



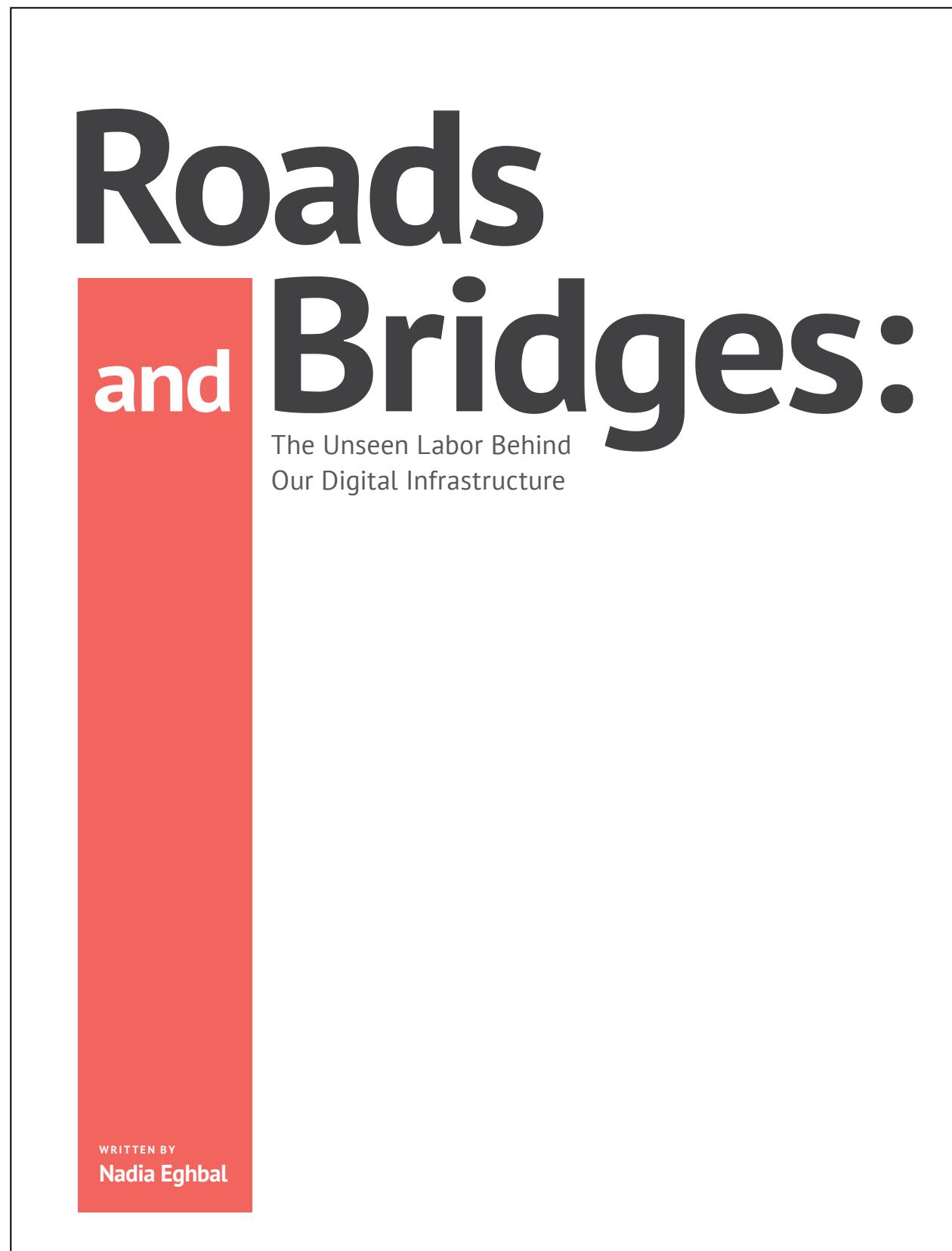
DALL·E 3 - "Networks of open-source software projects"

# The Strength of Weak Ties in Open-Source Software Development Networks

Bogdan Vasilescu  
At UC Irvine, October 3rd, 2024

# Open source software has become digital infrastructure

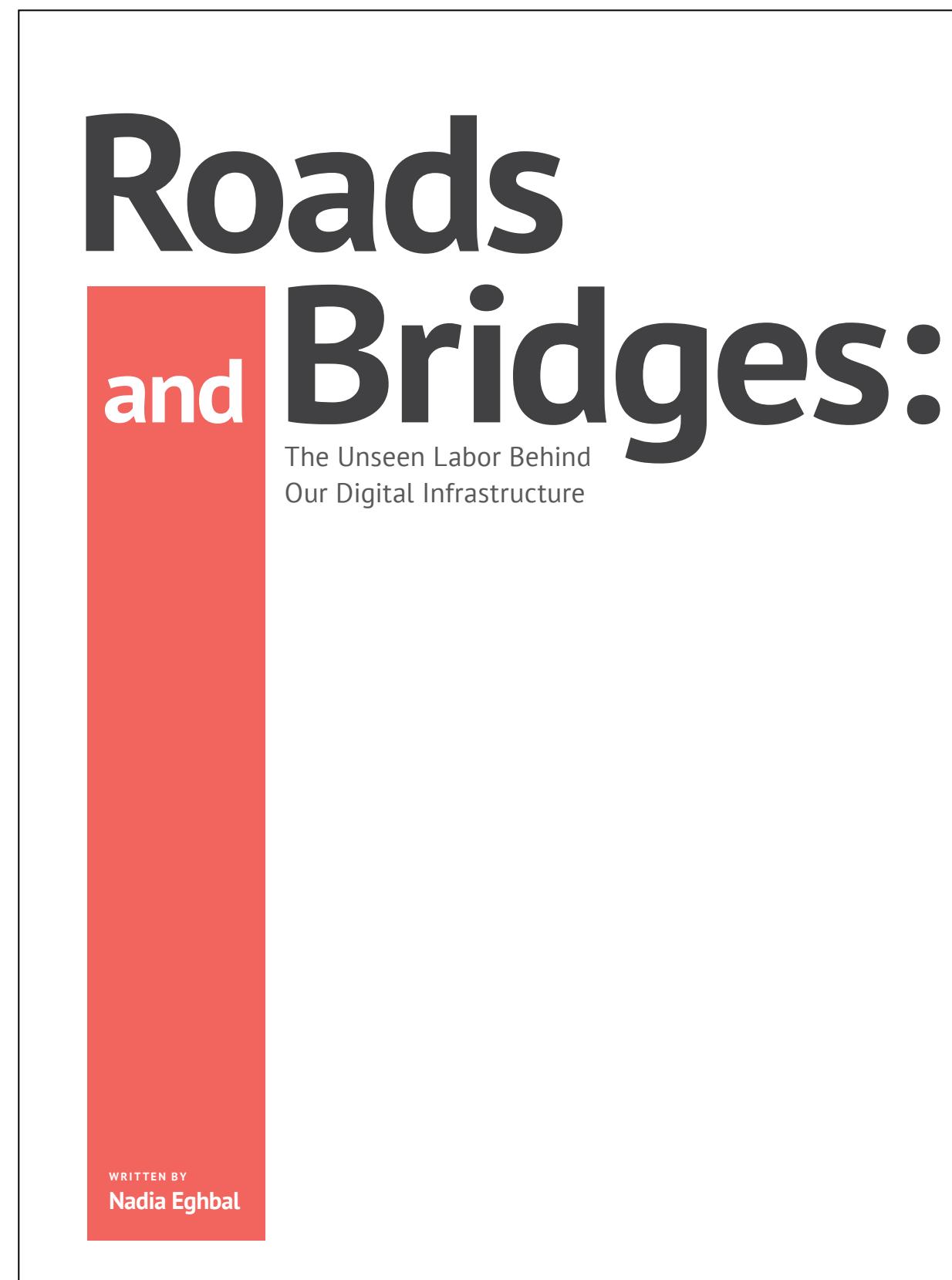
---



Everybody uses open source:

- Fortune 500 companies
- Major software companies
- Startups
- Government
- ...

# Like any infrastructure, it needs regular upkeep and maintenance



Everybody uses open source:

- Fortune 500 companies
- Major software companies
- Startups
- Government
- ...

If undermaintained:

- Brittle supply chains
- Risks for downstream users
- Slows down innovation
- ...

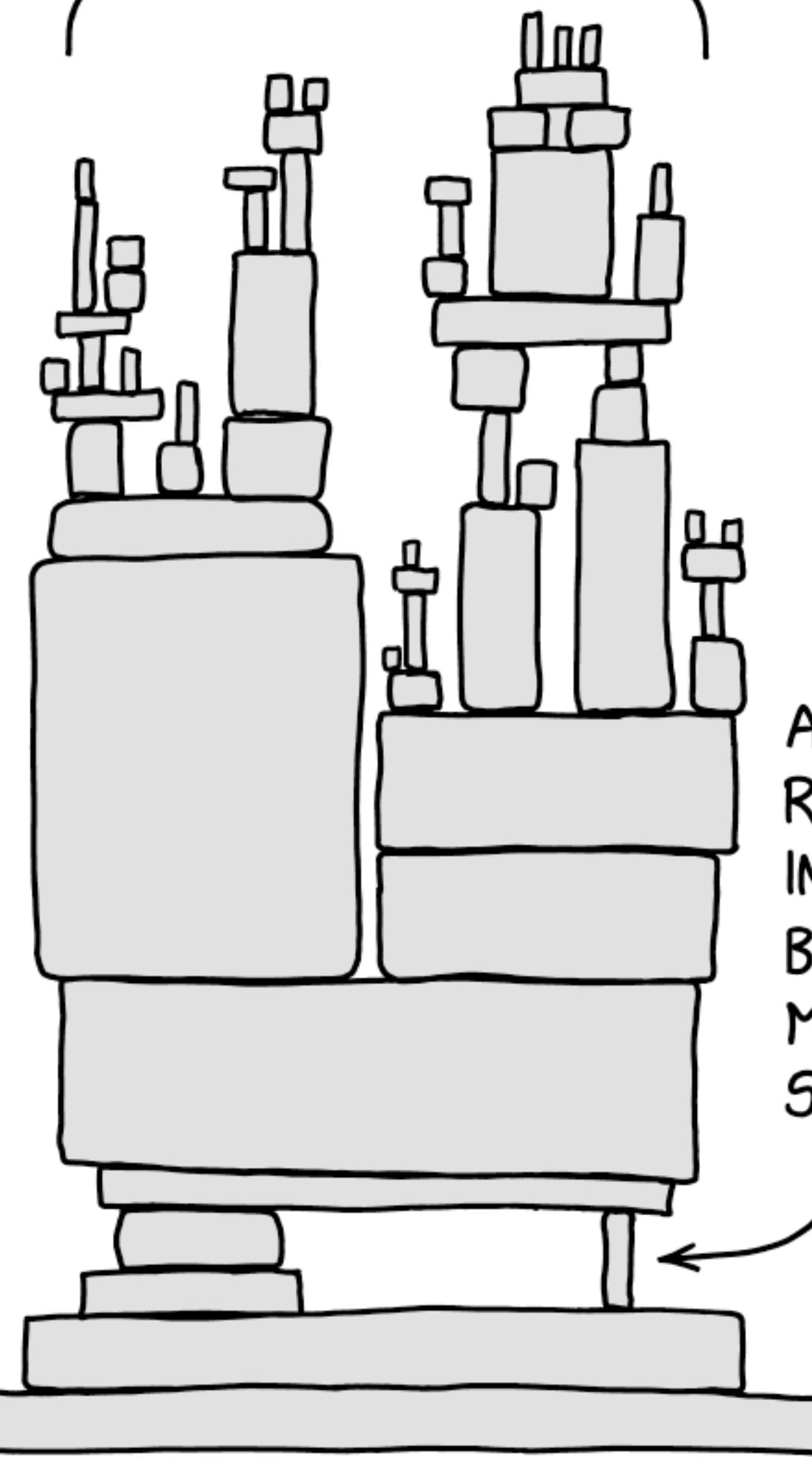
A screenshot of a news article from Quartz. The title is 'How one programmer broke the internet by deleting a tiny piece of code'. Below the title, it says 'By Kaith Collins • March 27, 2016'. The main content shows a code editor with a file named 'leftpad.js' containing the following code:

```
1 module.exports = leftpad;
2 function leftpad(str, len, ch) {
3     str = String(str);
4     var i = -1;
5     if (!ch && ch !== 0) ch = ' ';
6     len = len - str.length;
7     while (++i < len) {
8         str = ch + str;
9     }
10    return str;
11 }
```

<https://qz.com/646467/how-one-programmer-broke-the-internet-by-deleting-a-tiny-piece-of-code/>



ALL MODERN DIGITAL  
INFRASTRUCTURE



Sustaining  
open source  
is hard

# Ever more open source software is being created (and reused)

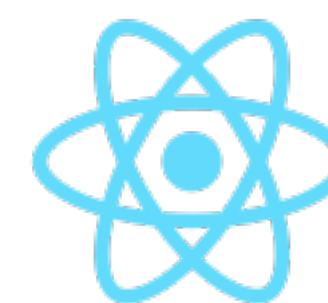
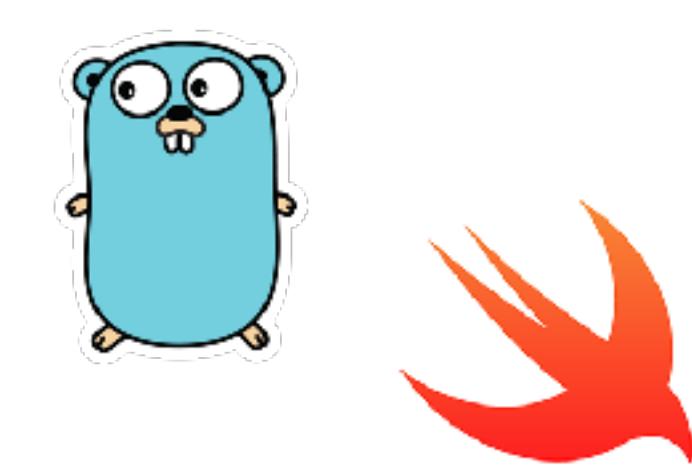
Explosion of production in the past 10 years



# There is increasing commercialization and professionalization

---

- Historically
  - Community-based projects (Python, RubyGems, Twisted)
- More recently, lots of commercial involvement
  - Companies (Go - Google, React - Facebook, Swift - Apple)
  - Startups (Docker, npm, Meteor)



- 23% of respondents to 2017 GitHub survey: job duties include contributing to open source

<http://opensourcesurvey.org/2017/>

# Expectations toward the quality, reliability, and security of open source infrastructure are high

Equifax (market cap \$14 billion) built products on top of open-source infrastructure, including Apache Struts

Equifax did not make any contributions to open source projects

A flaw in Apache Struts contributed to the breach (CVE-2017-5638)

Equifax publicly blamed (with national news coverage) Apache Struts for the breach

## Equifax confirms Apache Struts security flaw it failed to patch is to blame for hack

The company said the March vulnerability was exploited by hackers.

By Zack Whittaker | September 14, 2017 -- 01:27 GMT (18:27 PDT) | Topic: Security



<https://www.zdnet.com/article/equifax-confirms-apache-struts-flaw-it-failed-to-patch-was-to-blame-for-data-breach/>

# High level of demands & stress

Easy to report issues / submit PRs

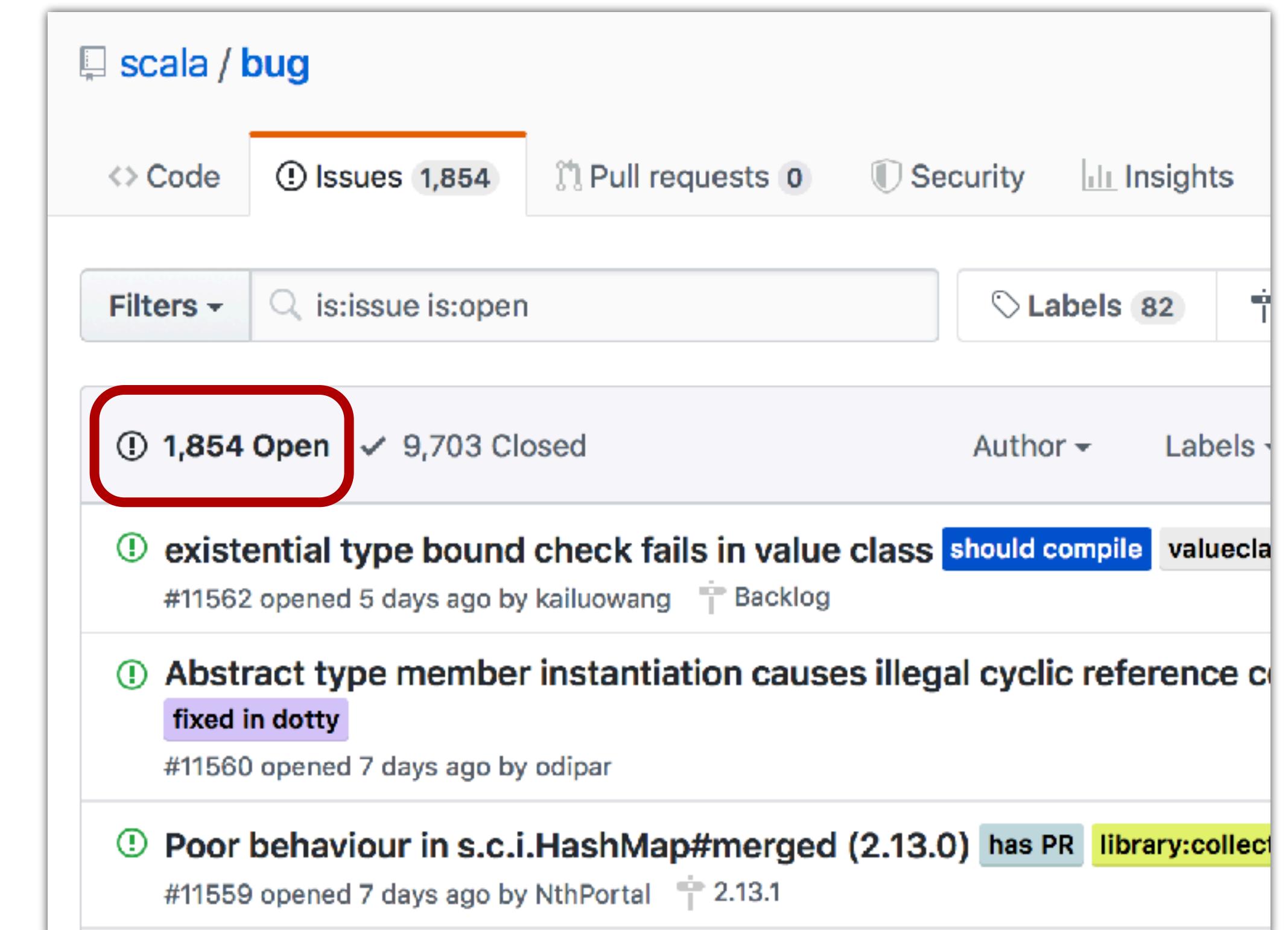
- Growing volume of requests

Social pressure to respond quickly

- Otherwise, off-putting to newcomers  
(Steinmacher et al. 2015)

Entitled, unreasonable users:

- *“I have been waiting 2 years for Angular to track the ‘progress’ event and it still can’t get it right?!?!”*
- *“Thank you for your ever useless explanations.”*



# The social platforms have won

Profile pages for users and projects

Rich inferences about people's expertise and level of commitment

Impacts collaboration, but also recruiting and hiring

- (Dabbish et al. 2012), (Marlow et al. 2013), (Marlow and Dabbish 2013)

The image shows two GitHub profile pages side-by-side.

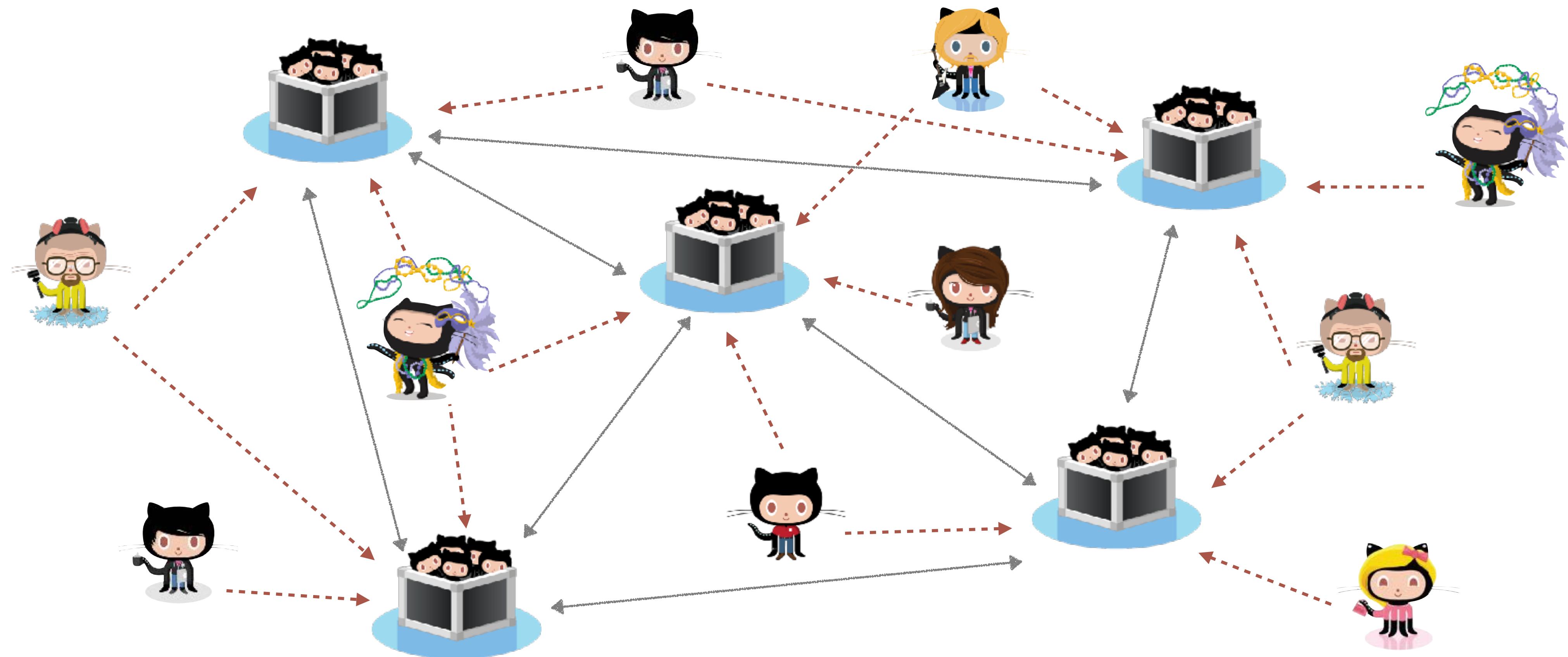
**Top Profile (npm, inc):**

- User Avatar:** A GitHub Octocat holding a laptop labeled "CV".
- Organizations:** Shows membership in four organizations: a hat icon, a blue cloud icon, a red star icon, and a green gear icon.
- Statistics:** Joined on Oct 31, 2011. 776 Followers, 38 Starred, 15 Following.
- Public contributions:** A heatmap showing activity from Feb to Jan. Summary: Contributions in the last year: 1,886 total (Jan 24, 2015 – Jan 24, 2016). Longest streak: 37 days (October 7 – November 12). Current streak: 7 days (January 18 – January 24).
- Repositories contributed to:** A list of repositories: npm/docs (44 stars), mozilla/publish.webmaker.org (2 stars), npm/marky-markdown (104 stars), artisan-tattoo/assistant-frontend (5 stars), and npm/npm-camp (1 star).

**Bottom Profile (caolan / async):**

- User Avatar:** A blue arrow pointing right next to the word "async".
- Statistics:** 1,629 commits, 11 branches, 72 releases, 206 contributors, MIT license.
- README.md:** Displays the content of the README file, featuring the "async" logo.
- Build Status:** build passing.
- Deployment:** npm v2.6.0, coverage 95%, gitter, join chat, examples 26348, jsDelivr 40.7k hits/month.
- Description:** Async is a utility module which provides straight-forward, powerful functions for working with asynchronous JavaScript. Although originally designed for use with Node.js and installable via npm install --save async, it can also be used directly in the browser.

# Contributors and projects form complex socio-technical networks!



# Science is needed for evidence-based recommendations

Rich socio-technical ecosystem  
Lots of changes  
Lots of challenges

Little evidence or theory

## Anecdotal evidence reliable? One man says “yes”.

A STUDY CONDUCTED YESTERDAY by a man on himself concluded that self-reported anecdotal evidence is, in fact, both reliable and relevant.

The landmark study, conducted by Mark Mattingly of Virginia Beach in his apartment, concluded with 100% accuracy that data collected from personal experience can disprove other data conducted by reputable scientific institutions, thereby proving once and for all that “statistics can’t be trusted”.

In a press release Mr. Mattingly took aim at his detractors saying that “...this study shows what I’ve been telling people on the internet for years: all your fancy evidence and statistics don’t mean nothing in the real world.”

A frequenter of internet forums, comment sections, and social media, Mr. Mattingly recounts that he was inspired to undertake the study when someone reportedly kept insisting that he provide evidence for his claims. “I think everyone’s entitled to an opinion, and that my opinion is worth just as much as anyone else’s” Mr. Mattingly said.

Academic types have criticised the study, and papers who are publishing it, saying that it lacks everything and makes no sense. When shown the study, Emeritus Professor James Albrecht of Carnegie Mellon University looked all confused and hopeless before making pining, guttural sounds.



Mr. Mattingly in his apartment looking all smug.

Mr. Mattingly has responded saying that this is just the first of many studies he intends to conduct, and that a meta-analysis of people who have opinions and anecdotal experiences independent of controls, methodological rigor, blinding and peer review are soon to be published, adding further weight to his initial findings.

# A great opportunity for research!

# A great opportunity for research!

... because (almost) everything being  
archived and public makes it possible  
to study the problem empirically

# A great opportunity for research!

... because (almost) everything being archived and public makes it possible to study the problem empirically



“The collection of public Git repositories as a whole [...] exceeds 1.5PB” (Ma et al, 2021)

# Today: Let's look at some concrete examples of network effects

---



Measuring innovation  
in software



Understanding how  
innovation emerges



Social capital



Social contagion

STRUDEL

# Today: Let's look at some concrete examples of network effects

---



Measuring innovation  
in software



Understanding how  
innovation emerges



Social capital



Social contagion

# Open-source software development is an avenue for innovation and creative expression.

---

(Lakhani & Wolf, 2005)

“How creative a person feels when working on the project is the strongest and most pervasive driver [of participation in open source]”

“Free software is directly responsible for today’s current startup renaissance.”

(Eghbal, 2016)

How to define innovation in software?

How to measure it?

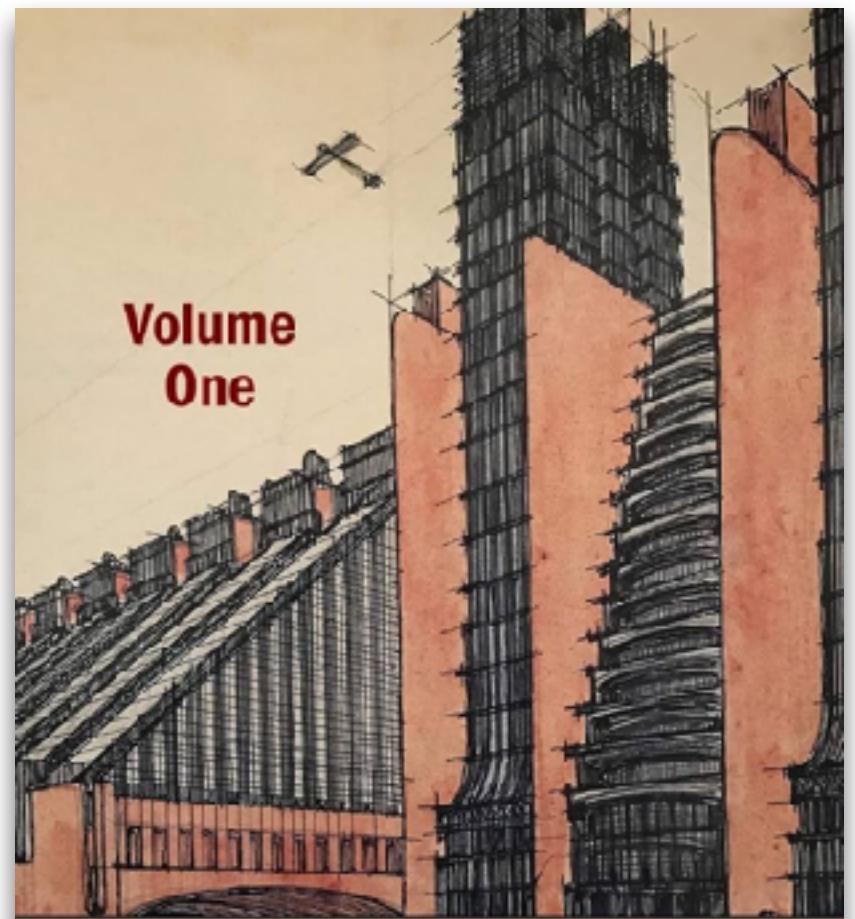
How does innovation emerge?

