# Prediction and optimization of the computing and storage resource usage in Grid infrastructure of CERN ALICE experiment

## CERN-ի ALICE Գիտափորձի Գրիդ ինֆրակառուցվածքում հաշվողական և կուտակային ռեսուրսների օգտագործման կանխատեսումը և օպտիմալացումը

### *Armenuhi Abramyan*

*Master of Informatics and Computer Science*

Supervisor:     *Prof. Ara Grigoryan,*
                *Leader of ANSL/ALICE team*

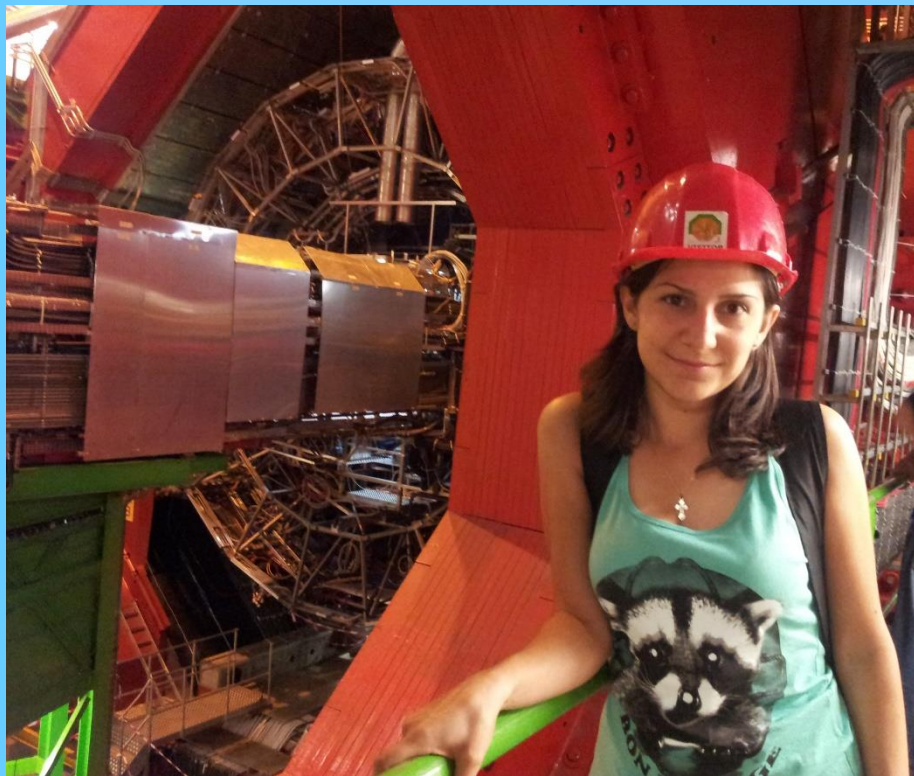*aabramya@{mail.yerphi.am, mail.cern.ch}*

# Problem statement

High Energy Physics experiments, such as the *Large Hadron Collider(LHC)* at *CERN*, pose a challenge to current big data handling methodologies. The biggest scientific communities produce petabytes of data each year, which must be efficiently stored, processed and analyzed. The modern approach to the solution of big data handling and processing is the exploitation of **Grid technologies**.

The exploitation of *Grid technologies* by various scientific communities has showed their high effectivity in the solution of various, even very complicated, problems. However, the continuous increase of the data volume and the complexity of the problems, with the financial limitations for the proportional increase of the storage and computing resources, revealed the following problems:

- *Need in maintenance of the continuous development process of Grid infrastructure software, in order to comply with modern requirements.*

- *Need of the detailed monitoring of the data usage in Grid infrastructure.*

- *Lack of the algorithms for the optimization of data distribution over the Grid infrastructure.*

- *Necessity of the prediction of the storage and computing resource requirements for the certain period of time.*

*My research and work is devoted to the solution of aforementioned problems within the Grid infrastructure of ALICE experiment of LHC at CERN.*

# ALICE and its Grid environment



The aim of *ALICE experiment*, one of *4* biggest *LHC* experiments at *CERN*, is to explore the primordial state of matter that existed in the first instants of our Universe, immediately after the initial hot *Big Bang*.

The *ALICE detector* has been built by a collaboration, which currently includes *1800* physicists and engineers from *176* institutes in *41* countries.

The *AliEn*, ALICE Grid Environment, is a set of *Grid* middleware and application tools and services which are exploited by *ALICE collaboration* to store and analyze the experimental data, as well as to perform Monte-Carlo simulations.

# *Outcomes and practical implementation*

# 1. Analysis of the functionality of AliEn services

In order to support the continuous process of the modernization of the *AliEn* services and the effective involvement of developers, especially of newcomers into the upgrade of the *AliEn* code, **ALICE** *Offline team* has decided to provide a transparent view of the *AliEn* functionality.
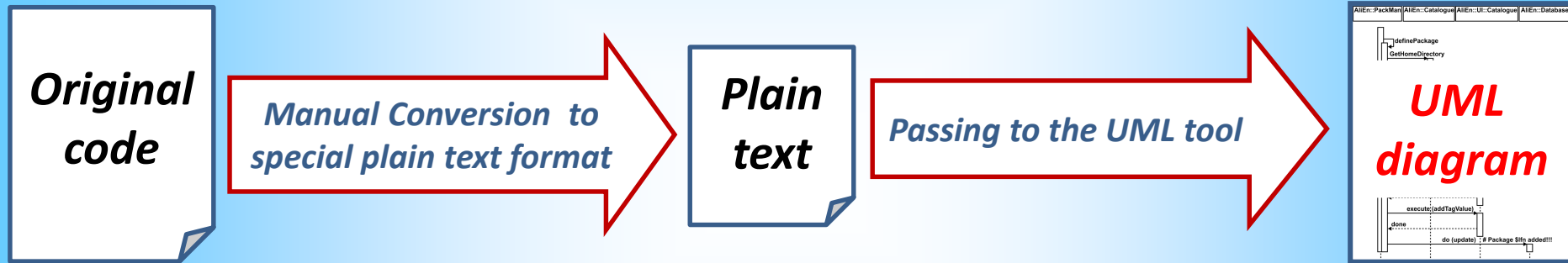
This is done by means of *Unified Modeling Language (UML) Sequence diagrams*.

The dynamic behaviour of system is described through the tracing and visualization of the sequence of interactions that occur between the objects or components of a system.

