# CS846 Project Infrastructure Overview

2025.01.30

**Xiang Chen**,
David R. Cheriton School of Computer Science

**Instructor: Mike Godfrey**
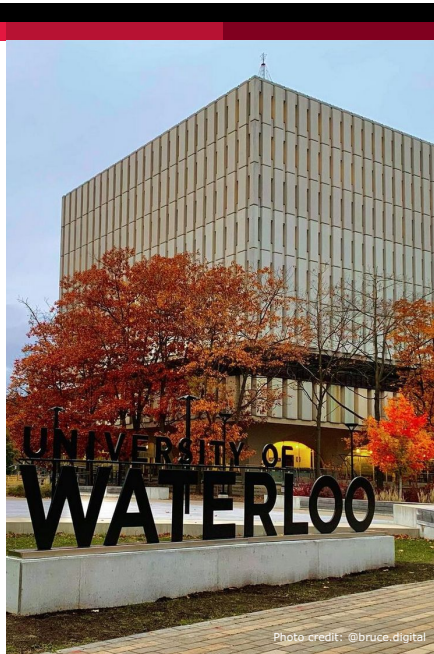
UNIVERSITY OF WATERLOO

Photo credit: @bruce.digital

---

# Taken From MSR 2025 Mining Challenge

https://2025.msrconf.org/track/msr-2025-mining-challenge#Call-for-Mining-Challenge-Papers

UNIVERSITY OF WATERLOO

---

# What is the MSR Mining Challenge?

- Annual challenge by the Mining Software Repositories (MSR) conference.

UNIVERSITY OF WATERLOO

---

# What is the MSR Mining Challenge?

- Annual challenge by the Mining Software Repositories (MSR) conference.
- Focuses on analyzing real-world software data.

UNIVERSITY OF WATERLOO

## What is the MSR Mining Challenge?

- Annual challenge by the Mining Software Repositories (MSR) conference.

- Focuses on analyzing real-world software data.

- Provides specific datasets and tools for researchers.

UNIVERSITY OF
WATERLOO

## What is the MSR Mining Challenge?

- Annual challenge by the Mining Software Repositories (MSR) conference.

- Focuses on analyzing real-world software data.

- Provides specific datasets and tools for researchers.

- Encourages innovative methods and insights.

UNIVERSITY OF
WATERLOO

## What is the MSR Mining Challenge?

- Annual challenge by the Mining Software Repositories (MSR) conference.

- Focuses on analyzing real-world software data.

- Provides specific datasets and tools for researchers.

- Encourages innovative methods and insights.

- Aims to advance software repository mining techniques.

UNIVERSITY OF
WATERLOO

## Theme For MSR 2025 Mining Challenge

- **Theme**: Dependency Analysis with **Goblin** framework

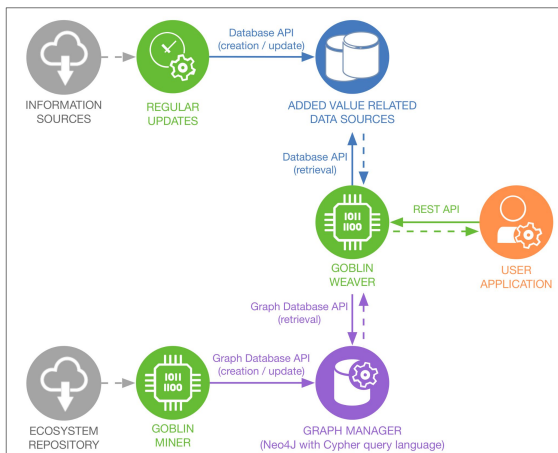UNIVERSITY OF
WATERLOO

## Theme For MSR 2025 Mining Challenge

- **Theme**: Dependency Analysis with **Goblin** framework

- **Goblin Framework Components**
  - Neo4J Maven Central dependency graph
  - **Miner**: generates the dependency graph
  - **Weaver**: Adds custom metrics to dependency graphs, and query the graph

UNIVERSITY OF
**WATERLOO**

---

## Theme For MSR 2025 Mining Challenge

- **Theme**: Dependency Analysis with **Goblin** framework

- **Goblin Framework Components**
  - Neo4J Maven Central dependency graph
  - **Miner**: generates the dependency graph
  - **Weaver**: Adds custom metrics to dependency graphs, and query the graph

