

# Heterogeneity in knowledge spillovers across regions: The effects of endogeneity, complexity and IPR environment

Ashwin Iyengar (1521001)  
ashwin.iyengar15@iimb.ernet.in

January 9, 2017

## Abstract

I line up empirical evidence to demonstrate the heterogeneity in the geographical distribution of knowledge spillovers across various regions. I then explore three potential mechanisms that may help explain this heterogeneity: endogenous aspects of the regions themselves, complexity of work, and the intellectual property rights environment of the location. While the empirical results are yet inconclusive and incomplete, the current work extends prior work on geographic spillovers of knowledge by integrating three hitherto alternative explanations.

Keywords: Knowledge Spillovers, Endogeneity, Complexity, IPR

## 1 Introduction

There has been a long and illustrious scholarly tradition highlighting the agglomeration characteristics of economic regions, going back at least as far as [Marshall \(2009\)](#), whose original work was published in 1890. More recently, scholars over the last three decades have demonstrated the paper trail of these knowledge spillovers through the study of patent citations (e.g., [Almeida and Kogut \(1999\)](#); [Jaffe et al. \(1993\)](#)). This tradition of scholarship has further

shaped our theoretical understanding of knowledge spillovers through mechanisms such as the effects of inventor mobility (e.g., Almeida and Kogut (1999)), differential Intellectual Property Rights environments across locations (e.g., Zhao (2006)) and of the role of international geography (e.g., Singh (2007)). The nature and extent of the geographical distribution of knowledge spillovers observed in practice is so highly heterogeneous across locations, firms and legal environments, that the understanding of the causal mechanisms leading to knowledge spillovers continues to intrigue the best of scholars. While this is, in no way dismissive of the enormous theoretical strides so far, the question assumes greater significance in the environment surrounding the second machine age as some scholars have begun to highlight (McAfee and Brynjolfsson, 2014)

Motivated by empirical evidence surrounding the heterogeneity in the nature of knowledge flows across the various regions, I intend to explore the three mechanisms ostensibly influencing knowledge spillovers. Complexity of patents invented as a potential mechanism influencing the extent of local knowledge spillovers. This approach is not to be construed as yet another mechanical departure from the current theory on spillovers. I argue so with the following reasons. First, from a human capital perspective, it is valuable to understand the impact of MNCs that dominate much of the cutting- and bleeding-edge innovation in emerging markets on the development of the talent pool in the host country. Does a significant group of local inventors develop? Is this affected by the strength of the IPR regime in the host country? Second, a specific flavor of this question is the investigation of the spillover effects of the

innovation process in emerging countries, or those known to have weaker IPR regimes. Specifically, do multinational firms that develop patentable technologies in emerging (or weaker IPR) countries create spillover effects in the host country talent pool? Or do the benefits remain localized to within multinational companies (MNCs) and their home country employees? Finally, the wide disparity in the extent of knowledge spillovers across locations, across firms and across IPR regimes is intriguing to the researcher and calls attention toward a creative response. a researcher to find the mechanisms that may lie behind such a phenomenon. Patents data allows us to ask these questions and to have them answered as has been in the tradition of [Jaffe et al. \(1993\)](#).

The choice of the three explanatory mechanisms is not arbitrary. Indeed, there has been a tradition of scholarly work in each of them<sup>1</sup>. First, several studies including [Almeida and Kogut \(1999\)](#) have conclusively demonstrated that the kind and extent of knowledge flows between firms in Silicon Valley is unparalleled in the rest of the world. Indeed, our analysis on a select chosen locations demonstrates this adequately in Figure 1 where the lines in purple and red stand out (Silicon Valley, or the broader San Francisco Bay Area is classified into two Metropolitan Statistical Areas (MSAs) by the United States census). Second, as evidenced by the respective scholarly traditions of [Baldwin and Henkel \(2015\)](#), [Ethiraj and Levinthal \(2004\)](#) and [Yayavaram and Ahuja \(2008\)](#), the complexity, intellectual property and organizational implications have been addressed by scholars in the context of patenting. In the spirit of [Ethiraj and Levinthal](#)

---

<sup>1</sup>I am however, unable to go in much of those details in the current article

