

// PHILIPP BOEING, LOREN BRANDT, RUOCHEN DAI, KEVIN LIM, AND BETTINA PETERS

The Anatomy of Chinese Innovation: Insights on Patent Quality and Ownership





The Anatomy of Chinese Innovation:

Insights on Patent Quality and Ownership^{*}

Philipp Boeing[†] Loren Brandt[‡] Ruochen Dai [§] Kevin Lim [¶] Bettina Peters ^{\parallel}

First Version: March 2024 This Version: May 2025

Abstract

We develop a new method to measure the importance of a patent for subsequent innovation, based on the use of a Large Language Model to process patent text data and a new model of the innovation process. We apply this method to study the evolution of patenting in China from 1985-2019, also classifying patent ownership using a comprehensive business registry. Our analysis yields seven novel facts about Chinese patenting. Among these are that patenting has become narrower and less innovative over time; that knowledge within China has become more important than knowledge outside of China for directing innovative activity in China; and that knowledge produced by Chinese entities within China has been more important than knowledge produced by foreign entities filing patents in China.

Keywords: innovation, patenting, China

JEL Codes: O30

^{*}We are grateful to Dietmar Harhoff, Josh Lerner, Georg Licht, Mark Roberts, Heiwai Tang, and Xiaodong Zhu for their helpful comments and discussions. We thank participants at the HKIMR-ECB-BOFIT joint conference on Europe, Asia and the Changing Global Economy, NBER-Chinese Economy Working Group Fall Meeting 2024, 10th ZEW/MaCCI Conference on the Economics of Innovation and Patenting, and 7th Annual Bank of Canada-University of Toronto Conference on the Chinese Economy for invaluable feedback. Boeing and Peters acknowledge funding by the German Federal Ministry of Education and Research under grant number 01DO21006A. Brandt thanks funding support from the Noranda Chair at the University of Toronto. Dai thanks funding support from the National Natural Science Foundation of China (grant numbers 72342031 and 72203252). The responsibility for the content of this study lies with the authors.

[†]Goethe University Frankfurt, Frankfurt am Main, and ZEW – Leibniz Centre for European Economic Research, Mannheim (philipp.boeing@zew.de).

[‡]University of Toronto (loren.brandt@utoronto.ca).

[§]Central University of Finance and Economics (r.dai@cufe.edu.cn).

[¶]University of Toronto (kvn.lim@utoronto.ca).

^IZEW – Leibniz Centre for European Economic Research, Mannheim (bettina.peters@zew.de).

1 Introduction

Patents are an important indicator of technological change and innovation (Griliches (1990)). In China, patenting activity has accelerated substantially over the last two decades, far outpacing growth in countries like the US. Figure 1.1 shows the number of invention patents applied for at the China National Intellectual Property Administration (CNIPA) in each year from 1985 to 2018 and the associated annual growth rates.¹ From 2000 to 2018, for example, the number of invention patents granted by the China National Intellectual Property Administration (CNIPA) grew at an average annual rate of 23.5%, in contrast with an average growth rate of 4.3% at the US Patent and Trademark Office (USPTO). In 2015, the number of patents granted at the CNIPA overtook the number of patents granted at the USPTO for the first time. A rapidly growing literature has emerged trying to explain this acceleration and its implications.² In this paper, we investigate how the quality of Chinese patenting has changed over time and quantify which sources of knowledge have been important for driving innovation in China.

Our analysis requires overcoming three important challenges. First, patents represent a vast and important source of codified knowledge, but one that is also difficult to incorporate systematically into quantitative analyses of patenting growth and quality because it is almost exclusively in the form of text data.³ To make progress on this front, we leverage the capabilities of an industry-leading Large Language Model (LLM) to generate embeddings for the abstracts of every Chinese patent filed at the CNIPA. Embeddings are vectors that represent the meaning of text in a vector space, where the mapping from text to vector is determined by an LLM that has been pre-trained on vast quantities of text data in multiple languages. These vector representations allow us to compute well-defined measures of similarity between patents (e.g., the cosine similarity between the patent text embeddings), which in turn offers new ways to integrate patent text data into an analysis of innovation.

The second challenge is that there is little existing theory about how patent text embeddings should be used to measure patent quality. To guide our analysis, we develop a simple model of innovation where the patent text embeddings are drawn from a dynamic stochastic process. The key idea is that the mean of the embedding distribution in the future is determined by the embeddings of patents in the past. This simple structure leads to novel insights about how to estimate the importance of a patent for innovation. Intuitively, we can now measure precisely the direction of innovation (i.e., the average direction in which patent embeddings are moving over time). Past patents whose embeddings are more similar to the direction of innovation are then inferred to be more important for the innovation process. As we show below, this new measure of patent importance is not only positively correlated with traditional measures of patent quality (e.g., citations and grant rates) but is also a stronger predictor of total

¹Invention patents in China are essentially equivalent to utility patents in the US.

²For discussions on the institutional reforms in China leading to an increase in patenting activity and the link between patenting and firm outcomes see Jefferson et al. (2003), Hu and Jefferson (2009), Li (2012), Boeing (2016), and Boeing et al. (2016). A more recent literature has focused on the possibly distorting effects of policy on the incentive to patent and its effect on patent quality. On this point see Dang and Motohashi (2015) and Hu et al. (2017).

 $^{^{3}}$ Earlier efforts based on counts of keywords include Younge and Kuhn (2016), Kelly et al. (2021), and Kalyani (2024) and in the Chinese context Fang et al. (2021)



Figure 1.1: Number and growth rate of patents in China

Notes: Panel (a) shows the number of invention patents applied for in each year at the CNIPA. Panel (b) shows the associated annual growth rates.

factor productivity (TFP) and output at the firm level than traditional quality measures. Our approach thus provides a new metric to assess the importance of different sources of knowledge for Chinese innovation. For instance, we use this to compare the importance of patents in China versus patents outside of China (filed at the USPTO).

The third and final challenge is that standard patent data provide very little information about patentees beyond basic names and addresses. Hence, even if one can precisely measure the importance of a patent, it remains challenging to determine which types of patentees have been producing more important patents. To address this problem, we leverage information on ownership of registered capital from a comprehensive business registry from China. Combining this with basic patent information allows us to differentiate between at least eight types of patenting entities: (i) privately-invested enterprises (PIEs); (ii) state-owned enterprises (SOEs); (iii) foreign-invested enterprises (FIEs); (iv) Chinese universities; (v) Chinese research institutes; (vi) individuals; (vii) other domestic patentees; and (viii) overseas patentees (those with an address outside of China). We can then compare the growth, composition, and importance of patenting across these patentee types.

Our analysis leads to seven novel facts about Chinese patenting. First, patents that are important for innovation have become less important on average, which is due in part to the fact that patenting in China has become less innovative overall and more crowded over time. This is consistent with other studies that have documented similar slowdowns in various measures of innovative activity in the US (for example, see Kelly et al. (2021)), Kalyani (2024), Bloom et al. (2020), Akcigit and Ates (2023), and Covarrubias et al. (2019)). Second, there is substantial heterogeneity across technologies in terms of how patent importance has changed over time, with technologies related to health, chemistry, and physics being the most dynamic. Third, the number of breakthrough (high-importance) patents has continued to grow in China, but growth rates have been falling since the mid-2000s. Fourth, knowledge within China has become more important than knowledge outside China for driving the direction of innovation in China. Fifth, knowledge produced by Chinese entities (PIEs, SOEs, universities, and research institutes) has been more important for Chinese innovation than knowledge produced by foreign entities (FIEs and overseas patentees). Sixth, new patentees have typically produced more important patents than experienced patentees, but this premium has dissipated steadily since the mid-2000s. Finally, patenting in China has become narrower, with the influence of knowledge in other related technologies on innovation diminishing over time.

The remainder of this paper is organized as follows. Section 2 describes the key sources of data used in our analysis, explains how we classify patentee types using novel information from the Chinese business registry, and discusses the generation of patent text embeddings. Section 3 develops our model of innovation and shows how the theory leads to a new measure of patent importance based on the use of text embeddings. Section 4 describes how we take the model to the data and section 5 provides validation of our new importance measure by showing how it correlates with patent citations and grant rates, as well as how it predicts TFP and output at the firm level. Section 6 then presents our main findings and section 7 concludes.

2 Data: basic patent statistics and text embeddings

We study an administrative database covering all invention patents applied for at the CNIPA from 1985 to 2019.⁴ This includes Patent Cooperation Treaty (PCT) invention patents, which are filed almost exclusively by overseas applicants.⁵ For brevity, we will henceforth refer to an invention patent application simply as a "patent".⁶

2.1 Classifying patentee types

Patent application documents provide information about the name of the patent applicant(s) and the address of the main applicant. Beyond this, little is known about who is responsible for patenting activity. For example, is the applicant of a patent an enterprise or a research institute? If it is an enterprise, who owns the firm? We define a taxonomy of what we refer to as the "patentee type" of a patent based on the intersection of three sources of information.

First, we determine from the address of the main patent applicant whether the *location* of the main patent applicant is domestic (reporting an address in China) or overseas (reporting an address outside of China). Second, for patents where the main applicant is domestic, we conduct a keyword search on the name of each applicant in a patent to identify the *entity*

⁴We obtain these data directly from CNIPA, which helps explain why they are more comprehensive and accurate than other databases for Chinese patents, (e.g., downloaded from the CNIPA website). We restrict attention to invention patents, since these are typically considered to be the most innovative. On average, 40% of annual patent applications are invention patent applications, with utility patent applications making up 40% and design patent applications the remaining 20%.

⁵China became a Patent Cooperation Treaty (PCT) contracting state on 1 January 1994.

⁶Several patent offices provide online resources on the Chinese patent system, including patenting procedures and fees (CNIPA), patent examination and laws/policies (WIPO), China's patent numbering system (EPO), and an IPR toolkit for China (USPTO).



Figure 2.1: Patent shares by patentee type

type of the applicant. This allows us to differentiate between enterprises, universities, research institutes, individuals, and other domestic patentees (e.g., military). Third, we merge the data for domestic enterprise patents with a comprehensive registry of businesses in China. This allows us to observe ownership of registered capital for each entity in the business registry and hence to determine the *ownership type* of each enterprise applicant.⁷ Finally, we combine the information about location, entity type, and ownership type to define eight mutually exclusive *patentee types*: (i) privately-invested enterprises (PIEs); (ii) state-owned enterprises (SOEs); (iii) foreign-invested enterprises (FIEs); (iv) universities; (v) research institutes; (vi) individuals; (vii) other domestic patentees; and (viii) overseas patentees. Online Appendix A provides a detailed explanation of this classification procedure.

This taxonomy allows us to shed light on the actors responsible for patenting activity in China at a level of detail that has not previously been possible to examine. For instance, Figure 2.1 shows the share of patents applied for in each year by patentee type. Before the mid-2000s, overseas patentees account for more than half of all patent applications in almost every year. However, from the early 2000s onward, there is rapid growth in the patenting activity by PIEs, which quickly become the dominant patentee type in terms of patent shares. By the late 2010s, PIEs account for around 40% of patent applications, whereas overseas patentees account for only around 10%. In contrast, the share of overseas patentees among granted USPTO utility patents actually increases from around 45% to 55% during the same time period (see Figure 1 in Online Appendix E for details). We also observe that the role of Chinese universities and SOE patentees increases slightly over time.

⁷Of all the invention patents that we identify as having an "enterprise" entity type, we are able to identify the ownership type of the patent from the business registry in around 88% of cases, while the remaining 12% cannot be matched to the business registry and are classified as "Other Domestic". Hence, the match between the patent data and the business registry is of reasonably high quality.

2.2 Technology classes

The most basic source of information about the knowledge that a patent represents is provided by a set of technology class codes reported by each patent, referred to as the International Patent Classification (IPC) codes. For example, IPC code "A01B 1/00" indicates that a patent is related to human necessities ("A"), specifically agriculture and forestry ("01"), more specifically soil working ("B"), and even more specifically hand tools ("1/00"). A patent can and often does report more than one IPC code but the main IPC code of each patent is also observed.

Figure 2.2 shows the share of patents applied for in each year by the section (first character) of each patent's main IPC code. Note that the technological composition of Chinese patents has been fairly stable over time. The share of patents by IPC section in the average year, listed in descending order, is as follows: G (Physics) (20.3%), B (Performing Operations; Transporting) (17.8%), H (Electricity) (17.2%), A (Human Necessities) (16.2%), C (Chemistry; Metallurgy) (15.4%), F (Mechanical Engineering) (7.6%), E (Fixed Constructions) (3.9%), and D (Textiles; Paper) (1.7%).

2.3 Other basic patent data features

Online Appendix B provides more detailed information about patent counts and shares by product versus process classification (based on the text of the claims of a patent), by location within China for domestic applicants, and shares by country for foreign applicants. We briefly mention several observations. First, just under half of all patents are product patents, a third are mixed (product and process), and the remaining are process patents. Second, domestic patent applications are highly concentrated in the coastal regions of China, with this concentration increasing over time, and the role played by provinces in central China also increasing. Third, applications from the top five overseas locations – Japan, USA, Germany, South Korea, and Taiwan – represent around 80% of all overseas patent applications after 2000.

2.4 Patent text embeddings

A patent's IPC codes are a coarse indicator of the knowledge that the patent embodies. Even at the finest level of disaggregation (e.g., "hand tools"), it is difficult if not impossible to discern from the IPC codes alone what the patent is really about. This information is contained, of course, in the textual descriptions of the patent – specifically, in the patent's abstract (a summary of the key elements of the patent) and claims (a detailed description of how the patent is claiming to be innovative). The text of a patent is hence a key source of codified knowledge, but this has traditionally been difficult to incorporate into quantitative analyses of patenting activity.

We overcome this problem by utilizing an industry-leading LLM to generate *text embeddings* for the abstract of every patent in the CNIPA database. Each embedding is a vector that represents the meaning of the text, where the mapping from text to vector is determined by the LLM based on vast quantities of training data (e.g., all text that is publicly available on the internet). The model that we use specifically is the *embed-multilingual-v2.0* model from Cohere, a Canadian technology company that specializes in natural language processing (NLP)



Figure 2.2: Patent shares by IPC section

and LLMs.⁸ The model represents each text with a 768-element vector, thereby reducing the extremely high-dimensional text data to a smaller but still high-dimensional object. In practice, we generate embeddings using the abstract for each patent. To draw comparisons with patenting activity outside of China, we also generate embeddings for all granted utility patents at the USPTO.⁹

To visualize the text embeddings of the Chinese patent abstracts, we first reduce each 768-vector to two dimensions using a process known as t-Distributed Stochastic Neighbor Embedding (t-SNE), which is a technique for dimensionality reduction that is particularly well-suited for visualizing high-dimensional data. Figure 2.3 shows heatmaps of the two-dimensional t-SNE reduction of each abstract embedding by the main IPC section for each CNIPA patent from 1985 to 2000.¹⁰ Evidently, patents belonging to the same IPC section tend to have abstract embeddings that exist in similar areas of the reduced vector space. For example, the embeddings for patents in sections G (Physics) and H (Electricity) – technology classes which are arguably more similar to each other than to other classes – tend to be clustered in the southwest corner of the two-dimensional space, whereas the embeddings for patents in section B (Performing Operations; Transporting) tend to be clustered in the southeast corner. This observed clustering based on IPC section is indicative that the embeddings of the patent abstracts are picking up differences and similarities in semantic meaning across patents in different technology classes.

⁸See https://txt.cohere.com/multilingual/ for a user-friendly introduction to the Cohere multilingual model.

