

## **Cover Sheet**

**Session Title: The Roots of Innovation**

**Paper Title: The Dynamics of Open Source Contributors**

Josh Lerner (corresponding author)

Rock Center

Harvard Business School

Boston, MA 02163 USA

Email: [josh@hbs.edu](mailto:josh@hbs.edu)

Phone: (617) 495-6065

Fax: (617) 495-3817

Parag A. Pathak

Baker Library 420E

Harvard Business School

Boston, MA 02163 USA

Email: [ppathak@fas.harvard.edu](mailto:ppathak@fas.harvard.edu)

Phone: (617) 864 6211

Jean Tirole

Universite des Sciences Sociales

Manufacture des Tabacs

Aile Jean-Jacques Laffont

Accueil MF 404

F-31000 Toulouse FRANCE

Email: [tirole@cict.fr](mailto:tirole@cict.fr)

Phone: 33(0) 5.61.12.85.89

Fax: 33(0) 5.61.12.86.37

# **The Dynamics of Open Source Contributors**

Josh Lerner, Parag A. Pathak, and Jean Tirole\*

There are substantial differences between open source projects and traditional innovative efforts in private firms. Private firms usually pay their workers, direct and manage their efforts, and control the output and intellectual property created. In an open-source project, however, a body of original material is made publicly available for others to use, under certain conditions. Contributions to open source projects are made by a diverse array of individual contributors and for-profit corporations, who in many cases must agree to make all enhancements to the original material widely available for nominal cost.

This paper empirically examines the dynamics of contributions to open source software projects. We show that the share of corporate contributions in a sample of approximately 100 open source projects tracked between 2001 and 2004 is greater in larger and growing projects.

## **I. Background<sup>1</sup>**

The decision to contribute without pay to freely available software may seem mysterious to economists. However, the unpaid programmer working on an open source software development project faces a variety of benefits and costs. The programmer incurs an opportunity cost of time, which can manifest itself in different ways. For example, a programmer who works as an independent on open source projects forgoes

the monetary compensation that could otherwise be earned by working for a commercial firm or a university. For a programmer with a commercial company, university or research lab affiliation, the opportunity cost of working on open source software comes from not focusing on other tasks. For example, the academic's research output may sag and the student's progress towards a degree may slow down.

Several short-or long-run benefits may counter these costs (see Josh Lerner and Jean Tirole, 2002, for a fuller discussion). First, open source programmers may improve rather than reduce their performance in paid work. This outcome is particularly relevant for system administrators looking for specific solutions for their company. Second, the programmer may find intrinsic pleasure if choosing a “cool” open source is more fun than a routine task set by an employer. Third, in the long run, open source contributions may lead to future job offers, shares in commercial open source-based companies, or future access to the venture capital market, and last (but not least) ego gratification from peer recognition. Of course, different programmers may put different values on monetary or personal payoffs, and on short-term or long-term payoffs.

Economic theory suggests that long-term incentives are stronger under three conditions: 1) the more visible the performance to the relevant audience (peers, labor market, and venture capital community); 2) the higher the impact of effort on performance; 3) the more informative the performance about talent (for example, Bengt Holmström, 1999). The first condition gives rise to what economists call “strategic complementarities.” To have an “audience,” programmers will want to work on software projects that will attract a large number of other programmers. This argument suggests the possibility of multiple equilibria. The same project may attract few programmers

because programmers expect that other programmers will not be interested; or it may flourish as programmers (rationally) have faith in the project.

