# ClearlyDefined

E. Lynette Rayle (GitHub)

Qing Tomlinson (SAP)

Nick Vidal (Open Source Initiative)

# Introduction

## Mission

ClearlyDefined's mission is to crowdsource a global database of licensing metadata for every software component ever published for the benefit of all

# Problem

With the move towards SBOMs everywhere for compliance and security reasons, organizations will face great challenges to generate these at scale for each stage on the supply chain, for every build or release.

Plus, multiple organizations will have to fix the same missing or wrongly identified licensing metadata over and over again.

# Solution

This is where ClearlyDefined comes in, by serving a cached copy of licensing metadata for each component through a simple API.

Organizations will also be able to contribute back with any missing or wrongly identified licensing metadata, helping to create a database that is accurate for the benefit of all.

# CLEARLYDEFINED

# Bringing clarity to Open Source Software licenses.

Use the Data

Curate Data

Contribute Data

Contribute Code

Add a Harvest

Adopt Practices

# Use the data

API: definitions, curations, harvest, attachments, notices


curl -X GET
"https://api.clearlydefined.io/definitions/npm/npmjs/-/lodash/4.17.21" -H
"accept: */*"



https://api.clearlydefined.io/api-docs/

# Curate the data

"contributionInfo": {

    "summary": "[Test] Update declared license",

    "details": "The declared license should be Apache as per the LICENSE file.",

    "resolution": "Updated declared license to Apache-2.0.",

    "type":"incorrect",

    "removeDefinitions":false

  },

# Contribute data

# Contribute code

# Add a harvest

**30,464,321**
Number of total definitions

**60**
Median licensed score

**30**
Median described score

| | | Total | Licensed | Described |
|---|---|---|---|---|
| npm | | 14,572,253 | 60 | 30 |
| gem | | 899,355 | 61 | 30 |
| pypi | | 2,110,095 | 60 | 100 |
| maven | | 3,015,557 | 60 | 100 |
| nuget | | 3,303,878 | 15 | 30 |
| git | | 2,938,206 | 62 | 100 |
| crate | | 46,504 | 65 | 30 |
| deb | | 379,653 | 4 | 100 |
| debsrc | | 25,115 | 13 | 100 |
| composer | | 696,538 | 60 | 100 |
| pod | | 10,016 | 75 | 100 |

# Adopt best practices



Join the ORT community at **github.com/oss-review-toolkit/ort**

# New developments

# LicenseRef

- [Beyond SPDX: expanding licenses identified by ClearlyDefined](#)

# Conda

- [Better identifying conda packages with ClearlyDefined](#)

# GUAC

- [GUAC adopts license metadata from ClearlyDefined](#)
  [Announcing GUAC v0.8.0 Enhancements](#)

Switch to Lynette

# Case Studies

# A Developer's Look at ClearlyDefined

and how GitHub is using it and why

# Why we're using ClearlyDefined…

ClearlyDefined houses **business critical data** for licenses and attributions.  We want to support the mission of making ClearlyDefined **THE** source of truth.

# Impact of ClearlyDefined at GitHub

GitHub added 17.5 million package licenses

- sourced from ClearlyDefined to our database,

- expanding the license coverage for packages that appear in

  - dependency graph
  - dependency insights
  - dependency review
  - repository's software bill of materials (SBOM)

Crawler == Harvester

# Overview of ClearlyDefined Harvesting Process

**website for ClearlyDefined**
Purpose: Simplified view of definitions

**cdcrawler queue**

**service for ClearlyDefined**
Purpose: Process license definitions
Purpose: API for accessing definitions

**crawler for ClearlyDefined**
Purpose: Run tools that find licenses

API GET /definitions

Add request to crawler queue

Pull from queue

- Read raw results
- Summarize
- Create definition

Read definition
Write definition

Notify harvest complete

Run
- reuse
- licensee
- scancode

Write raw results

**definition** store

**raw results** store

Curation == Human

# Overview of ClearlyDefined Curation Process

# Clearly Defined Change Notifications

**changes-notifications**
Azure blob storage
(public read access)

```
|- changes
|   |- index
|   |- 2024-02-29-22
|   |- 2024-02-29-23
|   |- 2024-03-01-00
|   |- ...
|- gem
|   |- rubygems
|   |   |- -
|   |   |   |- rspec
|   |   |   |   |- revisions
.   |   |   |   |   |- 3.13.0.json
|- ...
```

index - lists all changeset IDs past and present
changeset - lists coordinates of changes named for the
            date and hour when the changeset was created
       example: 2024-02-29-22 is
                    Feb 29, 2024 at 2200 hours

path based on coordinates

definition - package version as a json file holding the
definition

# GitHub & ClearlyDefined

# Handling High Volume of Requests

**changes-notifications**
Azure blob storage
(public read access)

hourly cron reads

local cache
store

Notify harvest complete

GH crawler queue

Package License Gateway
Purpose: Respond to License Queries

Write

Read

Add request to crawler queue

Pull from queue

GET /licenses

service for ClearlyDefined
Purpose: Process license definitions
Purpose: API for accessing definitions

GitHub Harvester for ClearlyDefined
Purpose: Prioritize GitHub Requests

Policy Service
Purpose: License Compliance

POST /requests

Run
- reuse
- licensee
- scancode

Write raw results

raw results
store

# Complying with Attribution Requirements

# How we setup a harvester at GitHub

- Used kubernetes config in [crawler.yaml](crawler.yaml) (in clearlydefined/crawler repo)
  - uses official crawler image in Docker Hub
    (may be moving to GitHub container registry)
  - lists configurations to be set

- Requires a token to write directly to ClearlyDefined's raw result store

- Requires significant hardware

  - Ex. Azure P3V2 (4 Virtual CPU's, 14GB Ram)
    (values for ClearyDefined's production crawler)

Switch to Qing

# ClearlyDefined Adoption @SAP

# SAP Open Source

We believe in co-innovation and collaboration with open source.



