# Science Reproducibility and Reusability with FutureGateway and a Zenodo-like repository: the PALMS experiment

**R. Bruno**[1*], R. Barbera[1,2], M. Fargetta[1], R. Rotondo[1], A. Anagnostou[3], S. J. Taylor[3]

1    Italian National Institute of Nuclear Physics, Division of Catania, Italy

2    Department of Physics and Astronomy "E. Majorana" of the University of Catania, Italy

3    Modelling & Simulation Group, Department of Computer Science, Brunel University London, UK

*    riccardo.bruno@ct.infn.it

# Driving considerations

INFN

"*Open Science refers to a scientific culture that is characterized by its openness. Scientists **share results** almost immediately and with a very wide audience*" [1]
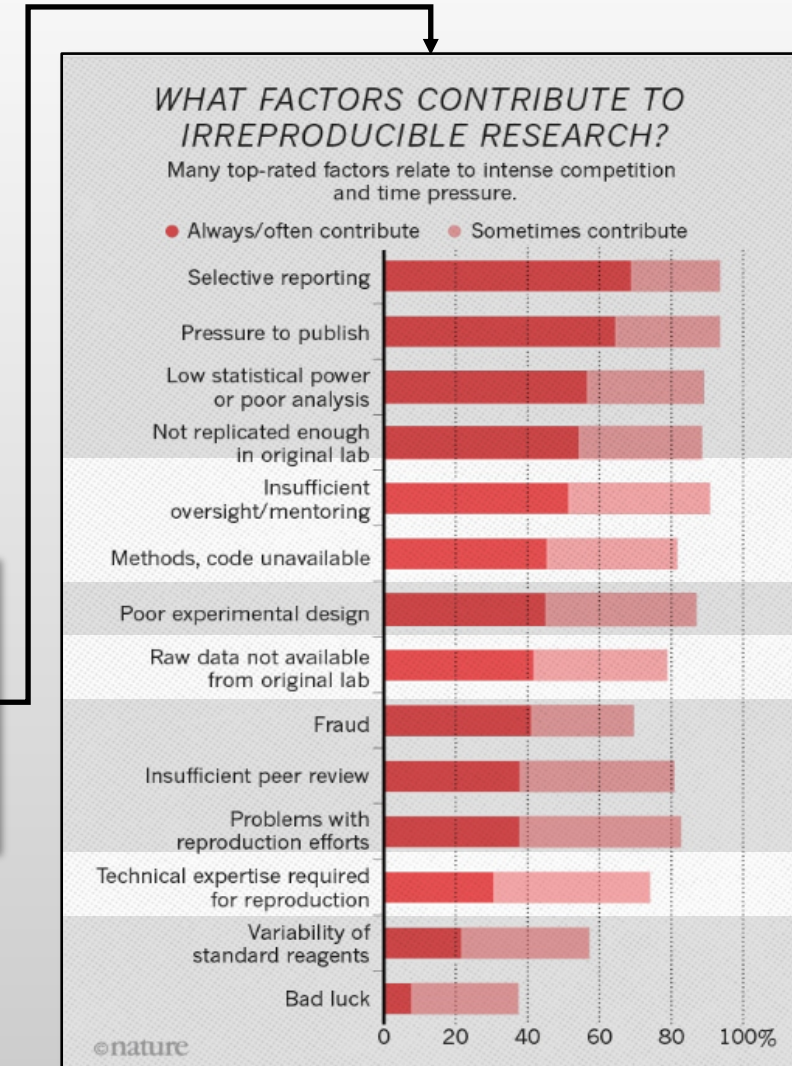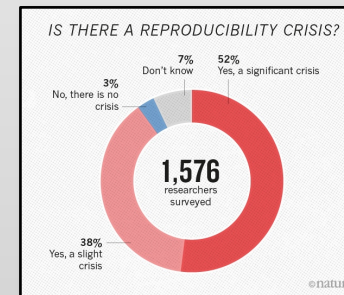
"*Open science is a means and not an end in itself and it is much more than just open access to publications or data; it includes many aspects and stages of research processes thus enabling full **reproducibility and re-usability** of scientific results.*" [2]

Reality check on reproducibility (Survey) [3]
- Insufficient oversight/mentoring
- Methods, code unavailable
- Raw data not available from original lab
- Technical expertise requires for reproducibility



IS THERE A REPRODUCIBILITY CRISIS?