 $^{^{9}}$ We focus on granted patents for the USPTO because data for these patents are more readily available. For instance, data for patent grants are available as far as back as 1976, whereas application documents are only available from the early 2000s onward.

¹⁰The t-SNE reduction of an vector depends on the set of all vectors that are being reduced. Since the computation time required for the t-SNE decomposition scales quickly with the number of vectors being reduced, we show this figure only for patents in the earlier years of our data.





Notes: This figure shows the two-dimensional t-SNE reductions of the abstract embeddings for all CNIPA invention patents from 1985-2000 by the main IPC section of each patent.

It is also important to emphasize that the use of embeddings to measure the "meaning" of a patent is fundamentally different from approaches to textual analysis that are based on counts of keywords, such as the Term Frequency-Inverse Document Frequency (TF-IDF) method (e.g., as used by Younge and Kuhn (2016) and Kelly et al. (2021)). To illustrate, consider two hypothetical scenarios. In the first, all patents in the past contain only the key phrase "automobile", while those in the future contain only the keyword "electric vehicle". In the second, all patents in the past contain only the key word "electric vehicle". In the future again contain only the keyword "electric vehicle". It is obvious that there is a greater leap in innovation in the second scenario compared with the first. However, keyword-based methods would not be able to discern this – one can only tell that the keywords in the past and future are different, but not *how different* they are. In contrast, the embeddings for "automobile", "electric vehicle", and "horse carriage" can be used to measure precisely the "distance" between these keywords.

In the context of text embeddings, a widely-used measure of (inverse) distance is cosine similarity, which is defined as follows. Consider two patents, i and j, and let $\vec{x}_i \equiv \left[x_i^1 \cdots x_i^K\right]'$ and $\vec{x}_j \equiv \left[x_j^1 \cdots x_j^K\right]'$ denote the embeddings of the abstracts for these two patents (where K = 768 is the dimension of each embedding). The cosine similarity of \vec{x}_i and \vec{x}_j is:

$$c\left(\vec{x}_{i}, \vec{x}_{j}\right) = \frac{\vec{x}_{i} \cdot \vec{x}_{j}}{||\vec{x}_{i}||||\vec{x}_{j}||} = \frac{\sum_{k=1}^{K} \vec{x}_{i}^{k} \vec{x}_{j}^{k}}{\left[\sum_{k=1}^{K} \left(\vec{x}_{i}^{k}\right)^{2}\right]^{\frac{1}{2}} \left[\sum_{k=1}^{K} \left(\vec{x}_{j}^{k}\right)^{2}\right]^{\frac{1}{2}}}$$
(2.1)

The cosine similarity of any two vectors lies in the interval [-1,1] and is closer to 1 when the

two vectors are pointing in more similar directions in the vector space.¹¹ This metric has the advantage of being invariant to the magnitudes of the embeddings, which tends to reflect the length of the text being embedded (in practice, we simply normalize each embedding \vec{x}_i to have unit length, so that $c(\vec{x}_i, \vec{x}_j)$ is simply the dot product of \vec{x}_i and \vec{x}_j). Now, for example, using the Cohere *embed-multilingual-v2.0* model, we find that c ("horse carriage", "electric vehicle") = 0.90 whereas c ("automobile", "electric vehicle") = 0.93, indicating that the model is able to recognize how different the keywords are. Note that the standard deviation of cosine similarity scores between pairs of patent abstract embeddings in our data is on the order of 0.01, so a difference of 0.03 is quite large.

3 A theory of patent importance

We now turn toward a central challenge: measuring the quality of a patent. Researchers studying patents have explored many different measures of patent quality, for example, the number of forward citations received by a patent (Hall et al. (2005), Kuhn et al. (2020)); the centrality of a patent in the citation network (Funk and Owen-Smith (2017), Park et al. (2023)); legal status changes (e.g., grants, unpaid renewal fees) (Hedge et al. (2023)); the timing of legal status changes (e.g., the time taken for a patent to be granted after application) (Chondrakis et al. (2021)); the number and length of patent claims (Marco et al. (2019)); and the existence of related patent filings at overseas patent offices (Harhoff et al. (2003)). Online Appendix C documents our findings about how each of these measures have evolved over time in China. Note, however, that none of these measures make direct use of the semantic content of patent text. Furthermore, the "quality" of a patent can mean many things: is a patent good for productivity, for example, or is it good for developing new products?

Since our interest in this paper is to investigate the sources of knowledge that drive innovation in China, we first focus on a specific notion of patent quality: a high-quality patent is one that is *important* for the direction of innovation. Note that important patents thus defined may not necessarily be the same as patents that are good for other economic outcomes, although below we will show that the measure of patent quality that we develop here is positively correlated with measures such as patent citations, firm TFP, and output.

3.1 A naive approach

How can patent text embeddings be used to determine which patents are important for innovation? A common conception in the literature is that a patent that is important for innovation should have two features: it should be similar to future patents ("impact") and different from past patents ("novelty"). Forward citations are by far the most widely used measure of impact, based on the assumption that if one patent cites another then the cited patent must have had some influence on the citing patent. Similarly, others have proposed that a lack of backward citations indicates that a patent is dissimilar to the past and hence is indicative of novelty (e.g., the so-called "radicalness" measure used by Ahuja and Lampert (2011) and Banerjee and Cole (2011)). More recently, Kelly et al. (2021) have also proposed

¹¹When K = 2, $c(\vec{x}_i, \vec{x}_j)$ is the cosine of the angle between the vectors \vec{x}_i and \vec{x}_j , hence the name.

that difference in the similarity of a patent's keywords to the future relative to the past is indicative of the patent's importance for innovation.

A candidate measure of the importance of a patent i applied for at time t based on the patent text embeddings would thus be:

$$\tilde{p}_{it} = \underbrace{c\left(\vec{x}_{it}, \vec{\mu}_{Ft}\right)}_{F_{it}} - \underbrace{c\left(\vec{x}_{it}, \vec{\mu}_{Bt}\right)}_{B_{it}} \tag{3.1}$$

where \vec{x}_{it} is the embedding of patent i, $\vec{\mu}_{Ft}$ is an average embedding of some patents applied for after t, $\vec{\mu}_{Bt}$ is an average embedding of some patents applied for at or before t, and we let F_{it} and B_{it} denote the "forward" and "backward" cosine similarity of patent i, respectively.

There are two major limitations to measuring patent importance in this way. First, there is no theoretical underpinning for why this one particular transformation of the information contained in the embeddings is the best measure of importance. For example, why should F_{it} and B_{it} be equally weighted? Second, this approach evaluates the importance of a patent without regard for the existence of other contemporaneous patents that might also be important for innovation. Arguably, the set of competing ideas needs to be considered jointly to determine which particular ideas win out in shaping the future. Yet, without additional structure, it is not exactly clear how one should do this. Hence, we next develop a simple theory of patent importance.

3.2 The innovation process

Consider patenting in a single IPC class. We begin with the assumption that the embedding of each patent applied for in the class at time t is a random vector drawn from some distribution with mean $\vec{\mu}_t$ (a $K \times 1$ column vector, where K is the dimension of the embeddings) and support on the (K - 1)-sphere (e.g., for K = 2, the support of the distribution is the unit circle). An example of such a distribution is the von-Mises Fisher distribution in \mathbb{R}^K , which is a multivariate normal distribution conditional on the associated random vectors having unit length, although the exact functional form of the distribution will play no role in what follows.

We refer to $\vec{\mu}_t$ as the state of knowledge at time t, which captures the current location of patenting activity in the embedding space. Our goal is to understand how this state changes over time and to quantify the importance of each patent for such changes. To introduce the key ideas of the theory, we start with a parsimonious model of the innovation process, which we will enrich below before connecting the theory with the empirical analysis. Suppose then that the state of knowledge at t + 1 is determined as follows:

$$\vec{\mu}_{t+1} = \rho_t \vec{\mu}_t + \sum_{i \in \Omega_t} p_{it} \vec{x}_{it} + \vec{\epsilon}_t \tag{3.2}$$

where \vec{x}_{it} is the embedding of patent *i* applied for at time *t*, Ω_t is the set of patents applied for at *t* (with $N_t \equiv |\Omega_t|$), and $\vec{\epsilon}_t$ is a mean-zero error vector that is orthogonal to both $\vec{\mu}_t$ and $\{\vec{x}_{it}\}_{i\in\Omega_t}$. The future state of knowledge thus depends on its current state as well as on the ideas that are embodied in new patent applications. We refer to ρ_t as the *memory* of the innovation process, which is a standard AR(1) coefficient except that it is allowed to vary over time. Memory allows the knowledge embodied in current and past patents to have an effect on the state of knowledge at any time in the future.¹² We also refer to p_{it} as the *importance* of patent *i*, which will be our key measure of how important a patent is for innovation.

3.3 Patent importance

Since the error vector $\vec{\epsilon}_t$ is assumed to be orthogonal to the patent embeddings $\{\vec{x}_{it}\}_{i\in\Omega_t}$, the moment condition $\mathbb{E}\left[X'_t\vec{\epsilon}_t\right] = 0$ must hold, where X_t is a $K \times N_t$ matrix with column *i* equal to \vec{x}_{it} . The sample analogue of this moment condition can be written as:

$$X_{t}'X_{t}\vec{p}_{t}' = X_{t}'\underbrace{(\vec{\mu}_{t+1} - \rho_{t}\vec{\mu}_{t})}_{\vec{d}_{t+1}}$$
(3.3)

where \vec{p}_t is the $1 \times N_t$ importance vector with element *i* equal to p_{it} . We define $\vec{d}_{t+1} \equiv \vec{\mu}_{t+1} - \rho_t \vec{\mu}_t$ as the *direction of innovation*, which captures the way in which the future is moving away from the present. Note that equation (3.3) is simply an ordinary least squares (OLS) moment condition, where \vec{d}_{t+1} is the "dependent variable", each patent embedding \vec{x}_{it} is an "independent variable" with coefficient p_{it} , and the "observations" are the *K* elements of the embeddings.

Now, note that since each patent embedding is normalized to have unit length, the matrix $X'_t X_t$ is simply the matrix of cosine similarities between the embeddings of patents applied for at date t. Furthermore, $X'_t \vec{\mu}$ for any unit vector $\vec{\mu}$ is simply a vector of cosine similarities between the patent embeddings in X_t and μ . Hence, equation (3.3) can also be expressed as:

$$C_t \vec{p}_t' = \vec{F}_t' - \rho_t \vec{B}_t' \tag{3.4}$$

where C_t is the $N_t \times N_t$ cosine similarity matrix with (i, j)-element equal to $c(\vec{x}_{it}, \vec{x}_{jt})$, \vec{F}_t is the $1 \times N_t$ forward similarity vector with element *i* equal to $c(\vec{x}_{it}, \vec{\mu}_{t+1})$, and \vec{B}_t is the $1 \times N_t$ backward similarity vector with element *i* equal to $c(\vec{x}_{it}, \vec{\mu}_t)$. The *i*th row of equation (3.4) in particular is:

$$p_{it} + \sum_{j \in \Omega_t \setminus \{i\}} c_{ijt} p_{jt} = F_{it} - \rho_t B_{it}$$

$$(3.5)$$

where $c_{ijt} \equiv c(\vec{x}_{it}, \vec{x}_{jt})$ is shorthand notation for the cosine similarity between patents *i* and *j*.

Equation (3.5) implies that patent i can be more similar to the direction of innovation if either patent i is more important or it is more similar to other important patents. In contrast with the "naive" measure of importance posited in equation (3.1), our measure now has two advantages. First, the relative weight on forward similarity F_{it} versus backward similarity B_{it} can be interpreted as depending on the memory of the innovation process (which we will estimate in the empirical implementation below). Second, equation (3.5) shows how it can be misleading to infer the importance of a patent based only on its own forward and backward

¹²Given that the right-hand side of equation (3.2) already includes the embeddings of patents at time t, one may wonder why we also include the state of knowledge at time t instead of at t-1. The reason for this modeling choice is the following. First note that a patent embedding \vec{x}_{it} at time t has a direct effect on the future state of knowledge $\vec{\mu}_{t+1}$ as long as $p_{it} \neq 0$. If we were to replace $\vec{\mu}_t$ with $\vec{\mu}_{t-1}$ on the right-hand side of equation (3.2), any change in $\vec{\mu}_{t+1}$ due to \vec{x}_{it} will have no effect on $\vec{\mu}_{t+2}$ but will have a non-zero effect on $\vec{\mu}_{t+3}$. To avoid these "holes" in the dynamic process, we include $\vec{\mu}_t$ on the right-hand side of equation (3.2).

similarity scores: a patent may look like it has high impact and novelty but may really just be similar to other patents that are truly impactful and novel. Our measure adjusts for this by taking into account the network of cosine similarities between patents.

To obtain a unique solution for the importance of each patent from equation (3.3), a necessary and sufficient condition is that the cosine similarity matrix $X_t'X_t$ is invertible. As we discuss below, invertibility will be a salient issue in taking the model to the data, but let us assume for now that invertibility holds in order to develop more intuition about what patent importance represents. With invertibility, the unique solution for the importance vector is:

$$\vec{p}_{t} = \left(X_{t}'X_{t}\right)^{-1}X_{t}'\vec{d}_{t+1}$$
(3.6)

Hence, the direction of innovation predicted by the patent embeddings X_t is:

$$\vec{d}_{t+1} = X_t \vec{p}_t^{\prime} \tag{3.7}$$

$$= X_t \left(X_t' X_t \right)^{-1} X_t' \vec{d}_{t+1}$$
(3.8)

Now, note that $X_t (X'_t X_t)^{-1} X'_t$ is the projection matrix for the subspace spanned by the columns of X_t , i.e., the embeddings of the patents applied for at t. $\hat{\vec{d}}_{t+1}$ is thus the projection of the direction of innovation into this subspace and the importance scores $\{p_{it}\}_{i\in\Omega_t}$ are simply the unique set of weights on each patent embedding that allow us to construct this projection.

To illustrate, consider Figure 3.1. In panel (a), we show an example where K = 2 and there are exactly two patent embeddings, \vec{x}_{1t} and \vec{x}_{2t} . In this case, the direction of innovation \vec{d}_{t+1} already lies in the subspace spanned by \vec{x}_{1t} and \vec{x}_{2t} . Hence, the projection $\hat{\vec{d}}_{t+1}$ is simply \vec{d}_{t+1} itself and the magnitude of the error vector $\vec{\epsilon}_t$ is zero. It is clear that in this example, there is a unique pair of weights p_{1t} and p_{2t} that allows us to construct \vec{d}_{t+1} using \vec{x}_{1t} and \vec{x}_{2t} . Furthermore, note that the direction of innovation is pointing toward \vec{x}_{1t} and away from \vec{x}_{2t} . Hence, patent 1 is assigned a positive importance score $(p_{1t} > 0)$ and patent 2 is assigned a negative importance score $(p_{2t} < 0)$.

In panel (b), we show a second example where K = 3 and we again have two patent embeddings, \vec{x}_{1t} and \vec{x}_{2t} . In this case, the dimension of the embedding space is greater than the number of patents and it is possible for the direction of innovation not to lie in the subspace spanned by \vec{x}_{1t} and \vec{x}_{2t} . The projection $\hat{\vec{d}}_{t+1}$, however, is the closest vector in this subspace to \vec{d}_{t+1} and there is again a unique pair of weights p_{1t} and p_{2t} that allows us to construct $\hat{\vec{d}}_{t+1}$ using \vec{x}_{1t} and \vec{x}_{2t} . The difference between \vec{d}_{t+1} and $\hat{\vec{d}}_{t+1}$ is the error vector $\vec{\epsilon}_t$, which is orthogonal to the subspace spanned by \vec{x}_{1t} and \vec{x}_{2t} by assumption. In sum, one can think of the patent importance scores as weights on each patent that allow us to get "as close as possible" to the direction of innovation using the patent embeddings.