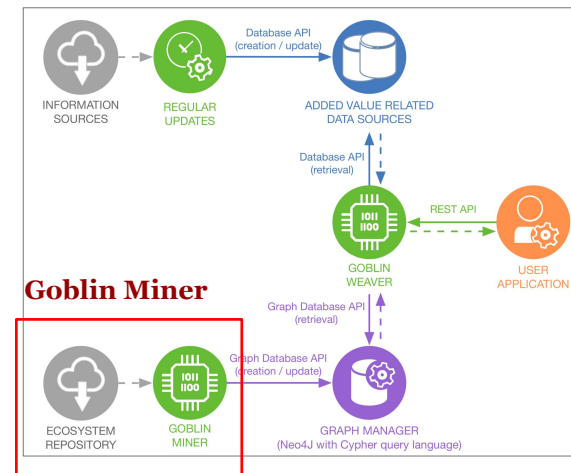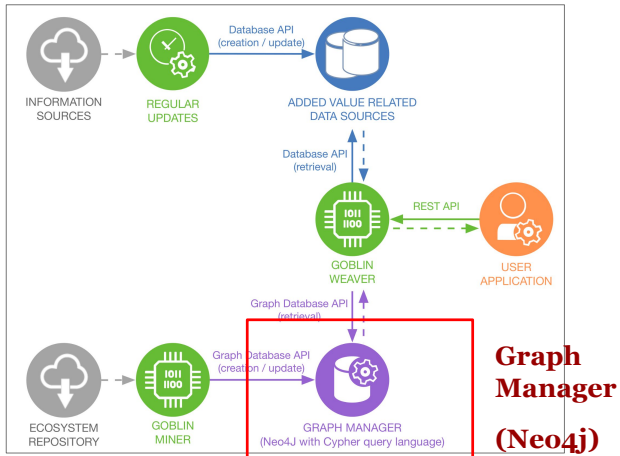- **Goblin**: A customizable tool for studying software ecosystems and dependencies

UNIVERSITY OF
**WATERLOO**

---

## Goblin Framework

UNIVERSITY OF
**WATERLOO**

---

## Goblin Framework

UNIVERSITY OF
**WATERLOO**

## Goblin Framework



**Graph Manager (Neo4j)**

UNIVERSITY OF WATERLOO

---

## Goblin Framework
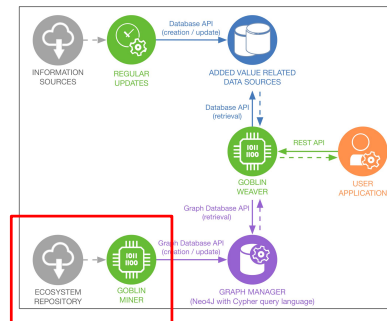


**Goblin Weaver**

UNIVERSITY OF WATERLOO

---

## Goblin Framework

### Goblin Miner

- Retrieve all releases in the Lucene Maven Central Index archive

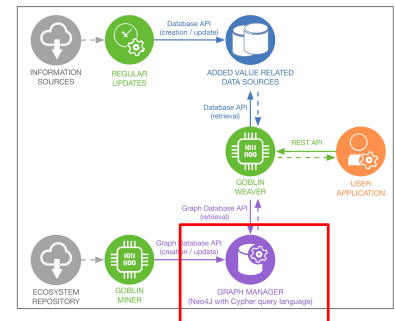- Retrieve their direct dependencies with the `org.eclipse.aether` library

UNIVERSITY OF WATERLOO
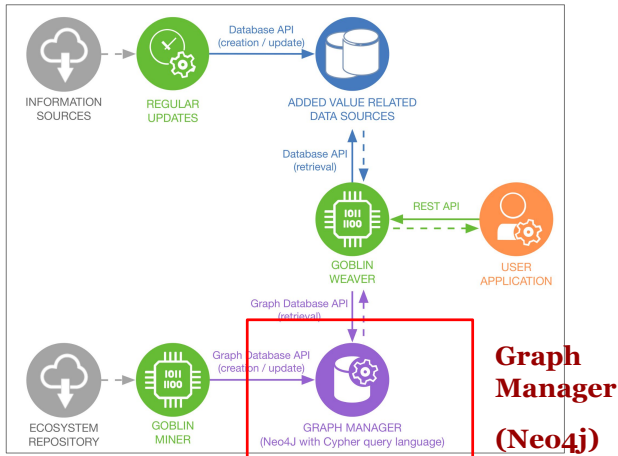
---

## Goblin Framework

### Goblin Miner

- Retrieve all releases in the Lucene Maven Central Index archive

- Retrieve their direct dependencies with the `org.eclipse.aether` library



**Dependency graph**

UNIVERSITY OF WATERLOO
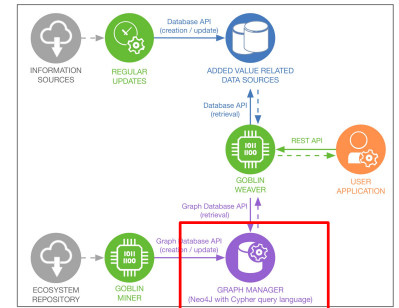
## Slide (Page 17)