7% Don't know
52% Yes, a significant crisis
3% No, there is no crisis
1,576 researchers surveyed
38% Yes, a slight crisis

WHAT FACTORS CONTRIBUTE TO IRREPRODUCIBLE RESEARCH?
Many top-rated factors relate to intense competition and time pressure.
● Always/often contribute   ● Sometimes contribute

- Selective reporting
- Pressure to publish
- Low statistical power or poor analysis
- Not replicated enough in original lab
- Insufficient oversight/mentoring
- Methods, code unavailable
- Poor experimental design
- Raw data not available from original lab
- Fraud
- Insufficient peer review
- Problems with reproduction efforts
- Technical expertise required for reproduction
- Variability of standard reagents
- Bad luck

0   20   40   60   80   100%

©nature

[1] Opening Science: The Evolving Guide on How the Internet is Changing Research (DOI: 10.1007/978-3-319-00026-8)
[2] Making Open Science a Reality: http://dx.doi.org/10.1787/5jrs2f963zs1-en
[3] Reality check on reproducibility: https://www.nature.com/news/1-500-scientists-lift-the-lid-on-reproducibility-1.19970

# Key actors (Overview)

Software

FGSG

Science Gateway

Portal

SSOD

**R**eprodiciblity and
**R**eusabiltiy platform

Experiment

**PALMS**

DAMS

**O**utput

**I**nput

**Software**

PALMS web application

?

Reproducibility and reusability of Physical Activity Lifelong Modelling & Simulations
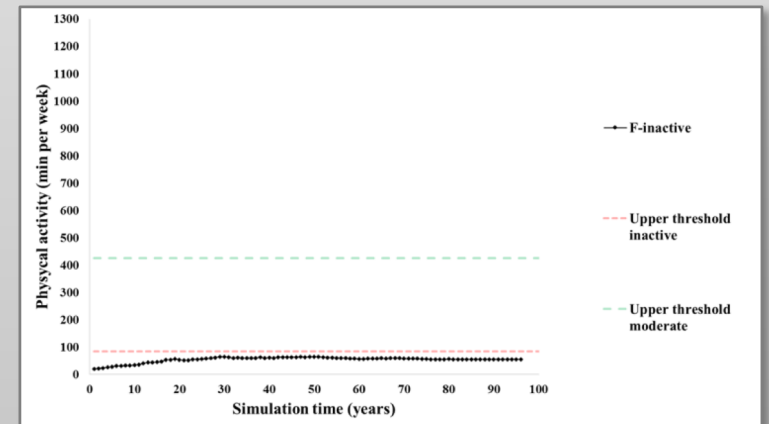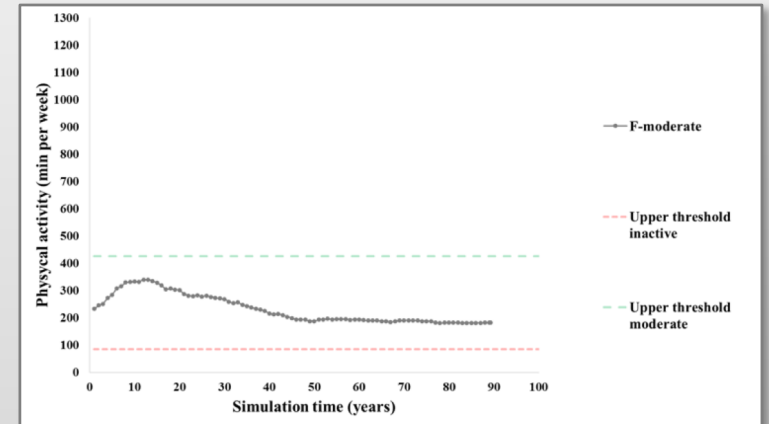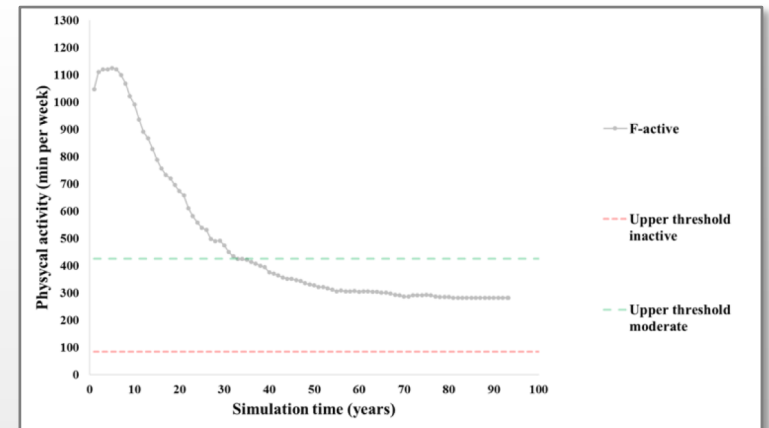
Digital Asset Management Systems (DAMSes) are *"intertwined structures incorporating both software and hardware that take care of management tasks and decisions surrounding the ingestion, annotation, cataloguing, storage, retrieval and distribution of digital assets"*

