# RA21 Academic Pilot Technical Evaluation

Editor: H. Flanagan, Academic Pilot Coordinator

# Executive Summary

Both RA21 pilots were successful in testing technical approaches to identity provider persistence, and we learned a great deal from both. Now that we have determined that taking RA21 forward will require the establishment and operation of at least some centralized infrastructure, we want to focus on just one option. After a thorough evaluation, we have chosen to move forward with the technical architecture prototyped by the P3W pilot, thus closing the WAYF Cloud. We are focusing on P3W principally because it minimizes the amount of data that is held in a single place, thus minimizing the potential for security or privacy breaches. It also has a lighter technical footprint for the central infrastructure and this is more attractive to potential operations partners.

A core capability of RA21 is persistence of the identity provider choice made by users. We believe a central identity provider discovery service will also be highly valuable, particularly

to offer an easy path to implementation for the "long tail" of smaller service providers, however it is not absolutely essential to the RA21 proposition. So, we propose to put forward a technical recommendation for identity provider persistence, but only best practices for how to present identity provider discovery options.

# Background

## What is RA21?

The RA21 project has been working to remove some of the barriers blocking the move to the use of federated identity for the management of access to scholarly information resources by testing new user experience concepts and flows through several pilot projects. Once the pilot phase is completed, RA21 intends to publish recommendations for best practice for future implementations via the NISO Recommended Practice process.

The project has identified the current identity provider (IdP) discovery process as a difficult-to-navigate step which currently drives many users away from authenticating using their home institution's credentials when accessing scholarly information resources despite the large installed base of SAML IdPs at relevant institutions. As a result, IP-based access control persists as the predominant mechanism of access control for these resources when accessing from with an institution's network, while an inconsistent patchwork of alternative solutions such as local service provider accounts with various provisioning mechanisms and URL-rewriting proxies have been employed to facilitate access to users when off network. These mechanisms present a large maintenance overhead for IdPs and SPs and provide for a confusing and inconsistent user experience.

RA21 began in 2016, and over the course of the project the following elements have emerged as the most promising components of an improved solution:

- A **common UI element** that participating service providers may deploy to their sites directing users to initiate the IdP discovery process
- An **improved, search-based IdP discovery interface** which makes use of enhanced IdP metadata to enable reliable selection of the appropriate IdP using institution name or email domain
- A **centralised IdP persistence service** which enables a user's previous choice of IdP to be remembered by their browser across participating service providers, thus decreasing the frequency with which the user has to choose their IdP

The future RA21 service will encompass all three of these elements while providing the highest possible protection for user privacy, and respecting all current best practices for information security.

## Who will use the service?

There will be three principal groups of users:

- Users at participating IdP organizations will interact with the service when making their IdP selection and having it remembered by their browsers.
- Participating IdP organizations will follow metadata guidelines issued by RA21 and train and educate their users about use of the service. The set of participating IdPs is envisaged to include the current set of principally academic identity federations as represented in edugain, as well as corporations who make use of scholarly information resources, either individually, or through yet-to-be-established corporate federations.
- Participating Service Providers (SPs) will incorporate the WAYF persistence service, and optionally the common UI elements, into their services, and may also make use of a centralised IdP discovery service. The set of participating SPs will include publishers, content aggregators, platform vendors, tool vendors and service providers to participating IdPs. They will range from volunteer publishers of individual publications to large enterprises, and will be drawn from the academic, non-profit and for-profit sectors.

RA21 is international in scope, and users, IdPs and SPs will come from around the world.

# Academic Pilot Descriptions

RA21's goal is to facilitate as seamless as possible a user experience while surmounting the technical and security limitations of IP address recognition, and strongly supporting network security and user privacy.

The user experience with federated identity can be broken down into two steps: the user's experience in initially selecting their Identity Provider, known as IdP discovery, and the persistence of that choice for future sessions involving federated identity, known as IdP persistence. The evaluation of the two academic pilots focused primarily on IdP persistence; how IdP discovery is implemented will vary more broadly among Service Providers, and RA21 has chosen to focus on offering best practices in the identity discovery space as opposed to proposing a ubiquitous central service for that component.

A diagram[1] of the data flows for both pilots is included below.

## Privacy Preserving Persistent WAYF (P3W)

The P3W pilot supports two models of integration. Level 1, the most basic integration model assumes that the Service Provider — generally a publisher in the RA21 use case — wants to completely externalize its federated identity discovery services. As such, it would use a common URL to point to a central discovery service that would allow the user to choose among a list of possible IdPs, and then record that choice in a user's browser so that it is available for future sessions.

In the more advanced level 2 scenario, the Service Provider would use a local IdP discovery service that could include any local accounts hosted by the SP, and use the central service only to store the user's choice of IdP in the browser. This would be accomplished by calling an API within the browser provided by JavaScript hosted by the central service on a trusted domain so that participating services providers can all access the same shared set of remembered IdPs.

