

Success and Geography: Evidence from open-source software

Gábor Békés¹ Julian Hinz² Miklós Koren³ Aaron Lohmann⁴

February 25, 2024

¹Central European University, KRTK, CEPR

²Bielefeld University, Kiel Institute for the World Economy

³Central European University, KRTK, CEPR, Cesifo

⁴Bielefeld University, Kiel Institute for the World Economy

Lugano, 2024

Introduction

Big Picture:

- How dispersed and unorganized developers can create great products.
- How and where good Open Source Software (OSS) is produced.

How and where good Open Source Software (OSS) is produced?

- Are there spatial frictions even though all online?
 - weightless economy – OSS no fixed costs, and no need for face-to-face interaction - pure online.
 - Geography may not significantly impact OSS development.
 - But: spread of talent globally
 - But: spatial frictions in search + matching

How and where good Open Source Software (OSS) is produced?

- Are there spatial frictions even though all online?
 - weightless economy – OSS no fixed costs, and no need for face-to-face interaction - pure online.
 - Geography may not significantly impact OSS development.
 - But: spread of talent globally
 - But: spatial frictions in search + matching

Data:

- Writing code together – Collaboration (Github)
- Using other people's code – imported dependencies (Libraries.io).

Open Source Software (OSS) is everywhere

Open Source Software (OSS) has a vast landscape, GitHub hosts over 330 million repositories.

OSS plays an important roles in

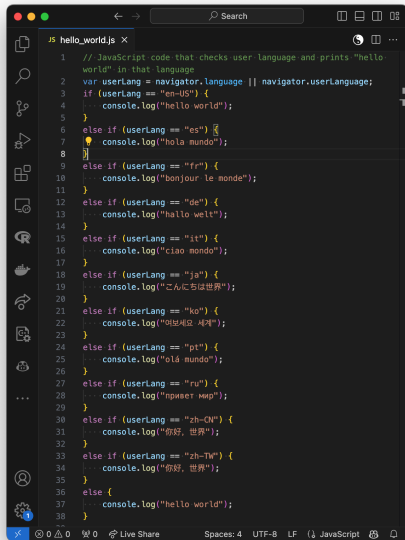
- Websites (JavaScript)
- Operating systems (Linux, Android)
- Data (R Tidyverse, Python Pandas, Julia)
- Machine Learning and AI (PyTorch, LLaMA)

OSS mostly free, but present in fee-based platforms

- Overleaf

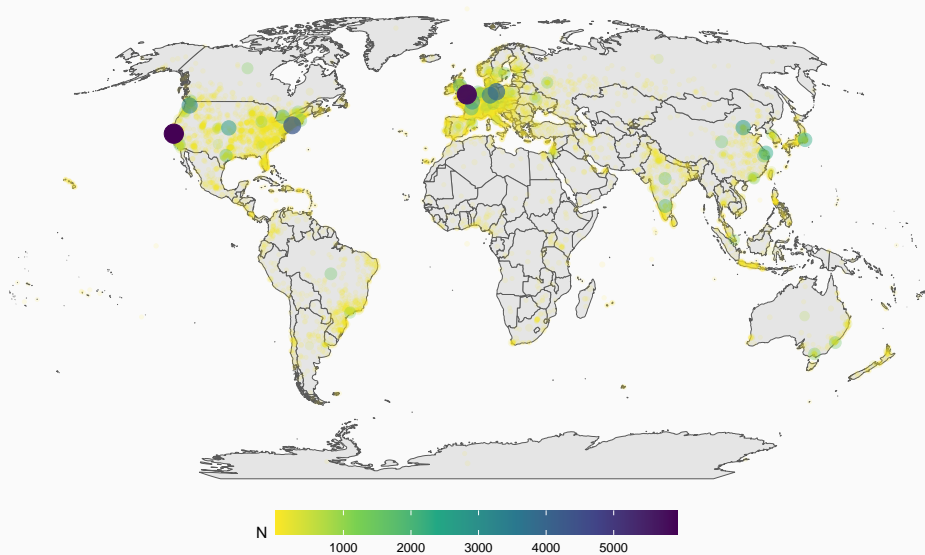
Focus on JavaScript

- JavaScript is one of the biggest programming languages
- used in web development and app development
- NPM is a package manager
- organizes packages and provides access

