

# CERN-Solid Code Investigation

Proof of Concept and Prospects

by

**Jan Schill**

schi@itu.dk

Supervisors:

**Philippe Bonnet** (ITU)

phbo@itu.dk

**Maria Dimou** (CERN)

maria.dimou@cern.ch

A thesis presented for the degree of  
Master of Science

IT UNIVERSITY OF COPENHAGEN

Computer Science  
IT University of Copenhagen  
Denmark  
01.06.2021

# Table of Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1	Context	3
2	Goal	3
<b>2</b>	<b>Related Work</b>	<b>4</b>
1	Background	4
2	Indico	4
2.1	Events	4
2.2	Storage Mechanisms	4
2.3	Conference Registration	4
3	Solid	4
3.1	Authentication With Solid	4
3.2	Reading and Writing Linked Data	4
3.3	Authorization Through WAC	4
3.4	Application Launcher	4
<b>3</b>	<b>Investigation</b>	<b>5</b>
1	Proof of Concepts	5
1.1	POC 1: Commenting Module for Events in Indico	5
1.1.1	Architectural Analysis and Synthesis	5
1.1.2	Screen Design	7
1.1.3	Design	7
1.1.4	Integration with Indico	12
1.1.5	Evaluation	12
1.1.6	Analysis	12
1.2	POC 2: Auto-Complete for Conference Registration in Indico	12
1.2.1	Architectural Analysis and Synthesis	12
1.2.2	Design	13
1.2.3	Integration With Indico	13
1.2.4	Evaluation	14
1.2.5	Analysis	14
1.3	Deployment of Indico Instance	14
2	Comparison of Solid and Indico Design Principles	14
3	Challenges, Advantages, and Gaps of Existing Solid Solutions versus CERN Ones	14
4	Proceedings in the CERN-Solid Collaboration	14
4.1	Servers	14
4.1.1	Solution	14
4.1.2	Hosting	14
4.2	Applications	14
<b>4</b>	<b>Conclusion</b>	<b>15</b>

# Introduction

- 1 Context
- 2 Goal

# Related Work

## 1 Background

## 2 Indico

### 2.1 Events

### 2.2 Storage Mechanisms

### 2.3 Conference Registration

## 3 Solid

### 3.1 Authentication With Solid

### 3.2 Reading and Writing Linked Data

### 3.3 Authorization Through WAC

### 3.4 Application Launcher

# Investigation

## 1 Proof of Concepts

One main part of the investigation into the CERN-Solid collaboration is the development of a proof of concept (POC). The POC contains the creation of two independent software modules in an existing system from European Organization for Nuclear Research (CERN). These software modules should show how it is to develop with the Solid principles in mind and to the Solid standard.

The goal of these modules is the symbiosis of decentralized stored data in a highly functional system without comprising its performance, security, or usability.

### 1.1 POC 1: Commenting Module for Events in Indico

The first POC is supposed to enrich the Indico system with some sort of Solid-based content. With the product owner and chief developer of Indico, the CERN-Solid project manager and a Solid developer it was decided a commenting module for Indico events is an adequate solution to include data from an external storage entity namely a data pod. The ability to allow users of Indico to leave a comment on an event, which then lives in a data pod completely controlled by the author of the comment was concluded to be an attractive feature for Indico.

#### 1.1.1 Architectural Analysis and Synthesis

In this section the architectural analysis will be looked into. The scope of the system and its attributes will be defined based on the system description and the quality attributes from the analysis of stakeholder needs.

#### *System Description*

The system aims at enabling commenting in web application with decentralized storage on data pods. The system in the context of Indico intends to allow Indico users with a data pod to comment on specific Indico events.

#### *Features*

The key features – but not limited to – for the module to succeed are the following:

1. A user with a Solid account can authenticate with their Solid identity provider (IDP)
2. A user can compose a text which is stored on their data pod
3. A user can see other users' comments

#### *User of the System*

There are two types of users in the system. The first group of users are the active authors of comments or observers of such. These will interact with the module by logging into their data pod and comment. The other subset of users are the administrator of the application (Indico). This type of user is interested in keeping the tone of comments to their desire and not allow any misuse.

