# TOSHIBA

Search GitLab, Redmine, and repositories with a single query

# Deploy an enterprise search server with Fess

## FOSDEM 2023

Takashi Kumagai, TOSHIBA Corporate Software Engineering Center
February 4, 2023

## Today's talk

- Our engineers use multiple content management tools

- **Searching** became a problem

  - Can't search laterally across multiple tools
  - Can't search texts inside binary files

- How we use Fess to solve the problem of search

- How we overcame (some of the) shortcomings of Fess

- Contributions (patches)

- Conclusion

# Contents

# Contents

# Who am I

- An engineer at Corporate Software Engineering Center in Toshiba
- 2006 – 2015
    - Develop applications for Toshiba's products (HDTV, etc.)
- 2015 – 2022
    - Maintain and improve the company's cloud infrastructure and automation processes.
        - Build in-house tools
        - Write automation scripts
        - Automate E2E tests
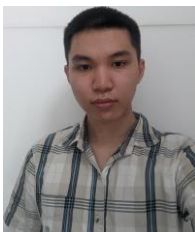        - Configure and deploy search servers to departments in other Toshiba companies

## Development team engineers

Vũ Thị Thanh Thanh

Đỗ Thành Trung

Hoàng Trung Hiếu

Lê Tiến

# Contents

# Providing software development tools



TOSHIBA CORPORATION

Engineers (users)

