2: Illustration of the network evolution where new  $k$ -shells emerge as new links are formed.

In Figure 2 we provide a simplified example of a growing network where similar to real R&D networks, new firms enter the network by creating new links either with existing firms or newcomers. In the first process firm,  $i$  plays an active role in forming new alliances, and its coreness may change accordingly. In the second process, firm  $i$  plays no direct role in forming new alliances. However, the network grows, and new shells emerge. Therefore, a firm  $i$ 's position may still change. This implies that the coreness of a particular firm  $i$  may change even without any new R&D alliances involving that firm.

All coreness values reported in this paper are based on the *weighted*  $k$ -core decomposition. This differs from the work by Powell et al. (1999), where the Katz-Bonacich centrality is used to operationalize embeddedness. This centrality is similar to the known eigenvector centrality, i.e. it considers the weight of direct and indirect neighbours in measuring the importance of nodes. We argue that the Katz-Bonacich centrality has substantial limitations in measuring embeddedness because it may reflect the degree of a node, which is not embeddedness.

To illustrate this argument, in Figure 3 we plot a network where nodes are characterized by three different common measures. The node's size indicates its Katz-Bonacich centrality, the colour of a node indicates its coreness value, and the number of links of a node indicates its

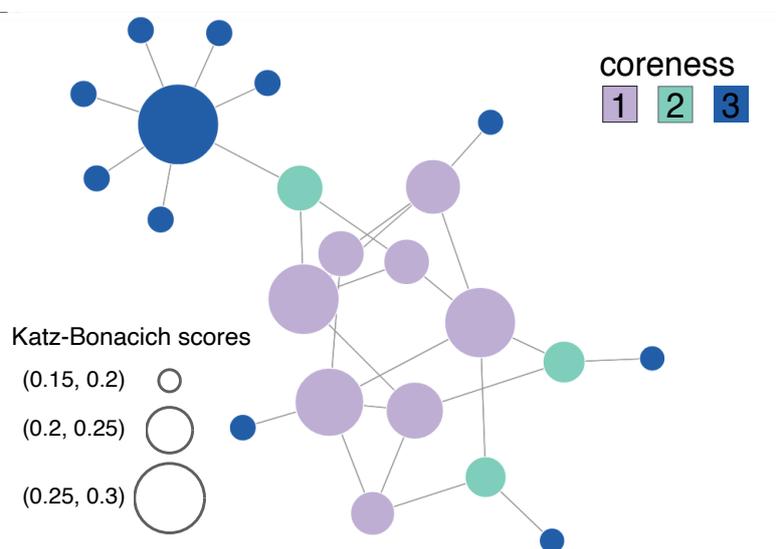


Figure 3: Differences between degree, Katz-Bonacich centrality, and coreness. We illustrate a undirected network where the Katz-Bonacich centrality fails to capture the embeddedness of nodes in a network as it mainly captures degree.

degree-centrality. We see that nodes with a high degree also have a large size, i.e., a higher Katz-Bonacich centrality. However, these nodes are not necessarily well embedded as the node in the upper left part of Figure 3 shows. Despite its high degree and high Katz-Bonacich centrality, this node can be easily disconnected from the network by removing one link; hence it is not well embedded. The coreness measure addresses this problem by assigning a low coreness value to nodes that remain connected after the link removal procedure.

## 2.4 Regression variables

For the regression analysis, we follow the approach of Schilling and Phelps (2007) who reconstructed the alliance network by aggregating alliances in time windows of three years. This aggregation yields an unbalanced panel because a firm can “disappear” if it does not form alliances in all time windows. Nevertheless, alternative approaches would result in other problems. For example, considering time-aggregated alliance networks with increasing time windows as done in (Tomasello et al., 2014; Vaccario et al., 2018) would have the drawback that new alliances have the same weight as much older alliances. This assumption is unreasonable to study how alliances affect innovation output. Empirical studies have shown that alliances end after about three years (Phelps, 2010). Therefore, in line with Schilling and Phelps (2007), we consider a time window of three years.

To work with the time windows, one could carry out a regression for each year separately and include only active firms. However, by adopting this approach would hinder the comparison of parameters across models. Alternatively, one could restrict the analysis to years with more firms forming alliances and filing patents. This approach, however, creates a bias because outcome variables would be used to prepare the data sample for the regression. For these reasons, we opt for a pool regression and use all firms' observations from different years. Moreover, we control for (i) sector and (ii) time fixed effects.

The *dependent* variable is the innovation output of a firm, which we quantify using the number of patents filed by a firm in the next year. The *independent* variable is the firms' embeddedness, which we quantify using the coreness value  $C$ . For each firm, we control for the number of patents published in the previous five years, the number of partners, betweenness centrality, local efficiency, local clustering coefficient, and local reach for the following reasons.

We control for the number of partners because Shan et al. (1994) have shown that innovation output is significantly affected by the number of cooperative relationships. Moreover, Powell and Brantley (1992) found that interfirm cooperation rises with their size, and hence, the number of partners is also a control for the firm size. Owen-Smith and Powell (2004) proxied a firm's ability to absorb information flows via betweenness centrality and shown that it affects firms' innovation. The local clustering coefficient measures how many neighbours of a focal node are also connected, this way forming cliques or clusters. Local reach quantifies the fraction of all nodes in a network reachable from a focal node. Both measures are the local versions of the variables investigated by (Schilling and Phelps, 2007) which argued for their importance. Also, they found that local efficiency, i.e., to which extent a firm's partners are non-redundant, has a significant effect. Finally, we control for industrial sectors as different sectors have different patenting practices.

