# From OpenRank to OpenPerf

## —— Enhancing Open Source Ecosystem Insights with Graph-Based Approaches

**Wei Wang**

**East China Normal University**
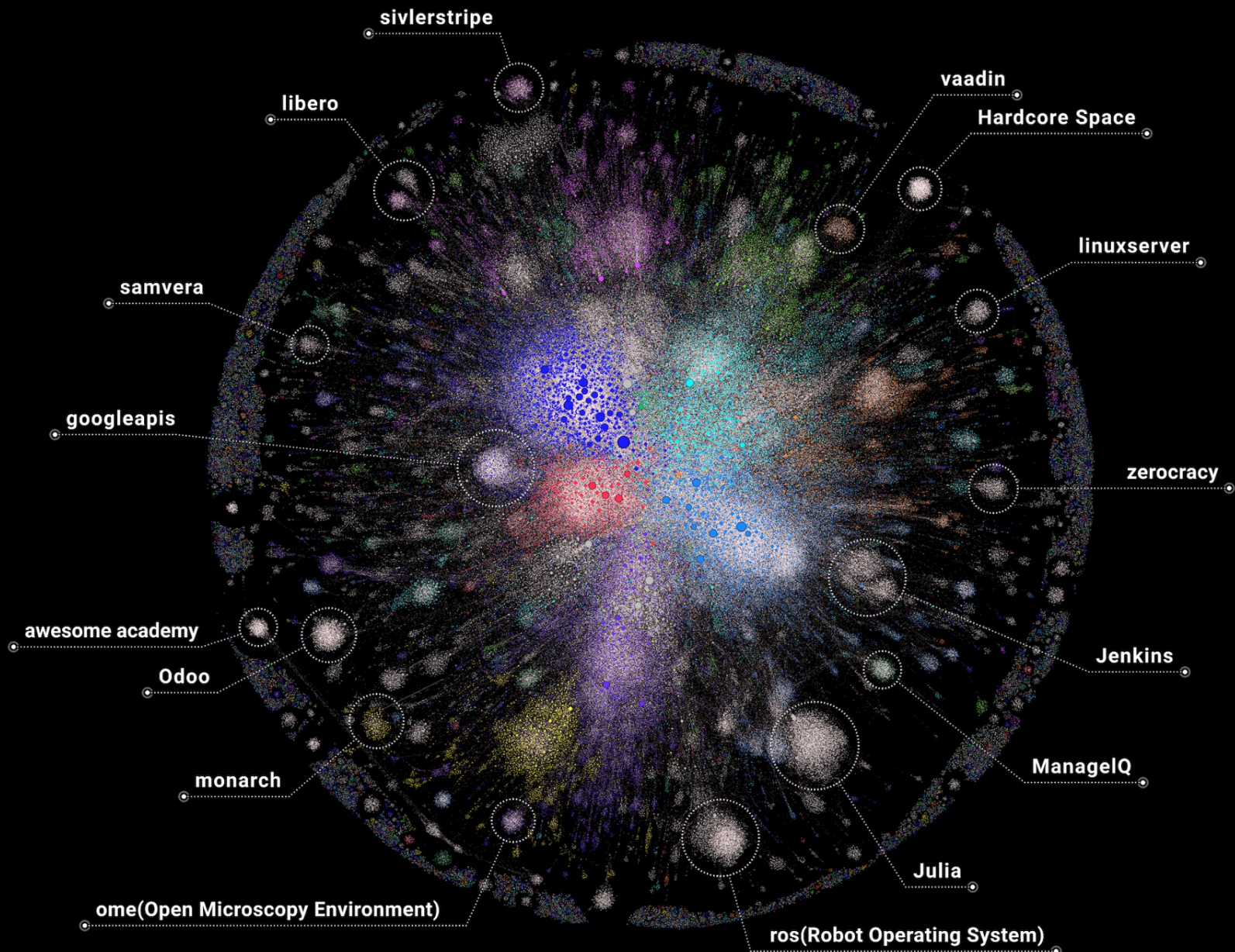
**X-lab Community**

**July 2024**

# OpenGalaxy 2019

OpenGalaxy is generated by collaboration network of all active GitHub repos in 2019. This graph contains 171,141 nodes and 2,811,489 edges. The generate method can be found in here [1] and the data is from GHArchive [2].