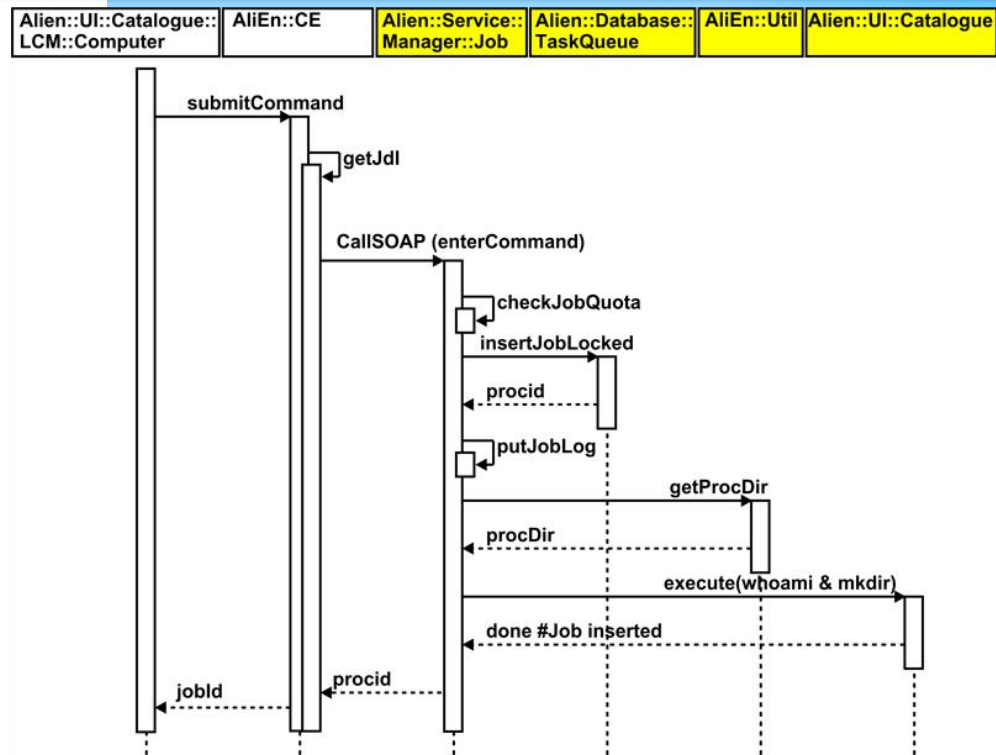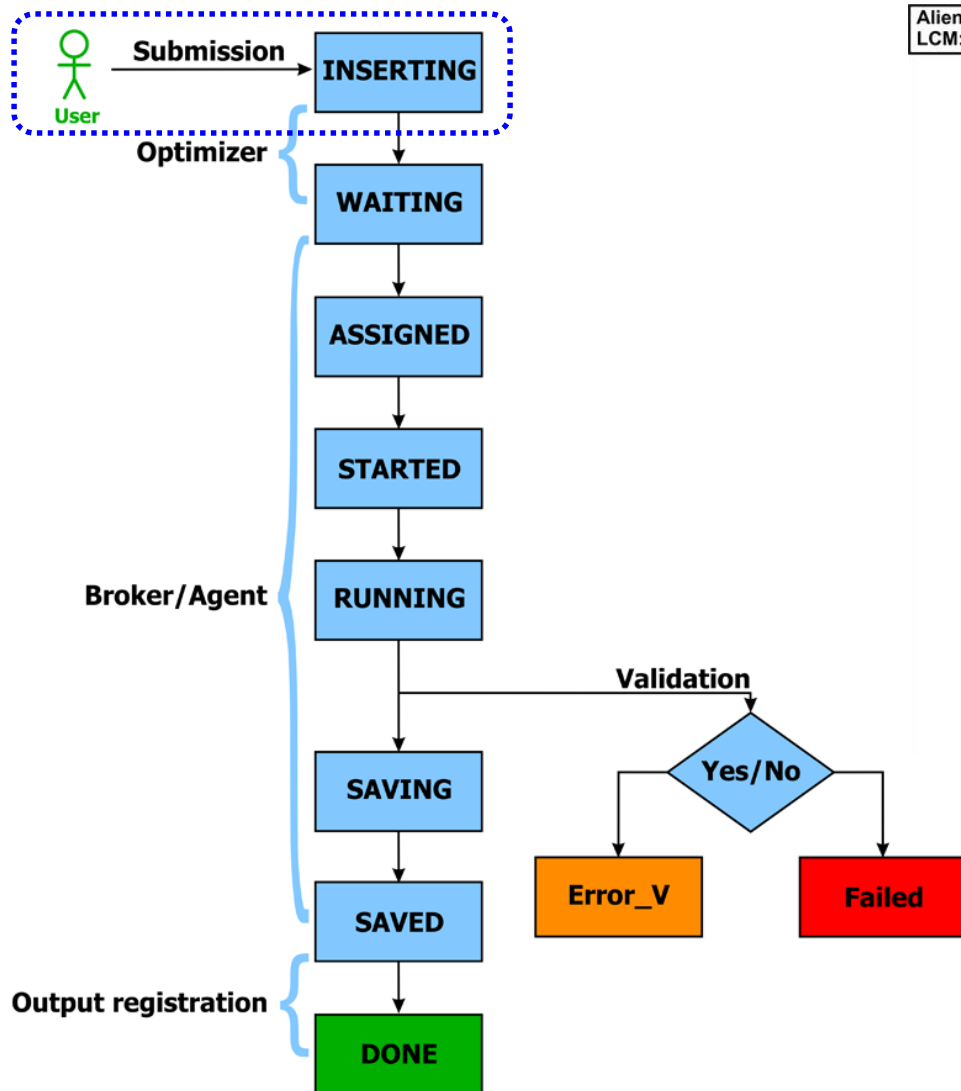**All created diagrams (30 in total) are included in the *AliEn* official documentation for developers**

http://alien2.cern.ch/index.php?option=com_content&view=article&id=84&Itemid=127
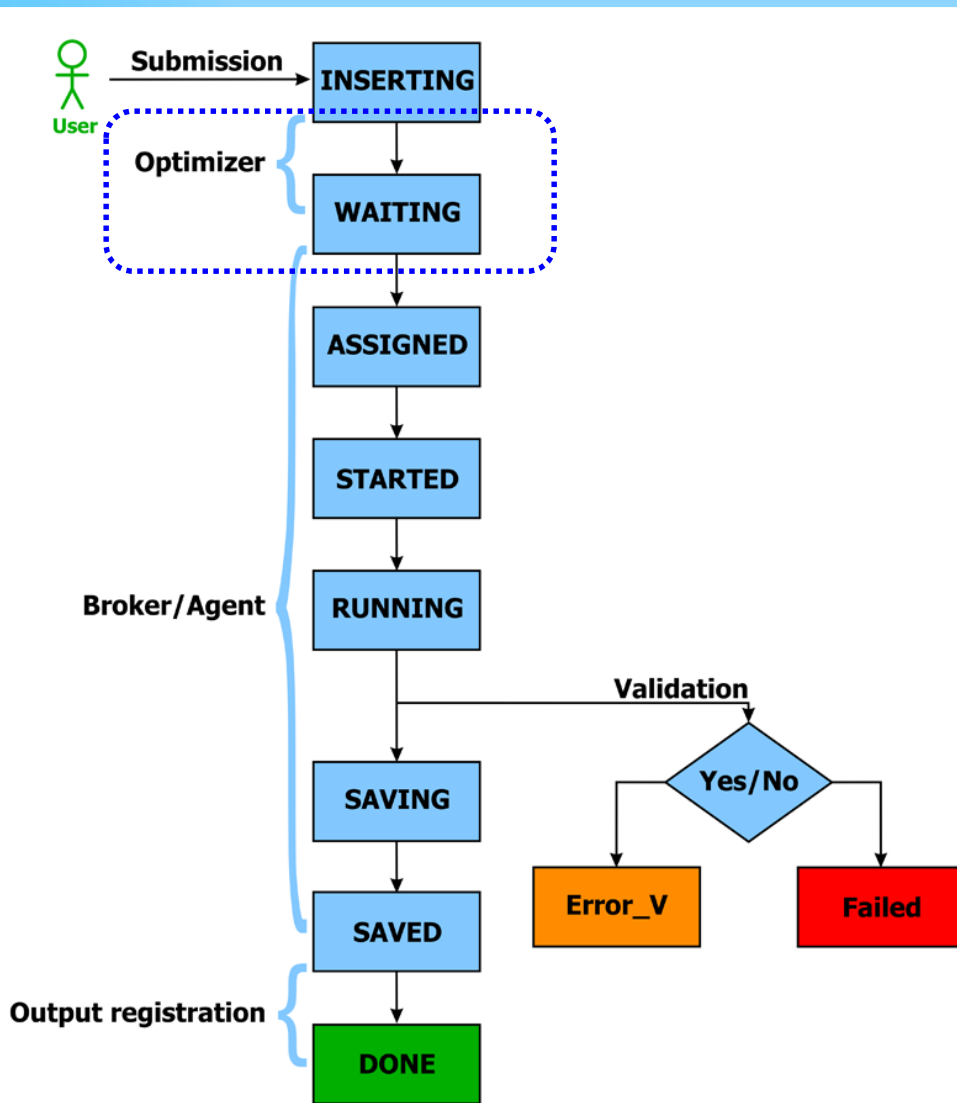
# Creation of Sequence diagrams