Finally, note that average value of patent importance within a given year is not a very meaningful statistic. To see why this is the case, suppose that the distribution of \vec{x}_{it} is symmetric around its mean $\vec{\mu}_t$ (as it is in the von-Mises Fisher distribution, for example). In two dimensions, the distribution of \vec{x}_{it} may then resemble the first diagram shown in Figure 3.2 below, where each patent embedding is represented by a square or circle. Now suppose that





Figure 3.2: Patent importance as a relative measure



the the state of knowledge $\vec{\mu}_t$ were to change to some new state $\vec{\mu}_{t+1}$, for example by rotating clockwise as illustrated in the second diagram of Figure 3.2. Suppose also for exposition that $\rho_t = 1$ so that the direction of innovation \vec{d}_t is simply the difference between $\vec{\mu}_{t+1}$ and $\vec{\mu}_t$. Note that the direction of innovation points in the same direction as exactly half of the patent embeddings (the circles in Figure 3.2) and in the opposite direction as exactly half of the patent embeddings (the squares). Hence, the former group of patents would be assigned positive importance scores while the latter group of patents would be assigned negative importance scores. Furthermore, this is true regardless of the way in which the state of knowledge moves. As such, the average importance of patents in a given IPC-year will typically be close to zero and hence is an uninteresting statistic. Instead, what will matter is which patents within a given IPC-year are important and which are not.

3.4 Allowing for interactions across IPC classes and foreign influence

The version of the model described above makes two strong assumptions: that innovation within each IPC class occurs in isolation and that only knowledge in China is relevant for innovation in China. There are at least two features of our data, however, that strongly suggest these to be poor assumptions. First, as described above, patents often report secondary IPC codes besides their main IPC code, indicating that these other technology classes are also relevant for the knowledge embodied in a patent. Second, from the information about backward citations provided by each patent application document, we often see patents in one IPC class citing patents in other IPC classes. Similarly, we often see Chinese patents citing patents filed at foreign patent offices as well. Hence, we will now extend the model to relax both of these assumptions.

Let us first explicitly index a patent's main IPC class by g. Now consider innovation in IPC class g and suppose that instead of equation (3.2), the innovation process is described by:

$$\vec{\mu}_{t+1}^{g} = \rho_t^g \vec{\mu}_t^g + \sum_{i \in \Omega_t^g} p_{it}^g \vec{x}_{it}^g + \underbrace{\sum_{g' \in \Gamma_g} \gamma_t^{gg'} \vec{\mu}_t^{g'}}_{\text{IPC interactions}} + \underbrace{p_{US,t}^g \vec{\mu}_{US,t}^g}_{\text{foreign influence}} + \vec{\epsilon}_t^g$$
(3.9)

Note that we now allow the state of knowledge in other IPCs, $\vec{\mu}_t^{g'}$, to play a role, where Γ_g is the set of *related IPCs* that matter for innovation in g and $\gamma_t^{gg'}$ is the importance of g' for innovation in g. Similarly, we proxy for the influence of foreign knowledge on innovation in China by including the current state of knowledge for granted US utility patents in the same IPC class, $\vec{\mu}_{US,t}^g$, where $p_{US,t}^g$ measures the importance of US knowledge for innovation in China.

4 Empirical implementation

4.1 OLS identification and a LASSO estimator

Our goal is to estimate the coefficients of the innovation process specified in equation (3.9) for every IPC and year: ρ_t^g , $\{p_{it}^g\}_{i\in\Omega_t}$, $\{\gamma_t^{gg'}\}_{g'\in\Gamma_g}$, and $p_{US,t}^g$. Given the discussion above, a natural approach would be to employ a simple OLS estimator. However, there are two issues with this method.

The first issue is that for identification to hold under OLS, all of the embeddings on the right-hand side of equation (3.9) must be linearly independent. This may fail to be true for two reasons. First, since $\vec{\mu}_t^g$ is defined as the mean of $\{\vec{x}_{it}^g\}_{i\in\Omega_t^g}$, for any realization of the latter there is likely to be a high degree of multicollinearity with the former. Hence, for the empirical model, we measure $\vec{\mu}_t^g$ as the average embedding in IPC g in years t, t-1, and t-2 instead (normalized to unit length), denoting this now by $\vec{\mu}_{Bt}^g$ to avoid confusion. For consistency, we also do the same for the secondary IPC embeddings $\vec{\mu}_t^{g'}$ and $\vec{\mu}_{US,t}^g$, replacing these with $\vec{\mu}_{Bt}^{g'}$ and $\vec{\mu}_{Bt}^g$, respectively. Similarly, we measure the future state of knowledge $\vec{\mu}_{t+1}^g$ as the average embedding in IPC g in years t+1, t+2, and t+3 instead (normalized to unit length), denoting this now by $\vec{\mu}_{Ft}^g$.

The second reason why identification may fail is that the dimension of the embedding space in the Cohere *embed-multilingual-v2.0* model is fixed at K = 768, whereas the number of patents N_t^g applied for in any given IPC-year may be much larger. This is a problem because a necessary condition for linear independence of the embeddings is that $\mathcal{M}_t^g \leq K$, where $\mathcal{M}_t^g = 2 + N_t^g + |\Gamma_g|$ is the number of coefficients being estimated (note that \mathcal{M}_t^g is the number of "regressors" and K is the number of "observations", so this is just the standard OLS condition that one cannot have fewer observations than regressors). For example, ignoring the other coefficients, the patent importance scores are only uniquely determined by equation (3.6) when the cosine similarity matrix is invertible, which requires $N_t \leq K$. Out of the 865,401 unique IPC-years in our data, only 987 IPC-years fail to meet this criterion, but these cases account for 12.7% of all patent applications.

To deal with this issue, we discard the OLS approach and instead make use of a Least Absolute Shrinkage and Selection Operator (LASSO) estimator. Specifically, the coefficients of interest, $\chi_t^g \equiv \left\{ \rho_t^g, \{p_{it}^g\}_{i \in \Omega_t}, \left\{ \gamma_t^{gg'} \right\}_{g' \in \Gamma_g}, p_{US,t}^g \right\}$, are estimated as follows:

$$\min_{\chi_t^g} \left\{ \|\vec{\mu}_{Ft}^g - \hat{\vec{\mu}}_{Ft}^g\|_2^2 + \alpha \|\vec{p}_t^g\|_1 \right\}$$
(4.1)

where $\hat{\mu}_{Ft}^g \equiv \rho_t^g \vec{\mu}_{Bt}^g + \sum_{i \in \Omega_t^g} p_{it}^g \vec{x}_{it}^g + \sum_{g' \in \Gamma_g} \gamma_t^{gg'} \vec{\mu}_{Bt}^{g'} + p_{US,t}^g \vec{\mu}_{US,Bt}^g$ is the predicted future state of knowledge. This adds to the standard OLS sum-of-squared-errors objective (the first term in equation (4.1)) an ℓ_1 -penalty (the second term) that promotes sparsity in the estimated patent importance scores, with the parameter α governing the strength of this penalty. This has a straightforward economic interpretation: the ℓ_1 -norm of the importance vector (i.e., the sum of the absolute values of $\{p_{it}^g\}_{i \in \Omega_t^g}$) can be viewed as a measure of the cost of *attention* that inventors have to pay to the current set of patents when choosing the direction of innovation. Hence, the ℓ_1 -penalty can be interpreted as a constraint that inventors have a limited attention span. Since only the embeddings $\{\vec{x}_{it}^g\}_{i \in \Omega_t^g}$ are at the individual patent level (whereas the other embeddings are averages across multiple patents), we apply the ℓ_1 -penalty only to the importance vector \vec{p}_t^g and not to the other coefficients in the estimation.

To select the value of α , a standard approach is to use *cross-validation*, where the dataset is first split into *n* groups, one of which is treated as test data and the remaining n - 1 are treated as training data. Given a candidate value of α , the LASSO model is estimated using the training data and the loss (the value of the objective function) is then evaluated on the test data given the estimated coefficients. This is done *n* times, with each iteration using a different group of data as the test. The value of α is then selected numerically to minimize the average loss across all iterations. While this can be implemented for each IPC-year separately in principle, this greatly increases the computational time for generating the LASSO estimates in practice. Hence, we use cross-validation on a small random sample of IPC-years to pick a value for α that we then use for all IPC-years. This process leads to a value of $\alpha = 10^{-9}$, so the LASSO estimator is in a sense very close to the OLS estimator. For the cases in which OLS is identified (i.e., when $\mathcal{M}_t^g \leq K$ and there is no multicollinearity), we also compute the OLS estimates of the coefficients and compare them to the LASSO estimates for robustness.

4.2 The relevant set of secondary IPC classes

The final issue that needs to be addressed before taking the model to data is defining the related set of IPCs, Γ_g , for each IPC g. While one may wish to impose as few assumptions about this set as possible, one should also note that there are 75,025 unique IPC codes in our data and hence some restriction on Γ_g is necessary given the discussion about identification and sparsity above. There are two sources of information that may be helpful in determining

 Γ_g : the set of secondary IPCs listed by a patent and the set of IPCs other than its own that a patent cites. As we discuss in Appendix A, the network of citations between IPCs is much sparser than the network defined by secondary IPCs. Hence, we choose to define Γ_g as the ten secondary IPCs reported most frequently by patents in IPC g^{13} .

One may be concerned that limiting the number of related IPCs to at most ten is imposing too strong a restriction on how technologies in different IPC classes can influence each other, especially since the maximum possible number of related IPC pairs is on the order of 5.6 billion $(75,025 \times 75,024)$ and we are allowing for only around 750,000. Nonetheless, even with this level of sparsity in the related-IPC network, around two-thirds of all IPCs are connected to each other either directly or indirectly in the so-called "giant component" of the network. Furthermore, these IPCs account for 97% of all patents. Hence, even with $|\Gamma_g| \leq 10$, almost all patents in our data are connected to each other at least indirectly in the related-IPC network.

5 Validation of patent importance

Before presenting our main findings, we validate our measure of patent importance by showing how it correlates with other economic outcomes. First, we investigate how patent importance correlates with two other commonly-used measures of patent quality that are used in the literature: the number of forward citations received by a patent and the grant status of a patent. While forward citations approximate the ex post importance of the cited patent to subsequent patents, the grant status is an ex ante measure of the inventive step of the patent relative to prior art. Second, we investigate which of these measures of patent quality – importance, forward citations, and grant status – best predict TFP and output at the firm level. One can view the first exercise as a form of validation internal to the context of patents and the second as a form of external validation.

5.1 Internal validation: correlations with forward citations and grant status

For the internal validation exercise, we first demean patent importance, number of forward citations received, and patent grant status by IPC-year effects.¹⁴ This effectively controls for IPC-specific time trends that might otherwise obfuscate any relationship between the patent outcomes. For instance, patents filed in earlier years may have more forward citations than patents filed in later years simply because the former have had a longer time to get cited. Similarly, the average patent grant rate may vary over time due to changes in the patent examination system that are unrelated to patent importance or patent citation behavior. Next, we assign patents into 50 percentile groups based on their demeaned importance score (so that each group accounts for 2% of all patents). We then regress the demeaned number of forward citations and grant status on dummy indicators for a patent's importance percentile group.

Figure 5.1 shows plots of the estimated coefficients for the demeaned number of forward citations and grant status, respectively, where the error bars denote 95% confidence intervals.

¹³If there are less than ten secondary IPCs reported by patents in IPC g, then $|\Gamma_g| < 10$.

 $^{^{14}}$ We exclude self-citations, where a self-citation is defined as one in which there is at least one common applicant name between the cited and citing patent.





Notes: These figures show the estimated coefficients from regressions of the number of forward citations received by a patent (panel (a)) and grant status (panel (b)) on dummy indicators for the patent's importance percentile group (two-percentile bins). All variables are residualized by IPC-year effects. Forward citations exclude self-citations. Bars around each point show the associated 95% confidence interval.

In panel (a), we observe a strikingly monotonic relationship between forward citations and a patent's importance percentile group: patents that have higher importance within an IPC-year tend to receive more forward citations than other patents within the same IPC-year. For instance, patents at the 90th percentile of the within-IPC-year distribution of patent importance (percentile group 45 on the x-axis) have, on average, 0.11 more citations than patents at the 10^{th} percentile of the distribution (percentile group 5). While this difference might seem small, the average number of forward citations per patent in our data is 0.44 with a standard deviation of 1.21, so a difference of 0.11 is sizable.

In panel (b) of Figure 5.1, we also observe a positive relationship between the grant status of a patent and its importance percentile group: patents that have higher importance within an IPC-year are more likely to be granted than other patents within the same IPC-year. Here, however, we observe an interesting pattern that is different from the case of forward citations – once a patent is "important enough", being more important is not associated with a higher grant rate. In particular, there appears to be a threshold around the median of the distribution (percentile group 25 on the x-axis), where patents with importance above the threshold have around a 2-3 percentage point higher grant rate than patents with importance below the threshold, but above the threshold, more important patents do not have a significantly different grant rate than less important patents. This is perhaps reflective of the patent examination process: to be granted, a patent simply has to be "innovative enough" to clear the examination bar and being more innovative beyond this does not lead to a higher probability of being granted.

5.2 External validation: predicting firm TFP and output

Figure 5.1 establishes that patent importance, forward citations, and grant status are positively correlated within an IPC-year. However, without a reference point external to the patent data, it is impossible to assess which of these measures is a "better" indicator of patent quality, if any. To investigate, we merge the patent data with data on above-scale firms from the Chinese National Bureau of Statistics (NBS) for the years 1998-2007. In this dataset, we have 1,784,156 firm-years and 438,707 firms with information on TFP and output.¹⁵ Of these, 52,033 firm-years and 16,183 firms have patents. The firms in these data account for around 25% of total output in the sample.

For each patent, we first determine whether it is in the top p% of the importance distribution within its application year (a "top-p% important patent") and in the top p% of the forward citation distribution within its application year (a "top-p% cited patent"), for a given threshold p. For each firm f and year t in the merged data, we then measure the following: (i) the stock of top-p% important patents applied for by f before and up to t, $N_{ft}^{topimp,p}$; (ii) the stock of top-p% cited patents applied for by f before and up to t, $N_{ft}^{topcit,p}$; and (iii) the stock of active patents owned by f at t, N_{ft}^{active} , i.e., patents applied for by f and granted before and up to tthat have not been terminated as of t.¹⁶ Appendix Table 1 provides summary statistics of these patent counts by year.

We then estimate the following regression via OLS:

$$y_{ft} = \beta^{topimp,p} N_{ft}^{topimp,p} + \beta^{topcit,p} N_{ft}^{topcit,p} + \beta^{active} N_{ft}^{active} + \gamma \delta_{ft} + \alpha_{h(f)t} + \epsilon_{ft}$$
(5.1)

where y_{ft} is a firm-year outcome of interest (e.g., log TFP or output), δ_{ft} is a dummy indicator equal to one if firm f has applied for at least one patent before and up to t, $\alpha_{h(f)t}$ is an industryyear fixed effect (where h(f) denotes the 4-digit industry code of firm f), and ϵ_{ft} is a residual. We will primarily be interested in comparing the regression coefficients $\beta^{topimp,p}$, $\beta^{topcit,p}$, and β^{active} , since this will be informative about which characteristics of a firm's patents best predict firm performance: having more important patents, having more highly-cited patents, or having more active patents.