When we consider the delayed rewards of working on an open source project, the ability to signal a high level of competence may be stronger in the open source mode for three reasons. First, in an open source project, outsiders can see the contribution of each individual, whether that component “worked,” whether the task was hard, if the problem was addressed in a clever way, whether the code can be useful for other programming tasks in the future, and so forth. Second, the open source programmer takes full responsibility for the success of a subproject, with little interference from a superior, which generates information about ability to follow through with a task. Finally, since many elements of the source code are shared across open source projects, more of the knowledge they have accumulated can be transferred to new environments, which makes programmers more valuable to future employers. To compare programmers' incentives in the open source and proprietary settings, we need to examine how the features of the two environments shape incentives. From the standpoint of the individual, commercial projects typically offer better current compensation than open source projects, because employers are willing to offer salaries to software programmers in the expectation that they will capture a return from a proprietary project.

Commercial companies may interact with an open source project in a number of ways. While improvements in the open source software are not appropriable, commercial companies can benefit if they also offer expertise in some proprietary segment of the market which is complementary to the open source program. Firms may temporarily encourage their programmers to participate in open source projects to learn about the

strengths and weaknesses of this development approach. For-profit firms may compete directly with open source providers in the same market. Firms may also be able to learn about potential employees when their staff interacts with open source programmers. Finally, commercial companies may interface with the open source world because it generates good public relations with programmers and customers.

A for-profit firm that seeks to provide services and products which are complementary to the open source product, but are not supplied efficiently by the open source community, can be referred to as “living symbiotically.” IBM, which has made open source software into a major focus for its consulting business, exemplifies this approach. A commercial company in this situation will want to have extensive knowledge about the open source movement, and may even want to encourage and subsidize open source contributions, both of which may cause it to allocate some programmers to the open source project. Because firms do not capture all the benefits of the investments in the open source project, however, the free-rider problem often discussed in the economics of innovation should apply here as well. Subsidies by commercial companies for open source projects should remain somewhat limited.

The code release strategy arises when companies release some existing proprietary code and then create a governance structure for the resulting open source development process. This strategy is to giving away the razor (the released code) to sell more razor blades (for instance, the related consulting services that IBM and HP hope to provide). In general, it will make sense for a commercial company to release proprietary code under an open source license if the increase in profit in the proprietary complementary segment offsets any profit that would have been made in the primary

segment, had it not been converted to open source. Thus, the temptation to go open source is particularly strong when the product is lagging behind the leader and making few profits, but the firm sees a possibility that if the released code becomes the center of an open source project and is utilized more widely, the profitability of the complementary segment will increase.

## **II. The Sample**

We built a panel data set of the contributors to approximately 100 open source projects (for full details on the dataset, see Josh Lerner, Parag A. Pathak and Jean Tirole (2006)). These projects are stratified to over-represent the largest open source projects. We extract the contributors to the project in each new official version of the program has been released, using a variety of text editing tools.

Table 1 summarizes the projects, and highlights that they differ considerably in their size and other characteristics. Open source projects periodically introduce new versions. The number of versions introduced between the beginning of data collection and July 2004 varies between one and twenty.<sup>2</sup> For each project we obtained information on the projects from SourceForge, press searches, and project websites. Key information includes the type of license of the project, whether venture capitalists had funded the company, and whether a corporate released some of its code as an open source project.

For each project, we opened the Tape Archive (known as “tarball”) to count the number of distinct references to each individual. The archive preserves information such as user and group permissions, dates, and directory structures. Open source projects are scrupulous about keeping track of contributors, which reflects the fact that giving credit

to authors is essential in the open source movement. This principle is included as part of the nine key requirements in the “Open Source Definition.”<sup>3</sup> This point is also emphasized by Eric Raymond (1999), who points out “surreptitiously filing someone’s name off of a project is, in cultural context, one of the ultimate crimes.” This point was also emphasized in our conversations with open source project managers and SourceForge officials. Each project release was then associated with a set of e-mails that appeared in the archive.<sup>4</sup>