**Original code** → **Manual Conversion to special plain text format** → **Plain text** → **Passing to the UML tool** → **UML diagram**

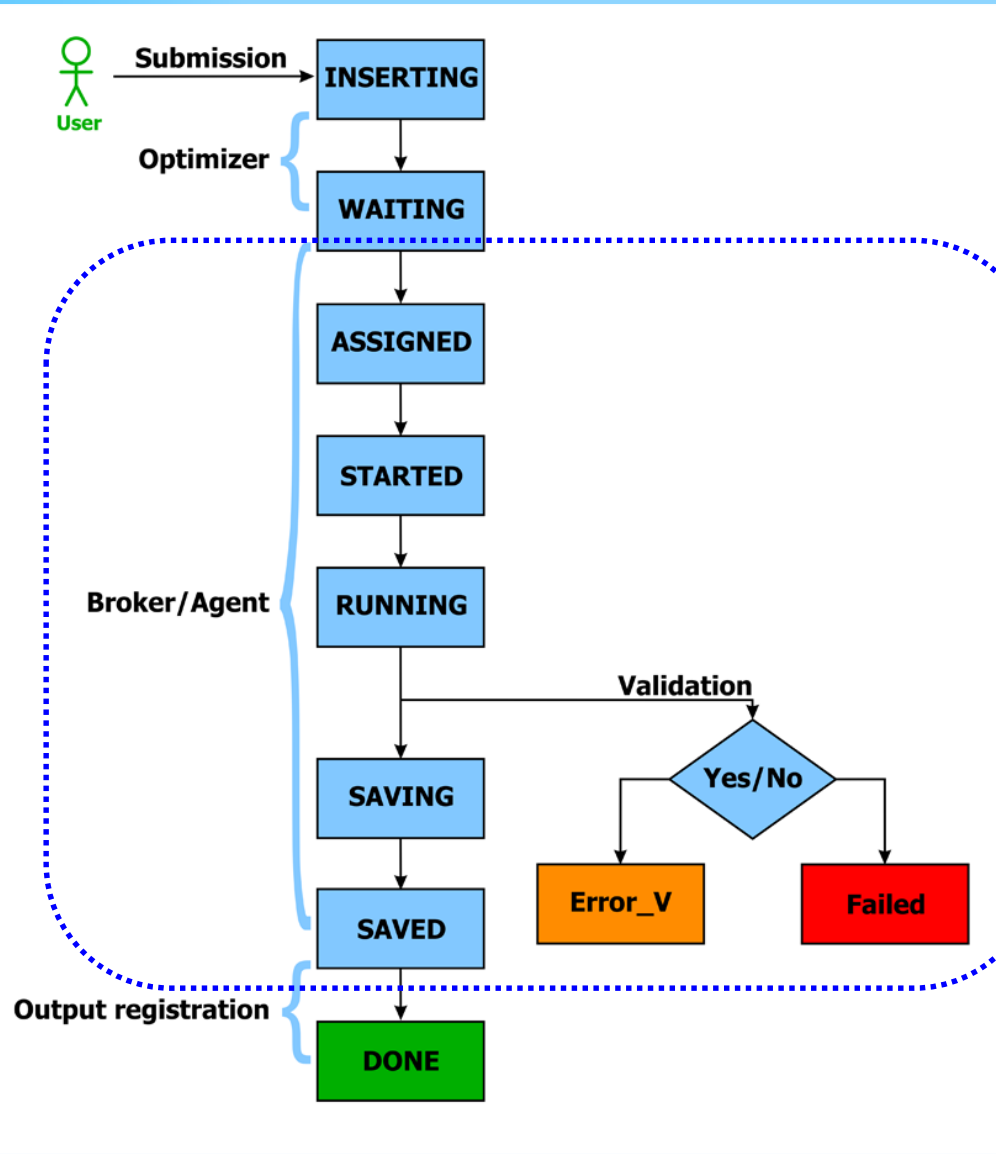# Sequence diagrams for AliEn Job Execution chain
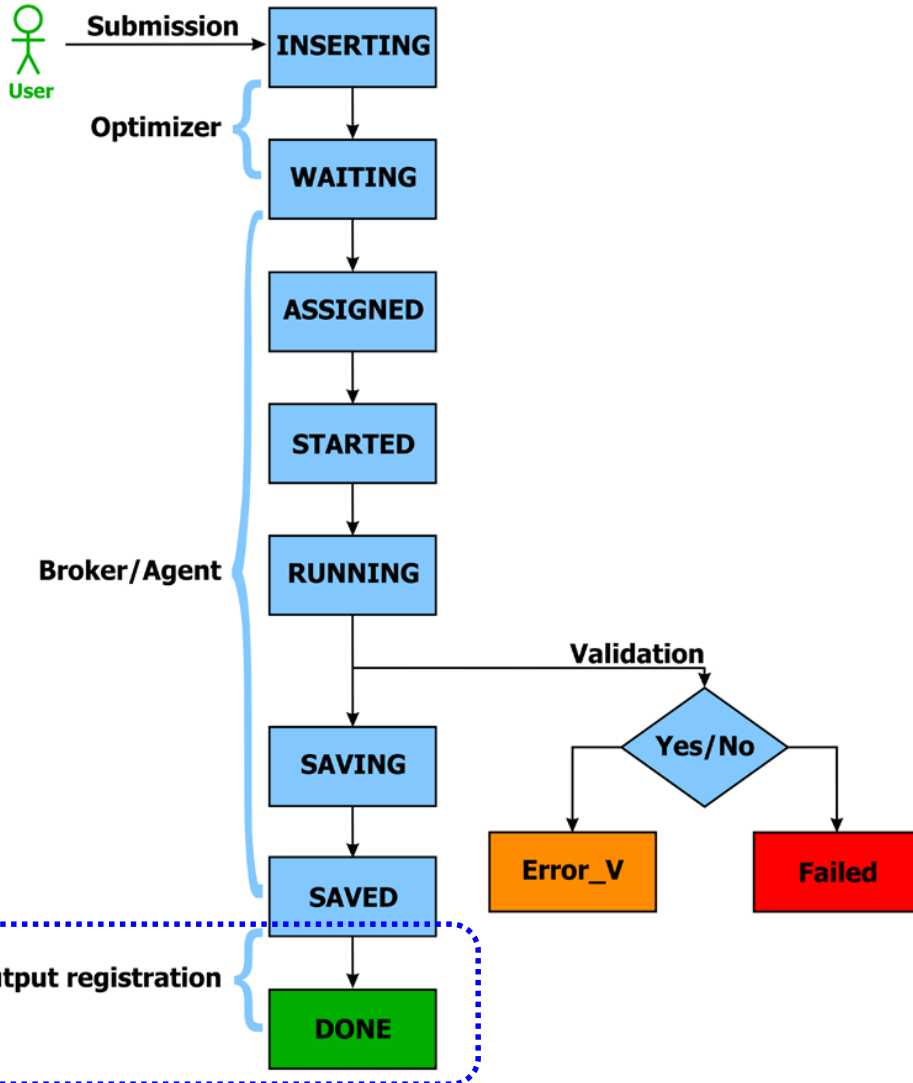## Submission stage

# Broker/Agent stage

# Output registration stage

# Sequence diagrams for the commands of AliEn PackMan service

The Package Manager (**PackMan**) service of *AliEn* provides automatic installation, removal, upgrade and configuration of application software packages needed for the execution of user Jobs.