Table 5.1 reports our estimated regression coefficients from equation (5.1) with *t*-statistics shown in parentheses, where the dependent variable y_{ft} is log TFP in panel (a) and log output in panel (b). Each column of the tables reports results for a different value of the top-patent threshold, *p*. Note that we also scale the dependent variable by a factor of 100 so that the magnitudes of the regression coefficients can be interpreted as percentage points (p.p.). We highlight four key observations.

First, the estimated coefficient on the stock of a firm's top-p% important patents, $\beta^{topimp,p}$, is always positive and significant at the 99% confidence level. For instance, each additional patent in the top 2% of the within-year importance distribution is associated with an increase in TFP of 0.40 p.p. (row i, column p = 2 of panel (a)) and an increase in output of 5.42

 $^{^{15}}$ The firm-level productivity estimates for 1998-2007 are taken from Brandt et al. (2025). Reporting problems in the NBS data after 2007 prevent analysis over a longer period.

¹⁶Patents may terminate because a firm fails to pay an annual renewal fee, for example.

(a) outcome: log TFP ($\times 100$); fixed effects: industry-year								
<i>p</i> =	2	5	10	25	50			
i. stock of top- $p\%$ important patents, $\beta^{topimp,p}$	0.40	0.31	0.17	0.05	0.03			
	(3.87)	(5.18)	(5.39)	(4.68)	(4.16)			
ii. stock of top- $p\%$ cited patents, $\beta^{topcit,p}$	0.14	0.08	0.05	0.02	-0.00			
	(2.88)	(3.17)	(3.22)	(2.21)	(-0.24)			
iii. stock of active patents, β^{active}	-0.03	-0.07	-0.09	-0.09	-0.09			
	(-1.48)	(-2.91)	(-3.55)	(-3.46)	(-3.25)			
iv. has patents, γ	3.66	3.63	3.63	3.66	3.68			
	(22.33)	(22.14)	(22.15)	(22.35)	(22.50)			
observations (m)	1.78	1.78	1.78	1.78	1.78			
R^2	0.73	0.73	0.73	0.73	0.73			
adjusted R^2	0.73	0.73	0.73	0.73	0.73			
(b) outcome: log output (×	100); fixed	effects: ind	lustry-year	r				
<i>p</i> =	2	5	10	25	50			
i. stock of top- $p\%$ important patents, $\beta^{topimp,p}$	5.42	5.11	2.69	0.76	0.29			
	(15.02)	(23.95)	(24.27)	(18.72)	(11.18)			
ii. stock of top- $p\%$ cited patents, $\beta^{topcit,p}$	0.59	0.47	0.34	-0.01	-0.13			
	(3.43)	(5.15)	(5.79)	(-0.19)	(-4.14)			
iii. stock of active patents, β^{active}	0.35	-0.28	-0.62	-0.32	-0.01			
	(4.90)	(-3.36)	(-6.91)	(-3.32)	(-0.10)			
iv. has patents, γ	110.64	110.09	110.17	110.64	111.01			
	(191.48)	(190.30)	(190.53)	(191.57)	(192.37)			
observations (m)	1.78	1.78	1.78	1.78	1.78			
R^2	0.15	0.15	0.15	0.15	0.15			
adjusted R^2	0.15	0.15	0.15	0.15	0.15			

Table 5.1: Regressions of firm TFP on patent stocks

Notes: This table reports the results of estimating equation (5.1) by OLS, where the dependent variable is log TFP in panel (a) and log output in panel (b). *t*-statistics are shown in parentheses.

p.p. (row i, column p = 2 of panel (b)). In other words, firms that own important patents are likely to be larger and more productive. Second, our estimate of $\beta^{topimp,p}$ is monotonically decreasing in the threshold percentile p. For example, each additional patent in the top 2%, 5%, 10%, 25%, and 50% of the within-year importance distribution is associated with an increase in TFP of 0.40, 0.31, 0.17, 0.05, and 0.03 p.p., respectively (row i of panel (a)). In other words, more important patents are more strongly correlated with firm TFP and output. These first two observations provide external validation that our measure of importance is capturing a characteristic of patents that is positively associated with measures of firm performance.

Third, in every specification, our estimate of $\beta^{topimp,p}$ is greater than our estimates of the coefficients on the other patent stock regressors, $\beta^{topcit,p}$ and β^{active} . For example, comparing rows i-iii of panel (a), we see that the estimated value of $\beta^{topcit,p}$ is around one-third the value of $\beta^{topimp,p}$, while the estimated value of β^{active} is negative. Similarly, comparing rows i-iii of panel (b), we observe that the estimates of $\beta^{topcit,p}$ and β^{active} are at least an order of magnitude smaller than our estimate of $\beta^{topimp,p}$. Finally, our estimate of $\beta^{topimp,p}$ decays much more slowly with the percentile threshold p compared with our estimate of $\beta^{topcit,p}$. For example, in panel (b), $\beta^{topimp,p}$ is positive and significant even for the p = 50 threshold, whereas $\beta^{topcit,p}$ is positive and significant only for the specifications where $p \leq 10$. In sum, our measure of patent importance is not only positively associated with firm performance, but it is also a better predictor of firm TFP and output than traditional measures of patent quality based on forward citations and the legal status of a patent.

In Online Appendix D, we provide results for three alternative versions of the regressions discussed above. In the first alternative, when constructing patent stocks in year t, we include only patents applied for between t-4 and t. In other words, we consider only "new" patents applied for by a firm. In the second alternative, we add firm fixed effects to the regressions. In both of these versions, we find that the four observations highlighted above continue to hold in almost all cases (the sole exception is that the coefficients on $N_{ft}^{topimp,p}$ are no longer statistically significant at the 99% confidence level when the outcome is log TFP and firm fixed effects are included in the regression). Finally, in the third alternative, we estimate the same baseline specification as in equation (5.1) but add as a regressor the stock of a firm's top-p% important patents where importance is based on the naive measure defined in equation (3.1). This allows us to assess whether our model-based measure of importance is superior to the naive measure in terms of predicting firm TFP and output. We find that in almost all cases, the estimated coefficient on our baseline measure of top-p% important patents is larger than the corresponding coefficient based on the naive measure of importance (the sole exception is for the specification where the outcome is log TFP and top patents are defined at the 2% threshold). Furthermore, in many specifications, the estimated coefficient on the naive measure of important patents is negative or statistically insignificant. The advantages that our baseline measure of patent importance have over the naive measure are thus meaningful for predicting firm performance.



Figure 6.1: The distribution of patent importance over time

Notes: Panel (a) shows the evolution of the distribution of patent importance across all patents. Panel (b) shows the evolution of the distribution of patent importance for patents with positive importance scores.

6 Results

We have now provided three forms of validation that text embeddings can be used to capture meaningful variation across patents and to measure the importance of patents: (i) patents in the same IPC section have embeddings that are more similar to each other (Figure 2.3); (ii) our measure of patent importance is positively correlated with other traditional measures of patent quality (Figure 5.1; and (iii) our measure of patent importance is a better predictor of firm TFP and output than traditional measures of patent quality (Table 5.1). We now discuss our main findings about how patent importance has evolved over time in China. We present these findings as a set of seven novel facts about innovation in China.

FACT 1. The importance of important patents has declined over time, due in part to falling innovativeness and increasing crowdedness of patenting in China.

Panel (a) of Figure 6.1 shows how the distribution of our estimated patent importance measure p_{it} changes over time. As expected given the discussion in section 3.3, the average value of patent importance is consistently close to zero in every year and is not a meaningful statistic. What is meaningful, however, is that the *dispersion* of the patent importance distribution is declining over time and, closely related, that the average importance of important patents, defined as those with $p_{it}^g > 0$, is also declining over time as shown in panel (b) of Figure 6.1. We propose two explanations for these trends.

First, the rate at which the state of knowledge in China is changing has slowed over time. To examine this formally, let us define the *innovativeness* of IPC g at time t as:

$$I_t^g = \frac{1}{2} \left[1 - c \left(\vec{\mu}_{Ft}^g, \vec{\mu}_{Bt}^g \right) \right]$$
(6.1)

Figure 6.2: Innovativeness and crowdedness of Chinese patenting



Notes: Panel (a) shows the evolution of the distribution of innovativeness defined in equation (6.1) across patents. Panel (b) shows the evolution of the distribution of crowdedness (number of patents applied for in the same IPC-year) across patents.

where $\vec{\mu}_{Ft}^{g}$ and $\vec{\mu}_{Bt}^{g}$ are the average patent embeddings (normalized to unit length) in IPC g in years $\{t + 1, t + 2, t + 3\}$ and $\{t - 2, t - 1, t\}$, respectively. Note that I_{t}^{g} is larger when $\vec{\mu}_{Ft}^{g}$ and $\vec{\mu}_{Bt}^{g}$ are less similar and hence innovativeness is a measure of how much the future state of knowledge differs from the past. Note also that $I_{t}^{g} \in [0, 1]$, since $c(\vec{\mu}_{Ft}^{g}, \vec{\mu}_{Bt}^{g}) \in [-1, 1]$. Panel (a) of Figure 6.2 shows how the distribution of innovativeness changes over time at the patent level (i.e., the distribution of $I_{it} \equiv I_{t}^{g(i)}$ where g(i) is the main IPC of patent i). Evidently, the innovativeness of Chinese patenting has declined steadily over time, indicating that the state of knowledge is changing at a slower pace. Recall that the patent importance scores are the weights assigned to the patent embeddings in order to "explain" the direction of innovation. Hence, when innovativeness is low, there is less that needs to be explained in the first place and the magnitudes of the patent importance scores tend to be small.

The second explanation for the decline in the dispersion of patent importance and the importance of important patents is that Chinese patenting has become much more *crowded* over time. More formally, for each patent *i* in IPC g(i) applied for at time *t*, we measure the total number of patents applied for in the same IPC-year, $N_{it} \equiv N_t^{g(i)}$. In panel (b) of Figure 6.2, we plot the distribution of N_{it} over time. Clearly, the crowdedness of Chinese patenting has increased substantially. For example, the average patent in 2005 has around 100 other patents in the same IPC-year, whereas this increases by about seven-fold by 2019. The crowdedness of innovation matters for the distribution of patent importance because when there are many candidate patents that can explain the direction of innovation, each individual patent is less likely to matter a lot.

To verify that the innovativeness and crowdedness of patenting are related to patent importance, we consider two regressions. First, at the patent-year level, we regress importance (p_{it})

	(1)	(2)	(3)	(4)
observation =	patent-year	patent-year	IPC-year	IPC-year
i. innovativeness	0.40	0.51	0.45	0.52
	(448.79)	(847.01)	(130.30)	(182.98)
ii. crowdedness	-0.13	-0.25	-0.10	-0.22
	(-91.45)	(-568.15)	(-40.85)	(-125.19)
IPC fixed effects	yes	no	yes	no
year fixed effects	yes	yes	yes	yes
observations (m)	6.01	6.01	0.42	0.42
R^2	0.36	0.28	0.39	0.24
adjusted \mathbb{R}^2	0.36	0.28	0.33	0.24
within \mathbb{R}^2	0.05	0.25	0.06	0.22

Table 6.1: Regressions of patent importance on innovativeness and crowdedness

Notes: Columns (1) and (2) report the results of regressing patent importance (p_{it}) on innovativeness (I_{it}) and crowdedness $(\log N_{it})$ at the patent-year level for the sample of important patents $(p_{it} > 0)$. Columns (3) and (4) report the results of regressing the average importance of important patents (\bar{p}_t^{g+}) on innovativeness (I_t^g) and crowdedness $(\log N_t^g)$ at the IPC-year level. Coefficients are reported as standardized beta coefficients and t-statistics are shown in parentheses.

on innovativeness (I_{it}) and crowdedness (measured as $\log N_{it}$) for the set of important patents $(p_{it} > 0)$. Columns (1) and (2) of Table 6.1 report the estimated standardized beta coefficients (*t*-statistics shown in parentheses) with and without the inclusion of IPC fixed effects, respectively. Second, at the IPC-year level, we compute the average importance of important patents, $\bar{p}_t^{g+} \equiv \frac{1}{|\Omega_t^{g+}|} \sum_{i \in \Omega_t^{g+}} p_{it}$, where $\Omega_t^{g+} \equiv \{i \in \Omega_t^g : p_{it} > 0\}$. We then regress \bar{p}_t^{g+} on innovativeness (I_t^g) and crowdedness ($\log N_t^g$) at the IPC-year level. Columns (3) and (4) of Table 6.1 report the estimated standardized beta coefficients, with and without the inclusion of IPC fixed effects, respectively.

Clearly, there is a close empirical relationship between patent importance and the innovativeness and crowdedness measures. For example, at the patent-year level within an IPC, an increase in innovativeness by one standard deviation is associated with an increase in patent importance by 43% of a standard deviation (row i, column (1) of Table 6.1). Similarly, an increase in crowdedness by one standard deviation is associated with a decrease in patent importance by 22% of a standard deviation (row ii, column (1) of Table 6.1). Hence, our finding that the average importance of important patents is declining over time is closely related to the fact that patenting in China is becoming less innovative and more crowded.

It is also important to note that the decline in innovativeness and increase in crowdedness of patenting is not unique to China. Similar patterns are observed since 1976 for granted USPTO utility patents (see Figure 2 in Online Appendix E for details). As a point of comparison, the average level of innovativeness for China in 2010 (around 0.02) is similar to the level of innovativeness reached in the US for granted utility patents in the mid 1990s. Similarly, the average level of crowdedness in China (around 200 patents) is close to the average level of crowdedness in the US in the early 2000s, although one should bear in mind that we are counting all patent applications for China but only granted patents for the US.

FACT 2. In some IPCs, innovativeness and the average importance of important patents have increased over time despite an increase in crowdedness. This is more likely to be the case for patents related to health, chemistry, and electricity; it is less likely to be the case for patents related to mining, printing, and apparel.

Despite the fact that the average importance of important patents is declining in the aggregate, there remains a great deal of heterogeneity in patent importance across IPCs. For instance, variation in innovativeness and crowdedness explain only 9% of the within-IPC variation in patent importance (last row of column (1) in Table 6.1).

To examine this heterogeneity further, panel (a) of Figure 6.3 shows how average innovativeness and crowdedness have changed between 2000-2010 and 2011-2019 for different IPCs, where each point in the scatter plot is an IPC and we include only IPCs that have at least an average of five patents per year in the two time windows.¹⁷ The area is divided into four cases – Case I: increasing innovativeness and increasing crowdedness; Case II: falling innovativeness and increasing crowdedness; Case III: increasing innovativeness and falling crowdedness; and Case IV: falling innovativeness and falling crowdedness. Evidently, most IPCs are in either case II or III, reflecting the negative correlation between changes in innovativeness and crowdedness on average. However, there are also some IPCs in case I, where innovativeness is increasing despite the increase in crowdedness. In panel (b) of Figure 6.3, we plot the average importance of important patents in the two time periods, where again each point in the plot is an IPC. We see that there are in fact some IPCs where the average importance of important patents is increasing over time and that this is more likely to be true for IPCs in case I versus those in case II. Hence, even though the average importance of important patents declines in China overall, this does not occur in all IPCs.