In Table 1, we summarize the variables that we will use to model firms' innovation output. Schilling and Phelps (2007) and Owen-Smith and Powell (2004) discuss these variables and their economic meaning exhaustively. The interested reader can refer to them.

## 2.5 Regression model

We use a zero-inflated negative binomial model. The negative binomial is a good model for firms' number of patents as this number is an over-dispersed count variable. Indeed, in one year, a single firm files between 0 to more than 2000 patents. The mean patent count is around 17.3, and the standard deviation is 111.24, i.e., almost 10 times bigger than the mean. However, we also have that more than 50% of firms file 0 patents in a year. We use a zero-inflated model to account for the excessive zero counts in the data. We assume that there are two processes driving firms' innovation output. Firms either have the capability to file patents, or they do not. If they do not

<b>Variables</b>	
<b>Dependent</b>	Number of patents between $t$ and $t + 1$ , $P_i(t + 1)$
<b>Independent</b>	Coreness $C_i(t)$
<b>Control</b>	
<b>Firm-level</b>	
Pre-sample of patents	# of patents in $[t - 5, t]$
Number of partners	$d_i(t)$
Betweenness centrality	$b_i(t)$
Local clustering coefficient	$2e_i/(d_i(d_i - 1))$
Local reach	$\sum_{j \neq i} 1/d(i, j)$
Local Efficiency	$\sum_j \left(1 - \sum_k \frac{A_{ik}}{\sum_l A_{il}} \delta_{jk}\right)$
<b>Network-level</b>	
Industrial sector	a dummy for each sector
<b>Time-level</b>	
Year of the alliance network	a dummy for each year
<b>Model specification</b>	Zero-inflated Negative Binomial

Table 1: The four types of variables used to explain firms' innovation output. For the details on how to compute them see Schilling and Phelps (2007); Newman (2018); Burt (2009).

have the ability, their patent count is 0, while in the other case, we model their patent count using a negative binomial. Note that by using a zero-inflated model, we assume that firms with the ability to file patents can still fail and then obtain zero as patent count.

We consider three different model specifications. In all models, the independent variables and controls are computed before time  $t$ , while the dependent variables are computed at  $t + 1$ . Hence, the general form of the models is:

$$P_i(t + 1) = \mathcal{F}[C_i(t), F_i(t), I_i(t)] \quad (2)$$

where  $P^i(t + 1)$  is number of patents filed by firm  $i$  between  $t$  and  $t + 1$ ,  $C_i(t)$  is the coreness of firm  $i$  at time  $t$ , and  $F^i(t)$  and  $I_i(t)$  are the controls at the firm and the network level. *Model 1* contains only the control variables used in (Schilling and Phelps, 2007). *Model 2* contains all the variables of Model 1, and we add the local reach and local clustering coefficient, i.e., the local versions of the variables studied in (Schilling and Phelps, 2007). In addition, we add the number of partners as a control for firm size. These two models are baseline models showing how well we can capture the innovation output of firms *without* using the firm embeddedness. *Model 3* contains all the variables of Model 2, and we introduce the measure of firms' embeddedness, i.e., their coreness  $C_i$ . For the regression models, we name our measure CORE.

### 3 Results

#### 3.1 Core-periphery structure

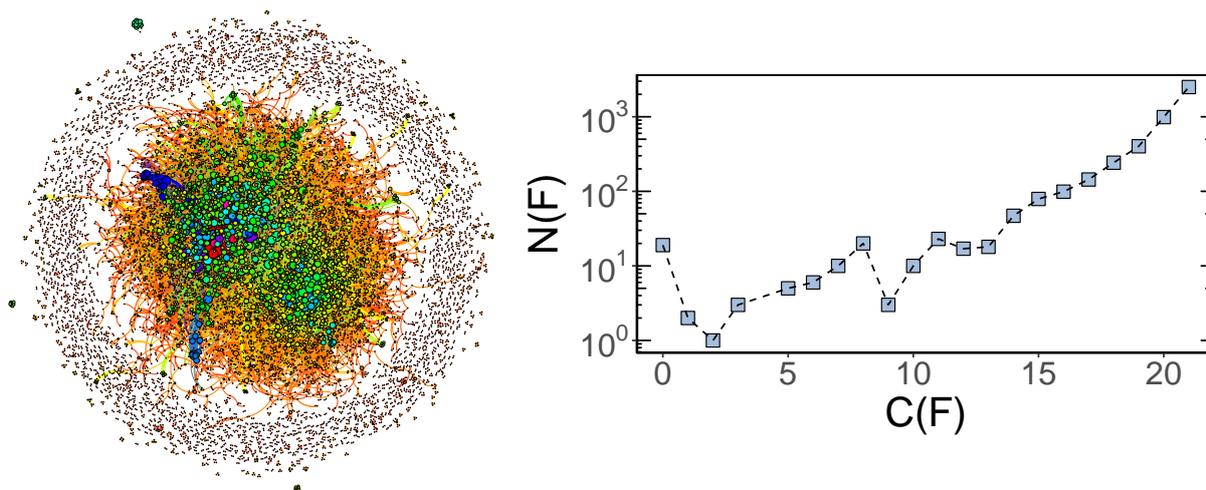


Figure 4: (*left*) Graphical representation of the cumulative R&D network at the end of 2009. This plot is made with Bastian et al. (2009) using the OpenOrd layout. The different colors represent different coreness values, with red assigned to the core nodes. (*right*) The network has a strong core-periphery structure (Borgatti and Everett, 2000), i.e. only a small number of nodes having small  $C$  values, while the majority of the nodes are located in the periphery and have large  $C$  values.

