# GENE-SWitCH

**The regulatory GENomE of SWine and CHicken: functional annotation during development**

## Deliverable D3.3
## Fully functional data coordination platform for production use

**Deliverable leader:** EMBL
**Authors:** Peter Harrison (EMBL-EBI)

**Version: 1.1**

| Due date of deliverable (as in DoA): | M18 |
| --- | --- |
| Actual submission date: | 24/06/2021, month M24 |

**Dissemination level:**

| | |
| --- | --- |
| **PU** Public | X |
| **CO** Confidential, only for members of the consortium (including the Commission Services) | |

Research and Innovation Action, SFS-30-2018-2019-2020 Agri-Aqua Labs
Duration of the project: 01 July 2019 – 30 June 2023, 48 months

# Table of contents

# 1 Summary

Within this website (DEC) report we detail the fully functional FAANG Data Coordination platform that has been developed for the production, submission and presentation of GENE-SWitCH datasets. This report will provide the website links, screenshots and descriptions of the different components that comprise the Data Coordination platform. The platform is in production use with, at the time of this report, eleven datasets in the GENE-SWitCH project page of the coordination platform (https://data.faang.org/projects/GENE-SWitCH) that comprise an array of different assay types for both chickens and pigs.

# 2 The FAANG Data Coordination Platform

The FAANG Data Coordination Centre was established at EMBL in 2015 to support the data coordination, archival and presentation of data generated by the global Functional Annotation of Animal Genomes (FAANG) project. The Data Coordination Platform has now been completely overhauled to support the data submission and presentation requirements of GENE-SWitCH. This report documents the different components that comprise the FAANG Data Coordination Platform (https://data.faang.org/) that is in production use by the project.

## 2.1 Metadata rulesets

**https://data.faang.org/ruleset/samples**
**https://data.faang.org/ruleset/experiments**
**https://data.faang.org/ruleset/analyses**

The FAANG metadata standards are a key component of interoperability in the FAANG project, ensuring that data recording is standardised across FAANG records regardless of which consortium member and in what location data is generated. The metadata standards have undergone an in-depth review and update in response to the new requirements of the EuroFAANG projects (see https://www.gene-switch.eu/eurofaang.html), with updated rulesets published for FAANG in the underlying GitHub and portal presentation pages (https://data.faang.org/ruleset/samples). The metadata pages provide a complete overview of the FAANG requirements for sample, experiment, and analysis submissions. Importantly this details whether the information is mandatory and whether controlled terminology or ontological limits are applied (Figure 1). This resource is crucial for consortium members to prepare data submissions to FAANG and ensuring that they collect and submit the required metadata information.

Figure 1. The metadata ruleset guide pages of the FAANG Data Coordination Platform, including different metadata rule groups, level of requirement and controlled terminology.

## 2.2 Validation and brokered submission

https://data.faang.org/validation/samples
https://data.faang.org/validation/experiments
https://data.faang.org/validation/analyses

Metadata rulesets are only useful if submitters are enabled and required to meet them. The FAANG Data Coordination Centre has released a new validation and submission process for the submission of GENE-SWitCH metadata, incorporating improved validation error handling and an updated user interface (https://data.faang.org/validation/samples; Figure 2). This makes the submission process and provision of rich FAIR (Findable, Accessible, Interoperable and Reusable) metadata easier for the EuroFAANG projects. The service takes a completed metadata spreadsheet from a GENE-SWitCH submitter, converts it to a JSON file format, and conducts a range of detailed validation checks. The service not only provides errors for metadata inaccuracies, but also provides warnings and suggestions for improvements that the users could make to their metadata submissions such as being more specific with ontologies. Once validated, the service takes the users archive submission credentials and submits the metadata to the underlying archive on the users behalf.

Figure 2. Validation and brokered submission component of the FAANG Data Coordination Platform.

## 2.3 Protocol upload and display

https://data.faang.org/protocol/samples

The EuroFAANG DCC also released a new protocol upload interface, to make it easier for the EuroFAANG projects to upload their sampling, sequencing, and analysis protocols to FAANG, in particular for sharing with the other EuroFAANG partners (https://data.faang.org/up-load_protocol). All of the protocols uploaded to FAANG are made available through the data portal (https://data.faang.org/protocol/samples). Protocols used within the GENE-SWitCH project are attached to the datasets so that the protocol can be obtained alongside the dataset.

Figure 3. Protocol browser within the FAANG Data Coordination Platform.


### 2.4 GENE-SWitCH presentation platform

https://data.faang.org/projects/GENE-SWitCH

The GENE-SWitCH project presentation platform is a fully customisable view of the FAANG data specific to the GENE-SWitCH project. It contains a project description, logo, funding logo and automated twitter stream. Data that is known to be part of GENE-SWitCH through the setting of the secondary project tag on data submission is automatically synced to this portal view. The page has interactive tables of GENE-SWitCH publications, datasets, files, organisms, and specimens.

