# Federated GMQL
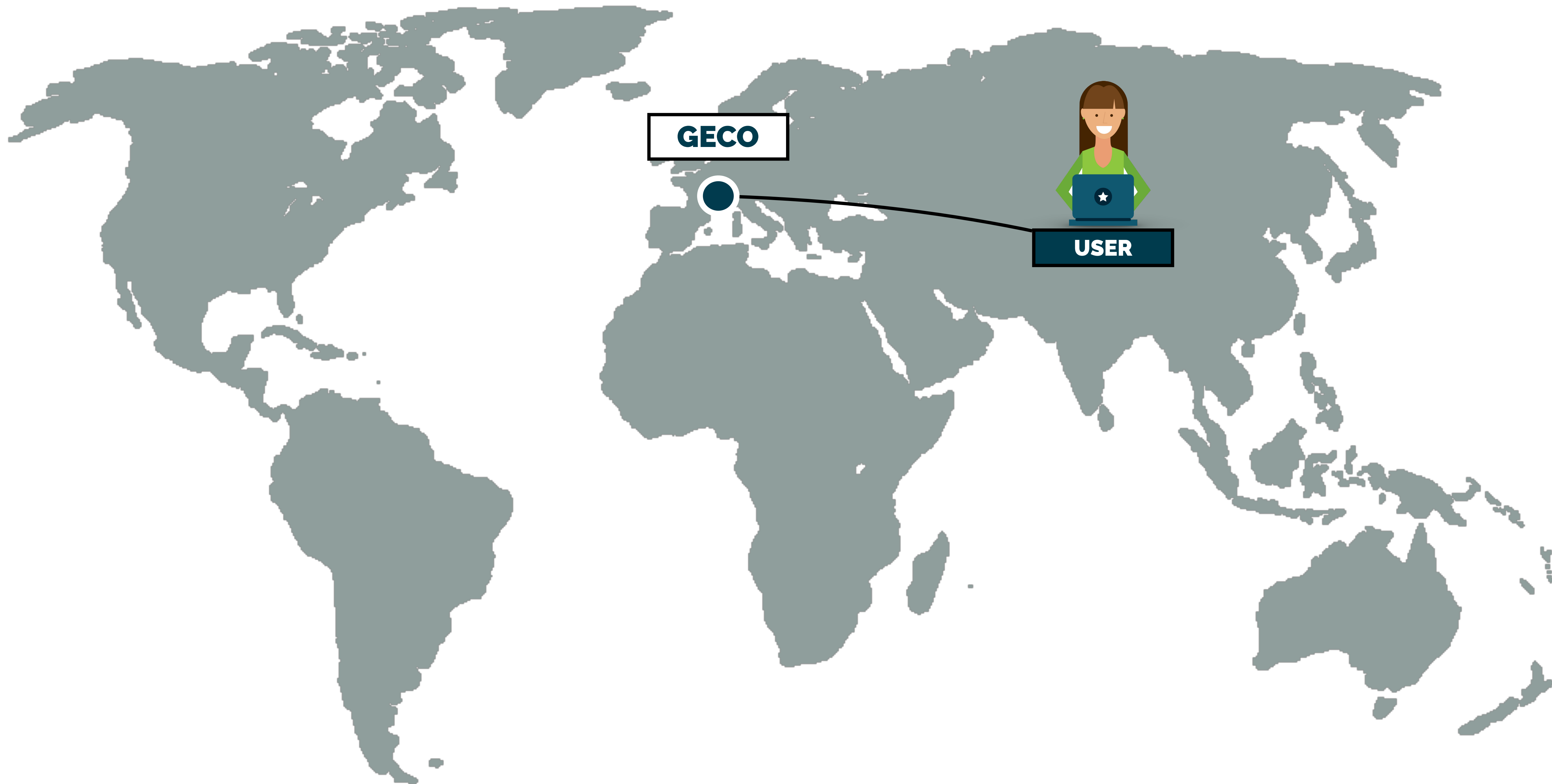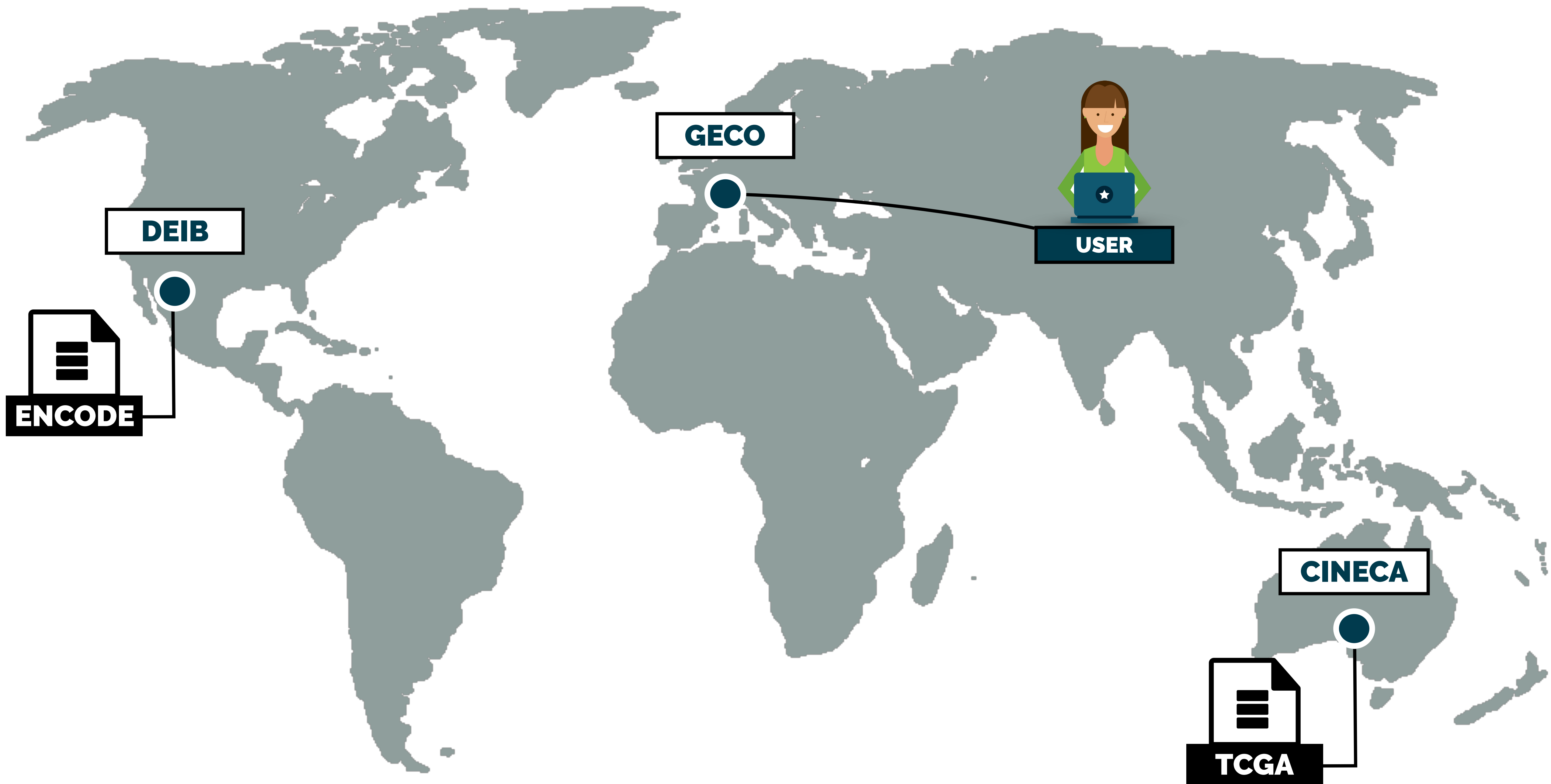
OVERVIEW

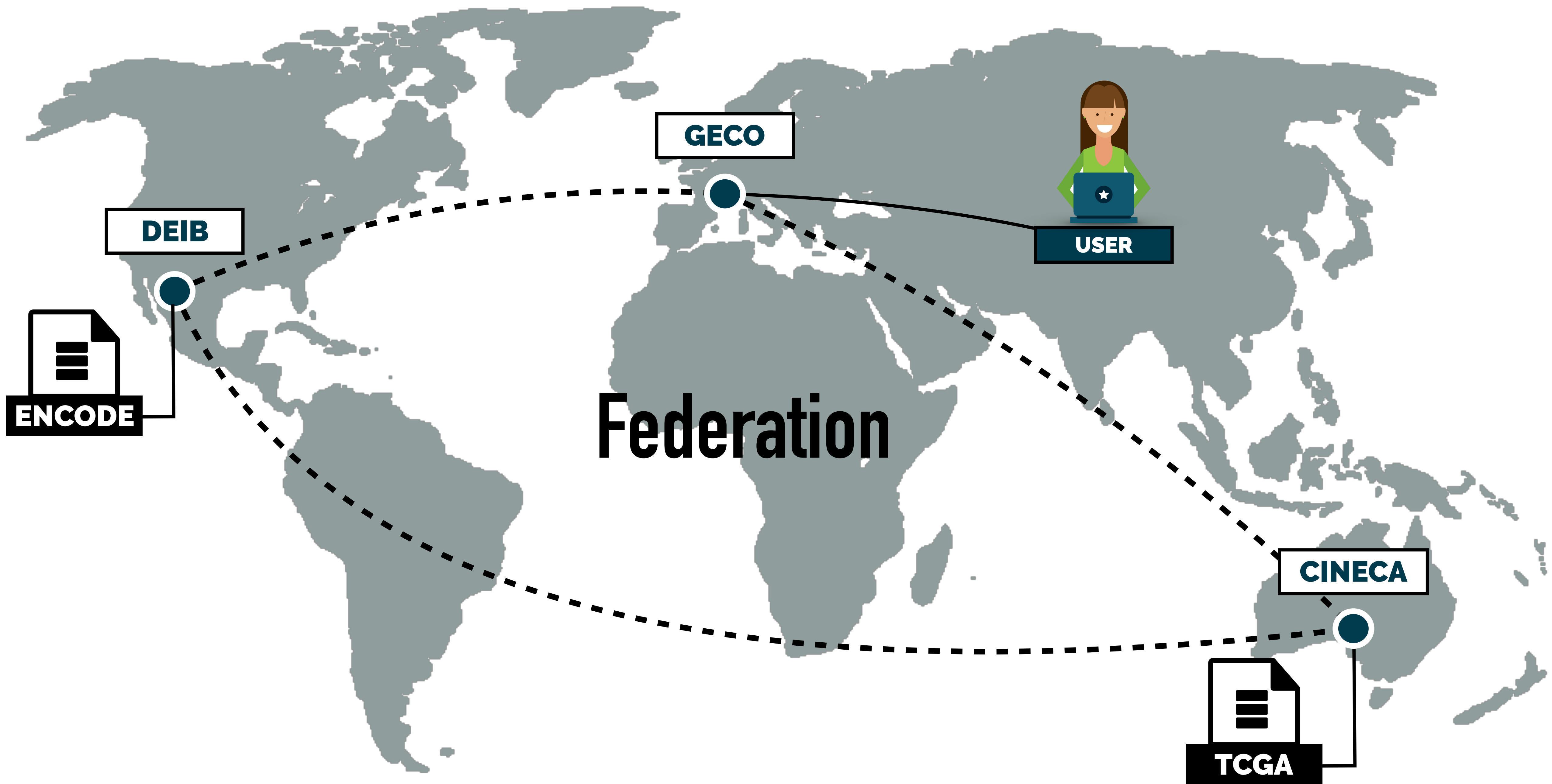# GenoMetric Query Language (GMQL)

- Open Source Project

- Enables querying hundreds of datasets and thousands of samples

- Runs on a cloud-computing system based on Apache Spark

- Publicly available through a user-friendly web interface

- Well maintained public repository which integrates genomic datasets from ENCODE, TCGA/GDC and Roadmap Epigenomics