# Goblin Framework



**Graph Manager (Neo4j)**

UNIVERSITY OF WATERLOO

## Slide (Page 18)

# Goblin Framework



**Graph Manager (Neo4j)**

- Node types
  - Libraries (type `Artifact`)
  - Releases (type `Release`)

UNIVERSITY OF WATERLOO
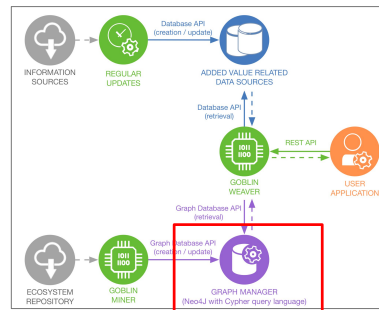
## Slide (Page 19)

# Goblin Framework



**Graph Manager (Neo4j)**

- Node types
  - Libraries (type `Artifact`)
  - Releases (type `Release`)

- Edge types
  - Dependencies (type `dependency`) are from `Release` to `Artifact`
  - Versioning (type `relationship_AR`) are from `Artifact` to `Release`
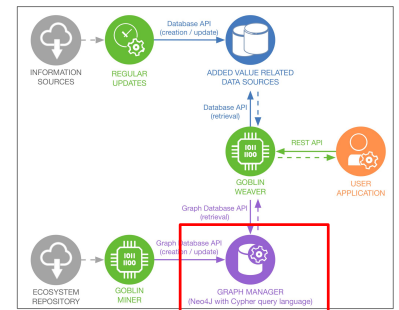
UNIVERSITY OF WATERLOO

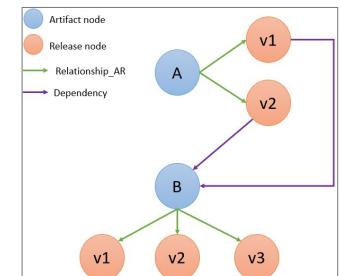## Slide (Page 20)

# Goblin Framework



**Graph Manager (Neo4j)**

- Node types
  - Libraries (type `Artifact`)
  - Releases (type `Release`)

- Edge types
  - Dependencies (type `dependency`) are from `Release` to `Artifact`
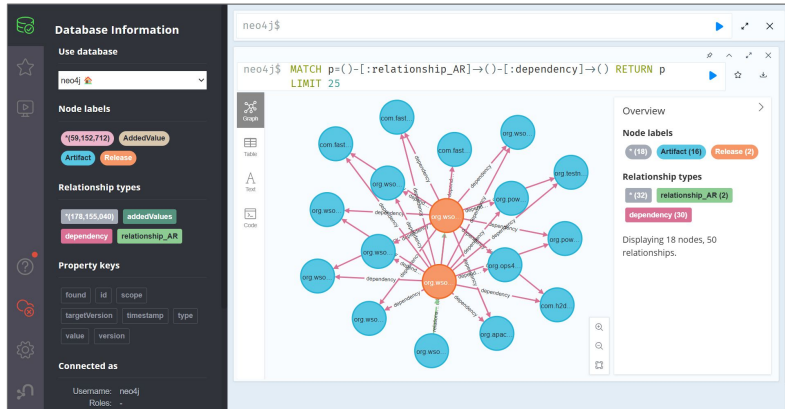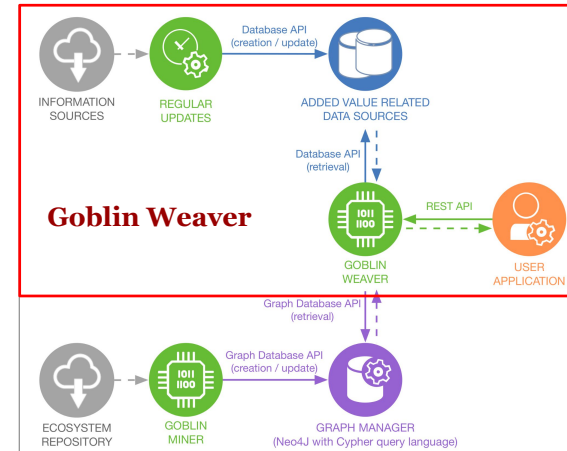  - Versioning (type `relationship_AR`) are from `Artifact` to `Release`