We aimed to distinguish individuals who were contributing code on their own behalf from those doing so as part of their employment. Our approach divided the contributors into five classes based on their e-mail addresses. These are corporate employees, individual hobbyists, and three classes of otherwise other contributors: unidentified international contributors, and those from organizations with top-level domains (TLDs) denoted “.org” and “.net,” which frequently denote non-profit and technical web sites. We included as corporate contributors all those with a “.com” address, excluding those sites used primarily as e-mail mailboxes, Internet Service Providers, or portals (e.g., “hotmail.com”). We also included overseas addresses that are associated with corporations (for instance, “.co.uk” and “.caldera.de”). We included as hobbyists contributions by individuals affiliated with universities and governments (again, employing both addresses with TLDs such as “.edu” and overseas domains like “.umontreal.ca”), as well as those who made contributions from addresses associated with portals, ISPs, and mailboxes.<sup>5</sup> The remaining categories—those from TLDs “.org” and “.net,” as well as the remaining international domains—were not classified in either

category, but rather treated separately, because we were not able to readily assign them (see Table 2).

### **III. Analysis of Project Contributions**

Our initial analysis seeks to understand the distribution of contributions to open source project by class of contributor, focusing on contributions by corporations and “hobbyists.” Table 3 presents some breakdowns, using the most direct measure: the number of contributions by each class of contributor for various classes of projects. This table presents the proportion of all contributions that are corporate.<sup>6</sup> The table also presents the result of F- and t-tests of the significance of the reported differences.

The table shows that corporate contributions are more frequent for larger projects. The share of corporate contributions is twice as larger in the largest quartile of projects than in the bottom. The pattern is similar, though somewhat less dramatic, when we compare the versions divided into quartiles based on their growth rates, defined here as the difference between the number of lines of code in the current and previous version. Both differences are highly statistically significant.

Patterns regarding license type and venture capital backing are less sharp. The share of corporate contributions is lowest among those projects with the most restrictive licenses (see Lerner and Tirole (2005b) for a discussion of our typology of license types), but there is no consistent relationship between license strength and corporate contributions. Corporate contributions are more common when venture capitalists have funded a company that is focusing on the open source project, but this difference is not



statistically significant. Finally, consistent with the results regarding project size above, corporate contributions are more common in later versions of projects.

#### **IV. Conclusions**

This paper presents only the beginnings of understanding of cross-sectional and time-series patterns of contributions to open source projects. In the approximately 100 software projects we track from 2001-2004, we have shown that the share of corporate contributions is much larger in large and growing projects. In the companion paper to this one, we develop a theoretical rational for these patterns and explore them in more depth.

## References

Holmström, Bengt, 1999, “Managerial Incentive Problems: A Dynamic Perspective,” *Review of Economic Studies*, 66, 169-182.

Lerner, Josh, and Jean Tirole, 2002, “Some Simple Economics of Open Source,” *Journal of Industrial Economics*, 52, 197-234.

Lerner, Josh, and Jean Tirole, 2005a, “The Economics of Technology Sharing: Open Source and Beyond,” *Journal of Economic Perspectives*, 19, 99-120.

Lerner, Josh, and Jean Tirole, 2005b, “The Scope of Open Source Licensing,” *Journal of Law, Economics, and Organization*, 21, 20-56.

Lerner, Josh, Parag A. Pathak, and Jean Tirole, 2006, “Open Source Contributors,” *Mimeo, Harvard Business School*.

Raymond, Eric, 1999, *The Cathedral and the Bazaar: Musings on Linux and Open Source by an Accidental Revolutionary*, Cambridge, O'Reilly.

**Table 1: Project Characteristics**

Sample:			
20 large projects	Started tracking:	05-2001	
78 randomly selected projects	Ended tracking:	07-2004	
98 total projects			
	min	median	max
Lines of Source Code <sup>1</sup>	1,253	81,671	4,032,921
	Wings 3D	jEdit	Linux
Absolute Change in Source Code <sup>1</sup>	-145,395	18,951	1,628,979
	AOLServer	Licq	Linux
Number of New Versions	1	8	20
	Dev-C++	Koffice	Wine
	Imprints	BZFlag	
	Kxlcq	glibc	
	KDE		
	Restrictive	Highly Restrictive	
License Type <sup>2</sup>	74%	51%	

**Notes:**

1. Measured at the end of the sample

2. BSD is an unrestrictive license, LGPL is restrictive, and GPL is highly restrictive. Three projects changed their license

during our sample period: Sendmail, PureFTPd, Wine. We take the license post-change.