Deploy tools

Our team

Use

Product development dept. at company A

Use

Product development dept. at company B

Use

R&D department at company B

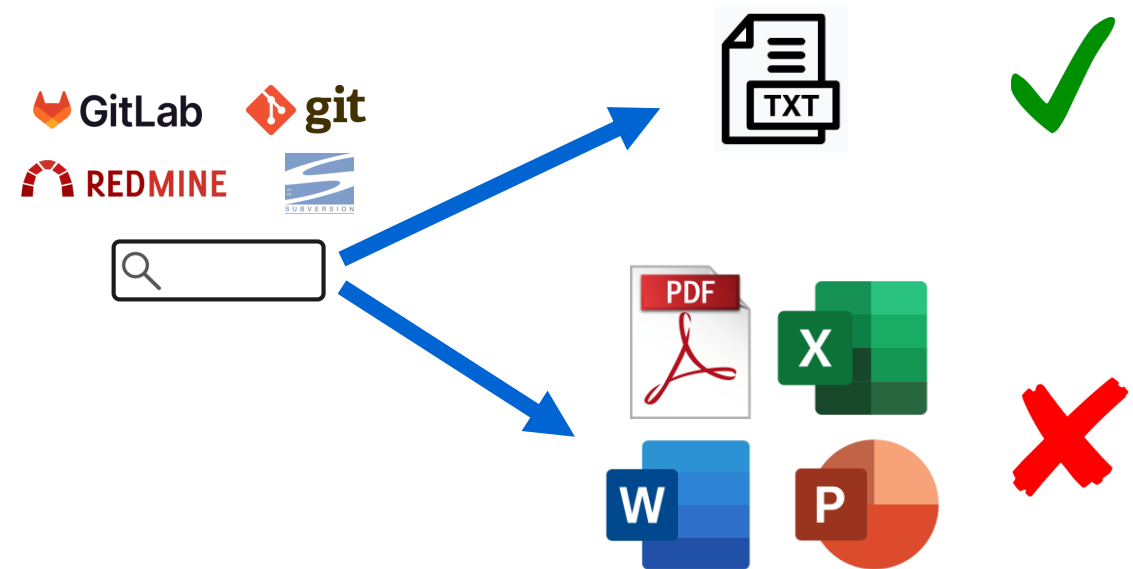about 200 deployments

# 2 major issues resulting from using multiple tools
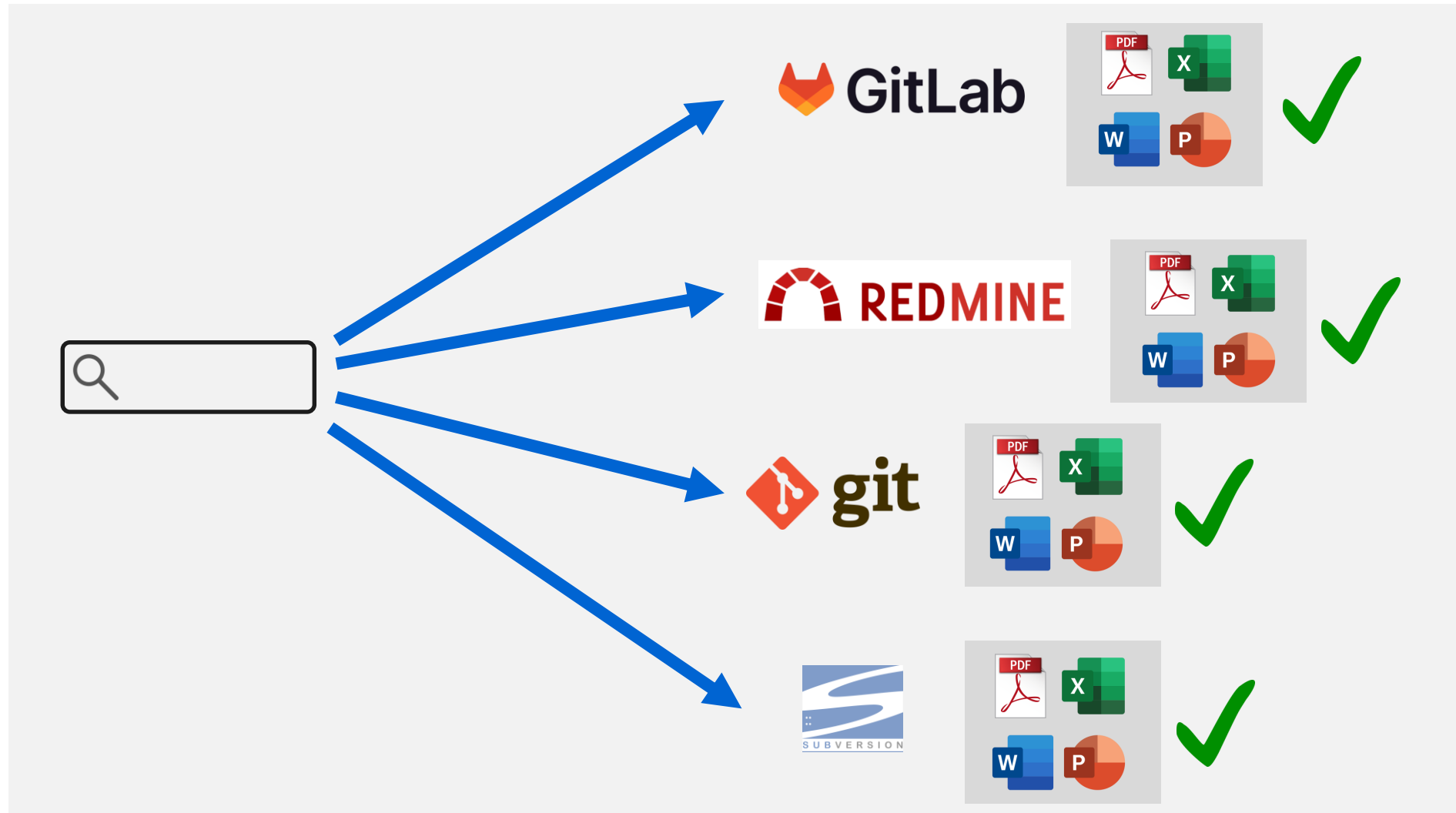
## Can't search across tools



Search

Resources, e.g. source code, documents, and wiki, are stored in multiple content management tools (fragmentation).

## Can't search texts inside binary files

# Goal

**One search box** to rule them all

# Approach

- Making a search engine by ourselves would be too expensive.
- Better to find an **OSS search server** and configure/customize it for our needs.

- Requirements:
  - Shall solve the aforementioned 2 problems.
  - Shall be complete with the essential features we need, e.g.
    - Search box UI for users
    - Features for admins (easily set up and run web crawlers, etc.).
  - *(A few more requirements; more on them later)*

- An OSS search server named **Fess** stood out as the closest to what we needed.
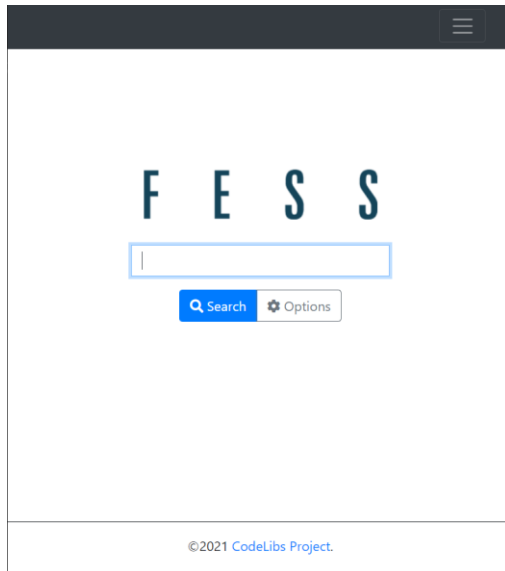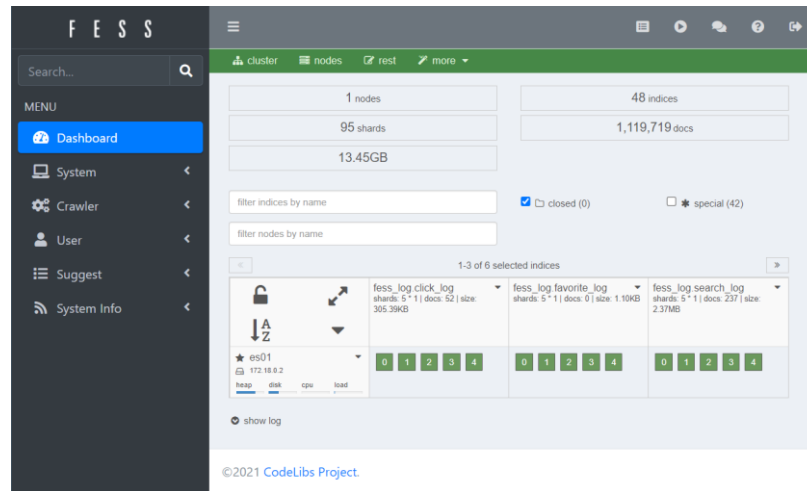
# Contents