We apply the weighted  $k$ -core decomposition to the cumulative R&D network to assign a coreness value to each firm. When the network evolves over time, a firms' position in the network may change. We take the cumulative R&D network at the very last year, which is the year 2009, and indicate the firms' coreness value in this maximal network as  $C_i(F)$  ( $F$  stands for final).

Figure 4 (left) shows a network plot of the cumulative R&D network in 2009. The nodes are coloured according to their coreness value  $C_i(F)$ , and their size is proportional to their cumulative degree  $d_i$ , i.e., the cumulative number of allied partners. The figure highlights that several firms, despite their many alliances, are not part of the core but of the periphery.

Figure 4 (right) shows the frequency of coreness values,  $N_{C(F)}$ , for the cumulative R&D network in 2009. Taking into account that  $N = 13,936$ , we see that the total number of firms with small coreness values,  $0 \leq C(F) \leq 5$ , i.e. firms that are part of, or very close to, the core are relatively small, but there is a large number of firms in the periphery,  $C(F) > 5$ . This topological property indicates a very pronounced core-periphery structure, and hence the R&D network exhibits a

large variation of coreness values. Given the broad distribution of coreness values, in our analysis, we can explore heterogeneous firm embeddedness.

### 3.2 Correlations between network position and innovation output

We explore to what extent the position of firms in the R&D network is indicative of innovation output. As explained in Section 2.1, we measure the firms' innovation output via the number of patents they file. We first show how this number correlates with the firms' coreness values. The results are shown in Figure 5. According to the distribution of coreness values shown in Figure 4 (right), each firm belongs to one out of 21 coreness classes indicated by  $C(F)$ . In Figure 5 we show the average number of patents for each of these classes. Although the patent data is somewhat scattered, there is a clear indication that the number of patents *increases* with a better network position, i.e., with smaller  $C(F)$  where  $C(F) = 0$  indicates the core.

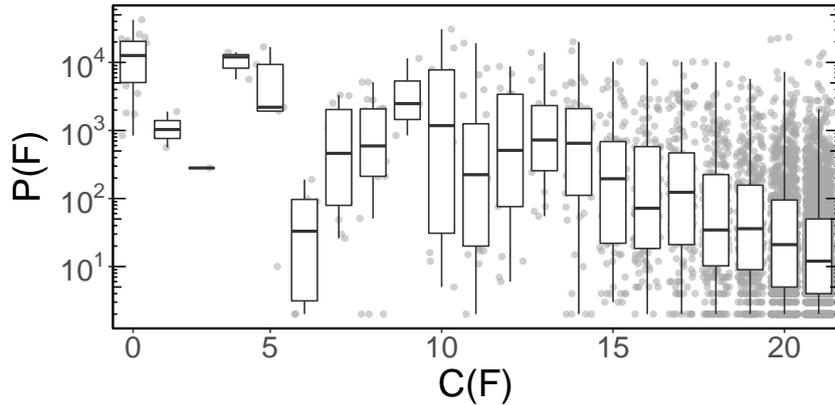


Figure 5: Box plot of the number of patents  $P$  against, coreness  $C(F)$ . The black middle line represents the median. The top and bottom of hinges denote the 25 and 75 percentiles, respectively. The whiskers represent the 95% confidence interval. The scatterplot shows the actual data points.

More precisely, by defining with  $\bar{P}(F)$  the average number of patents filed by firms with coreness  $C(F)$ , the pearson correlation between  $\bar{P}(F)$  and  $C(F)$  is  $-0.581$ . This means that the weighted coreness of a firm — i.e. a topological measure — becomes a very strong indicator of its performance in R&D activities, as measured by the number of patents.

After showing that the coreness of a firm correlates positively with innovation output in the cumulative network, we look at its evolution over time. To make coreness values comparable at different times, we define relative coreness as  $c_i(t) = C_i(t)/C_{\max}(t)$ , i.e. as the ratio between the