Table 6.2 also shows the share of patents in the pre-period (2000-2010) based on which of the above cases a patent's IPC belongs to, reporting this separately for 22 technology groups based on IPC class (the first three characters of the IPC code). We see, for instance, that 14.6% of patents in "Health; Life-saving" are in IPCs where both innovativeness and crowdedness are increasing whereas this is only true for 1% of patents in "Earth or Rock Drilling; Mining", suggesting that the former group of technologies is in some sense a more dynamic area of innovation than the latter. More generally, among the larger IPC groups, "Health; Life-saving", "Chemistry", and "Electric elements, Power, and Techniques" stand out as being the most dynamic, whereas "Earth of Rock Drilling; Mining", "Printing", and "Personal or Domestic Articles" appear to be least dynamic.

FACT **3.** The number of breakthrough patents has been growing in China, but growth rates have fallen sharply since the mid-2000s.

The decline in the average importance of important patents occurs in parallel with rapid growth in the number of patents overall (Figure 1.1). Hence, the stock of highly important patents may well be increasing. To investigate, we define a *breakthrough* patent as one that

 $^{^{17}\}mathrm{These}$ IPCs account for around 75% of all patents in 2000-2010 and 70% of all patents in 2011-2019.

	No. of patents, % of patents is				atents in:	
IPC group	Classes	2000-2010	Case I	${\rm Case~II}$	Case III	Case IV
Health; Life-saving	A61-A62	$173,\!846$	14.6	68.0	13.0	4.3
Chemistry	C01-C14	296,833	13.7	74.1	8.7	3.5
Foodstuffs; Tobacco	A21-A24	$50,\!518$	12.1	82.8	3.7	1.4
Electric Elements, Power, and Techniques	H01-H05	459,606	11.7	60.3	22.5	5.6
Engines or Pumps	F01-F04	47,231	10.1	82.0	5.3	2.7
Non-nuclear Physics	G01-G16	381,420	9.6	69.2	18.2	3.0
Paper	D21	5,732	9.1	69.7	13.6	7.7
Lighting; Heating	F21-F28	60,099	8.6	84.7	4.0	2.7
Agriculture	A01	36,793	8.4	81.7	8.7	1.1
Weapons; Blasting	F41-F42	2,014	6.5	74.3	8.8	10.4
Textiles or Flexible Materials	D01-D07	$30,\!159$	6.0	80.9	8.6	4.6
Metallurgy	C21-C30	$45,\!924$	5.9	89.6	3.3	1.2
Shaping	B21-B33	98,568	5.7	88.9	3.3	2.1
Separating; Mixing	B01-B09	59,323	5.5	90.4	2.9	1.1
Amusement	A63-A63	8,702	5.2	69.9	22.1	2.8
Building	E01-E06	46,301	5.1	83.7	5.4	5.9
Transporting	B60-B68	$90,\!536$	5.0	87.2	5.5	2.3
Personal or Domestic Articles	A41-A47	37,289	4.9	91.3	2.3	1.5
Nuclear Physics	G21	$1,\!615$	4.8	78.2	6.2	10.8
Engineering in General	F15-F17	43,717	4.6	87.9	4.1	3.4
Printing	B41-B44	20,387	4.4	73.4	17.0	5.2
Earth or Rock Drilling; Mining	E21	10,859	1.0	95.6	0.9	2.4

Table 6.2: Chan	ge in innovativenes	s and crowdedness.	, 2000-2010 to	2011-2019,	by IPC group

Notes: This table shows the share of patents in each IPC group in 2000-2010 with IPCs belonging to each of the four cases shown in Figure 6.3.





Notes: Panel (a) is a scatter plot of the change in innovativeness against the change in log crowdedness, where each point is an IPC and changes are defined as the average from 2011-2019 minus the average from 2000-2010. Panel (b) is a scatter plot of the average importance of patents with positive importance scores from 2011-2019 against the average from 2000-2010. The plots exclude IPCs with fewer than five patents per year in either of the two time periods.

has an importance score above the k^{th} percentile of the importance distribution over the entire sample of patents (so that the threshold does not vary across years).¹⁸ Figure 6.4 shows the number of breakthrough patents in each year for $k \in \{99, 95, 90\}$ and the associated annual growth rates. We indeed find that the number of breakthrough patents is increasing over time, with higher rates of growth in the early 2000s of around 15-20% per year. However, these growth rates decline steadily from 2005 onward, which is consistent with the decline in other measures of economic activity (e.g., GDP, exports, and TFP) in China during this period.

FACT 4. The reliance of Chinese innovation on knowledge within China versus outside China was constant from 1985-2000, increasing from 2000-2010, and constant again from 2010 onward. The increase from 2000-2010 occurred in a broad range of technology classes.

We now turn our attention to investigating which sources of knowledge are the most important for shaping the direction of innovation in China. The first margin that we consider is the relative importance of knowledge embodied in patents filed in China versus outside of China. Panel (a) of Figure 6.5 shows how the distribution of US importance, $p_{US,t}^g$ varies over time. Much like the distribution of importance for important Chinese patents, the importance of US knowledge is steadily declining. This decline may be explained in part by the same forces as discussed above – i.e., less innovative and more crowded patenting in China – but some of the decline may also reflect changes in the *relative* importance of knowledge embedded in Chinese versus US patents.

 $^{^{18}{\}rm This}$ is similar to the definition of a break through patent in Kelly et al. (2021) based on frequency counts of key words.



Figure 6.4: Number and growth rate of breakthrough patents

Notes: Panel (a) shows the number of breakthrough patents in each year, where a breakthrough patent is defined as one with an importance score above the x^{th} -percentile of the importance distribution across patents in all years, where $x \in \{99, 95, 90\}$. Panel (b) shows the corresponding growth rates in the number of breakthrough patents.

To investigate the latter, we measure for each IPC g and year t the share of Chinese patents with an importance score greater than the importance of US patents:

$$\mathcal{D}_t^g \equiv \frac{1}{N_t^g} \sum_{i \in \Omega_t^g} \mathbb{1}_{\left[p_{it}^g > p_{US,t}^g\right]}$$
(6.2)

We henceforth refer to this measure as *domestic reliance*. Panel (b) of Figure 6.5 plots the distribution of domestic reliance over time at the patent level (i.e., the distribution of $\mathcal{D}_{it} \equiv \mathcal{D}_t^{g(i)}$). We observe that average domestic reliance is roughly constant from 1985 to 2000 at around 30%. Throughout the 2000s, however, there is a sustained increase in the measure, so that by 2010, the average is just under 60%. From 2010 onward, average domestic reliance levels off and again remains roughly constant at around 60%. The 2000s were thus a key period during which innovation in China became increasingly directed by knowledge within China as opposed to knowledge outside.

Figure 6.5 also shows that there is substantial heterogeneity across IPCs in domestic reliance. To examine how the relative importance of Chinese versus US knowledge differs across technology classes, Table 6.3 reports how domestic reliance has changed over time for the 22 technology groups discussed above. Column (1) reports the share of all patents accounted for by each IPC group G, S^G . Columns (2) and (3) report domestic reliance for each IPC group G pooling all patents across 2000-2009 (\mathcal{D}_{00-09}^G) and 2010-2019 (\mathcal{D}_{10-19}^G), respectively. Column (4) reports the log change in domestic reliance across these two time periods, $\Delta \log \mathcal{D} \equiv \log \mathcal{D}_{10-19}^G - \log \mathcal{D}_{00-09}^G$, and column (5) reports the log change in domestic reliance for each IPC group G relative to the same log change across all patents.



Figure 6.5: US importance and domestic reliance

Notes: Panel (a) shows the evolution of the distribution of USPTO patent importance, $p_{US,t}^g$, across patents. Panel (b) shows the evolution of the distribution of domestic reliance (defined as the share of CNIPA patents in an IPC-year that are more important than USPTO patents).

Two observations are noteworthy. First, there is an increase in domestic reliance for all 22 IPC groups from 2000-2009 to 2010-2019 (all values in column (4) are positive). The relative shift in the direction of innovation toward Chinese versus foreign knowledge is not limited to only a few technology classes but rather is a widespread phenomenon in China. Second, even though there is an increase in domestic reliance for all IPC groups, there is also substantial heterogeneity across groups in the rate of this increase. Domestic reliance grows the fastest in technologies such as Paper (D21), Health; Life-saving (A61-A62), Metallurgy (C21-C30), and Chemistry (C01-C14), whereas it grows the slowest in Personal or Domestic Articles (A41-A47), Engines or Pumps (F01-F04), Printing (B41-B4), Transporting (B60-B68), and Engineering in General (F15-F17). The IPC groups in which domestic reliance grows the fastest also tend to be those in which both innovativeness and crowdedness increase over time (Case I in Table 6.2) – Health; Life-saving (A61-A62) and Chemistry (C01-C14), for instance, rank highly in both dimensions.

FACT 5. SOE, university, research institute, and PIE patents have tended to be the most important within an IPC-year and have been over-represented among breakthrough patents, whereas FIE and overseas patents have been less important and under-represented among breakthrough patents. Differences in importance between these two groups of patents have been falling since the early 2000s.

We now examine which sources of knowledge within China have been more important for driving the direction of innovation in China. We base our analysis on the patentee types described in section 2.1. Recall that there are important changes over time in the shares of patents accounted for by different patentee types, with rapid growth in the PIE share together with a marked decline in the share for overseas patentees filing patents in China.

		(1)	(2)	(3)	(4)	(5)
IPC group	Classes	S^G	\mathcal{D}^G_{00-09}	\mathcal{D}^G_{10-19}	$\Delta \log \mathcal{D}^G$	$\Delta \log \hat{\mathcal{D}}^G$
Paper	D21	0.2	42.8	57.1	28.8	11.3
Health; Life-saving	A61-A62	7.1	46.3	60.8	27.2	9.8
Metallurgy	C21-C30	2.1	47.3	60.2	24.1	6.7
Chemistry	C01-C14	13.3	48.5	60.6	22.4	4.9
Weapons; Blasting	F41-F42	0.1	41.8	51.0	19.9	2.5
Amusement	A63	0.5	48.0	58.2	19.2	1.7
Nuclear Physics	G21	0.1	41.0	49.3	18.5	1.0
Foodstuffs; Tobacco	A21-A24	3.4	49.7	59.7	18.2	0.8
Electric Elements and Techniques	H01-H05	17.3	47.0	56.0	17.7	0.2
Earth or Rock Drilling; Mining	E21	0.7	50.2	59.0	16.2	-1.3
Separating; Mixing	B01-B09	4.0	51.7	60.0	14.9	-2.5
Non-nuclear Physics	G01-G16	20.2	47.8	55.3	14.7	-2.8
Lighting; Heating	F21-F28	3.3	48.2	55.6	14.4	-3.0
Shaping	B21-B33	6.9	50.7	58.4	14.2	-3.2
Agriculture	A01	2.8	49.4	56.7	13.9	-3.5
Textiles or Flexible Materials	D01-D07	1.5	48.2	54.7	12.6	-4.8
Building	E01-E06	3.1	50.8	57.6	12.5	-4.9
Engineering in General	F15-F17	2.3	52.1	58.7	12.0	-5.4
Transporting	B60-B68	6,0	50.1	56.0	11.0	-6.5
Printing	B41-B44	0.7	49.8	55.0	10.0	-7.4
Engines or Pumps	F01-F04	1.9	52.0	56.7	8.5	-8.9
Personal or Domestic Articles	A41-A47	2.5	50.9	55.0	7.7	-9.7

Table 6.3: Domestic reliance by IPC group

Notes: Column (1) reports the share of patents across all years accounted for by each IPC group G, S^G . Columns (2) and (3) report the domestic importance share pooling all patents across 2000-2009 and 2010-2019, respectively. Column (4) reports the log difference between columns (3) and (2). Column (5) subtracts from column (4) the log change in the domestic importance share from 2000-2009 to 2010-2019 for all patents. Values in columns (1)-(3) are reported as percentages. Rows are sorted in descending order of the values in column (5).





Notes: This figure shows the evolution of average patent importance residualized by IPC-year effects for different patentee types. Values shown for year t are computed as the average between years t - 4 and t for visual clarity.

To compare the importance of patents across patentee types, we first compute the importance of each patent parsed of IPC-year effects. We then average this measure over all patents applied for by a given patentee type in a given year. This measures the *relative* importance of patents by each patentee type within an IPC-year. Figure 6.6 shows how this relative measure changes over time. For visual clarity, we omit the "Individual" and "Other Domestic" patentee types and plot five-year rolling averages of the relative importance scores.

We highlight three key observations. First, state-related patents (SOE, university, and research institute patents) and PIE patents are more important on average than FIE and overseas patents in almost every year. Before 2000, university and research institute patents tend to be the most important, but by 2005, SOE patents assume the top position in terms of relative importance. Second, in the mid-1990s, overseas patents had slightly higher than average importance. However, the relative importance of these patents declines steadily, so that from 2000 onward, overseas patents are less important than patents applied for by enterprises, universities, and research institutes in China. A potential explanation for this is that overseas patentees may be filing patents in China mainly to protect the intellectual property of goods that are being exported to the Chinese market and hence these patents are not important for innovation in China, although one would need patentee-level export data to investigate this thoroughly. We also find that overseas patents at the USPTO have also been less important for US innovation than domestic patents since the mid-1980s, although the gap in importance is much smaller than it is for China (see Figure 3 in Online Appendix E for details).¹⁹ Finally, we observe that differences in average importance between state-related and PIE patents on the one hand and

¹⁹These estimates are generated by estimating the same model specified in equation (3.9) for granted USPTO utility patents, replacing the average USPTO embedding on the right-hand side with the average CNIPA embedding.

FIE and overseas patents on the other hand have been falling since the early 2000s, which is due in part to the fact that the overall dispersion of patent importance has been falling over time (Panel (a) of Figure 6.1).

In Appendix B, we show that patent importance can also be decomposed as:

$$p_{it}^g \approx \sum_s S_{Ft}^{gs} p_{it}^{gs} \tag{6.3}$$

where p_{it}^{gs} is the importance of patent *i* for patenting by patentee type *s* and S_{Ft}^s is the share of future patents attributed to patentee type *s* in IPC *g*.²⁰ This decomposition allows us to further examine how patent importance varies by patentee type pairs. In Figure 6.7, each matrix shows, for a particular five-year window, the average importance of patents by patentee type *s* (rows) for patenting by patentee type *s'* (columns), defined as $\tilde{p}_t^{ss'} = \frac{1}{|\Omega_t^s|} \sum_g \sum_{i \in \Omega_t^{gs}} p_{it}^{gs'}$ where Ω_t^{gs} is the set of patents in IPC *g* applied for by patentee type *s* and $\Omega_t^s \equiv \bigcup_g \Omega_t^{gs}$. The column on the right of each matrix shows the average importance of patents for each patentee type, $\frac{1}{|\Omega_t^s|} \sum_g \sum_{i \in \Omega_t^{gs}} p_{it}^g$. For visual clarity, all importance scores are multiplied by 1000.

As expected, we find that patents by a given patentee type are always the most important for themselves (the diagonals of each matrix are almost always larger than the off-diagonals). In particular, we observe that overseas patents are always important for themselves (and in fact are increasingly more so from 2000-2019) but are almost never important for patenting by domestic patentees. This sheds light on why the average relative importance of overseas patents falls overall: these patents are only important for themselves and account for a smaller share of total patents over time. We also observe that among domestic patentees, the average importance of overseas patents is highest for FIEs, which again is to be expected. Patents by PIEs and SOEs are also important for each other, but the importance of SOE patenting for PIE patenting tends to be higher than the other way around. FIE patents tend to be important for PIE patenting but less so for SOE patenting. University and research institute patents are important for each other as well as for PIE and SOE patenting, but tend not to be important for FIE patenting.