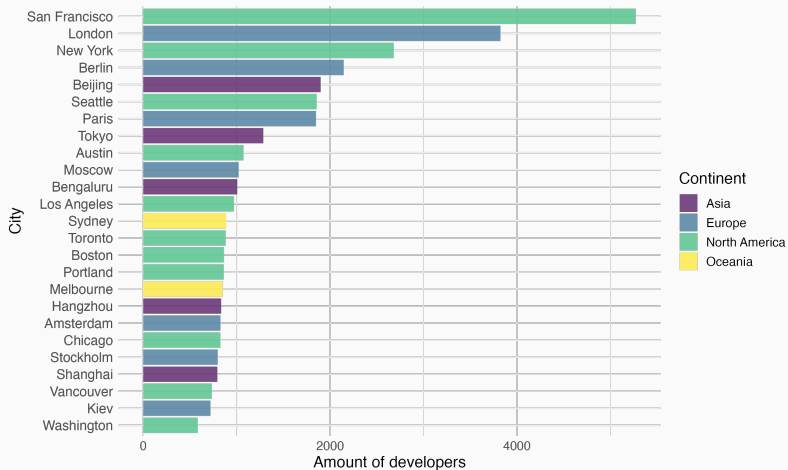


```
JS hello_world.js X
1 // JavaScript code that checks user language and prints "hello
  world" in that language
2 var userLang = navigator.language || navigator.userLanguage;
3 if (userLang == "en-US") {
4   console.log("hello world");
5 }
6 else if (userLang == "es") {
7   console.log("hola mundo");
8 }
9 else if (userLang == "fr") {
10  console.log("bonjour le monde");
11 }
12 else if (userLang == "de") {
13  console.log("hallo welt");
14 }
15 else if (userLang == "it") {
16  console.log("ciao mondo");
17 }
18 else if (userLang == "ja") {
19  console.log("こんにちは世界");
20 }
21 else if (userLang == "ko") {
22  console.log("안녕하세요 세계");
23 }
24 else if (userLang == "pt") {
25  console.log("olá mundo");
26 }
27 else if (userLang == "ru") {
28  console.log("привет мир");
29 }
30 else if (userLang == "zh-CN") {
31  console.log("你好，世界");
32 }
33 else if (userLang == "zh-Tw") {
34  console.log("你好，世界");
35 }
36 else {
37  console.log("hello world");
38 }
```


Global industry: Number of JavaScript developer per city



Dispersion and concentration: top cities per number of developers












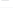

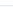
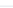


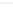

Collaboration is done mostly online

 **git-extras** Public

[Watch](#) 214 [Fork](#) 1.2k [Star](#) 16.6k

[main](#) 3 Branches 53 Tags [Add file](#) [Code](#)

 **vanpip** test(browse-ci): add unit tests (#1130) ✓ 5f19424 · 3 weeks ago 1,764 Commits

	test(git-browse): add unit tests (#1127)	last month
	feat: add reverse option to git-brv (#1123)	2 months ago
	feat: add reverse option to git-brv (#1123)	2 months ago
	fix: No longer pollute env with GREP_OPTIONS	last year
	feat: add reverse option to git-brv (#1123)	2 months ago
	test(browse-ci): add unit tests (#1130)	3 weeks ago
	Improve defaults for testing suite (#1104)	3 months ago
	Improve defaults for testing suite (#1104)	3 months ago
	test(git-authors): add unit test (#1098)	3 months ago
	maintenance: Add my name as maintainer in AUTHORS (#11...	3 months ago
	chore: add poetry to handle the tests of the git extras (#1121)	3 months ago
	feat: add reverse option to git-brv (#1123)	2 months ago
	Version 7.1.0 (#1097)	4 months ago
	Add more comprehensive dependencies (#1111)	3 months ago
	Mention initial copyright year and add contributors to copyr...	9 years ago
	makefile: Allow bypassing conflict check (#1080)	5 months ago

About

GIT utilities -- repo summary, repl, changelog population, author commit percentages and more

[git](#)

- [Readme](#)
- [MIT license](#)
- [Activity](#)
- [16.6k stars](#)
- [214 watching](#)
- [1.2k forks](#)

Report repository

Releases 22

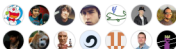
[7.1.0 \(Hauyne\)](#) Latest on Oct 29, 2023

[+ 21 releases](#)