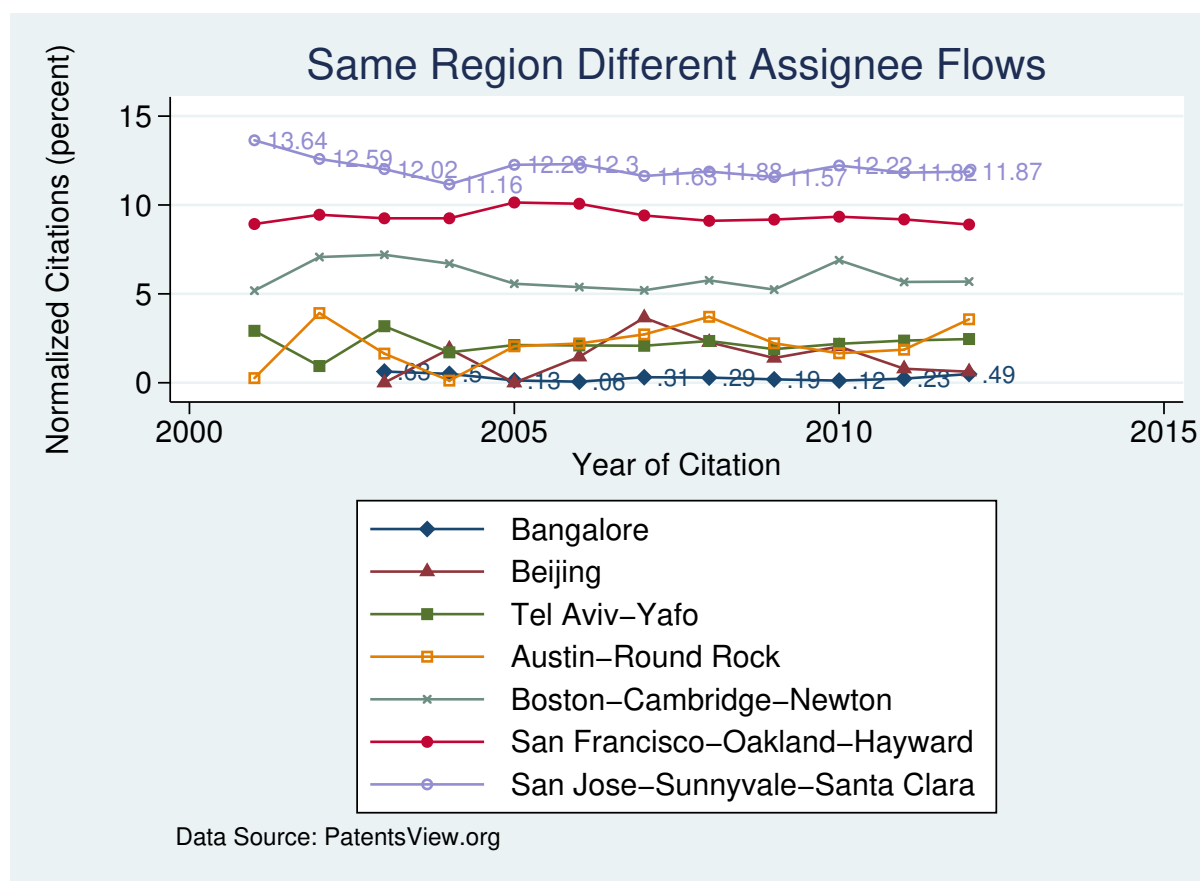


Figure 1: Local and External Flows by Region

(2004), I propose a definition of complexity that is rooted in the question of knowledge spillovers. Specifically, I suggest that complexity may be seen as either an attribute of usage, or as an attribute of invention. A patent that is used (cited) by several patents belonging to distinct and different patent technology classes maybe seen as modular by virtue of it being able to be plugged into multiple, diverse applications. Alternatively, a patent that is constructed with few dependencies may also be seen as being modular by virtue of its capacity to be developed standalone, or with minimal intervention from other modules. For the purposes of this study, I use a definition of Complexity that captures both the effects above. Finally, the scholarly tradition in the international business area has extensively analyzed the relationship between economic geography

Singh (2007), intellectual property environments (Zhao, 2006) and political geography (Singh and Marx, 2013).

The current work is placed at the confluence of these three traditions, with the focus on implications for beneficial knowledge spillovers. A second objective of the current work is to understand the local impact of inventing activity by multinationals in emerging nations. I attempt to answer the following questions. First, how does the nature of the geographic distribution of citations made by inventions from a region affect the quantum of citations received. Second, how do complexity of inventions and cross border differences in intellectual property environments affect the previous relationship.

The benefits of understanding geographic and multinational collaboration in invention is that we may seek to inform both managers and firms about the potential opportunities of tapping into or creating spillover effects in the host country talent pool. Does a significant group of local inventors who develop due to spillovers? Do they then move around to cross-pollinate to other firms? How do domestic firms integrate and appropriate rents in this context. These are some of the many interesting and valuable directions spawned by the current approach.

The rest of this article is organized as follows. In the following section, I present the motivation for empirically designing this study from a selected sampling