What are its consequences?



DALL-E 3 - "An old-looking map with uncharted territory, here be dragons style"

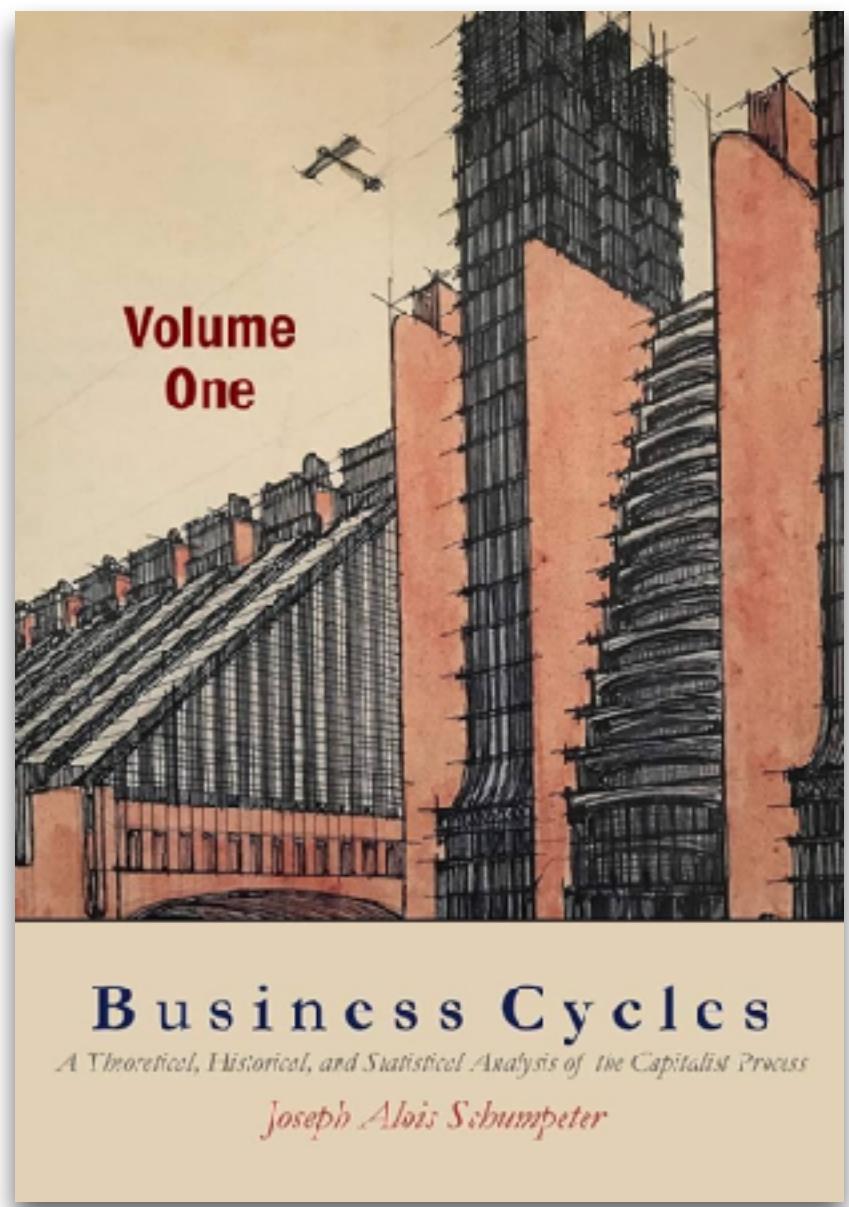
# Key idea: Innovation as novel recombination



“[We may say] that innovation combines factors in a new way, or that it consists in carrying out new combinations.”

(Schumpeter, 1939)

# Key idea: Innovation as novel recombination



(Schumpeter, 1939)

“[We may say] that innovation combines factors in a new way, or that it consists in carrying out new combinations.”

“... how scientists search for ideas is premised in part on the idea that teams can span scientific specialties, effectively combining knowledge that prompts scientific breakthroughs.”

## Atypical Combinations and Scientific Impact

Brian Uzzi,<sup>1,2</sup> Satyam Mukherjee,<sup>1,2</sup> Michael Stringer,<sup>2,3</sup> Ben Jones<sup>1,\*</sup>

Novelty is an essential feature of creative ideas, yet the building blocks of new ideas are often embodied in existing knowledge. From this perspective, balancing atypical knowledge with conventional knowledge may be critical to the link between innovativeness and impact. Our analysis of 17.9 million papers spanning all scientific fields suggests that science follows a nearly universal pattern: The highest-impact science is primarily grounded in exceptionally conventional combinations of prior work yet simultaneously features an intrusion of unusual combinations. Papers of this type were twice as likely to be highly cited works. Novel combinations of prior work are rare, yet teams are 37.7% more likely than solo authors to insert novel combinations into familiar knowledge domains.

**S**cientific enterprises are increasingly concerned that research within narrow boundaries is unlikely to be the source of the most fruitful ideas (1). Models of creativity emphasize that innovation is spurred through original combinations that spark new insights (2–10). Current interest in team science and how scientists search for ideas is premised in part on the idea that teams can span scientific specialties, effectively combining knowledge that prompts scientific breakthroughs (11–15).

<sup>1</sup>Kelogg School of Management, Northwestern University, 2001 Sheridan Road, Evanston, IL 60228, USA. <sup>2</sup>Northwestern Institute on Complex Systems, Northwestern University, 600 Foster, Evanston, IL 60208, USA. <sup>3</sup>Datascope Analytics, 180 West Adams Street, Chicago, IL 60603, USA. <sup>4</sup>National Bureau of Economic Research, 1050 Massachusetts Avenue, Cambridge, MA 02138, USA.

\*Corresponding author. E-mail: bjones@kellogg.northwestern.edu

468

25 OCTOBER 2013 VOL 342 SCIENCE www.sciencemag.org

(Uzzi et al, 2013)

between exten-  
binations of k  
advantages of  
ing is critical to  
and impact. Ho  
composition o  
can achieve it

In this stud

search articles  
see how prior v  
that indicate (i)

pers referenc

ations of prio

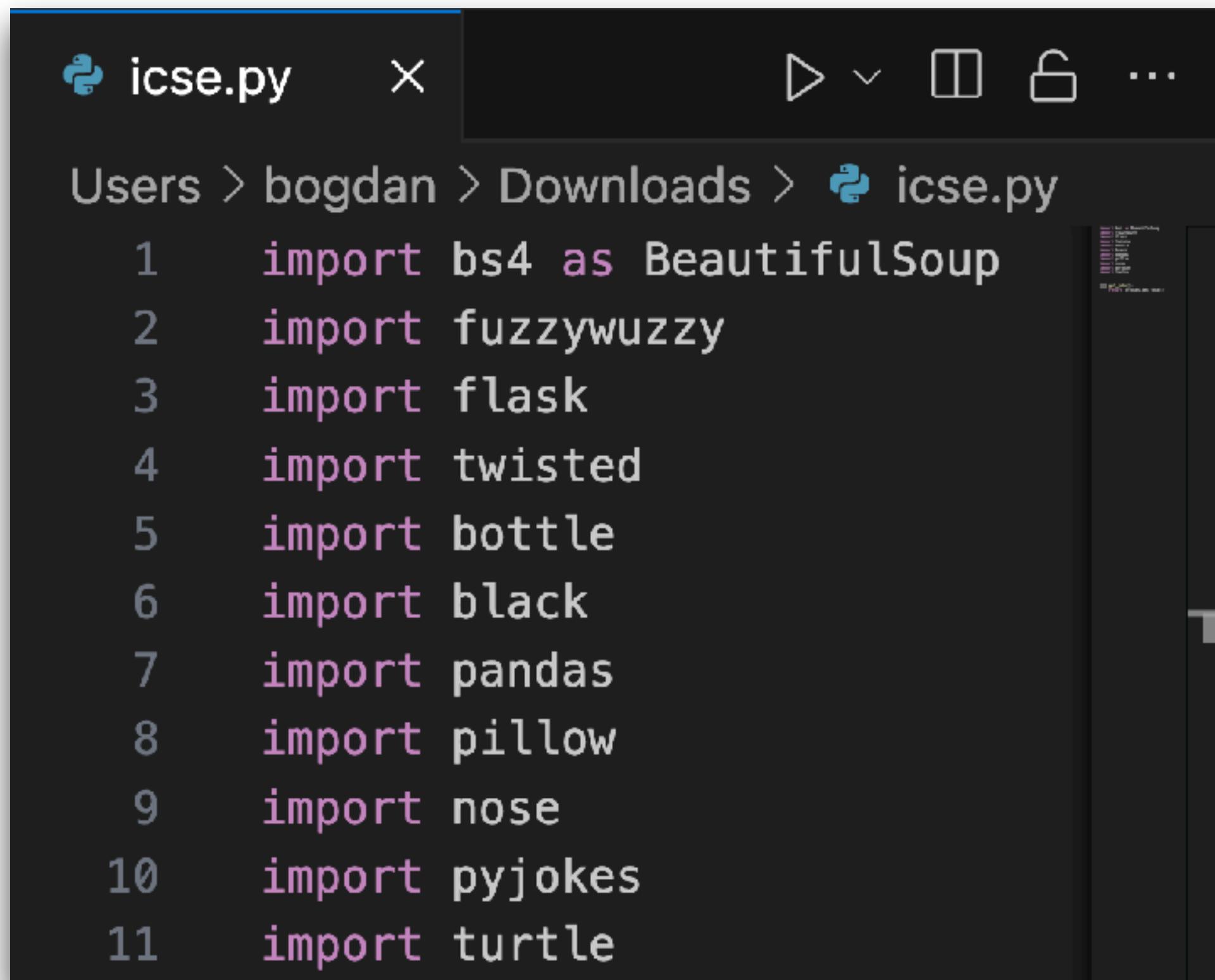
papers based

upon, and (ii)

collaboration.

We consider  
ences in the bil  
We counted th  
pair across all  
WOS and com  
to those expec  
citation netwo  
networks, all i  
in the WOS w  
Carlo algorith  
serves the total  
paper and the  
counts forward  
that a paper (c  
observed new  
randomized ne  
the randomized  
we aggregated  
respective joun  
combinations (t  
over 122 milli  
by the 15,613

# Software innovation as novel recombination of software libraries



A screenshot of a terminal window titled 'icse.py'. The window shows the file path: 'Users > bogdan > Downloads > icse.py'. The code content is as follows:

```
1 import bs4 as BeautifulSoup
2 import fuzzywuzzy
3 import flask
4 import twisted
5 import bottle
6 import black
7 import pandas
8 import pillow
9 import nose
10 import pyjokes
11 import turtle
```

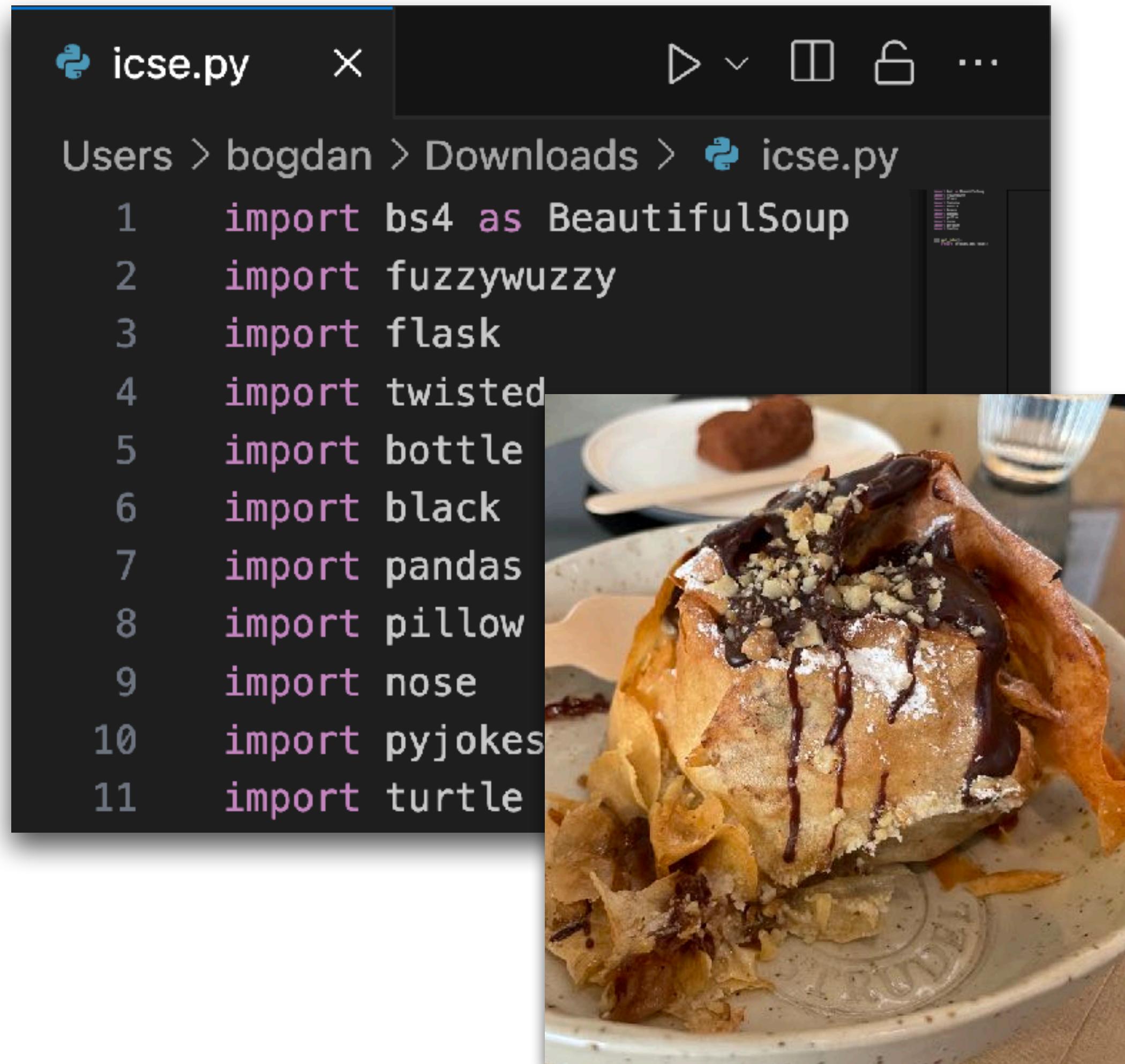
Lots of combinations:

- **(twisted, bottle)**
- **(turtle, nose)**
- **(black, pandas)**
- **(fuzzywuzzy, pillow)**
- ...

$C(n,2)$  unique pairs of packages.

Some of these may be highly innovative because they are atypical.

# Software innovation as novel recombination of software libraries



Lots of combinations:

- **(twisted, bottle)**
- **(turtle, nose)**
- **(black, pandas)**
- **(fuzzywuzzy, pillow)**
- ...

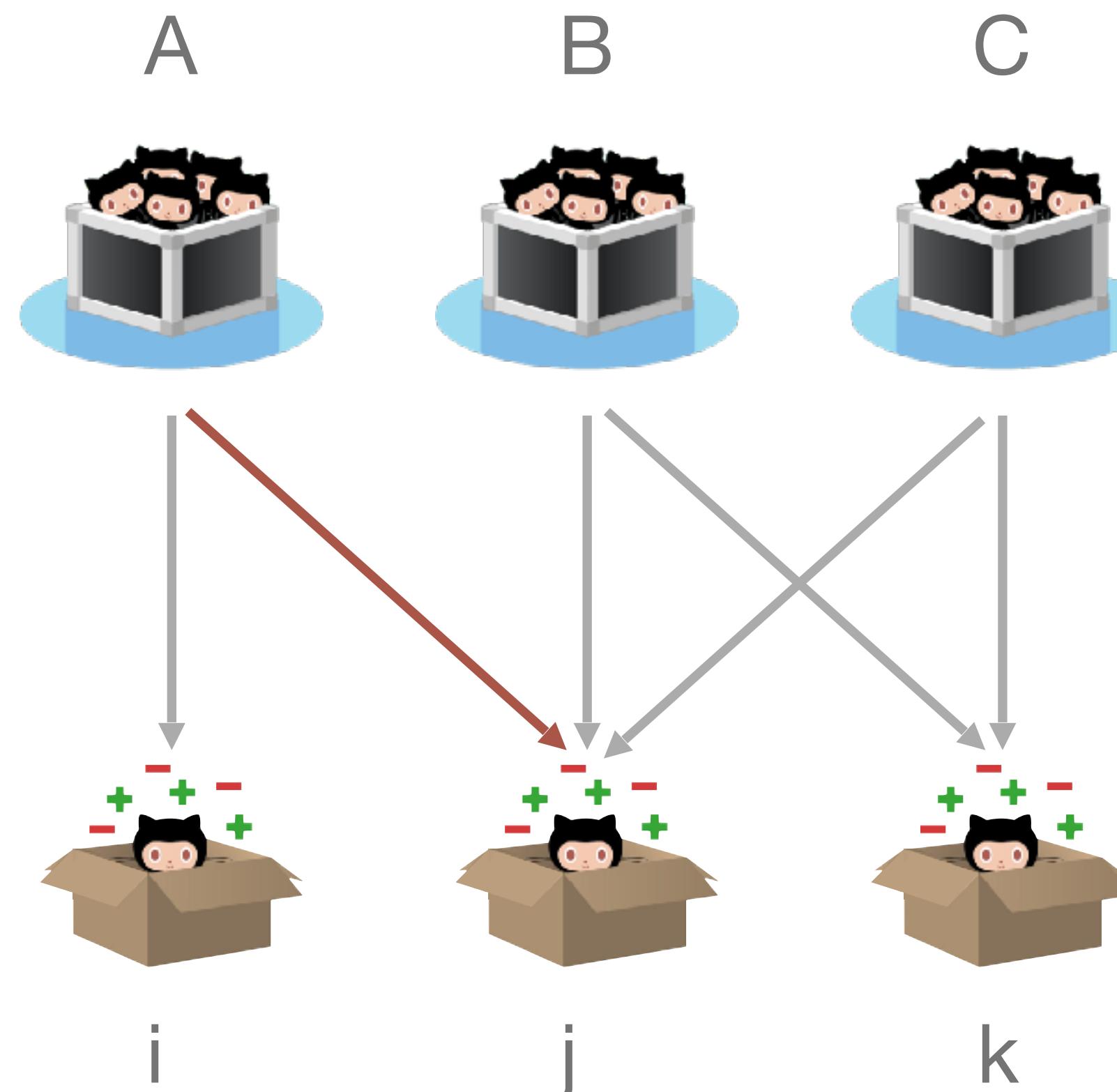
$C(n,2)$  unique pairs of packages.

Dark chocolate + apple strudel is arguably innovative because it is atypical.

# Key idea from network science: Comparison to null (random) model

Observed reality:

Projects:

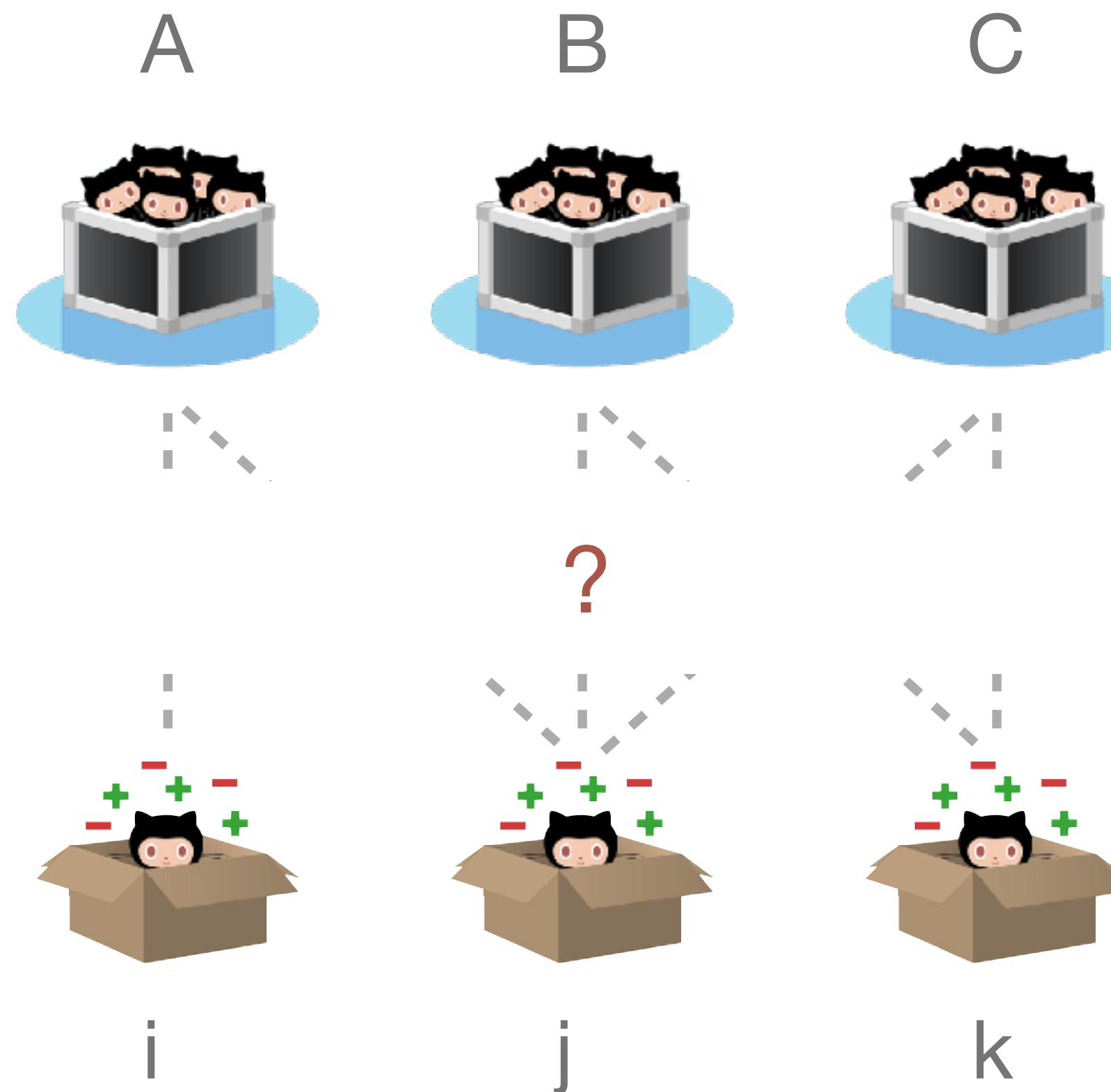


Project A adds a dependency on package j.  
New combinations are formed, e.g., (i, j).  
How atypical is (i, j)?

# Key idea from network science: Comparison to null (random) model

Counterfactual:

Projects:



Preserve:

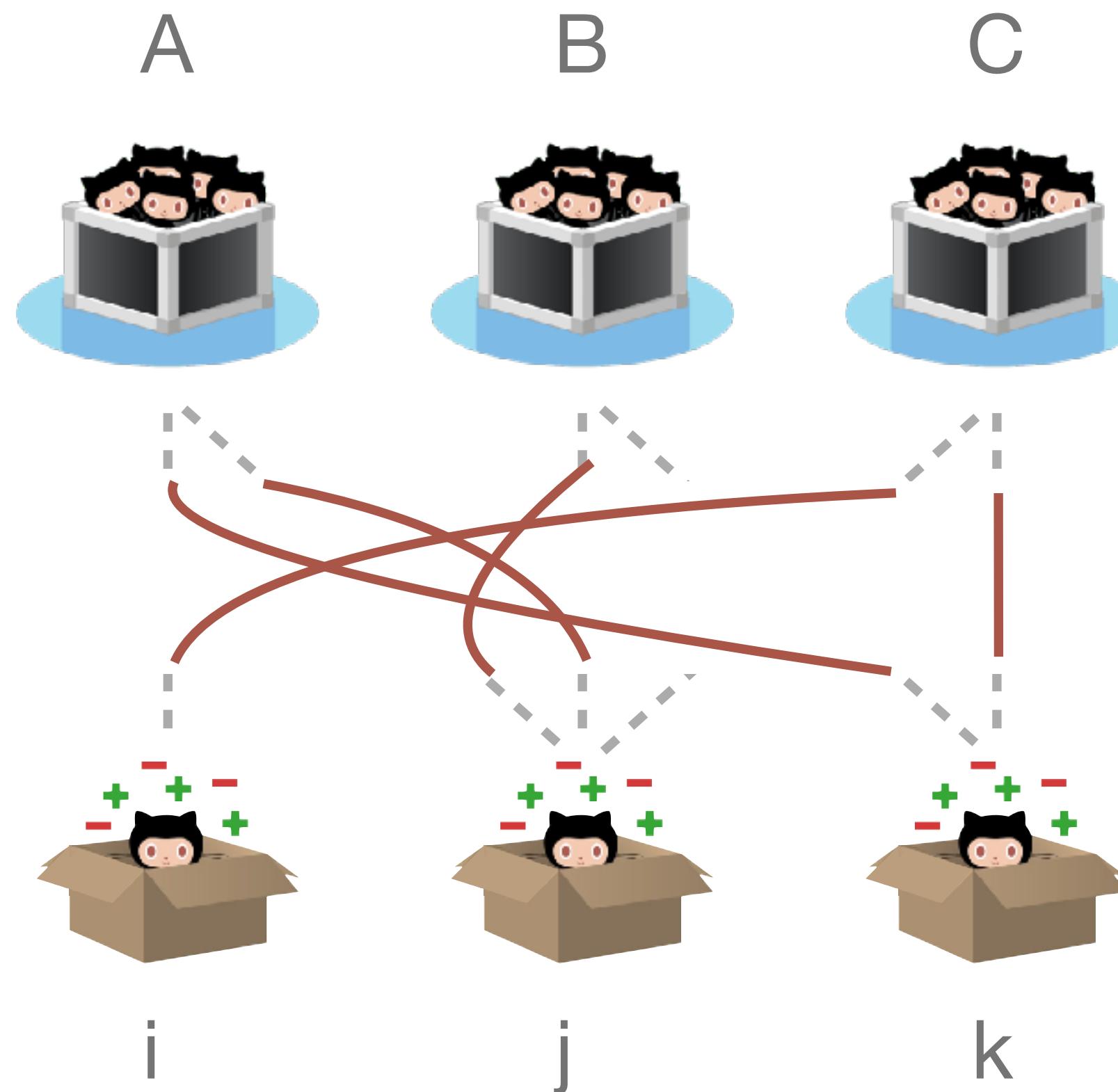
- all the projects
- all the libraries
- the distribution of imports per project
- the distribution of imports per library

Libraries:

# Key idea from network science: Comparison to null (random) model

Counterfactual:

Projects:



Libraries:

Preserve:

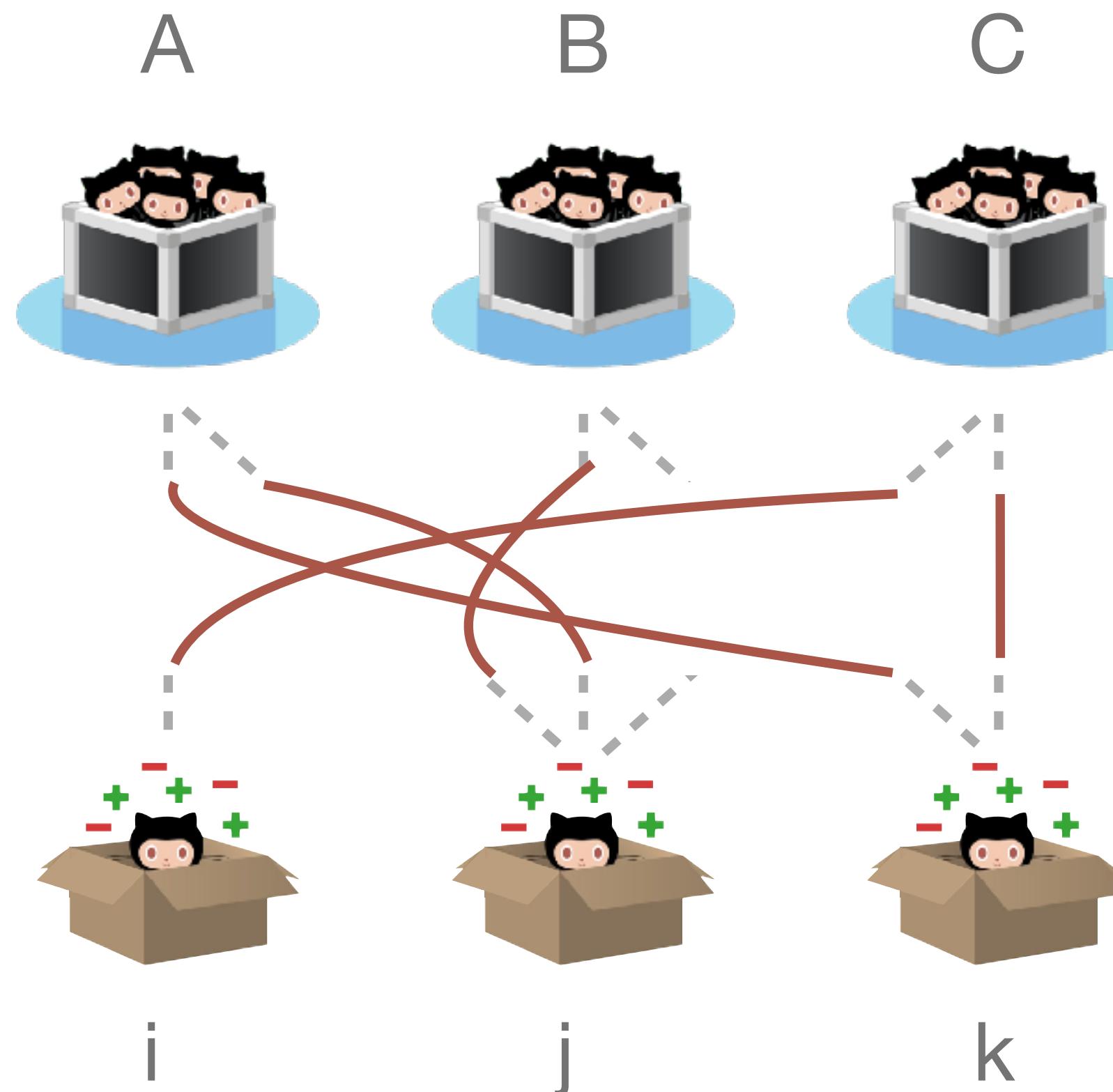
- all the projects
- all the libraries
- the distribution of imports per project
- the distribution of imports per library

But randomly rewire the network.

