# USAGE OF THE JINR SSO AUTHENTICATION AND AUTHORIZATION SYSTEM WITH DISTRIBUTED DATA PROCESSING SERVICES

## D.I. Gavrilov[1], A.A. Iachmenev[1], I.A. Matveev[1], D.A. Oleynik[2,3], A.Sh. Petrosyan[2,3,a]

[1] *Dubna State University, 19 Universitetskaya st., 141980, Dubna, Russia*

[2] *Joint Institute for Nuclear Research, 6 Joliot-Curie st., 141980, Dubna, Russia*

[3] *Plekhanov Russian University of Economics, 36 Stremyanny per., 117997, Moscow, Russia*

E-mail: [a] artem.petrosyan@jinr.ru

The amount of data produced by the scientific community is already measured in tens and hundreds of petabytes and will significantly grow in the future. High-throughput computing systems have proven to be a solution for handling such data streams. Technologies and a set of products that enable the deployment and use of important components of a distributed system already exist, are quite stable and supported. Meanwhile, building a fully functional system, even with existing components, is not trivial. The Unified Resource Management System is under development at the Meshcheryakov Laboratory of Information Technologies of JINR. This system uses technologies and solutions that were developed during the evolution of middleware platforms for the distributed processing of data from the LHC experiments, but oriented to JINR-based experiments. The basis for the integration of components to the system is the usage of a common authentication and authorization service. This article presents the experience of integrating the CRIC information system and the Airflow-based workflow management system with the JINR SSO authentication service.

Keywords: SSO, distributed computing, workflow management, grid, CRIC, Apache Airflow, PanDA, Python

Artem Petrosyan, Andrey Iachmenev, Ivan Matveev, Danila Oleynik, Dmitry Gavrilov

## 1. Introduction

Distributed data processing systems consist of many components interacting with each other. A classic example of the distributed computing environment is the World Wide LHC Computing Grid (WLCG) built for processing experimental data at the Large Hadron Collider. In a grid environment, authorization and authentication systems allow connecting distributed components into a single computing environment and providing user access control. There are X.509 certificates used to authenticate users, services, and compute nodes, as well as a specially developed authorization system, VOMS [1,2].

Unlike the WLCG, when building a data processing management system on heterogeneous resources within an organization, the problem of not only combining geographically distributed resources into a single computing environment, but also simplifying work with available resources of various types is being solved. Since the construction of the WLCG, authentication and authorization methods that are more convenient to use than the X.509 standard have emerged and become widespread. All this has led to the fact that the WLCG is currently working on a gradual transition from X.509 to modern industry standards, for example, to the use of tokens [3].

Using token-based authentication (OAuth2) makes authentication easier for users and enables integration with third-party applications that also support these standards. In addition, since the industry has not explicitly used X.509 certificates for user authentication en masse, there are significantly fewer resources with support for authentication mechanisms according to this standard than with support for modern standards [4].

All JINR employees receive a login and password to use electronic resources such as mail. Some time ago, the Institute's networking service began developing a Single Sign-On (SSO) service, which allows other services to use user information and thus provides a single entry point for users. This service also makes it possible to store additional information about the user, for example, his affiliation to a particular department and an experiment. Therefore, the use of this service allows eliminating the need to use X.509 and VOMS certificates for users at the organization level. The interaction of services can be carried out, as before, using X.509 certificates.

## 2. Unified Resource Management System

The processing of data, collected by modern physics experiments, is performed using software and hardware systems created for each specific experiment. As a rule, each complex solves the problem of data processing only for the experiment for which it was built and does not imply use in other experiments. However, in recent years, preconditions for changing this situation have appeared. Many middleware services of the grid environment have moved from the category of systems developed individually for a specific experiment into the category of software products that can be installed, configured and used in various experiments. Moreover, some of these software products already allow using one installation to work with several experiments.

On the other hand, there is a serious demand for the unification of software products used in large scientific organizations, such as JINR, which implements a large number of scientific projects, and MLIT JINR, which supports the IT infrastructure of these experiments. This is due to the fact that it is very difficult to provide high-quality support for a wide variety of software systems and products used by experiments.

To solve this problem, the Unified Resource Management System (URMS) is being developed at MLIT JINR, it consists of a set of basic services that support the possibility of using several experiments simultaneously and provide a unified interface to data storage and processing resources. The system is based on components that are used in large centers and scientific infrastructures, such as CERN, BNL, WLCG, and have proven their reliability, flexibility and scalability. In order not to register a user in each of the subsystems, it is proposed to use a single entry point, which will serve as a source of user data for various components of the system and will be responsible for authentication and authorization. In the URMS, this entry point will be the JINR SSO [5]. The URMS architecture is presented in Figure 1.

The system consists of the following components:

• Workflow management system – controls the process of data processing on each step of processing. Produces tasks, which is required for processing of a certain amount of data, manages task execution.

• Workload management system – processes task execution by splitting a task into small jobs, where each job processes a small amount of data. Manages the distribution of jobs across a set of computing resources. Takes care of generating the proper number of jobs until the task is completed (or failed).

• Data management system – responsible for the distribution of all data across computing facilities, managing data (storing, replicating, deleting, etc.).