#### *Context Diagram*

Other users of the system are highlighted in 1, which includes developers of either the system/module itself or the internal developer maintaining the application.

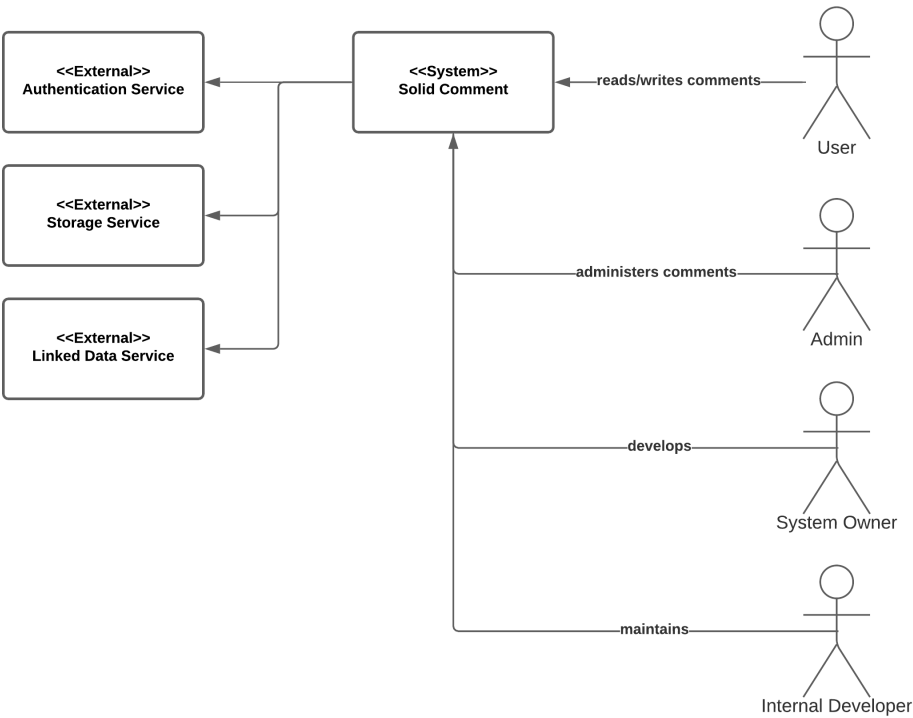


Fig. 1: Context diagram showing users and external services of system.

Sequence Diagram

A sequence diagram is suitable overview for any software architecture, but especially useful for decentralized or systems containing several separate services. It gives a clear understanding of each service’s tasks and their relationship when passing messages around in the overall system.