# Key idea from network science: Comparison to null (random) model

Counterfactual:

Projects:



Libraries:

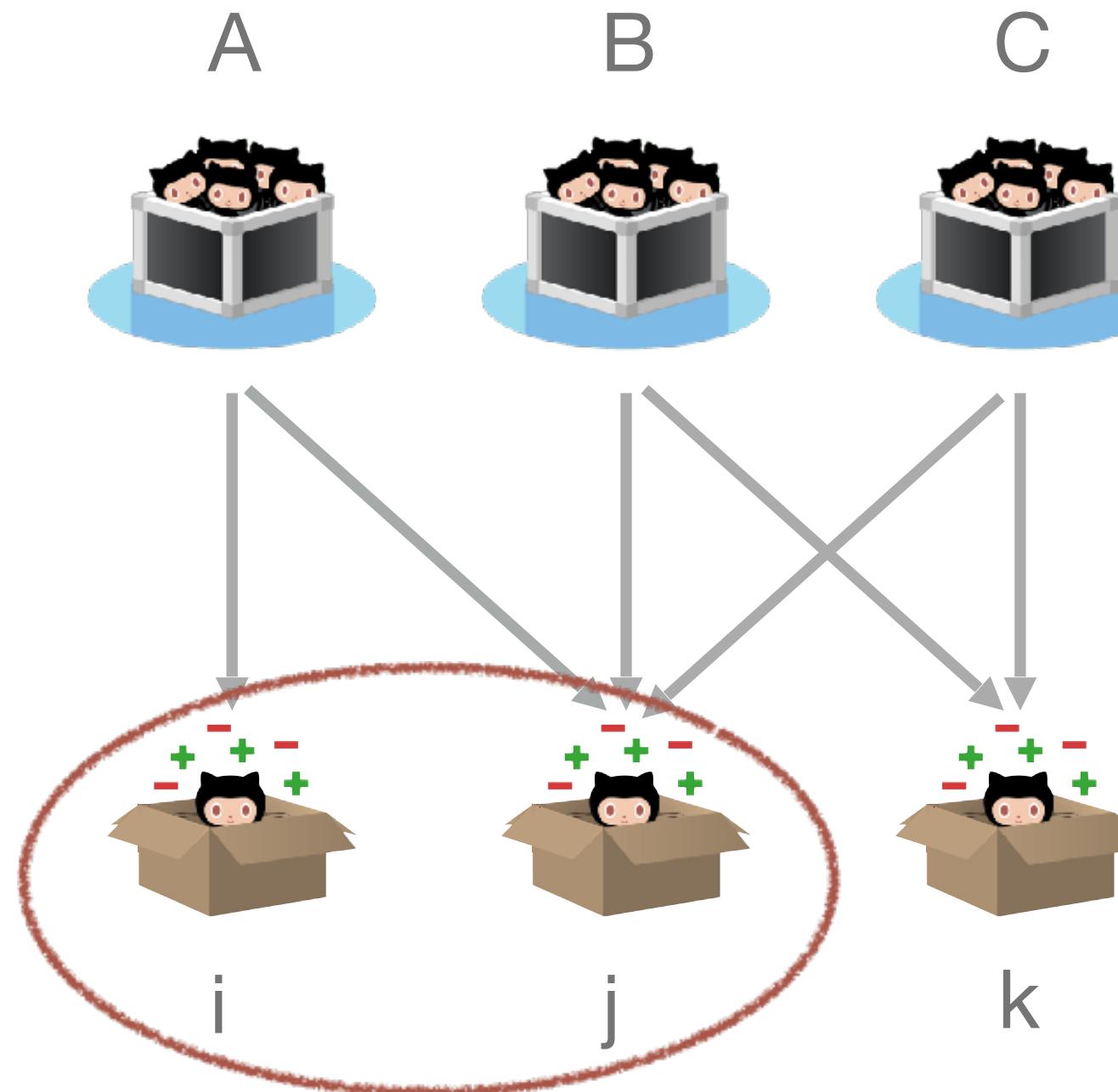
Preserve:

- all the projects
- all the libraries
- the distribution of imports per project
- the distribution of imports per library

But randomly rewire the network.

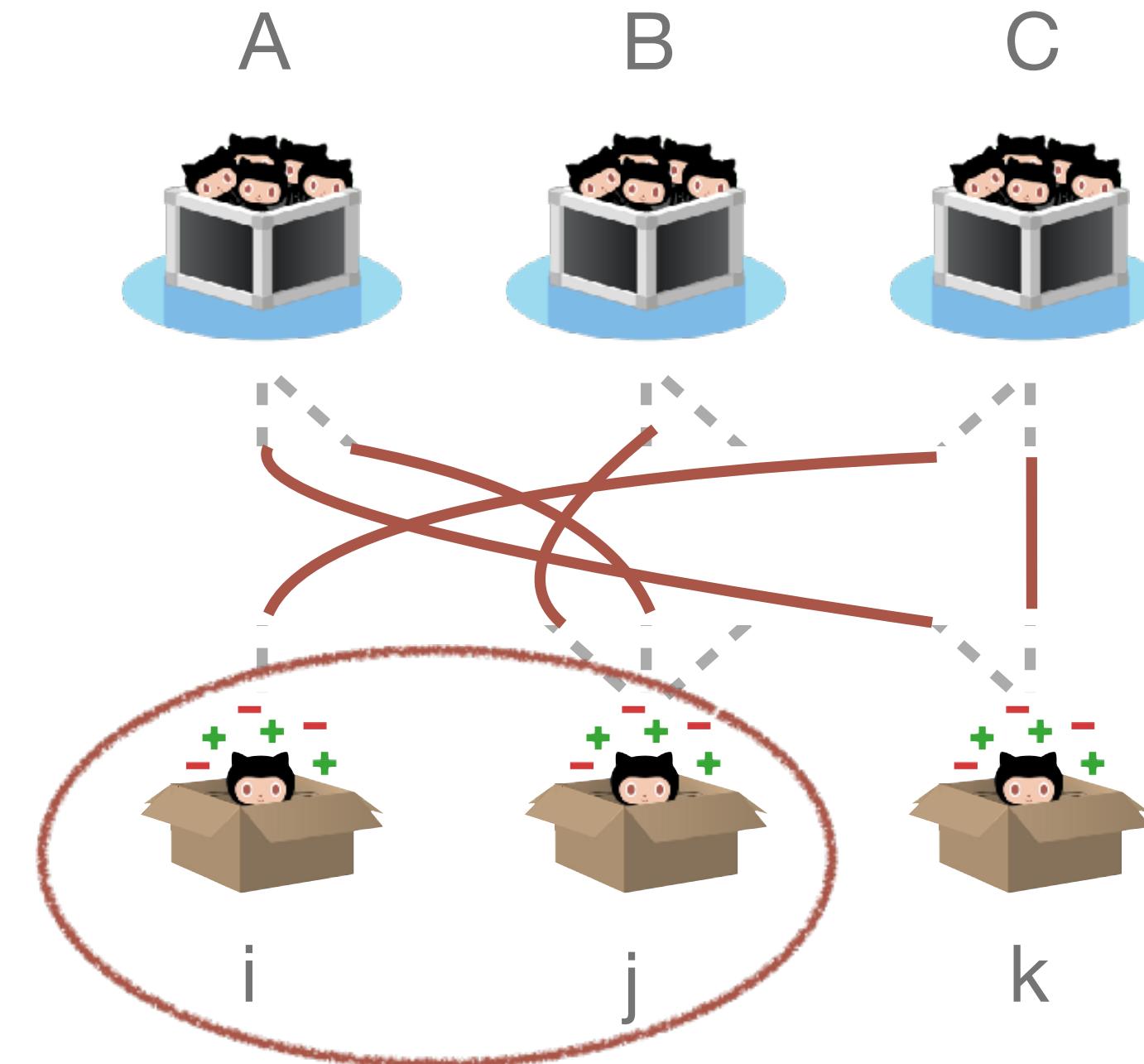
And repeat many times.

This z-score estimates if two packages are used together more, less, or about as much as could be expected by chance.



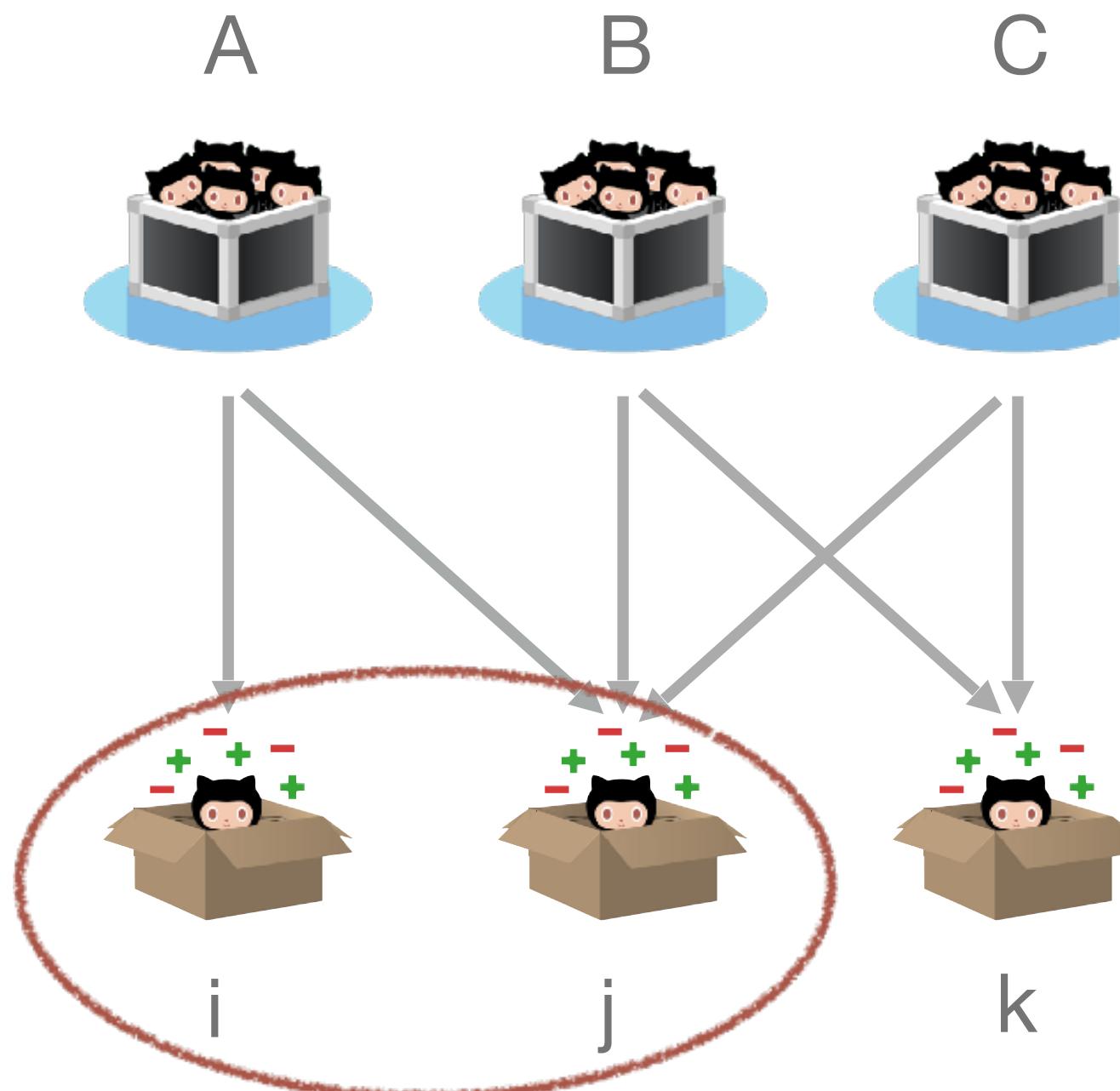
Observed number of times packages  $i$  and  $j$  appeared together until year  $t$ .

$$z_{ijt} = (obs_{ijt} - exp_{ijt}) / (\sigma_{ijt})$$



Average (i.e., expected) number of times packages  $i$  and  $j$  appeared together over  $N$  simulations.

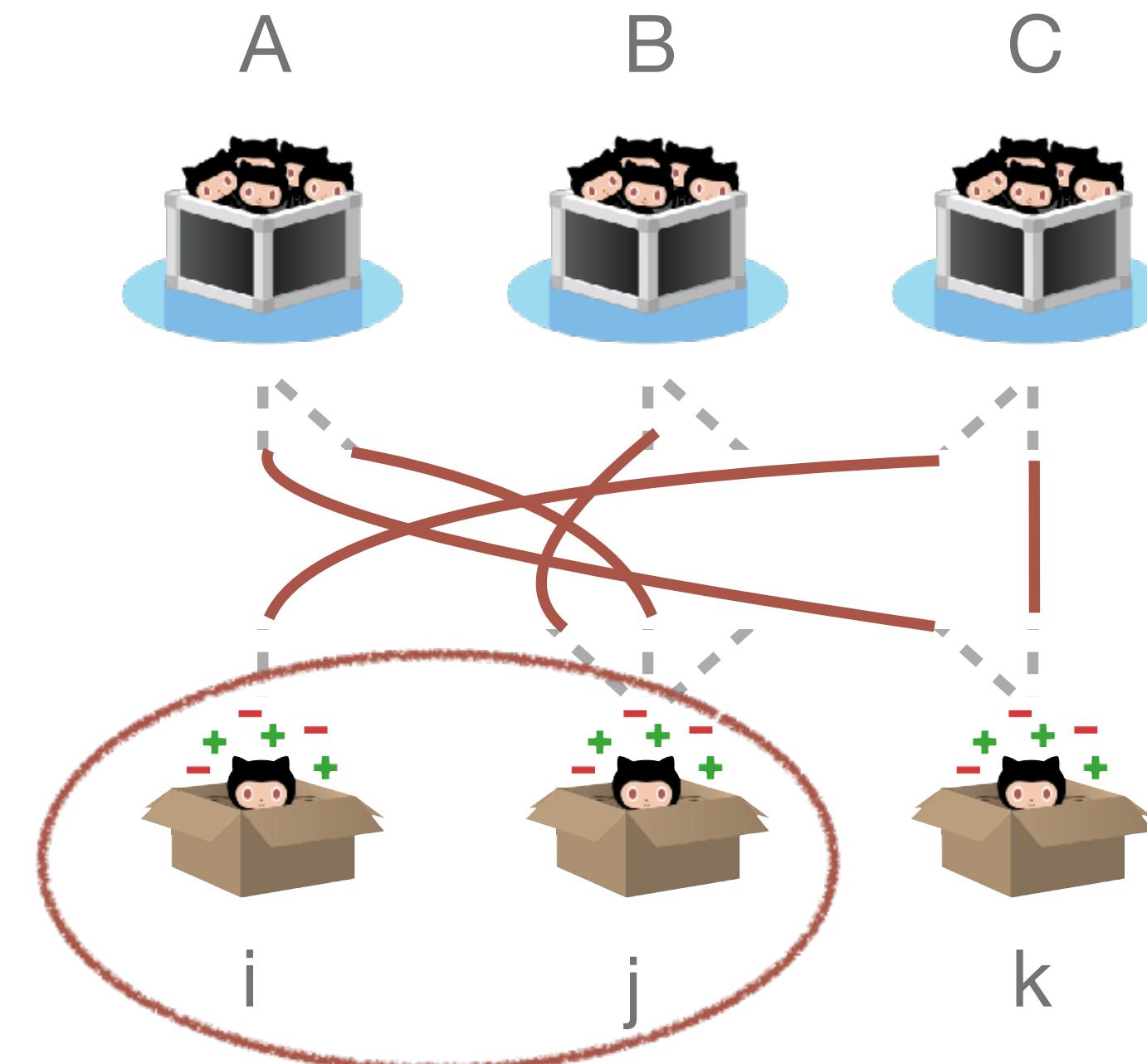
This z-score estimates if two packages are used together more, less, or about as much as could be expected by chance.



Observed number of times packages  $i$  and  $j$  appeared together until year  $t$ .

$$z_{ijt} = \frac{obs_{ijt} - exp_{ijt}}{\sigma_{ijt}}$$

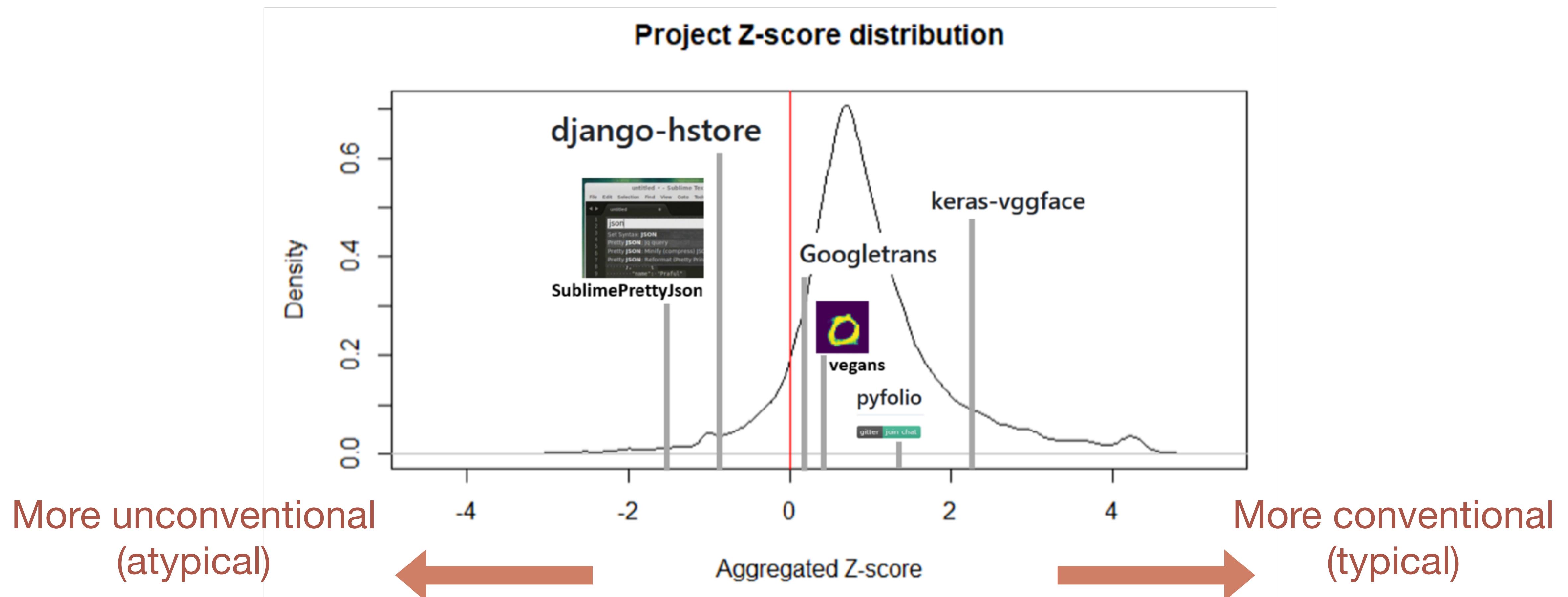
low ↘      high ↗  
⇒ atypical combination



Average (i.e., expected) number of times packages  $i$  and  $j$  appeared together over  $N$  simulations.

# Project-level aggregation is the average of pairwise atypicality z-scores

On average, projects are quite conventional.

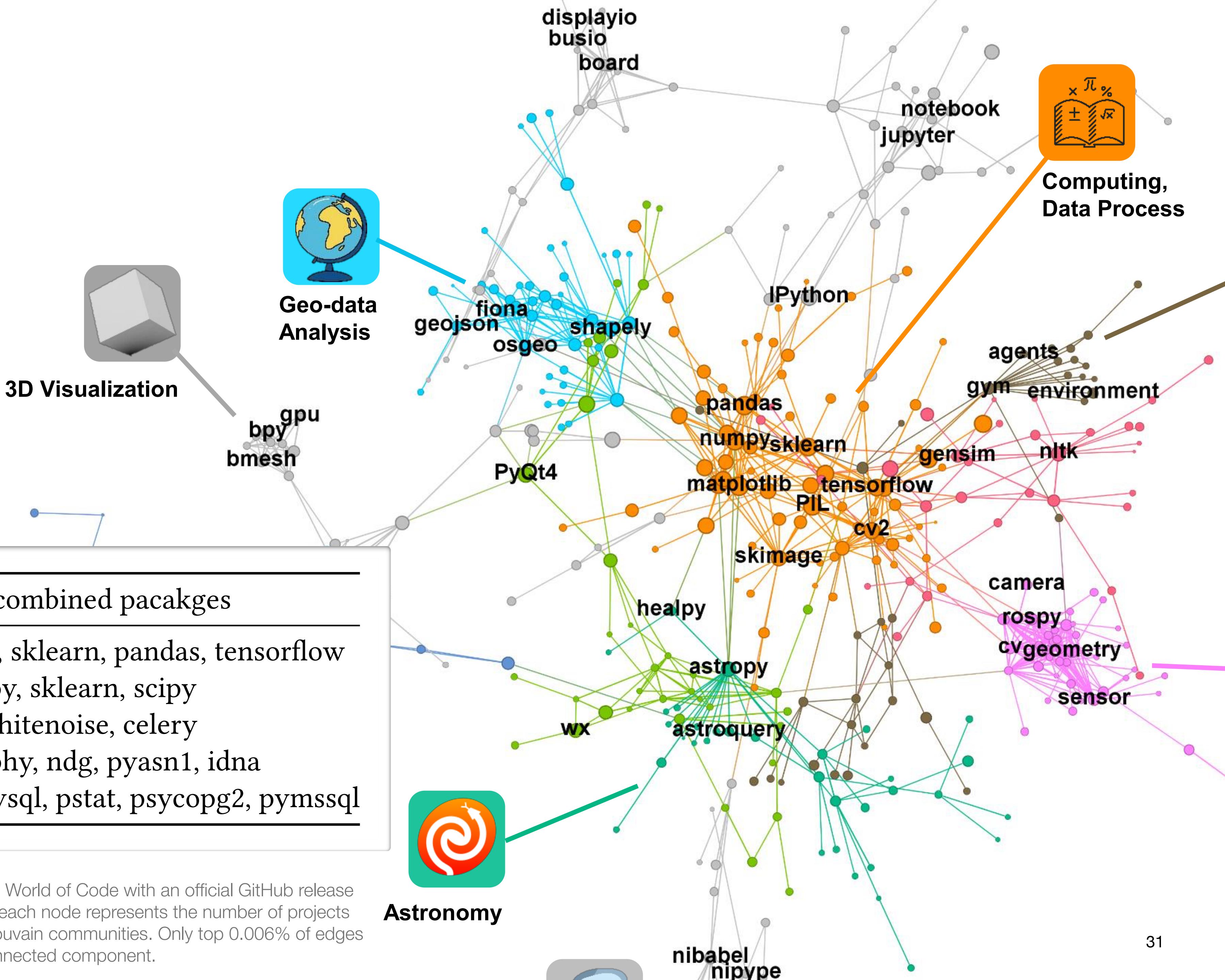


# Sanity checking

No ground truth on **atypical** package combinations, but at least the **typical** combinations should be meaningful!

Focal package	Top five mostly combined packages
numpy	matplotlib, scipy, sklearn, pandas, tensorflow
tensorflow	keras, cv2, numpy, sklearn, scipy
django	rest, dj, south, whitenoise, celery
OpenSSL	ntlm, cryptography, ndg, pyasn1, idna
pymysql	MySQLdb, aiomysql, pstat, psycopg2, pymssql

Fine print: Starting data consists of all Python projects in World of Code with an official GitHub release (75,388 projects and 7,728 packages total). The size of each node represents the number of projects that imported the package by 2019. Colors represent Louvain communities. Only top 0.006% of edges with the highest z-score shown, and only the largest connected component.



# Software innovation as novel recombination of software libraries

---

Combining software libraries that are not often used together is like using unusual ingredients in your cooking.

- People may be impressed by your culinary creativity.
- Serving unusual dishes can be risky if the chefs are unable to perfect the recipes and the customers are unwilling to try new things.



<https://www.tasteofhome.com/recipes/chocolate-peanut-butter-pizza/>

# Software innovation as novel recombination of software libraries

---

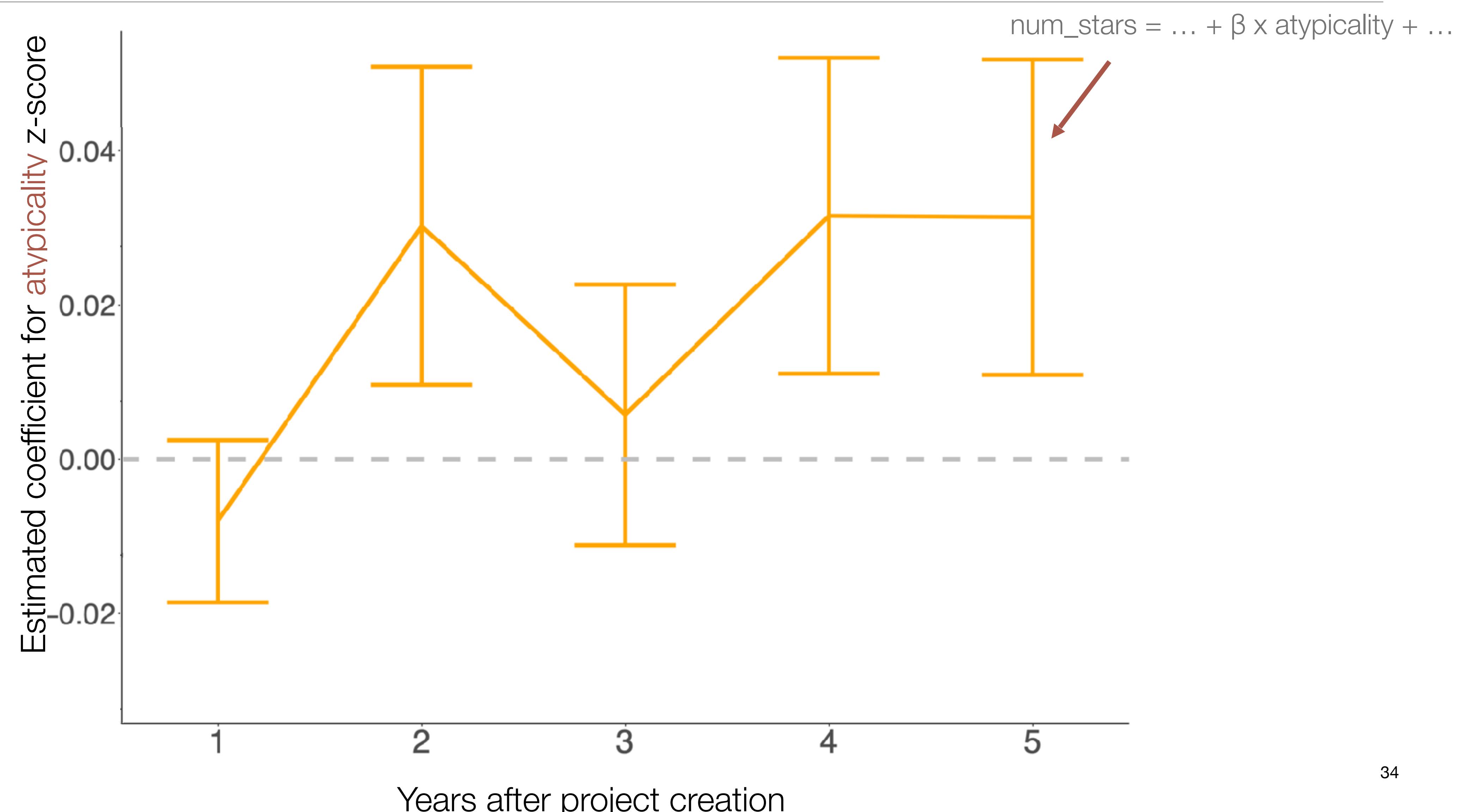
Combining software libraries that are not often used together is like using unusual ingredients in your cooking.