# Fess

- *An enterprise search\* server*
- Developed by a company named CodeLibs
- Uses Elasticsearch as the search engine
    - Supports indexing of certain binary files, e.g. Microsoft Office, PDF, and zip.
- Comes with several types of crawlers, which can crawl
    - documents on a web server
    - file system
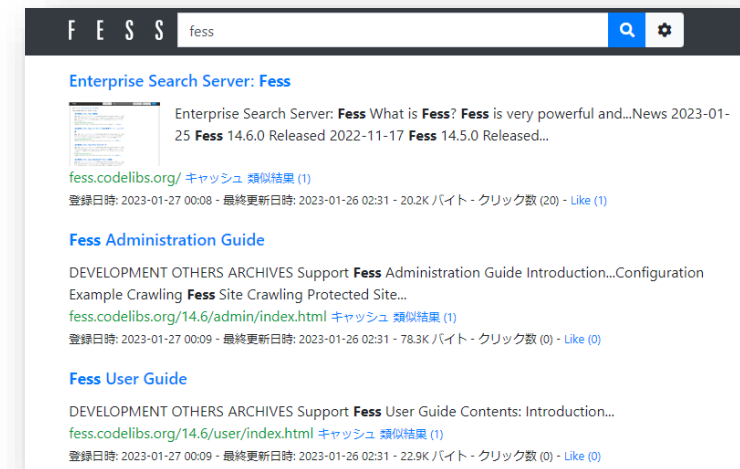    - data store (such as a CSV or database).