# Goblin Framework

## Graph Manager (Neo4j)

UNIVERSITY OF WATERLOO

---

# Goblin Framework



**Goblin Weaver**

UNIVERSITY OF WATERLOO

---

# Goblin Framework

## Goblin Weaver

- On-demand enrichment of the dependency graph

UNIVERSITY OF WATERLOO

---

# Goblin Framework

## Goblin Weaver

- On-demand enrichment of the dependency graph

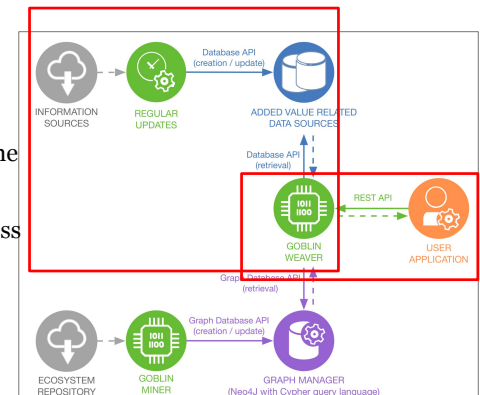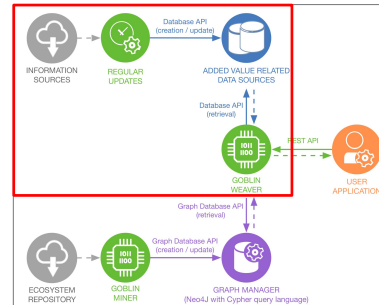- An alternative for direct access to the Neo4j database

UNIVERSITY OF WATERLOO

# Goblin Framework

## Goblin Weaver

- On-demand enrichment of the dependency graph

- An alternative for direct access to the Neo4j database
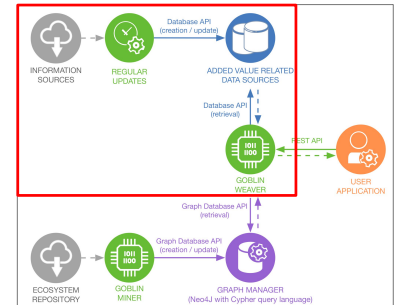
UNIVERSITY OF WATERLOO

---

# Goblin Framework

## Goblin Weaver

- On-demand enrichment of the dependency graph

❖ Release nodes added values

UNIVERSITY OF WATERLOO

---

# Goblin Framework

## Goblin Weaver

- On-demand enrichment of the dependency graph

❖ Release nodes added values

   a.  CVE (Common Vulnerabilities and Exposures): use the osv.dev dataset
- Name, cwe (type of vulnerability) and severity (low, moderate, high, critical)

UNIVERSITY OF WATERLOO

---

# Goblin Framework

## Goblin Weaver

- On-demand enrichment of the dependency graph

❖ Release nodes added values

   a.  CVE (Common Vulnerabilities and Exposures): use the osv.dev dataset
- Name, cwe (type of vulnerability) and severity (low, moderate, high, critical)

   b.  FRESHNESS
- The number of more recent releases available
- The time elapsed between it and the most recent release
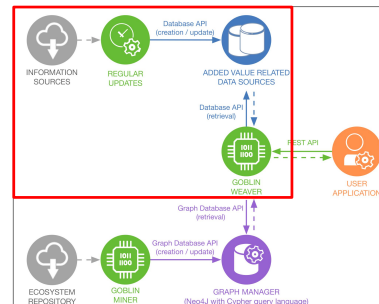
UNIVERSITY OF WATERLOO

# Goblin Framework

## Goblin Weaver

- On-demand enrichment of the dependency graph
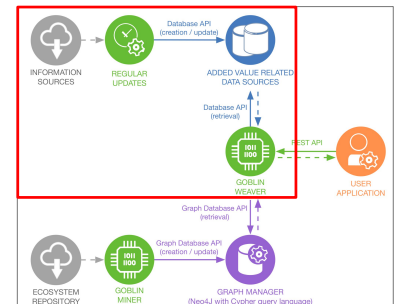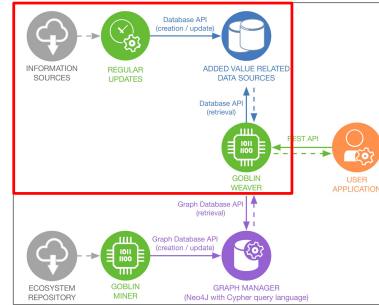