- Hyp: Projects that use more atypical combinations of libraries tend to be **more popular**.
- Hyp: More innovative projects tend to be **less sustainable** in the long term.

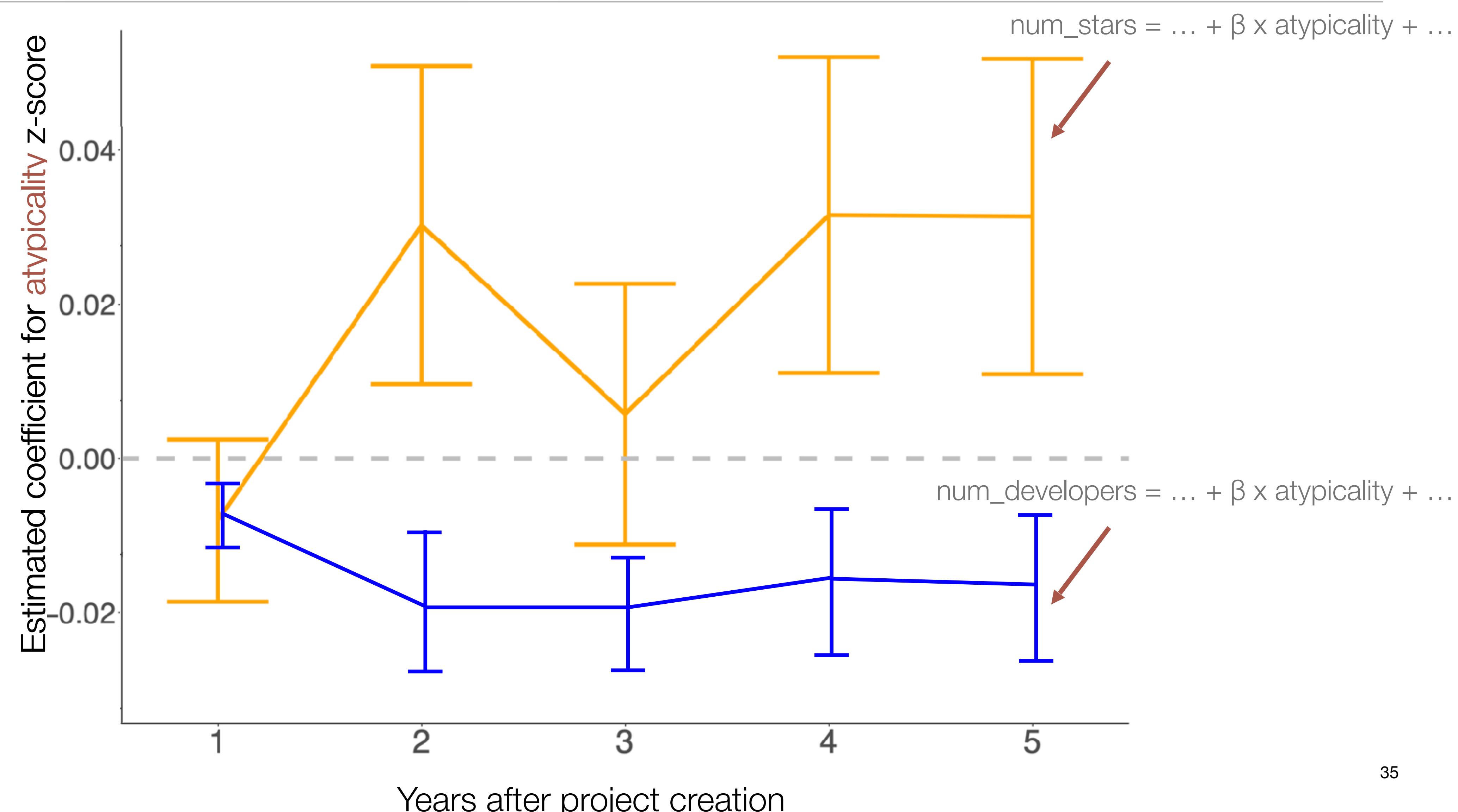


<https://www.tasteofhome.com/recipes/chocolate-peanut-butter-pizza/>

Atypical (novel) projects tend to have more stars.



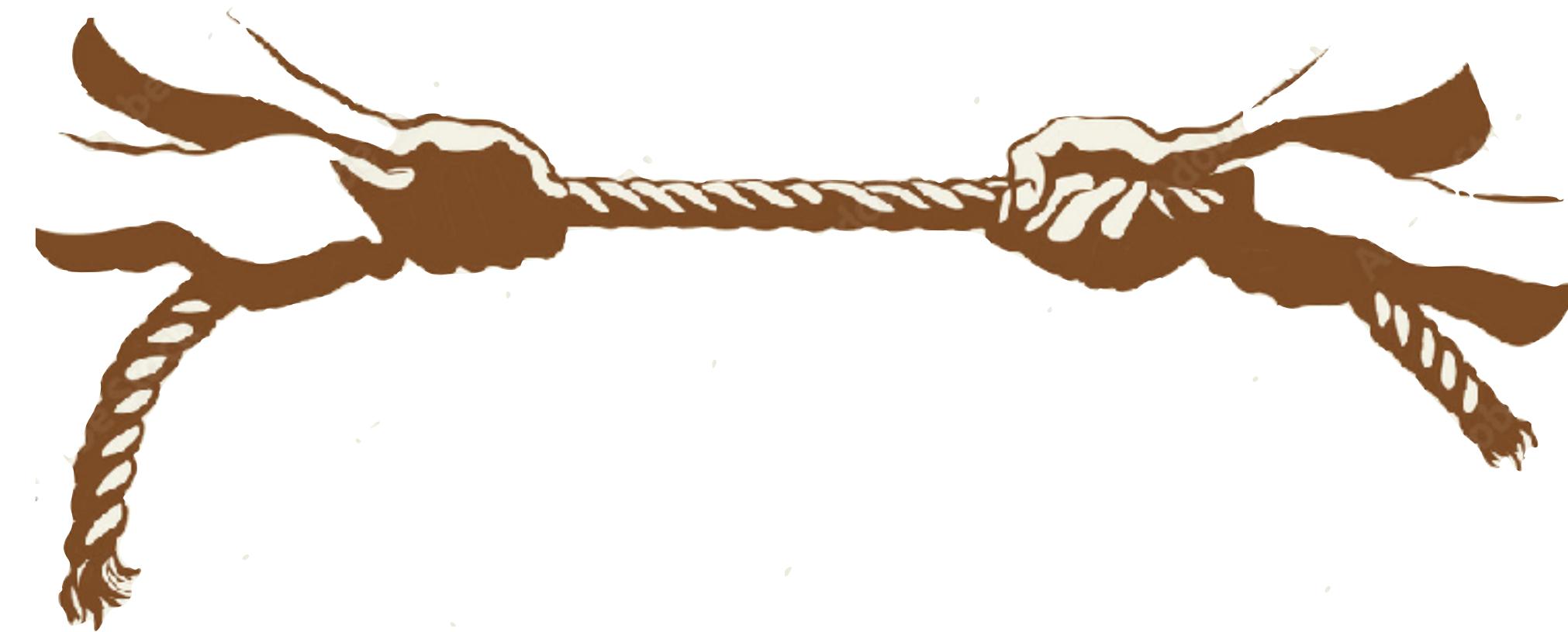
Atypical (novel) projects tend to have smaller teams (and higher probability of becoming abandoned).



# Tension between innovation and open source sustainability?

---

Incentive to create  
ever-new things



The “grunt work”  
of maintaining  
existing systems

- Creative expression is a main driver of contributing to open source
- Innovation seems to be rewarded with increased popularity

Will it become increasingly harder to ensure that sufficient maintenance attention (developers, funding, etc) is being allocated to the projects that need it the most?

# Today: Let's look at some concrete examples of network effects

---



Measuring innovation  
in software



Understanding how  
innovation emerges



Social capital



Social contagion

# Once upon a time, a PhD student at Harvard University was writing their dissertation ...

---

**Stanford**  
Sociology  
SCHOOL OF HUMANITIES AND SCIENCES

## Mark Granovetter

Joan Butler Ford Professor  
in the School of Humanities  
and Sciences; Professor of  
Sociology

A.B. Princeton University 1965  
Modern European and American  
History  
Ph.D. Harvard University 1970  
Sociology



<https://sociology.stanford.edu/people/mark-granovetter>

## The Strength of Weak Ties<sup>1</sup>

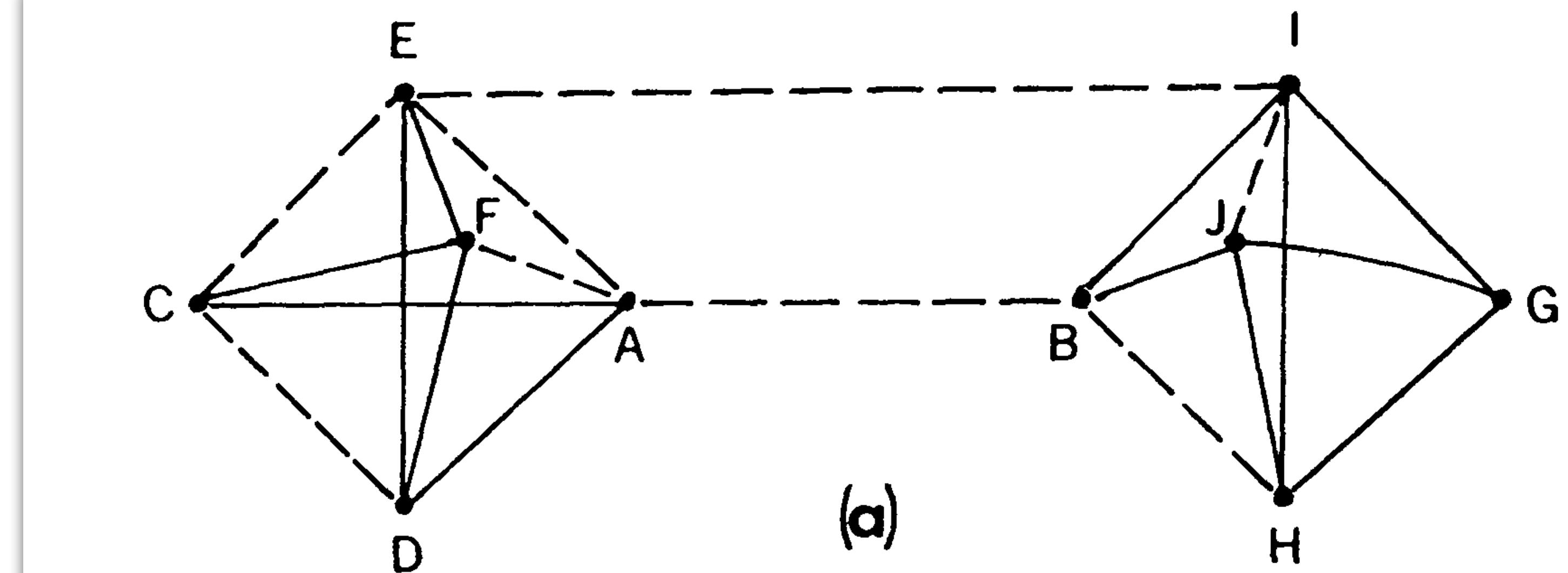
Mark S. Granovetter  
*Johns Hopkins University*

Analysis of social networks is suggested as a tool for linking micro and macro levels of sociological theory. The procedure is illustrated by elaboration of the macro implications of one aspect of small-scale interaction: the strength of dyadic ties. It is argued that the degree of overlap of two individuals' friendship networks varies directly with the strength of their tie to one another. The impact of this principle on diffusion of influence and information, mobility opportunity, and community organization is explored. Stress is laid on the cohesive power of weak ties. Most network models deal, implicitly, with strong ties, thus confining their applicability to small, well-defined groups. Emphasis on weak ties lends itself to discussion of relations *between* groups and to analysis of segments of social structure not easily defined in terms of primary groups.

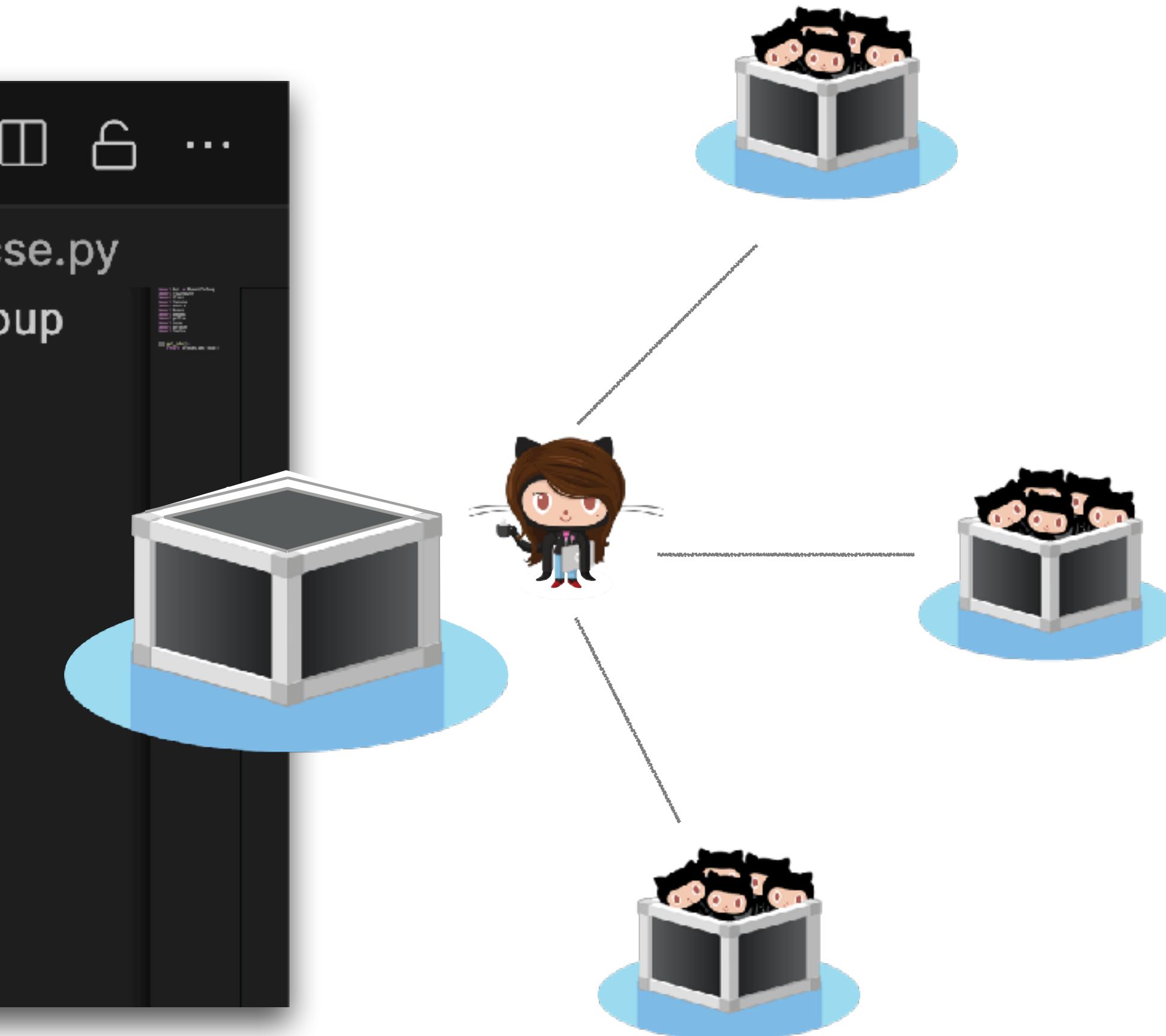
Weak ties are more effective in job searches because they act as bridges.

The majority of people found their jobs through acquaintances (weak ties) rather than close friends or family (strong ties).

In a random sample of recent professional, technical, and managerial job changers living in a Boston suburb, I asked those who found a new job through contacts how often they *saw* the contact around the time that he passed on job information to them. I will use this as a measure of tie strength.<sup>15</sup> A natural a priori idea is that those with whom one has strong ties are more motivated to help with job information. Opposed to this greater motivation are the structural arguments I have been making: those to whom we are weakly tied are more likely to move in circles different from our own and will thus have access to information different from that which we receive.



# Do OSS developers also find their new ideas through weak ties?

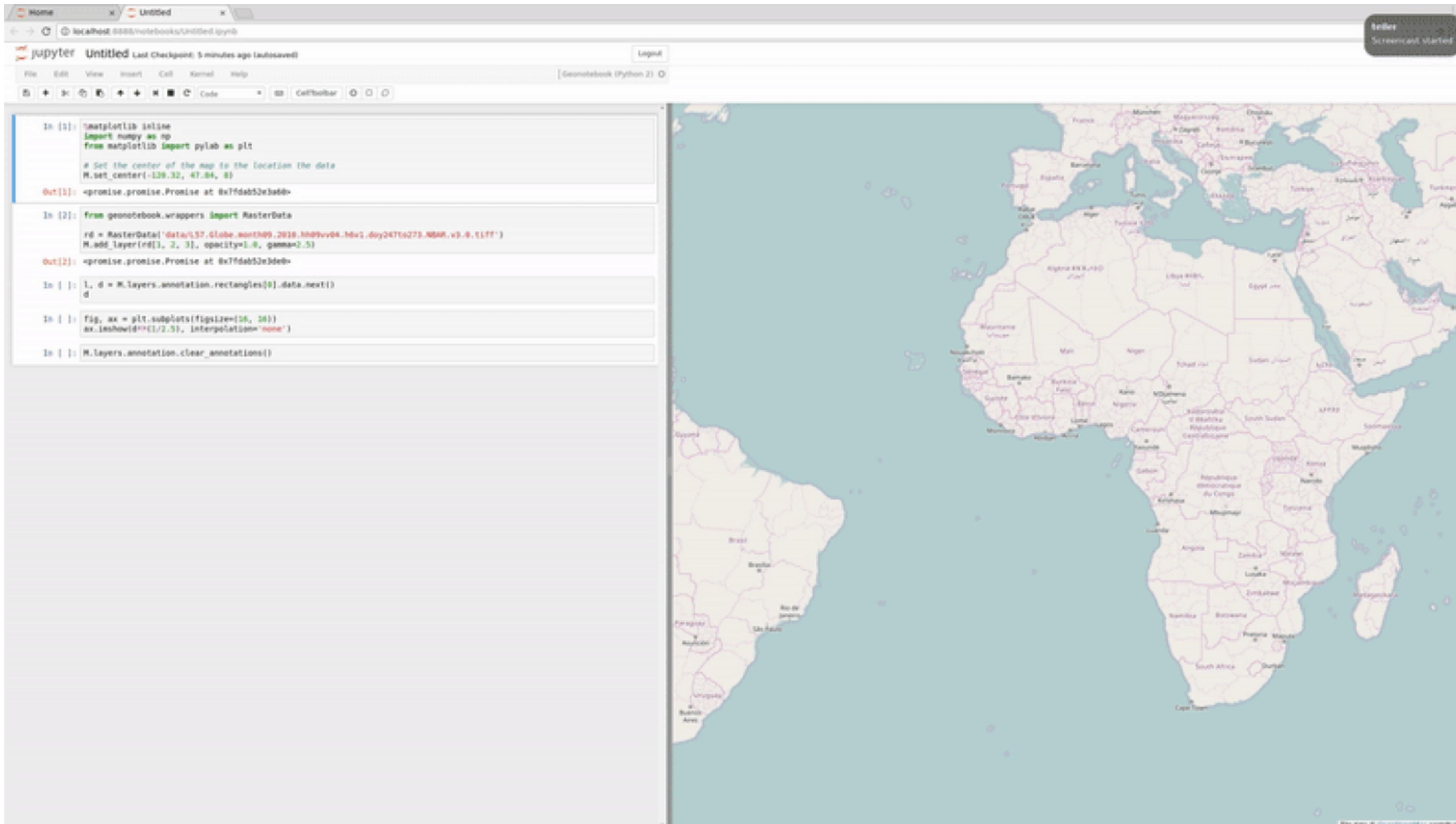


```
icse.py  X  ▶ ▾ ⏺ 🔒 ...
```

Users > bogdan > Downloads > icse.py

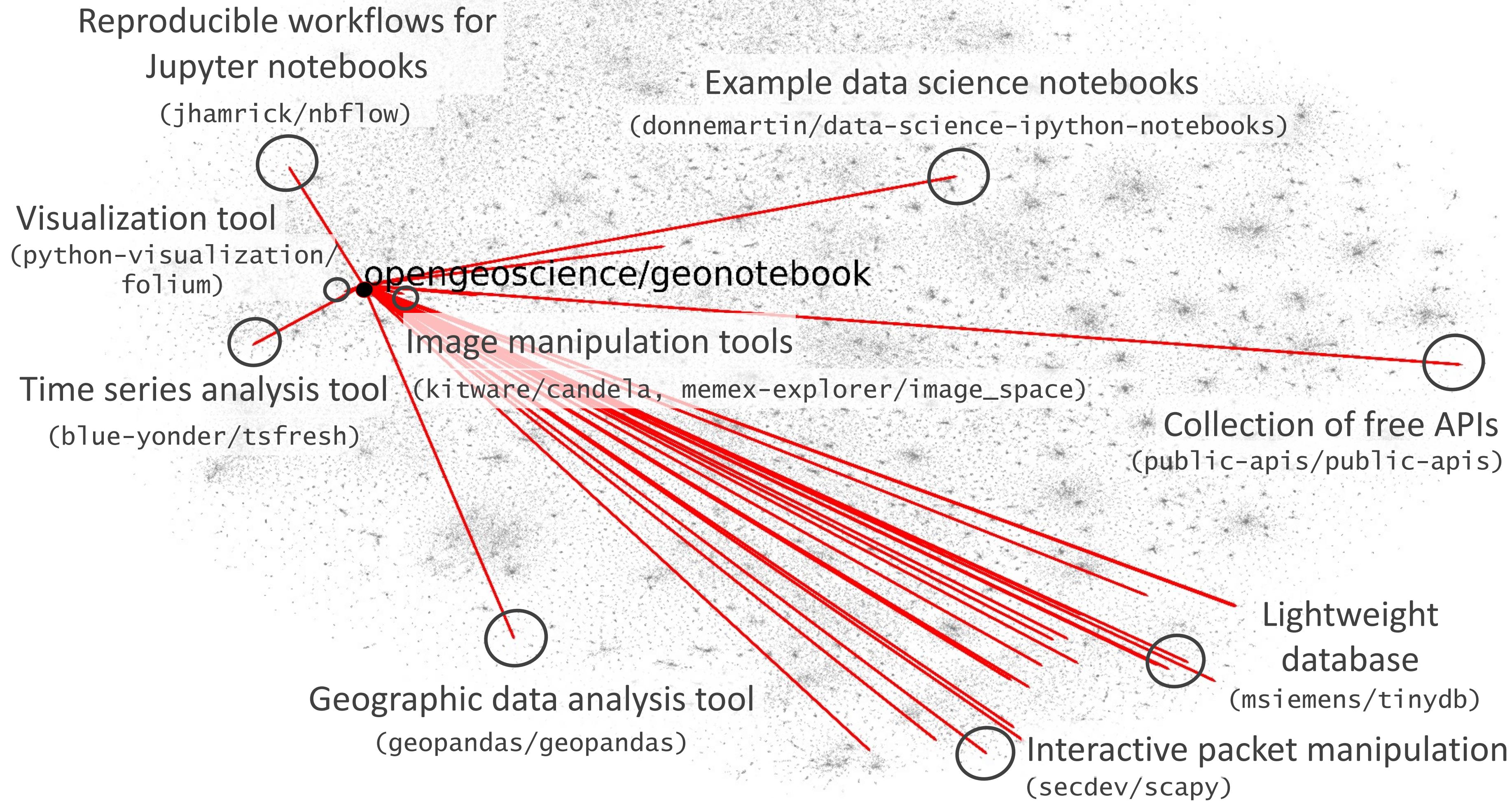
```
1 import bs4 as BeautifulSoup
2 import fuzzywuzzy
3 import flask
4 import twisted
5 import bottle
6 import black
7 import pandas
8 import pillow
9 import nose
10 import pyjokes
11 import turtle
```

# Do OSS developers also find their new ideas through weak ties?



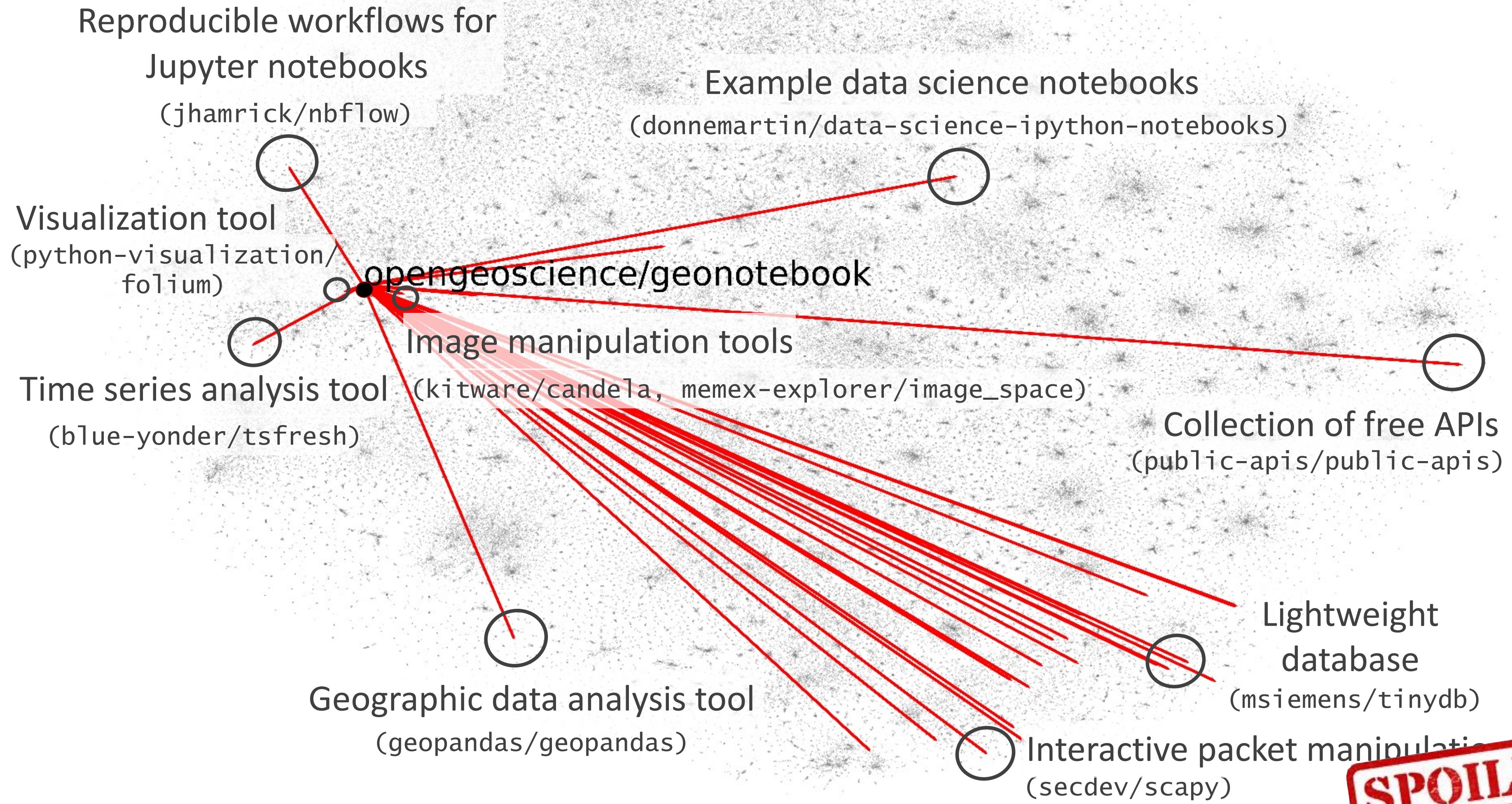
<https://github.com/opengeo-science/geonotebook>

# Do OSS developers also find their new ideas through weak ties?



Anecdotally, yes

# Do OSS developers also find their new ideas through weak ties?

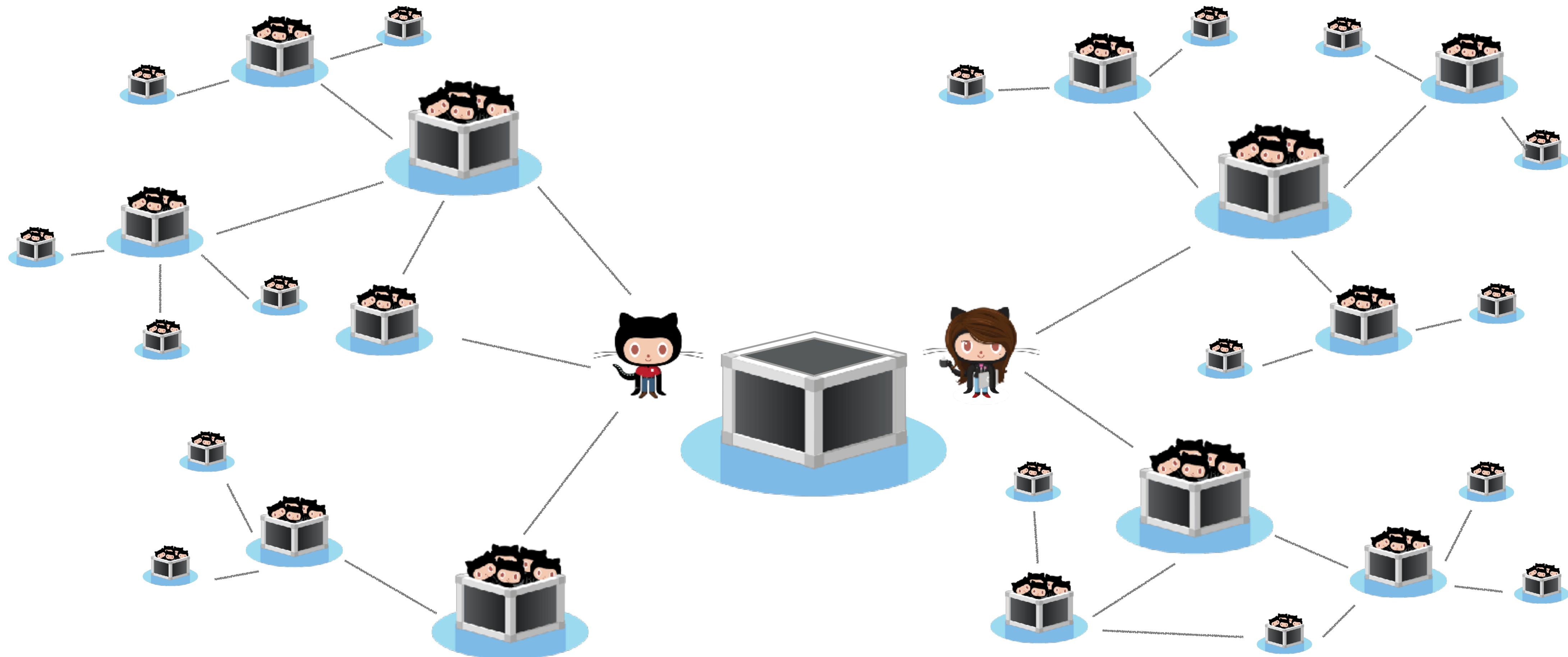


**SPOILER ALERT**

Amazingly, statistically also yes!