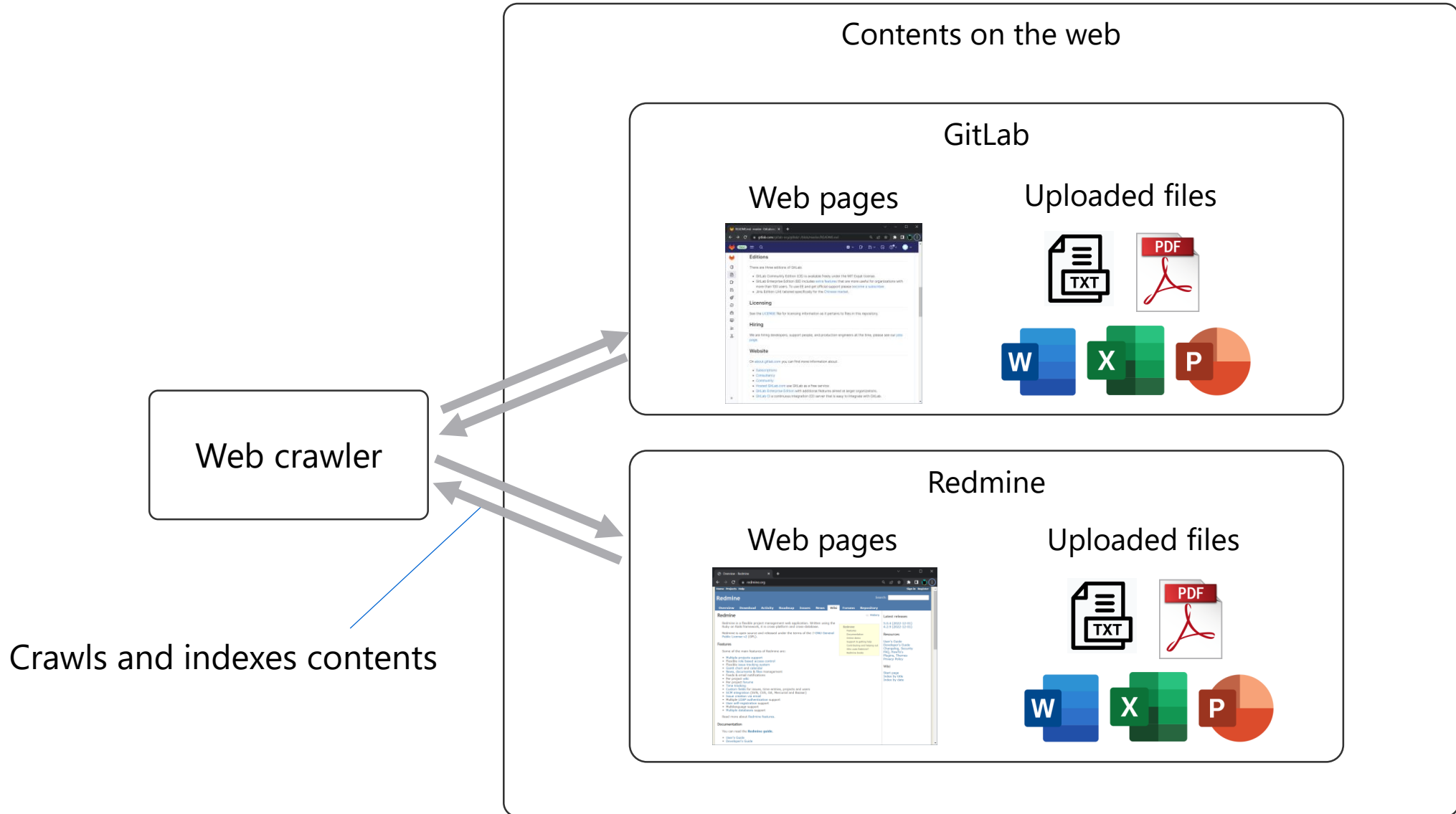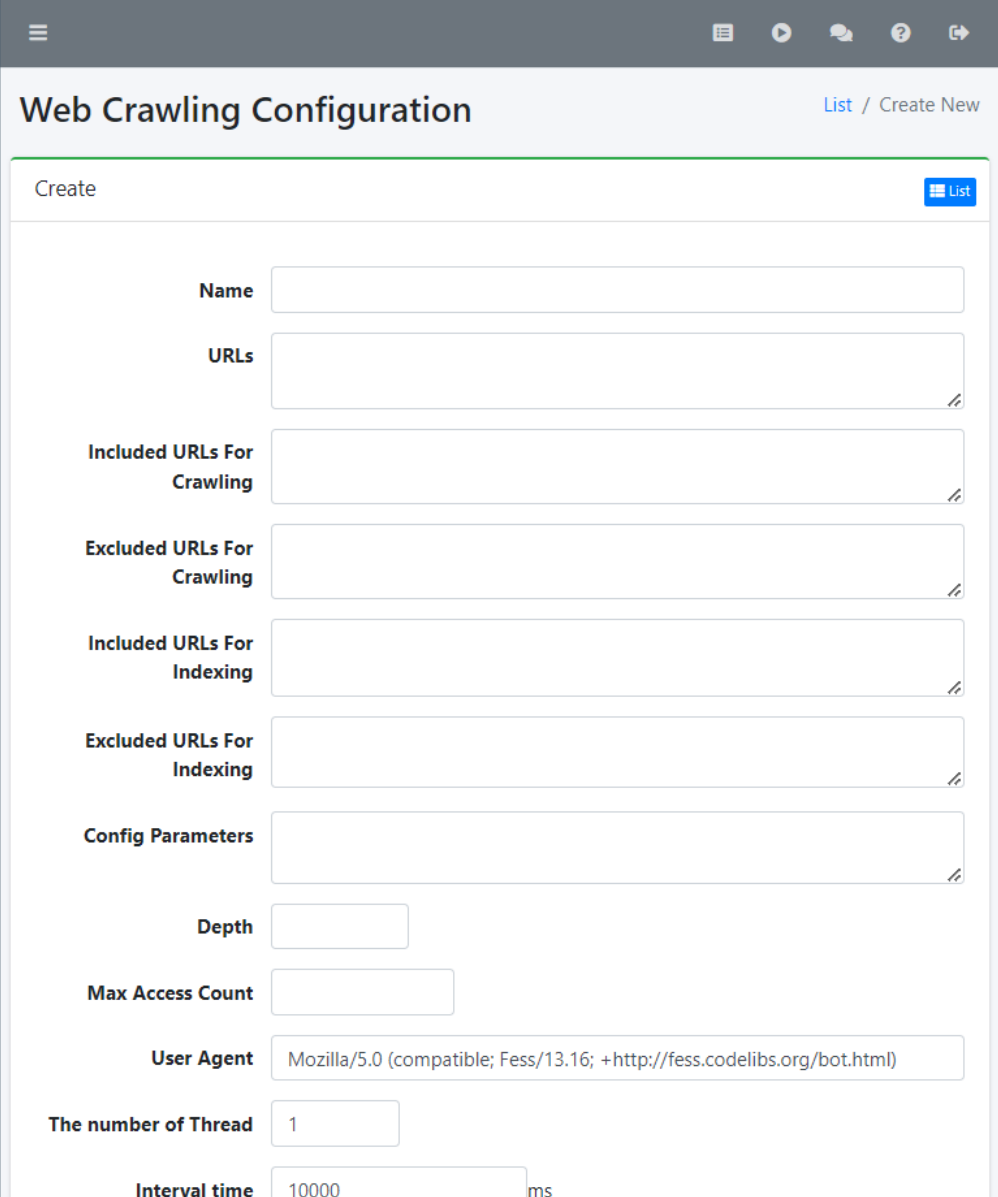https://github.com/codelibs/fess



Search box GUI



Admin console



Search engine results page (SERP)

# Fess web crawler



Web crawler

Crawls and indexes contents

Contents on the web

GitLab

Web pages

Uploaded files

Redmine

Web pages

Uploaded files

# Web crawler configuration basics (1 of 3)

- Crawls and indexes web page contents (texts, uploaded files, etc.) by following links recursively.
- Fess provides an admin console GUI to create a crawler.