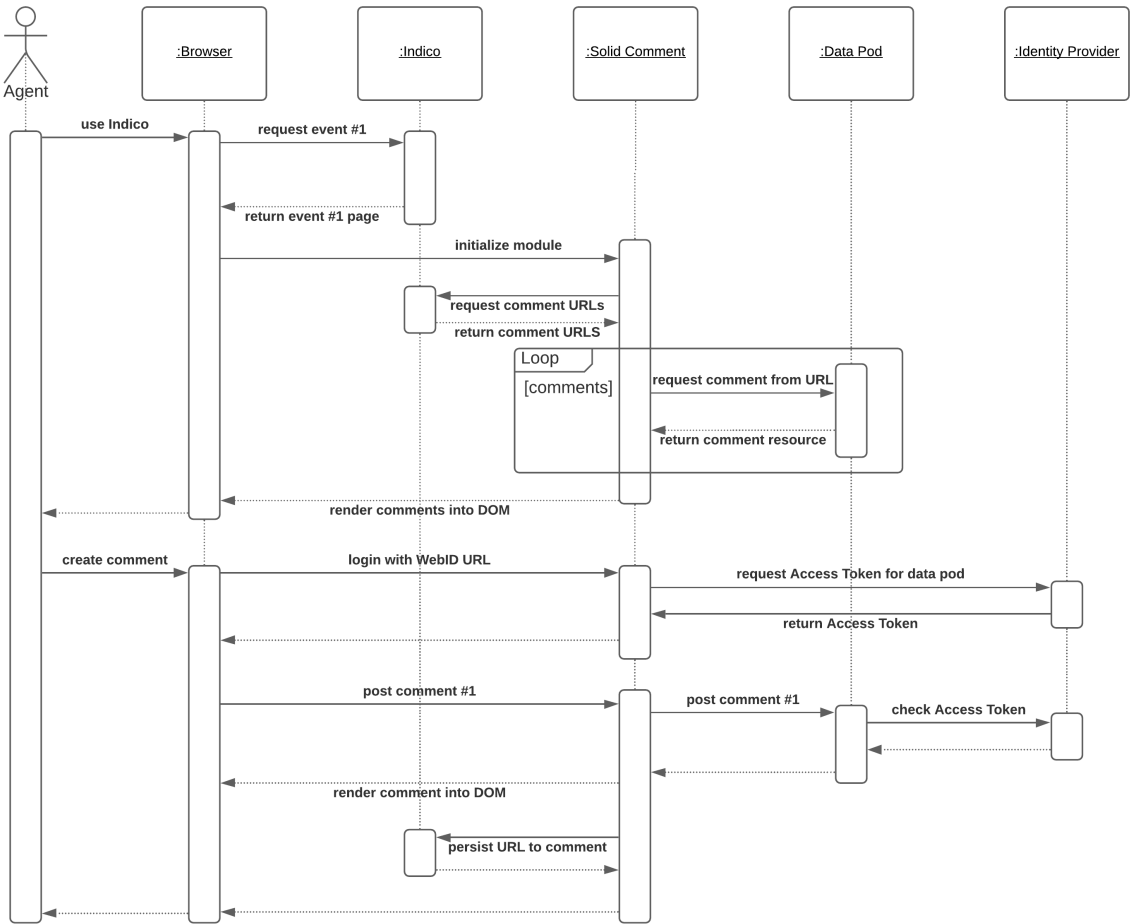


Fig. 2: Sequence diagram showing the sequential process through posting a comment.

Stakeholders

This section covers the various stakeholders in relation to the system. This both includes active users, but also various external and internal stakeholders who are impacted by the system or have an important say in the development process.

## Drivers

### 1.1.2 Screen Design

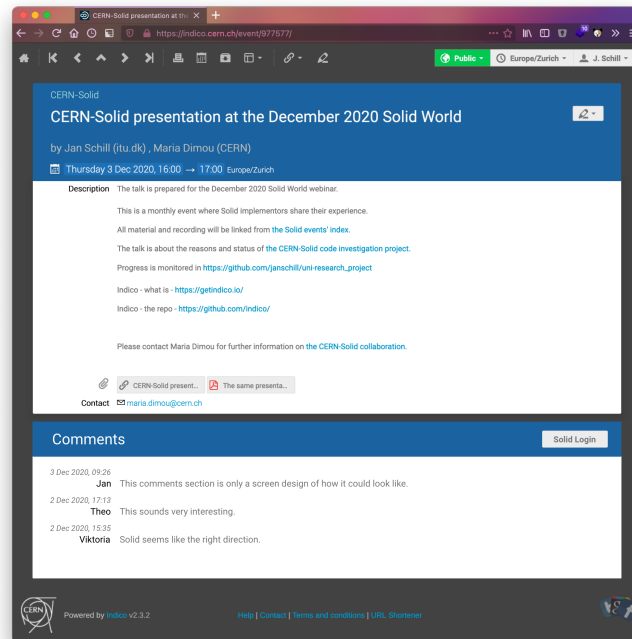


Fig. 3: User interface showing the comment module.

### 1.1.3 Design

For the implementation of this module several design decisions had to be made. From the fundamental choice of the module running on the client device or be computed on the server and then propagated to the client afterwards or even with a microservice proxying all traffic through it to enable Solid without changing Indico. Other design challenges were around how to protect the resources holding the comment information. These resources reside on the external data pod and need to be fetched from the application and read by other agents. Can access control lists (ACLs) be configured to allow the specific use-case?