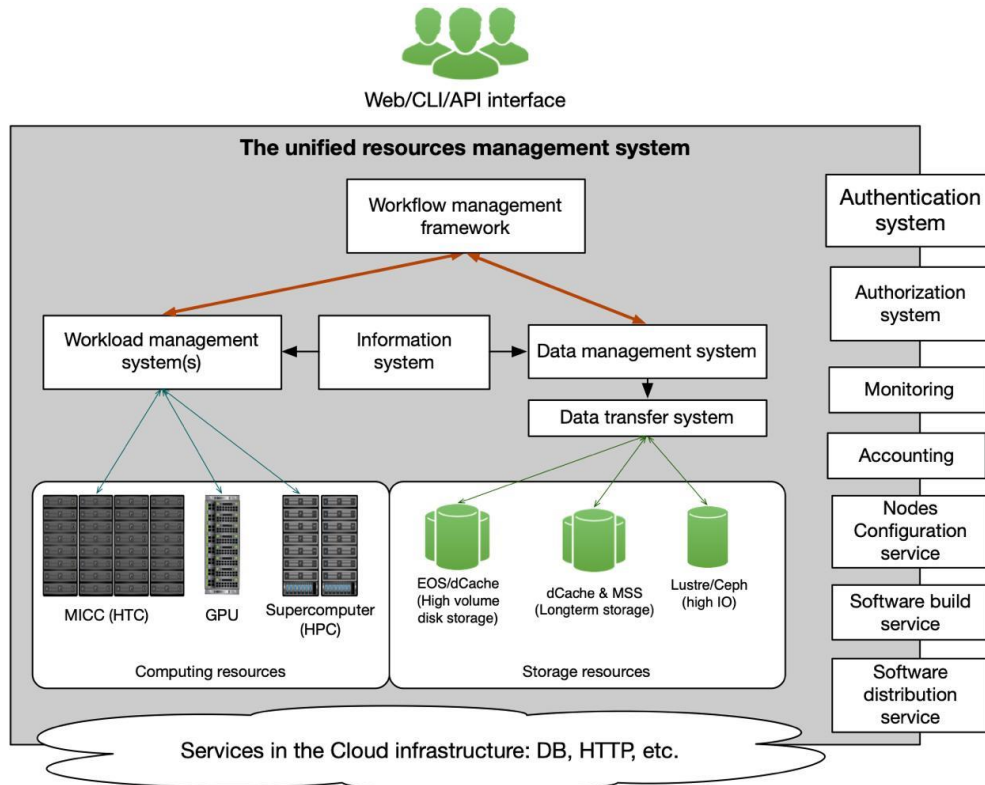• Data transfer service: takes care of major data transfers. Allows asynchronous bulk data transfers.



Figure 1. URMS architecture

## 3. JINR SSO

As mentioned above, JINR is a large scientific organization, and it is absolutely impossible to manage a modern organization of such a scale without the widespread introduction of electronic systems. It is clear that the number of such systems is very significant, given the number of departments and areas of activity of the Institute. Since the main activity of the Institute is physics research, IT solutions in the field of information support are usually selected from the most suitable ones available on the market. Usually, each of them has its own built-in authentication mechanism. Naturally, at some point, the task arose to streamline access to IT products used by the Institute and integrate them with the personnel accounting system to simplify administration and ensure information security.

The JINR SSO service is a key system that provides a single entry point for a user of the JINR network and information infrastructure. This service is integrated with the accounting system in the personnel department of JINR, which makes it possible to authenticate and authorize users based on the data entered into this system.

The JINR SSO operates according to the OAUTH2.0 standard [6]. Communication with the service takes place using POST and GET requests. A user wishing to use the system UI must go through the authentication and authorization procedure. When this procedure is initiated, the user is redirected to a page with a form for entering data from the JINR account in the JINR SSO service. If the entered data is correct, the service will return a confirmation code. Next, the web application exchanges the code for a token from the JINR SSO service. With a valid token, a third-party application can obtain information about the user for the authorization procedure (Figure 2).
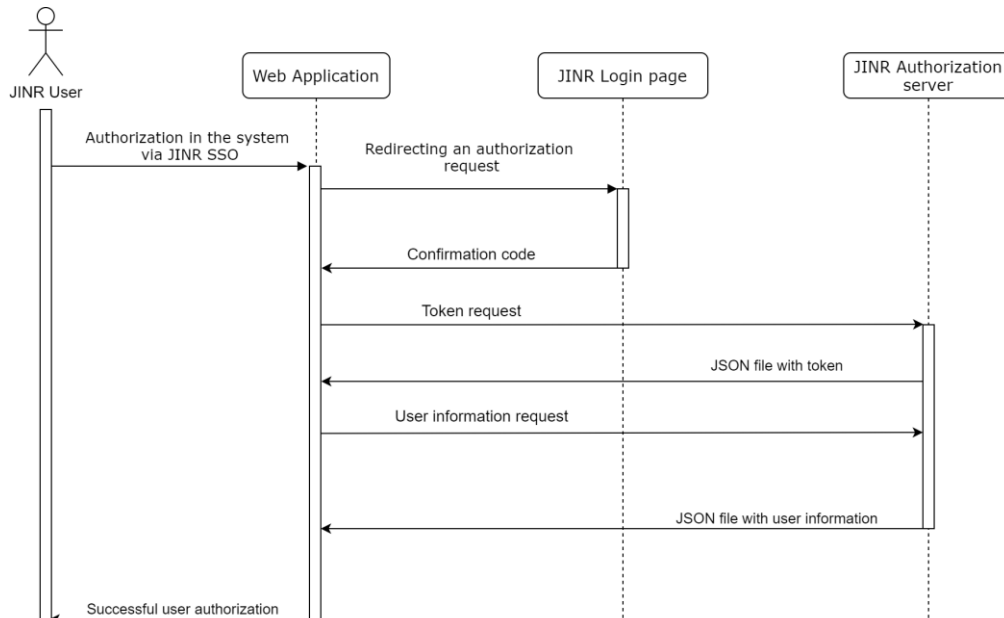
Figure 2. Scheme of the interaction of the application with the JINR SSO