Finally, we note that measures of average patent importance may mask differences across patentee types among the most important patents. Hence, we now examine the representation of each patentee type among the breakthrough patents discussed above. For patentee type s in year t, we measure this representation using the following statistic:

$$\bar{r}_{st} = S_{st}^* - S_t^* \tag{6.4}$$

where S_{st}^* is the share of patents by patentee type s in all IPCs that are breakthrough patents, and S_t^* is the share of all patents in all IPCs that are breakthrough patents. Hence, \bar{r}_{st} measures whether patentee type s is over- or under-represented among breakthrough patents overall.

Figure 6.8 plots the representation measure \bar{r}_{st} over time, where breakthrough patents are defined as either those in the top 10% (panel (a)) or top 1% (panel (b)) of the importance

²⁰This decomposition is exact if patent importance is estimated via OLS and hence is an approximation only because we estimate importance using the LASSO estimator. Nonetheless, the decomposition is almost exact: the correlation between p_{it} and the right-hand side of (6.3) is 0.98.



Figure 6.7: Patent importance by patentee type pairs

Figure 6.8: Representation among breakthrough patents by patentee type



Notes: This figure shows the representation of different patentee types among breakthrough patents as defined in equation (6.4). Panels (a) and (b) define breakthrough patents as those with importance above the 90^{th} -percentile and 99^{th} -percentile of the importance distribution across patents in all years, respectively. Values shown for year t are computed as the average between years t - 4 and t for visual clarity.

distribution across all years. Again, for visual clarity, we omit the "Individual" and "Other Domestic" patentee types and plot five-year rolling averages. Consistent with the patterns discussed above, we see in Figure 6.8 that overseas patents are under-represented among breakthrough patents within IPCs, although differences in representation across patentee types are smaller when breakthrough patents are defined at the 1% threshold. There are also interesting differences between the patterns based on average importance and those based on representation among breakthrough patents. For instance, even though FIE patents have higher-than-average importance scores within IPC-years in Figure 6.6, we find these patents are under-represented among breakthrough patents. Furthermore, SOE, university, and research institute patents stand out much more clearly in terms of being over-represented among breakthrough patents.

FACT 6. New patentees have accounted for around one-fifth of patenting in China and have produced more important patents than experienced patentees, although this gap has fallen steadily since the mid-to-late 2000s.

We next consider the role played by new patentees. We define a patent to be owned by a new patentee if none of the applicants on the patent have applied for a patent in China before. Panel (a) of Figure 6.9 shows the share of patents in each year applied for by new versus experienced patentees.²¹ In the early 1990s, around half of all patents are applied for by new patentees. This share then declines steadily until the mid-2000s, when it levels off at around 20%. To examine differences in patent importance between new and experienced patentees, we regress patent importance parsed of IPC-year effects on a dummy indicator for whether the

 $^{^{21}}$ We show this from 1991 onward because our data begin in 1985 and hence in the earlier years of our sample, almost all patentees are considered "new".

patent is owned by a new patentee, interacted with the application year of the patent. We can then interpret the coefficient on the regressor as the new patentee importance premium. This is shown in panel (b) of Figure 6.9, where the y-axis values are shown in terms of the standard deviation of importance across all patents. Beginning in the early 2000s, new patentees tend to have slightly more important patents than experienced patentees (with a difference of around 5% of the importance standard deviation on average). However, the new patentee importance premium declines steadily over time, so that by the end of our sample, new and experienced patentees have patents that are almost identical in terms of average importance.



Figure 6.9: Representation among breakthrough patents by patentee type

Notes: Panel (a) shows the share of patents accounted for by new patentees in each year, where a patent is defined as being owned by a new patentee if none of the patentees on the application have applied for a patent with the CNIPA before. Panel (b) shows the estimated coefficients from a regression of patent importance parsed of IPC-year effects on a dummy indicator for whether the patent is owned by a new patentee, interacted with the application year of the patent. Bars around each point show 95% confidence intervals.

FACT 7. Patenting in China has become narrower over time, with the importance of secondary IPCs declining and the importance of own-IPC memory increasing over time.

Finally, we examine how the persistence of the innovation process and the importance of interactions across IPCs have evolved over time. Panel (a) of Figure 6.10 shows our estimates of the memory coefficient ρ_t^g , while panel (b) shows our estimates of secondary IPC importance $\gamma_t^{gg'}$. Note that in panel (b), secondary IPCs g' are ordered based on their rank of $S_{sec}^{gg'}$ for each main IPC g, where recall that $S_{sec}^{gg'}$ is the share of patents with main IPC g that list g' as a secondary IPC. We observe that memory is roughly constant in the initial years of our sample, but trends steadily upward beginning in the early 2000s. Furthermore, the average importance of secondary IPCs is not only lower than own-IPC importance as expected, but also trends steadily downard throughout our sample. This is suggestive that patenting in China is becoming "narrower", where knowledge generated within an IPC is more important for innovation in the IPC as opposed to knowledge generated in other related IPCs. We see similar patterns at the



Figure 6.10: Memory and importance of secondary IPCs

Notes: Panel (a) shows the evolution of the distribution of memory, ρ_t^g . Panel (b) shows the evolution of the average importance of secondary IPCs, where secondary IPCs are ranked according to how frequently they are reported as secondary IPCs.

USPTO (see Appendix Figure 4), although the upward trend in memory begins about a decade earlier (in the 1990s). Lastly, note from panel (b) that the ranking of average importance across secondary IPCs is almost always perfectly correlated with the ranking of secondary IPCs by $S_{sec}^{gg'}$, i.e., secondary IPCs that are reported more frequently by patents in an IPC are estimated to have higher importance for innovation in that IPC. This once again provides validation that our importance measure is correlated with other patent characteristics in exactly the way that one would expect.

7 Conclusion

In this paper, we have leveraged frontier methods in machine learning, a new model of innovation, and novel data on patent ownership to quantify which sources of knowledge have been important for driving the direction of innovation in China. We have developed a new measure of patent importance based on the embeddings of patent text and have shown that this measure is correlated with traditional measures of quality and outperforms these measures in terms of predicting TFP and output at the firm level.

Our findings about the evolution of patent importance in China can be summarized as follows. Overall, the importance of important patents in China has declined, due in part to falling innovativeness and increasing crowdedness of Chinese patenting, although some technologies such as those related to health, chemistry, and electricity have been more dynamic than others. Despite this, the number of breakthrough patents in China has continued to rise, but growth rates have been falling steadily since the mid-2000s. This has been accompanied by increasing reliance of Chinese innovation on domestic rather than foreign knowledge. State-related and PIE patents have also been more important than FIE and overseas patents on average, although these differences have been dissipating since the mid-2000s. Furthermore, new patentees initially were the source of more important patents than experienced patentees, but this premium has declined steadily since the mid-2000s. Finally, Chinese patenting has become narrower, with within-IPC knowledge becoming relatively more important than knowledge in other IPCs.

Several directions for future research are promising. First, an important question that arises from our analysis is why the role of overseas patents has sharply diminished in China. Some of this may be expected given improvements in China's own innovative capabilities, but the extent of the decline appears out of line with the growing importance of the Chinese market in a global context more broadly and the experience of the US. Second, our approach to measuring patent importance in China can also be applied to other contexts, such as patenting at the USPTO. This will not only be helpful for examining the forces underlying the decline in the growth of overseas patents in China, but also to provide context for the trends in Chinese patent importance that we have documented here. Third, we have conducted the analysis in this paper largely at the patent level, but aggregating patent information to the applicant level will also be important for identifying the firms, universities, and institutes that are the most key for Chinese patenting activity. Fourth, innovation is a key focus of Chinese industrial policy and hence it is essential to understand how patenting and innovative activity more broadly are being directed by policies such as those targeting strategic and emerging industries. Are these policies putting a priority on indigenous innovation, import substitution and leap-frogging helping to reverse the secular decline in innovativeness, or hastening it? Finally, our method of employing LLMs to study patent abstracts can also be extended to incorporate information contained in patent claims. This source of text data is much larger in scale and may be helpful in constructing more refined measures of patent importance and similarity.

A Defining the relevant set of secondary IPC classes

Formally, let $S_{sec}^{gg'}$ and $S_{cit}^{gg'}$ denote the shares of patents with main IPC g that list g' as a secondary IPC and that cite at least one patent with main IPC g', respectively. One might then deem $g' \in \Gamma_g$ if either $S_{sec}^{gg'}$ or $S_{cit}^{gg'}$ is above a certain threshold. Note that the secondary IPCs and cited IPCs partially overlap. For example, on the set of all IPC pairs gg' where both $S_{sec}^{gg'} > 0$ and $S_{cit}^{gg'} > 0$, the correlation between the two shares is 0.46. In other words, IPCs that are listed as secondary classes also tend to get cited.²² However, the network of citations between IPCs is much sparser than the network defined by the secondary IPCs. For instance, $S_{sec}^{gg'}$ is strictly positive for around 4.5 million unique pairs of IPCs, whereas this number is only 1.8 million for $S_{cit}^{gg'}$. Similarly, out of all 75,025 unique IPCs, 93% (69,503 IPCs) have patents that report at least one secondary IPC, whereas only 75% (56,513 IPCs) have patents that cite patents in at least one other IPC besides their own. In light of this, we choose to define Γ_g on the basis of the secondary IPC network.

Note that most patents report only a small number of secondary IPCs. For example, if we define Γ_g to be the set of IPCs where $S_{sec}^{gg'}$ is greater than 10%, which is already a fairly low threshold, then the average patent is in an IPC with $|\Gamma_g| = 2$ and the patent at the 99th percentile of this measure has $|\Gamma_g| = 15$. Hence, we adopt a simpler rule and define Γ_g to be the top ten secondary IPCs for g on the basis of $S_{sec}^{gg'}$.

B Decomposition of patent importance by patentee type pair

Consider the OLS estimate of the coefficient vector \vec{P}_t^g for the model specified in equation (3.9) with the modifications described in section (4.1). This is given by:

$$\left(\vec{P}_t^g\right)' = \left(\bar{X}_t^{g'} \bar{X}_t^g\right)^{-1} \bar{X}_t^{g'} \vec{\mu}_{Ft}^g \tag{B.1}$$

where \bar{X}^{g}_{t} contains the embeddings of all right-hand side regressors. The average forward embedding $\vec{\mu}_{Ft}^{g}$ can be decomposed into the average forward embedding for patents owned by each patentee type s, $\vec{\mu}_{Ft}^{gs}$:

$$\vec{\mu}_{Ft}^{g} = \sum_{s} S_{t+1}^{gs} \vec{\mu}_{Ft}^{gs} \tag{B.2}$$

where S_{t+1}^s is the share of future patents attributed to patentee type s. Hence, the OLS coefficient estimate can be decomposed as:

$$\left(\vec{P}_t^g\right)' = \sum_s S_{Ft}^{gs} \left(\vec{P}_t^{gs}\right)' \tag{B.3}$$

where $\left(\vec{P}_t^{gs}\right)' = \left(\bar{X}_t^{g'}\bar{X}_t^g\right)^{-1} \bar{X}_t^{g'}\vec{\mu}_{Ft}^{gs}$ contains measures of the importance for patenting by patentee type s.

²²We also find a positive correlation of 0.48 between $S_{sec}^{gg'}$ and the share of patents with main IPC g' that get at least one citation from a patent with main IPC g.

References

- Ahuja, G. and C. M. Lampert (2011). Entrepreneurship in the large corporation: A longitudinal study of how established firms create breakthrough inventions. <u>Strategic Management</u> Journal 22(6-7), 521–543.
- Akcigit, U. and S. Ates (2023). What happened to u.s. business dynamism? Journal of Political Economy 131(8), 2059–2124.
- Banerjee, P. M. and B. M. Cole (2011). Globally radical technologies and locally radical technologies: The role of audiences in the construction of innovative impact in biotechnology. IEEE Transactions on Engineering Management 58(2), 262–274.
- Bloom, N., C. I. Jones, J. V. Reenen, and M. Webb (2020). Are ideas getting harder to find? American Economic Review 110(4), 1104–1144.
- Boeing, P. (2016). The allocation and effectiveness of china's rd subsidies: Evidence from listed firms. Research Policy 45(9), 1774–1789.
- Boeing, P., E. Mueller, and P. Sandner (2016). China's rd explosion—analyzing productivity effects across ownership types and over time. Research Policy 45(1), 159–176.
- Brandt, L., J. Van Biesebroeck, L. Wang, and Y. Zhang (2025). Where has all the dynamism gone? productivity growth in china's manufacturing sector, 1998-2013. Working paper.
- Chondrakis, G., E. Melero, and M. Sako (2021). The effect of coordination requirements on sourcing decisions: Evidence from patent prosecution services. <u>Strategic Management</u> Journal 43, 1141–1169.
- Covarrubias, M., G. Gutiérrez, and T. Philippon (2019). From good to bad concentration? u.s. industries over the past 30 years. Working paper.
- Dang, J. and K. Motohashi (2015). Patent statistics: A good indicator for innovation in china? patent subsidy program impacts on patent quality. China Economic Review (35), 137–155.
- Fang, H., Z. M. Song, and Y. Zhang (2021). An anatomy of the patent quality: China vs. us. Working paper.
- Funk, R. J. and J. Owen-Smith (2017). A dynamic network measure of technological change. Management Science 63(3), 791–817.
- Griliches, Z. (1990). Patent statistics as economic indicators: A survey. <u>Journal of Economic</u> Literature 28(4), 1661–1707.
- Hall, B. H., A. Jaffe, and M. Trajtenberg (2005). Market value and patent citations. <u>RAND</u> Journal of Economics 36, 16–38.
- Harhoff, D., F. M. Scherer, and K. Vopel (2003). Citations, family size, opposition and the value of patent rights. <u>Research Policy</u> 32(8), 1343–1363.

- Hedge, D., K. Herkenhoff, and C. Zhu (2023). Patent publication and innovation. <u>Journal of</u> Political Economy 131(7), 1845–1903.
- Hu, A. G. Z. and G. H. Jefferson (2009). A great wall of patents: What is behind china's recent patent explosion? Journal of Development Economics (90), 57–68.
- Hu, A. G. Z., P. Zhang, and L. Zhao (2017). China as number one? evidence from china's most recent patenting surge. Journal of Development Economics (124), 107–119.
- Jefferson, G., A. G. Z. Hu, X. Guan, and X. Yu (2003). Ownership, performance, and innovation in china's large-and medium-size industrial enterprise sector. <u>China Economic Review</u> <u>14</u>(1), 89–113.
- Kalyani, A. (2024). The creativity decline: Evidence from us patents. Working paper.
- Kelly, B., D. Papanikolaou, A. Seru, and M. Taddy (2021). Measuring technological innovation over the long run. American Economic Review: Insights 3(3), 303–320.
- Kuhn, J., K. Younge, and A. Marco (2020). Patent citations reexamined. <u>RAND Journal of</u> Economics 51(1), 109–132.
- Li, X. (2012). Behind the recent surge of chinese patenting: An institutional view. <u>Research</u> Policy (41), 236–249.
- Marco, A. C., J. D. Sarnoff, and C. A. W. deGrazia (2019). Patent claims and patent scope. Research Policy 48(9).
- Park, M., E. Leahey, and R. J. Funk (2023). Papers and patents are becoming less disruptive over time. Nature 613(7942), 138–144.
- Younge, K. A. and J. M. Kuhn (2016). Patent-to-patent similarity: A vector space model. Working paper.

Online Appendix to

"The Anatomy of Chinese Innovation: Insights on Patent Quality and Ownership"

by

Philipp Boeing, Loren Brandt, Ruochen Dai, Kevin Lim, and Bettina Peters

A Identifying entity and patentee types

We first assign each unique applicant name *appname* in the patent database an entity type *enttype* based on the following sequential keyword search.