# SAP's Open Source Vision

**Driving innovation through open collaboration**

**Empowerment through openness**

**Trust and transparency**

**Community engagement**

# Why ClearlyDefined?

- ClearlyDefined supports an impressive and growing number of harvest sources
- Crowd-sourcing license review process saves time and effort

# ClearlyDefined Adoption – Key Milestones

**2018** — Officially became a member of the community

**2019** — Internal policy developed, completion of pilot

**2020** — Large scale harvest automation

**2022** — Further automation, emerging use cases, contribute to codebase

# Overview of Automatic Harvesting

# ClearlyDefined Benefits

Less time scanning, more time focusing on compliance

Better data quality

Community centered

# Giving Back – Contributions to Data



The ClearlyDefined Lifecycle

1 — Component Harvest
2 — Build Base Definition
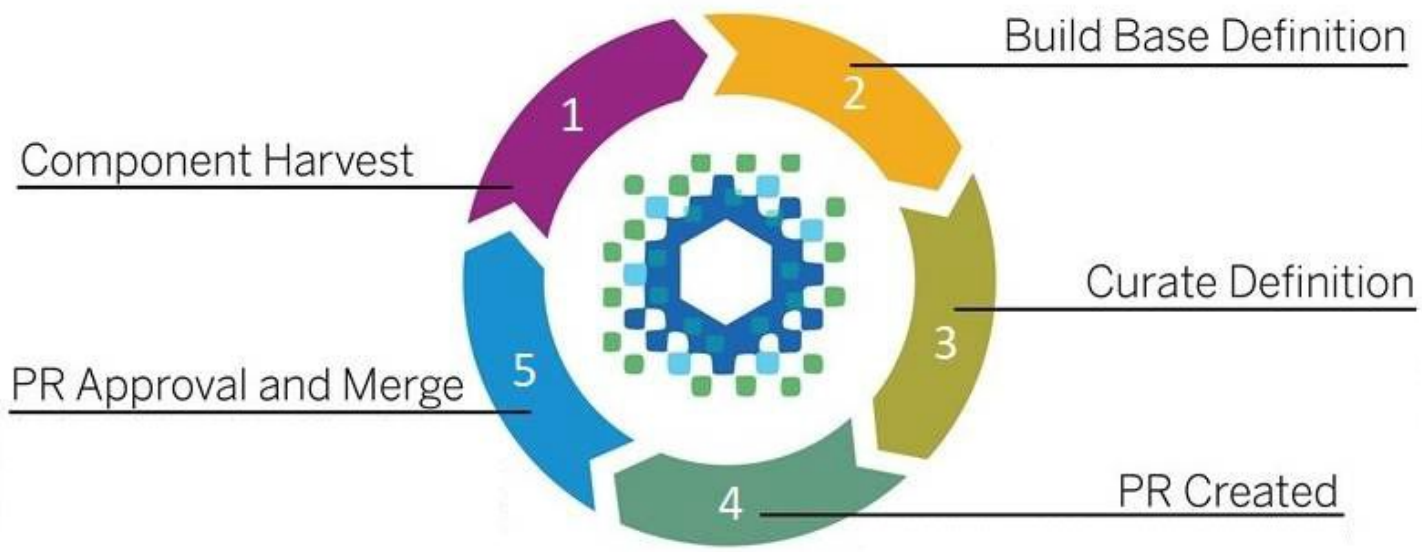3 — Curate Definition
4 — PR Created
5 — PR Approval and Merge

# Giving Back – Contributions to Code

Harvest
Gradle Plugin

Optimizations

Analyses and
Improvements

Collaborations

Switch to Nick

# Governance

# Governance

ClearlyDefined Governance:

- [Charter](Charter)

Governing Board:

- the Executive Director of the Open Source Initiative;
- the Steering Committee Chair; and
- the Outreach Committee Chair.

# Governance

Steering Committee:

- E. Lynette Rayle (GitHub) – Chair

- Qing Tomlinson (SAP)

- Jeff Mendoza (GUAC / Kusari)

# Governance

Outreach Committee:

- Nick Vidal (OSI) - Chair

- Jeff Luszcz (GitHub)

- Brian Duran (SAP)

- Alyssa Wright (Bloomberg)

# Conclusion

# How to get involved...

- [Weekly dev meetings](Weekly dev meetings)
- [https://docs.clearlydefined.io/docs/community/meetings](https://docs.clearlydefined.io/docs/community/meetings)
- [Hang out in Discord](Hang out in Discord)
- Paired programming
- Understand the inner working of ClearlyDefined
- Influence the priorities of development
- Help sustain and keep ClearlyDefined strong
- Set up a harvester on your hardware

Open Collaboration

Thank you

Join us:
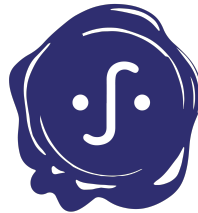https://clearlydefined.io/

ClearlyDefined