REPAST

FTP
HTTPD

docker compose

3

# PALMS

- **P**hysical **A**ctivity **L**ifelong **M**odelling & **S**imulations
  - Is an agent based micro-simulation that predicts the lifelong physical activity behaviour of a population taking into account individual characteristics and their effect on physical activity over time
  - Produces individual and aggregated quantitative outputs for quality of life and health conditions related costs

- The software
  - Uses REPAST [4] an open source agent-based modeling and simulation platform
  - A specific dockerhub image exists for PALMS executions (osabuoun/repast) [5]
  - Two inputs necessary: model file (REPAST) and a parameters' file

[4] https://repast.github.io
[5] https://hub.docker.com/r/osabuoun/repast

# FutureGateway Framework

INFN software project aiming to build secure and reliable Science Gateways [6]
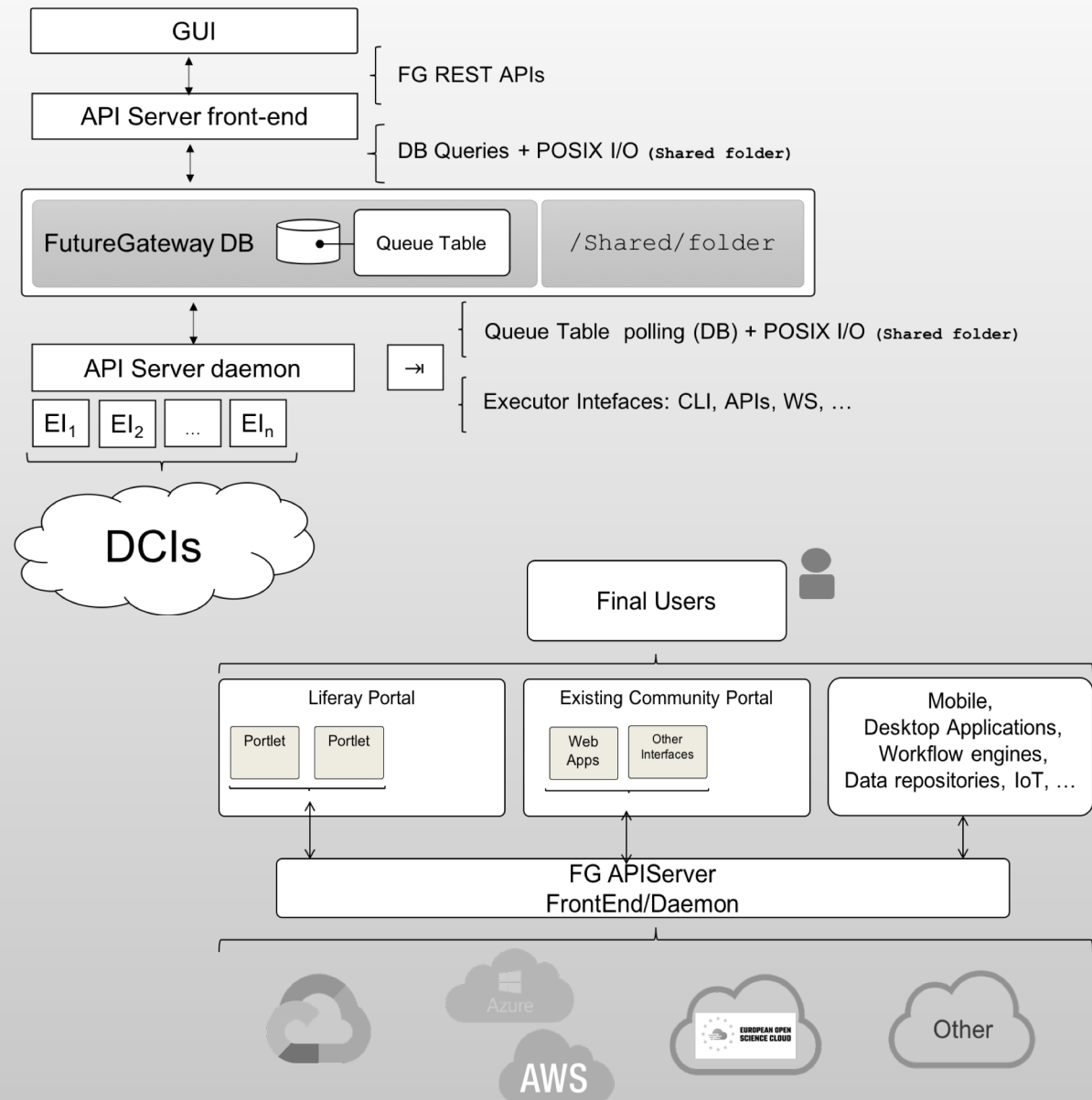
**Three core components:**