## Web crawler parameters

| Parameter name | Description | Example |
|---|---|---|
| URLs | Starting point of the web crawling/indexing | https://mygitlab.io |
| Include URLs (regex) | Web crawler will crawl the page only if the URL matches any of the listed regex. | ^https://mygitlab\.io/myproject/-/issues/\d+$ <br> ^https://mygitlab\.io/myproject/-/wikis(?!.*/edit$).*$ |
| Exclude URLs (regex) | Prevent crawlers from indexing irrelevant pages. | ^https://mygitlab\.io/help($\|[/?].*$) <br> ^https://mygitlab\.io/(?!.*/-/.*/-/).+/-/new/.*$ |
| Depth | Maximum recursion depth | 1000 |
| Max Access Count | The maximum number of pages to crawl. | 10000000 |
| Permissions | Determines who can find the indexed contents | {group}mygroup <br> {user}user1 <br> {user}user2 <br> {user}user3 |

Permission settings and access control

| Parameter name | Description | Example |
|---|---|---|
| Permissions | Determines who can find the indexed contents | {group}mygroup<br>{user}user1<br>{user}user2<br>{user}user3 |

- Allows admin to implement per-user access control
- Permission settings are **per-web crawler**
    - e.g. 100 projects on GitLab → create 100 web crawlers, each of which has its own permission settings
- **username** in {user}**username** can be either
    - Users created on Fess
    - Users authenticated via LDAP directory service

# Contents

# Customization (patching Fess)

- Patches delineated in this segment:
  - Authentication for web crawlers
  - Crawling and indexing repository contents

• Patches merged upstream
• Patch submitted but not merged
• Patches kept proprietary

| | Patch | Link* / patched file(s) |
|---|---|---|
| 1 | Disable default PostConstruct | codelibs/fess-crawler/pull/56 |
| 2 | Make webdriver quit after it finishes crawling | codelibs/fess-crawler/pull/58 |
| 3 | (webdriver) add option for specifying chrome arguments | codelibs/fess-crawler/pull/57 |
| 4 | (webdriver) add option for chrome.prefs | codelibs/fess-crawler/pull/57 |
| 5 | (webdriver) add authentication | - |
| 6 | (webdriver) support crawl files | - |
| 7 | (webdriver) add basic auth | - |
| 8 | Add saml auth | - |
| 9 | map filesystem paths to URLs | - |
| 10 | (webdriver) customize last modified | - |
| 11 | (webdriver) prioritize fess cookies | - |
| 12 | (webdriver) add waiting configurations | - |

*paths of Link URLs following https://github.com/

# Authentication for web crawlers

- Web crawlers need authentication.
- Fess's **web authentication** supports certain types of login pages but not all.



Directory service (LDAP)

Web Authentication    List / Create New

Create

| Hostname | |
| Port | |
| Realm | |
| Scheme | Form |
| Username | user |
| Password | •••••••••••••• |

Parameters
```
encoding=UTF-8
token_method=GET
token_url=https://mygitlab.io/login
token_pattern=name="authenticity_token"[^>]+value="([^"]+)"
token_name=authenticity_token
login_method=POST
login_url=https://mygitlab.io/login
login_parameters=username=${username}&password=${password}
```

token_pattern=name="authenticity_t oken"[^>]+value="([^"]+)" token_name=authenticity_token

Keycloak/SAML

Not supported

- **We wrote a patch to authenticate the web crawler through our SAML-enabled GitLab sign-in page**
  - **Patched webdriver client and http form**
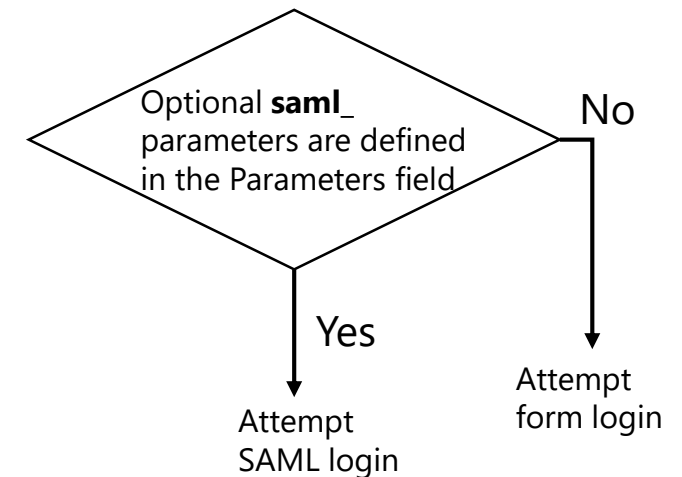
# Patching webdriver client and HTTP form

## New optional parameters on Fess admin console

Fess console:
(web authentication)

| Scheme | Form |
|---|---|
| Username | user |
| Password | ****** |
| Parameters | encoding=UTF-8<br>token_method=GET<br>token_url=https://tccloud2.toshiba.co.jp/ditc/gitlab/users/sign_in<br>token_pattern=name="authenticity_token" +value="([^\"]+)"<br>token_name=authenticity_token<br>login_method=POST<br>login_parameters=username=${username}&password=${password}<br>**saml_url=https://mygitlab.io/users/auth/saml**<br>**saml_login_url_pattern=action="([^"]+)"**<br>**saml_response_pattern=name="SAMLResponse" +value="([^"]+)"**<br>**saml_response_name=SAMLResponse**<br>**saml_callback_url=https://mygitlab.io/users/auth/saml/callback** |