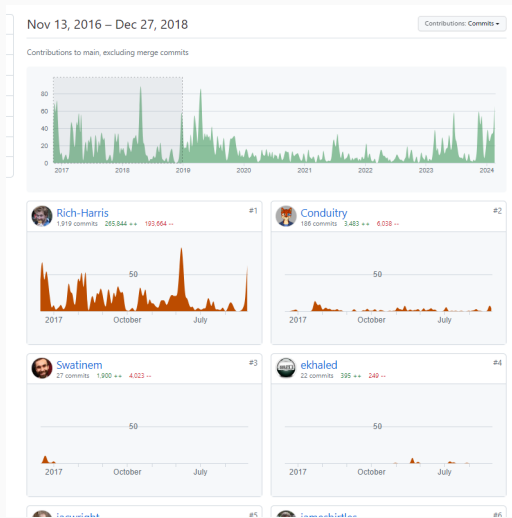
Packages

No packages published

Contributors 224



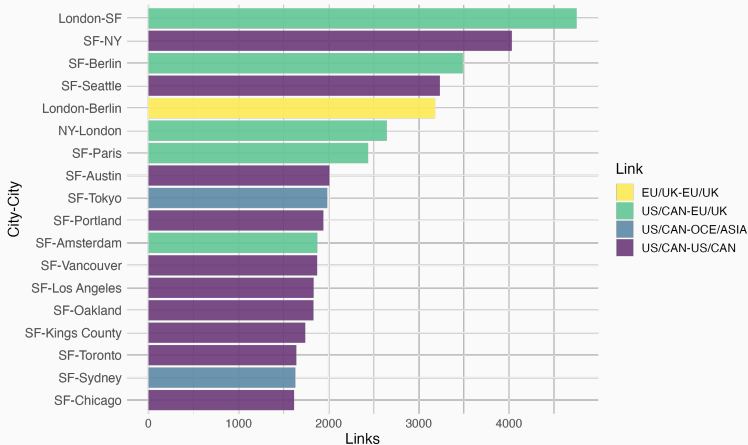
Collaboration is done mostly online



... but personal contacts still matter

- Personal meeting, esp. workplace (CEU, Oracle)
- Local community events, science parks (Xaccelerator)
- Regional events (R Ladies Auckland, VDSG Meetup, PyData Berlin)
- Conferences 1: dozens of events every month such CityJS Berlin, React Summit US,
- Conferences 2: developers directly such as Node-js fwdays23 in Kyiv, where new packages are presented.
- Learn about packages, devs: online forums, Stack Overflow, Twitter

Collaboration across cities is mostly North-North



Most frequent city-pairs for repos developed from 2 cities

- **Geographical Distance / Network formation / Agglomeration:** Chaney (2014) Bernard et al. (2019) Davis and Dingel (2019) Bailey et al. (2021)
- **Gravity: Digital:** Blum and Goldfarb (2006) Anderson et al. (2018)
- **Frictions in services:** Stein and Daude (2007) Bahar (2020)
- **Patents and science:** Bircan et al. (2021), Head et al. (2019), Jaffe et al. (1993), Singh (2008) AlShebli et al. (2018)
- **OSS:** Lerner and Tirole (2002) , Laurentsyeve (2019) Wachs et al. (2022) Fackler et al. (2023)

Open source software vs patents and academia

- R&D and patenting
 - Need machines, secrecy, often top-down
 - Distance matters in collaboration
 - More cited patents – geographically focused authors
- Science (math, academic papers)
 - Similar, but often longer projects, not open, F2F important to think and discuss
 - Distance matters in collaboration
 - Major role of top Universities / Centers

- OSS and data
- The role of space in collaboration
 - Gravity
 - Success

Open source software data

Open Source vocabulary

- Package: A unit of software, provision of a (bundle of) functionality
- Project: A software project offering solution to a use case. Typically one package, but may be more.
- Repository: A storage for one project (what we observe)
- Commit: The smallest unit of contribution
- Git: Distributed version control system for software projects
- GitHub: A platform to collaboratively work on software projects
- Dependency: An imported package that provides a functionality

Collaboration — Working on the same code with others

- GHTorrent: Tracks metadata on GitHub usage
- Commits, locations and user organisations
- Row: One commit from a developer to a repository
- Focus on links: binary if a developer committed at all to a repository

Dependencies — Sourcing of intermediate inputs

- Libraries.io: Tracks data on single software repositories
- dependency linkages
- Row: An imported dependency (package) to repo 1 from repo 2
- Can be mapped to repositories on GitHub

Scope of data

- Data coverage: 2013 – 2019
- We know location as city for developers
- Contributions by 217K developers,
- 300K repos
- 17% of repos have multiple developers (ie have collaboration)
- 70K organizations, with 120K developers

Sample design: exclude later arrival, bug-fixers

We focus on collaborating partners, who are likely to have interaction, joint decisions. Exclude

1. Bugfixers – as external "consultants" who come in help solve a problem
 - Less than 4 commits or 1% of commits — less than 10 commits total
2. Late arrivals – developers who take over maintenance or add important extensions late
 - Developers who first commit 730 days after the first commit

As we look at dynamics, we focus on projects we see the first commit, ie after 2013.