The **UML Sequence diagrams** describing the functionality of each of the commands of **PackMan** service have been created.

### PackMan commands

- define
- list
- listInstalled
- Install
- dependencies
- installLog
- remove
- undefine
- synchronize
- test

# 2. Upgrades for PackMan service (with Narine Manukyan)

After analysis of the functionality of the *AliEn* **PackMan** service on the base of **UML Sequence diagrams**, a model of the upgrade of this service, concerning especially the removal of the bottleneck, caused by the slow execution of commands of this service, has been suggested and implemented.

| PackMan performance statistics before and after its upgrade | |
|---|---|
| **PackMan commands** | **Ratio of execution times Before and After upgrade** |
| list | 32.6 |
| listInstalled | 27.8 |
| undefine | 2.3 |
| installLog | 1.4 |
| install | 1.3 |
| remove | 1.2 |
| define | 1.04 |
| test | 1.02 |

**The upgraded code has been implemented in *AliEn* software, since 2012.**

# 3. FAMoS - File Access Monitoring Service
## (with Narine Manukyan)

The purpose of the *FAMoS* is to monitor (in offline mode) the calls/accesses to the files in AliEn and to provide their storage in an organized manner.

*FAMoS* records the values of the following attributes to a special database, called *accesses*.

| Attribute | Description |
|---|---|
| **File name** | *Name of the* file in *AliEn* file system |
| **SE name** | Name of SE from where the file was accessed |
| **User name** | Name of user by whom the file was accessed |
| **Access Time** | Time and date when the file was accessed |
| **Operation result** | Successful or failed access |

*The monitoring of the accesses to the files is performed not only on the level of the files themselves, but also on the level of their (meaningful) combinations, called* **Categories.**
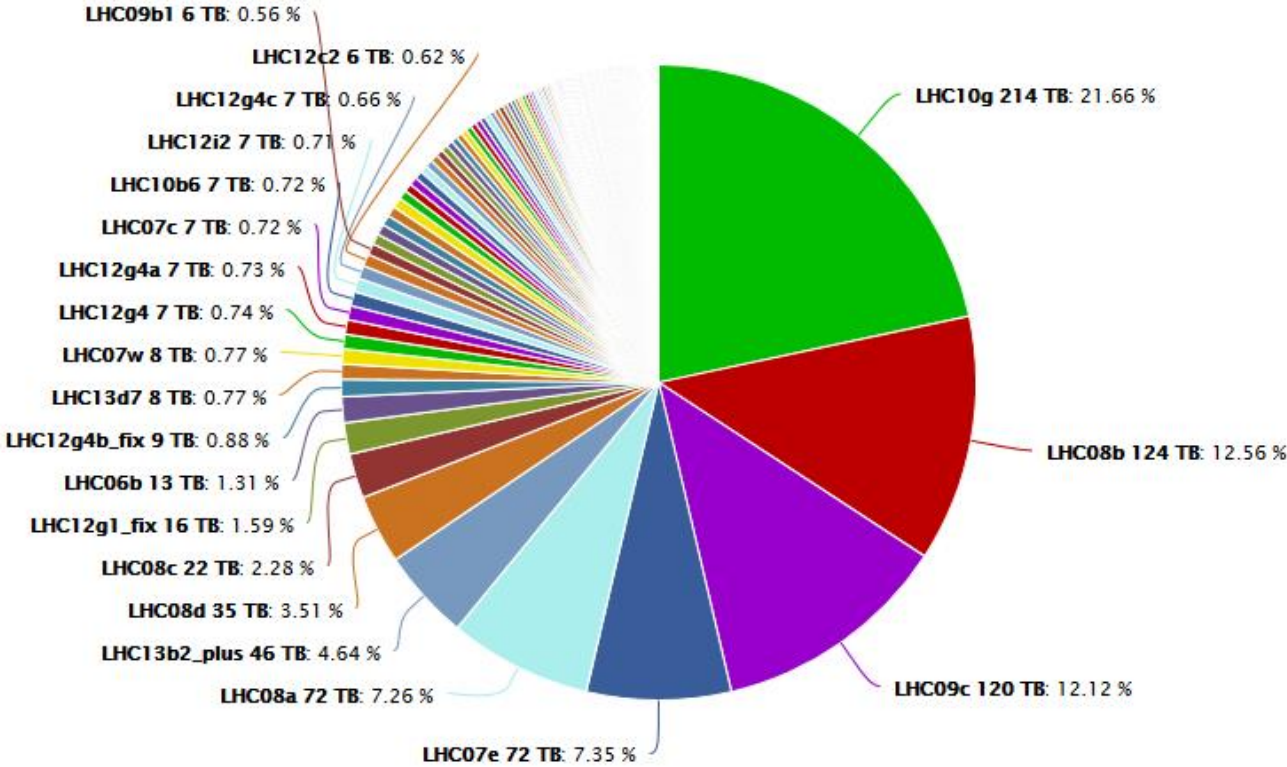
# FAMoS components

FAMoS consists of the following components:

1. **Collector** -A Bash script runs on daily basis on the API and Authen servers and uploads the log file of the last day accesses to the FAMoS server.
2. **Parser** - Perl module running on the FAMoS server and processing the log files of last day; aggregates the file attributes by categories, access time, user, etc. The Parser inserts the processed data into a database, called Accesses.
3. **Accesses** - MySQL database, which contains access information with the time granularity of one day.
4. **Web interface** for statistical representation of categories access information, volumes, access time, user and SE.
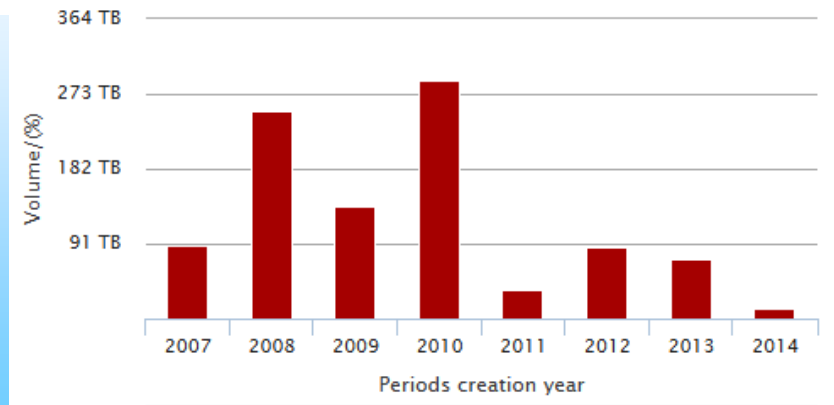
*The FAMoS service has been put into exploitation by ALICE collaboration since August 2013.*

