# The Functional Annotation of Animal Genomes Data Portal

Current and future perspectives for data reuse

Peter Harrison

Genome Analysis Team Leader

EMBL-European Bioinformatics Institute
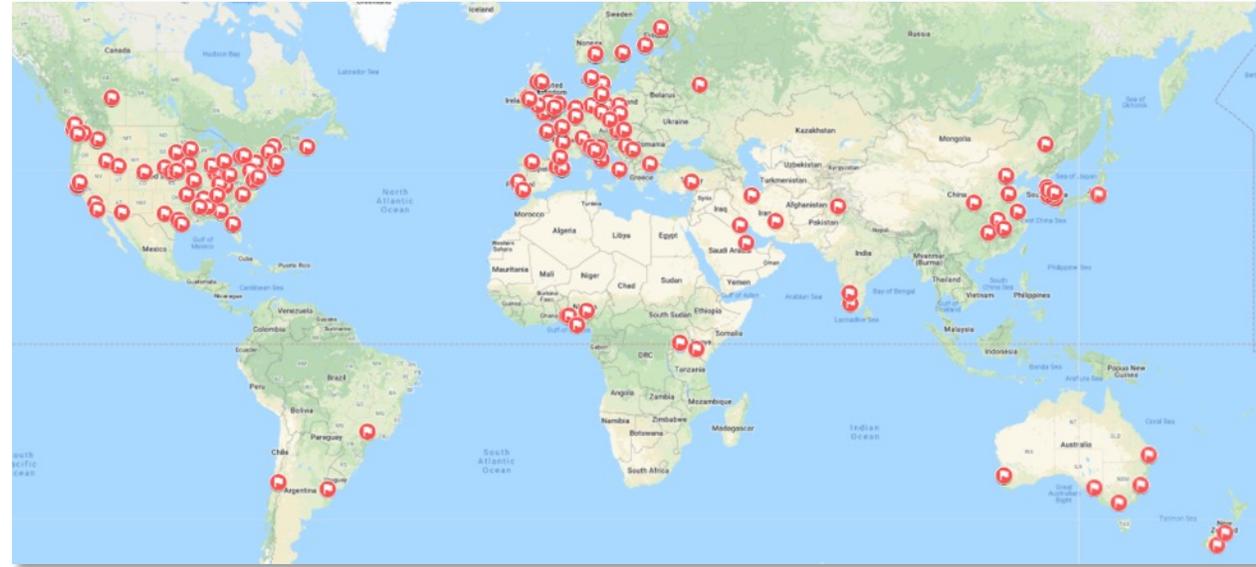
peter@ebi.ac.uk

@peterwharrison

EMBL-EBI

A coordinated international action to accelerate genome to phenome research

- Coordinated international effort to provide **high quality functional annotation** of animal genomes, with a focus on **livestock and aquaculture** communities, but extending to other animals.

- Core aims of data **openness, reusability, rich metadata and standardisation to** create a harmonised **rich genome to phenome resource**.

- EuroFAANG is a particularly coordinated effort within Europe to standardize our research processes to collectively improve animal production and welfare.

EMBL-EBI

# A global collaboration

- Steering committee, scientific advisory board, and working groups focus on range of key data reuse aspects.

- Comprised of multiple individual funded grants -> why coordination and standardisation is so important.

- FAANG projects collaborate and reuse existing genomic data to accelerate research and provide value to funders. Strong H2020 and USDA project data reuse.



**Functional Annotation of the Porcine Genome**

USDA NIFA *Grant 2018-67015-27501* [**$2,500,000**] (2018-01-01~2021-12-31)

OBJECTIVES: *(1) Maximize functional annotation in healthy adult porcine tissues relevant to phenotypes important for genetic improvement; (2) Create functional annotation of important tissues during fetal development and the association of allele-specific expression with allele-specific chromatin modification; (3) Identify the functional components of the immune system as a resource for improving resilience in pigs; (4) Integrate all public and project data to develop a higher-order regulatory understanding of the porcine genome, including a predicted chromatin state map.*

**Genome-wide annotation of cis-regulatory elements in the chicken genome**

USDA NIFA *Grant 2018-67015-27499* [**$1,000,000**] (2018-10-01~2021-09-30)

OBJECTIVES: *(1) Identify coding and long noncoding RNA transcripts in the chicken genome. (2) Identify promoter, enhancer, silencer and insulator elements in the chicken genome. (3) Characterize activity states and tissue specificity of cis-regulatory elements. (4) Map distal regulatory elements to their interacting target promoters.*

PROJECT DIRECTORS:

EMBL-EBI

# The FAANG Data Portal

- A single access point to all FAANG metadata, data and publications.