OpenGalaxy 是通过 GitHub 2019 年全域所有活跃项目的协作网络生成的。本图共包含 171,141 个节点和 2,811,489 条边。具体生成方法请参见这里 [1]，数据来自于 GHArchive [2]。
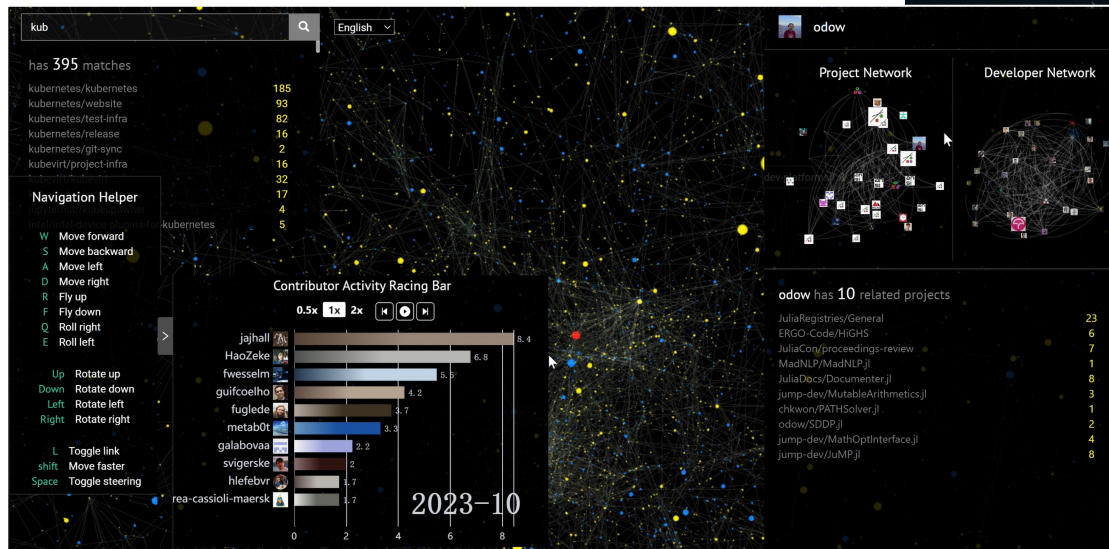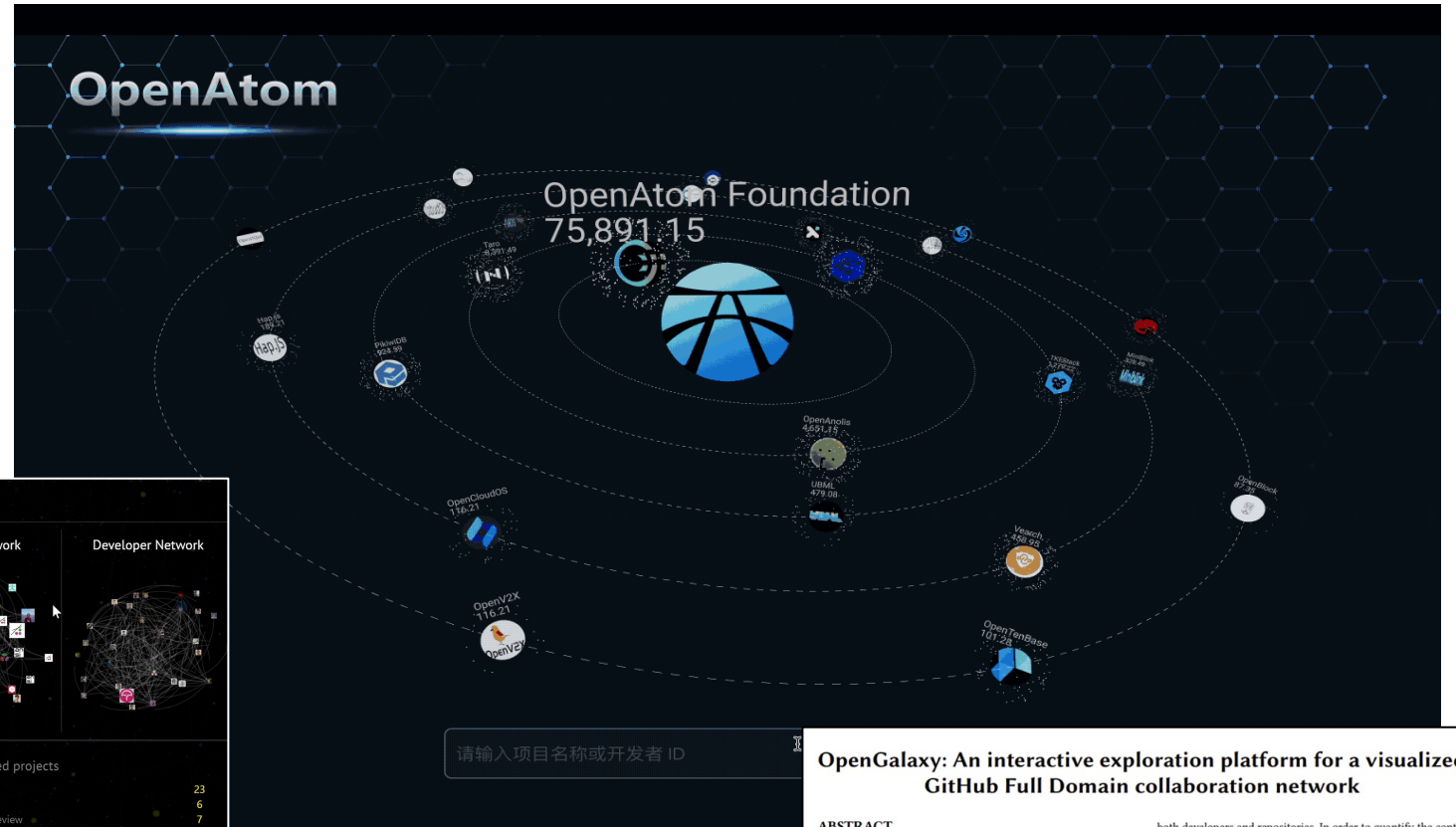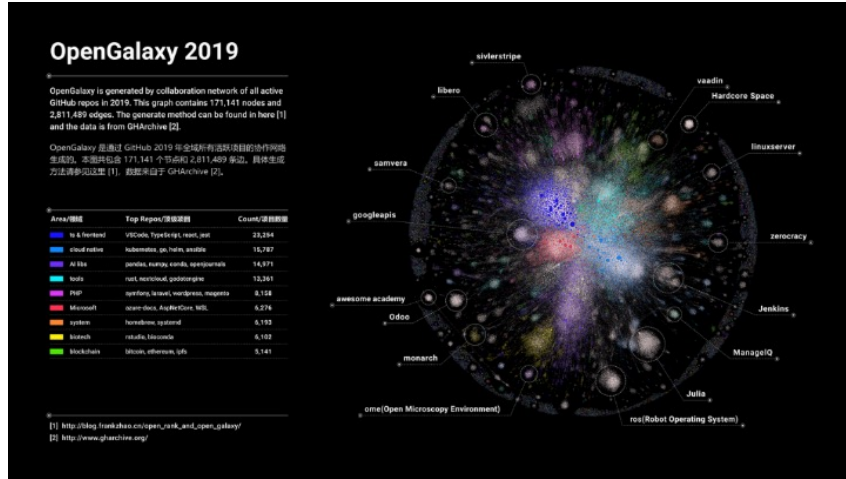
| Area/领域 | Top Repos/顶级项目 | Count/项目数量 |
|---|---|---|
| ts & frontend | VSCode, TypeScript, react, jest | 23,254 |
| cloud native | kubernetes, go, helm, ansible | 15,787 |
| AI libs | pandas, numpy, conda, openjournals | 14,971 |
| tools | rust, nextcloud, godotengine | 13,361 |
| PHP | symfony, laravel, wordpress, magento | 8,158 |
| Microsoft | azure-docs, AspNetCore, WSL | 6,276 |
| system | homebrew, systemd | 6,193 |
| biotech | rstudio, bioconda | 6,102 |
| blockchain | bitcoin, ethereum, ipfs | 5,141 |