#### *Client- Versus Server-Side Versus Microservice*

When an agent browses to a running instance of Indico most of the functionality is being prepared on the server hosting Indico. It retrieves the specific request, builds the Hypertext Markup Language (HTML), and sends it to the user. For Indico most of the functionality is built with Python and the web framework Flask. Sometimes functionality needs to be closer to the user, an example is dynamic rendering of Document Object Model (DOM) elements. This is useful when new data needs to be shown right away without getting the blank white screen on a page reload. Indico does send JavaScript (JS), which is used for client-side features, but it focuses on keeping most its features on the server.

To make the right decision if the module should be primarily developed for the client- or server-side or even as a microservice, a list of requirements to the module had to be defined. With the defined requirements in place it had to be figured out how much functionality can be extracted from existing libraries and how much needed to be implemented with the new module. Implementing existing functionality for a new programming language would defeat the POC's purpose of showing how an existing software could work with the Solid principles.

The rudimentary set of features to enable commenting for users in Indico while saving the data in a data pod includes:

1. Authentication with a Solid IDP
2. (Authenticated) Requests to a data pod
3. Parsing of structured data (Linked Data)

#### *Client Approach*

The module runs in the browser and is therefore written in JS. A programming language which compiles to JS, such as TypeScript (TS), is also possible. This means Indico remains mostly untouched, but would have to serve the needed JS to the client on traffic to an event endpoint where the comment module is integrated.

The communication flow with the data pod and the module would happen primarily from the browser.

Problem	Solution
Language	JS or TS
Framework	Native JS
Client	solid-client-js [1]
Authentication	solid-client-authn-browser [2]
RDF	solid-common-vocab-js [3], rdflib.js [4]

Table 1: Existing solutions to problems for a client approach.

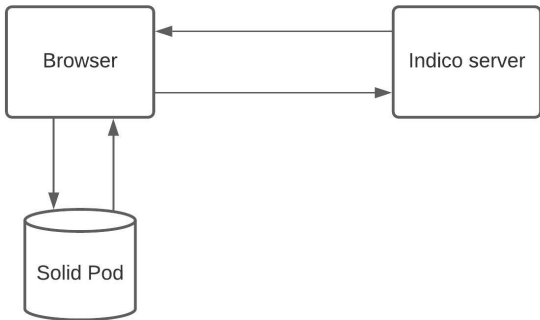


Fig. 4: Communication flow for a module developed on the client.

*Microservice Approach*

The microservice approach would allow developing the needed Solid logic on a separate service, which proxies all Solid related traffic through it an enable the Solid functionality. Most of the libraries from the client implementation can be used as well, as both developments would be written in JS. Only the authentication flow would work a bit different.

Problem	Solution
Language	JS or TS
Framework	Node.js
Client	solid-client-js [1]
Authentication	solid-client-authn-node [5]
RDF	solid-common-vocab-js [3], rdflib.js [4]

Table 2: Existing solutions to problems for a microservice approach.

The microservice module would handle take all requests aimed at the data pod and make it compliant with the Solid server. It would provide the client with the proper Solid OpenID Connect (Solid OIDC) flow to attach the access token to all authenticated requests.



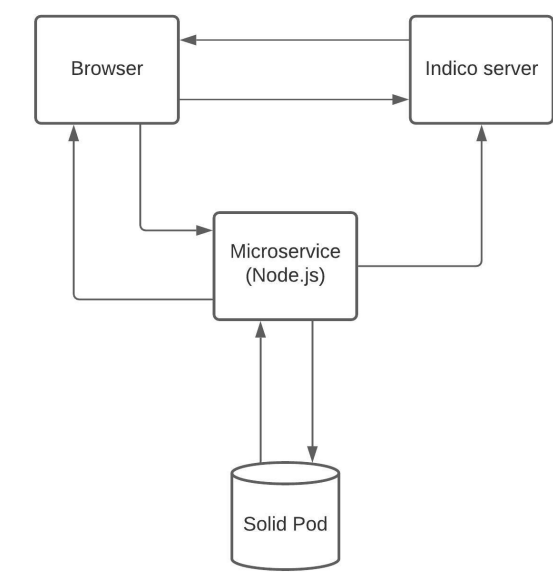


Fig. 5: Communication flow for a module developed as a microservice.