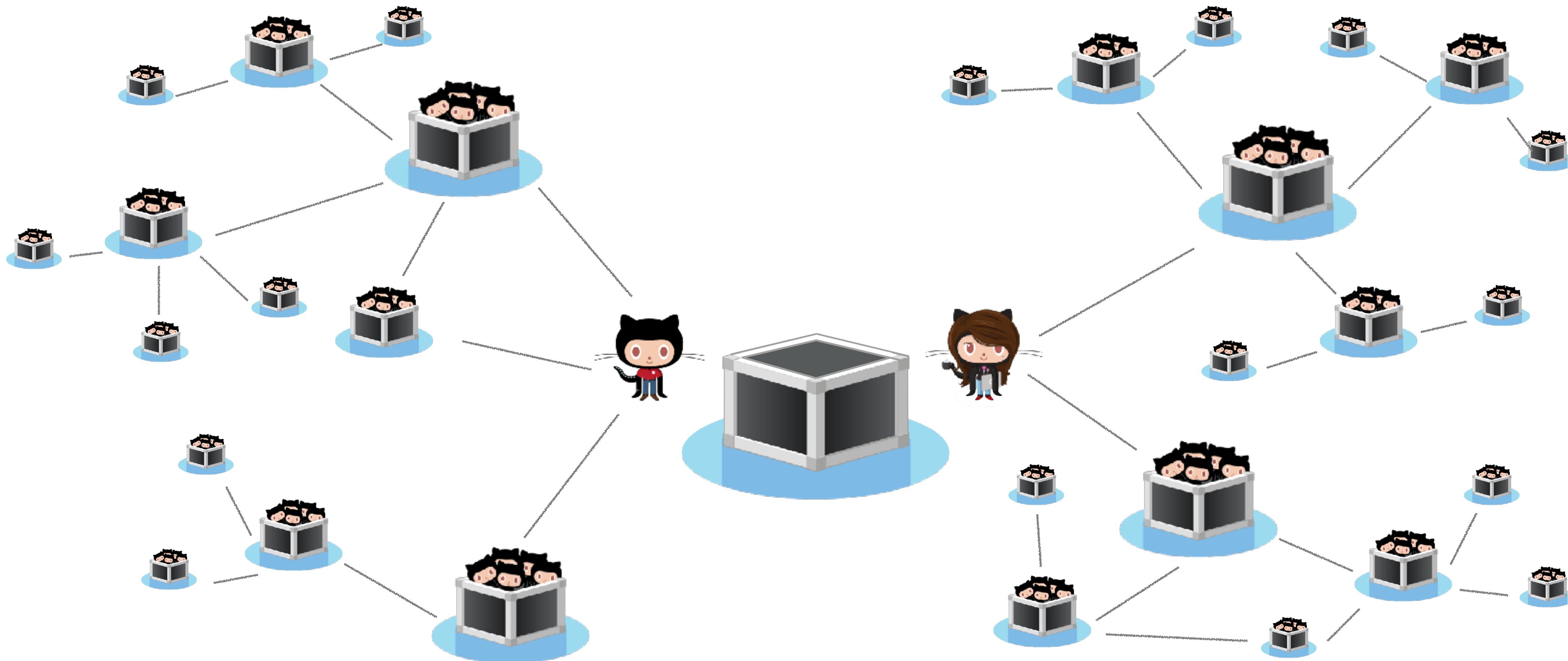
People interact with artifacts and with each other. This creates ties.

---

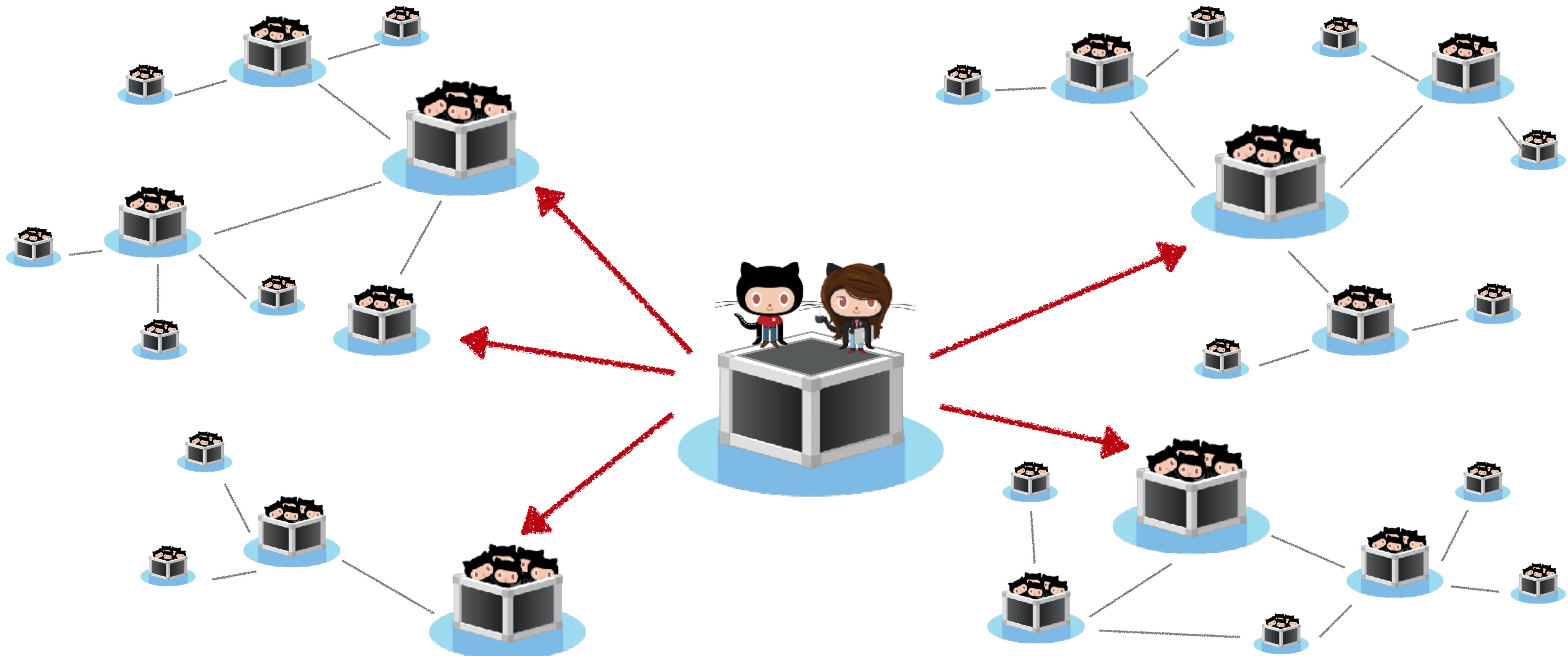


Hypothesis 1: The bigger developers' networks are, the better informed they are, and the more innovative their projects are.

---

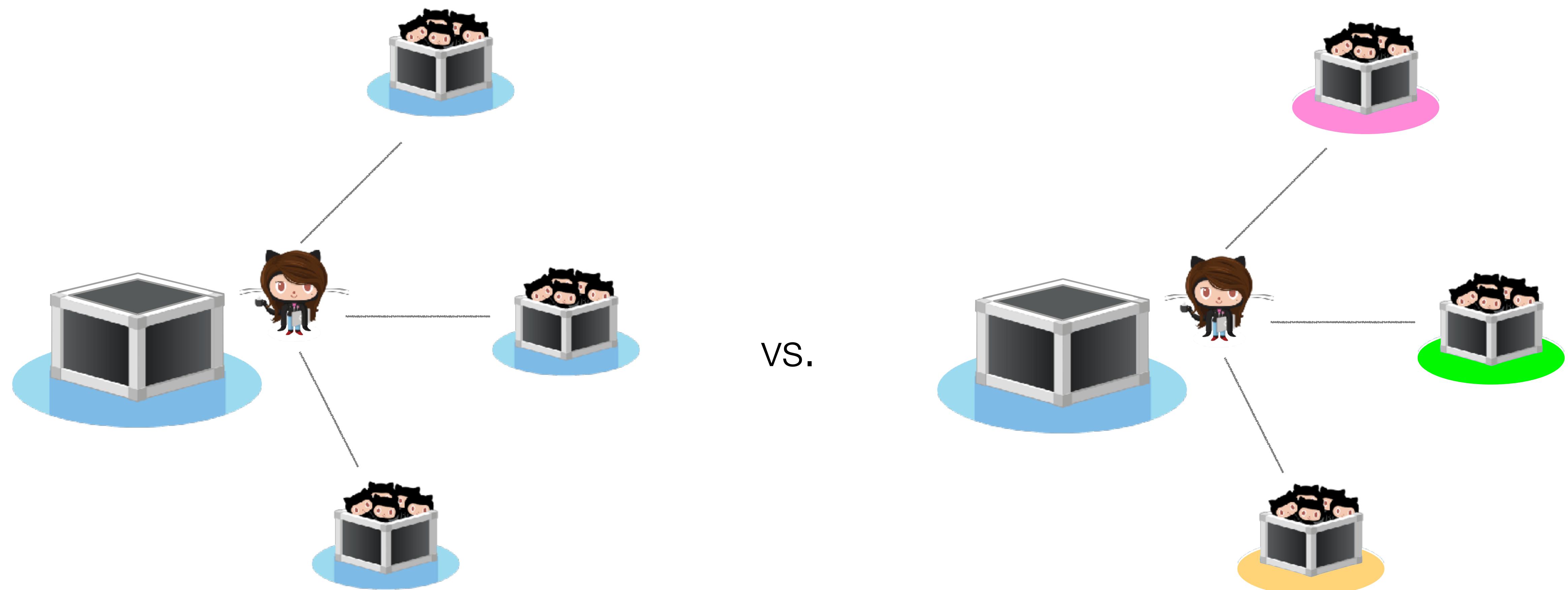


# Measure: Out-degree centrality



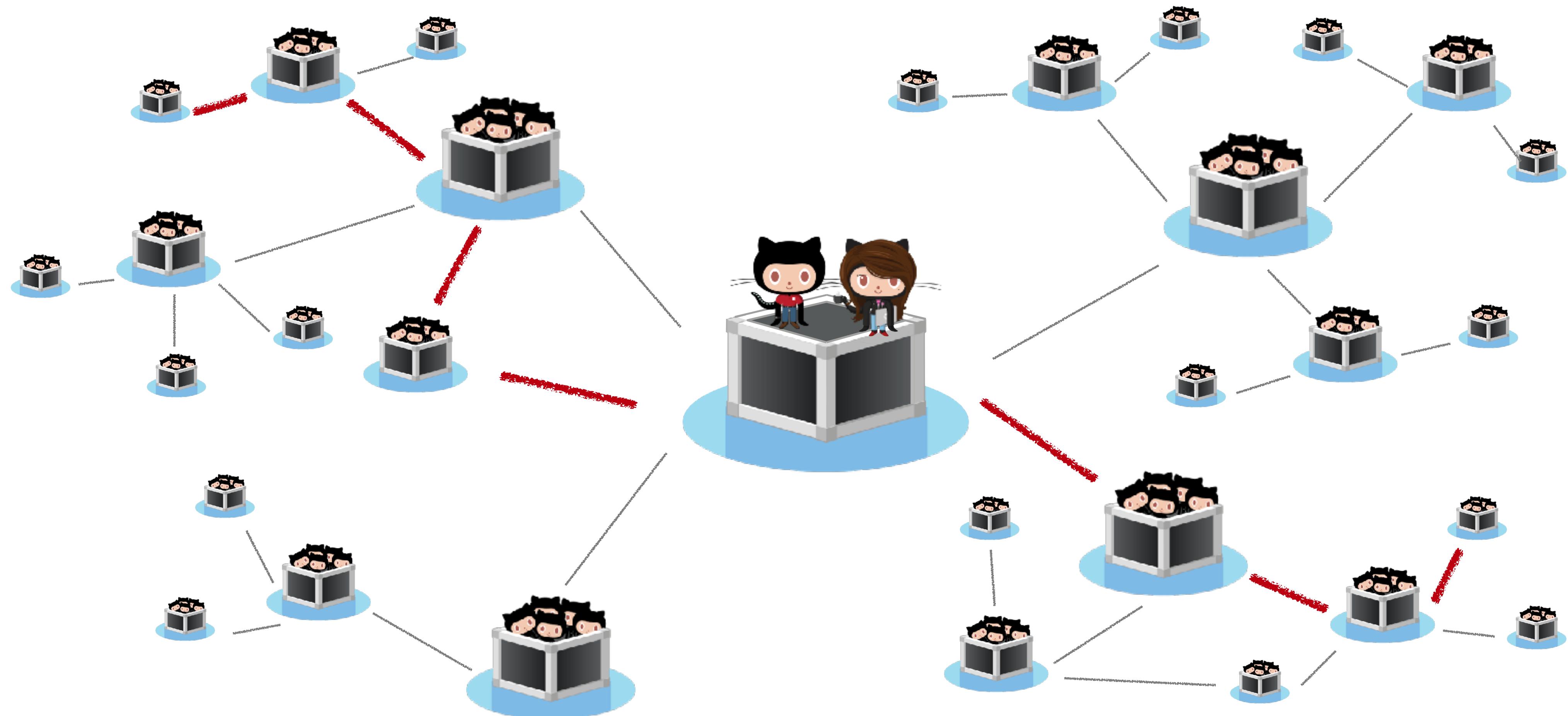
Hypothesis 2: The greater the informational diversity of developers' networks, the more innovative their projects are.

---



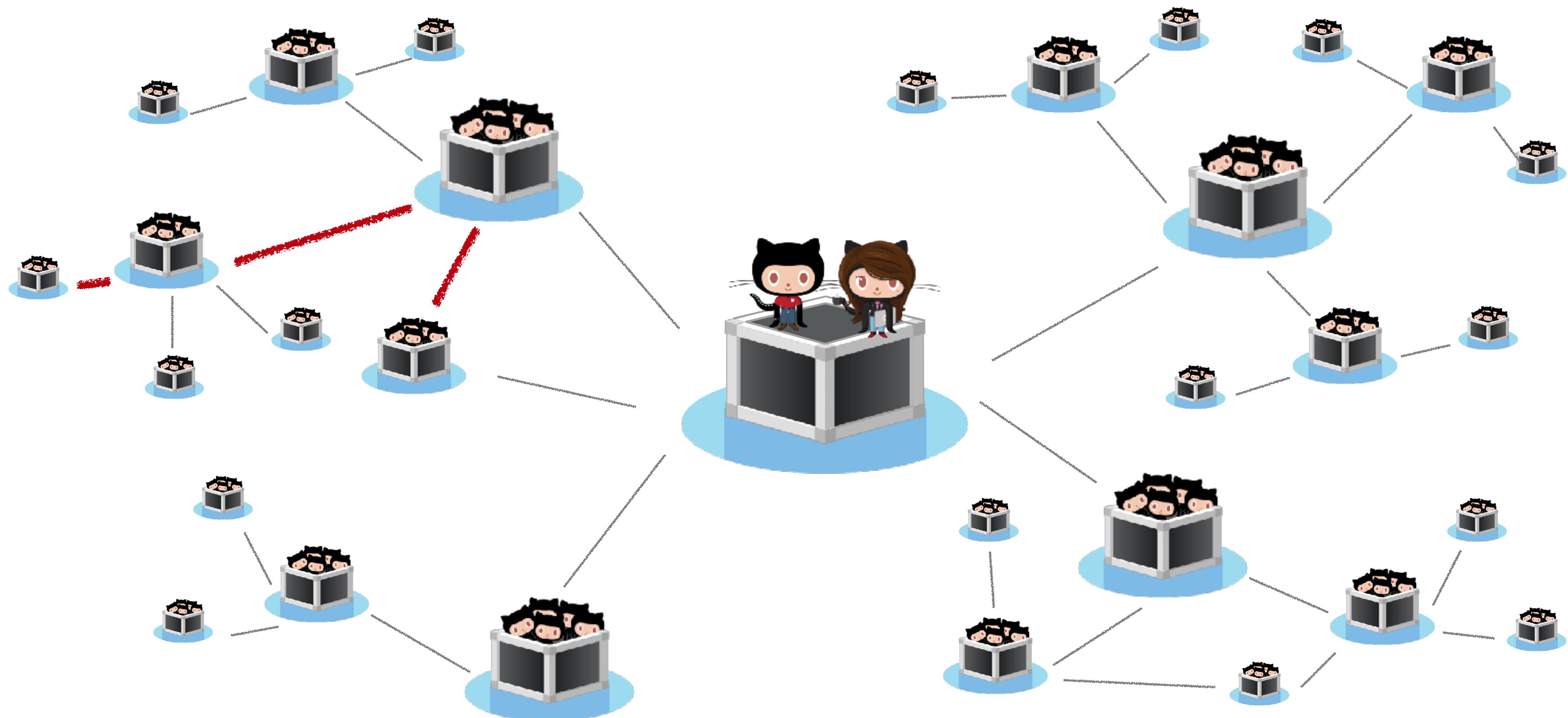
Measure: First, we generate Node2Vec embeddings for each project

---



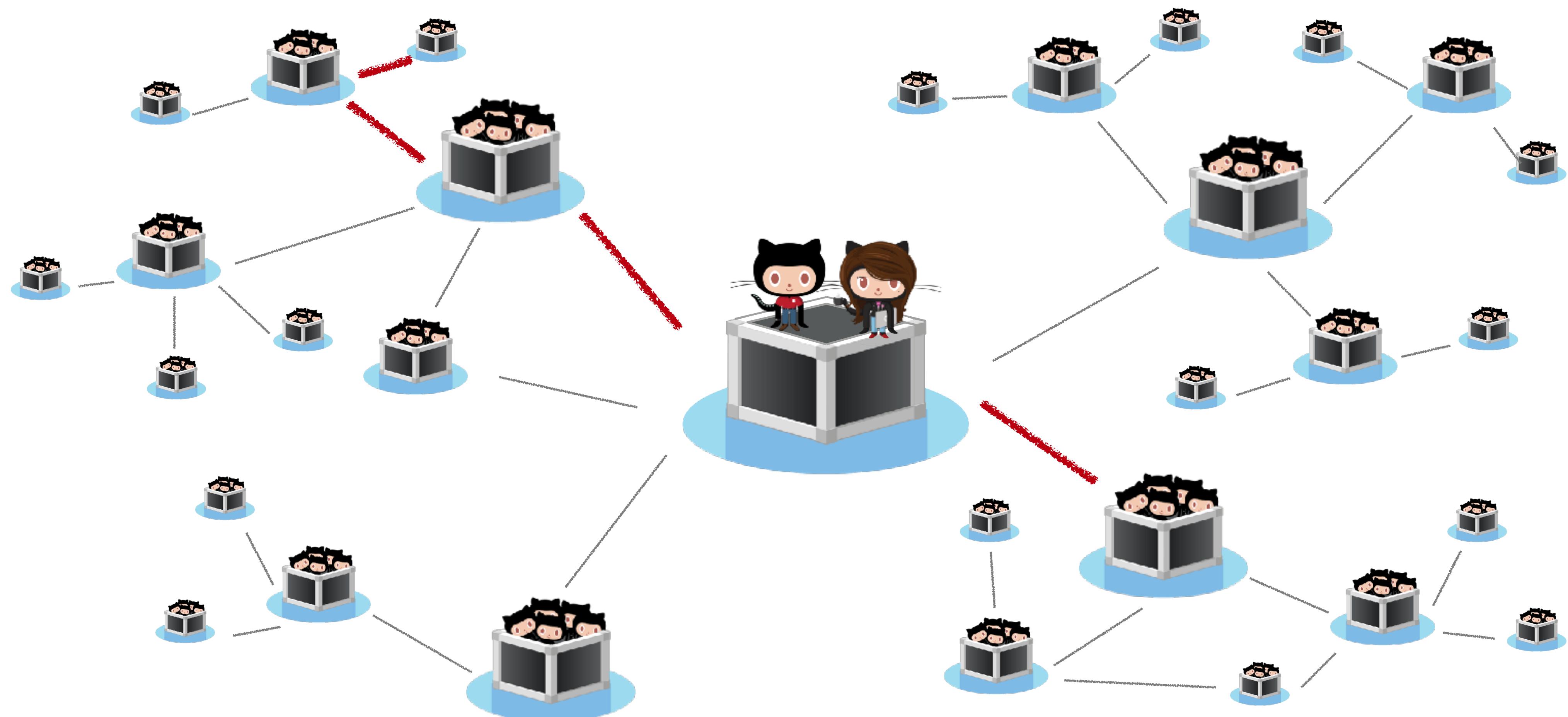
Measure: First, we generate Node2Vec embeddings for each project

---



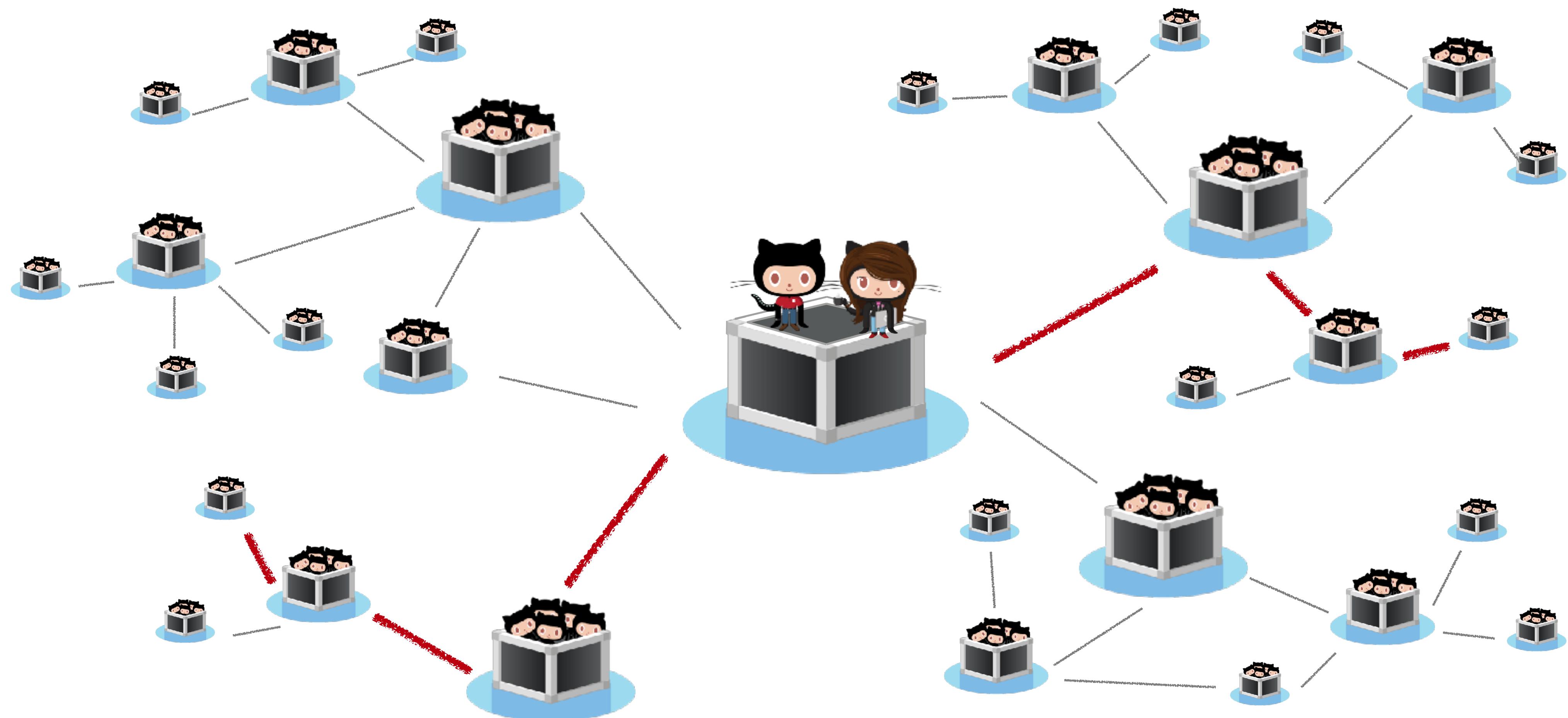
Measure: First, we generate Node2Vec embeddings for each project

---



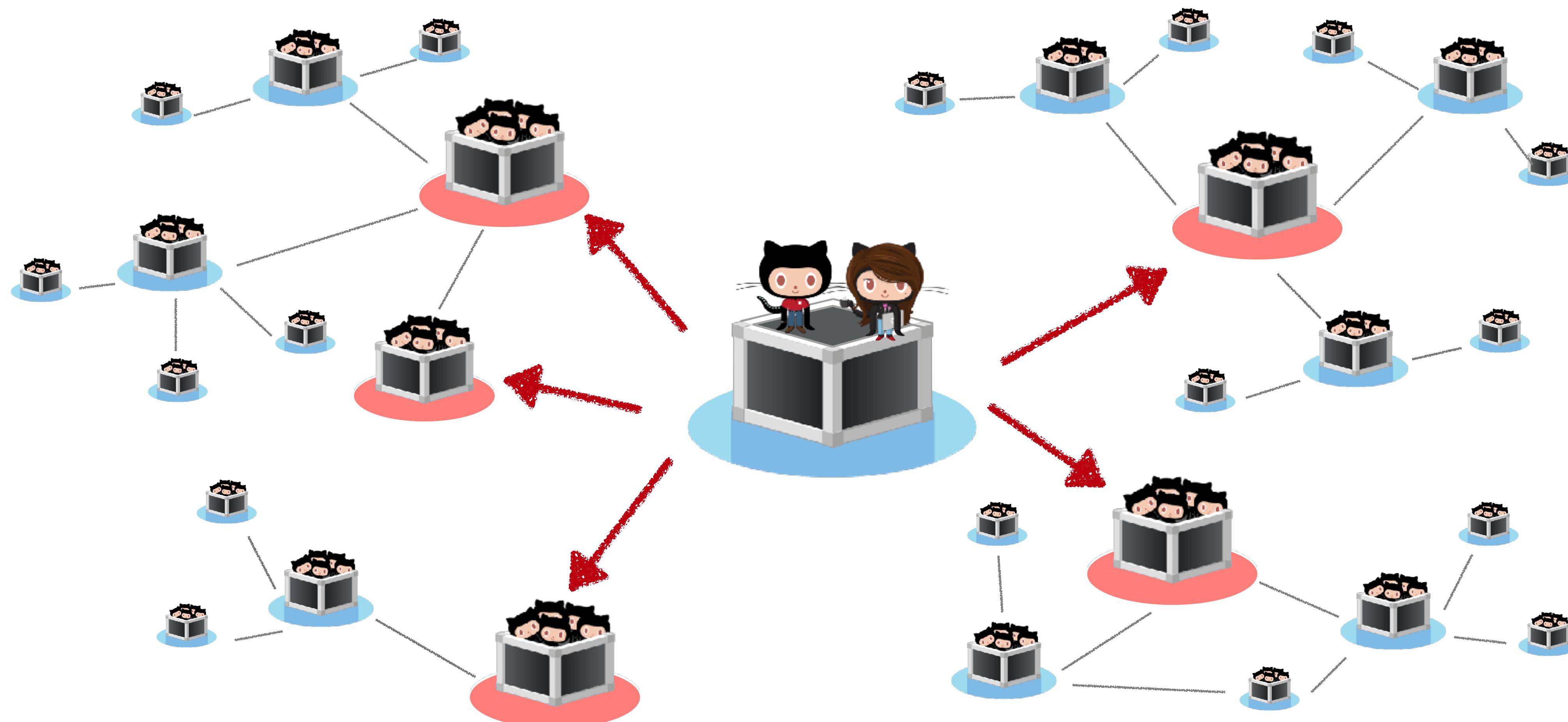
Measure: First, we generate Node2Vec embeddings for each project