# Volumes of accessed/not accessed periods
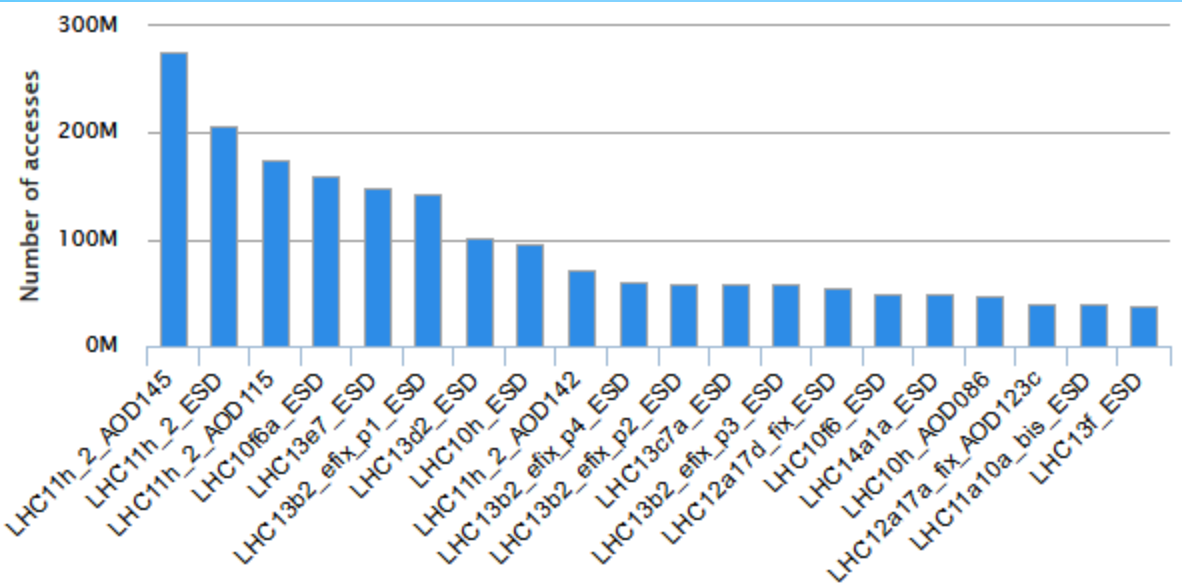


**Volumes of not accessed periods**

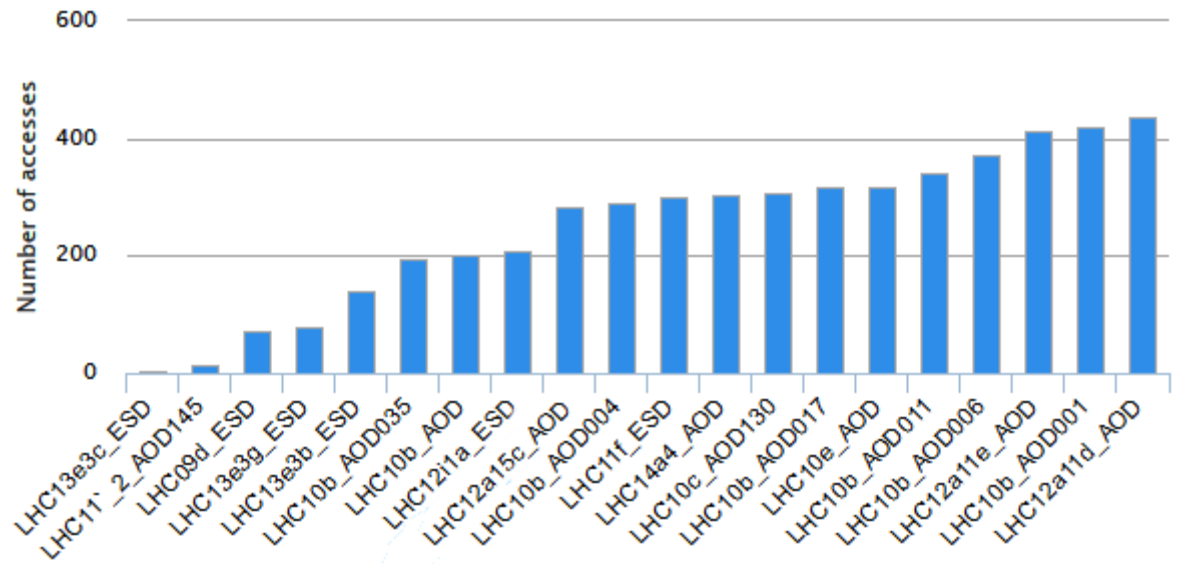**Volumes of not accessed periods grouped by their creation year**

# Most and least accessed categories



**20 most accessed categories**

**20 least accessed categories**

# Usage of the ESD and AOD files over time

# Volumes of data versus number of accesses



**ALICE number of accesses in time X**

Legend: ● 3 months  ● 6 months  ● 1 year

Y-axis: Aggregated data size / PB
X-axis: Number of accesses

(younger than X)

Chart generation startdate: June 2018

Volumes of data versus number of accesses in 3-, 6- and 12- month periods. For each period X, data created in that period but not accessed is in the second bin. The first bin is for data created before the period began.

# *FAMoS* web interface

# 4. ALICE Computing Model simulation software (*with Narine Manukyan*)

The software is to perform discrete-event simulations (DES) of ALICE data taking process for certain period of time for a given computing model 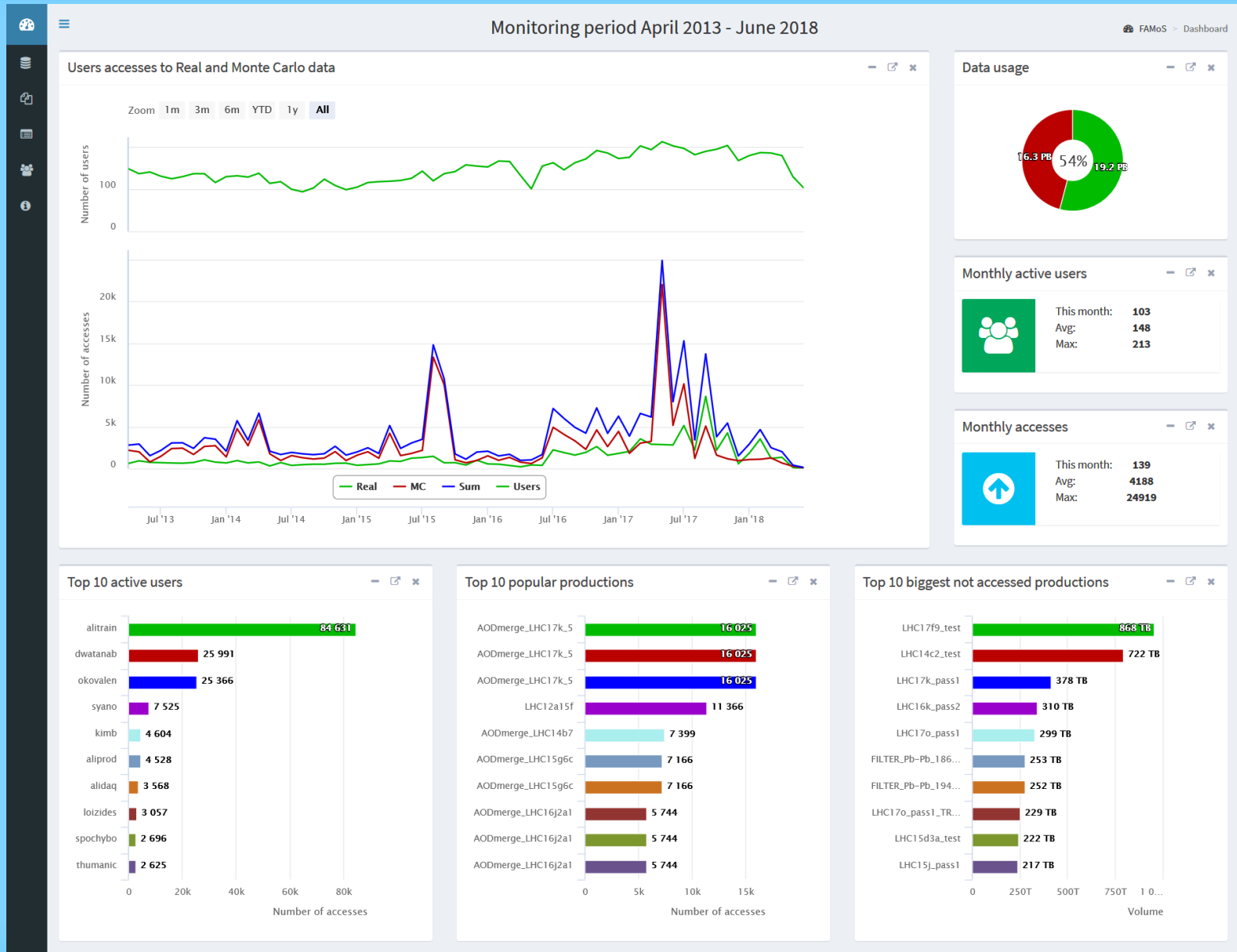layout* with the aim to estimate the usage of ALICE resources required to process and store the data during Run 3 (2021-2023) and Run 4(2026-2028).