*Server Approach*

Goal of the server approach would be just like with the microservice approach to decouple the logic needed to work with Solid from the client and have it run on a server instance. The attractiveness for the server approach would be it could be fully integrated within Indico and be part of its Python code base. The major drawbacks are no direct Solid libraries written in Python exist to allow a seamless integration into the ecosystem.

Problem	Solution
Language	Python
Framework	Flask
Client	-
Authentication	pyoidc [6] missing DPoP
RDF	solid-common-vocab-js [3], rdflib.js [4]

Table 3: Existing solutions to problems for a server approach.

The authentication library pyoidc allows authenticating with OpenID Connect (OIDC) systems, but is missing a mandatory feature called Demonstrating Proof-of-Possession (DPoP), which is needed to make requests to protected resources on a data pod.

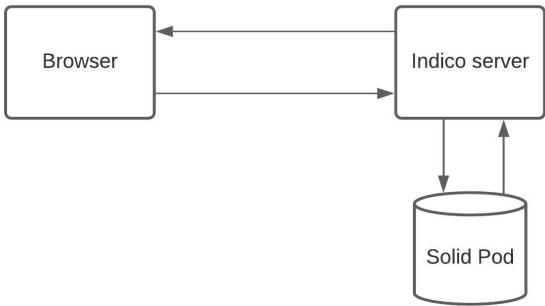


Fig. 6: Communication flow for a module developed on the server.

Comparison of the Different Approaches

Benefits from developing the module for the client:

- Necessary libraries exist (Major release for all basic Solid flows exist)
- Community support
- Programming effort for an minimum viable product (MVP) lowest
- Documentation on developing Solid apps in JS exist

Library	Description
solid-client	A client library for accessing data stored in Solid Pods.
solid-client-authn	A set of libraries for authenticating to Solid identity servers:solid-client-authn-browser for use in a browser.solid-client-authn-node for use in Node.js.
vocab-common-rdf	A library providing convenience objects for many RDF-related identifiers, such as the Person and familyName identifiers from the Schema.org vocabulary from Google, Microsoft and Yahoo!
vocab-solid-common	A library providing convenience objects for many Solid-related identifiers.
vocab-inrupt-common	A library providing convenience objects for Inrupt-related identifiers.

Table 4: Existing solutions to problems for a server approach.

Single Versus Multiple Resource(s) for Comments

Storing the comments in Resource Description Framework (RDF) can be done in two ways: storing it in one file as a graph with a list of comments, or creating a file for every comment.

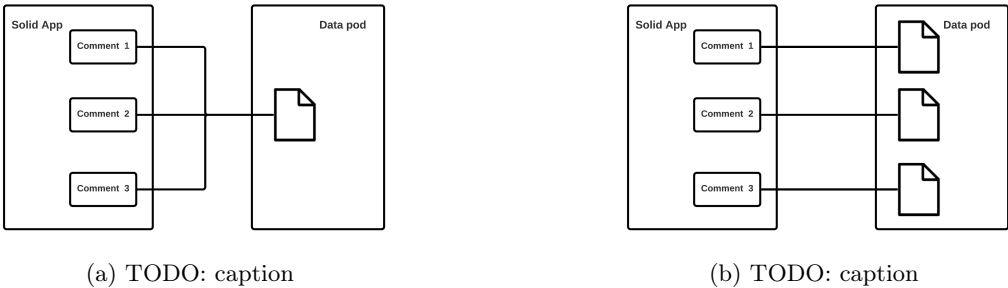


Fig. 7: TODO: two approaches

When fetching the container with all the comment resources the request returns a Turtle file describing the container, but not the actual content of the contained resources. The misconception of receiving the content as well, was thought to be a problem, as it was assumed an initial request to the container reading its child resources had to be made and then for each resource a request needed to be built to retrieve the resource. This overhead was later disproved as the application where this module is embedded maintains the list of resources to be fetched. Therefore, no manual building of requests to those resources had to be done. The next paragraph also lays out how this design would not work with the protection of the resources.