- Database, APIServer front-end, APIServer daemon + Executor Interfaces

**The framework:**

- Core components are enriched with a suite of tools, APIs and installation + maintenance scripts

- Open Source code available on GitHub

**Targets:**

- Desktop and Mobile applications, Workflow Engines, IoT and **Open Science**

[6] https://www.xsede.org/ecosystem/science-gateways



**Architecture**
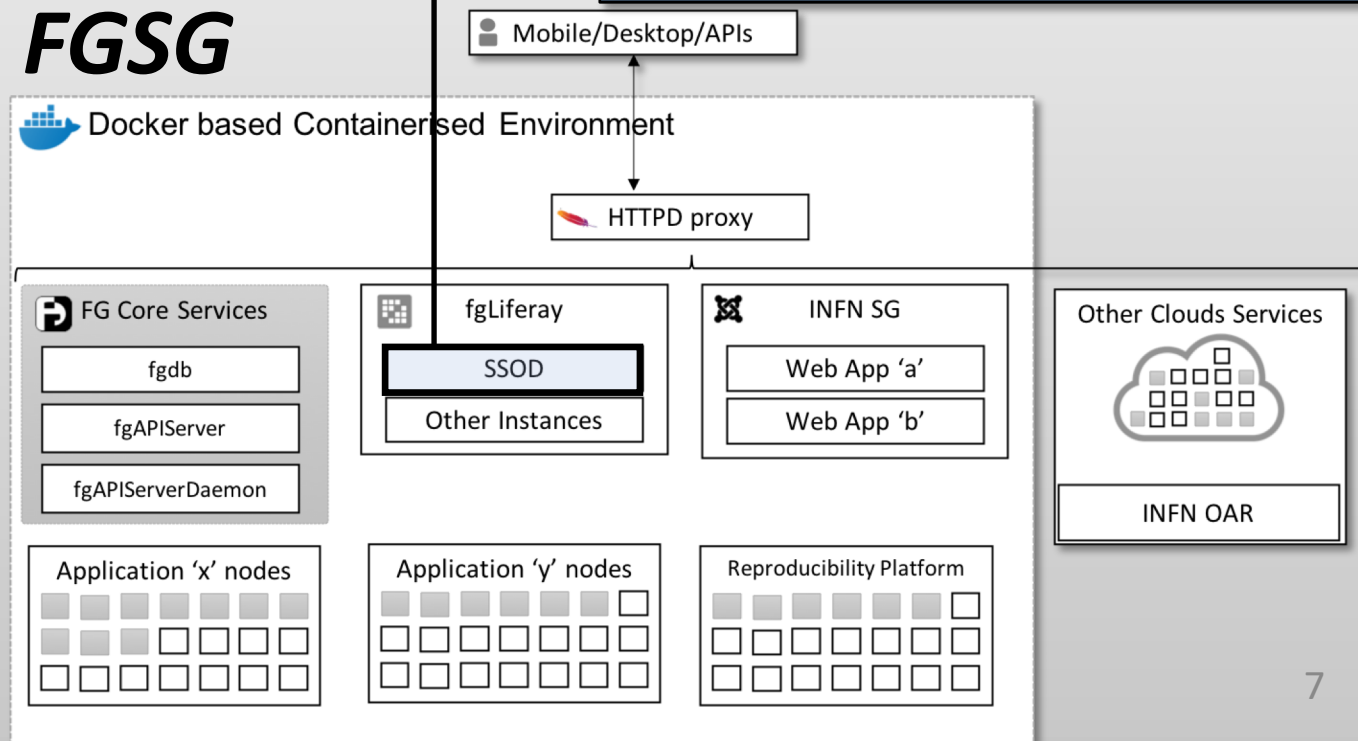
**Usage scenarios**

5

# FG APIs

- FG uses a set of RESTful APIs to perform operations on Distributed Computing Infrastructures (DCI)
- APIs are splitted in three families:
  - **IAT**, Specifications [7]
    - Infrastructures – Specify how to target the DCI
    - Applications – Specify which kind of activity perform on the DCI
    - Tasks – Application instances, relative statuses and output retrieval
  - **UGR**, Documented
    - Users – FG manages its own set of API users
    - Groups – Membership can be grouped to assign different set of Roles
    - Roles – Are assigned to groups and specify the operations allowed to perform
  - **AAA**, Design
    - Authorisation Auditing and Accounting (not yet available, db level functionality)

[7] https://fgapis.docs.apiary.io/#

# *FGSG –* FutureGateway based Science Gateway

Fully docker containerised environment built in the context of the EOSC-hub project [8], to provide a General purposes Science Gateway: the EGI [9] **Science Software on Demand** (SSOD) [10]

- The system allows to dynamically instantiate and destroy docker containers (it supports docker compose as well as docker swarm)

- FG core services + SSOD service

- SSOD service powered by an enterprise portal framework (Liferay)