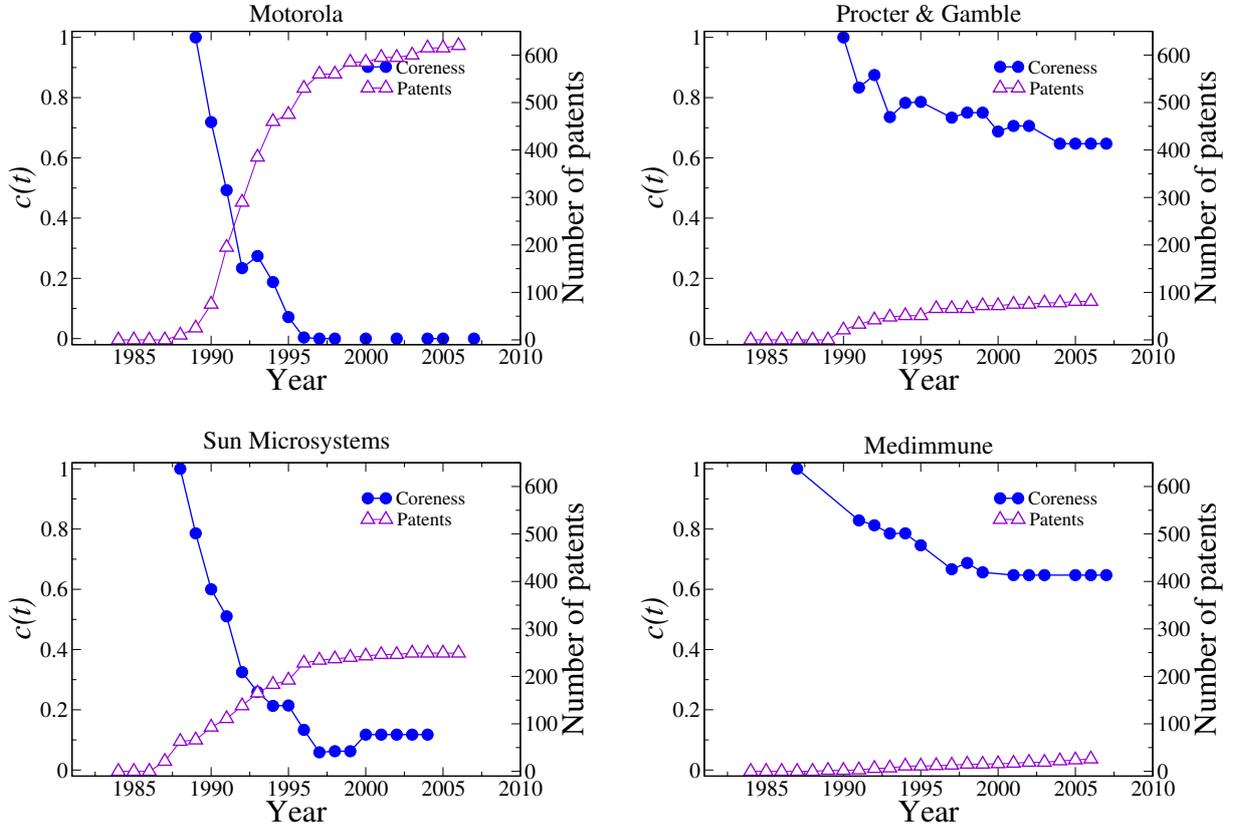


Figure 6: Coreness change and number of patents: (top-left) Motorola, (bottom-left) Sun Microsystems, (top-right) Procter & Gamble, and (bottom-right) Medimmune.

current coreness  $C_i(t)$  and the maximum coreness  $C_{\max}(t)$  at the same time. The variable  $c_i(t)$  lies between 0 and 1 where 0 is the very *core* and 1 to the outermost *periphery*.

In Figure 6, we present two firms, **Motorola** and **Sun Microsystems** that reached the core of the R&D network and two that did not: **Procter&Gamble** and **Medimmune**. All four firms start with high relative coreness values. **Motorola** and **Sun Microsystems** have a declining relative coreness which reach small values. In other words, these two firms move closer to the core thanks to the actively formed alliances and the general evolution of the network. This dynamic distinguishes them from **Procter&Gamble** and **Medimmune** that remain in the periphery. Looking at the number of patents filed, we see that the two firms moving closer to the core also file more patents, while the two firms in the periphery filed fewer. A more detailed investigation of the coreness evolution and the possible mechanisms to reach the core are discussed in Schweitzer et al. (2021).

### 3.3 Regression results

In Table 2, we report a summary of the regressions performed. The first result is that the past number of filed patents,  $\log(\text{PAT} + 1)$ , is an important control. Precisely, this variable controls for firms' heterogeneity (Blundell et al., 1995), and hence, for their ability to either file or not file patents. In other words, a firm's momentum in filing patents positively affects its innovation output. This result is in line with the one obtained by Schilling and Phelps (2007).

Our hypothesis is that coreness has a negative effect on firms' innovation output. This hypothesis is supported by the significant and negative effect of **CORE** in *Model 3* (see Table 2). If we compare *Model 2* and *Model 3*, we find that the two control variables change their significance: The effect size of **LOCAL\_REACH** becomes indistinguishable from zero, while **EFF** (local efficiency) becomes statistically significant. Their change in significance means that these variables have some correlations with **CORE**. Such correlations are partially expected because all these centrality measures tend to be correlated (Freeman, 1979). To check that these correlations do not affect our results, we performed a robustness analysis.

We create a fourth model, *Model 4*, in which we remove **EFF**, the local efficiency, but keep **CORE**. This way, we check if the effect of **CORE** is driven by the potential correlation with this variable. For this model, the parameter of **CORE** loses about 10% of its value (see Table 2). However, it stays negative and significant.

In Figure 7, we also illustrate this change by plotting the effect size of the analyzed network measures for *Model 3* and *Model 4*. From this visualization, we see that the effect size of **CORE** remains almost unchanged. This indicates that the correlation with **EFF** increases the importance of **CORE** but does not change the sign and significance of its effect. In the Appendix, we report the result for a fifth model where we keep only **CORE** and remove all the other centrality measures. Also, in this model the significance and sign of **CORE** do not change (see Table 5).

Finally, we note that *Model 3* and *Model 4* have about the same explanatory power (Log-Likelihood), but both are better than *Model 2*, i.e., the model without coreness. To make test this statement more rigorously, we use the Vuong closeness test (Vuong, 1989). This test allows us to compare different models, perform model selection, and tell us which model captures the data better. In Table 3, we report the results of the Vuong closeness test and find that *Model 3* and *Model 4* are both significantly better than *Model 2*. However, we cannot decided between *Model 3* and *Model 4*. Overall, our analyses confirm that **CORE** is a good explanatory variable for firms' innovation output.