[1]  http://blog.frankzhao.cn/open_rank_and_open_galaxy/

[2]  http://www.gharchive.org/

# OpenGalaxy 3D

Xingran Zhang, Xiaoya Xia, Shengyu Zhao, Wei Wang, **OpenGalaxy: An interactive exploration platform for a visualized GitHub Full Domain collaboration network**, ICPC, 2024.

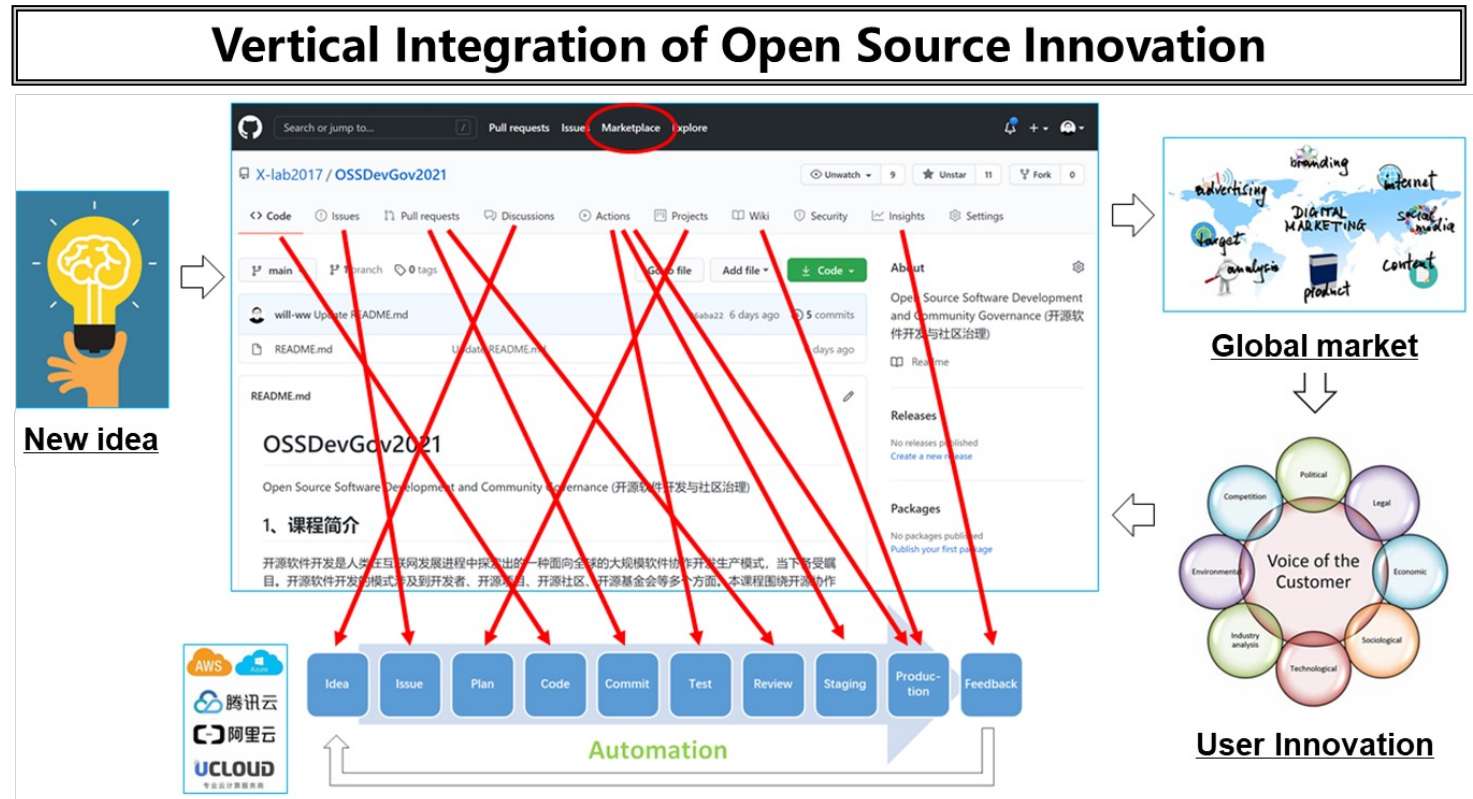# Motivation: GitHub Platform and Open Data

## GitHub

- The world's largest code hosting platform
- Acquired by Microsoft in May 2019 for $7.5 billion

## Not Just Code

- Issue Management
- Distributed Collaboration
- Continuous Integration
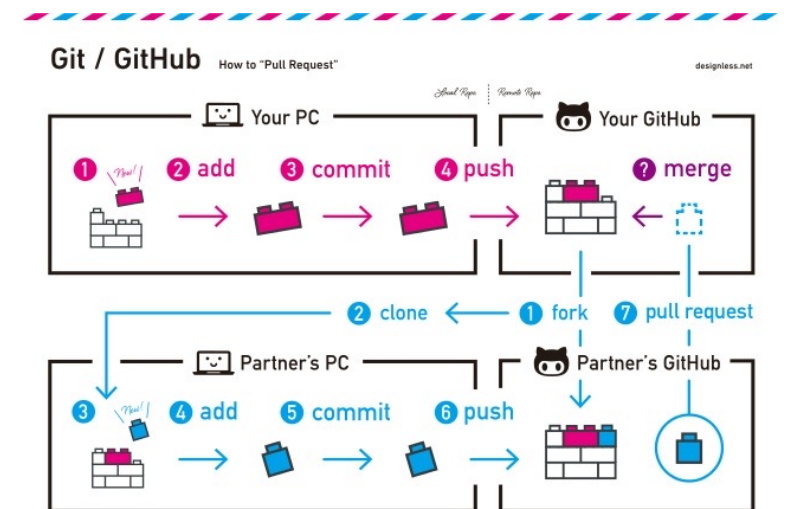- Project Management
- Security Risk Management

## Open data

- GitHub Restful API
- GitHub Event log

# Method: Data-Driven Developer Behavioral Science

**Activity data** in open source software development and ecosystem evolution is a very broad concept. Any data generated in the process of **software development**, **maintenance**, **operation**, as well as ecosystem **governance**, **evolution**, etc., can be called open source software activity data, including but not limited to:

- Git/GitHub Log Event
- Source Code
- Documentation
- Configuration Files
- Changes

- Development Process
- Developers
- Package Hosting Platforms
- Social Data
- Software Ecosystem Network