❖ <u>Release nodes added values</u>

   a. CVE (Common Vulnerabilities and Exposures): use the <u>osv.dev</u> dataset
- Name, cwe (type of vulnerability) and severity (low, moderate, high, critical)

   b. FRESHNESS
- The number of more recent releases available
- The time elapsed between it and the most recent release

   c. POPULARITY_1_YEAR
- Number of dependants over a one year window

---

# Goblin Framework

## Goblin Weaver

---

# Goblin Framework

## Goblin Weaver

- On-demand enrichment of the dependency graph



❖ <u>Artifact nodes added values</u>
- ➢ SPEED
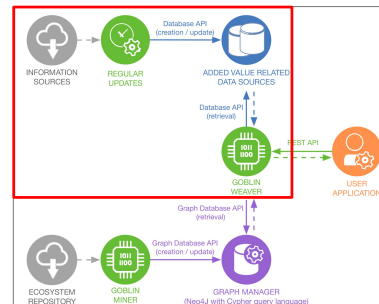  - Average number of releases per day of a library

---

# Goblin Framework

## Goblin Weaver

# Goblin Framework

## Goblin Weaver

- An alternative for direct access to the Neo4j database

UNIVERSITY OF WATERLOO

---

# Goblin Framework

## Goblin Weaver

- An alternative for direct access to the Neo4j database

**Cypher**: Neo4j's built-in query language

UNIVERSITY OF WATERLOO

---

# Goblin Framework

## Goblin Weaver

- An alternative for direct access to the Neo4j database

**Cypher**: Neo4j's built-in query language

https://neo4j.com/docs/cypher-manual/current/queries/basic/

UNIVERSITY OF WATERLOO

---

# Goblin Framework

## Goblin Weaver

- An alternative for direct access to the Neo4j database

**Cypher**: Neo4j's built-in query language

https://neo4j.com/docs/cypher-manual/current/queries/basic/

MATCH (r:Release) WHERE r.id='org.jgrapht:jgrapht-core:1.5.2' RETURN r

UNIVERSITY OF WATERLOO

# Goblin Framework

**Cypher**: Neo4j's built-in query language

---

# Goblin Framework



## Goblin Weaver

- On-demand enrichment of the dependency graph

- An alternative for direct access to the Neo4j database

---

# Goblin Framework

## Goblin Weaver

---

# Goblin Framework



## Goblin Weaver Example

- Which are the latest versions available after `jgrapht-core 1.5.0`?
- Add their CVE, freshness and popularity

# Goblin Framework

## Goblin Weaver Example

- Which are the latest versions available after `jgrapht-core 1.5.0`?
- Add their CVE, freshness and popularity



```
Method: POST
ROUTE: /release/newVersions
Body:
{
    "groupId": "org.jgrapht",
    "artifactId": "jgrapht-core",
    "version": "1.5.0",
    "addedValues": ["CVE", "FRESHNESS", "POPULARITY_1_YEAR"]
}
```

---

# Goblin Framework

## Goblin Weaver Example



POST  /release/newVersions   Get newer versions of a release from GAV

Get newer versions of a release from groupId:ArtifactId:Version with added values

Parameters                    Cancel      Reset

No parameters

Request body required                     application/json

```
{
    "groupId": "org.jgrapht",
    "artifactId": "jgrapht-core",
    "version": "1.5.0",
    "addedValues": ["CVE", "FRESHNESS", "POPULARITY_1_YEAR"]
}
```
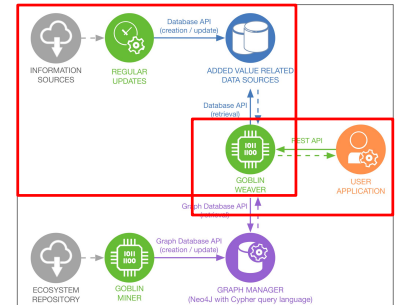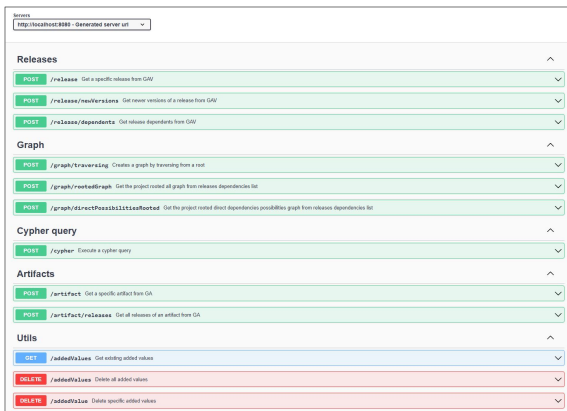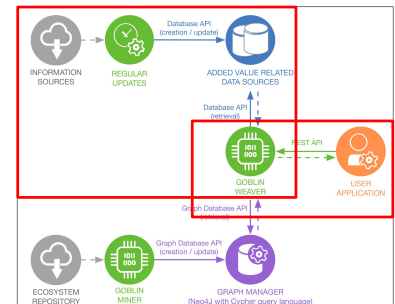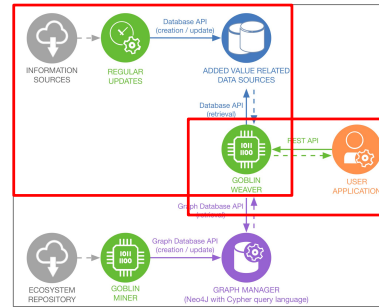
---