  - One section dedicated to the **Reproducibility & Reusability platform**
  - The platform exploits the FutureGateway and the INFN Open Access Repository

[8] https://www.eosc-hub.eu/
[9] https://www.egi.eu/
[10] https://fgsg.egi.eu/egissod/web/ssod/

**SSOD**

**FGSG**



7

# R&R for PALMS in SSOD

- Web application making use of FG REST APIs to reproduce and reuse PALMS results.

- Usage:
  1. Specify a PALMS output DOI registered in the INFN Open Access Repository -> 'Prepare'
  2. Specify PALMS input DOIs (model+prameters). They can also be HTTP URLs
  3. Execute PALMS with input specified at 2.
  4. User can access PALMS reproduced outputs; User can upload modified input files (reusability)

# User Auditing and Accounting in SSOD

- Only registered users can access the application





- SSOD uses EGI Check-in Proxy [11] (Federated Authentication)

  - It supports many authentication sources available in EU and supports Social Networks accounts as well

  - User activity performed by users is tracked by the FG with AAA API set (currently querying the DB)

  - Once logged the access to the application has to be authorized (not at the moment)

[11] https://www.egi.eu/services/check-in/

# INFN OAR
## Open Access Repository



- INFN joined the Plan S [12] initiative to promote open access

- INFN OAR is a DAMS hosted at INFN Catania [13] running on a dedicated Kubernets cluster

- It uses Zenodo [14] open source DAMS software

- PALMS input, output, software and papers files are registered with a referencing DOI

- INFN OAR allows to define references among registered DOIs and supports DOI versioning

- Software used by *FGSG* to run PALMS simulations is published as well on the INFN OAR (linked from GitHub)