Protection on Resource

Every container and resource in Solid is protected with Web Access Control (WAC), which determines if specific agents, groups, or the world can have read, write, append, or control access. These control access modes are defined in ACL files. The Solid ACL inheritance algorithm looks for an ACL file attached to a specific resource, if it cannot find one it goes recursively up the file hierarchy and looks for ACLs on the containers. Indico allows two general types of protection *private* and *public* on its events. Public means open to everyone, no Indico account or any type of authorization is needed to see the event. Whereas private can be as fine-grained as only to specific agents or groups. A comment module is only valuable if the comments can be read by anyone and be written by authorized users.

In order for visitors of a private or public event in Indico to be able to see the comment, the comment’s ACL needs to allow the public to read the resource. This can be achieved by using the *public* container, which comes with public-read by default on the Node Solid Server (NSS) or by creating a new container and setting the ACL with:

Listing 3.1: TODO: Label caption

```
1 @prefix acl: <http://www.w3.org/ns/auth/acl#>.
2 @prefix foaf: <http://xmlns.com/foaf/0.1/>.
3
```

```
4 # ... Definition for owner
5
6 <#example-container-name>
7   a acl:Authorization;
8   acl:agentClass foaf:Agent;
9   acl:accessTo <./>;
10  acl:mode acl:Read.
```

Every resource in this container is by definition readable by the public – if not otherwise stated in a more detailed resource ACL. The above definition even allows the reading of the container’s content, meaning a request to the container would yield a list of resources in the container. This becomes unpleasant if the Indico event is private and the comments for this Indico event should not be read by the world, which is entirely possible, when browsing to the location of a specific data pod and then looking into the public container.

To prevent a random agent to see the contents of a container, the container can be set to private, with the container’s resources to still be public. This would allow everyone provided they have the Uniform Resource Locator (URL) to browse to the public resource and read it, but not look into the resource’s parent container. To achieve this behavior with ACL, the container needs to just define its owner and no specific rules for the public, as WAC comes with a default private access control. Each child resource needs to define an ACL now, allowing public read. The container’s ACL would look like the following with just a owner defined:

Listing 3.2: TODO: Label caption

```
1 @prefix acl: <http://www.w3.org/ns/auth/acl#>.
2
3 <#owner>
4   a acl:Authorization;
5   acl:agent <https://janschill.net/profile/card#me>;
6   acl:accessTo <./>;
7   acl:default <./>;
8   acl:mode acl:Read, acl:Write, acl:Control.
```

A child resource would allow public read with:

Listing 3.3: TODO: Label caption

```
1 @prefix : <#>.
2 @prefix acl: <http://www.w3.org/ns/auth/acl#>.
3 @prefix foaf: <http://xmlns.com/foaf/0.1/>.
4
5 # ... Definition for owner
6
7 :Read
8   a acl:Authorization;
9   acl:accessTo <test.txt>;
10  acl:agentClass foaf:Agent;
11  acl:mode n0:Read.
```

Another approach and the one implemented after iterating through the previous ones is to have the container’s ACL resource define a default access mode for its child resources. This way one ACL only needs to be created on the container and all resources have proper access modes for public read and are not listed publicly in the container’s description.