The patched form parser checks for optional **saml_** parameters and stores them if they are present
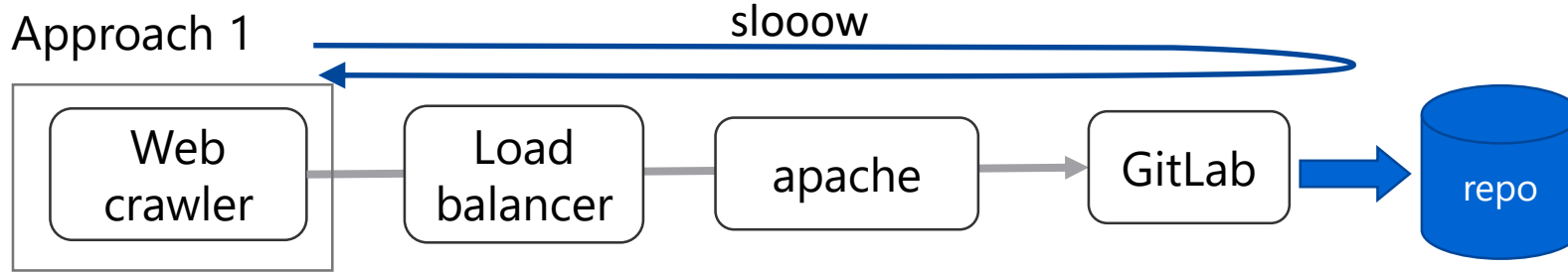
The patched web driver client attempts SAML login if optional parameters are defined

Optional **saml_** parameters are defined in the Parameters field

No → Attempt form login

Yes → Attempt SAML login
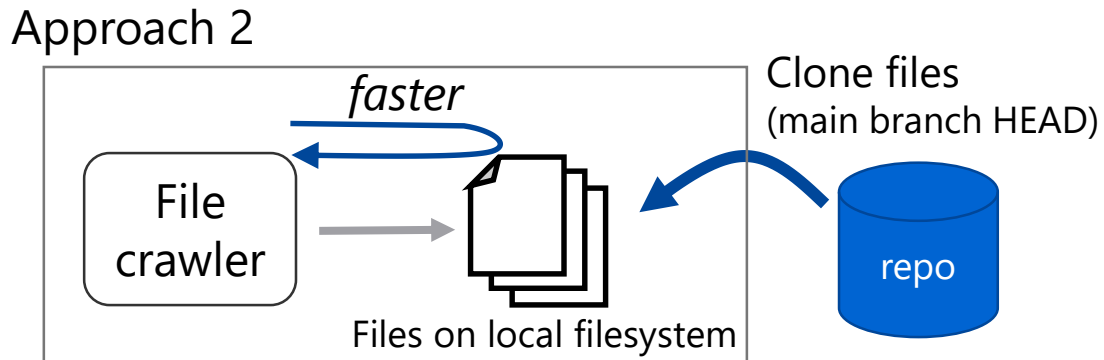
**Crawler supports SAML login**

# Crawling and indexing repository contents

- Many of our repositories are several GBs in size.
- Indexing the repo contents using web crawlers turned out to be too slow.
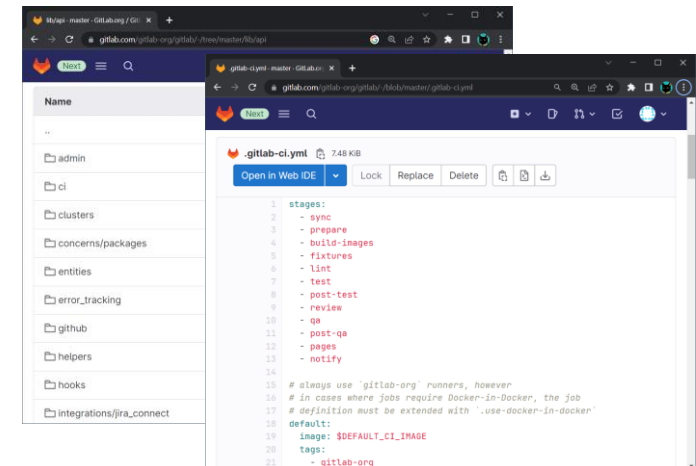
Approach 1

slooow



Web crawlers can index files in repo but it takes much longer
- HTTP request to GitLab → GitLab fetches repo file contents →
GitLab renders the file contents on the page*

*Repo files/folders as seen on the browser (GitLab)

Approach 2

Clone files
(main branch HEAD)

faster

Files on local filesystem

Open the file locally and read its content → much faster
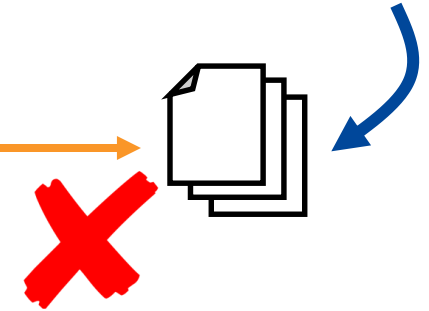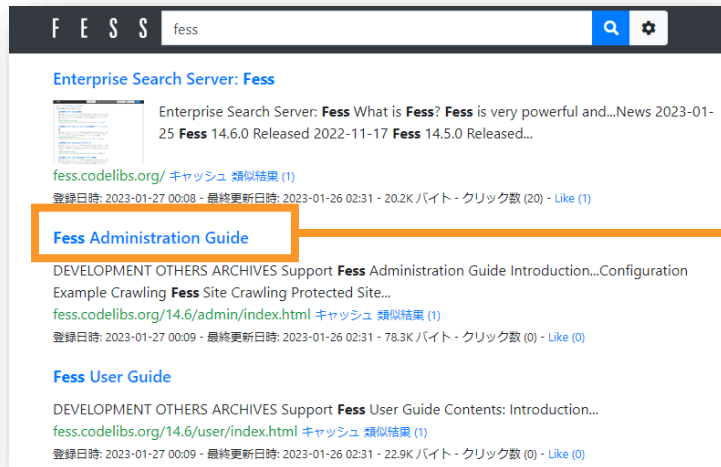
# Problem with approach 2