Note that the P3W architecture only supports storing the user's choice (or choices) of IdP in their browser, no usernames, passwords or other personally identifying information is stored. If a user uses a private browsing mode, any choices made will not be stored after that browser windows is closed.
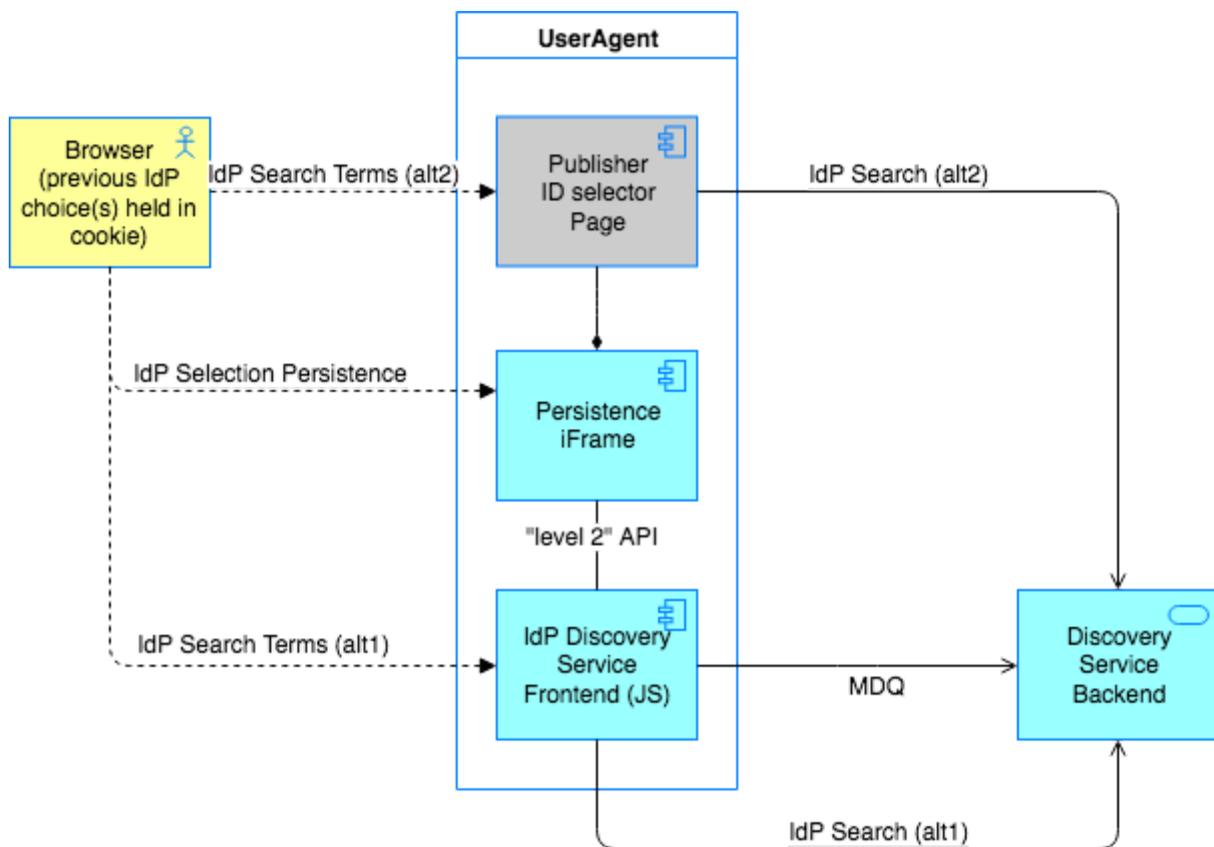


Figure 1: P3W Data Flow Diagram

## WAYF Cloud

The WAYF Cloud service assumes that IdP discovery is handled by the Service Provider. The focus of this service is to persist that choice such that any participating Service Provider can use the WAYF Cloud to present the user's choice back to the user. The WAYF Cloud does this by using both JavaScript and back-end API calls to store a mapping between each

service provider's local unique device identifier and a common global unique device identifier in a central database. The central database then maps the global shared device IDs to the IdPs the user has successfully logged into.

Note that the WAYF Cloud architecture does not store username or passwords or other personally identifying information. The data trail will exist during the life of the session when the user is in Incognito mode, and in the browser until the cache is cleared.
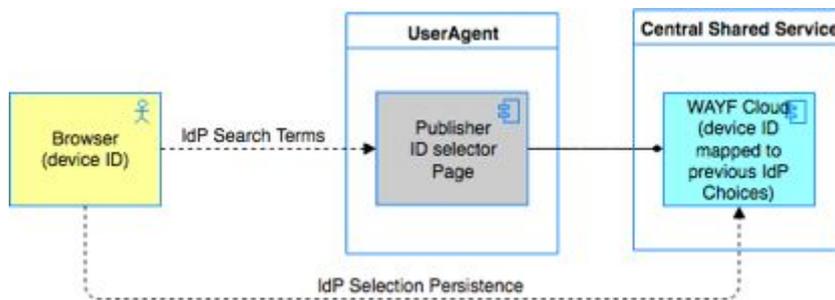


Figure 2: WAYF Cloud Data Flow Diagram

# Evaluation Process

The evaluation process for the pilots included an in-depth security analysis, a privacy review, and a detailed technical architecture comparison.