Figure 4. GENE-SWitCH project view of the FAANG Data Coordination Platform.

## 2.5 FAANG dataset portal pages

https://data.faang.org/dataset

The data portal pages provide interactive tables of the organisms, specimens, datasets, and files of FAANG data. Preconfigured filters are provided to identify data of interest by filtering by fields such as species, assay type, breed, and instrument. Each record has links to a full page with additional metadata detail and links to datasets in the underlying archive. The pages provide summary files of the tabulated metadata that can be downloaded. The datasets can be sorted by different metadata columns.
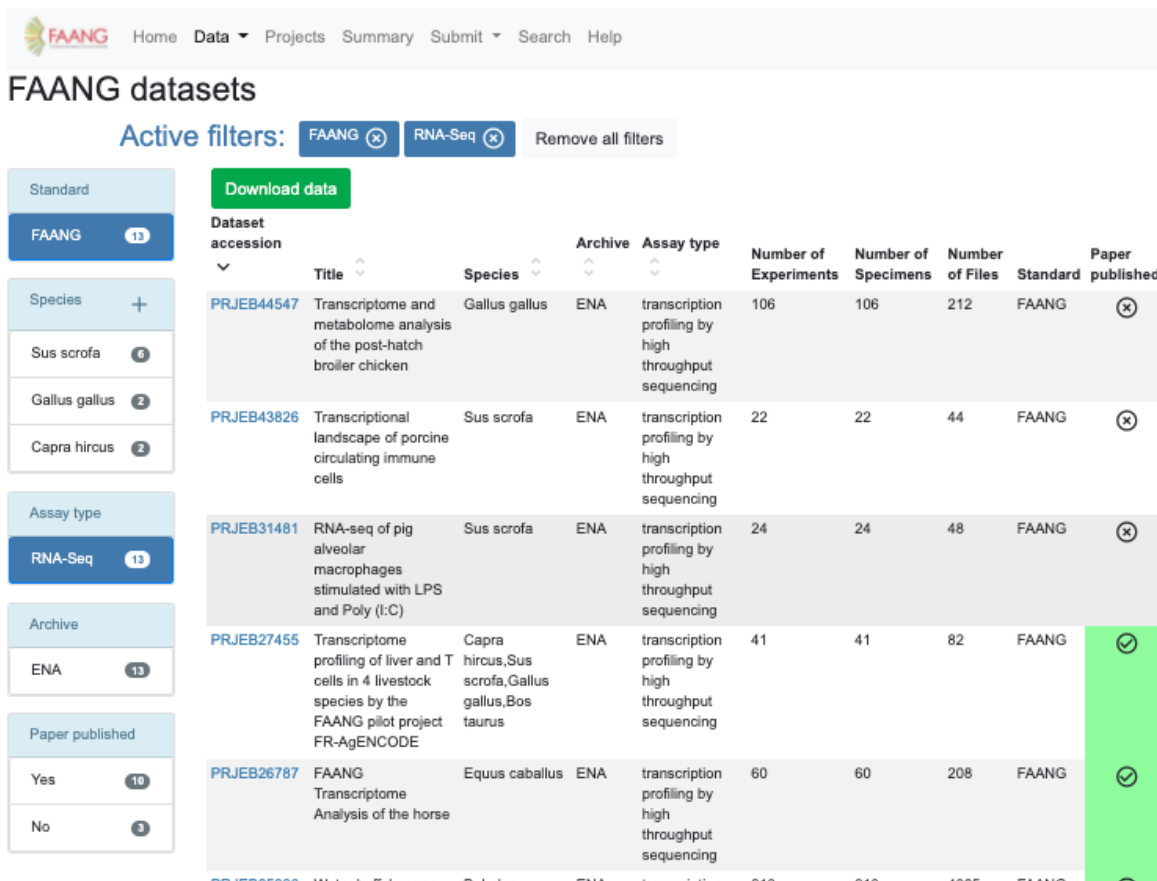
Figure 5. Data portal tables of the FAANG Data Coordination Platform.

## 2.6 Publication browser

https://data.faang.org/article

Publications on any FAANG dataset are automatically identified using the EMBL Literature services EuropePMC literature scraping technology (https://europepmc.org/). This identifies INSDC FAANG identifiers used within the main text of any publication, these publications are then imported into the FAANG data portal and linked to the dataset quoted. Links to the full text of the publication are included, and in the main data portal tables datasets that have publications associated with them are clearly labelled.

Figure 6. Publication browser of the FAANG Data Coordination Platform.

## 2.7 FAANG search interface

https://data.faang.org/search

The search interface provides a keyword search that utilises ElasticSearch tokenisation. It simultaneously searches across metadata from organisms, specimens, files and datasets. It has the option to exclude legacy data that does not met the full FAANG standards. Search results provide the number of hits, an interactive table with a limited summary view of the data and links to view the full information within the data portal pages.