## 4. CRIC Integration

The Computing Resource Information Catalog (CRIC) information system (at that time AGIS – ATLAS Grid Information System) was originally developed by the ATLAS collaboration at the LHC for their own needs. Later, the system began to be used in other projects. At present, the CRIC system is the main information system of the WLCG project [7]. The system has a robust web interface and allows one to describe any combination of systems and services of the computing infrastructure of any experiment.

CRIC is used in the URMS as an information system. The information system, which is one of the crucial components of the URMS, stores and provides a description of resources in a consistent way. CRIC implements a CERN-specific authentication and authorization mechanism. Therefore, this mechanism has been redeveloped to work with the JINR SSO service. The CRIC authorization system provides a wide set of rights, which allows it to be applied to the JINR authorization system.

## 5. Apache Airflow Integration

Airflow is a platform for programmatically creating, scheduling and monitoring task chains. In the URMS, Airflow is used as a workflow management system [8]. Airflow is assumed to operate at the top level of the data processing system. The responsibilities of the workflow management system include defining the tasks to be processed, launching these tasks in the PanDA workload management system, monitoring the execution of tasks. Airflow does not work directly with computing resources and only interacts with the PanDA system. To work with Airflow, users use a Web UI, in which they create tasks and control their execution. At this stage, there is a need for user authentication and authorization. Airflow supports plugins for various authentication and authorization methods. For its integration with the JINR SSO, a plugin was developed.

## 6. PanDA

PanDA is a task execution system [9]. In the URMS, it is used as a workload management system. Tasks are created in the workflow management system (Airflow) and submitted to the workload management system (PanDA) for execution. The workload management system sends tasks to computing resources, data about which comes from the CRIC information system. Thus, a person without a certificate and other encumbrances can start their tasks for processing in a distributed computing environment using a workload management system that is integrated with Airflow using a special operator.

To interact with PanDA, the Airflow operator uses the PanDA Client. This is a set of classes that implements interaction with the PanDA API. To send a task to the PanDA workload management system, a special object, which describes all the parameters of the task, is created. This object is serialized, and a request is made to the workload management system using the Panda API.

## 7. Conclusion

We are building a system that consists of different subsystems, each of which requires authentication and authorization. As part of the work on the development of the URMS, it is shown that it is possible to provide simple and understandable access to the computing resources of a large scientific organization through an application, which is widely used in the industry. In principle, in the same way, through a middleware layer in the form of PanDA, one can send tasks to another scientific computing infrastructure. All the complexity of such infrastructures remains hidden from the user and does not rise above the middleware with this approach. We managed to create a single entry point for them using the JINR SSO and leave certificates only at the service level, which is much more convenient for users. Basically, not only PanDA, but also DIRAC can be used as an integrating middleware, via describing it as a plugin in Apache Airflow [10].

## 8. Acknowledgement

## References

[1] X.509, available at https://www.itu.int/rec/T-REC-X.509 (accessed 28.09.2021)

[2] VOMS, available at https://italiangrid.github.io/voms/ (accessed 28.09.2021)

[3] B. Bockelman et al., WLCG Authorisation from X.509 to Tokens, EPJ Web Conf., Vol. 245, 2020

[4] D. Dykstra et al., Secure Command Line Solution for Token-based Authentication, EPJ Web Conf., Vol. 251, 2021

[5] JINR SSO, available at https://noc.jinr.ru/en/service/sso-serv.php (accessed 28.09.2021)

[6] OAuth 2.0, available at https://oauth.net/2/ (accessed 28.09.2021)

[7] Anisenkov et al, CRIC: Computing Resource Information Catalogue as a unified topology system for a large scale, heterogeneous and dynamic computing infrastructure, EPJ Web Conf., Vol. 245, 2020

[8] Apache Airflow, available at https://airflow.apache.org/ (accessed 28.09.2021)

[9] PanDA, available at https://panda-wms.readthedocs.io/en/latest/index.html (accessed 28.09.2021)

[10] DIRAC, available at https://dirac.readthedocs.io/en/latest/ (28.09.2021)