- 1. If appname contains 解放军, assign enttype "military".
- If appname contains {公司, 会社, 独立行政法人, 董事会, 集团, 有限责任, 股份, 企业, 公室, 工司, 公社, 商店, 工业}, contains 厂 and has length greater than three characters, or ends in 站, assign *enttype* "enterprise".
- 3. If *appname* contains {部, 厅, 局} and has length greater than three characters, assign *enttype* "government".
- 4. If appname contains {大学, 学校}, assign enttype "university".
- 5. If *appname* contains {研究所, 技术院, 设计院, 检定所, 验所, 技术协会, 科院, 科所, 科研 所, 检定所}, assign *enttype* "institute".
- 6. If appname contains 学院, assign enttype "university".
- 7. If appname contains 研究, assign enttype "institute".
- 8. If *appname* contains {中心, 医院, 协会, 制造院, 基金会, 工技协, 学会, 技术委, 技术室, 服务处, 委员会}, assign *enttype* "organization".
- 9. If appname has length no greater than three characters, assign enttype "individual".
- 10. If appname does not satisfy any of the above, assign enttype "other".

Note that the keyword search is sequential in the sense that if an applicant name has already been assigned an entity type at some step in the process, it is removed from the set of applicant names that remain to be assigned at later steps in the process. For example, $\psi \equiv \Lambda \neq \Re \oplus \Pi$ (Institute of Microbiology, Chinese Academy of Sciences) is identified as an institute at step 5 while 福建工程学院 (Fujian College of Engineering) is identified as a university at step 6, even though both names contain the characters $\neq \Re$.

The next step is to assign each patent an entity type *penttype* based on the applicant names listed in the patent application. Note that entity types at the patent level (*penttype*) and applicant level (*enttype*) can differ because some patents have multiple applicants. The assignment of *penttype* is based on the following procedure.

- 1. If at least one entity type is "enterprise" and no entity types are in {"university", "institute", "organization", "government", "military"}, assign *penttype* "enterprise".
- 2. If at least one entity type is "university" and no entity types are in {"company", "institute", "organization", "government", "military"}, assign *penttype* "university".
- 3. If at least one entity type is "institute" and no entity types are in {"company", "university", "organization", "government", "military"}, assign *penttype* "institute".
- 4. If at least one entity type is "organization" and no entity types are in {"company", "university", "institute", "government", "military"}, assign *penttype* "organization".
- 5. If at least one entity type is "government" and no entity types are in {"company", "university", "institute", "organization", "military"}, assign *penttype* "government".
- 6. If at least one entity type is "military" and no entity types are in {"company", "university", "institute", "organization", "government"}, assign *penttype* "military".

- 7. If at least one entity type is "individual" and no entity types are in {"company", "university", "institute", "organization", "government", "military"}, assign *penttype* "individual".
- 8. If no entity types are in {"company", "university", "institute", "organization", "government", "military", "individual"}, assign *penttype* "other".
- 9. Assign all remaining patents *penttype* "mixed".

Finally, we combine information about the location, entity type, and ownership type of each patent to define eight mutually exclusive patentee types:

	location	entity type	ownership
private-invested enterprise (PIE)	domestic	enterprise	private
state-owned enterprise (SOE)	$\operatorname{domestic}$	enterprise	state
foreign-invested enterprise (FIE)	$\operatorname{domestic}$	enterprise	foreign
university	$\operatorname{domestic}$	university	-
institute	$\operatorname{domestic}$	institute	-
individual	$\operatorname{domestic}$	individual	-
other domestic	$\operatorname{domestic}$	other	-
overseas	overseas	-	-

Table 1: Classification of patentee types

B Basic patent statistics

Technology classes. Panel (i) of Table 1 categorizes patents by their main International Patent Classification (IPC) section (1-digit). The composition of invention patents by technology class changes little over time. Sections A (Human Necessities), B (Performing Operations; Transporting), C (Chemistry; Metallurgy), G (Physics) and H (Electricity) account for the largest shares of invention patents. Patenting in sections A (Human Necessities) and C (Chemistry; Metallurgy) become slightly less prevalent after 2000, whereas the opposite is true for patenting in section G (Physics). Patenting in section H (Electricity) was also relatively more prevalent from 2000-2010 compared with after 2010.

Product versus process patents. Panel (ii) of Table 1 categorizes patents into process versus product patents based on the text of the claims of a patent. First, a set of keywords referring to prior claims is used to identify dependent claims in a patent and separate them from independent claims. Second, each independent claim is searched for a set of 60 positive keywords that are associated with process innovation. This search is carried out in the claim preamble (identified by the position of the so-called transitional phrase, a comma, or a colon, in that order) or in the entire claim text if no preamble is found. If there are several keyword matches, the keyword coming last in the selected text is used. If at least one such keyword is found, the claim is labeled as "process"; if not, the claim is labeled as "product". We then define process (product) patents as patents for which every independent claim is a process (product) claim. Patents that have both process and product claims are labeled "mixed". We observe that the majority of patents in China are product or mixed patents, i.e., they have at least one product claim. These account for more than 80% of all patents in the average year.

Geography. Panels (i)-(iii) of Table 2 categorize patents by the location of the primary patent applicant (based on the reported address). The shift in patenting activity away from overseas applicants (those with an address outside of China) toward domestic applicants (those with an address in China) is clear from panel (i): overseas applicants account for more than half of all

	'85-'89	'90-'94	'94-'99	'00-'04	'05-'09	'10-'14	'15-'19	all years	
all invention patents	41	75	182	482	1,203	2,181	$7,\!129$	11,292	
(i) invention patents by main IPC section									
A: Human Necessities	8	21	41	85	179	341	$1,\!155$	1,831	
	(18.3)	(28.2)	(22.8)	(17.6)	(14.9)	(15.6)	(16.2)	(16.2)	
B: Performing Operations	9	14	28	62	158	357	$1,\!379$	2,005	
	(21.6)	(17.8)	(15.4)	(12.8)	(13.1)	(16.4)	(19.3)	(17.8)	
C: Chemistry; Metallurgy	9	15	32	82	199	362	1,035	1,735	
	(22.7)	(20.5)	(17.7)	(17.0)	(16.6)	(16.6)	(14.5)	(15.4)	
D: Textiles; Paper	1	2	4	9	21	39	117	193	
	(2.5)	(2.4)	(2.0)	(1.9)	(1.7)	(1.8)	(1.6)	(1.7)	
E: Fixed Constructions	2	3	4	13	32	78	303	435	
	(3.7)	(3.3)	(2.4)	(2.6)	(2.7)	(3.6)	(4.3)	(3.9)	
F: Mechanical Engineering	4	5	12	31	92	185	530	859	
	(8.8)	(7.3)	(6.6)	(6.6)	(7.7)	(8.5)	(7.4)	(7.6)	
G: Physics	5	8	27	93	230	398	1,529	2,290	
	(12.6)	(10.7)	(14.8)	(19.3)	(19.1)	(18.2)	(21.4)	(20.3)	
H: Electricity	4	7	33	107	291	422	1,080	1,945	
	(9.9)	(9.9)	(18.3)	(22.3)	(24.2)	(19.3)	(15.2)	(17.2)	
	(ii) invent	tion pater	nts by pro	ocess/pro	duct type)			
process	7	15	27	78	217	405	1,243	1,990	
	(16.6)	(20.0)	(14.7)	(16.1)	(18.0)	(18.6)	(17.4)	(17.6)	
product	17	35	79	190	474	1,003	3,356	5,152	
	(40.3)	(46.7)	(43.3)	(39.4)	(39.4)	(46.0)	(47.1)	(45.6)	
mixed	11	25	76	214	503	751	2,507	4,086	
	(27.1)	(33.3)	(41.9)	(44.4)	(41.8)	(34.4)	(35.2)	(36.2)	
unknown	7	0	0	1	10	23	23	63	
	(15.9)	(0.0)	(0.0)	(0.1)	(0.8)	(1.1)	(0.3)	(0.6)	

Table 1: Basic patent statistics: IPC sections and product vs. process

Notes: All patent counts are reported in thousands. Parentheses report shares out of all invention patents.

patent applications from 1985-1989, more than two-thirds by the late 1990s, but only 11.2% from 2015-2019. Domestic patent applications are highly concentrated in the coastal regions of China and this concentration has been increasing over time. For example, the top five locations in China ranked by the total number of patent applications from 1985-2019 are Beijing and four coastal provinces: Jiangsu, Guangdong, Zhejiang, and Shandong.²³ These five locations accounted for around 35% of all invention patent applications in China before 2005 and around 55% after 2005. The role played by the next top five locations in China – which are mainly provinces in central China – has also been increasing over time. For overseas patent applications, concentration is even higher. For instance, applications from the top five overseas locations – Japan, USA, Germany, South Korea, and Taiwan – represent around 80% of all overseas patent applications after 2000. As we document in detail below, there is also a significant shift over time away from patenting by overseas applicants and toward patenting by domestic applicants.

C Traditional measures of patent quality in China

Number of forward citations. Figure 1 shows how the average number of forward citations that CNIPA invention patents receive (from other CNIPA patents) changes over time. From 2000 onward, the majority of a patent's forward citations are received within 6 years of appli-

²³Locations within China are classified as administrative divisions at the province level. Most of these are provinces (e.g., Jiangsu, Guangdong) but some are direct-administered municipalities (e.g., Beijing, Shanghai).

	'85-'89	'90-'94	'94-'99	'00-'04	'05-'09	'10-'14	'15-'19	all years
	(ii) inven	tion pater	nts by do	mestic vs	. overseas	5		
domestic	18	40	53	183	669	1,661	6,330	8,953
	(44.2)	(52.5)	(28.9)	(37.9)	(55.6)	(76.2)	(88.8)	(79.3)
overseas	23	36	129	299	534	519	796	2,336
	(55.8)	(47.5)	(71.1)	(26.1)	(44.4)	(23.8)	(11.2)	(20.7)
	(ii) do	mestic in	vention p	atents by	region			
top 1-5 Chinese locations	6	14	19	81	378	946	3,407	4,851
	(34.6)	(35.9)	(36.5)	(44.2)	(56.4)	(56.9)	(53.8)	(54.2)
top $6-10$ Chinese locations	4	7	10	40	132	330	$1,\!408$	1,929
	(19.6)	(17.7)	(18.3)	(22.0)	(19.7)	(19.8)	(22.2)	(21.5)
other Chinese locations	8	18	24	62	160	386	1,514	2,172
	(45.8)	(46.4)	(45.2)	(33.8)	(23.9)	(23.2)	(23.9)	(24.3)
	(iii) ov	erseas inv	vention p	atents by	region			
top 1-5 overseas locations	15	26	98	235	429	421	632	1,857
	(67.1)	(72.2)	(76.3)	(78.7)	(80.4)	(81.0)	(79.4)	(79.5)
top $6-10$ overseas locations	5	6	20	40	63	55	82	271
	(19.7)	(17.0)	(15.4)	(13.5)	(11.7)	(10.6)	(10.3)	(11.6)
other overseas locations	3	4	11	23	42	43	82	208
	(13.3)	(10.8)	(8.3)	(7.8)	(7.9)	(8.3)	(10.3)	(8.9)

Table 2: Basic patent statistics: geography

Notes: All patent counts are reported in thousands. Parentheses report shares out of all invention patents (panel (i)), domestic invention patents (panel (ii)), and overseas invention patents (panel (iii)) in the same time period. In panels (ii) and (iii), locations are ranked based on total invention patent applications, 1985-2019. Top Chinese locations in order: Jiangsu, Guangdong, Beijing, Zhejiang, Shandong, Anhui, Shanghai, Sichuan, Hubei, and Shaanxi. Top overseas locations in order: Japan, USA, Germany, South Korea, Taiwan, France, Netherlands, Switzerland, UK, and Sweden.



Figure 1: Average number of forward citations per patent

Notes: This figure shows the average number of forward citations that CNIPA invention patents receive from other CNIPA invention patents in each year. The bars are colored according to the time between the application date of the cited patent and the date that the citation is made.

cation and almost all citations are received within 8 years. The average patent applied for in 2000 receives around one citation within 6 years and this average is fairly constant until 2010. After 2010, however, this average declines steadily. This is largely due to censoring, since we only observe backward citations made by patents that are granted up to 2019, which highlights one of the key limitations of using forward citation counts as a measure of patent quality.²⁴

Citation centrality. A common approach to measuring patent quality is to compute measures of centrality based on citation linkages. To do this, we first define a directed graph where a node in the graph is a 3-digit IPC category and the strength of the link from node A to node B is the share of backward citations by patents in node B that cite patents in node A. Based on this graph, we then compute the eigenvector centrality of each node.²⁵ Figure 2 shows how these centrality measures compare for patents from 2000-2010 versus patents from 2011-2019. We observe that eigenvector centralities are very constant over time – technology classes that are central from 2000-2010 also tend to be central from 2011-2019. Examples of IPC categories with high centrality scores are G01 (Measuring), G06 (Computing), H01 (Electric Elements), H04 (Electric Communication Techniques), and A61 (Medical/Veterinary Science). On the other hand, examples of IPC categories with low centrality scores are A22 (Butchering), A42

 $^{^{24}}$ There is also censoring in the earlier years of our sample, as we observe fewer backward citations per patent application for patents that are applied for before 2000. Hence, patents applied for in earlier years only receive citations much further into the future.

²⁵The eigenvector centrality of a node *i* is defined recursively as $\lambda c_i = \sum_{j=1}^N w_{ij}c_j$, where *N* is the number of nodes in the graph, w_{ij} is the weight of the link from node *i* to *j*, and λ is the largest eigenvalue of the weighted adjacency matrix (a matrix with *ij*-element equal to w_{ij}).



Figure 2: Eigenvector centralities by IPC 3-digit category

Notes: This figure shows the eigenvector centralities in a directed graph where a node is a 3-digit IPC code and the strength of the link from node A to node B is the share of backward citations by patents in node B that cite patents in node A. The x-axis displays centrality scores for patents from 2000-2010 and the y-axis displays centrality scores of patents from 2011-2019.

(Headwear), C13 (Sugar), and D07 (Ropes).

Legal status changes. After a patent application is submitted, various events can occur that alter its legal status. We can then categorize patents into eight mutually exclusive cases based on the most common legal status change events. First, patents that have not been granted fall into one of four categories: (i) those that have not been examined; (ii) those that have been examined but that have been rejected; (iii) those that have been examined, were not rejected, but were withdrawn before being granted; and (iv) those that were examined, have not been rejected or withdrawn, and have not been granted. Second, patents that have been granted also fall into one of four categories: (i) those that were terminated because of the failure to pay a renewal fee; (ii) those that were terminated for reasons other than unpaid renewal fees and expiration; (iii) those that were terminated due to expiration; and (iv) those that were still active as of the latest date in our data sample.

Figure 3 shows the share of patents that fall into each of these eight categories in each year, by the main IPC section of each patent. Patent grant rates are between 40-60% in most cases. We observe increasing grant rates in all IPC sections throughout the 1990s, followed by falling grant rates from 2000 to 2015. Many ungranted patents never request examination and many are withdrawn without being rejected. The patent rejection rate also increases steadily from 2000 to 2015. For example, across all IPC sections, the rejection rate was less than 5% in 2000 but around 20% in 2015. Amongst granted patents, the primary reason for termination of patent rights before the expiration of a patent is the non-payment of renewal fees. For example, more than two-thirds of the patents that were applied for in 2000 and that were eventually



Figure 3: Patent shares by legal status

Notes: This figure shows the shares of invention patent applications within each IPC section and year based on eight mutually-exclusive legal status events.

granted were terminated before expiration for this reason. This share falls steadily from 2000 to 2015, although this is due in part to censoring, since the average non-payment event occurs around ten years after the application date of a patent.