# OpenDigger project



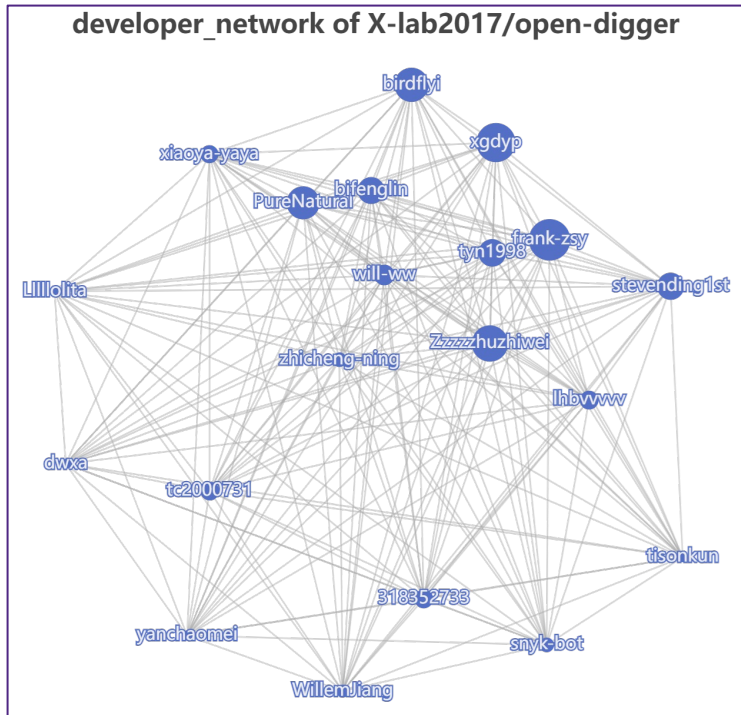- **GitHub action log entries: 5.8 billion**
- **Artifact repository data from NPM/PyPI, etc: 6.2 million entries**
- **CVE security vulnerability data: 160,000 entries**
- **StackOverflow Q&A posts: 25 million**
- **Labeled data, including 413 GitHub orgs, covering 89,427 repositories**
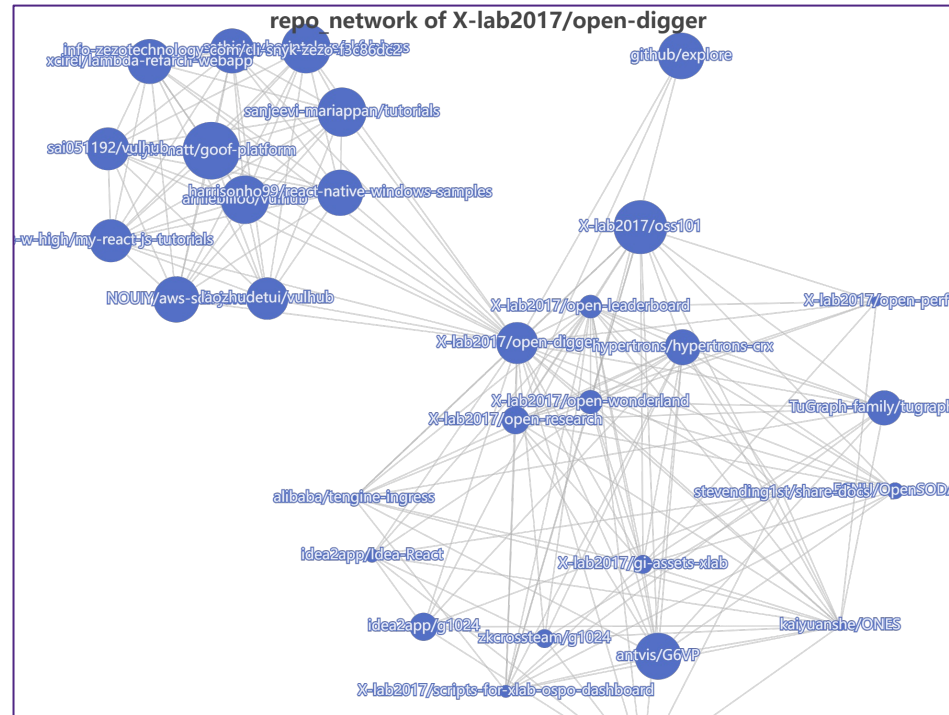
**Research chanllenges: Identification, Recognition, Accounting, and Rewarding of Open Source Contributions**

