# State of play of the Data Acquisition Service (DAS) of the Web Intelligence Hub (WIH)
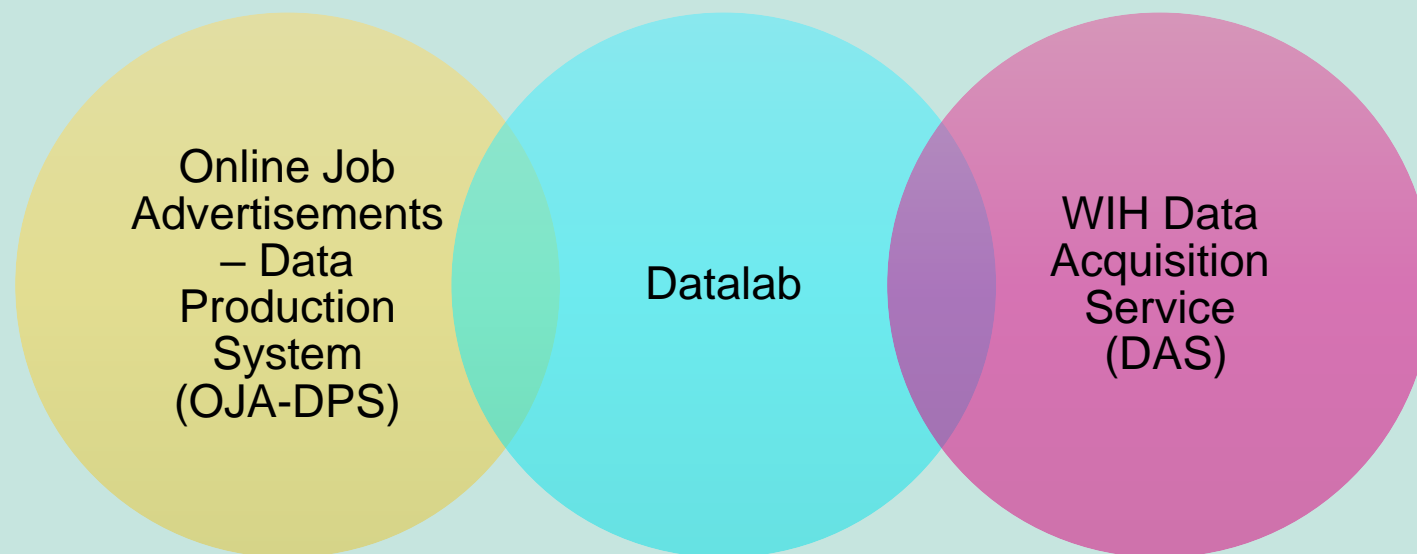
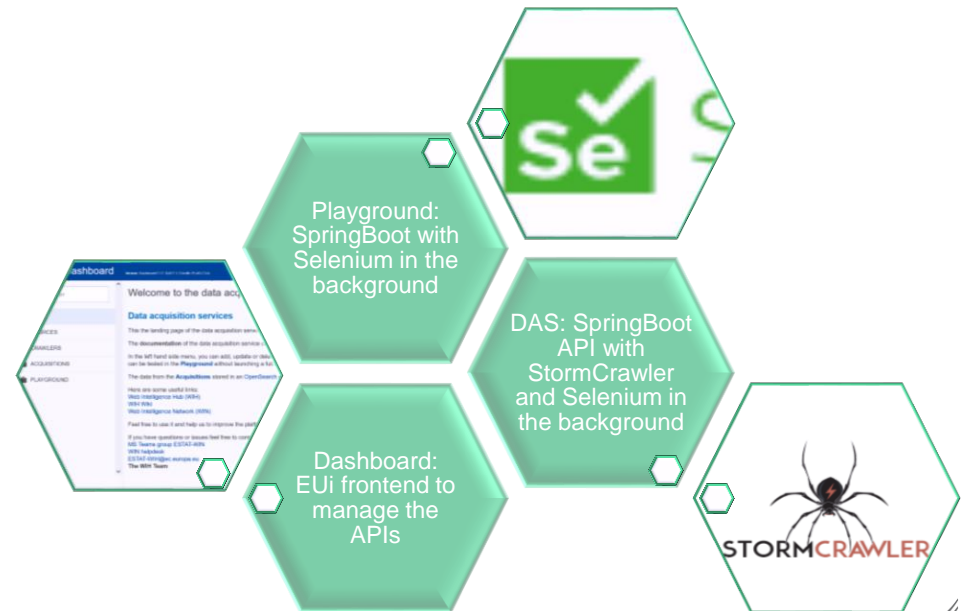Mészáros Mátyás