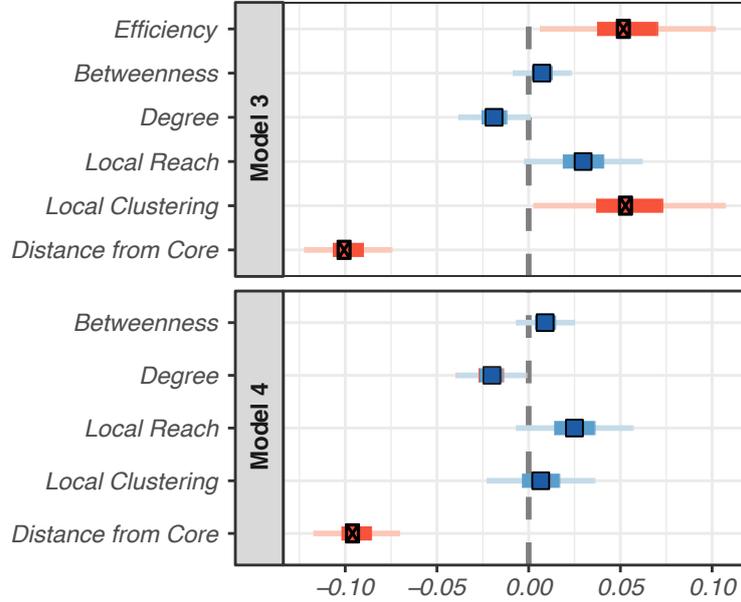


Figure 7: Coefficient plots for *Model 3* and *Model 4*. In blue we plot the effect size that are not indistinguishable from zero, while in red we plot the statistically significant effects. Note that in *Model 4* we remove the (local) efficiency EFF and find that the effect size for the coreness CORE remains almost unchanged.

	Model 1	Model 2	Model 3	Model 4
Zero model: (Intercept)	1.25 (0.07)***	1.25 (0.07)***	1.28 (0.21)***	1.28 (0.21)***
Zero model: log(PAT + 1)	-1.13 (0.04)***	-1.13 (0.04)***	-1.13 (0.07)***	-1.14 (0.07)***
(Intercept)	-0.38 (0.12)**	-0.32 (0.12)**	-0.31 (0.12)**	-0.30 (0.11)**
log(PAT + 1)	0.90 (0.01)***	0.89 (0.01)***	0.88 (0.01)***	0.87 (0.01)***
EFF	-0.02 (0.01)	0.02 (0.03)	0.05 (0.02)**	
BETWEENNESS_NORM	0.02 (0.01)***	0.01 (0.01)	0.01 (0.01)	0.00 (0.01)
DEGREE		0.02 (0.01)	-0.02 (0.01)	-0.02 (0.01)
LOCAL_REACH		0.05 (0.02)**	0.03 (0.02)	0.02 (0.02)
LOCAL_CLUSTERING		0.03 (0.03)	0.05 (0.05)	0.01 (0.03)
CORE			-0.10 (0.03)***	-0.09 (0.02)***
AIC	48419.62	48410.93	48344.97	48347.99
Log Likelihood	-24180.81	-24173.47	-24139.48	-24143.99
Num. obs.	12649	12649	12649	12649

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$ . The and year fixed effect are reported in Appendix.

The errors have been clustered at CORE (21 classes).

Table 2: Zero-inflated negative binomial models. For this regression, we have standardized the variable in order to be able to compare the effect sizes. For CORE, the average value is around 16.53, and the standard deviation is around 2.82.

	Vuong z-statistic		$H_A$		$p$ -value
Raw	-3.357278	<i>Model 3</i>	>	<i>Model 2</i>	0.00039357
AIC-corrected	-3.258480	<i>Model 3</i>	>	<i>Model 2</i>	0.00056005
BIC-corrected	-2.890686	<i>Model 3</i>	>	<i>Model 2</i>	0.00192201
Raw	-3.24855	<i>Model 4</i>	>	<i>Model 2</i>	0.00057997
AIC-corrected	-3.24855	<i>Model 4</i>	>	<i>Model 2</i>	0.00057997
BIC-corrected	-3.24855	<i>Model 4</i>	>	<i>Model 2</i>	0.00057997
Raw	1.1754242	<i>Model 3</i>	>	<i>Model 4</i>	0.11991
AIC-corrected	0.6598649	<i>Model 3</i>	>	<i>Model 4</i>	0.25467
BIC-corrected	-1.2593904	<i>Model 3</i>	<	<i>Model 4</i>	0.10394

Table 3: Results of the Vuong closeness test.

## 4 Discussion

**Firms’ embeddedness and its evolution.** In this paper, we argue that the embeddedness of firms in the R&D network is indicative of their ability to innovate. The latter is proxied by the number of patents firms file in a given time window. To proxy embeddedness, we introduce a new measure, coreness. It is a relative measure to compare firms’ positions in the network at different stages of the evolution. To calculate coreness, we use the *weighted k-core decomposition* (Garas et al., 2012), which takes into account multiple R&D alliances between firms. The decomposition applies a sequence of node removals to prune the network and assigns a value to each firm that indicates its current distance from the core. In this respect, it reflects the embeddedness of firms better than previously used network centrality measures.

As the R&D network evolves both as new firms enter and new alliances are formed, firms’ coreness values change, reflecting changes in their network positions. We observe the emergence of a clear core-periphery structure characterized by a dense core containing a smaller number of firms and a sparse periphery containing the majority of less integrated firms. Firms that have reached the core of the alliance network are shown to be more successful in their innovation output than firms in the periphery. Analyzing the relation between the coreness values of firms and the number of patents, we find a strong correlation.

**Embeddedness and innovation.** To better quantify the higher innovation output coming from higher embeddedness, we have performed a regression. We have used a zero-inflated negative binomial regression to model firms’ innovation output, measured by the number of patents filed in the subsequent year. We find that decreasing the coreness by one unit increases the logarithm of the patent count by  $\sim 0.09/2.82 = 0.03$  (see Table 2). This effect size implies that, for example, given a firm with coreness  $C$  filing 100 patents, an identical copy of this firm with coreness  $C - 1$  would file about 103 patents. Another example: Given a firm with coreness  $C$  filing 10

patents, an identical copy of this firm with coreness  $C - 1$  would file about 13.5 patents. Our results are significant even after controlling for several factors, including firms' past innovation output (Blundell et al., 1995), ability to absorb knowledge flows in the alliance network (Phelps, 2010), the number of partners, network clustering and structural holes (Ahuja, 2000; Baum et al., 2000). Overall, our results suggest that embeddedness has a positive and significant effect on innovation output in R&D activities.