# The Essence of Open Source Contribution: A Graph Perspectives
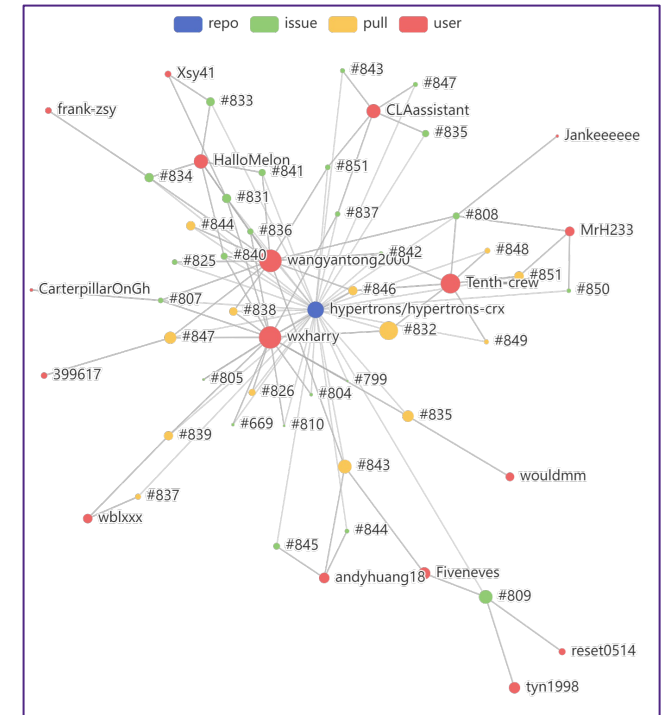


**Deveoper network**

**Repo network**

**Value unit network**

From the Perspective of Digital Economics: Not only is the Internet a network, but all economic phenomena are networks

# OpenRank Developer Contribution Evaluation

The **OpenRank algorithm** is an evaluative method that generalizes the PageRank algorithm to accommodate directed, weighted, heterogeneous networks with initial values that are not necessarily strongly connected. When applied to the assessment of contributions in the open-source context.

- **Significant positive correlation with traditional metrics. Developers endorse the OpenRank results.**
- **Noticeable impact on developers already in the projects, particularly in issue discussions, PRs submission and emojis.**
- **Observation of desirable developers' behavior impact and improvement of community collaboration.**



Research Questions

RQ1
OR Effectiveness
Comparison
Survey

RQ2
Impact on Projs
RDD over
CHAOSS metrics

RQ3
Devs Perception
Semi-structured
Interview

Shengyu Zhao, Xiaoya Xia, Brian Fitzgerald, et al., **Motivating Open Source Collaborations Through Social Network Evaluation: A Gamification Practice from Alibaba**, ICSE 2024.

# The Applications of OpenRank



Corporate Open Source Governance Dashboard

Community Multi-Project Governance Dashboard

Community Single Project Governance Dashboard

Open Source Community Incentive Dashboard

# From OpenRank to OpenPerf



OpenPerf enhances the sustainability and growth of the OSS ecosystem by providing tools for measuring and evaluating project metrics, enabling data mining for research, offering methodologies for ranking projects, and assessing contribution levels.

# OpenPerf Suite Architecture



**Table 1: Benchmarking Tasks**

| Benchmarking Task | Data Type | Problem Type | Scene | Research Field |
|---|---|---|---|---|
| Behavior Data Completion and Prediction[11] | Time Series | Regression | Enterprise Governance | Data Flow |
| OSS Bot Identification and Classification[6] | Time Series | Classification | Software Development | Data Flow |
| Community Sentiment Classification[54, 62] | Text Data | Classification | Community Operations | NLP |
| Software Supply Chain Risk Prediction[32] | Time Series | Regression | Ecosystem Strategy | Complex Networks |
| Project Influence Ranking[68] | Graph & Network | Ranking | Community Operations | Complex Networks |
| Archived Project Prediction[60] | Time Series | Regression | Enterprise Governance | Web Mining |
| Network Metric Prediction[59] | Graph & Network | Regression | Enterprise Governance | Data Flow |
| Community Anomalous Detection[10] | Time Series | Anomaly Detection | Enterprise Governance | Complex Networks |
| OSS Project Recommendation[56] | Graph & Network | Recommendation | Community Operations | Recommendation |