2025/02/04

European Commission

# DAS development principles

- Scalable
- Build on open-source tools
- Try to use the state of the art
- Can handle static and dynamic content
- No coding, only configuration by the user
- Universal, can be used for several use cases (OJA, MNE, price, etc.)
- Separation of use cases with possible collaboration in the same use case

# The beginning of the DAS

- **Version 1** was released in 2021 Nov
  - Generic data acquisition service with API access for static and dynamic web pages
  - Using StormCrawler and Selenium
  - EUi frontend (Dashboard) with user authentication (EU Login and AWS Cognito)
  - Deployment using Infrastructure as Code (IaC)

- **Version 2** was released in 2022 Apr
  - Adding the playground for data acquisition to test filters and dynamic web pages

- **Version 3** was released in 2022 Sep
  - Additional Selenium filters based on the needs of tourism websites

# First testing by the WIN

Based on the feedback

- **Version 4** was released in 2022 Dec
  - Multitenancy was introduced
  - Moving authentication from AWS Cognito to Keycloak

- **Version 5** was released in 2023 Feb
  - 1st Security testing and update
  - Adding new functionalities like advanced search, copy configuration and acquisition action history

- **Version 6** was released in 2023 Mar
  - Introduction of user roles (guest, developer and admin)
  - Possibility to use sitemap discovery

# Continuous improvement of the DAS

**Version 7** was released in 2023 Oct
- 2nd Security testing and fixing issues
- Handling large number of sources for the OBEC use case
- Moving from Elastic to OpenSearch
- Automatic stopping of acquisitions
- Verifying bot IPs

**Version 8** was released in 2024 May
- 3rd Security testing and update
- Exporting data to S3 and access it through API
- Scheduling of acquisition
- Automatic resizing of resources

# Latest version

**Version 9** was released in 2025 Jan

- Introduction of URLfrontier and incremental crawling
- Improved scheduling
- Availability of crawler history
- Possibility to store only extracted text without source html
- No parallel run of the same crawler

# Current state

# The way forward

- 4th security testing and corrections
- Migration of OJA-DPS crawlers on the DAS
- Updating the components
- Possible use by the new individual grants

Frontier for : GUS_STAT_JOBS

Crawler has **46** Sources Frontier

Q Search

| | Next Fetch Date ⇅ |
|---|---|
| | 03/02/2025 22:14:20 |
| | Fetched |
| | Fetched |
| ...dziale-badan-i-strategii-regionalnych-w-dep... | Fetched |
| ...dziale-budzetu-i-analiz-w-departamencie-adm... | Fetched |
| ...dziale-czasopism-naukowych-w-departamencie-... | Fetched |
| ...dziale-finansowo-ksiegowym-w-departamencie-... | Fetched |
| ...dziale-metodologii-badan-statystycznych-w-d... | Fetched |
| ...dziale-projektow-dofinansowywanych-ze-srodk... | Fetched |
| ...dziale-projektow-dofinansowywanych-ze-srodk... | Fetched |
| https://bip.stat.gov.pl/ogloszenia/wolne-stanowiska-pracy/glowny-urzad-statystyczny/glowny-specjalista-w-wydziale-rozwoju-kompetencji-innowacyjnych-w-... | Fetched |
| https://bip.stat.gov.pl/ogloszenia/wolne-stanowiska-pracy/glowny-urzad-statystyczny/glowny-specjalista-w-wydziale-rozwoju-kompetencji-innowacyjnych-w-... | Fetched |
| https://bip.stat.gov.pl/ogloszenia/wolne-stanowiska-pracy/glowny-urzad-statystyczny/glowny-specjalista-w-wydziale-zarzadzania-bezpieczenstwem-informac... | Fetched |
| https://bip.stat.gov.pl/ogloszenia/wolne-stanowiska-pracy/glowny-urzad-statystyczny/glowny-specjalista-w-wydziale-zarzadzania-bezpieczenstwem-informac... | Fetched |

*Questions? Comments?*

# Thank you