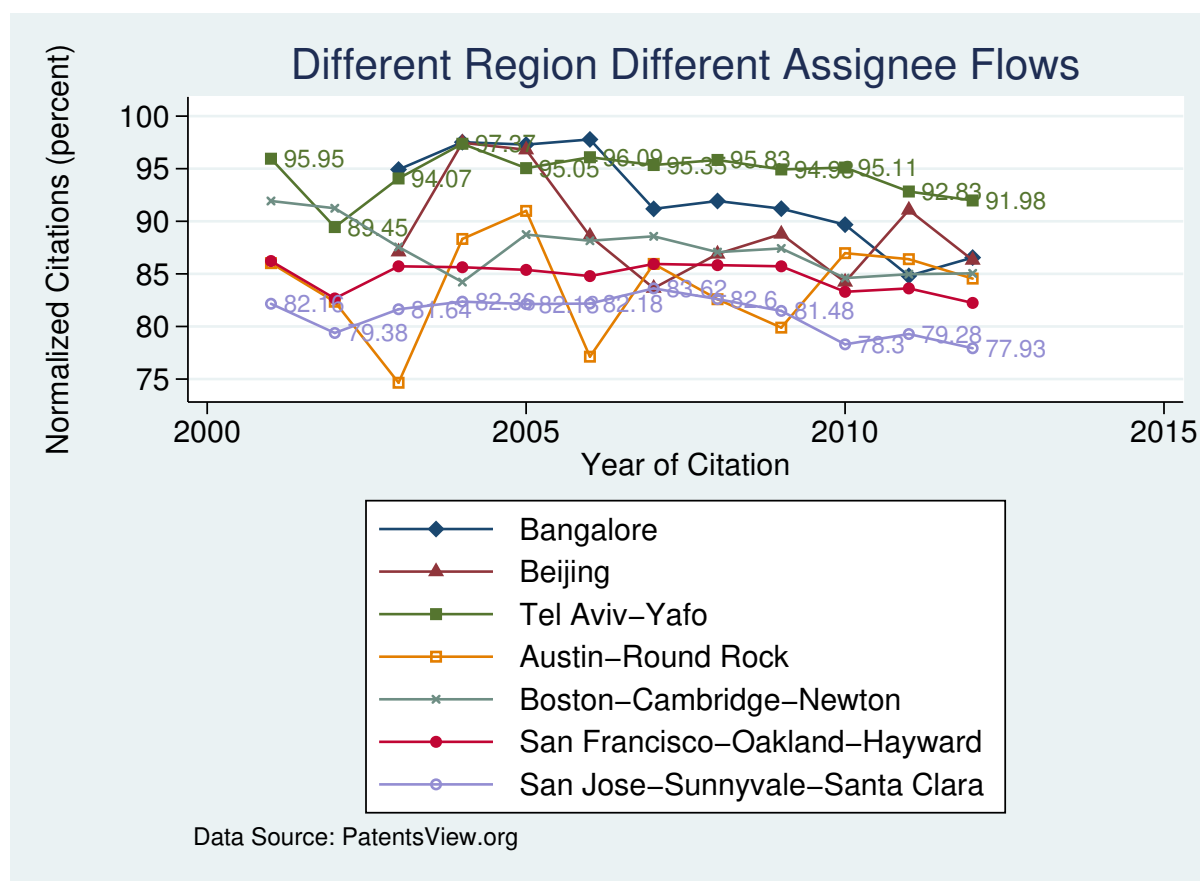


Figure 2: Non-local and External Flows by Region

study. The section following that proposes the hypotheses I intend to test empirically. The approach I take is to assume prior scholarly results about spillovers as a given, and nuance new insights building on top of these giants. The following section on research design presents the constructs created for computing complexity, as well as a discussion on the methodology of the research. My initial, incomplete results are then presented, and a conclusion is drawn of this article as very much a work in progress.

## 2 Motivation

I motivate this study with a small sample analysis of the knowledge spillover characteristics of a selected sample of seven regions across continents, and tech-

nological specialization and IPR strength. In each of Figure 1, Figure 2, Figure 3 and Figure 4, citation counts of regions are expressed as a normalized percentage number (so as to be able to fairly compare across regions with vastly different pools of knowledge and inventors). Figure 1 already demonstrated that the northern California regions stood out in terms of the extent of local knowledge spillovers to other firms. Scholars have explained this using the mechanism of employee mobility (Almeida and Kogut, 1999). Figure 2, on the other hand demonstrates that Tel Aviv-Yafo in Israel stands out as accounting of the highest proportion of flows to external firms in external locations. The employee mobility explanation may not be able explain the phenomenon here. Figure 3 suggests that the Bangalore region sees very little local knowledge spillovers at all, while Figure 4 suggests that the Bangalore region flows are dominated by those to the multinational parent location. These disparate spillover behaviors across locations have been attempted to be explained by scholars (e.g., Singh (2007); Zhao (2006)) by slicing the problem in a specific context (e.g., of an MNC Parent - MNC Subsidiary). This rather wide disparity between a selection of inventing locations, provides us with the context to dive into the understanding the mechanisms that underlie this divergence in knowledge spillover patterns across regions in the aggregate without making simplifying assumptions.