*Resources required to process and store data during Run3/Run4*

* A combination of FLPs, EPNs and other resources, their role as well as the network topology by which these resources are connected.

**Under development**

Flexible and highly configurable tool that gives possibility to estimate (via DES) the **CPU** and **storage resource** usage required to process and store each type of ALICE data (during Run3 and 4), by taking into account *LHC running schedule*, *Conference calendar*, *data management/removal policies* and any other criteria.

*The works on the simulation with SIM.JS are done with Tim Hallyburton.*

**For each collision type we specify:**

✓ Planned number of collisions (for year) – $N_{collisions}$

✓ Collision rate (Number of collisions per second) - $C_{rate}$

✓ CTF size per event - $E_{size}$

✓ Data taking efficiency factor– Efficiency (%)

**CTF_size_per_day** = ($C_{rate}$ * $E_{size}$ * Efficiency/100)
* Seconds_in_a_day

**Parameters that we can play with:**

- $C_{rate}$
- **Efficiency**

**ALICE running scenario for the LHC Run3 and 4**

| | Year | Collision type | $N_{collisions}$ | $E_{size}$ (kB) |
|---|---|---|---|---|
| **Run3** | 2020 | pp<br>Pb-Pb | $2.7 * 10^{10}$<br>$2.3 * 10^{10}$ | 50<br>1600 |
| | 2021 | pp<br>Pb-Pb | $2.7 * 10^{10}$<br>$2.3 * 10^{10}$ | 50<br>1600 |
| | 2022 | pp<br>pp | $2.7 * 10^{10}$<br>$4 * 10^{11}$ | 50<br>50 |
| **Run4** | 2025 | pp<br>Pb-Pb | $2.7 * 10^{10}$<br>$2.3 * 10^{10}$ | 50<br>1600 |
| | 2026 | pp<br>Pb-Pb<br>p-Pb | $2.7 * 10^{10}$<br>$1.1 * 10^{10}$<br>$10^{11}$ | 50<br>1600<br>100 |
| | 2027 | pp<br>Pb-Pb | $2.7 * 10^{10}$<br>$2.3 * 10^{10}$ | 50<br>1600 |

**LHC Schedule (2017)**

Legend:
- Technical stop
- Recommissioning with beam
- Scrubbing (indicative - dates to be established)
- Machine development
- Physics runs
- Special physics runs (indicative - schedule to be established)

e $_{O2}$  j $_{T0/T1s}$

| Site Type | Site Name | CPU resources (N of CPU cores) | Storage Resources | | Tape resources | | ⊕ |
|-----------|-----------|--------------------------------|-------------------|---|---------------|---|---|
| O2 | O2 | 5000 | 90000 | GB ▼ | 80000 | GB ▼ | ⊖ |
| T1 | T1 | 7000 | 80000 | GB ▼ | 60000 | GB ▼ | ⊖ |
| T2 | T2 | 25000 | 10000 | GB ▼ | 0 | GB ▼ | ⊖ |

| | CpuTransformations (HS06 - consumed CPU seconds per event) | | | | CpuShare (%) | | | |
|---|---|---|---|---|---|---|---|---|
| | pp | pPb | PbPb | pp-ref | O2 | T1 | T2 | AF |
| RAW -> CTF | 0 | 0 | 0 | 0 | 100 | 0 | 0 | 0 |
| CTF -> ESD -> AOD | 300 | 710 | 3800 | 300 | 67 | 33 | 0 | 0 |
| MC -> MCAOD | 1000 | 3000 | 45000 | 1000 | 0 | 0 | 100 | 0 |
| AOD -> HISTO | 200 | 700 | 3700 | 200 | 0 | 0 | 0 | 100 |
| ⊕ | | | | | | | | |

**CTF**
**ESD** (15% of CTF size)
**AOD** (10% of CTF size)
**MC** (100% of CTF size)
**MCAOD** (30% of CTF size)
**HISTO** (1% of ESD size)

| Data Types | Replication factor | | Storage Sharing (%) | | | | | LifeTime on Disk (days) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Disk | Tape | O2 | T0 | T1 | T2 | AF | O2 | T0 | T1 | T2 | AF |
| CTF | 1 | 1 | 66.6 | 0.0 | 33.3 | 0.0 | 0.0 | 270 | 0.0 | 270 | 0.0 | 0.0 |
| ESD | 1 | 0 | 75.0 | 0.0 | 25.0 | 0.0 | 0.0 | 10 | 15 | 25 | 0.0 | 0.0 |
| AOD | 1 | 1 | 15 | 0.0 | 0.0 | 1.0 | 0.0 | 150 | 250.0 | 100 | 0.0 | 100 |
| MCAOD | 1 | 1 | 0 | 25 | 75 | 0.0 | 1.0 | 100 | 100.0 | 100.1 | 5.0 | 100 |
| HISTO | 1 | 0 | 10.0 | 5.0 | 75.0 | 0.0 | 10.0 | 10.0 | 100.0 | 150.0 | 0.0 | 50.0 |
| ⊕ | | | | | | | | | | | | |

# ALICE Calendar as an additional input for simulations



Under development and discussion stage

Amount of created data by their type

With the LHC 2017[th] schedule and Run3/Run4 configuration parameters (without data removal), at the end of the 2017 we expect (only on T1 sites):

SUM – **6.1 PB**
CTF – **4.9 PB**
ESD – **0.7 PB**
AOD – **0.5 PB**

# *Presentation of the work*

| Authors | # | Reference |
|---|---|---|
| A. Abramyan | 1 | "**Unified Modeling Language Diagrams for Grid Middleware of CERN ALICE Experiment**". Proc. 8[th] Int. Conf. "Computer Science and Information Technologies",  Yerevan, Armenia, pp. 352-355, 2011. |
| A. Abramyan, N. Manukyan | 1 | "*AliEn File Access Monitoring Service – FAMoS* ". Proc. 9[th] Int. Conf. "Computer Science and Information Technologies",  Yerevan, Armenia, pp. 311-313, 2013. |
| A. Abramyan, N. Manukyan and 11 members of ALICE Offline software team | 2 | "**AliEn Extreme JobBrokering**". Proceedings of Computing in High Energy and Nuclear Physics (CHEP) conference, New York City, NY, USA (2012) |
|  |  | "*ALICE Environment on the GRID*". Proceedings of Computing in High Energy and Nuclear Physics (CHEP) conference, New York City, NY, USA (2012) |
| **A. Abramyan**, **N. Manukyan** and 4 members of ALICE Offline software team | 1 | "*A simulation tool for ALICE storage and computing resource usage*" abstract has been accepted by the CHEP 2018 Conference, 9-13 July 2018 Sofia, Bulgaria |
| A. Abramyan, et al. (ALICE Collaboration) | 9 | … |

# Thank you

# Backup slides

# *Practical relevance*

The problem of prediction and optimization of the computing and storage resource usage is known not only in Grid networks, but also in all contemporary scientific cyber-infrastructures, where the big data transfer and caching operations are handled. For example, collaborations in meteorology, seismology, ecology, biology, physics, astronomy and medicine need to store huge amounts of data (of the order of hundreds of Petabytes), as well as to conduct massive computations. However, there is no common solution, which can be applied to all networks. Consequently, every private network requires its unique approach.