**Comparison with previous works.** Similar to (Schilling and Phelps, 2007), our regression analysis uses negative binomial as firms' patent count is an over-dispersed variable. Different from (Schilling and Phelps, 2007), we use a zero-inflated version of the negative binomial as 50% of the firms in our sample file 0 patents in a year. Similar to (Powell et al., 1999), we find that embeddedness has a significant effect on firms innovation output. However, different from (Powell et al., 1999), we do not quantify firms' embeddedness using the Katz-Bonacich centrality but use a new indicator, *coreness*. It is based on the  $k$ -core decomposition as in (Al-Laham and Bort, 2011) but takes multiple relations between firms into account. This decision was motivated by the fact that Katz-Bonacich centrality may fail to measure embeddedness (see Sect. 2.3).

Following (Tomasello et al., 2014, 2016; Vaccario et al., 2018), we consider an alliance network with firms across industrial sectors. That means, in contrast to the results of Powell et al. (1999); Al-Laham and Bort (2011), our findings are not restricted to the biotechnology sector. Firms' innovation might depend on complementary capabilities coming from any industrial sector, and hence, firms' embeddedness should be quantified in a network containing all the industrial sectors. By doing this, the presented results are more general.

**Limitations and outlooks.** The analysed data is limited since it contains only collaborations until 2009 and patent data until 2006. At the same time, to our knowledge, we have performed the largest analysis to quantify the link between the firms' embeddedness and their innovation output, with more than 13,000 firms over a time horizon of 25 years. We have performed various statistical analyses to ensure the robustness of our analysis.

Given the evolution of the alliance network (Tomasello et al., 2016) and firms' embeddedness (see Sect. 3.2), further studies could investigate how quickly the innovation output of firms change after they have obtained a more central position. To extend our setup, we could add a regression analysis also considering the patents filed after two and three years and investigate how the effect size changes.

## References

- Ahuja, G. (2000). Collaboration networks, structural holes, and innovation: A longitudinal study. *Administrative Science Quarterly* 45(3), 425–455.
- Al-Laham, A. and S. Bort (2011). Chapter 13 the innovation outcomes of mnc subsidiaries' local embeddedness: Evidence from the german "bioregion rhein-neckar-dreieck" local network". *Entrepreneurship in the Global Firm (Progress in International Business Research, Volume 6)*. Emerald Group Publishing Limited, 291–323.
- Bastian, M., S. Heymann, and M. Jacomy (2009). Gephi: An open source software for exploring and manipulating networks. In *International AAAI Conference on Weblogs and Social Media*.
- Baum, J. A., T. Calabrese, and B. S. Silverman (2000). Don't go it alone: Alliance network composition and startups' performance in canadian biotechnology. *Strategic management journal* 21(3), 267–294.
- Blundell, R., R. Griffith, and J. V. Reenen (1995). Dynamic count data models of technological innovation. *The Economic Journal* 105(429), 333–344.
- Bollobás, B. (1988). *Graph theory and combinatorics 1988*. Elsevier.
- Borgatti, S. P. and M. G. Everett (2000). Models of core/periphery structures. *Social Networks* 21(4), 375 – 395.
- Burt, R. S. (2009). *Structural holes: The social structure of competition*. Harvard university press.
- Carmi, S., S. Havlin, S. Kirkpatrick, Y. Shavitt, and E. Shir (2007, July). A model of Internet topology using k-shell decomposition. *Proceedings of the National Academy of Sciences of the United States of America* 104(27), 11150–4.
- Freeman, C. (1991). Networks of innovators: a synthesis of research issues. *Research policy* 20(5), 499–514.
- Freeman, L. C. (1979). Centrality in social networks conceptual clarification. *Social Networks* 1(3), 215–239.
- Garas, A., P. Argyrakis, C. Rozenblat, M. Tomassini, and S. Havlin (2010, November). Worldwide spreading of economic crisis. *New Journal of Physics* 12(11), 113043.
- Garas, A., F. Schweitzer, and S. Havlin (2012, August). A  $k$ -shell decomposition method for weighted networks. *New Journal of Physics* 14(8), 083030.
- Gilsing, V., B. Nooteboom, W. Vanhaverbeke, G. Duysters, and A. van den Oord (2008). Network embeddedness and the exploration of novel technologies: Technological distance, betweenness centrality and density. *Research Policy* 37(10), 1717–1731. Special Section Knowledge Dynamics out of Balance: Knowledge Biased, Skewed and Unmatched.
- Granovetter, M. (1985). Economic action and social structure: The problem of embeddedness. *American journal of sociology* 91(3), 481–510.
- Gulati, R., N. Nohria, and A. Zaheer (2000). Strategic networks. *Strategic management journal* 21(3), 203–215.
- Hall, B. H., A. B. Jaffe, and M. Trajtenberg (2001). The nber patent citation data file: Lessons, insights and methodological tools.
- Kitsak, M., L. K. Gallos, S. Havlin, F. Liljeros, L. Muchnik, H. E. Stanley, and H. A. Makse (2010, August). Identification of influential spreaders in complex networks. *Nature Physics* 6(11), 888–893.
- Newman, M. (2018). *Networks*. Oxford university press.
- Newman, M. E. J. (2010). *Networks: an introduction*. Oxford; New York: Oxford University Press.