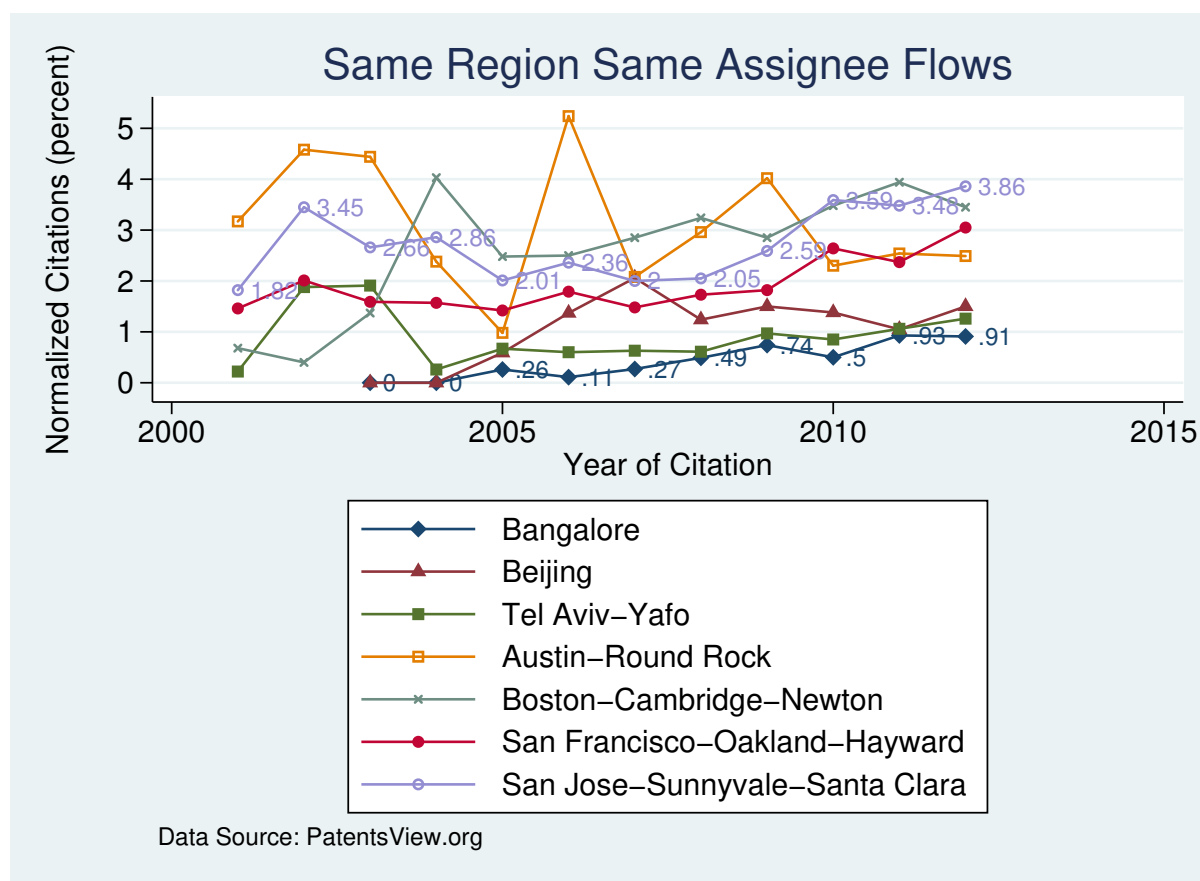


Figure 3: Local and Internal Flows by Region

### 3 Theory

Our approach toward theory building is assume the premise presented by prior scholars as valid to start, and then proceed step by step to nuance those arguments based on our arguments about the interplay of path dependence of the priors of the region, the level of complexity of the inventions produced and the relative strength of the intellectual property rights environment. Building off on [Jaffe et al. \(1993\)](#), I propose hypothesis 1 consistent with the priors supported in northern California.

*Hypothesis 1: The higher the number of local citations made by inventors from a region, the higher the number of overall citations received by inventors from*