In the case of ALICE experiment, this problem is relevant, because its solution will make it possible to maintain and efficiently process the data accumulated in the collisions of LHC, the world's largest proton and ion accelerator, containing the most important and fine-grained information about the structure of the matter and the origin of our universe.

# ALICE - A Large Ion Collider Experiment

ALICE is one of the four large experiments running at *Large Hadron Collider (LHC)*. ALICE detector is designed to study both bulk and microscopic properties of the matter created in PbPb, pPb and pp collisions, through the identification and analysis of numerous signals emanated by the matter.

Early stages of PbPb collisions at *LHC* energies - an extremely dense (more than $10^{15}$ g/cm³) and hot (more than $10^{12}$ K) medium of free quarks and gluons - *Quark-Gluon Plasma (QGP)* is formed.

**Our Universe was in *QGP* state at the age of a few microseconds after the Big Bang.**

The ALICE apparatus - **18** different subsystems: detectors, magnets, absorber.

Overall dimensions **16×16×26 m³**,

Total weight ~ **10 000 t**.

The apparatus provides full identification of several thousands of particles produced in the central barrel acceptance



## ALICE Computing environment

Computing in ALICE - an experiment-specific, complex Object-Oriented framework, called *AliRoot.* The Analysis of very large sets of experimental data collected in ALICE and massive simulations - ALICE *Grid environment*, called *AliEn*.

# *ALICE Virtual Organization (ALICE VO)*

*VO = Computing environment + Authenticated and Authorized Users*

*ALICE VO* enables *ALICE* user community to effectively perform their simulations and analysis on the distributed computing centers.



CERN
France-Switzerland

TIER 0
(Volume ≈ 30000 TB)

IN2P3 France | CNAF Italy | FZK Germany | KISTI South Korea | RAL England | SARA Netherlands

TIER1 x 6
(Total volume ≈ 33400 TB)

OSC USA | Wuhan China | UNAM Mexico | SaoPaulo Brasil | YerPhI Armenia

TIER 2 x 53
(Total volume ≈ 113000 TB)

A transparent access to the resources of the *ALICE* Collaboration is provided by the *AliEn Grid* middleware.

# AliEn² - ALICE Environment on the Grid

*AliEn* was the first *Grid* put into the exploitation at CERN (in the end of 2001)

The core services of *AliEn* infrastructure are open source *web services*, which communicate with each other via *JSON (JavaScript Object Notation)* messages using *RPC (Remote Procedure Call)* protocol and thus create a network of collaborating services.

**2 categories of the *AliEn* services:**

- **Central services:** *Job Manager, Job Broker, Task Queue and File Catalogue databases, Job Optimizer, …*

- **Site services:** *Computing Element, Cluster Monitor, Package Manager and Job Agent.*

**AliEn functionality provides:**
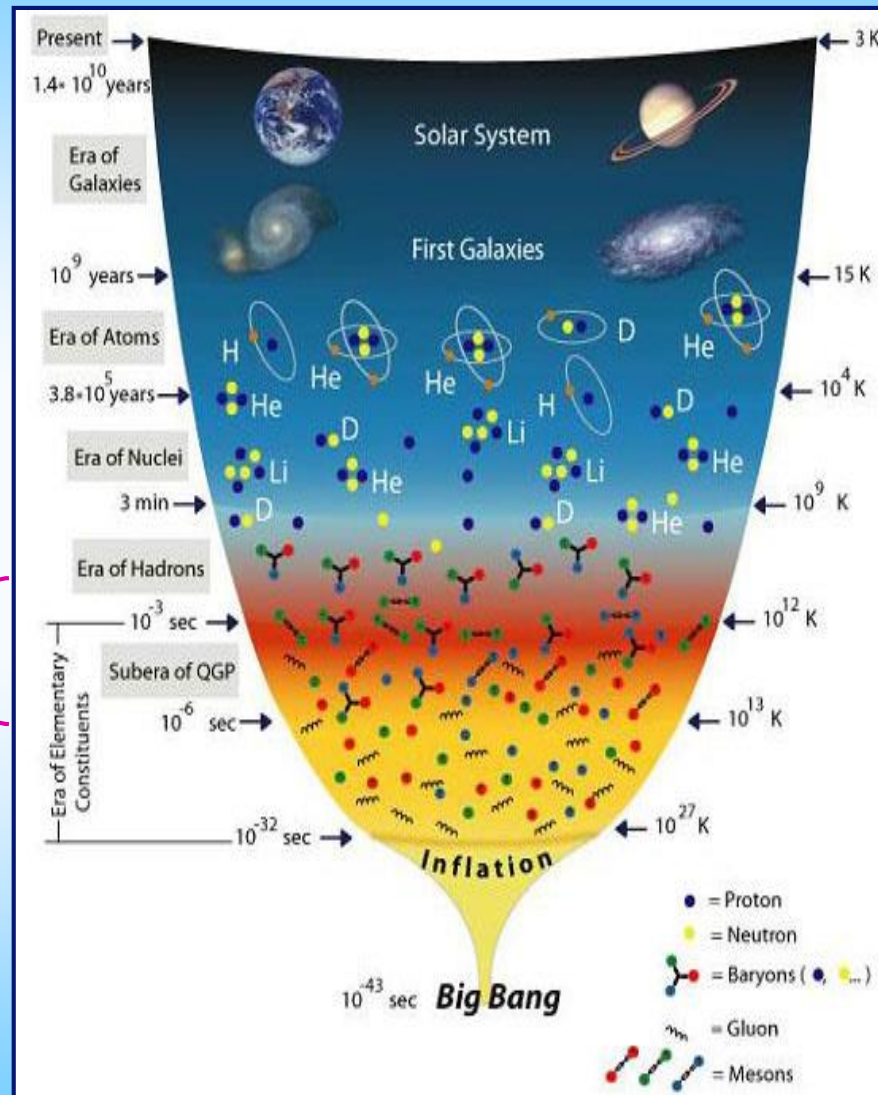
- **Security**;

- **Global file system;**

- **File management, transmission and replication;**

- **Job execution**;

- **Management of application software packages**;

- **Log recording/keeping;**

- **Monitoring**.

The native *AliEn* code is written (by Predrag Buncic and Pablo Saiz) predominantly in *Perl*.