## Influence Ranking Comparison Results

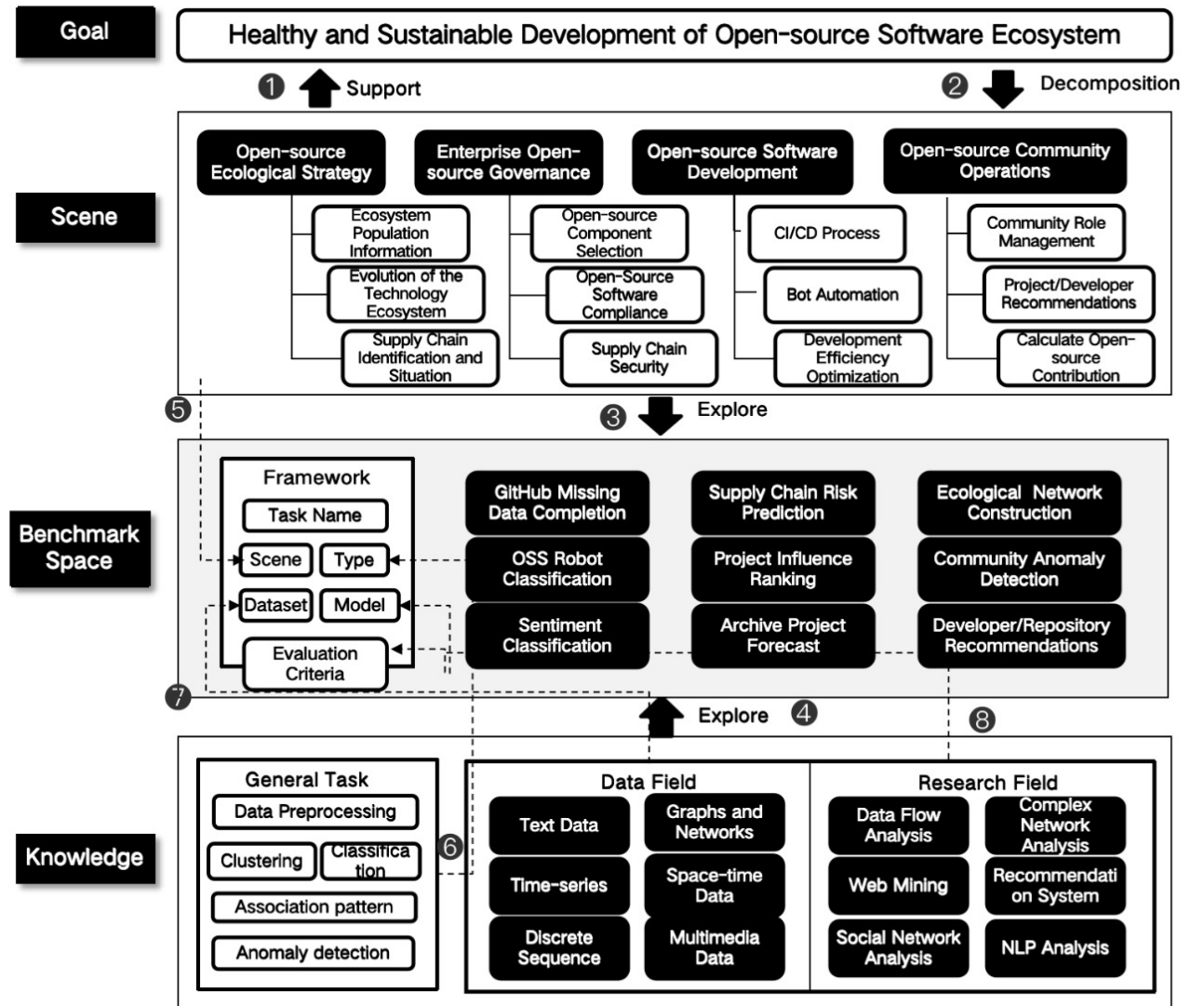| Repository | Degree Centrality | PageRank | OpenRank |
|---|---|---|---|
| home-assistant/core | 0.015660 | 0.0035 | 2393.86 |
| NixOS/nixpkgs | 0.008743 | 0.0008 | 2207.5 |
| microsoft/vscode | 0.015247 | 0.003 | 1960.39 |
| flutter/flutter | 0.012138 | 0.002 | 1460.34 |
| pytorch/pytorch | 0.009624 | 0.0012 | 1421.18 |
| azure-docs | 0.239616 | 0.08 | 1216.01 |
| dotnet/runtime | 0.004141 | 0.0006 | 1181.12 |
| winget-pkgs | 0.061954 | 0.0075 | 1106.3 |
| godotengine/godot | 0.203330 | 0.045 | 1105.51 |
| odoo/odoo | 0.175534 | 0.043 | 907.97 |

# OSGraph

**OSGraph (Open Source Graph)** is an open-source graph-based analytics tool that leverages the comprehensive graph of GitHub open-source data to provide insights into developer behavior and project community ecosystems. It offers developers, project owners, DevRel advocate, and community operators a clear and intuitive view of open-source data, helping you and your project to create a personalized open-source business card, find compatible development partners, and unearth deep community value.





Project Contribution Graph

Project Ecosystem Graph

Project Community Graph

Developer Activity Graph

Open-source Partner Graph

Open-source Interest Graph

OSGraph: A Data Visualization Insight Platform
for Open Source Community

Wenrui Huang[1][0009−0009−8645−3193], Xiaoya Xia[1], Aoying Zhou[1], Xuan Zhou[1],
Wei Wang[1], Shengyu Zhao[2], Zhiyong Wang[3], and Sikang Bian[3]

[1] East China Normal University, Shanghai, China
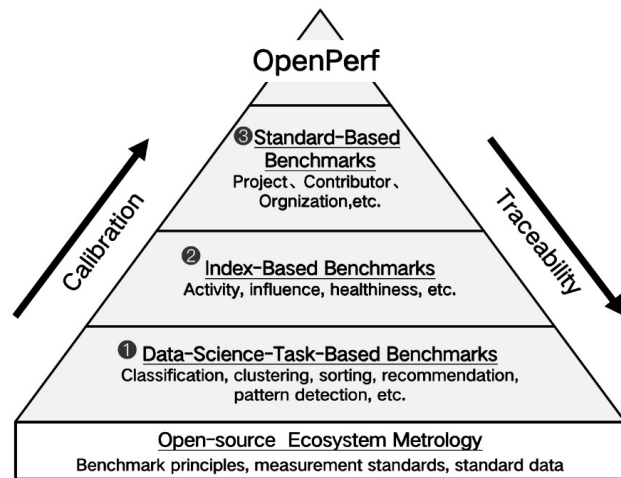hwr0577@gmail.com, xiaoya@stu.ecnu.edu.cn, ayzhou@sei.ecnu.edu.cn,