---

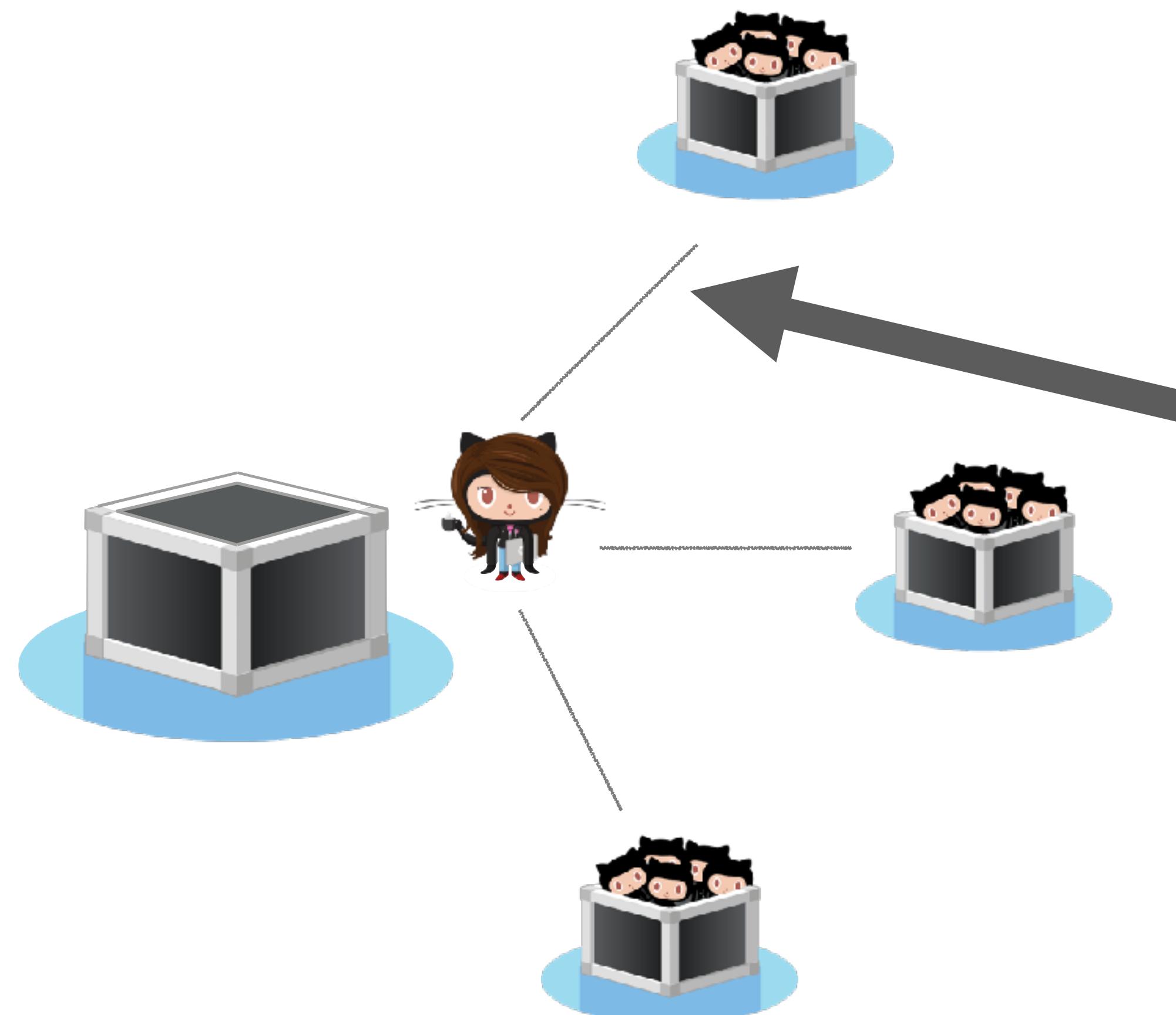


Measure: Then, we compute the average pairwise distance (inverse cosine similarity) between a focal project's direct neighbors

---



# From interactions to ties of varying strength



 1 file changed +1 -1 lines changed 

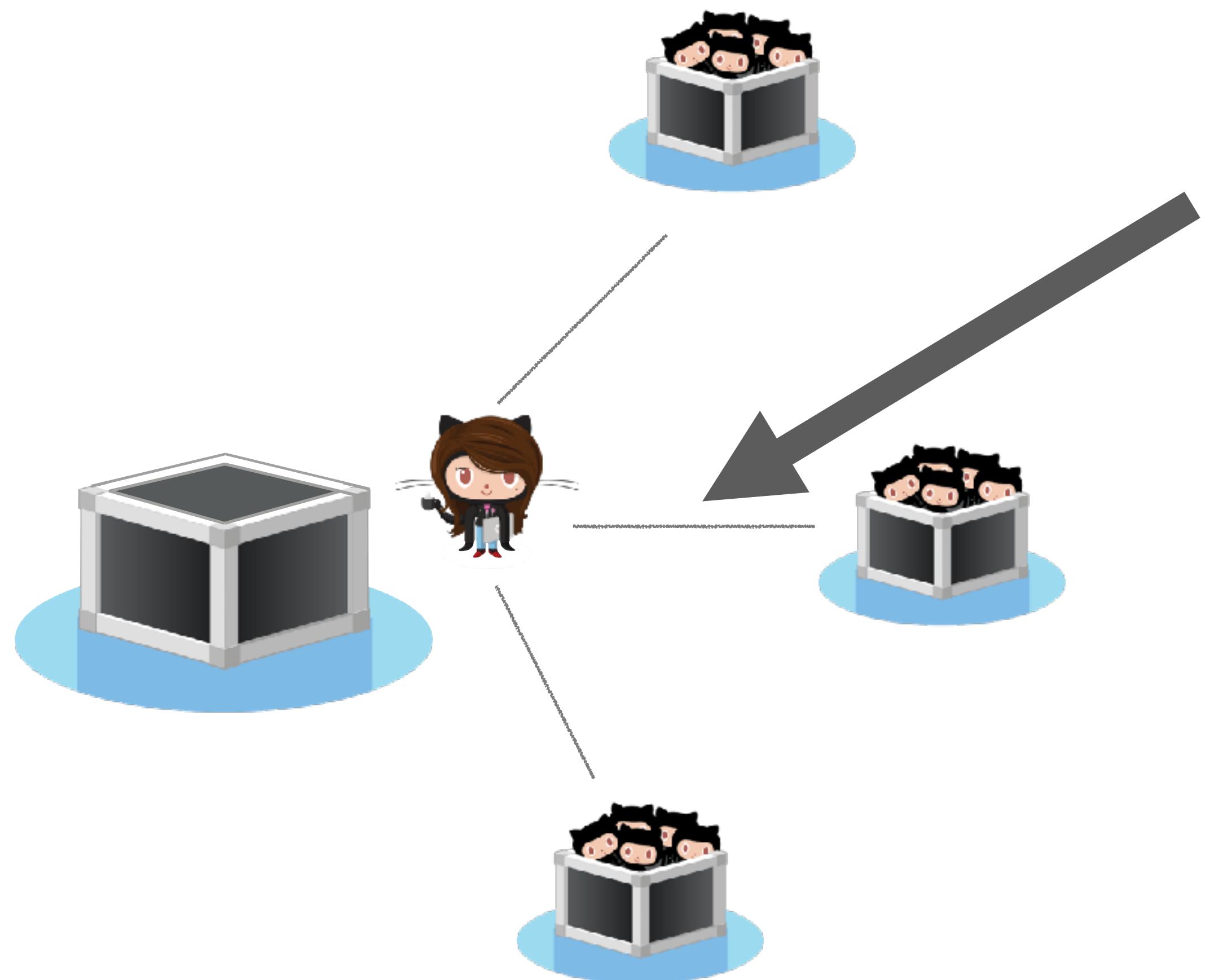
js/config/resolve.js   +1 -1  ⋮

... @@ -1,6 +1,6 @@

1 var path = require('path');	1 var path = require('path');
2	2
3 - var renderer = process.env.GEONOTEBOOK_MAP_RENDERE R    'geojs';	3 + var renderer = process.env.GEONOTEBOOK_MAP_RENDERE R    'ol';
4	4
5 module.exports = {	5 module.exports = {
6 alias: {	6 alias: {
....	

# Commits to the codebase (relatively deep understanding of the codebase)

# From interactions to ties of varying strength



commented on Dec 7, 2017

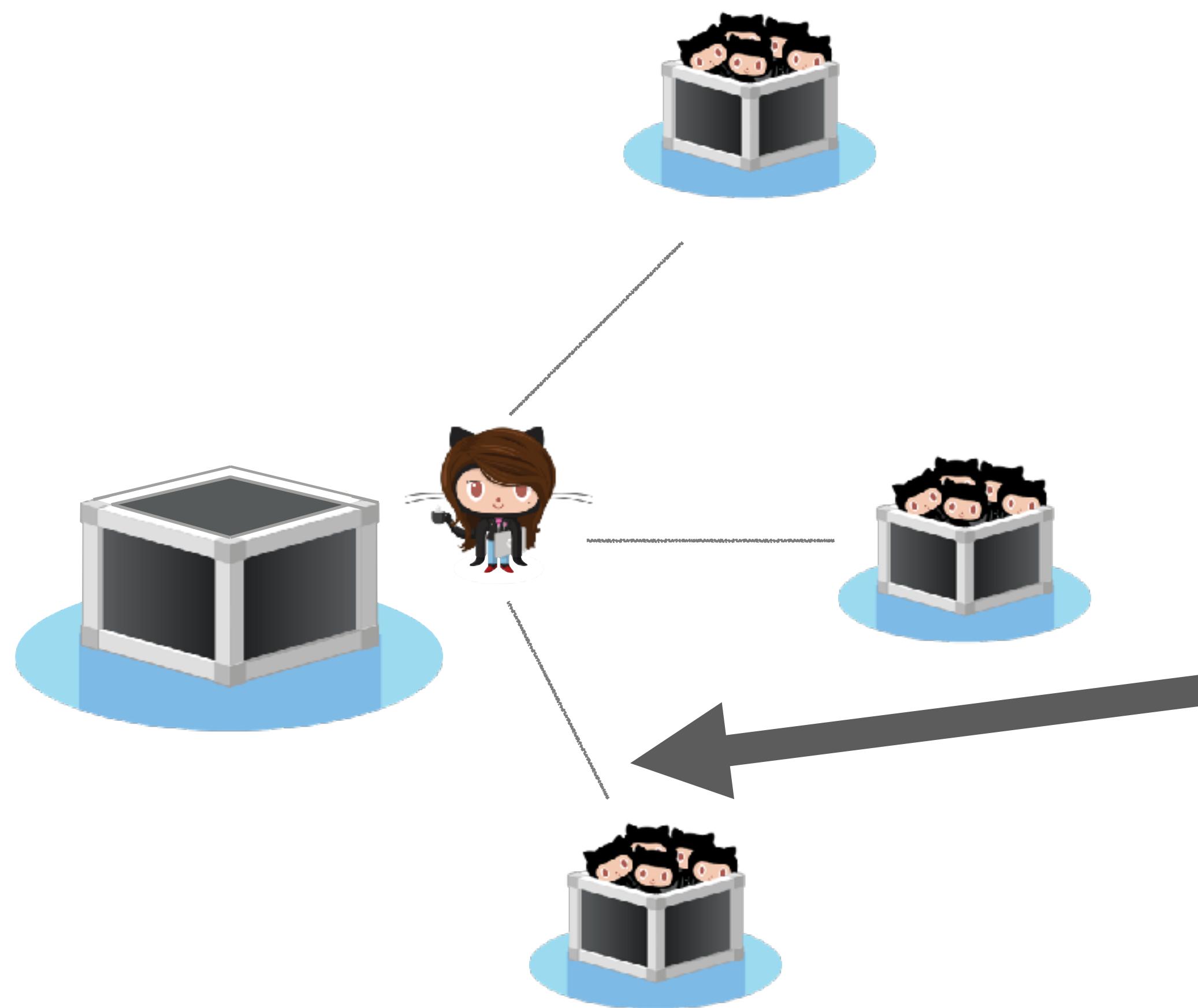
Hello,  
I am new to GeoNotebook, I am at the stage where I try to understand how GeoNotebook works, or more precisely what each of the python libraries that are used in GeoNotebook do.  
  
What I didn't understand is how I can change the projection of the rasters overplayed in Mapnik? What is the library that does this, is it Mapnik or Rasterio? For the vectors, is Shapely, if I am not mistaken.

Smiley face icon

Assignees	No one assigned
Labels	None yet
Projects	None yet
Milestone	No milestone
Development	No branches or pull requests

Issue reports  
(some understanding of the project)

# From interactions to ties of varying strength



Stars

Search stars  Search Type: All Language Sort by: Recently starred

[OpenGeoscience / geonotebook](#) Starred

A Jupyter notebook extension for geospatial visualization and analysis

Python 1,081 141 Updated on Jan 21, 2019

Stars  
(awareness of the project)

Many interactions are possible, these were just three examples.

Stars

Search stars  Search Type: All Language Sort by: Recently starred

[OpenGeoscience / geonotebook](#)  Starred 

A Jupyter notebook extension for geospatial visualization and analysis

 Python  1,081  141 Updated on Jan 21, 2019

	<p>[REDACTED] commented on Dec 7, 2017</p>	<p><b>Assignees</b></p> <p>No one assigned</p> <hr/> <p><b>Labels</b></p> <p>None yet</p> <hr/> <p><b>Projects</b></p> <p>None yet</p> <hr/> <p><b>Milestone</b></p> <p>No milestone</p> <hr/> <p><b>Development</b></p> <p>No branches or pull requests</p>
	<p>Hello,</p> <p>I am new to GeoNotebook, I am at the stage where I try to understand how GeoNotebook works, or more precisely what each of the python libraries that are used in GeoNotebook do.</p> <p>What I didn't understand is how I can change the projection of the rasters overplayed in Mapnik? What is the library that does this, is it Mapnik or Rasterio? For the vectors, is Shapely, if I am not mistaken.</p>	

1 file changed +1 -1 lines changed

js/config/resolve.js

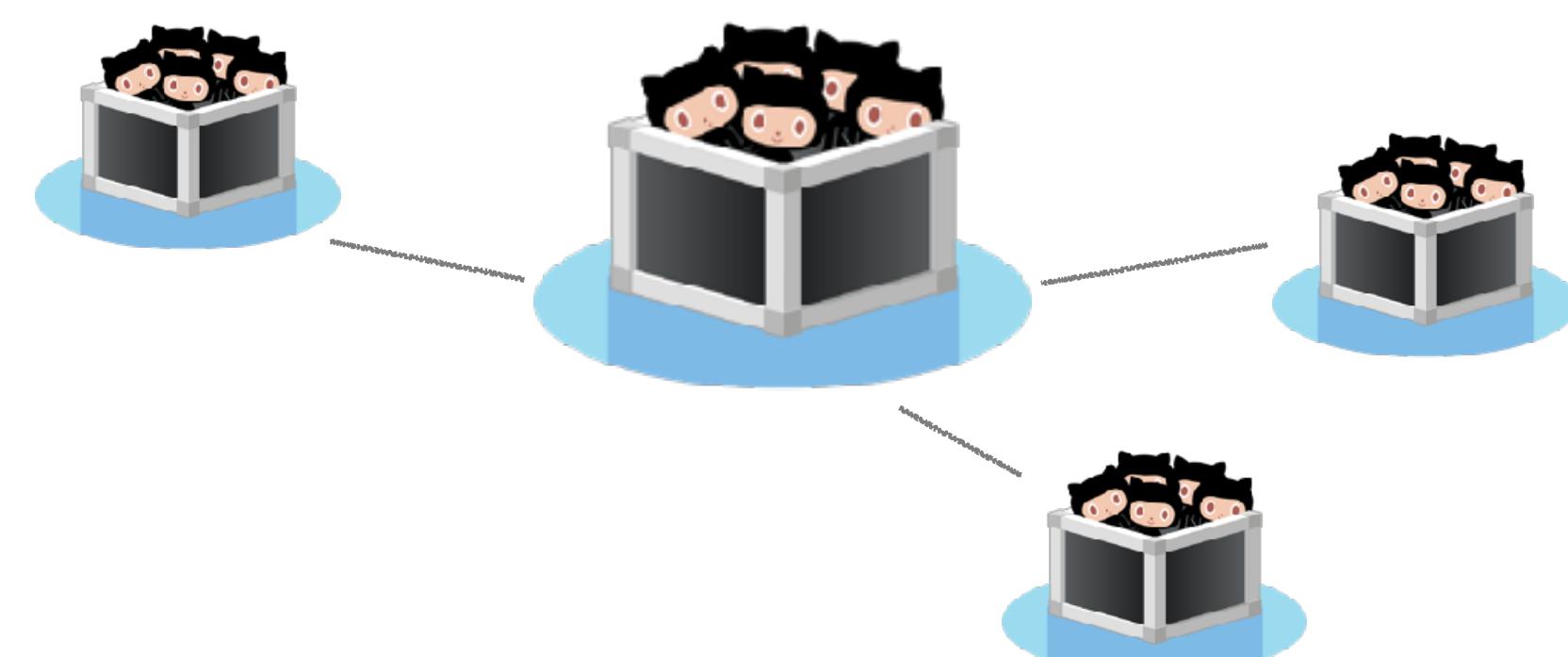
+1 -1

```
@@ -1,6 +1,6 @@
1 var path = require('path');
2
3 -var renderer =
4   process.env.GEONOTEBOOK_MAP_RENDERE
5   R || 'geojs';
6
7 module.exports = {
8   alias: {
9     ....
```

# Stars

## Issues

# Commits

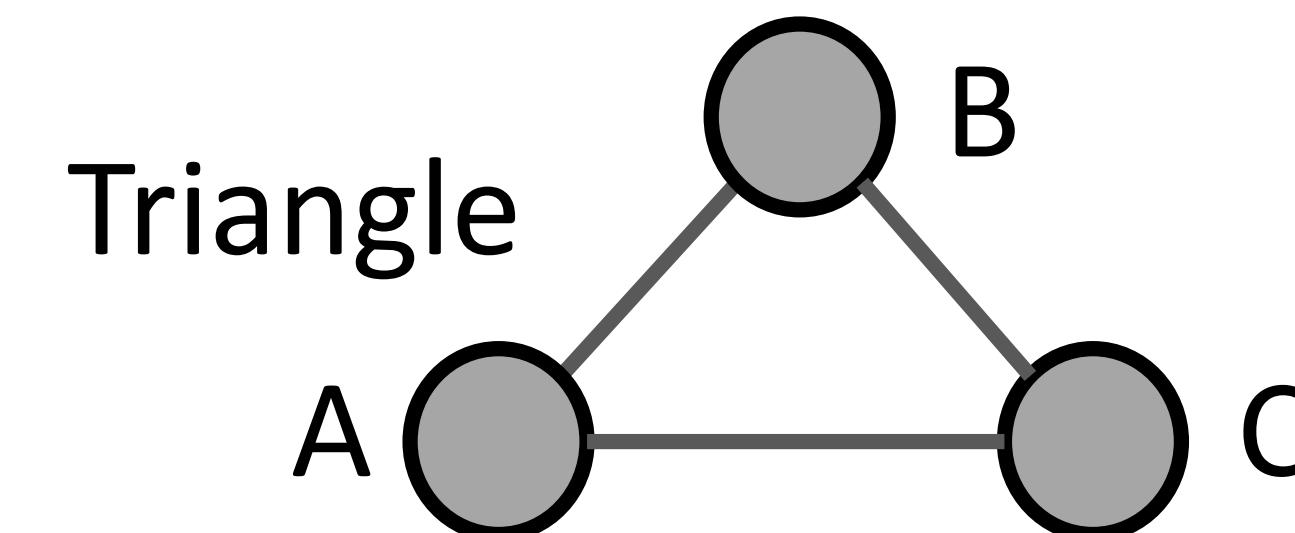
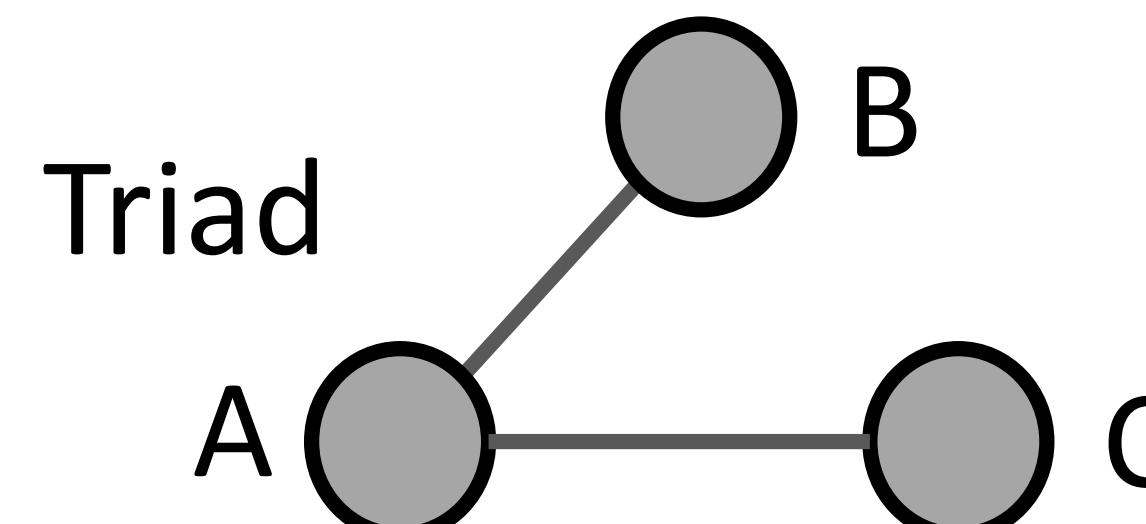


# Weaker ties

# Stronger ties

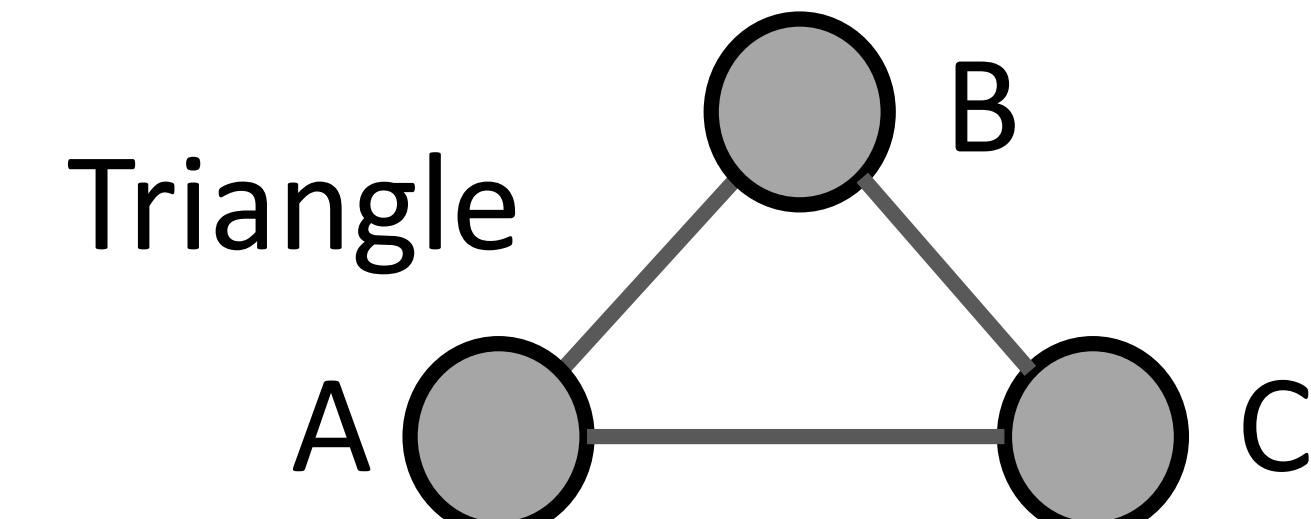
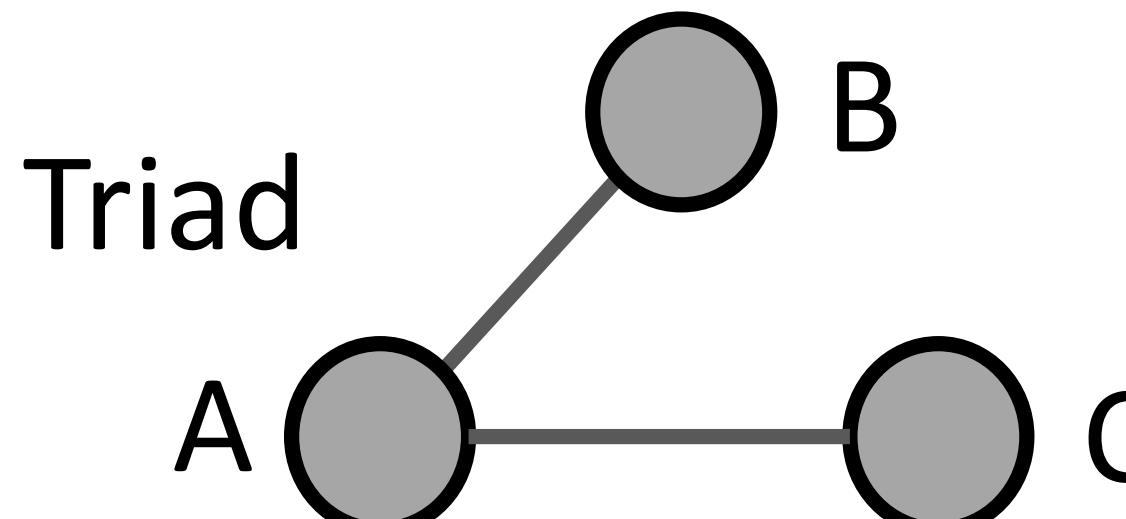
In strongly-tied social networks, triads are unlikely.

---



There is ~an order of magnitude (10 $\times$ ) difference in transitivity values between each pair of networks.