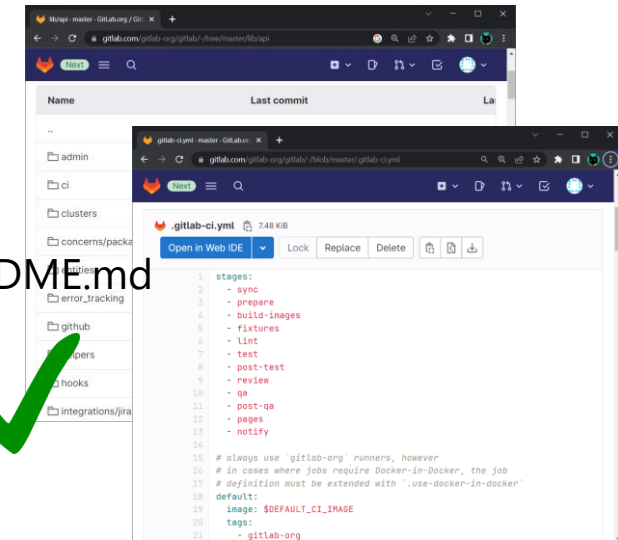File path pointed to the local file on the Fess instance

**file:///home/user**/repo/myproject/README.md

Patched web driver client re-maps the filesystem path to URL

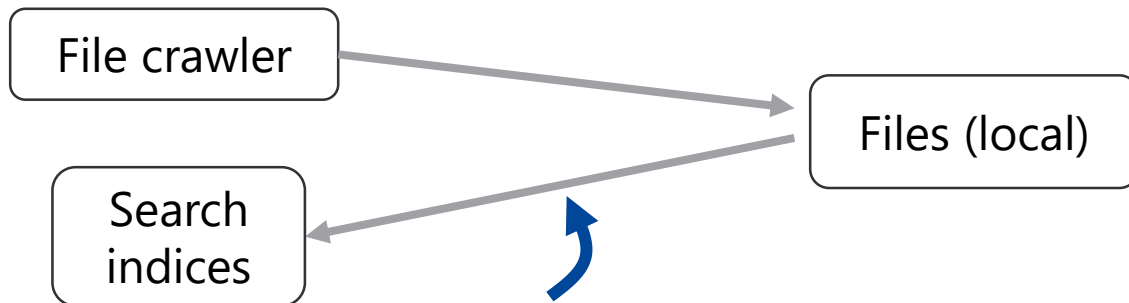**https://mygitlab.io/**myproject/-/blob/main/README.md

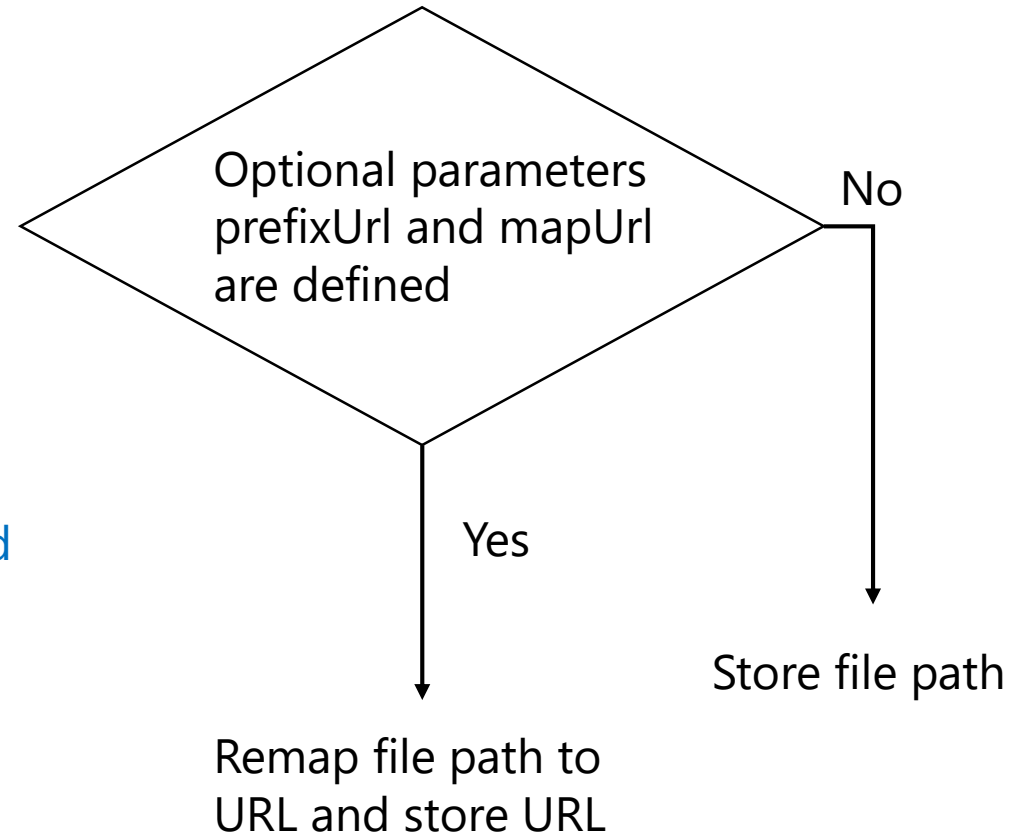# Patching file system client (re-mapping file system paths to URLs)