Timing of legal status change events. Figure 4 shows the average number of years that it takes for examination, withdrawal, granting, and rejection to occur after the application date of a patent. An important observation is that the average time taken for a patent to be granted falls steadily from an average of around 6 years in 1995 to around 3 years in 2015, a trend that is consistent across all IPC sections. This is due in part to the fact that the time taken for a patent to receive examination also falls after 2005.

Number and length of claims. Panel (a) of Figure 5 shows the average number of independent claims per patent application by patent application year, while panel (b) shows the average number of Chinese characters per claim for the average patent application. Throughout the 1990s, the average number of independent claims per patent rises in all IPC sections. This is particularly prominent in sections G (Physics) and H (Electricity), where the average number of independent claims grows from just under two in 1991 to around three in 2000. Beginning in the early 2000s, however, we observe a steady decline in this average across all IPC sections. On the other hand, the average length of each patent claim (number of Chinese characters) rises steadily throughout our sample period in all IPC sections, increasing by between 50-100% from 1990 to 2019.

Priority filings at the USPTO. Figure 6 shows the share of CNIPA invention patents that have a related priority filing at the USPTO. This share is typically very small – below 1% in most IPC sections throughout our sample – although there are significant increases in sections G (Physics) and H (Electricity). In these two sections, the share of patents with USPTO filings grows from around 0.3% in 2000 to 2.2% and 3.0% respectively by 2012. There is a smaller increase in this share in section F (Mechanical Engineering) as well.



Figure 4: Timing of legal status events

Notes: This figure shows the average time between the application date of a patent and the date that the patent is examined, withdrawn, granted, and rejected (conditional on each of these events occurring for a patent). Average durations are shown by IPC section and application year.

D Supplementary tables for firm TFP and output regressions

This section provides supplementary information about data and results for the firm-level regressions discussed in section 5.2 of the main text. Table 1 provides summary statistics about the regression sample. Table 2 provides results from estimating equation (5.1) where patent stocks at t include only patents applied for between t-4 and t ("new patents"). Table 3 provides results from estimating equation (5.1) where firm fixed effects are added to each regression.



Figure 5: Number and length of claims per patent

Notes: Panel (a) shows the average number of independent claims associated with a patent. Panel (b) shows the average number of Chinese characters per independent claim for the average patent. All results are shown by IPC section and application year.



Figure 6: Share of patents with USPTO priority filings

Notes: This figure shows the share of CNIPA patents within each IPC section and application year that also have a priority filing at the USPTO.

	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007
no. of firms	110,982	112,421	116,779	128,742	141,451	159,778	235,210	231,200	257,313	290,280
no. of firms with:										
patents	1,605	1,789	2,089	2,444	3,308	4,338	6,303	7,492	9,790	$12,\!875$
active patents	188	198	260	381	495	738	1,391	1,961	2,692	3,562
top-1% important patents	29	30	36	52	77	127	197	278	420	612
top-2% important patents	51	52	69	102	163	247	415	554	795	$1,\!149$
top-5% important patents	117	138	175	233	379	571	951	1,278	1,789	2,492
top-10% important patents	200	241	321	442	681	1,025	1,637	2,137	2,984	4,105
top-25% important patents	429	525	692	913	1,348	1,919	3,002	3,823	5,217	7,112
top-50% important patents	737	863	1,073	1,355	1,958	2,684	4,172	5,221	7,021	9,417
top-1% cited patents	28	40	53	73	106	147	235	304	400	516
top-2% cited patents	51	69	90	123	185	259	432	542	717	921
top-5% cited patents	118	150	187	244	353	501	819	1,032	1,390	1,883
top-10% cited patents	207	246	319	421	630	895	1,343	1,627	2,081	2,780
top-25% cited patents	383	454	614	809	1,221	1,697	2,577	3,112	4,073	5,416
top- 50% cited patents	663	759	967	1,212	1,748	2,368	3,573	4,319	5,653	7,490
mean (conditional on positive) of:										
patents	2.12	2.29	2.52	3.19	4.15	5.50	6.47	7.89	8.82	9.57
active patents	1.94	1.93	2.01	1.85	1.87	2.71	3.47	3.84	4.26	4.72
top-1% important patents	1.10	1.17	1.14	1.12	1.34	1.34	1.33	1.41	1.44	1.91
top-2% important patents	1.24	1.27	1.29	1.30	1.40	1.49	1.50	1.56	1.61	1.96
top-5% important patents	1.36	1.37	1.39	1.54	1.67	1.81	1.88	2.00	2.12	2.34
top-10% important patents	1.42	1.48	1.52	1.71	1.93	2.18	2.33	2.54	2.69	2.91
top-25% important patents	1.66	1.75	1.87	2.24	2.64	3.16	3.45	3.87	4.15	4.41
top-50% important patents	1.80	1.95	2.17	2.87	3.67	4.70	5.22	6.04	6.55	6.94
top-1% cited patents	1.32	1.33	1.51	2.34	2.91	3.28	3.11	3.35	3.32	3.25
top-2% cited patents	1.37	1.41	1.61	2.45	3.17	3.49	3.40	3.72	3.80	3.71
top-5% cited patents	1.54	1.63	1.82	2.58	3.21	3.57	3.61	3.99	4.11	4.18
top-10% cited patents	1.67	1.76	2.00	2.72	3.33	3.85	3.94	4.29	4.54	4.60
top-25% cited patents	1.75	1.85	2.12	2.86	3.51	4.25	4.42	4.99	5.36	5.57
top-50% cited patents	1.78	1.91	2.16	2.87	3.62	4.52	4.85	5.57	6.08	6.41

Table 1: Summary statistics for the firm regression sample

Notes: This table provides summary statistics for the sample of firms in the NBS data that have information on both TFP and output from 1998-2007.

(a) outcome: log TFP ($\times 100$); fixed effe	ects: indust	try-year		
<i>p</i> =	2	5	10	25	50
i. stock of new top- $p\%$ important patents, $\beta^{topimp,p}$	0.47	0.37	0.18	0.05	0.03
	(4.55)	(5.85)	(5.65)	(4.61)	(3.82)
ii. stock of new top- $p\%$ cited patents, $\beta^{topcit,p}$	0.15	0.08	0.05	0.02	-0.00
	(2.64)	(2.72)	(2.77)	(1.59)	(-0.43)
iii. stock of new active patents, β^{active}	-0.05	-0.11	-0.15	-0.13	-0.12
	(-1.16)	(-2.52)	(-3.13)	(-2.68)	(-2.41)
iv. has patents, γ	3.66	3.63	3.64	3.66	3.68
	(22.31)	(22.13)	(22.19)	(22.36)	(22.49)
observations (m)	1.78	1.78	1.78	1.78	1.78
R^2	0.74	0.74	0.74	0.74	0.74
adjusted R^2	0.74	0.74	0.74	0.74	0.74
(b) outcome: log output (×10	0); fixed ef	fects: indu	stry-year		
<i>p</i> =	2	5	10	25	50
i. stock of new top- $p\%$ important patents, $\beta^{topimp,p}$	5.03	4.86	2.54	0.72	0.30
	(13.72)	(21.87)	(22.12)	(17.34)	(11.04)
ii. stock of new top- $p\%$ cited patents, $\beta^{topcit,p}$	0.31	0.18	0.13	-0.12	-0.21
	(1.57)	(1.74)	(0.97)	(-3.44)	(-6.18)
iii. stock of new active patents, β^{active}	1.01	0.16	-0.33	0.268	0.67
	(7.34)	(1.03)	(-2.00)	(1.62)	(3.84)
iv. has patents, γ	110.66	110.20	110.29	110.64	110.94
	(191.51)	(190.54)	(190.77)	(191.54)	(192.18)
observations (m)	1.78	1.78	1.78	1.78	1.78
R^2	0.15	0.15	0.15	0.15	0.15
adjusted R^2	0.15	0.15	0.15	0.15	0.15

Table 2: Regressions of firm TFP and output on new patent stocks

Notes: This table reports the results of estimating equation (5.1) by OLS where the dependent variable is log TFP in panel (a) and log output in panel (b). Patent stocks at t include only patents applied for between t - 4 and t. t-statistics are shown in parentheses.

a) outcome: log TFP ($\times 100$); fixed effects: industry-year, firm								
<i>p</i> =	2	5	10	25	50			
i. stock of top- $p\%$ important patents, $\beta^{topimp,p}$	0.09	0.07	0.02	0.01	0.01			
	(0.96)	(1.26)	(0.73)	(0.59)	(0.93)			
ii. stock of top- $p\%$ cited patents, $\beta^{topcit,p}$	-0.02	-0.01	-0.00	-0.00	-0.01			
	(-0.33)	(-0.32)	(-0.27)	(-0.35)	(-0.77)			
iii. stock of active patents, β^{active}	0.01	0.00	0.00	0.00	0.00			
	(0.42)	(0.12)	(0.10)	(0.13)	(0.06)			
iv. has patents, γ	0.78	0.78	0.78	0.79	0.79			
	(3.32)	(3.28)	(3.31)	(3.33)	(3.33)			
observations (m)	1.69	1.69	1.69	1.69	1.69			
R^2	0.89	0.89	0.89	0.89	0.89			
adjusted R^2	0.86	0.86	0.86	0.86	0.86			
(b) outcome: log output ($\times 100$);	fixed effec	ets: indus	try-year, i	firm				
<i>p</i> =	2	5	10	25	50			
i. stock of top- $p\%$ important patents, $\beta^{topimp,p}$	0.71	0.83	0.48	0.14	0.07			
	(3.57)	(6.39)	(6.74)	(5.45)	(3.89)			
ii. stock of top- $p\%$ cited patents, $\beta^{topcit,p}$	0.22	0.15	0.10	0.03	-0.01			
	(1.78)	(2.38)	(2.49)	(1.41)	(-0.51)			
iii. stock of active patents, β^{active}	-0.05	-0.15	-0.20	-0.21	-0.19			
	(-1.11)	(-3.01)	(-3.99)	(-3.88)	(-3.52)			
iv. has patents, γ	22.94	22.82	22.80	22.87	22.92			
	(43.87)	(43.60)	(43.57)	(43.72)	(43.85)			
observations (m)	1.69	1.69	1.69	1.69	1.69			
R^2	0.86	0.86	0.86	0.86	0.86			
adjusted R^2	0.83	0.83	0.83	0.83	0.83			

Table 3: Regressions of firm TFP and output on patent stocks with firm fixed effects

Notes: This table reports the results of estimating equation (5.1) by OLS where the dependent variable is log TFP in panel (a) and log output in panel (b), with firm fixed effects added to each regression. *t*-statistics are shown in parentheses.

(a) outcome: log TFP ($\times 100$); fixed effects: industry-year year								
<i>p</i> =	2	5	10	25	50			
i. stock of top-p% important patents, $\beta^{topimp,p}$	0.27	0.26	0.16	0.06	0.08			
	(2.55)	(3.96)	(4.58)	(3.88)	(4.32)			
ii. stock of top- $p\%$ important patents (naive), $\tilde{\beta}^{topimp,p}$	0.62	0.15	0.03	-0.01	-0.05			
	(3.94)	(2.51)	(0.86)	(-0.60)	(-2.95)			
iii. stock of top- $p\%$ cited patents, $\beta^{topcit,p}$	-0.27	-0.08	0.02	0.03	0.03			
	(-2.35)	(-1.14)	(0.47)	(1.43)	(2.13)			
iv. stock of active patents, β^{active}	-0.05	-0.08	-0.10	-0.09	-0.08			
	(-2.56)	(-3.33)	(-3.64)	(-3.25)	(-2.89)			
v. has patents, γ	3.65	3.64	3.63	3.64	3.60			
	(22.31)	(22.17)	(22.07)	(22.06)	(21.85)			
observations (m)	1.78	1.78	1.78	1.78	1.78			
R^2	0.74	0.74	0.74	0.74	0.74			
adjusted R^2	0.74	0.74	0.74	0.74	0.73			
(b) outcome: log output ($\times 100$); fixed effects: industry-year								
<i>p</i> =	2	5	10	25	50			
i. stock of top- $p\%$ important patents, $\beta^{topimp,p}$	5.31	5.94	3.72	1.80	2.41			
	(14.07)	(26.05)	(30.78)	(33.21)	(35.64)			
ii. stock of top- $p\%$ important patents (naive), $\tilde{\beta}^{topimp,p}$	0.59	-2.22	-2.34	-1.73	-2.17			
	(1.07)	(-10.29)	(-21.26)	(-28.90)	(-33.95)			
iii. stock of top- $p\%$ cited patents, $\beta^{topcit,p}$	0.20	2.81	3.24	2.11	1.17			
	(0.49)	(11.47)	(21.83)	(26.46)	(23.55)			
iv. stock of active patents, β^{active}	0.32	-0.12	0.04	0.29	0.37			
	(4.37)	(-1.39)	(0.44)	(2.98)	(3.76)			
v. has patents, γ	110.61	109.72	109.02	108.65	108.39			
	(191.40)	(189.39)	(187.82)	(186.92)	(186.29)			
observations (m)	1.78	1.78	1.78	1.78	1.78			
R^2	0.15	0.15	0.15	0.15	0.15			
adjusted R^2	0.15	0.15	0.15	0.15	0.15			

Table 4: Regressions of firm TFP and output on patent stocks, baseline vs. naive importance

Notes: This table reports the results of estimating equation (5.1) by OLS where the dependent variable is log TFP in panel (a) and log output in panel (b), with the stock of top patents based on the naive measure of importance in equation (3.1) added to each regression. *t*-statistics are shown in parentheses.

E Results for granted USPTO utility patents

This section presents results for granted USPTO utility patents. Figure 1 shows the share of patents by patentee type, which is identified based on the address of the main applicant. Patents identified as belonging to "US applicants" have only US applicants, those identified as belonging to "overseas non-CN applicants" have at least one non-China overseas applicant and no applicants from China, and those identified as belonging to "CN applicants" have at least one applicant from China. Figure 2 shows how the innovativeness and crowdedness of patenting has changed over time, as discussed under Fact 1. Figure 3 shows how relative patent importance within IPC-years by patentee type has changed over time, as discussed under Fact 5. Figure 4 shows our estimates of memory and secondary IPC importance, as discussed under Fact 7



Figure 1: Patent shares by patentee type — granted USPTO patents



Figure 2: Innovativeness and crowdedness of patenting – granted USPTO patents

Figure 3: Relative patent importance by patentee type – granted USPTO patents





Figure 4: Memory and importance of secondary IPCs – granted USPTO patents



↓

Download ZEW Discussion Papers:

https://www.zew.de/en/publications/zew-discussion-papers

or see:

https://www.ssrn.com/link/ZEW-Ctr-Euro-Econ-Research.html https://ideas.repec.org/s/zbw/zewdip.html

IMPRINT

ZEW – Leibniz-Zentrum für Europäische Wirtschaftsforschung GmbH Mannheim

ZEW – Leibniz Centre for European Economic Research

L 7,1 · 68161 Mannheim · Germany Phone +49 621 1235-01 info@zew.de · zew.de

Discussion Papers are intended to make results of ZEW research promptly available to other economists in order to encourage discussion and suggestions for revisions. The authors are solely responsible for the contents which do not necessarily represent the opinion of the ZEW.