## Security Review

The security review followed the Microsoft STRIDE Threat Classification Model, commonly used for analyzing threats from six standard categories:

· **S**poofing of user identity
· **T**ampering
· **R**epudiation
· **I**nformation disclosure (privacy breach or data leak)
· **D**enial of service
· **E**levation of privilege

A detailed report describing the full results of this analysis is available on the RA21 website under Results. In summary, however, the security analysts determined that both pilots are very low risk. Both pilots would need to follow best common security practices if deployed, such as regular auditing and penetration testing, security of the server and API components, etc.

## Privacy Review

For data privacy risks, a data protection impact assessment ("DPIA") was performed to analyze potential issues based on GDPR requirements. The DPIA is an analysis of expected processing activities related to assessments and covers details of the processing activity itself and an assessment of the risks associated with the processing activities including any measures that need to be taken to mitigate those risks.

Full detail of the privacy review is included in the [Security & Privacy report](#). In summary, however, the analysis indicated that the data collected by either pilot during the discovery and persistence processes is of very low value from a privacy perspective. A key difference however is that in P3W architecture, the IdP persistence information for any given device is stored locally within the browser to which is pertains, whereas in the WAYF Cloud model, the data for all devices for which a preference has been set is stored in a centrally hosted database.

## Service Architecture

The technical architecture required to support a service offering for either of the pilots is fairly small in terms of infrastructure, but there are key differentiators.

The WAYF Cloud has a heavier footprint in terms of the ask for the infrastructure operator in that it works on a more traditional client-server model. The infrastructure would require a high-availability database that holds the device-specific keys and IdP mappings, and appropriate processes for an entity to handle the legal overhead of reviewing, asserting, and maintaining GDPR compliance. While there are benefits of trained staff managing such a service, there are also concerns about the level of support required.

The P3W is a much lighter service from an infrastructure perspective. Any central service offering would only be required to make the necessary JavaScript API highly available. If the central service host also offered a discovery service, maintenance of any necessary metadata query tools would also be required.

In both WAYF Cloud and P3W, a central service provider would need to maintain and offer the appropriate security audits going forward.

# Academic Pilot Findings

Overall, the P3W and WAYF Cloud technologies offer similar functionality for supporting IdP persistence. Key differentiators for the services included:

- The minimum data required to make the service work
  - WAYF Cloud
    - Publisher local device IDs
    - WAYF Cloud service global device IDs
    - IdP preferences (represented as IdP entity IDs) for every user in the database in a central location

- ○ P3W
  - ■ Identity Provider choices (represented as IdP entity IDs), for one device, in browser local storage
  - ■ Data needed to secure access to the Javascript API

- ● Identified (minor) risks
  - ○ WAYF Cloud
    - ■ There is a potential for information leakage via the referrer logs of the central infrastructure.
    - ■ The organization that runs the service (or a man in the middle) has access to the data.
    - ■ WAYF Cloud still leaves a short data trail that exists for the length of the browser session of what publishers were accessed, even when user is Incognito mode.
  - ○ P3W
    - ■ Depends on some mechanism to secure the JavaScript API; failure to secure the JavaScript API may result in third party access and manipulation of the IdP choice information stored in browser local storage.

In both cases, the JavaScript (WAYF Cloud) and JavaScript API (P3W) must be secured in order to consider the overall service secure.

The primary differentiator in the pilots is where the data regarding the user's choice of IdP is stored. As mentioned, in the WAYF Cloud, that information is stored in a central database. In P3W, that information is stored in a user's browser. While the value of the information in the central database is very low from a security and privacy perspective, given another option that does not require the consolidation of such data, the choice was made to go with P3W. The P3W architecture offers a functionally equivalent service to the WAYF Cloud architecture, but with no central collection of information thus adhering to the privacy principle of data minimisation.  On the other hand, the WAYF Cloud architecture would obligate the party hosting the central infrastructure to define data retention policies and be subject to GDPR data access requests.

# Acknowledgements

- Dan Ayala (Proquest)
- Meltem Dincer (Wiley)
- Ken Ferris (Taylor and Francis)
- Peter Reid (Bath Spa University)
- Todd Carpenter (NISO)
- Christian Pruvost (Elsevier)
- Andy Sanford (EBSCO)
- Ralph Youngen (ACS)
- Heather Ruland Staines (Hypothes.is)
- Adam Snook (OpenAthens)
- Richard Northover (Elsevier)
- Chris Shillum (Elsevier)
- Joe Greene (CAS)
- Jos Westerbeke (University of Rotterdam)
- Will Simpson (ORCID)
- Paul Dixon (LibLynx)
- Phil Leahy (OpenAthens)
- David Orr (OpenAthens)
- Diane Cogan (Ringgold)
- Laird Barrett (Springer Nature)