Fess admin console
(file crawling configuration)

| URLs | https://mygitlab.io |
|------|---------------------|
| Included URLs | https://mygitlab\.io/.+ |
| … | … |
| Config Parameters | **client.prefixUrl=/home/fess/repo/mygroup/**<br>**client.mapUrl="https://mygitlab.io/mygroup/" + project +**<br>**"/-/blob/main/" + relativeFilePath** |

The patched config parser checks for optional **client.prefixUrl** and **client.mapUrl** parameters and stores them if they are present

Optional parameters prefixUrl and mapUrl are defined

No

Yes

Store file path

Remap file path to URL and store URL

File crawler

Files (local)

Search indices

```
If( client.prefixUrl && client.mapUrl ) {
    url = remap(filepath)
  }
```

# Contents

# Automating configurations

- Manual configuration using GUI became impractical as the number of configurations increased.
- For each Fess instance,
    - More than 10 configurations
    - Several hundred web crawlers to create


- Without automation, the administrator would have to
    - manually edit lots of config files.
    - do repetitive yet complicated GUI operations to create web crawlers, web auth objects, schedulers, etc. on admin console.

# Why automate



Search results page

Project A
Project B

user1 — Search

user2 — Search → Project A

Search results should show only the resources (pages, files, etc.) in projects the user has access to.

## GitLab

**Project A**

members:
user1
user2

**Project B**

members:
user1
user3

### Web Crawling Configuration
List / Create New

Create                    List

**Name**

**URLs**

**Permissions**

{user}user1
{user}user2

### Web Crawling Configuration
List / Create New

Create                    List

**Name**

**URLs**

**Permissions**

{user}user1
{user}user3

Web crawler for project A          Web crawler for project B

Need to create as many web crawlers as the number of the GitLab projects. → …

# Automation example using Fess APIs

Creating web crawlers for each GitLab project:



```python
# Get all the projects on GitLab (GitLab API)
projects = gitlab.get_all_projects()

for project in projects:
  # Get all members of project
  users = gitlab.get_project_members(project)

  # Create web crawler for project (Fess API)
  webcrawler = fess.create_webcrawler(project.url, users, ...)

  # Create web authentication for the web crawler
  webauth = fess.create_webauth(project.url + '/login', webcrawler,
  username, password)

  # Create job scheduler for the web crawler
  ...
```
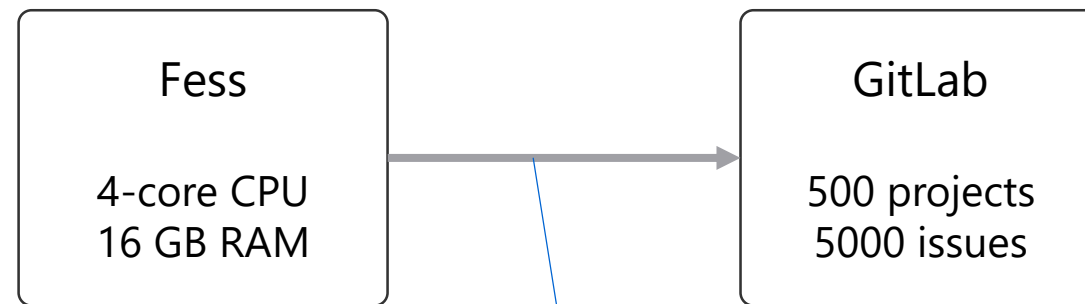
# Contents

- Did Fess solve our problems?
  - Can't search across tools
    - → Solved. Our users can now search content management tools laterally.
  - Can't search texts inside binary files
    - → Solved. Our users can find PDF and Office documents by texts contained in them.

- Web crawler performance
  - What our preliminary deployments revealed



| Fess<br><br>4-core CPU<br>16 GB RAM | → | GitLab<br><br>500 projects<br>5000 issues |

Takes about a couple of days
to index everything

# Committed to People, Committed to the Future.

• Patches merged upstream
• **Patches kept proprietary**

1. Fess's webdriver-based crawler client can index most of the dynamically loaded contents (Fess 13)
   - You need to change the crawler client from the default to a webdriver-based one, which requires editing config files (XML) and downloading and copying jar files to a specific directory (but no patching is required).
   - Web driver client tends to be less tested/mature; we fixed memory leaks and the PRs were accepted and merged upstream.

2. Our dev team wrote a patch to index contents on GitLab 15
   - The stock webdriver client didn't support issue pages on GitLab 15
   - **Patched web driver client to wait for the specified elements to appear on the page.**

> - Starting Fess 14, webdriver client is no longer supported.
> - A new crawler client powered by playwright is on its way for Fess 14, but no roadmap has been given by CodeLibs