- Collaboration – link developers who contribute to the same repo.
- Dependencies – link developers from one package using another
- One observation is one link
- Aggregated at city (city pair) level

- Start with the developer's link to a repository (via commits)
- Directed but (mostly fully) symmetric
- Transform it to developer to developer links
- Aggregate at city level

Links in the contribution network

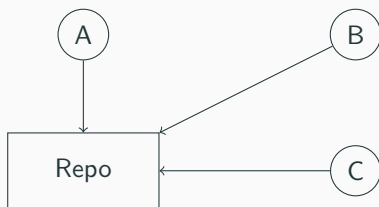


Figure 1: Developers committing to a repository.

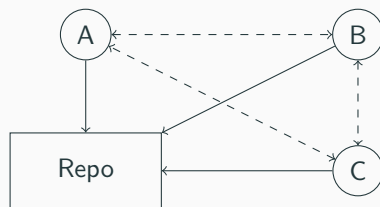


Figure 2: Developers committing to a repository including implied contributor to contributor links.

Solid lines are what we **observe**. Dashed lines is what we **infer**.

- In a repo, all developers create links with each other
- If two people have 3 repo together, will generate 3 links
- Also look at *intensive* margin – weighted by commits

- Github collaboration system
- Mostly amateurs (like CEU Econ)
- Includes corporations (like Oracle)
- Today: mostly focus outside organizations

Estimating gravity

Gravity: finding a partner

- The role of distance in finding a partner
- Search and maintenance
- Each coder can choose any partner: logit
- Aggregate + transform: Poisson at city pair level: number of links as function of distance
- (*Yes, like structural gravity: PPML, FEs*)

MORE: ▶ From logit to Poisson

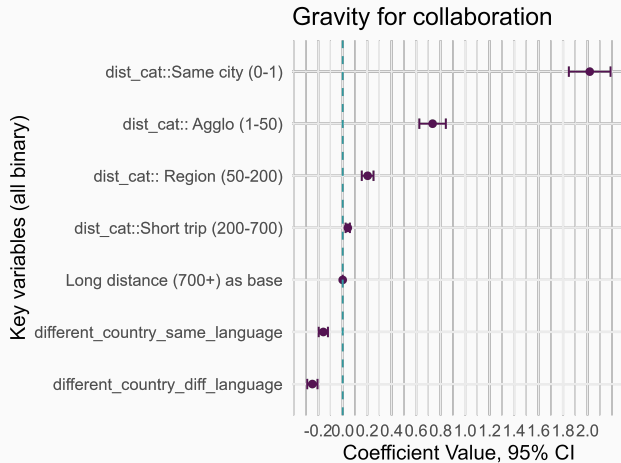
$$\Pr(Y_{od}|x_o, x_d, d_{od}) \approx \text{Poisson}[N_o \times N_d \times \exp(\beta_1 x_o + \beta_2 x_d + \beta_3 d_{od})]$$

- Outcome: Number of links between cities o, d
- d_{od} Distance measured as a set of indicators / log-linear
- Origin and destination city FE
- $N_o \times N_d$ -Exposure: Number of developers in city $o \times d$

Modelling search and maintenance costs

- Meeting – distance in terms of travel
 - Same city – e.g. universities, office parks
 - Agglomeration (1-50km) – regional events
 - Regional (50-200km) – national conferences
 - Short trip (200-700km) – big conferences
 - Beyond 700km (*as base*) – global events
- Travel difficulty
 - Crossing borders
 - Crossing borders — different language