# Goblin Framework

## Goblin Weaver Example

```
{
    "nodes": [
        {
            "cve": [],
            "id": "org.jgrapht:jgrapht-core:1.5.1",
            "nodeType": "RELEASE",
            "freshness": {
                "numberMissedRelease": "1",
                "outdatedTimeInMs": "66882028000"
            },
            "version": "1.5.1",
            "popularity_1_year": 105,
            "timestamp": 1616171280000
        },
        {
            "cve": [],
            "id": "org.jgrapht:jgrapht-core:1.5.2",
            "nodeType": "RELEASE",
            "freshness": {
                "numberMissedRelease": "0",
                "outdatedTimeInMs": "0"
            },
            "version": "1.5.2",
            "popularity_1_year": 953,
            "timestamp": 1683053308000
        }
    ]
}
```

---

# Goblin Framework

## Goblin Weaver Example

```
{
    "nodes": [
        {
            "cve": [],
            "id": "org.jgrapht:jgrapht-core:1.5.1",
            "nodeType": "RELEASE",
            "freshness": {
                "numberMissedRelease": "1",
                "outdatedTimeInMs": "66882028000"
            },
            "version": "1.5.1",
            "popularity_1_year": 105,
            "timestamp": 1616171280000
        },
        {
            "cve": [],
            "id": "org.jgrapht:jgrapht-core:1.5.2",
            "nodeType": "RELEASE",
            "freshness": {
                "numberMissedRelease": "0",
                "outdatedTimeInMs": "0"
            },
            "version": "1.5.2",
            "popularity_1_year": 953,
            "timestamp": 1683053308000
        }
    ]
}
```

# Goblin Framework

```
Code     Details

200

Response body
    "nodes": [
        {
            "cve": [],
            "id": "org.jgrapht:jgrapht-core:1.5.1",
            "nodeType": "RELEASE",
            "freshness": {
                "numberMissedRelease": "1",
                "outdatedTimeInMs": "66882028000"
            },
            "version": "1.5.1",
            "popularity_1_year": 105,
            "timestamp": 1616171280000
        },
        {
            "cve": [],
            "id": "org.jgrapht:jgrapht-core:1.5.2",
            "nodeType": "RELEASE",
            "freshness": {
                "numberMissedRelease": "0",
                "outdatedTimeInMs": "0"
            },
            "version": "1.5.2",
            "popularity_1_year": 953,
            "timestamp": 1683053308000
        }
    ]
}

Response headers

connection: keep-alive
content-type: application/json
date: Wed,29 Jan 2025 03:27:20 GMT
keep-alive: timeout=60
transfer-encoding: chunked
```

PAGE  45

UNIVERSITY OF
WATERLOO

---

# Goblin Framework

```
Code     Details

200

Response body
    "nodes": [
        {
            "cve": [],
            "id": "org.jgrapht:jgrapht-core:1.5.1",
            "nodeType": "RELEASE",
            "freshness": {
                "numberMissedRelease": "1",
                "outdatedTimeInMs": "66882028000"
            },
            "version": "1.5.1",
            "popularity_1_year": 105,
            "timestamp": 1616171280000
        },
        {
            "cve": [],
            "id": "org.jgrapht:jgrapht-core:1.5.2",
            "nodeType": "RELEASE",
            "freshness": {
                "numberMissedRelease": "0",
                "outdatedTimeInMs": "0"
            },
            "version": "1.5.2",
            "popularity_1_year": 953,
            "timestamp": 1683053308000
        }
    ]
}

Response headers

connection: keep-alive
content-type: application/json
date: Wed,29 Jan 2025 03:27:20 GMT
keep-alive: timeout=60
transfer-encoding: chunked
```

PAGE  46

UNIVERSITY OF
WATERLOO

---

# Some Useful Links

1. **Challenge preprint**:
https://hal.science/hal-04777703

### Navigating and Exploring Software Dependency Graphs using Goblin

Damien Jaime
*Sorbonne Université, CNRS, LIP6,
Université Paris Nanterre*
F-75005, Paris, France
damien.jaime@lip6.fr