**Table 2: Characteristics of Contributions**

	min	median	max
Number of Contributors <sup>1</sup>	1	67	3,521
	CsvJdbc	Miranda ICQ client	Linux
	EverQuest	Gstreamer	
	JFS		
Absolute Growth in Contributors	-343	16	1,174
	JFS	Common C++ Libraries	Linux
		Gstreamer	
		Jext	
% Growth in Contributors	-100%	36%	4200%
	JFS	Gabber	PPTP Client
Number of Contributions <sup>1</sup>	2	374	52,607
	JFS	Cluster Infrastructure	GCC
Absolute Growth in Contributions	-1,208	80	24,611
	XFree86	ROX Desktop	GCC
% Growth in Contributions	-100%	50%	5800%

JFS

Licq

AWStats

---

---

Notes:

1. Measured at the end of the sample

**Table 3: Distribution of Corporate Contributors As a Share of All Contributors**

Size of Code Base		Growth of Code Base		License Type	
Smallest size quartile	21.4%	Smallest growth quartile	29.9%	Unrestrictive licenses	32.0%
Mid-small size quartile	22.2%	Mid-small growth quartile	26.3%	Restrictive licenses	37.1%
Mid-large size quartile	33.1%	Mid-large growth quartile	26.6%	Highly restrictive licenses	29.0%
Largest size quartile	44.2%	Largest growth quartile	43.3%		
p-Value, F- (or t-)test	0.000		0.000		0.093
Venture Backing		Version		-	
Venture-backed					
projects	35.0%	Less than version 4	5.5%		
Non-venture backed	31.6%	Version 4 to 6	24.8%		
		Version 7 to 11	38.4%		
		More than version 12	43.9%		
p-Value, F- (or t-)test	0.398		0.000		

---

<sup>\*</sup>Harvard University and NBER; Harvard University; and University of Toulouse and Massachusetts Institute of Technology. We thank for research assistance Vakha Elmurzev, Lee Gao, James Hunter, Payal Loungani, Susanna Kim, Qicheng Ma, and John Sheridan. Helpful comments were received from Tim Bresnahan, participants in the 2006 American Economics Association meetings and the Toulouse Network on Information Technology's Fall 2005 meeting, and a number of practitioners. Jeff Bates was extremely helpful in regard to SourceForge access. We thank Harvard Business School, the National Science Foundation, and the Toulouse Network on Information Technology for financial support.

<sup>1</sup>This section is based on Josh Lerner and Jean Tirole (2005a).

<sup>2</sup>We did an initial analysis using 20 SourceForge projects beginning in May 2001. In January 2002, we expanded the data collection to include the entire sample, which was tracked until July 2004.

<sup>3</sup>[http://www.opensource.org/docs/definition\\_plain.php](http://www.opensource.org/docs/definition_plain.php) (accessed December 4, 2005).

<sup>4</sup>The database went through an extensive cleaning process to remove invalid email addresses and to deal with situations where there were two email addresses from the same individual. Examples of the decisions made are in Lerner, Pathak, and Tirole (2006).

<sup>5</sup>One complication was posed by sites such as "aol.com," which are used by both corporate employees and as an e-mail service. We treat these cases as corporate contributors. We have experimented with further portioning the corporate contributors into subcategories, where cases like "aol.com" will be considered to be separate. With this further breakout of the corporate sample, the qualitative results are similar.



---

<sup>6</sup>Results looking at the ratio of contributions by corporate contributors and hobbyists generate similar results.