Results 1: More work together when closer



Gravity 1: N of links between cities declines with distance

Dependent Variable: Model:	N of links between contributors		
	(1)	(2)	(3)
ln_dist	-0.1650*** (0.0083)		
same_org	5.609*** (0.0871)	5.556*** (0.0855)	
dist_cat = Samecity(0-1)		1.746*** (0.0772)	2.018*** (0.0858)
dist_cat = Agglo(1-50)		0.6351*** (0.0724)	0.7344*** (0.0873)
dist_cat = Region(50-200)		0.1905*** (0.0319)	0.2039*** (0.0307)
dist_cat = Shorttrip(200-700)		0.0245* (0.0127)	0.0416*** (0.0101)
different_country_same_language		-0.1749*** (0.0215)	-0.1581*** (0.0184)
different_country_diff_language		-0.2856*** (0.0369)	-0.2476*** (0.0322)
Pseudo R ²	0.84163	0.84385	0.86084
Observations	3,634,243	3,636,122	3,478,716

Origin, destination city FE, Clustered (city_destination & city_origin) standard-errors in parentheses

Results 2: Commits as kinda intensive margin

- Look at commits – number of changes in code
- Bit like extensive margin

Gravity 2: Co-location = more intensive work

Dependent Variables: Model:	N links (1)	commit share (2)
<i>Variables</i>		
dist_cat = Samecity(0-1)	2.018*** (0.0858)	0.7564*** (0.1309)
dist_cat = Agglo(1-50)	0.7344*** (0.0873)	0.1838 (0.1410)
dist_cat = Region(50-200)	0.2039*** (0.0307)	0.0906 (0.0795)
dist_cat = Shorttrip(200-700)	0.0416*** (0.0101)	-0.0192 (0.0399)
<i>Fixed-effects</i>		
city_destination	Yes	Yes
city_origin	Yes	Yes
<i>Fit statistics</i>		
Pseudo R ²	0.86084	0.52444
Observations	3,478,716	451,423

Origin, destination city FE, Clustered (city_destination & city_origin) standard-errors in parentheses

- Maybe a few very large repositories dominate and flatten the curve. No
- Also no huge difference excluding few largest cities

Estimating success and dispersion

Success (popularity) and spatial dispersion

- Popularity = measures the number of other packages which declare a dependency on a the repository in NPM
- Measures on spatial dispersion
- Controls

Success (popularity) and spatial dispersion

$$\Pr(Y_i|.) \approx \text{Poisson}[\exp(\beta_1 \text{cities}_i + \beta_2 \text{countries}_i) + \gamma \mathbf{Z}]$$

- Outcome: Number of repos importing this repo i
- countries_i number of countries
- cities_i number of cities
- **Z: f(number of developers), f(age of project) + other controls**

Mechanic controls

- Number of developers (as categorical)
- Age of project, ys (as categorical)

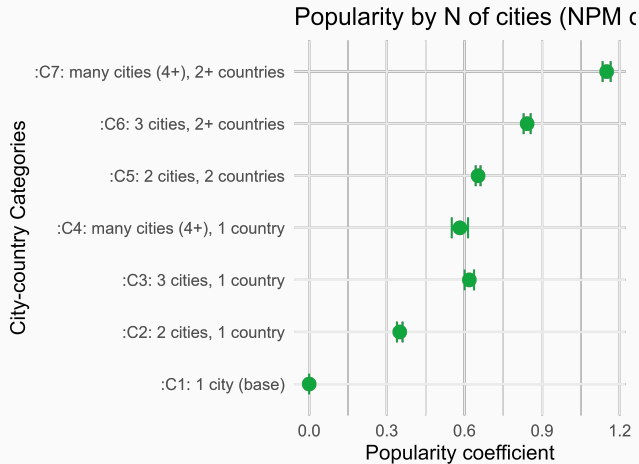
Coder quality

- Size of the city of coder(s) – top 3 coders
- Quality of coder(s) as measured by number of stars– top 3 coders

Commits

- sum commits (squared)

Results 1: More popular dependency - higher spatial dispersion



Results 1: More popular dependency - higher spatial dispersion

Dependent Variable: Model:	N Dependents (NPM)		
	(1)	(2)	(3)
Constant	1.745*** (0.0548)	1.148*** (0.0857)	1.673*** (0.0674)
Count of cities	0.4075*** (0.0403)	0.2306*** (0.0491)	
Count of countries	0.3431*** (0.0628)	0.3057*** (0.0637)	
City cat × CI2 × 2cities			0.3851*** (0.0813)
City cat × CI3 × 3cities			0.4925*** (0.1242)
City cat × CI4 × manycities(4+)			0.6543*** (0.1642)
Country cat × CO2 × 2countries			0.2461*** (0.0813)
Country cat × CO3 × manycountries(3+)			0.6269*** (0.1462)
Age, N.Dev	No	Yes	Yes
Commits	No	No	No
Coders	No	No	No
Pseudo R ²	0.05100	0.11532	0.11586
Observations	36,491	36,491	36,491

What's behind this correlation?