[12] https://www.coalition-s.org
[13] https://www.openaccessrepository.it
[14] https://github.com/zenodo/zenodo

# Summary and conslusions

- PALMS use case has a quite simple structure
- The R&R platform developed for PALMS can be easily used to host other use cases
- OAR + R&R Platform are open to any interested community
- FutureGateway is a mature product successfully used by:
  - Desktop and Mobile Applications, Workflow engines, IoT and **OpenScience**
- Investigations are in progress to extend this work to a more general and widely usable solution

"Reusing" Reprodiciblity and Reusabiltiy platform

# Further info

# FAIR Principles in Open Science

**As expressed by GOFAIR , to be compliant with the Open Science paradigm, research data should be FAIR:**

- **F**indable

  **F1. (Meta)data are assigned a globally unique and persistent identifier**
  **F2. Data are described with rich metadata (defined by R1 below)**
  **F3. Metadata clearly and explicitly include the identifier of the data they describe**
  **F4. (Meta)data are registered or indexed in a searchable resource**

- **A**ccessible

  **A1.1 The protocol is open, free, and universally implementable**
  **A1.2 The protocol allows for an authentication and authorisation procedure, where necessary**

- **I**nteroperable

  **I1. (Meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.**
  **I2. (Meta)data use vocabularies that follow FAIR principles**
  **I3. (Meta)data include qualified references to other (meta)data**

- **R**eusable

  **R1.1. (Meta)data are released with a clear and accessible data usage license**
  **R1.2. (Meta)data are associated with detailed provenance**
  **R1.3. (Meta)data meet domain-relevant community standards**

# The «pillars» of the Scientific Method

- **Repeatability**
  - The closeness of agreement between independent results obtained with the same method on identical test material, under the same conditions (same operator, same apparatus, same laboratory and after short intervals of time)
  - Affected by *random errors*

- **Reproducibility**
  - The closeness of agreement between independent results obtained with the same method on identical test material but under different conditions (different operators, different apparatus, different laboratories and/or after different intervals of time)
  - Affected by *systematic errors*

# The (Re-)'s of Science

**Repeat**

**Replicate**

Same experiment
Same laboratory

Same experiment
Different laboratory

**Test**

**Reproduce**

**Reuse**

Same experiment
Different setup

Different experiment
Some of same

# Open Science (a receipe rather than a definition)

Open science and research refers to effort to promote open procedures in scientific research activities. The key objective is, in the context set by research ethics and legal frameworks, to publish research outputs (research publications, research data, research methods) so the can be examined and user by any interested party.

Open science and research involves practices, such as promoting open access to research publications, open availability fo research data, harnessing open source software and open standards, and open documentation of research process.



Source: Davis, F. D. (1989), "Perceived usefulness, perceived ease of use, and user acceptance of information technology», MIS Quarterly, **13(3)**: 319-340

# Concepts and definitions

- **Open Access Repositories** (OARs) are powered by **Digital Asset Management Systems** (DAMSes), which are "intertwined structures incorporating both software and hardware that take care of management tasks and decisions surrounding the ingestion, annotation, cataloguing, storage, retrieval and distribution of **digital assets**"

- A **digital asset** in essence is "anything that exists in a binary format and comes with the right to use"

- "Types of digital assets include, but are not exclusive to, photography, logos, illustrations, animations, audio-visual media, presentations, spreadsheets, Word and/or PDF documents, data and a multitude of other digital formats and their respective metadata"

# FutureGateway Architecture

# API Server front-end



Front-End

REST

GUI

Command

Response

AuthN/Z check

Process the command

Enqueue command

Prepare response

Id

Queue

UGR configuration or external AAI systems (PTV)

Shared Folder

- Operations
    - GUIs send a command via REST APIs
    - The 'command' may contain a JSON stream specifying command parameters
    - The Front-End first check for requestor Authorization and Authentication eventually using UGR configuration or external AAI mechanisms using PTV information
    - The command is processed querying and/or updating the DB accordingly and/or updating the shared folder
    - Commands to be finalized by the APIServer daemon, are stored in the queue table
    - Command output is returned back into the response as a JSON stream