- Nooteboom, B. (1999). *Inter-firm alliances: Analysis and design*. Psychology Press.
- Owen-Smith, J. and W. W. Powell (2004). Knowledge networks as channels and conduits: The effects of spillovers in the boston biotechnology community. *Organization science* 15(1), 5–21.
- Paier, M. and T. Scherngell (2011). Determinants of collaboration in european r&d networks: Empirical evidence from a discrete choice model. *Industry and Innovation* 18(1), 89–104.
- Phelps, C. C. (2010). A longitudinal study of the influence of alliance network structure and composition on firm exploratory innovation. *Academy of Management Journal* 53(4), 890–913.
- Polanyi, K. and R. M. MacIver (1944). *The great transformation*, Volume 2. Beacon press Boston.
- Powell, W., K. Koput, L. Smith-Doerr, and J. Owen-Smith (1999, 01). Network position and firm performance: Organizational returns to collaboration in the biotechnology industry. *Research in the Sociology of Organizations* 16.
- Powell, W. W. and P. Brantley (1992). *Networks and Organizations: Structure, Form, and Action (Chapter 14)*. Harvard Business School Press.
- Schilling, M. A. and C. C. Phelps (2007). Interfirm collaboration networks: The impact of large-scale network structure on firm innovation. *Management science* 53(7), 1113–1126.
- Schweitzer, F., A. Garas, M. V. Tomasello, G. Vaccario, and L. Verginer (2021). The role of network embeddedness on the selection of collaboration partners: An agent-based model with empirical validation. *Advances in Complex Systems* 24(7-8), xxxx.
- Seidman, S. B. (1983). Network structure and minimum degree. *Social Networks* 5(3), 269 – 287.
- Shan, W., G. Walker, and B. Kogut (1994). Interfirm cooperation and startup innovation in the biotechnology industry. *Strategic Management Journal* 15(5), 387–394.
- Thomson-Reuters (2013). Sdc platinum dataset. Date of access: 07/04/2014.
- Tomasello, M. V., M. Napoletano, A. Garas, and F. Schweitzer (2016). The Rise and Fall of R&D Networks. *Industrial and Corporate Change* dtw041.
- Tomasello, M. V., N. Perra, C. J. Tessone, M. Karsai, and F. Schweitzer (2014). The role of endogenous and exogenous mechanisms in the formation of r&d networks. *Scientific Reports* 4.
- Uzzi, B. (1997). Social structure and competition in interfirm networks: The paradox of embeddedness. *Administrative science quarterly*, 35–67.
- Vaccario, G., M. V. Tomasello, C. J. Tessone, and F. Schweitzer (2018, August). Quantifying knowledge exchange in r&d networks: A data-driven model. *Journal of Evolutionary Economics* 28(3), 461–493.
- Vuong, Q. H. (1989). Likelihood ratio tests for model selection and non-nested hypotheses. *Econometrica* 57(2), 307–333.

## Appendix

### A Industry and time effects.

In Table 4 we report all the control variable of *Model 1*, *Model 2*, *Model 3*, and *Model 4*.

In Table 5 we report all the control variable of *Model 3*, *Model 4*, and *Model 5*.

	<i>Model 1</i>	<i>Model 2</i>	<i>Model 3</i>	<i>Model 4</i>
(Intercept)	-0.38 (0.12)**	-0.32 (0.12)**	-0.35 (0.16)*	-0.35 (0.16)*
log(PAT + 1)	0.90 (0.01)***	0.89 (0.01)***	0.87 (0.02)***	0.87 (0.02)***
I_automotive bodies and parts	-0.21 (0.06)***	-0.18 (0.06)**	-0.16 (0.12)	-0.17 (0.12)
I_chemicals	-0.60 (0.05)***	-0.57 (0.05)***	-0.50 (0.03)***	-0.50 (0.03)***
I_computer and office equipment	-0.29 (0.05)***	-0.29 (0.05)***	-0.34 (0.07)***	-0.34 (0.07)***
I_household audiovisual equipment	-0.12 (0.07)	-0.10 (0.07)	-0.19 (0.13)	-0.20 (0.13)
I_measuring and controlling devices	-0.46 (0.06)***	-0.44 (0.06)***	-0.40 (0.05)***	-0.40 (0.06)***
I_medical equipment	-0.42 (0.05)***	-0.39 (0.05)***	-0.34 (0.06)***	-0.34 (0.06)***
I_petroleum refining and products	-0.80 (0.12)***	-0.76 (0.12)***	-0.72 (0.13)***	-0.72 (0.13)***
I_pharmaceuticals	-0.40 (0.04)***	-0.39 (0.04)***	-0.39 (0.07)***	-0.38 (0.07)***
I_aerospace equipment	-0.24 (0.07)***	-0.22 (0.07)***	-0.22 (0.14)	-0.22 (0.14)
I_telecommunications equipment	-0.33 (0.05)***	-0.33 (0.05)***	-0.31 (0.08)***	-0.30 (0.08)***
1988	-0.03 (0.14)	-0.03 (0.14)	-0.08 (0.09)	-0.07 (0.09)
1989	-0.10 (0.14)	-0.09 (0.14)	-0.14 (0.12)	-0.14 (0.12)
1990	-0.15 (0.13)	-0.14 (0.13)	-0.19 (0.10)*	-0.19 (0.10)*
1991	-0.12 (0.12)	-0.12 (0.12)	-0.14 (0.12)	-0.14 (0.12)
1992	-0.15 (0.12)	-0.18 (0.12)	-0.16 (0.13)	-0.16 (0.13)
1993	-0.00 (0.12)	-0.06 (0.12)	-0.01 (0.12)	-0.01 (0.12)
1994	0.17 (0.12)	0.08 (0.12)	0.16 (0.16)	0.16 (0.17)
1995	-0.22 (0.12)	-0.33 (0.12)**	-0.23 (0.12)	-0.23 (0.12)*
1996	-0.15 (0.12)	-0.25 (0.12)*	-0.16 (0.12)	-0.15 (0.12)
1997	-0.34 (0.12)**	-0.40 (0.12)***	-0.31 (0.12)**	-0.32 (0.12)**
1998	-0.32 (0.12)**	-0.37 (0.12)**	-0.29 (0.10)**	-0.29 (0.10)**
1999	-0.40 (0.12)***	-0.42 (0.12)***	-0.36 (0.09)***	-0.36 (0.09)***
EFF	-0.02 (0.01)	0.02 (0.03)	0.05 (0.02)**	
BETWEENNESS_NORM	0.02 (0.01)***	0.01 (0.01)	0.01 (0.01)	0.01 (0.01)
Log(theta)	0.56 (0.03)***	0.57 (0.03)***	0.58 (0.03)***	0.58 (0.03)***
Zero model: (Intercept)	1.25 (0.07)***	1.25 (0.07)***	1.28 (0.21)***	1.28 (0.21)***
Zero model: log(PAT + 1)	-1.13 (0.04)***	-1.13 (0.04)***	-1.13 (0.07)***	-1.13 (0.07)***
DEGREE		0.02 (0.01)	-0.02 (0.01)	-0.02 (0.01)
LOCAL_REACH		0.05 (0.02)**	0.03 (0.02)	0.02 (0.02)
LOCAL_CLUSTERING		0.03 (0.03)	0.05 (0.05)	0.01 (0.03)
CORE			-0.10 (0.03)***	-0.09 (0.02)***
AIC	48419.62	48410.93	48344.97	48347.53
Log Likelihood	-24180.81	-24173.47	-24139.48	-24141.76
Num. obs.	12649	12649	12649	12649