- Packages written by *more* dispersed people will be used more.
 - Is it because a diverse group reaches diverse regions generating use? (reverse causality-ish)
 - Let us look at other packages using ours as dependency

What's behind this correlation?

- Packages written by *more* dispersed people will be used more.
 - Is it because a diverse group reaches diverse regions generating use? (reverse causality-ish)
 - Let us look at other packages using ours as dependency
 - Is this generated by a sorting model?
 - Think about source of sorting
 - Condition on selection

Preparation: Aggregating dependencies to city level

- We observe a repository importing another one as dependency.
- Directed, not symmetric
- Transform it to developer to developer links
 - Use knowledge of producers of the dependency as well
- Aggregate at city level

Links in the dependency network

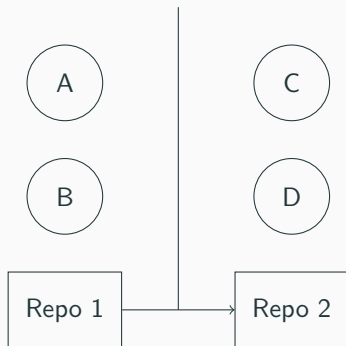


Figure 3: Dependency of repository 1 on repository 2 with the respective developers.

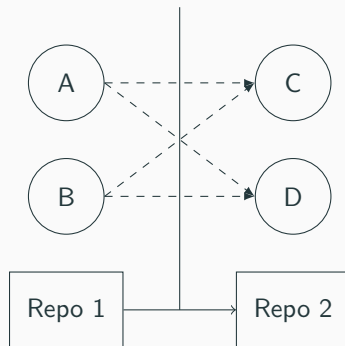


Figure 4: Dependency of repository 1 on repository 2 with the respective developers. Dashed lines indicate implied links between developers.

Again, solid lines are what we **observe**. Dashed lines is what we **infer**.

1. Not reverse causality - dependency use just mildly spatial

Dependent Variables: Model:	contr_n.links (1)	dep_value (2)
<i>Variables</i>		
dist_cat = Samecity(0-1)	2.018*** (0.0858)	0.0754*** (0.0138)
dist_cat = Agglo(1-50)	0.7344*** (0.0873)	0.0805*** (0.0127)
dist_cat = Region(50-200)	0.2039*** (0.0307)	0.0254*** (0.0095)
dist_cat = Shorttrip(200-700)	0.0416*** (0.0101)	0.0045 (0.0036)
different_country_same_language	-0.1581*** (0.0184)	-0.0222*** (0.0082)
different_country_diff_language	-0.2476*** (0.0322)	-0.0499*** (0.0115)
Pseudo R ²	0.86084	0.98866
Observations	3,478,716	3,202,202

Origin, destination city FE, Clustered (city-destination & city-origin) standard-errors in parentheses

2. Models of sorting

- Coders vary by skill, match within and across cities randomly, large cities have better coders. Opposite correlation
- Coders vary by many types of skill. FC of matching across cities. Better coders self-select into diverse teams producing better code. Selection on coder skills.

MORE: [▶ More on a sketch of a theory](#)

Selection?

- Coder quality: Number of repos, number of followers
- Location: Size of the city (N of total coders)
- Commits: number of commits to repo

Results 2: Selection? Partialing out coder quality and commits

Dependent Variables: Model:	N Dependents (NPM)			N commits
	(1)	(2)	(3)	(4)
Count of cities	0.2306*** (0.0491)	0.1906*** (0.0488)	0.2818*** (0.0509)	-0.1087*** (0.0358)
Count of countries	0.3057*** (0.0637)	0.3243*** (0.0634)	0.2904*** (0.0637)	0.0355 (0.0583)
Sum of commits (ln)			0.3330*** (0.0202)	
Constant	1.148*** (0.0857)	-0.0641 (0.1140)	-1.746*** (0.1507)	5.090*** (0.0854)
Age, N_Dev	Yes	Yes	Yes	Yes
Coders	No	Yes	Yes	Yes
Commits	No	No	Yes	—
Pseudo R ²	0.11532	0.15056	0.18345	0.26074
Observations	36,491	36,491	36,491	36,491

- Compare coders of similar quality based in similar locations
- Exclude success driven by bigger spatial reach of developers
- Group of diverse coders will create more successful projects
- Even conditioning on intensity of collaboration (commits)

- Organization
 - Results robust, but something is going on in orgs.
- Missing city info
 - Half of city information is missing – could be non-random
 - For many we know the organisation they work for – l8r
 - Results robust to filter on 2-3 member teams with known location

Discussion

Summary

- Location matters even for coding

Summary

- Location matters even for coding
- Will the best coders congregate in big cities to create best code?