# The evolution of our Universe



**Investigated by ALICE experiment**

# Types of UML diagrams

*Unified Modeling Language (UML) serves for the specification, visualization, construction and documentation of software systems".*

**UML diagrams are grouped into two categories:**

- ***Structure diagram** – Describing the static structure of a system and the relationships among the objects.*

- ***Behaviour diagram** – Describing the dynamic behaviour of a system through the behavior of the objects in a system, including their methods, collaborations, activities and states. The dynamic behavior of a system can be described as a series of changes to the system over time.*

# Creation of Sequence diagrams
## First step. Manual conversion of the Perl code to a special plain text format

### A snippet of original code

```
sub definePackage
{
 ...
 my $lfnDir=lc($self->{CATALOGUE}->{CATALOG}->
GetHomeDirectory()."/packages");
 ...
 while (my $arg=shift){
   if ($arg=~ /^-?-se$/ )     { ... }
   ...
 }
 ...
 $self->{CATALOGUE}->{CATALOG}->isFile($lfn)
  and return;
  ...
}
```

check arguments

### Plain text format

```
AliEn::PackMan.definePackage


AliEn::Catalogue.GetHomeDirectory -> $lfnDir



AliEn::Catalogue.check arguments



AliEn::Catalogue.isFile
 ...
```

# Creation of Sequence diagrams
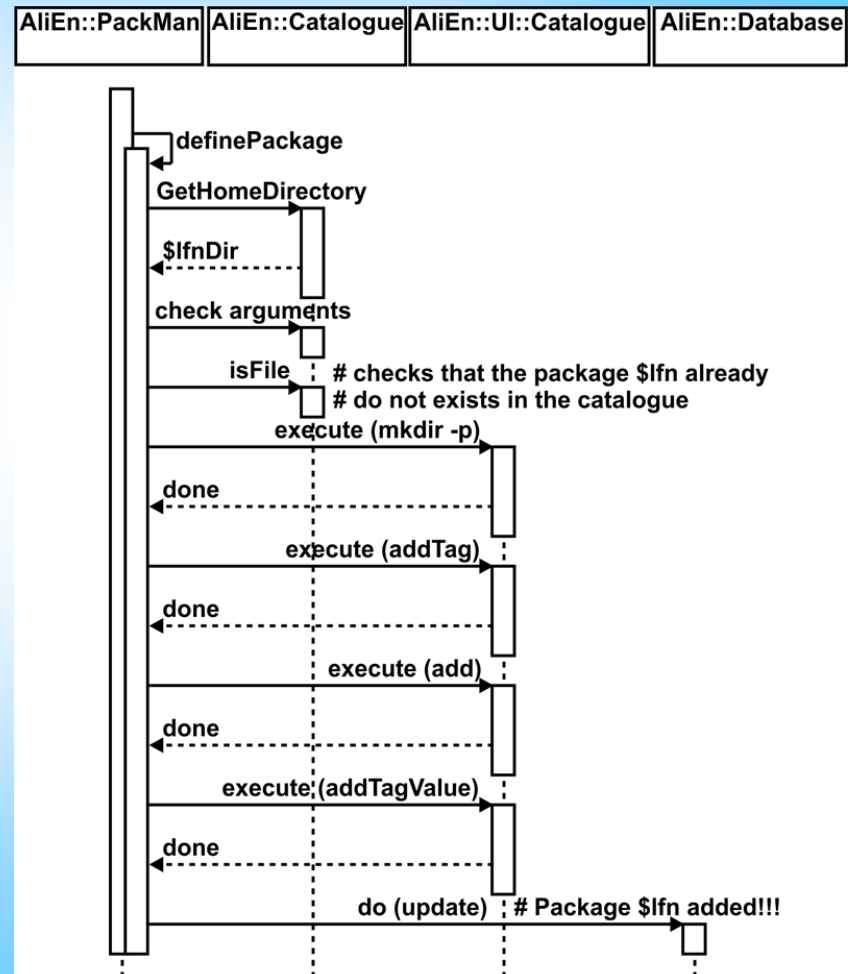## Second step. Passing created file to the UML::Sequence tool

### Plain text

AliEn::PackMan.definePackage
   AliEn::Catalogue.GetHomeDirectory -> $lfnDir
   AliEn::Catalogue.check arguments
   AliEn::Catalogue.isFile
   …

### Diagram generation commands

### UML diagram

1. Bachelor diploma thesis "**Study and presentation of the CERN ALICE experiment Grid infrastructure on the base of UML diagrams**" in two languages - Armenian and English.

   The thesis has been defended **<u>magna cum laude</u>** in the end of **May 2012** and it has been presented by the Administration of SEUA to the National competition for the Armenian President Award in the category of Information Technologies (IT).

The awarding ceremony took place on **4 October 2012**.



***1st prize of the Armenian President as the Best Bachelor Student in the field of IT***

# Student activity (Master study)

MS thesis
*"Development of a service for the optimization of the distribution of data and software packages over Grid infrastructure of ALICE experiment at CERN"*

*The thesis defense will be held in the end of May of this year*

# Publications

| | |
|---|---|
| 1 | A. Abramyan, "**Unified Modeling Language Diagrams for Grid Middleware of CERN ALICE Experiment**". Proc. Int. Conf. "Computer Science and Information Technologies", Yerevan, Armenia, pp. 352-355, 2011. |
| 2 | A. Abramyan et al., "**AliEn Extreme JobBrokering**". Proceedings of Computing in High Energy and Nuclear Physics (CHEP) conference, New York City, NY, USA (2012) |
| 3 | A. Abramyan et al., "*ALICE Environment on the GRID*". Proceedings of Computing in High Energy and Nuclear Physics (CHEP) conference, New York City, NY, USA (2012) |
| 4 | Miguel Martinez Pedreria, ..., A. Abramyan et al. "*Creating a simplified global unique file catalogue*". 20th International Conference on Computing in High Energy and Nuclear Physics (CHEP2013), Amsterdam, Beurs van Berlage (2013) |
| 5 | A. Abramyan, N. Manukyan "*AliEn File Access Monitoring Service – FAMoS* ". Proc. 9th Int. Conf. "Computer Science and Information Technologies", Yerevan, Armenia, pp. 352-355, 2013. |
| 6 | ..., A. Abramyan et al. (ALICE Collaboration), "*Charged-particle multiplicity measurement in proton–proton collisions at √s = 0.9 and 2.36 TeV with ALICE at LHC*". Eur. Phys. J. C (2010) 68: 89–108 |
| 7 | Betty Bezverkhny Abelev,..., A. Abramyan *et al.* (ALICE Collaboration), "*Measurement of quarkonium production at forward rapidity in pp collisions at S√= 7 TeV*". arXiv:1403.3648 [nucl-ex] CERN-PH-EP-2014-042 |
| 8 | Betty Bezverkhny Abelev,..., A. Abramyan *et al.* (ALICE Collaboration), "*K\*(892)^0 and PHI(1020) production in Pb-Pb collisions at sqrt(sNN) = 2.76 TeV*". arXiv:1404.0495 [nucl-ex] CERN-PH-EP-2014-060 |
| 9 | ..., A. Abramyan et al. (ALICE Collaboration), "*Charged-particle multiplicity measurement in proton–proton collisions at √s = 7 TeV with ALICE at LHC*". Eur. Phys. J. C (2010) 68: 345–354 |
| 10 | K. Aamodt, ..., A. Abramyan, *et al.* (ALICE Collaboration), "*Two-pion Bose-Einstein correlations in pp collisions at √s=900 GeV*". Phys. Rev. D 82, 052001 (2010), Volume: 82, Issue: 5, Publisher: American Physical Society, Pages: 1-14 |
| 11 | K. Aamodt, ..., A. Abramyan, et al. (ALICE Collaboration), "*Midrapidity Antiproton-to-Proton Ratio in pp Collisons at √s=0.9 and 7 TeV Measured by the ALICE Experiment*" Phys. Rev. Lett. 105, 072002 (2010) [12 pages] |
| 12 | K Aamodt, ..., A Abramyan, et al. (ALICE Collaboration), "*Transverse momentum spectra of charged particles in proton–proton collisions at √s=900 GeV with ALICE at the LHC*". Physics Letters B 693 (2010) 53–68 |