\*\*\* $p < 0.001$ ; \*\* $p < 0.01$ ; \* $p < 0.05$

The errors have been clustered at CORE (21 classes).

Table 4: Zero-inflated negative binomial models.

	<i>Model 3</i>	<i>Model 4</i>	<i>Model 5</i>
(Intercept)	-0.35 (0.16)*	-0.35 (0.16)*	-0.35 (0.16)*
log(PAT + 1)	0.87 (0.02)***	0.87 (0.02)***	0.87 (0.02)***
I_automotive bodies and parts	-0.16 (0.12)	-0.17 (0.12)	-0.17 (0.12)
I_chemicals	-0.50 (0.03)***	-0.50 (0.03)***	-0.51 (0.03)***
I_computer and office equipment	-0.34 (0.07)***	-0.34 (0.07)***	-0.34 (0.07)***
I_household audiovisual equipment	-0.19 (0.13)	-0.20 (0.13)	-0.20 (0.12)
I_measuring and controlling devices	-0.40 (0.05)***	-0.40 (0.06)***	-0.41 (0.05)***
I_medical equipment	-0.34 (0.06)***	-0.34 (0.06)***	-0.35 (0.05)***
I_petroleum refining and products	-0.72 (0.13)***	-0.72 (0.13)***	-0.74 (0.14)***
I_pharmaceuticals	-0.39 (0.07)***	-0.38 (0.07)***	-0.38 (0.07)***
I_aerospace equipment	-0.22 (0.14)	-0.22 (0.14)	-0.22 (0.14)
I_telecommunications equipment	-0.31 (0.08)***	-0.30 (0.08)***	-0.31 (0.08)***
1988	-0.08 (0.09)	-0.07 (0.09)	-0.07 (0.09)
1989	-0.14 (0.12)	-0.14 (0.12)	-0.14 (0.13)
1990	-0.19 (0.10)*	-0.19 (0.10)*	-0.19 (0.10)
1991	-0.14 (0.12)	-0.14 (0.12)	-0.14 (0.12)
1992	-0.16 (0.13)	-0.16 (0.13)	-0.16 (0.13)
1993	-0.01 (0.12)	-0.01 (0.12)	0.00 (0.12)
1994	0.16 (0.16)	0.16 (0.17)	0.17 (0.17)
1995	-0.23 (0.12)	-0.23 (0.12)*	-0.20 (0.12)
1996	-0.16 (0.12)	-0.15 (0.12)	-0.13 (0.12)
1997	-0.31 (0.12)**	-0.32 (0.12)**	-0.30 (0.12)**
1998	-0.29 (0.10)**	-0.29 (0.10)**	-0.29 (0.10)**
1999	-0.36 (0.09)***	-0.36 (0.09)***	-0.36 (0.09)***
EFF	0.05 (0.02)**		
BETWEENNESS_NORM	0.01 (0.01)	0.01 (0.01)	
DEGREE	-0.02 (0.01)	-0.02 (0.01)	
LOCAL_REACH	0.03 (0.02)	0.02 (0.02)	
LOCAL_CLUSTERING	0.05 (0.05)	0.01 (0.03)	
CORE	-0.10 (0.03)***	-0.09 (0.02)***	-0.09 (0.02)***
Log(theta)	0.58 (0.03)***	0.58 (0.03)***	0.58 (0.03)***
Zero model: (Intercept)	1.28 (0.21)***	1.28 (0.21)***	1.28 (0.21)***
Zero model: log(PAT + 1)	-1.13 (0.07)***	-1.13 (0.07)***	-1.14 (0.07)***
AIC	48344.97	48347.53	48345.42
Log Likelihood	-24139.48	-24141.76	-24144.71
Num. obs.	12649	12649	12649

\*\*\* $p < 0.001$ ; \*\* $p < 0.01$ ; \* $p < 0.05$

The errors have been clustered at CORE (21 classes).

Table 5: Zero-inflated negative binomial models.