- Location matters even for coding
- Will the best coders congregate in big cities to create best code?
- No. Spatially dispersed developers create code that is more widely adopted.
- Sorting matters: good coders write good code used by more. But not explains
- It's also no more effort – when looking at commits, we see the opposite...

- Location matters even for coding
- Will the best coders congregate in big cities to create best code?
- No. Spatially dispersed developers create code that is more widely adopted.
- Sorting matters: good coders write good code used by more. But not explains
- It's also no more effort – when looking at commits, we see the opposite...
- There is something else. Maybe higher fixed costs of starting a diverse project generates a selection on project ideas.

References

- AlShebli, Bedoor K., Talal Rahwan, and Wei Lee Woon**, “The preeminence of ethnic diversity in scientific collaboration,” *Nature communications*, 2018, 9, 1–10.
- Anderson, James E, Ingo Borchert, Aaditya Mattoo, and Yoto V Yotov**, “Dark costs, missing data: Shedding some light on services trade,” *European Economic Review*, 2018, 105, 193–214.
- Bahar, Dany**, “The hardships of long distance relationships: time zone proximity and the location of MNC’s knowledge-intensive activities,” *Journal of International Economics*, 2020, 125, 103311.

- Bailey, Mike, Abhinav Gupta, Sebastian Hillenbrand, Theresa Kuchler, Robert Richmond, and Johannes Stroebe**, “International Trade and Social Connectedness,” *Journal of International Economics*, Mar 2021, 129, 103418.
- Bernard, Andrew B, Andreas Moxnes, and Yukiko U Saito**, “Production networks, geography, and firm performance,” *Journal of Political Economy*, 2019, 127 (2), 639–688.
- Bircan, Cagatay, Beata Javorcik, and Stefan Pauly**, “Creation and Diffusion of Knowledge in the Multinational Firm,” 2021. Working Paper.
- Blum, Bernardo S and Avi Goldfarb**, “Does the internet defy the law of gravity?,” *Journal of international economics*, 2006, 70 (2), 384–405.
- Chaney, Thomas**, “The network structure of international trade,” *The American Economic Review*, 2014, 104 (11), 3600–3634.
- Davis, Donald R and Jonathan I Dingel**, “A spatial knowledge economy,” *American Economic Review*, 2019, 109 (1), 153–170.

- Fackler, Thomas, Michael Hofmann, and Nadzeya Laurentsyevea**, “Defying Gravity: What Drives Productivity in Remote Teams?,” Technical Report, LMU CRCT Discussuion Paper 427 2023.
- Head, Keith, Yao Amber Li, and Asier Minondo**, “Geography, Ties, and Knowledge Flows: Evidence from Citations in Mathematics,” *The Review of Economics and Statistics*, 10 2019, 101 (4), 713–727.
- Jaffe, Adam B, Manuel Trajtenberg, and Rebecca Henderson**, “Geographic localization of knowledge spillovers as evidenced by patent citations,” *Quarterly journal of Economics*, 1993, 108 (3), 577–598.
- Laurentsyevea, Nadzeya**, “From friends to foes: National identity and collaboration in diverse teams,” Technical Report, CESifo Discussion Paper 2019.
- Lerner, Josh and Jean Tirole**, “Some simple economics of open source,” *The journal of industrial economics*, 2002, 50 (2), 197–234.

Stein, Ernesto and Christian Daude, “Longitude matters: Time zones and the location of foreign direct investment,” *Journal of International Economics*, 2007, 71 (1), 96–112.

Wachs, Johannes, Mariusz Nitecki, William Schueller, and Axel Polleres, “The Geography of Open Source Software: Evidence from GitHub,” *Technological Forecasting and Social Change*, 2022, 176, 121478.