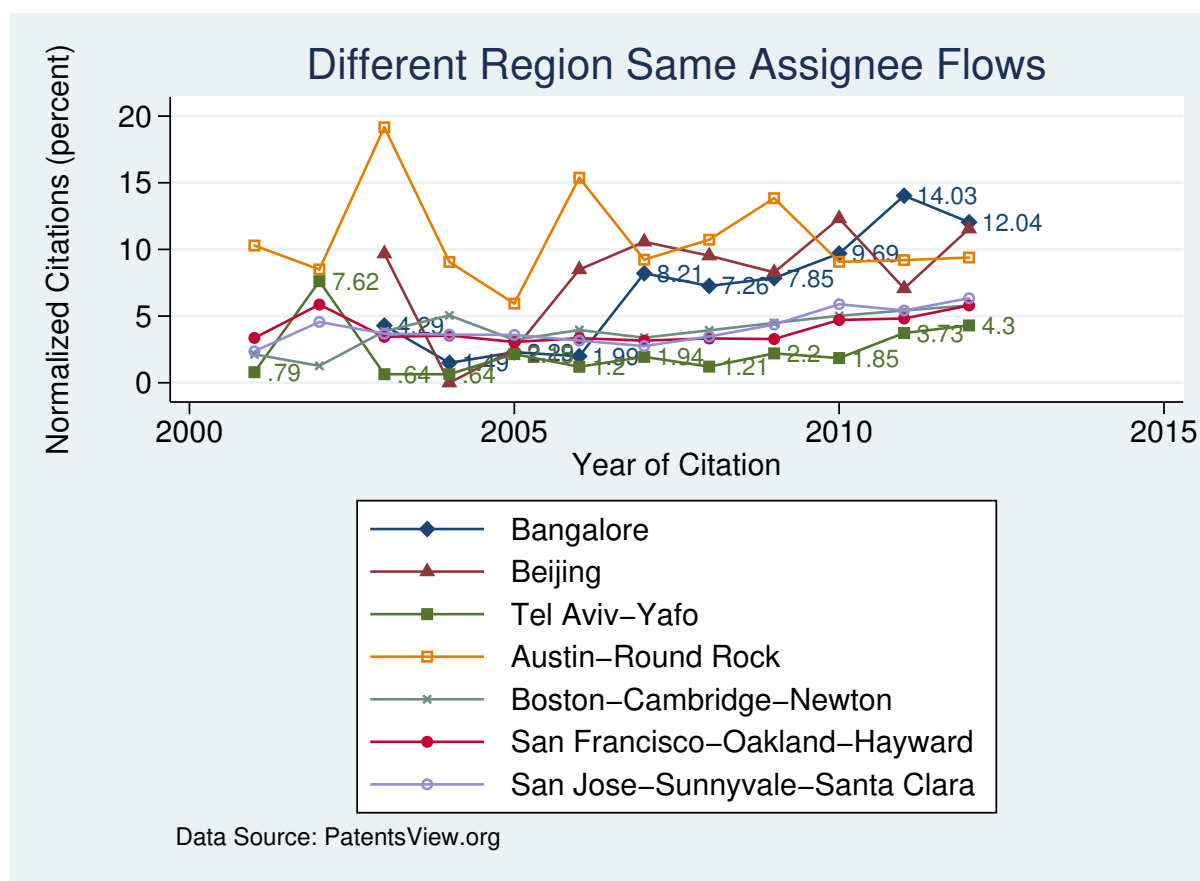


Figure 4: Non-local and Internal Flows by Region

*that region*

I then build on top of Zhao (2006) and Singh (2007) to propose hypothesis 2.

*Hypothesis 2: The number of overall citations received by inventors from a region is negatively related with the strength of the IPR environment of the region*

Scholars (Baldwin and Henkel, 2015; Yayavaram and Ahuja, 2008), have argued that increased interaction with a larger number of components creates organizational impediments to an increase in reusability of prior work. In the presence of a stronger differential in the IPR environments between inventing locations, Zhao (2006) suggests that organizational mechanisms may stand to counter the treat posed by weaker property rights. In a similar vein, I argue

that a differential in the IPR rights environment creates the organizational response to increase complexity of the inventions shared across country and IPR boundaries.

*Hypothesis 3: The number of overall citations received by inventors from a region is positively related with the complexity of the invention*

Finally, departing from prior scholarship, I suggest that path dependence would require that the inventing characteristics of locations are deeply embedded in their past. This would therefore suggest that past knowledge flow patterns should strongly predict future knowledge flow patterns.

*Hypothesis 4a: The number of overall citations received by inventors from a region in a particular period is positively related to the number of citations received by inventors from that region in prior period*

Hypothesis 4b is the logical conclusion from applying hypotheses 3 and 4a, suggesting therefore that the requirements of developing inventions of high complexity may lead firms to strengthen the prior attributes of the patenting region.

*Hypothesis 4b: The effect of Hypothesis 4a is stronger for inventions of high complexity*

## 4 Research Design

### 4.1 Complexity

I construct my measure of complexity based interactions between the different patent sub-classes. Since each of the interactions between patent sub-classes may introduce a new interaction, I model interactions on a binomial function. Specifically, when `subclass` represents the number of distinct patent sub-classes, I define `interaction(subclass)` as follows:

$$interaction(subclass) = \begin{cases} 1 & : subclass \leq 2 \\ \binom{subclass}{2} & : subclass > 2 \end{cases}$$

I would expect, from a user perspective that the more number of contexts in which the patent is valuable, the lower should be the complexity. If `complexity` represents my measure of the complexity of the patent, and `usage contexts` represents the number of distinct contexts where the patent is found valuable, I should expect the following relationship to hold:

$$Complexity \propto \frac{1}{usage\ contexts}$$

Similarly, from an inventor perspective, the more the number of contexts that the patent is built on, the higher should be the complexity. A patent that is developed without citing any other patents is an extreme case of lowest complexity, while one that requires to be built upon several `source contexts` is properly understood as being more complex.

The relationship between `source contexts` and `complexity` is therefore a normal one as depicted below.

$$\text{complexity} \propto \text{source contexts}$$

Using the principles above, I therefore develop the following definition of complexity.

$$\text{complexity} = \frac{\text{interaction}(\text{subclass}_{\text{cited}})}{\text{interaction}(\text{subclass}_{\text{patent}})}$$

By the definition above, a patent that cites no patents (and hence has  $\text{subclass}_{\text{cited}} = 0$ ) but is itself assigned to 4 sub-classes (and hence has  $\text{subclass}_{\text{patent}} = 4$ ) will have a raw Complexity score of  $\frac{1}{\binom{4}{2}} = 0.16$ . If the patent itself had been assigned onto 2 sub-classes, the raw complexity score would have been just 1. Therefore, the more the number of patent sub-classes a patent is assigned to, the lower its complexity score (by a square term). A similar but inverse relationship would hold for sub-classes arising out of cited patents. Here, I take a set union of patent sub-classes assigned to each cited patent, and use that count to determine the value of the `interaction` function.