Figure 7. Search interface of the FAANG Data Coordination Platform.

## 2.8 Summary statistics

https://data.faang.org/summary/organisms

The summary statistics page of the FAANG Data Coordination Platform provides overall summary statistics on FAANG data for organisms, specimens, datasets and files (Figure 8). It provides statistic plots on sex, whether data is associated with a publication, organisms, legacy status, breeds, cell type and assay type.
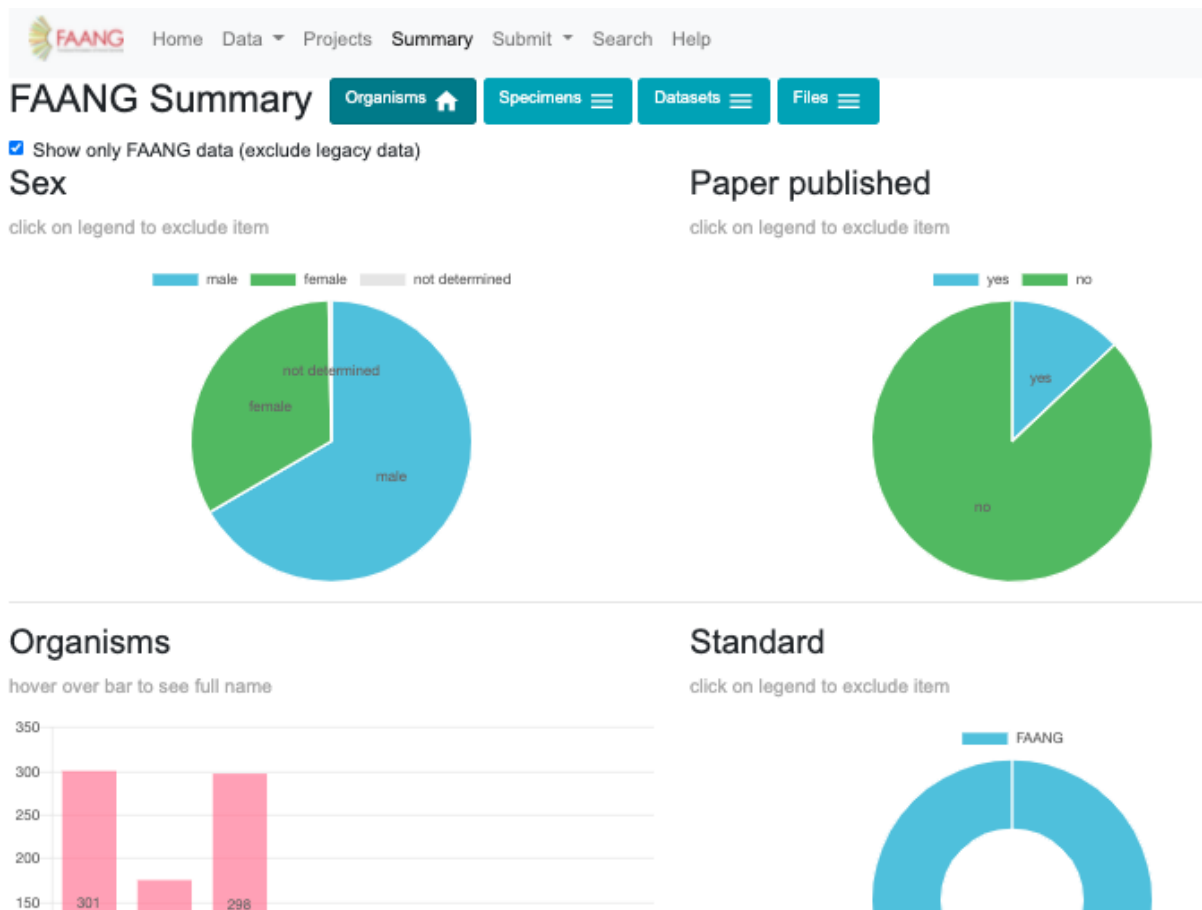
Figure 8. Summary statistics pages of the FAANG Data Coordination Platform.

## 3    Conclusion

The FAANG Data Coordination Platform has been developed and deployed to meet the data coordination, archiving and presentation requirements of the GENE-SWitCH project. The platform is in full production used with eleven datasets already coordinated, archived and on display in the data portal.

## 4    Deviations or delays

This deliverable was delayed by six months due to the secondment of EMBL DCC staff from March to September 2020 to support the emergency European COVID-19 scientific response to develop the European COVID-19 data portal (https://www.covid19dataportal.org/).

## 5    Glossary

| | |
|---|---|
| DCC | Data Coordination Centre |
| EuroFAANG | The European node of the Functional Annotation of Animal Genomes Data Coordination Centre |
| EMBL | European Molecular Biology Laboratory |
| FAANG | Functional Annotation of Animal Genomes |
| FAIR | Findable, Accessible, Interoperbale and Reusable |