Behind Poisson 1: Individual matching decision

Collaboration or dependency link between developer i and j ,

$$\Pr(Y_{ij} = 1 | x_i, x_j, d_{ij}) = \Pi(\beta_1 x_i + \beta_2 x_j + \beta_3 d_{ij})$$

with

$$\Pi(z) = e^z / (1 + e^z)$$

the logistic function

Assumption: Independence across links, add fixed effects

Behind Poisson 2: Aggregate to Poisson

In practice, distance only varies at the city level. Take origin city o and destination city d .

$$Y_{od} := \sum_{i \in o} \sum_{j \in d} Y_{ij}$$

$$\Pr(Y_{od} | x_o, x_d, d_{od}) = \text{Binomial}[N_o \times N_d, \Pi(\beta_1 x_i + \beta_2 x_j + \beta_3 d_{ij})]$$

Here $N_o \times N_d$ is the total number of *potential* links between cities o and d .

When Π is small, we aggregate i into cities o , and j into cities d

$$\Pr(Y_{od} | x_o, x_d, d_{od}) \approx \text{Poisson}[N_o \times N_d \times \exp(\beta_1 x_o + \beta_2 x_d + \beta_3 d_{od})]$$

Behind Poisson 3: Having exposure is key

We may also look at a subsample (like users not in the same GitHub organization)

$$Y_{od, \text{not org}} := \sum_{i \in o} \sum_{j \in d, j \notin \text{org}(i)} Y_{ij}$$

This changes the *exposure variable*,

$$\Pr(Y_{od, \text{not org}} | x_o, x_d, d_{od}) \approx \text{Poisson}[N_{od, \text{not org}} \times \exp(\beta_1 x_o + \beta_2 x_d + \beta_3 d_{od})],$$

with $N_{od, \text{not org}}$ the number of user pairs in city o, d , *not sharing* an organization.

Important: $N_{od, \text{not org}}$ may be zero.

What is a Poisson regression?

First-order conditions for Maximum Likelihood:

$$\sum_o \sum_d x_i [Y_{od} - N_{od} \exp(\beta_1 x_o + \beta_2 x_d + \beta_3 d_{od})] = 0$$

$$\sum_o \sum_d x_j [Y_{od} - N_{od} \exp(\beta_1 x_o + \beta_2 x_d + \beta_3 d_{od})] = 0$$

$$\sum_o \sum_d d_{ij} [Y_{od} - N_{od} \exp(\beta_1 x_o + \beta_2 x_d + \beta_3 d_{od})] = 0$$

- Level (not log) error terms are orthogonal to RHS variables.
- Exposure variable has fixed exponent of 1 (\approx weighting).
- Standard errors computed from GMM, not ML. E.g., we allow for two-way city clustering.

What is an observation?

Two interpretations:

1. 10 billion potential developer pairs
2. 3.7 million city pairs

Model sketch

- Production of code is driven by utility gains of creating code used by many people
- Coders are heterogeneous in coding quality.
- Coders collaborate with others when
 - Task is too complex for a single person. Economies of scale.
 - ...
- There is selection into projects: best coders write most complex packages.

Model sketch 2: The role of geography

- Coders are dispersed geographically – located in a discrete set of N_c cities
 - City size (number of coders) Pareto distributed
 - Size may be driven by first geography (later), such as proximity to University, tech firms or the beach.
- Heterogeneity of coders: at every location, their distribution is Pareto
- Random matching: simple random selection of collaborators
- Assortative matching: Coders match with coders of same quality

Model: self selection of coders

- If best programmers are in big cities (Pareto with different k across cities): size and quality correlated
- Top coders coming from large cities will produce best code – \hat{c} more popular code.
- Best code will come more than proportionally from large cities
- Assortative matching reinforces this aspect, as big city coders will only work with big city coders
- Best code written by people in top cities (like SF) – homogeneity

Model: There are search costs

- Costs of setting up a partnership and maintaining it
- Search costs of inputs (code chunks)
 - Written together – finding a collaborator
 - Using already published code – finding a package
- Search costs vary with distance – lower inside the city

Model: Coder heterogeneity

- There is a set of possible coding skills, S
- Coders randomly vary in each skill, $s = 1, 2, 3 \dots S$
- Two coders who are on average same quality still have difference and can benefit from collaboration, where the pair's skill is max

Model: Dispersion forces

- Coders differ to some extent, and so search is needed
- There is a search cost, higher for other cities
- Better coders pay higher search cost and hence can search a larger pool across cities

Model: Additional aspects

- Face to face matters when creating complex projects.
- Some cities specialize in some tasks