DEIB

ENCODE

GECO

USER

CINECA

TCGA

GECO

DEIB

USER

ENCODE

Federation

CINECA

TCGA

**Federated Datasets**

# Federated Query

```
# SELECT TCGA AT CINECA
TCGA = SELECT(tumor_tag == "BRCA"; at:CINECA) CINECA.TCGA;

# SELECT ENCODE AT DEIB
ENCODE = SELECT(cell_line == "H1-hESC"; at:DEIB) DEIB.ENCODE;

# JOIN DS1 AND DS2 AT DEIB
JOINED = JOIN(dist < 0; output: left_distinct; at:LOCAL) TCGA ENCODE;

# SELECT MUTATION AT GeCo
MYMUTATION = SELECT(at:LOCAL) MUTATION;

# Map MUTATION AT GeCo
RES = MAP(count_name: mut_count; at:LOCAL) JOINED MYMUTATION;

# Materialize at GeCo
MATERIALIZE RES INTO ResGenes;
```

# Federated Query

```
# SELECT TCGA AT CINECA
TCGA = SELECT(tumor_tag == "BRCA"; at:CINECA) CINECA.TCGA;

# SELECT ENCODE AT DEIB
ENCODE = SELECT(cell_line == "H1-hESC"; at:DEIB) DEIB.ENCODE;

# JOIN DS1 AND DS2 AT DEIB
JOINED = JOIN(dist < 0; output: left_distinct; at:LOCAL) TCGA ENCODE;

# SELECT MUTATION AT GeCo
MYMUTATION = SELECT at:LOCAL MUTATION;

# Map MUTATION AT GeCo
RES = MAP(count_name: mut_count at:LOCAL) JOINED MYMUTATION;

# Materialize at GeCo
MATERIALIZE RES INTO ResGenes;
```

SELECT(cell_line == "H1-hESC") ENCODE

GECO

DEIB

SELECT(tumor_tag == "BRCA") TCGA

USER

CINECA

POLITECNICO DI MILANO

JOIN ... SELECT ... MAP ... MATERIALIZE

GECO

DEIB

USER

CINECA

POLITECNICO DI MILANO

GMQL

GMQL    GMQL-REST    Demo Video    Documentation    Example Queries ▾    GeCo

Hello Guest    Logout

## Datasets ⚙

0.0%

☑ ➖ Private
➕ Public
➕ Federated

⊕ Add    🗑 Delete    ⊘ Download    ☁ UCSC

## Query editor ⚙    Select query ▾    🗑

```
1   # SELECT TCGA AT CINECA
2   TCGA = SELECT(tumor_tag == "BRCA"; at:CINECA) CINECA.HG19_TCGA_rnaseqv2_gene;
3
4   # SELECT ENCODE AT DEIB
5   ENCODE = SELECT(cell_line == "H1-hESC"; at:DEIB)  DEIB.HG19_ENCODE_NARROW_2019_01;
6
7   # JOIN DS1 AND DS2 AT DEIB
8   JOINED = JOIN(dist < 0; output: left_distinct; at:LOCAL) TCGA ENCODE;
9
10  # SELECT MUTATION AT GeCo
11  MYMUTATION = SELECT(at:LOCAL) MUTATION;
12
13  # Map MUTATION AT GeCo
14  RES= MAP(count_name: mut_count; at:LOCAL) JOINED MYMUTATION;
15
16  # Materialize at GeCo
17  MATERIALIZE RES INTO ResGenes;
18  |
```
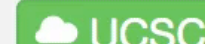
**Query name**    queryname

**Output format**    ◉ Tab delimited    ○ GTF

≣ Show jobs    RUNNING ■    ⚙ Compile    ▶ Execute

## Metadata browser

Copy

```
DATA_SET_VAR = SELECT()
    queryname_20190913_130954_ResGenes;
```

➕ New condition    ▤ Test    ▤ Download

## Sample metadata

## Schema

DEIB - POLITECNICO DI MILANO

# Federated GMQL

SETUP

# Installing a GMQL Instance

- ## Using **Docker**

  Follow the  guide at : https://github.com/DEIB-GECO/GMQL-Docker

- ## Full Installation

  Follow the  guide at  https://github.com/DEIB-GECO/GMQL-WEB/wiki/installation

# Name Server

Centralized component of the Federated System that:

- Allows joining the Federation

- Manages:

  - Instances and their authentication mechanisms

  - Federated Datasets and their privacy

  - Groups

Publicly available at : http://genomic.elet.polimi.it/nameserver/

# Joining the Federation

- Register your instance on the Name Server
- Login
- Add the required configuration to your repository.xml file

Name Server

# Name Server

## Register a new Instance

Username (instancename)    MyInstance

Email    myemail@myinstitute.com

Description    My Instance

Password    ••••••••

Repeat Password    ••••••••

GMQL API URL    http://myserver.com/gmql-rest/

Login    **Register**

Name Server

# Name Server

## Login

Username    myInstance

Password    ••••••••

Register    **Login**

**Name Server**   Home   Datasets   Instances   Groups   Logout

# 👤 myInstance

Your *token* for communication with the NameServer is:

| 7cadeca22c38b9321f1029c267938e7342e74510 | 👁 |
|---|---|

You should have the following configuration in your repository.xml file:

```
<property name="GF_ENABLED">true</property>
<property name="GF_NAMESERVER_ADDRESS">http://genomic.elet.polimi.it/nameserver</property>
<property name="GF_INSTANCENAME">myInstance</property>
<property name="GF_TOKEN">7cadeca22c38b9321f1029c267938e7342e74510</property>
```

Your instance URL is the following:

| http://myserver.com/gmql-rest/ | Save |
|---|---|

## Name Server    Home   Datasets   Instances   Groups

# 👤 myInstance

Your *token* for communication with the NameServer is:

```
7cadeca22c38b9321f1029c267938e7342e74510
```
👁

You should have the following configuration in your repository.xml file:

```xml
<property name="GF_ENABLED">true</property>
<property name="GF_NAMESERVER_ADDRESS">http://genomic.elet.polimi.it/nameserver</property>
<property name="GF_INSTANCENAME">myInstance</property>
<property name="GF_TOKEN">7cadeca22c38b9321f1029c267938e7342e74510</property>
```

Your instance URL is the following:

```
http://myserver.com/gmql-rest/
```
Save

# Federating a Dataset

- **Login to the Name Server**

- **Go to the Datasets section**

- **Click on Add a new dataset**

- **Set dataset details and privacy**
    - **The dataset name must be identical to the name of the dataset as it appears in the public repository of the instance**
- **Save the new dataset**

Name Server

Name Server    Home    Datasets    Instances    Groups    Logout

## Add a new dataset

| | |
|---|---|
| **Name** | MYDATASET |
| **Author** | ENCODE Project |
| **Description** | ENCODE data mapped to HG19 human genome. Broad (or Regions) Peaks format is selected and archived/revoked data are avoided. The release date of this dataset is August 2017. https://www.encodeproject.org/ |

**Privacy**

Available instances:

GMQL-ALL (Group)
admin
genomic
canakoglu
imn_gqml_instance
andreagulino
luca_portatile

>

Selected instances:

myInstance
GMQL-ALL                    remove

**Repositories**

Available instances:

admin
genomic
canakoglu
imn_gqml_instance
andreagulino

>

Selected instances:

myInstance

Name Server    Home    Datasets    Instances    Groups    Logout

## Add a new dataset

**Name**         MYDATASET

**Author**       ENCODE Project

**Description**  ENCODE data mapped to HG19 human genome. Broad (or Regions) Peaks format is selected and archived/revoked data are avoided. The release date of this dataset is August 2017. https://www.encodeproject.org/

**Privacy**

Available instances:

GMQL-ALL (Group)
admin
genomic
canakoglu
imn_gqml_instance
andreagulino
luca_portatile

Selected instances:

myInstance
GMQL-ALL                    remove

**Repositories**

Available instances:

admin
genomic
canakoglu
imn_gqml_instance
andreagulino

Selected instances:

myInstance

DEIB - POLITECNICO DI MILANO

# Federated GMQL

SETUP