Joyce El Haddad
*Université Paris Dauphine-PSL,
CNRS, LAMSADE*
CF-75016, Paris, France
joyce.elhaddad@lamsade.dauphine.fr

Pascal Poizat
*Sorbonne Université, CNRS, LIP6,
Université Paris Nanterre*
F-75005, Paris, France
pascal.poizat@lip6.fr

*Abstract*—Using package managers is a simple and common method for reusing code through project dependencies. However, these, direct, dependencies can themselves rely on additional packages, resulting in indirect dependencies. It may then become complex to get a grasp of the whole set of dependencies of a project. Beyond studying individual projects, a deep understanding of software ecosystems is also a critical prerequisite for achieving sustained success in software development. This paper presents the 2025 edition of the MSR conference mining challenge. This year's mining challenge focuses on dependencies and dependency ecosystem analysis using the Goblin framework that has been presented at the previous edition of the MSR conference. Goblin is composed of a Neo4j Maven Central dependency graph and a tool called Weaver for on-demand metric weaving into dependency graphs. As a whole, Goblin is a customizable framework for ecosystem and dependency analysis.

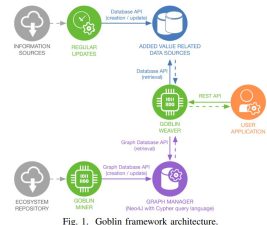*Index Terms*—software ecosystem, dependency graph, dataset, mining software repositories, maven central

Fig. 1.  Goblin framework architecture.

UNIVERSITY OF
WATERLOO

---

# Some Useful Links

2. **Goblin framework**:
https://dl.acm.org/doi/abs/10.1145/3643991.3644879?casa_token=bTRwoEEBLooAAAAA:9v4cP65Fwic4JmVCEg0mGGA4XYSlbg6vaMJ-TAbbw84kXgJrt3sFTGrm1UTRGCtSTbFLV7ySOjKy

### GOBLIN: A Framework for Enriching and Querying the Maven Central Dependency Graph

Damien Jaime
Sorbonne Université, CNRS, LIP6
F-75005, Paris, France
SAP France S.A
F-92300, Levallois-Perret, France
damien.jaime@lip6.fr

Joyce El Haddad
Université Paris Dauphine-PSL,
CNRS, LAMSADE
F-75016, Paris, France
joyce.elhaddad@lamsade.dauphine.fr

Pascal Poizat
Sorbonne Université, CNRS, LIP6
F-75005, Paris, France
Université Paris Lumières, Université
Paris Nanterre
F-92000, Nanterre, France
pascal.poizat@lip6.fr

**ABSTRACT**

Dependency graphs support software maintenance and software ecosystem analysis. Several metrics can be used on top of these graph models but the set of such metrics is able to evolve over time. Further, some metrics have a dynamic nature, requiring being able to "rewind" dependency graphs at some point in time. To address these issues we propose the GOBLIN framework. It is composed of a dependency graph metamodel with time-related information, a miner to retrieve the graph from Maven Central, and a tool for on-demand metric weaving into dependency graphs. As a whole, GOBLIN is a customizable framework for ecosystem and dependency analysis. This is illustrated with a set of complementary experiments. Our tools, datasets, and experiments are freely available online.

dependency relations. On top of these models, it is needed to compute different metrics before being able to measure the quality of a project (from a dependency perspective) and to support or even suggest updates. There are several metrics of interest here. The set of CVEs (Common Vulnerabilities and Exposures) concerning a dependency is central for security. Other metrics incorporate time, e.g., freshness [3] or rhythm [5].

These metrics could be part of a DG metamodel. Yet, they are numerous and evolve over time (and research). We believe that a better solution is to weave them on-demand over DG models. Some of these metrics have a dynamic nature, e.g., CVEs, freshness, and rhythm, are all not constant in time given some package. Dependencies may support ranges (instead of requiring version g:a:1.0

PAGE  48

UNIVERSITY OF
WATERLOO

# Some Useful Links

3. **Maven Central Neo4j dependency graph datasets:**

https://zenodo.org/records/13734581

# Some Useful Links

4. **Goblin Weaver:**

https://github.com/Goblin-Ecosystem/goblinWeaver

# Some Useful Links

5. **Goblin Miner:**

https://github.com/Goblin-Ecosystem/goblinDependencyMiner

# Some Useful Links

6. **Goblin tutorial**:

https://github.com/Goblin-Ecosystem/goblinTutorial?tab=readme-ov-file

# Taken From MSR 2025 Mining Challenge

https://2025.msrconf.org/track/msr-2025-mining-challenge#Call-for-Mining-Challenge-Papers