---

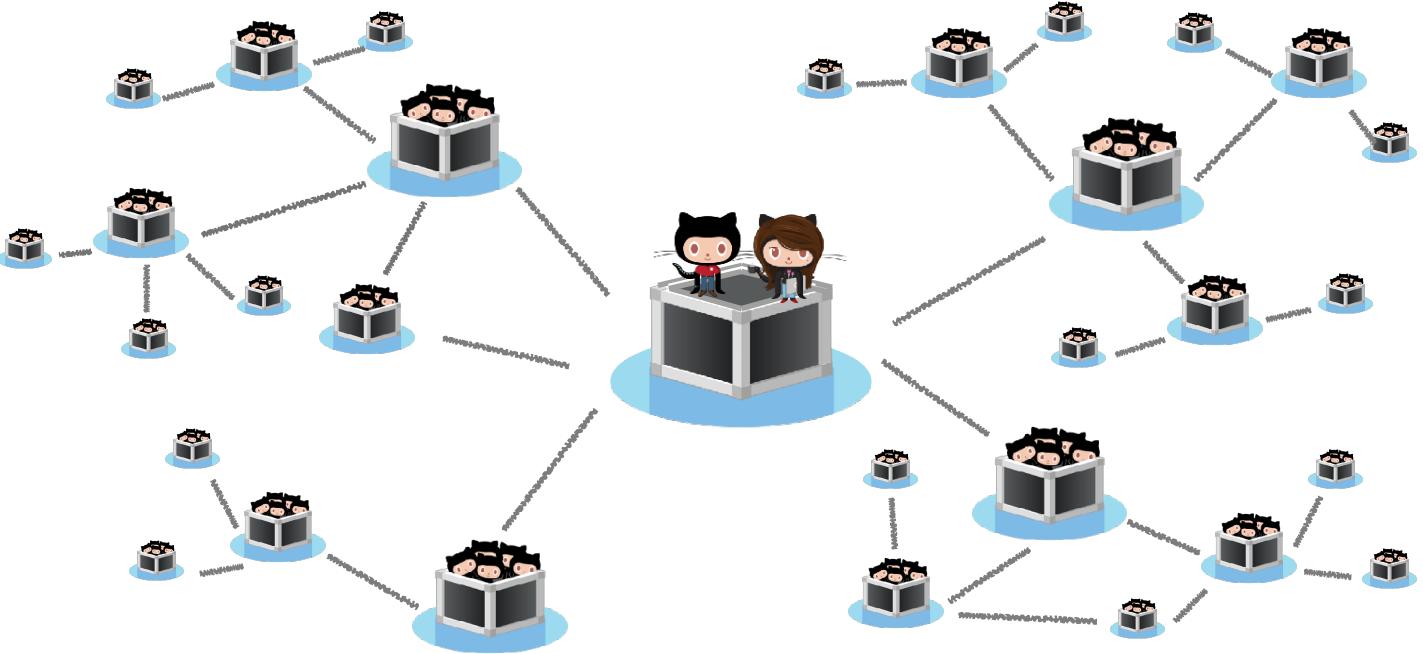


Interaction	#Nodes	#Edges	Transitivity ( $\times 10^{-2}$ )
Commits	763,062	1,926,978	30.04
Issues	278,945	727,255	3.42
Stars	480,394	3,658,543	0.23

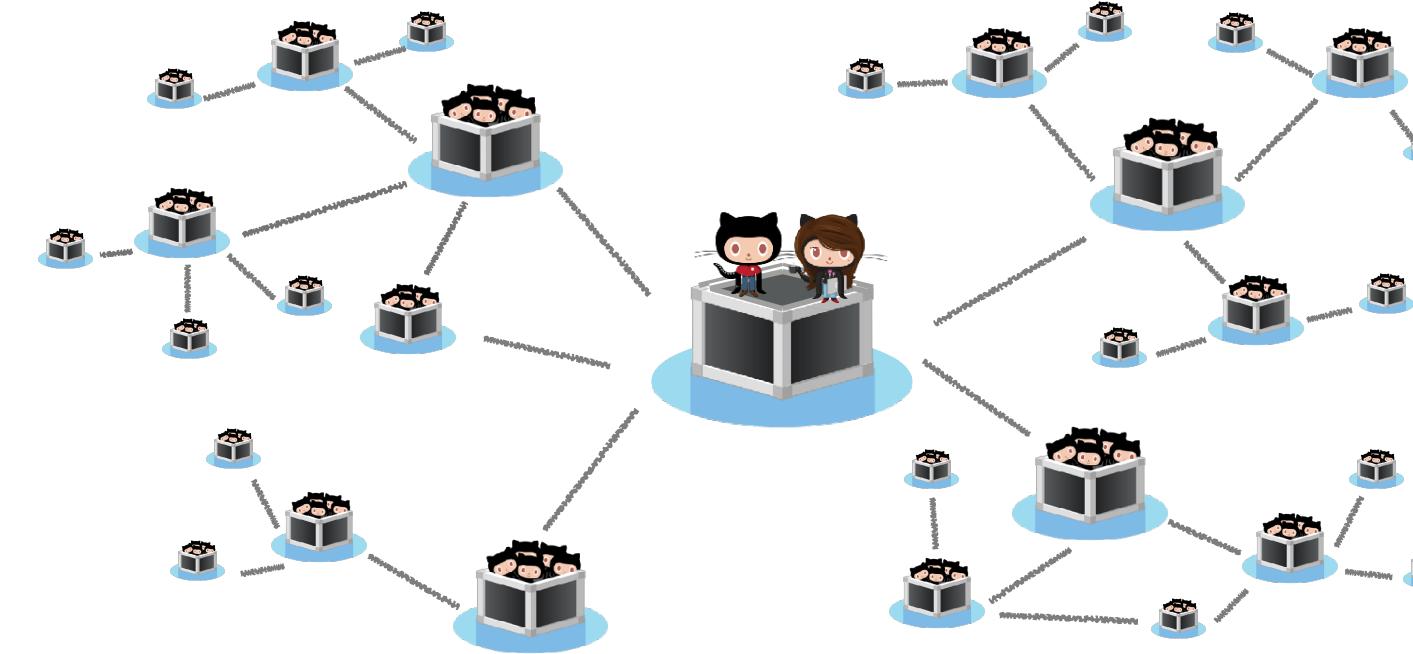
$$\text{Transitivity} = 3 * \frac{\text{N}_{\text{triangles}}}{\text{N}_{\text{triads}}}$$

Commits >> Issues >> Stars

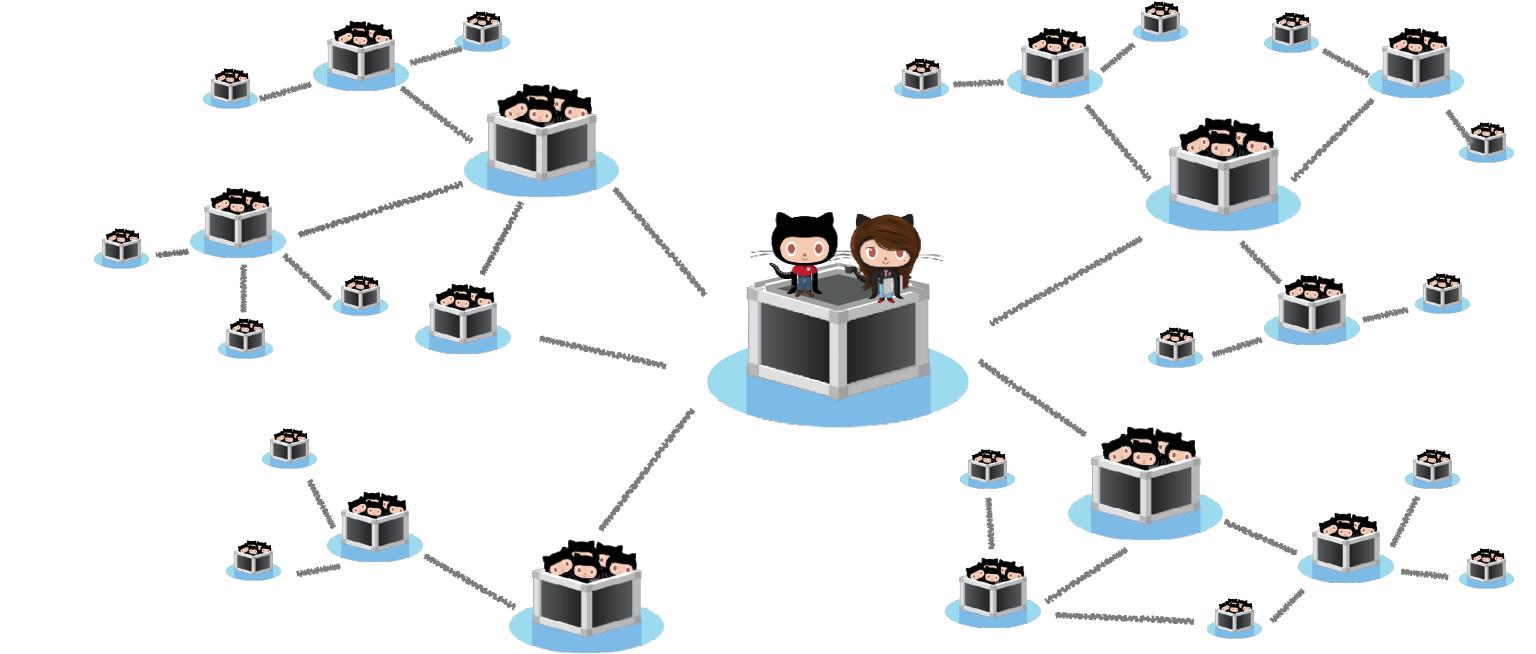
# Now what?



Stars

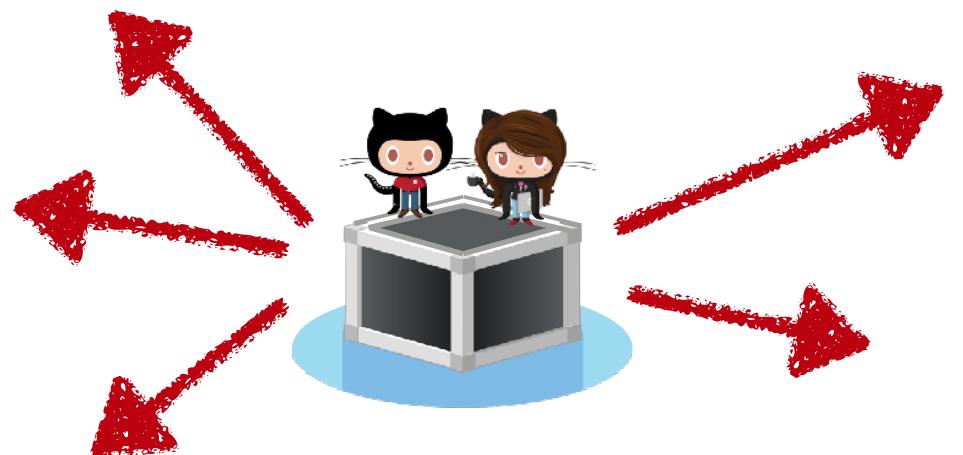


Issues

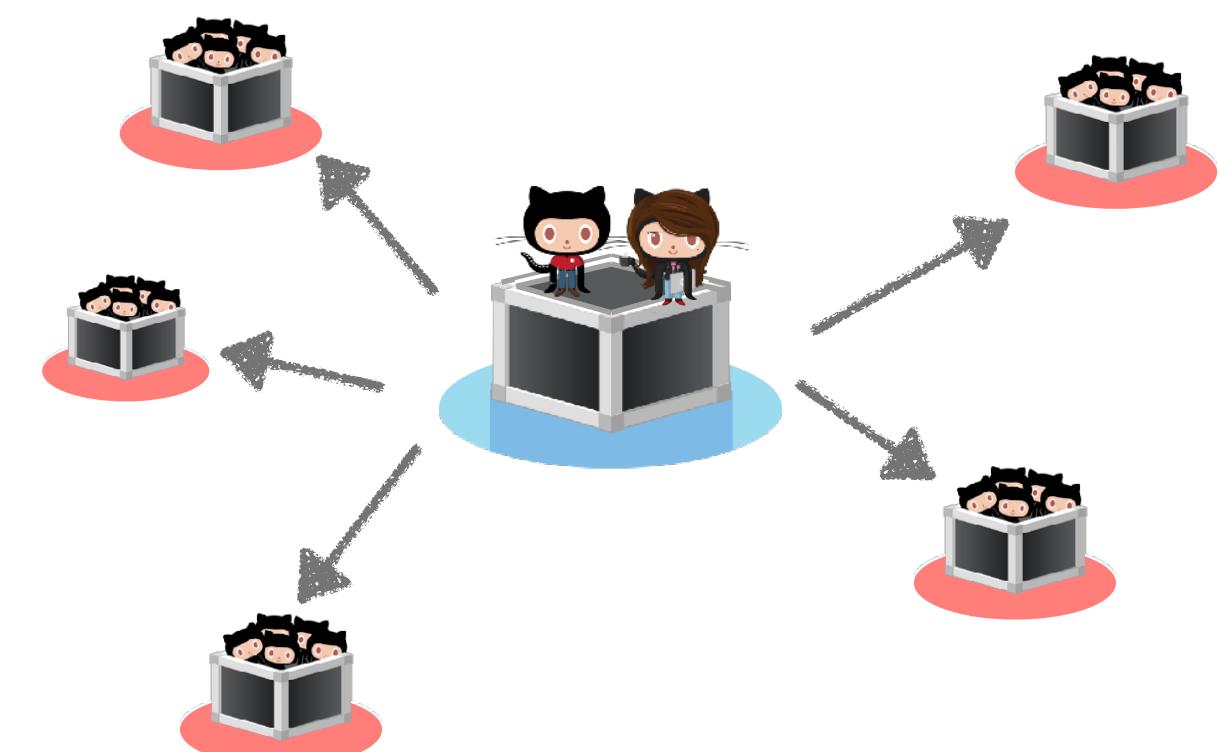


Commits

Out-degree centrality  $\times 3$ ?



Information diversity index  $\times 3$ ?



The first two PCs cumulatively explain over 80% of the variance.

---

	Out-deg. centrality			Diversity index		
	PC1	PC2	PC3	PC1	PC2	PC3
D <sub>commit</sub>	0.60	-0.45	0.67	0.63	-0.36	0.69
D <sub>issue</sub>	0.61	-0.28	-0.74	0.64	-0.24	-0.72
D <sub>star</sub>	0.52	0.85	0.11	0.43	0.90	0.08

PC1: Average volume of information available /  
Average diversity of the knowledge space (hyp 2)

The first two PCs cumulatively explain over 80% of the variance.

---

	Out-deg. centrality			Diversity index		
	PC1	PC2	PC3	PC1	PC2	PC3
D <sub>commit</sub>	0.60	-0.45	0.67	0.63	-0.36	0.69
D <sub>issue</sub>	0.61	-0.28	-0.74	0.64	-0.24	-0.72
D <sub>star</sub>	0.52	0.85	0.11	0.43	0.90	0.08

PC2: Where the connectivity / diversity comes from  
**(The strength of weak ties)**

Hypothesis 3: The more the informational diversity can be attributed to weak ties, the more innovative the projects are.

---

	Out-deg. centrality			Diversity index		
	PC1	PC2	PC3	PC1	PC2	PC3
D <sub>commit</sub>	0.60	-0.45	0.67	0.63	-0.36	0.69
D <sub>issue</sub>	0.61	-0.28	-0.74	0.64	-0.24	-0.72
D <sub>star</sub>	0.52	0.85	0.11	0.43	0.90	0.08

PC2: Where the connectivity / diversity comes from

(The strength of weak ties)

# Finally, the novelty regression:

- Hypothesis 1 (greater connectivity): weak/inconsistent effects
- Hypothesis 2 (greater info diversity): small but clear effects (25–75 percentile: 4% change in the distribution)
- Hypothesis 3 (strength of weak ties): clear effects, comparable size

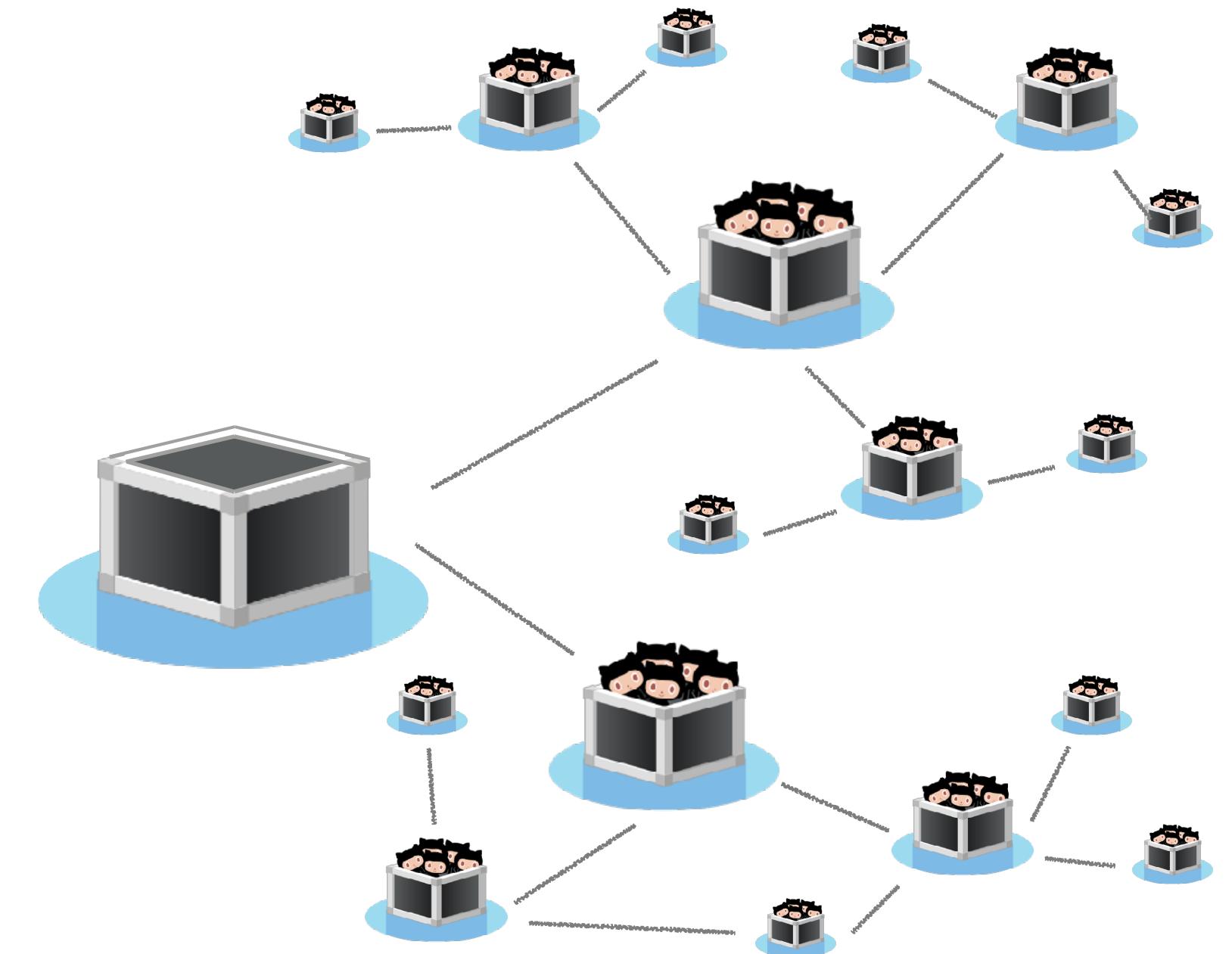
	Model III	Model IV
<b>Variables of interest</b>		
$Deg_{ave}$ ( $H_1$ )		-0.002*** (0.001)
$Deg_{weakness}$		-0.005*** (0.001)
$Div_{ave}$ ( $H_2$ )	0.007*** (0.001)	0.008*** (0.001)
$Div_{weakness}$ ( $H_3$ )	0.005*** (0.001)	0.007*** (0.001)
Observations	38,164	38,164

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

# Exposure to diverse ideas through weak ties predicts novel combinations of packages.

---

- Lurking on the GitHub platform seems to have quantifiable benefits. Redesign the Trending page?
- Automated project recommendation tools may be counterproductive?
- Well-informed but not necessarily highly active developers may also be experts at their craft?
- How to track and give credit to ideas?
- Surface-level vs deep-level diversity?
- AI-generated code: novel or regression to the mean?



# Today: Let's look at some concrete examples of network effects

---



Measuring innovation  
in software



Understanding how  
innovation emerges



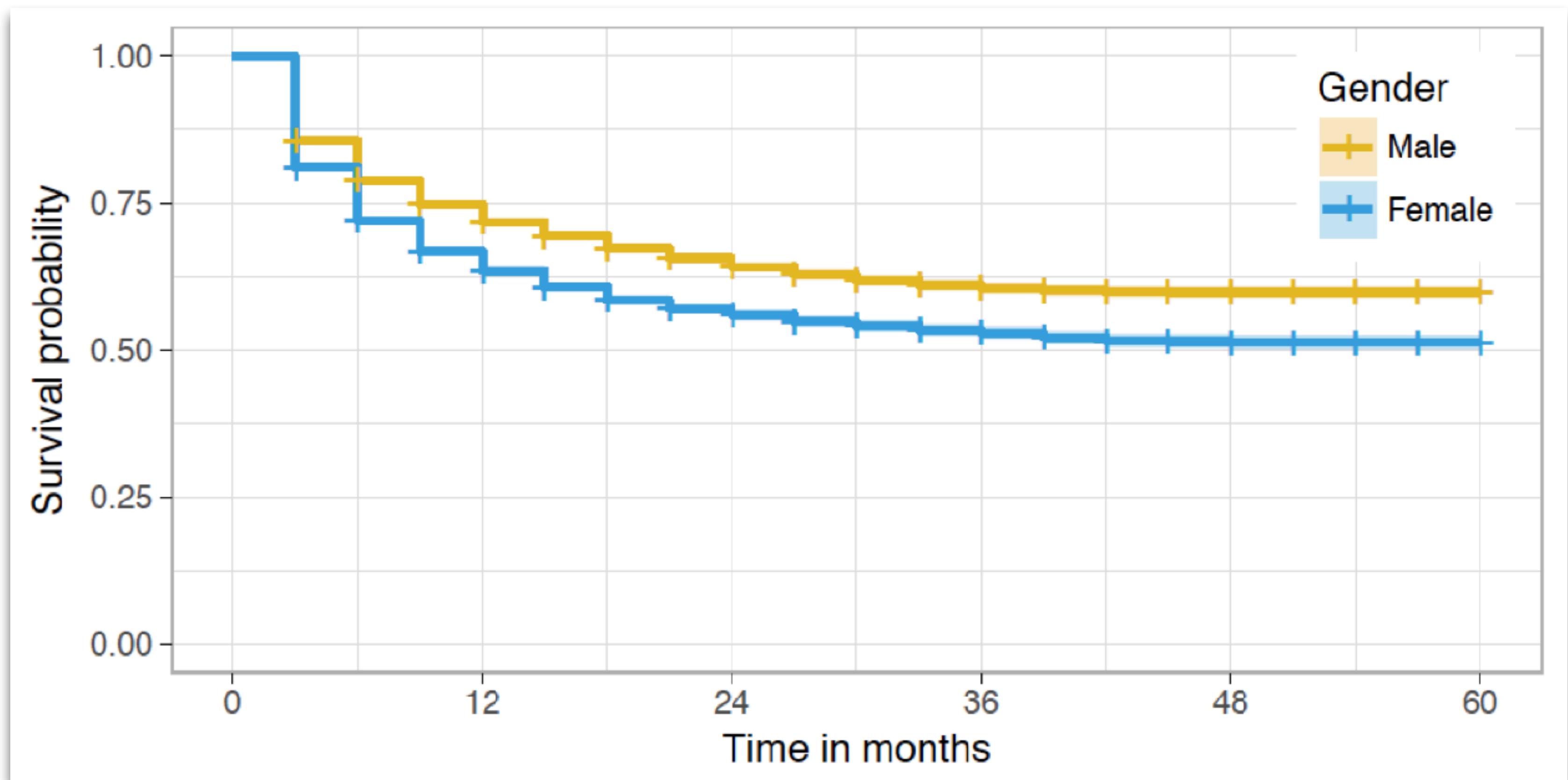
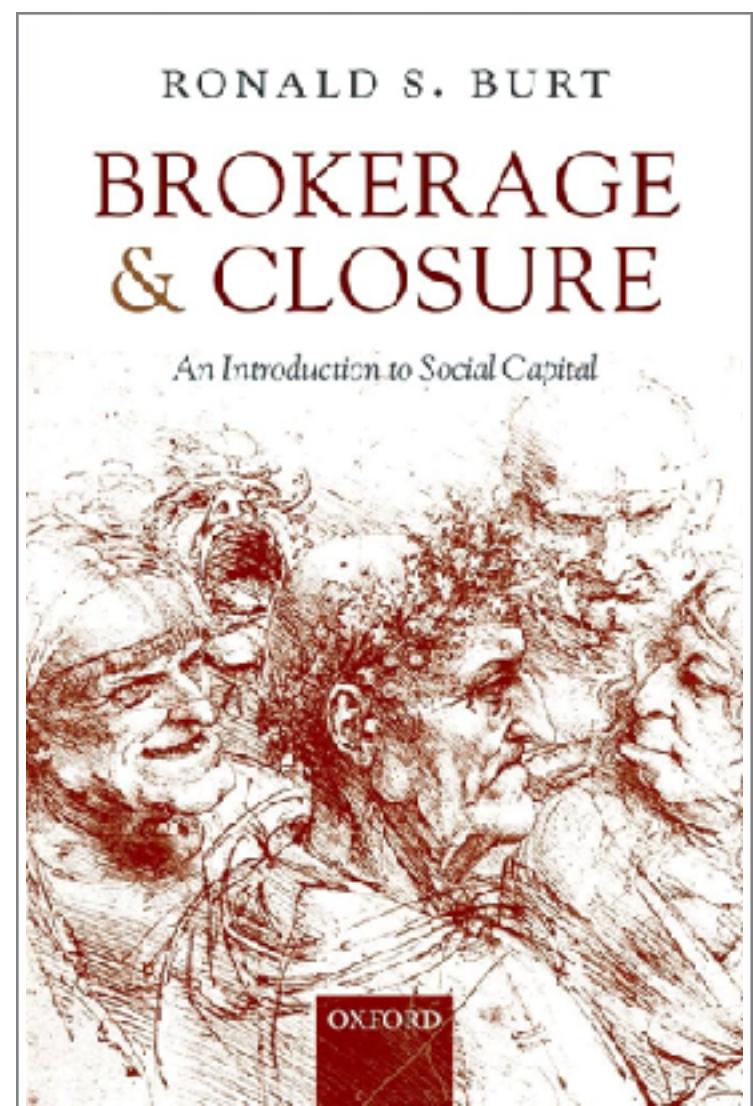
Social capital



Social contagion

# Weak ties predict longer-term participation for women

Social capital mechanism



# Weak ties predict the spread of tools

Diffusion of innovations mechanism



## 12 popular quality assurance tools

Continuous integration

Dependency management

build passing

dependencies up to date

Travis  
Circle  
Appveyor  
Codeship

Code coverage reporters

coverage 94%

Coveralls  
Codeclimate  
Codecov  
Codacy

Cross browser testers

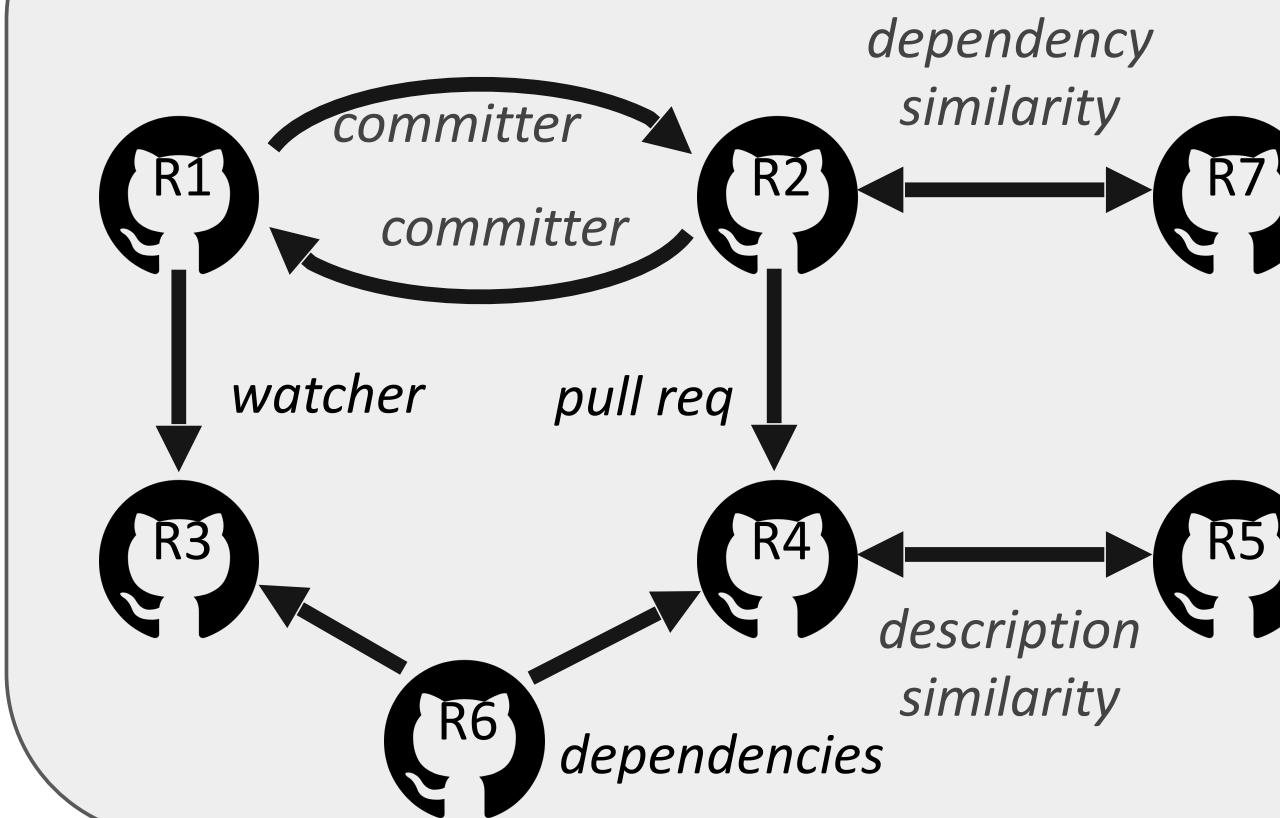
Firefox 82  
Chrome 86  
Windows 7  
Mac OS X 10.14  
Linux 5.3