# API Server daemon



- Operations
  - Commands (Tasks=Command(Action,EI)) are extracted from the front-end queue
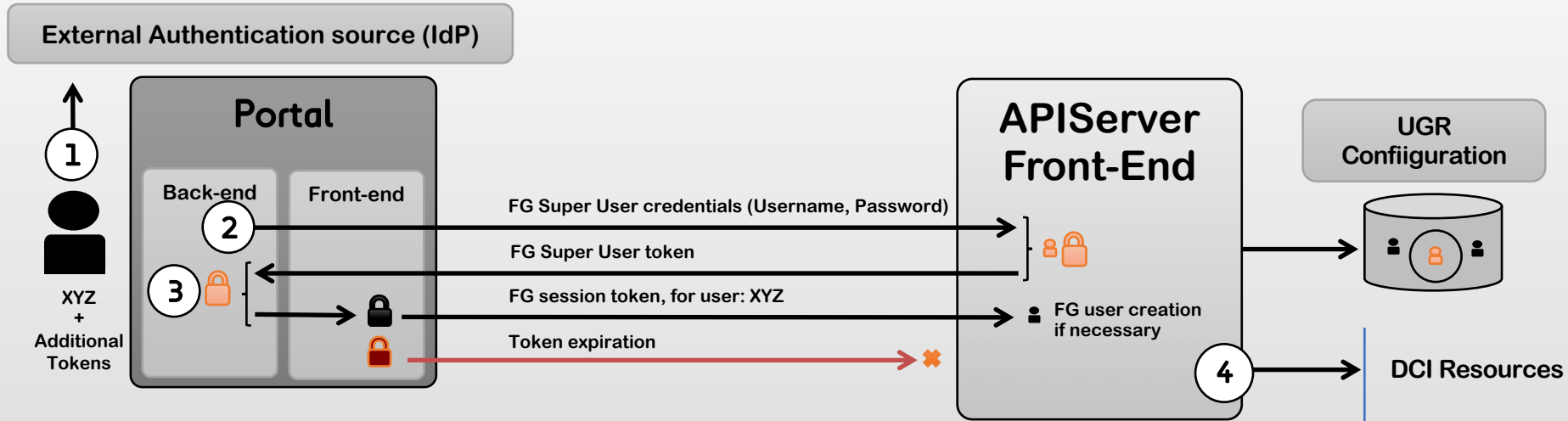  - Each 'command' contains the 'Target Executor' field which specifies the Executor Interface name
  - Executor interfaces are dynamically instantiated by the APIServer by its name, applying the specified action on DCI
  - Other queue daemons may extract commands from the queue having their own EIs implemented.
  - New EIs can be easily developed just inheriting a standard Executor Interface class
  - Current available EIs:
    - GridEngine (A core component of the CSGF using JSAGA and targeting: ssh, rOCCI and wms)
    - ToscaIDC (INDIGO-dc orchestrator with IAM)

# FutureGateway AuthN/Z
## (UGR: user-behalf)



**1** **Log-In**
The portal authenticate users using an external Authorization source (IdP). Once authenticated it receives a unique user name (XYZ) + eventually additional tokens.

**2** **FG Super user access token**
With UGR API calls, it is possible to obtain a session token having FG super user capabilities. This operation should be done in a protected environment such as a portal back-end environment.