## MSR 2025
### 22nd International Conference on Mining Software Repositories
April 28-29, Ottawa, Canada

co-located with ICSE 2025

Attending ▾   Tracks ▾   Organization ▾   🔍 Search   Series ▾          Sign in   Sign up

🏠 ICSE 2025 (series) / 🏠 MSR 2025 (series) /

### Mining Challenge

MSR 2025

Accepted Papers   **Call for Mining Challenge Papers**   Call for Mining Challenge Proposals

#### Call for Mining Challenge Papers

**NEW: Challenge preprint available** here

Using package managers is a simple and common method for reusing code through project dependencies. However, these direct dependencies can themselves rely on additional packages, resulting in indirect dependencies. It may then become complex to get a grasp of the whole set of dependencies of a project. Beyond individual projects, a deep understanding of how software ecosystems work and evolve is also a critical prerequisite for achieving sustained success in software development.

This year's mining challenge focuses on dependencies and dependency ecosystem analysis using the Goblin framework that has been presented at the previous edition of the MSR conference. Goblin is composed of a **Neo4J Maven Central dependency graph** and a tool called **Weaver** for on-demand metric weaving into dependency graphs. As a whole, Goblin is a customizable framework for ecosystem and dependency analysis.

**Important Dates** 🌐🕐 AoE (UTC-12h)

Wed 5 Feb 2025
Camera Ready Deadline

Sun 12 Jan 2025
Author Notification

Fri 6 Dec 2024
Paper Deadline

Tue 3 Dec 2024
Abstract Deadline

Mining Challenge - Program Committee

UNIVERSITY OF WATERLOO

---

# Possible Research Questions

https://2025.msrconf.org/track/msr-2025-mining-challenge#Call-for-Mining-Challenge-Papers

1. Ecosystem evolution
   i. What are the patterns in the growth of the Maven Central graph across different time periods?
   ii. Do libraries tend to use more dependencies than in the past?
   iii. Is the rhythm of library releases higher than in the past, and how has this rhythm evolved over time?
   iv. Does the emergence of project management methods (e.g., agile methods) have any impact on the release rhythm of libraries?
   v. To what extent does the ecosystem contain unmaintained libraries?
   vi. How do projects with unmaintained dependencies cope with the challenges they face?
2. Clustering
   i. Can we deduce different clusters from Maven Central's comprehensive dependency graph? How do these clusters interact with one another?
   ii. Can dependency-based clustering reveal domain-specific groupings, and how well do they align with known categorizations of projects?
   iii. How can clustering be used to identify high-risk clusters in the Maven Central ecosystem?
   iv. Which artifacts serve as the most crucial dependencies for the ecosystem (i.e., most depended upon)?
   v. How do these central nodes affect the overall health and stability of the ecosystem?
3. Dependency update
   i. How often do projects update their dependencies, and what factors influence this frequency (e.g., project size, popularity, type)?
   ii. Whenever an artifact releases a new version, how do its dependents react?
   iii. How does the removal or failure of certain projects affect the overall network (e.g., log4j Vulnerability)?
   iv. How do major versus minor dependency updates differ in frequency and impact?
   v. Do projects tend to avoid major updates due to the potential for breaking changes?
4. Trends
   i. How has the adoption of new frameworks (e.g., Spring Boot, Microservices) changed the dependency structures in Maven Central?
   ii. What impact do modern dependency management tools (e.g., Dependabot) have on the ecosystem?
   iii. How does the adoption of newer Java versions influence dependency graphs?
   iv. Does an artifact's number of dependents correlate with other popularity metrics such as GitHub stars?
5. Graph theory
   i. How do metrics such as degree distribution, clustering coefficient, and average path length characterize the dependency graph?
   ii. Is the graph scale-free, small-world, or does it exhibit other known graph structures?
   iii. Are certain types of projects more likely to be central (hubs) or peripheral (leaves) in the graph structure?
   iv. Is the graph made up of connected components with no relationship between them?
   v. How do shortest path lengths between projects vary, and what does this tell us about the overall connectivity of the ecosystem?
6. Vulnerability
   i. How do vulnerabilities propagate through the dependency network, and which projects are most affected?
   ii. What proportion of releases have vulnerabilities? What is the proportion of releases directly and transitively impacted?
   iii. What is the average time taken to patch a vulnerability in a dependency?
   iv. How do users of an artifact react to the discovery of a vulnerability in that artifact?
7. Licensing and Compliance
   i. Are there dominant license types, and how do they influence the usage and distribution of projects?
   ii. How does the choice of licenses affect the artifact graph structure?
   iii. What percentage of projects have conflicting licenses within their dependency trees?

UNIVERSITY OF WATERLOO