~86,000 npm package repositories

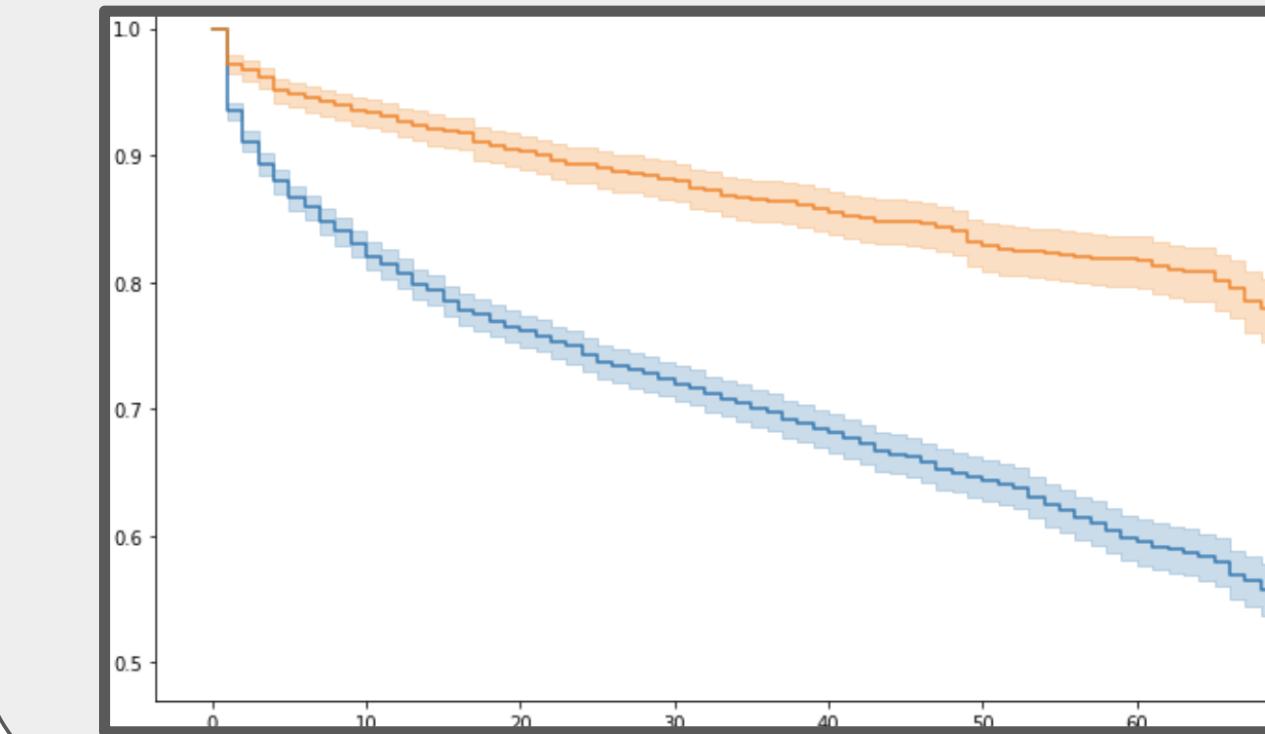


For each tool:

## Heterogeneous network



## Hazard modeling (Cox regression)



# Acknowledgements

---



Courtney Miller



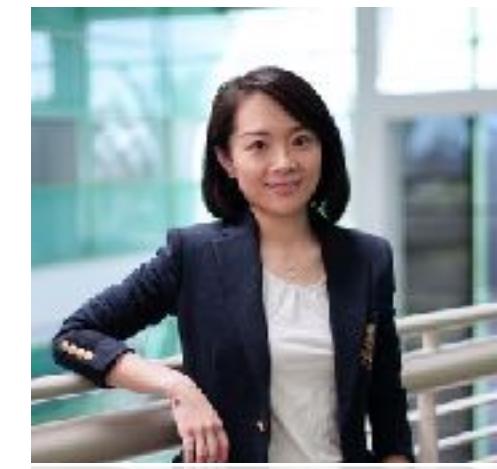
Anita Brown



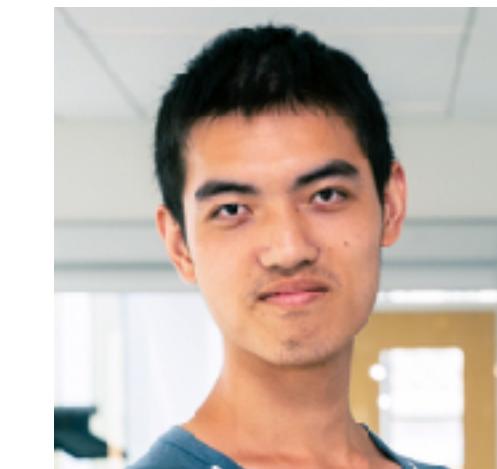
Asher Trockman



Jim Herbsleb



Shurui Zhou



Hongbo Fang



Anita Sarma



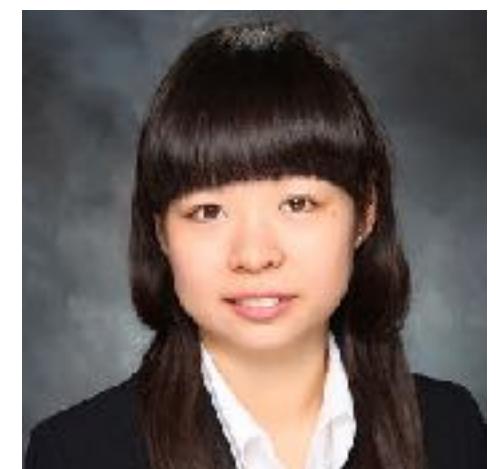
Cassandra Overney



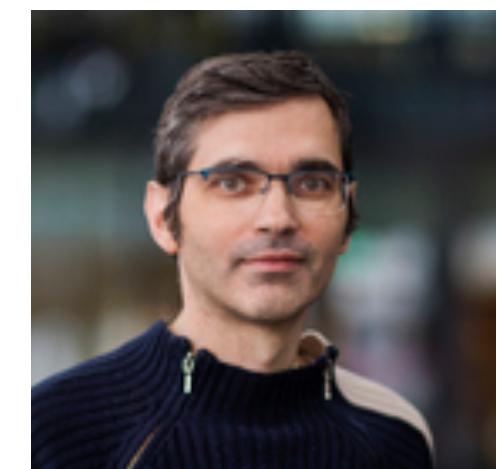
Audris Mockus



Alex Nolte



Sophie Qiu



Alex Serebrenik



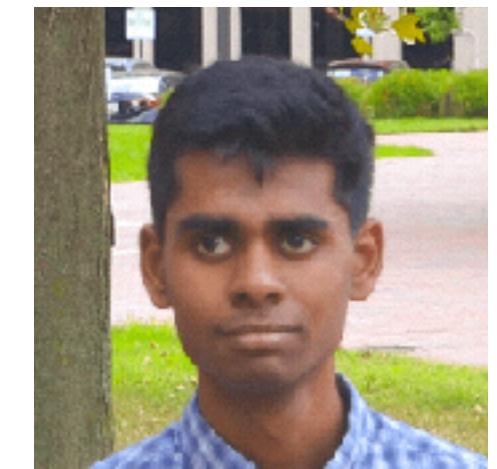
Marat Valiev



Laura Dabbish



Lily Li



Naveen Raman



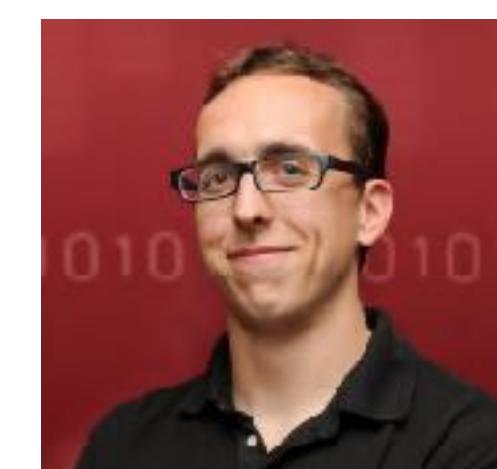
Hao He



Christian Kästner



Hemank Lamba



Emerson  
Murphy-Hill



Alfred P. Sloan  
FOUNDATION



FORDFOUNDATION

# STRUDEL sustainability research on ...

## Project practices

- [CHASE 2023](#) (social media)
- [ICSE 2020](#) (forking)
- [ESEC/FSE 2019](#) (forking)
- [ESEC/FSE 2018](#) (abandonment factors)

## Funding models

- [ICSE 2020](#) (donations)

## Sunsetting

- [ESEC/FSE 2023](#) (dealing with abandonment)

## Attracting contributors

- [ICSE 2022](#) (Twitter)
- [MSR 2020](#) (Twitter)
- [CSCW 2019](#) (signals)
- [ESEC/FSE 2015](#) (social connections)

## Transparency and signaling

- [ESEC/FSE 2020](#) (diffusion of practices)
- [CSCW 2019](#) (signals)
- [ICSE 2018](#) (badges)

## Stress, burnout, disengagement

- [ICSE 2022](#) (toxicity theory)
- [ICSE SEIS 2022](#) (toxicity vs pushback)
- [ICSE NIER 2020](#) (toxic language)
- [ICSE 2019](#) (overwork)
- [OSS 2019](#) (dropout, survival analysis)

## Diversity and inclusion

- [CHI 2023](#) (ClimateCoach)
- [ICSE SEIS 2023](#) (census)
- [ICSE 2019](#) (social capital)
- [CHI 2015](#) (gender & tenure)
- [CHASE 2015](#) (survey)

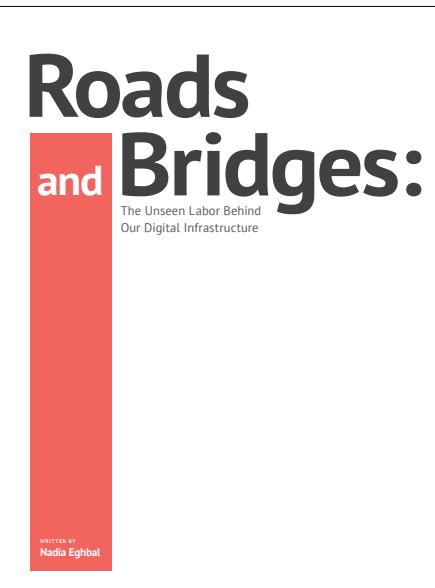
## Novelty and innovation

- [ICSE 2024](#) (atypical combinations)

## Network effects

- [ICSE 2024](#) (innovation)
- [ESEC/FSE 2023](#) (labor pools)
- [ICSE 2022](#) (Twitter)
- [ESEC/FSE 2020](#) (diffusion of practices)
- [ICSE 2019](#) (social capital)
- [ESEC/FSE 2018](#) (abandonment factors)

Like any infrastructure, it needs regular upkeep and maintenance

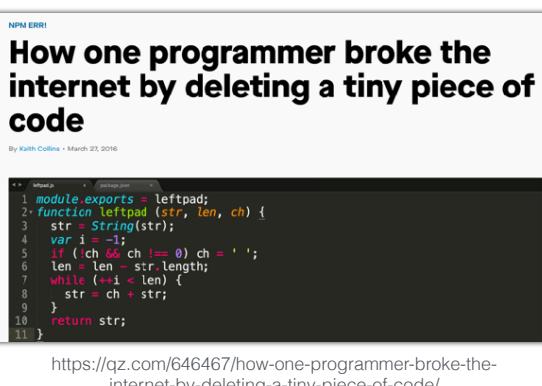


Everybody uses open source:

- Fortune 500 companies
- Major software companies
- Startups
- Government
- ...

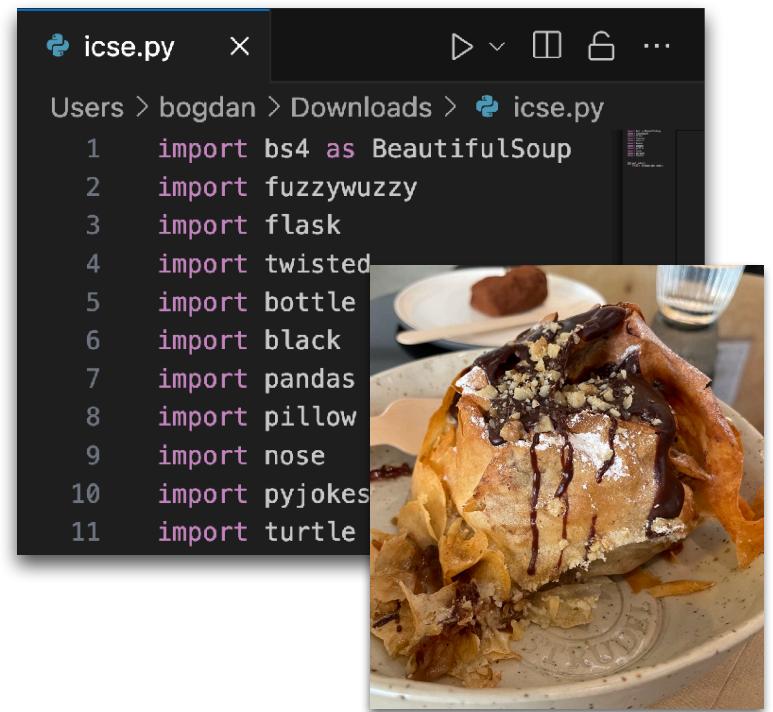
If undermaintained:

- Brittle supply chains
- Risks for downstream users
- Slows down innovation
- ...



3

Software innovation as novel recombination of software libraries



Lots of combinations:

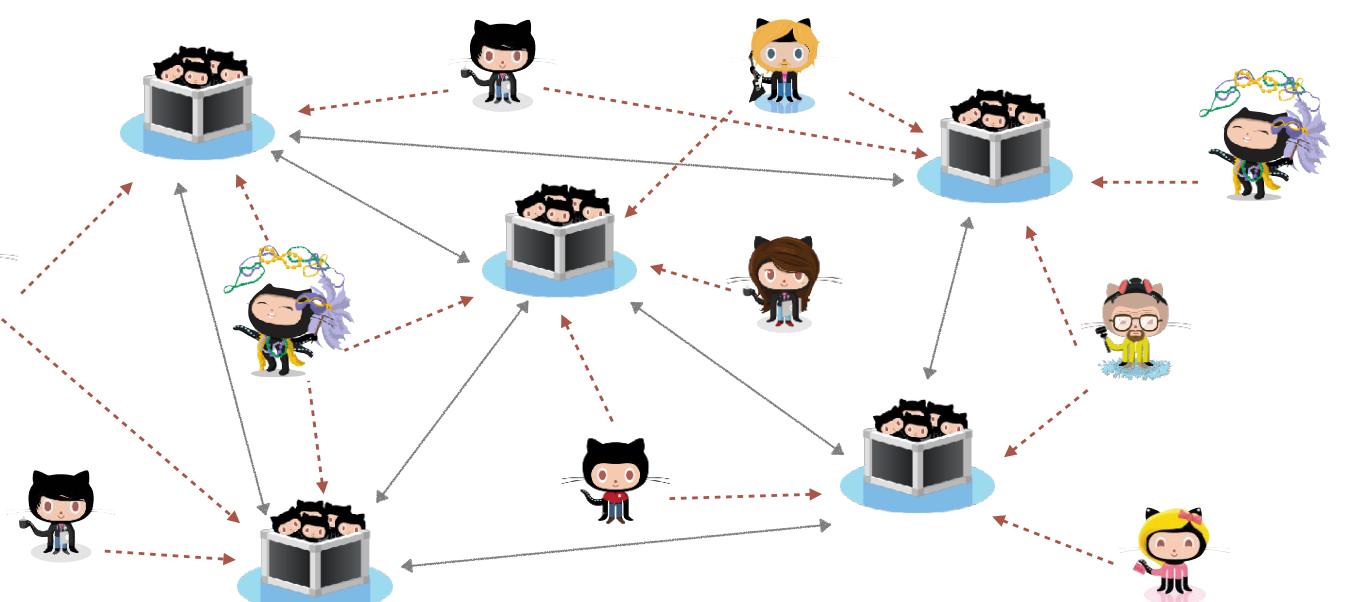
- (twisted, bottle)
- (turtle, nose)
- (black, pandas)
- (fuzzywuzzy, pillow)
- ...

C( $n, 2$ ) unique pairs of packages.

Dark chocolate + apple strudel is arguably innovative because it is atypical.

24

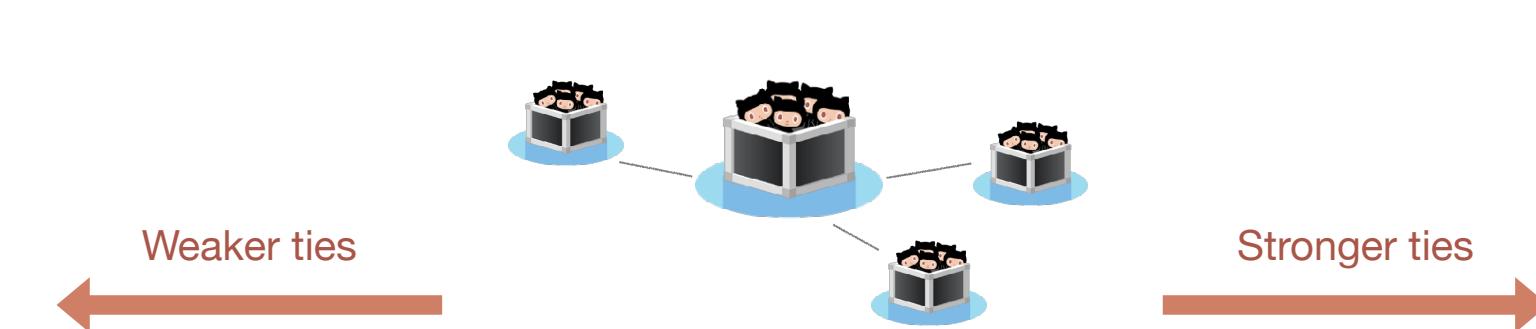
Contributors and projects form complex socio-technical networks!



Can we measure the network effects?

11

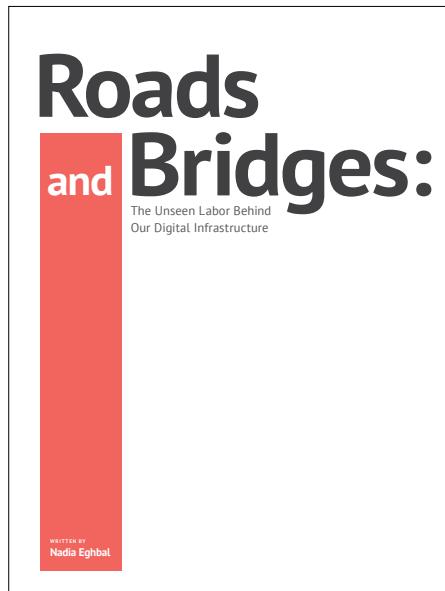
Exposure to diverse ideas through weak ties predicts novel combinations of packages.



# The Strength of Weak Ties in Open-Source Software Development Networks

Bogdan Vasilescu  
At UC Irvine, October 3rd, 2024

Like any infrastructure, it needs regular upkeep and maintenance



**Everybody** uses open source:

- Fortune 500 companies
- Major software companies
- Startups
- Government
- ...

If undermaintained:

- Brittle supply chains
- Risks for downstream users
- Slows down innovation
- ...

**How one programmer broke the internet by deleting a tiny piece of code**

```
1 module exports = leftpad;
2 function leftpad(str, len, ch) {
3     var i, l;
4     if (l <= 0) return str;
5     len = len - str.length;
6     while (len < 0) {
7         str = ch + str;
8         len++;
9     }
10    return str;
11 }
```

<https://qz.com/64647/how-one-programmer-broke-the-internet-by-deleting-a-tiny-piece-of-code/>



Software innovation as novel recombination of software libraries

```
icse.py
Users > bogdan > Downloads > icse.py
1 import bs4 as BeautifulSoup
2 import fuzzywuzzy
3 import flask
4 import twisted
5 import bottle
6 import black
7 import pandas
8 import pillow
9 import nose
10 import pyjokes
11 import turtle
```

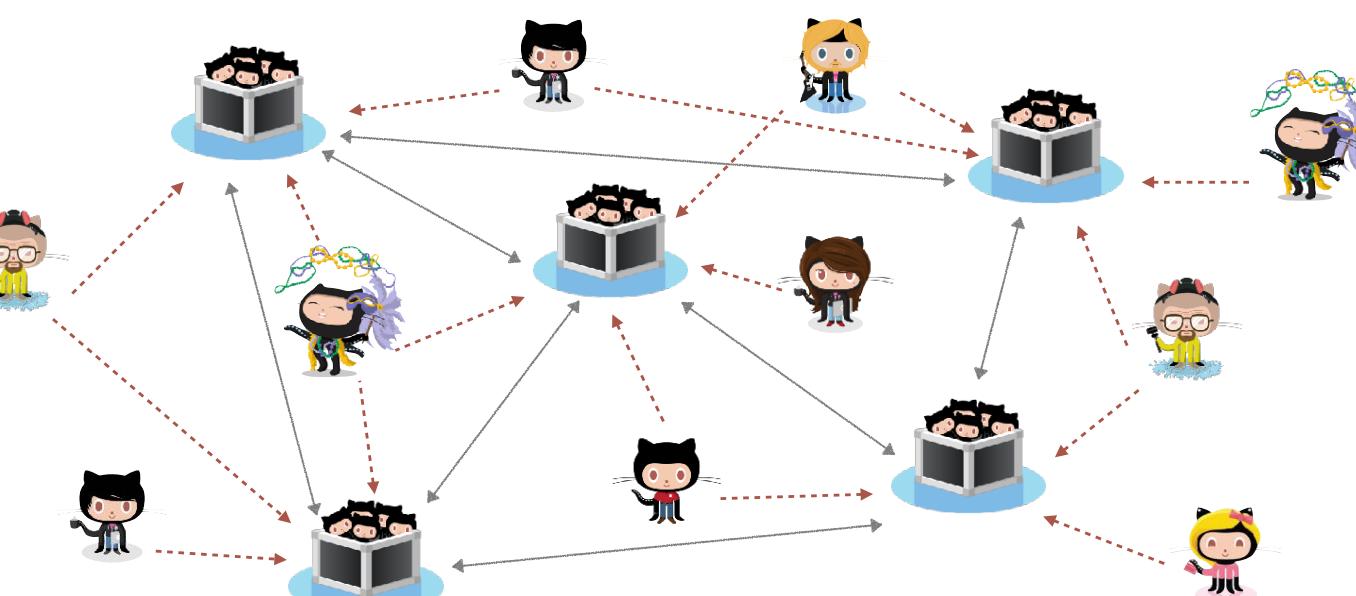
Lots of combinations:

- (twisted, bottle)
- (turtle, nose)
- (black, pandas)
- (fuzzywuzzy, pillow)
- ...

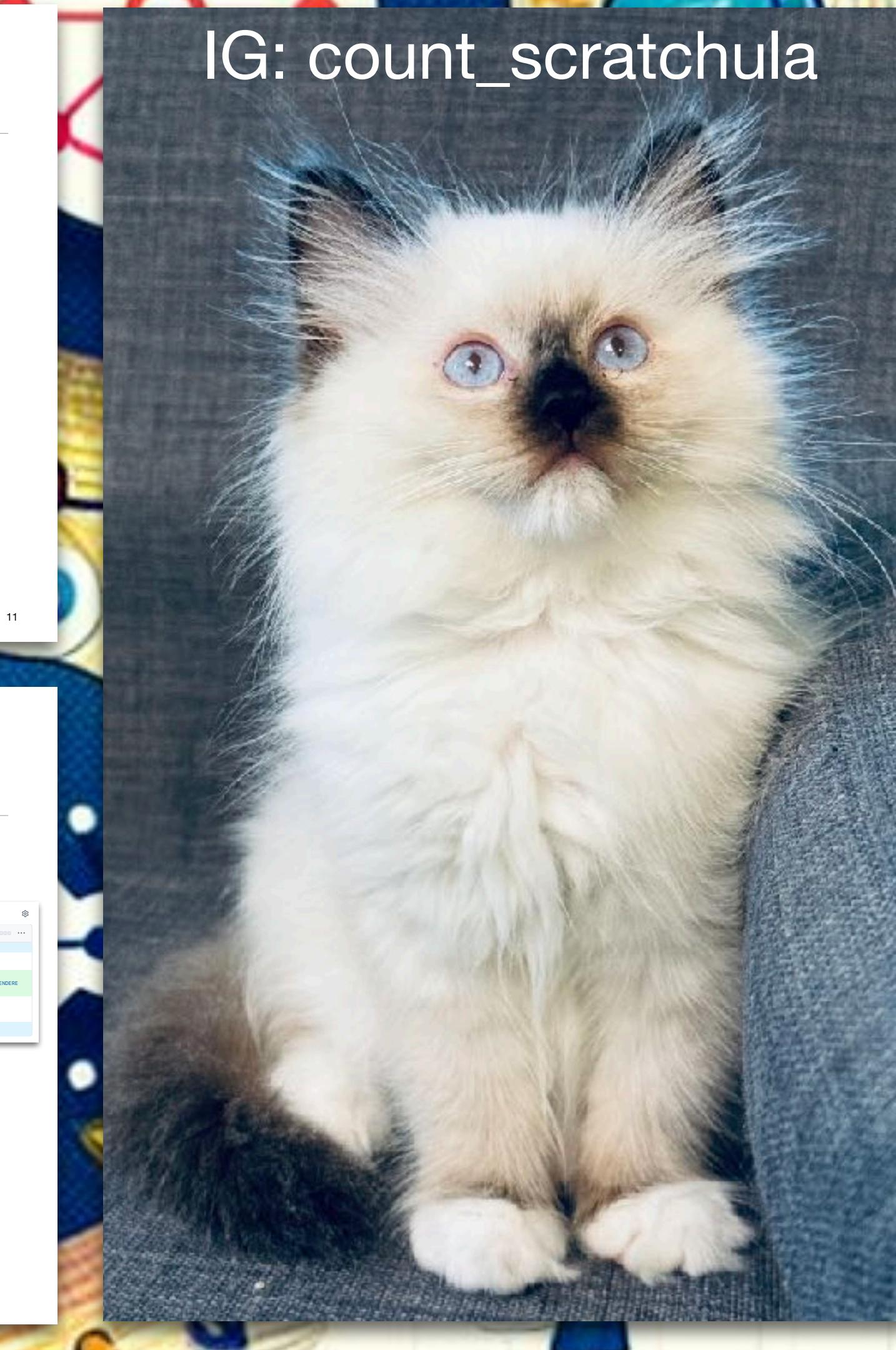
$C(n,2)$  unique pairs of packages.

Dark chocolate + apple strudel is arguably innovative because it is atypical.

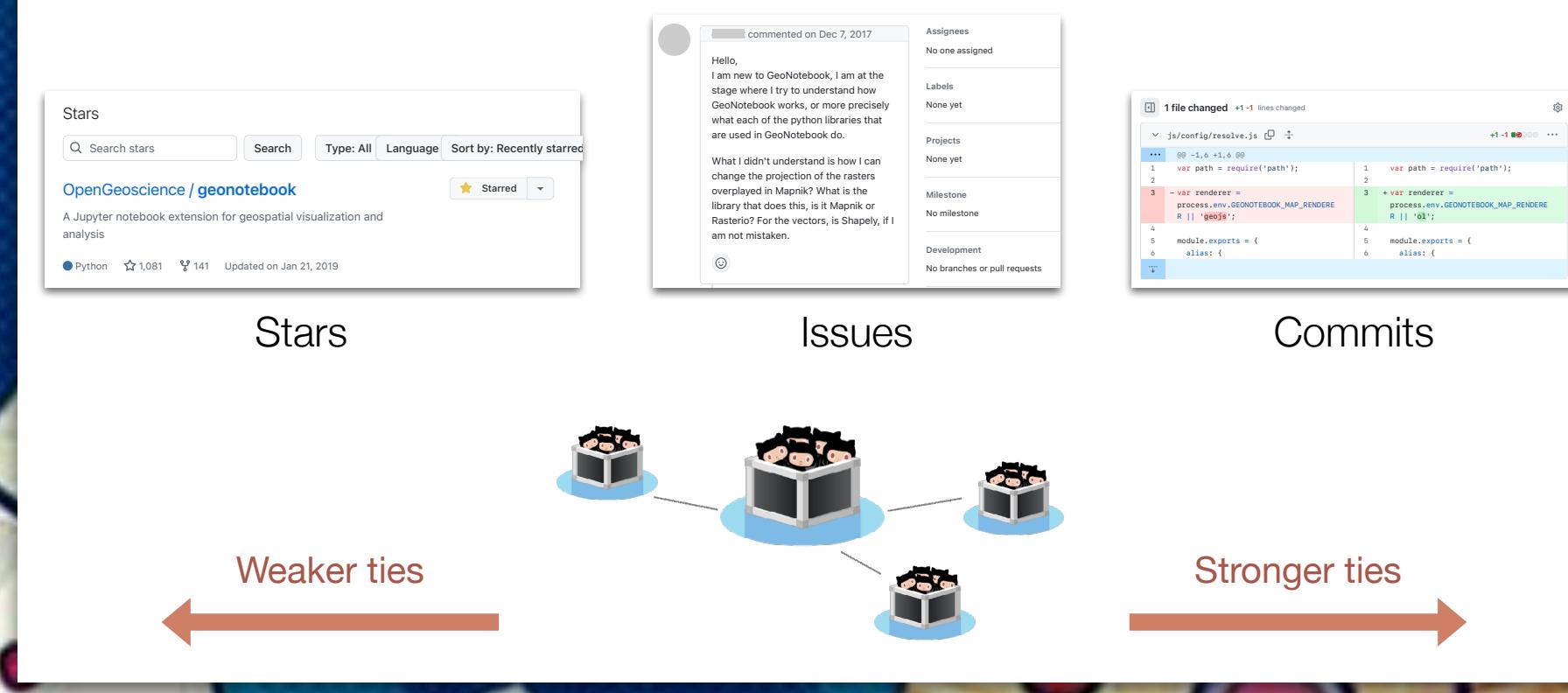
Contributors and projects form complex socio-technical networks!



Can we measure the network effects?



Exposure to diverse ideas through weak ties predicts novel combinations of packages.



# The Strength of Weak Ties in Open-Source Software Development Networks

Bogdan Vasilescu  
At UC Irvine, October 3rd, 2024