## 4.2 IPR Classification

A review of the academic literature surrounding the construction of IPR indexes indicated that there were several, as was also evident in Zhao (2006) constructing a composite measure for the purposes of her article. Lesser (2010) provides an alternative, composite scoring system that includes the following components: protectable subject matter, membership in convention, enforcement, administration and duration of protection. I have therefore used the scores generated by Lesser (2010) for the purposes of this study. The extensive table of IPR scores has not been presented here to adhere to the page restriction, but

can be made available on request. The listing has several countries for which scores have not been provided. However none of the top patenting nations were among them, and I therefore chose to go along with this scale.

### 4.3 Data Source

I derive all patents data for this study from [patentsview.org](https://patentsview.org). The dataset considered is for all USPTO patents filed in the period 1976 to 2015. For the IPR Scores, I rely on the scores generated by [Lesser \(2010\)](#). For country definitions, I use the resources provided by [Thematic Mapping](#). To determine if spillovers are local, I use a composite data source as described in the following. For locations in the United States, it has been standard to use Metropolitan Statistical Areas (MSA) for analyses related to economic geography. Such standardized data is unavailable for non-US locations. Urban areas are a close substitute for economic centers, and I therefore determine to use one such definition for non-US locations. My data source for MSA of US locations is [the US census](#) and that for urban areas for world wide locations is [Natural Earth Data](#).

This automatically raises conflicting definitions for locations in the United States. So that the MSA definitions take precedence, I eliminated all data pertaining to US locations from the Natural Earth urban centers data and integrated this with the MSA information. With this I generated a single database of location information for economic centers around the world. A sample region definition is depicted in Figure [5](#). Here all points in the yellow region are considered to fall within the San Jose-Sunnyvale-Santa Clara, CA MSA. In Figure [6](#), I present a non-MSA example of a geographic definition based on the urban centers data

from [Natural Earth Data](#). As will be noticed in Figure 6, the Bangalore urban center is seen to include parts of Hosur as well.

#### 4.4 Unit of Analysis

The unit of analysis for this study is the `region - year`

#### 4.5 Dependent Variable

My primary dependent variable is the number of citations received by a region in a year. Being a count variable, I use a negative binomial estimation method (`xtnbreg` in Stata)

#### 4.6 Explanatory Variables

##### 4.6.1 Citations Made

My primary explanatory variables are the four counts of Citations Made along the two dimensions of same/different region and same/different assignee. While most patents have multiple inventors, and some patents also have multiple assignees, my question requires us to associate a single location to the inventor of a patent, and a single location for the assignee of the patent. For the inventor location, I tabulate the count of each of the regions that each inventor is a resident of at the time of the filing of the patent application. In doing so, I treat all inventors equally and allocate the most frequently occurring location as the

location of the inventor for that patent. In case of a tie, I assign the location of the first inventor (given by the sequence number of the inventor on the patent) as the location of the inventor of the patent.

For the assignee location, I treat multiple assignees as having been granted separate patents. I do this since the number of patents with multiple assignees is small, and so as to not lose potentially valuable information.

## **4.7 Control Variables**

### **4.7.1 Size of the patent pool, and Number of patents generated in a region - year**

Since the priors of the region may themselves explain the extent of citations received, I control for both the number of patents generated that year, as well as the aggregate pool of patents invented within that region. The reason both variables are logged is because of the exponential nature of the estimation method used.

## **5 Results**

### **5.1 Spillover Effects**

Table 1 presents the preliminary results from my negative binomial regression. The first model uses year dummies but does not include region fixed effects, while the second includes both. As evident, the effects do not seem stable and much empirical work will be needed to determine the right identification method. A couple of approaches that are planned to be tried are to look at

Table 1: Effect of Geographic Distribution of Citations Made on Citations Received

	(1) Citations Received	(2) Citations Received
Citations Received		
Citations Made to [Same Region, Same Assignee]	0.000178 (.)	-0.00393*** (-6.54)
Citations Made to [Same Region, Different Assignee]	0.0000721*** (9.14)	-0.00594*** (-7.73)
Citations Made to [Different Region, Same Assignee]	0.0000548 (1.73)	0.00600*** (9.10)
Citations Made to [Different Region, Different Assignee]	-0.0000127*** (-9.60)	0.000711*** (9.65)
Citations Made to [Other]	0.0000790*** (17.44)	-0.00124*** (-7.62)
Log (Num Patents)	-0.539*** (-33.15)	-0.0456 (-0.10)
Log (Patent Pool Size)	-0.636*** (-46.17)	-0.154 (-0.28)
Constant	8.944*** (212.68)	-43.51 (-0.00)
ln_r		
Constant	0.0000135 (0.00)	
ln_s		
Constant	-0.0000123 (.)	
Year Dummy	Yes	Yes
Region Fixed Effects	No	Yes
Observations	2624	2624

*t* statistics in parentheses\*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$



controlling for technology class, as regions may vary widely on this aspect. As indicated at the outset, the empirical aspect of this article remains work in progress.

## **6 Limitations and Looking Ahead**

I started this study attempting to understand if I could bring three mechanisms to bear together in explaining the heterogeneity in knowledge flows across regions. While there seems to be theoretical promise to exploring this question, this was a study too big to have been completed within the constraints of a term. Specifically, the endeavor has exposed me to the challenges to demonstrating empirically driven work with the objective of building theory. I intend to continue to pursue this further and integrate the IPR level data and complexity data to the flows. In addition, I plan to explore the prospect of controlling for technology classes, and see if that may lead to a strong result.