Listing 3.4: TODO: Label caption

```
1 @prefix acl: <http://www.w3.org/ns/auth/acl#>.
2 @prefix foaf: <http://xmlns.com/foaf/0.1/>.
3 @prefix target: <./>.
4
5 :ReadDefault
6   a acl:Authorization;
7   acl:default target;;
8   acl:agentClass foaf:Agent;
9   acl:mode acl:Read.
```

**Preventing Resources From Unwanted Discovery**

With the resources having proper access modes but being publicly readable a simple naming convention of taking the International Organization for Standardization (ISO) 8601 string and using it as a filename for the resources created on the data pod does not suffice – even though it is a good strategy when looking for a reliable naming convention to prevent duplication. Considering performance improvements such as pagination for future iterations of the module, which would require some sort of iterative indication, a combination of randomness, but also a order indicator it was settled for using universally unique identifier (UUID) plus the ISO 8601 string to form a filename.

Other ideas included hashing a random string with the timestamp to generate non-guessable filenames. The filename would need to use the same hash function to decipher the filename to figure out when the comment was generated. UUID is a reliable and easy to use system to generate *truly* globally unique strings.

**Modification of Resource From Data Pod**

## *Mitigation of Spam*

Enabling user input in form of comment module without application authentication is a gateway to spam. Even though authentication with a Solid IDP is necessary, it does not hinder a malicious actor to create a multitude of Solid accounts and spam into the application. In the first iteration of the module only Solid authentication was integrated into the module and would therefore allow anyone with a Solid account to post comments and thus also spam the Indico event theoretically. In a second iteration of the module another authentication layer was added to mitigate spam from outside of Indico. For CERN's use-case an authenticated Indico session was enough. This adds an extra step to the comment process and it is ensured only registered Indico users can comment.

## *Giving Application Full Control of Data Pod*

To create or change programmatically ACLs requires *control* access to the container. By default NSS asks for the permissions when authenticating for the first time in the Solid OIDC flow between the *solid-comment* module and Solid IDP. The permissions granted to the application are on the root container of the data pod. Meaning, giving an application control access allows the application to read, write, change ACLs on the entire data pod. This is obviously troubling, as a simple application as a commenting module need to have control access to set the needed ACLs for the containers and resources it creates.

The current implementation has not a built-in solution, but one way of solving it, is the use of an application launcher, which is an application itself with full control access and then limits the access controls of the *solid-comment* module by creating a dedicated container for it and setting the needed ACL for this specific container only.

**1.1.4 Integration with Indico** The need for integration with Indico is twofold: serving the module to the client and being the provider to the list of references to the comments that have been posted on a specific event.

## *Storing Reference to Comments in Indico*

## *Enforce Authenticated Session For Posting Comments*

### **1.1.5 Evaluation**

This section focuses on evaluating the module. This shall be done by iterating through the modules requirements from the stakeholders and how the architecture lives up to those expectations. The evaluation is done with the help of the architectural Software Quality Assurance (aSQA) framework to ensure continuous quality assessment and prioritizing with a light-weight technique [7]. aSQA

## *Metrics*

1. Security
2. Performance
3. Usability

## *Levels*

## *Components*

## *Scenarios*

1. A user who has commented deletes their comment in their data pod
2. A user who has commented deletes their data pod
3. A user uses a script to generate a large number of comments
4. A user creates a comment with a cross-site scripting (XSS) attack

**1.1.6 Analysis** \* EventProxy is not scalable, when comments are sent at the same time could be overwritten \* log IP and blog on spam \* extra layer api to cache comments, performance and security against IP logging \* one pod sends different content to users depending on their location

## **1.2 POC 2: Auto-Complete for Conference Registration in Indico**

### **1.2.1 Architectural Analysis and Synthesis**

## *System Description*

## *Features*

## *Context Diagram*

## *Sequence Diagram*

Stakeholders

Drivers

1.2.2 Design

TODO: 1st iteration, save data in pod 2nd iteration, only pull data from pod  
TODO: Include these somehow:

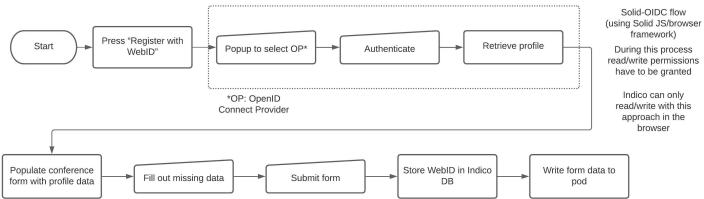


Fig. 8: TODO:

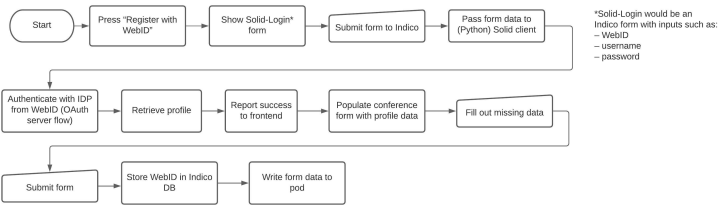


Fig. 9: TODO:

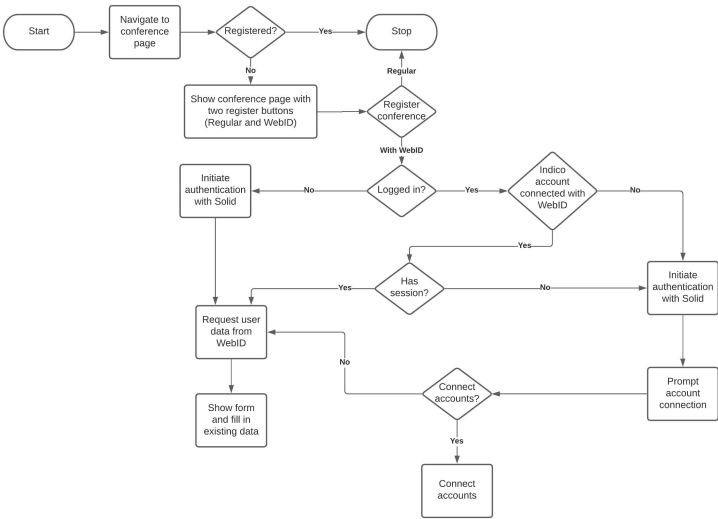


Fig. 10: TODO:

Modification of Resource From Data Pod

Payment on Input Fields

Performance of Large Conference

Availability of Crucial User Data

1.2.3 Integration With Indico

Bind to Dynamically Created Form

#### **1.2.4 Evaluation**

*Metrics*

*Levels*

*Components*

#### **1.2.5 Analysis**

#### **1.3 Deployment of Indico Instance**

### **2 Comparison of Solid and Indico Design Principles**

### **3 Challenges, Advantages, and Gaps of Existing Solid Solutions versus CERN Ones**

### **4 Proceedings in the CERN-Solid Collaboration**

#### **4.1 Servers**

##### **4.1.1 Solution**

##### **4.1.2 Hosting**

#### **4.2 Applications**

# Conclusion







## References

- [1] *Solid JavaScript Client - solid-client*. URL: <https://github.com/inrupt/solid-client-js/>. (Accessed: 08.05.2021).
- [2] *Solid JavaScript Authentication - solid-client-authn*. URL: <https://github.com/inrupt/solid-client-authn-js/>. (Accessed: 08.05.2021).
- [3] *The Solid Common Vocab library for JavaScript*. URL: <https://github.com/inrupt/solid-common-vocab-js>. (Accessed: 08.05.2021).
- [4] *rdflib.js*. URL: <https://github.com/linkedata/rdflib.js/>. (Accessed: 08.05.2021).
- [5] *Solid JavaScript Authentication - solid-client-authn*. URL: <https://github.com/inrupt/solid-client-authn-js/>. (Accessed: 08.05.2021).
- [6] *A Python OpenID Connect implementation*. URL: <https://github.com/OpenIDC/pyoidc/>. (Accessed: 08.05.2021).
- [7] Henrik Bærbak Christensen, Klaus Marius Hansen, and Bo Lindstrøm. *aSQA: Architectural Software Quality Assurance: Software Architecture at Work - Technical Report 5*: English. WorkingPaper. Department of Computer Science, Aarhus University, 2010.