**3** **User session token**
The UGR APIs can obtain a user session token on behalf of another user, only super-user can do this operation (see role: user_impersonate).
During this opertation, the back-end Portal should create a new FG user if not yet present.
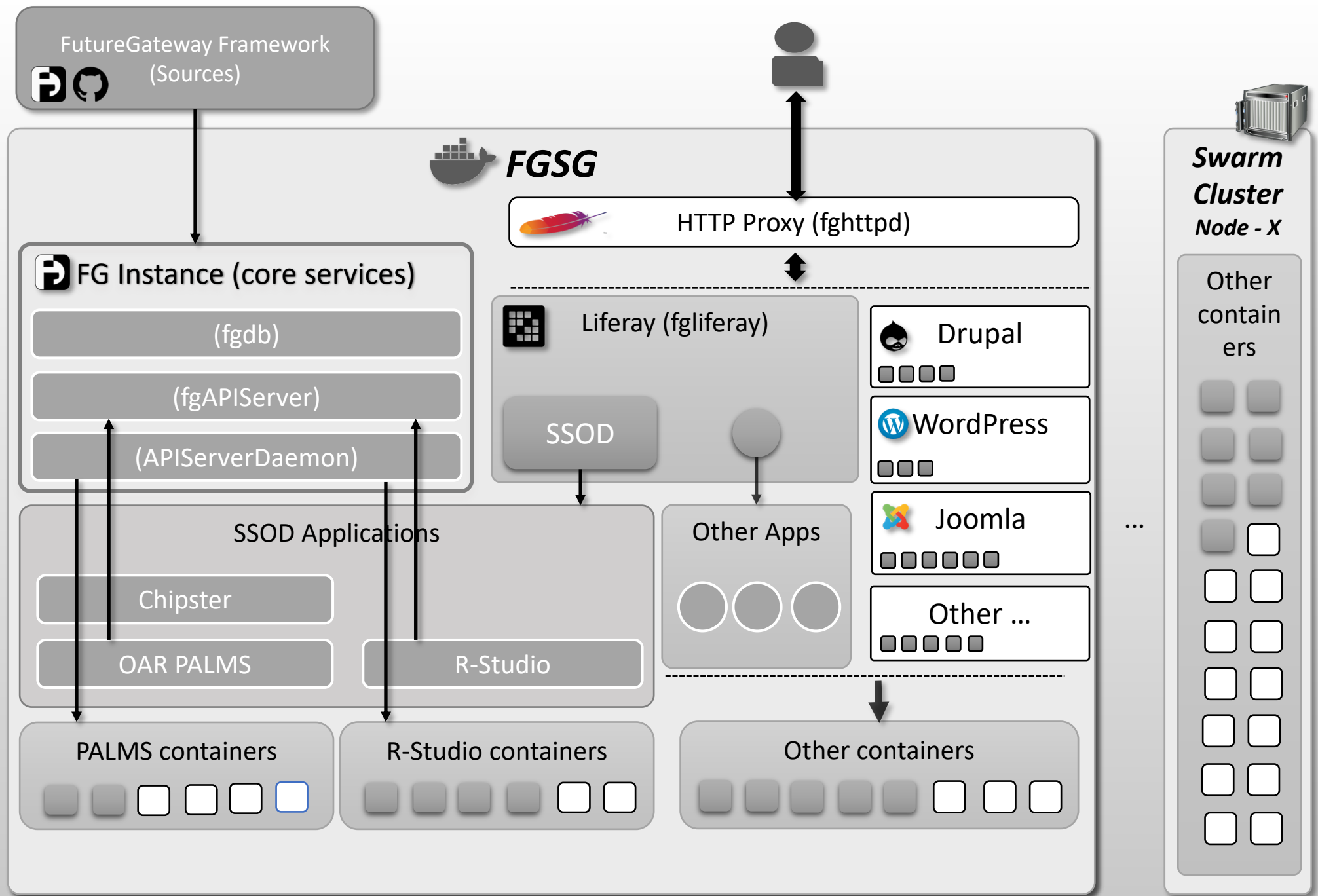
**4** **Additional tokens**
During the external authentication, additional tokens may be retrieved. This information can be used as FG username or saved using URG APIs. This kind of tokens may be needed by FG applications to access DCIs resources (See INDIGO-dc IAM).

**INFN**

# FGSG
## Structure

IBM Blade center-X
One slot with 32 cores
128 GB RAM
Ubuntu 18.04
Docker version
17.12.1-ce
150 GB on board
RAID-2
5TB ext. Storage

FutureGateway Framework (Sources)

**FGSG**

HTTP Proxy (fghttpd)

FG Instance (core services)

(fgdb)

(fgAPIServer)

(APIServerDaemon)

Liferay (fgliferay)

SSOD

Drupal

WordPress

Joomla

Other ...

SSOD Applications

Chipster

OAR PALMS

R-Studio

Other Apps

PALMS containers

R-Studio containers

Other containers

*Swarm Cluster Node - X*

Other containers

...

# INFN OAR Deployment Structure