## **7 Acknowledgements**

I am greatly indebted to Pranav Garg for having emphasized the importance of picking up skills in using Stata while still in the first year. While I might have not done as much justice to it, I cannot imagine having made as much empirical progress on this project if not for that early start. I am also indebted to him for having pushed me to see the theoretical relevance and contribution of empirical work.

I am also grateful to Sai Yayavaram for having introduced me to the literature

on innovation, and for having hand held me with working on the patents data. Indeed many of the ideas underlying this article owe their origin to him. All mistakes though, remain entirely mine.

## References

- Almeida, P. and Kogut, B. (1999). Localization of knowledge and the mobility of engineers in regional networks. *Management Science*, 45(7):905–917.
- Baldwin, C. Y. and Henkel, J. (2015). Modularity and intellectual property protection. *Strategic Management Journal*, 36(11):1637–1655.
- Ethiraj, S. K. and Levinthal, D. (2004). Modularity and innovation in complex systems. *Management Science*, 50(2):159–173.
- Jaffe, A. B., Trajtenberg, M., and Henderson, R. (1993). Geographic localization of knowledge spillovers as evidenced by patent citations. *The Quarterly Journal of Economics*, 108(3):577–598.
- Lesser, W. (2010). Measuring intellectual property strength and effects: An assessment of patent scoring systems and causality. *J. Bus. Entrepreneurship & L.*, 4:345.
- Marshall, A. (2009). *Principles of Economics: Unabridged Eighth Edition*. Cosimo, Inc.
- McAfee, A. and Brynjolfsson, E. (2014). *The second machine age*. NY: WW Norton & Company.
- Singh, J. (2007). Asymmetry of knowledge spillovers between mncs and host country firms. *Journal of International Business Studies*, 38(5):764–786.
- Singh, J. and Marx, M. (2013). Geographic constraints on knowledge spillovers: Political borders vs. spatial proximity. *Management Science*, 59(9):2056–2078.
- Yayavaram, S. and Ahuja, G. (2008). Decomposability in knowledge structures and its impact on the usefulness of inventions and knowledge-base malleability. *Administrative Science Quarterly*, 53(2):333–362.
- Zhao, M. (2006). Conducting r&d in countries with weak intellectual property rights protection. *Management Science*, 52(8):1185–1199.