DASFAA 2024
Gifu Japan

Wenrui Huang, Xiaoya Xia, Shengyu Zhao, Wei Wang, **OSGraph: A Data Visualization Insight Platform for Open Source Community**, DASFFA, 2024

# OpenPerf × OSGraph

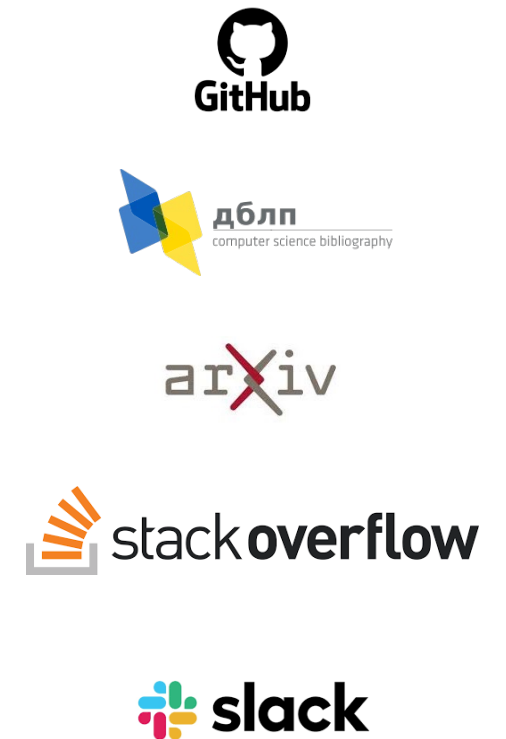| Scenario / Task | Graph Data Mining Tasks | Graph Neural Network Analysis Tasks | Network Science Tasks |
|---|---|---|---|
| **Project Contribution Graph**<br>**Project Ecosystem Graph**<br>**Project Community Graph**<br>**Developer Activity Graph**<br>**Open-source Partner Graph**<br>**Open-source Interest Graph** | • **Attribute Analysis**<br>• **Graph Matching**<br>• **Graph Retrieval**<br>• **Graph Clustering**<br>• **Graph Classification**<br>• **Frequent Subgraph Mining**<br>• **Graph Pattern**<br>• **Link Prediction**<br>• **Anomaly Detection** | • **Representation Learning**<br>• **Node Classification**<br>• **Graph Classification**<br>• **Link Prediction**<br>• **Graph Matching**<br>• **AutoML**<br>• **Dynamic Graphs**<br>• **Heterogeneous Graphs** | • **Properties Analysis**<br>• **Models Analysis**<br>• **Evolution**<br>• **Degree Correlation**<br>• **Robustness Analysis**<br>• **Communities Analysis**<br>• **Propagation** |

# Connected World = OpenPerf × OSGraph × ...

# References

1. Shengyu Zhao, Xiaoya Xia, Brian Fitzgerald, et al., **Motivating Open Source Collaborations Through Social Network Evaluation: A Gamification Practice from Alibaba**, *International Conference on Software Engineering (ICSE),* 2024.

2. Liang Chen, Wei Wang, Yun Yang, **Temporal Autoregressive Matrix Factorization for High-dimensional Time Series prediction of OSS**, *IEEE Transactions on Neural Networks and Learning Systems*, 2024.

3. Yenan Tang, Shengyu Zhao, Xiaoya Xia, et al., **HyperCRX: A Browser Extension for Insights into GitHub Projects and Developers**, *International Conference on Program Comprehension* (*ICIP*)**,** 2024.

4. Xinran Zhang, Shengyu Zhao, Yenan Tang, et al., **OpenGalaxy: An Interactive Exploration Platform for a Visualized GitHub Full Domain Collaboration Network**, *International Conference on Program Comprehension* (*ICIP*), 2024.

5. Wenrui Huang, Xiaoya Xia, Aoying Zhou, et al., **OSGraph: A Data Visualization Insight Platform for Open Source Community**, *International Conference on Database Systems for Advanced Applications* (*DASFAA*), 2024.

6. Xiaoya Xia, Wei Wang, Shengyu Zhao, **Understanding the Archived Projects on GitHub**, *IEEE SANER*, 2023.

7. **OpenDigger**: https://github.com/X-lab2017/open-digger

8. **OpenGalaxy:** https://github.com/X-lab2017/open-galaxy

9. **OpenPerf**: https://arxiv.org/abs/2311.15212

10. **OSGraph**: https://github.com/TuGraph-family/OSGraph