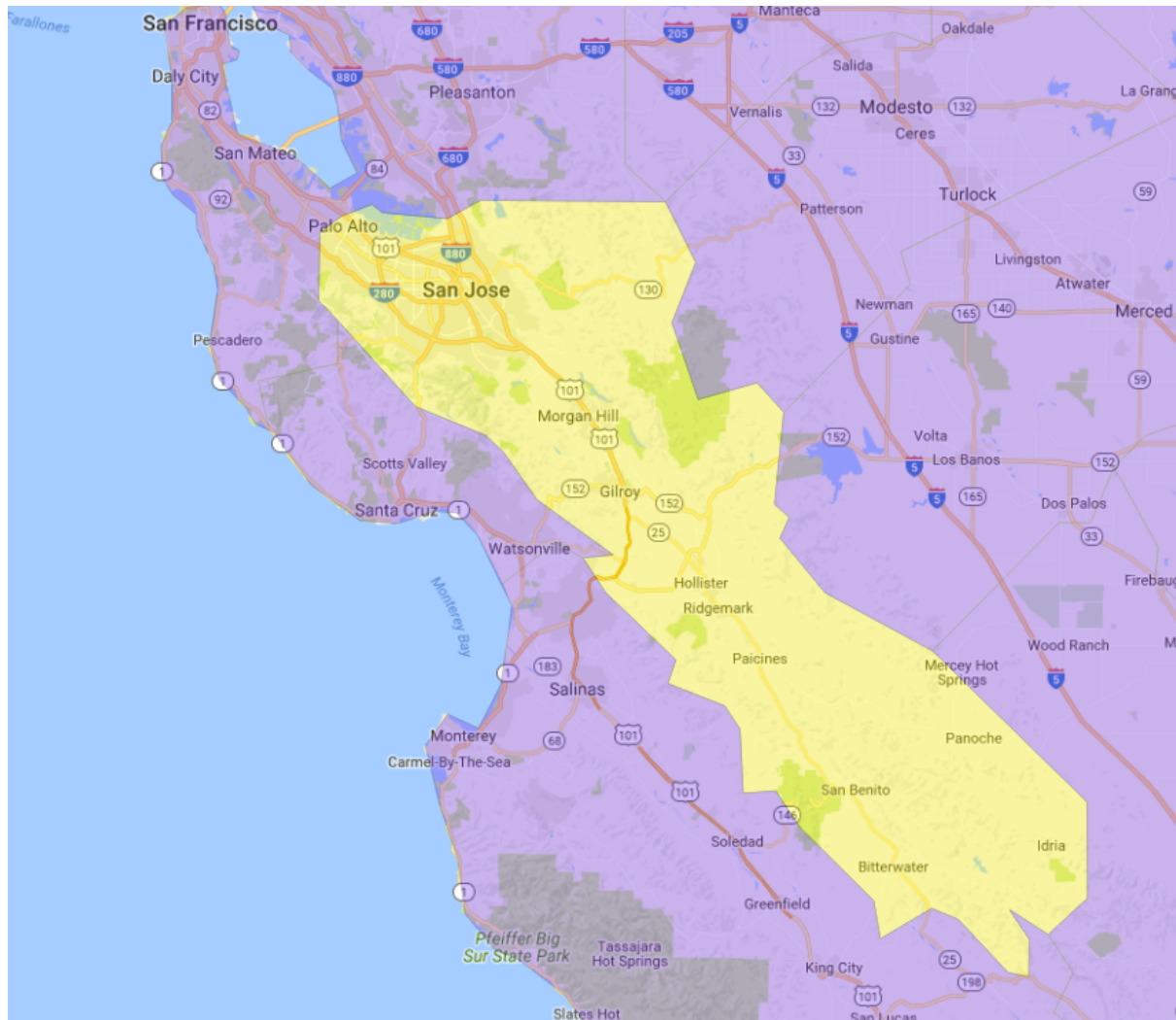


Figure 5: Geographic Definition of San Jose-Sunnyvale-Santa Clara, CA

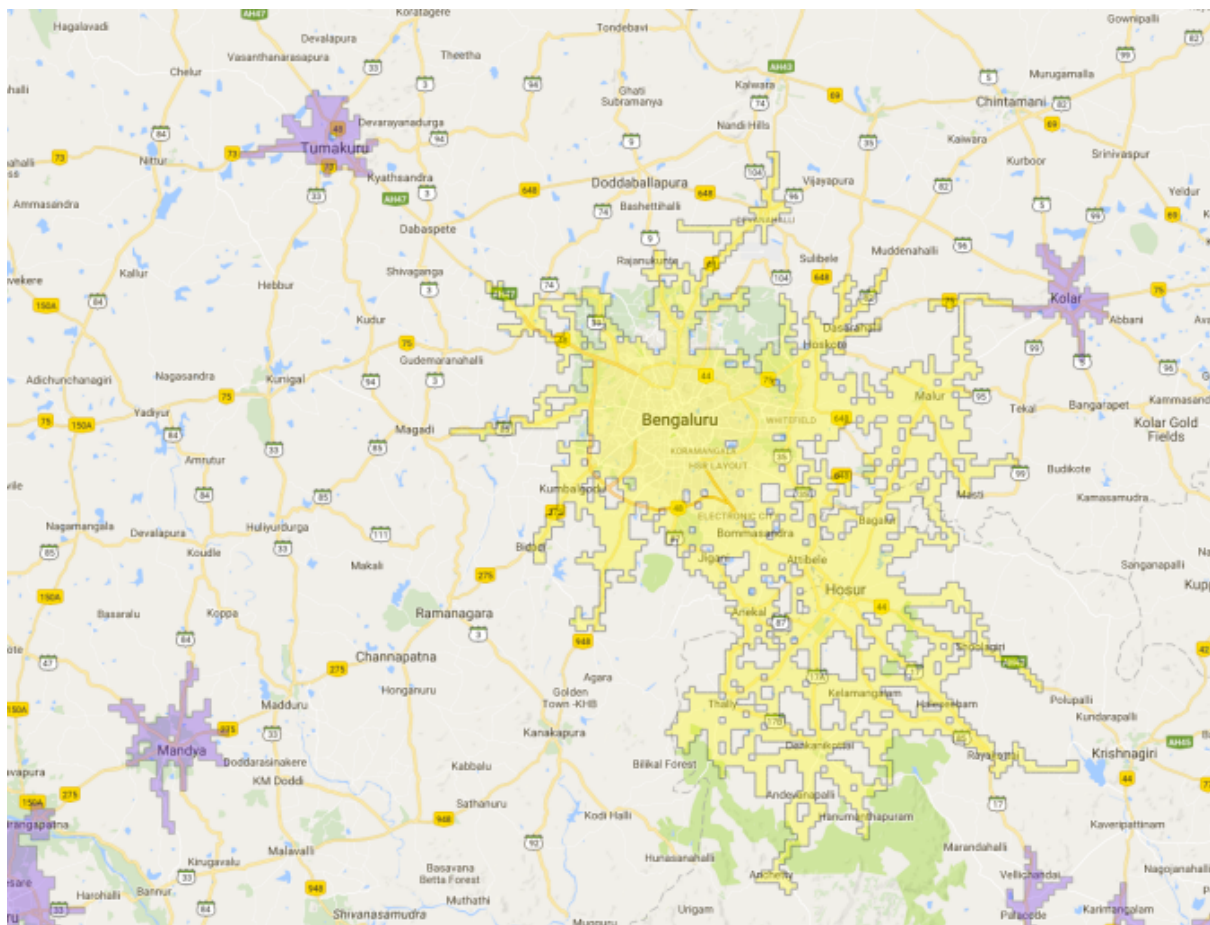


Figure 6: Geographic Definition of Bengaluru