- Providing direct access to download all data from various underlying public archives.

- Automatically identifies dataset (re)use in publications, and links these publications to datasets.

- Intuitive search and filtering.

- Has sub-project pages to access just that projects data.



https://data.faang.org/          https://data.faang.org/projects

# What makes the FAANG datasets special

- Rich, consistent and validated metadata descriptions.

- Standardised set of core assays from contributing projects.

- Mandatory sampling, experiment and analysis protocols connected to each dataset and available with the datasets for download.

- Many projects using standardised analysis pipelines attached to each dataset.

- A data platform and community drive that ensures data is open and FAIR.

F indable   A ccessible   I nteroperable   R eusable

*To accelerate genome to phenome research*

Photo: CODATA

EMBL-EBI

# A full metadata solution for FAANG

- Requires **>200** different metadata questions for different studies.

- Constantly evolving with the community, recent changes include aquaculture, single cell sequencing and focus on developmental timepoints.

- Terminology controlled through standardised ontologies to make downstream search and analysis more powerful. Drives portal data filters.

| Name | Description | Type | Required? | Allow multiple? | Valid values | Valid units | Valid terms | Condition |
|------|-------------|------|-----------|-----------------|--------------|-------------|-------------|-----------|
| Organism | NCBI taxon ID of organism. | ontology id | mandatory | No | | | NCBITaxon:1 | |
| Sex | Animal sex, described using any child term of PATO_0000047. | ontology id | mandatory | No | | | PATO:0000047 | |
| birth date | Birth date, in the format YYYY-MM-DD, or YYYY-MM where only the month is known. For embryo samples record 'not applicable'. | string | recommended | No | | YYYY-MM-DD, YYYY-MM, YYYY | | |
| breed | Animal breed, described using the FAANG breed description guidelines (http://bit.ly/FAANGbreed). Should be considered mandatory for terrestiral species, for aquatic species record 'not applicable'. | ontology id | recommended | No | | | LBO:0000000 | |
| health status | Healthy animals should have the term normal, otherwise use the as many disease terms as necessary from EFO. | ontology id | recommended | Yes | | | PATO:0000461 EFO:0000408 | |
| diet | Organism diet summary, more detailed information will be recorded in the associated protocols. Particularly important for projects with controlled diet treatements. Free text field, but ensure standardisation within each study. | string | optional | No | | | | |
| birth location | Name of the birth location. | string | optional | No | | | | |

EMBL-EBI

# Validation and brokered submission

- Rich metadata rulesets are only useful if they are met by all submissions.

- All FAANG data goes through pre-submission validation, that blocks submission till compliant.

- Validation service not only highlights errors, it warns on suggested improvements such as being more specific in ontologies.

- Brokered submission to underlying archives to ensure standard presentation.

# Detailed protocols mandatory with every submission

- Enhances reproducibility, reuse and comparative study possibilities.

- Our protocol browser shows all past protocols. A useful reference of methodologies for future studies.

- Encourages standardisation across future studies.

EMBL-EBI

# FAANG data analysis

- Shared development of a complete set of open pipelines across FAANG.

- Development based on the principles of open science, open source code and reproducible workflows.

- Researchers reuse and improve a common set of pipelines.

- EuroFAANG projects are exemplifying this shared development approach, focussing around nf-core.

- Will never be one pipeline that fits all, so we capture metadata on pipeline parameters for full reproducibility.

# How FAANG promotes highly reusable standardised datasets



Automated Literature services

Rich, validated metadata

Accessible Analysis pipeline and parameters

Detailed sampling and analysis protocols

FAANG dataset

Legacy contextual datasets

Secondary analyses and annotations

EMBL-EBI

# Intuitive identification of relevant reusable data



Google style search



APIs for external access -> FAANGmine.

Ontology and

metadata filters



FAANG data drives improved Ensembl annotations

EMBL-EBI

# FAANG Ontology improvement service   **Beta**

- Frequently ontologies are not appropriate for use in agriculture and G2P, as are medical or model organism based.

- We are developing a FAANG Ontology improvement service that will community source improvements to ontologies of importance in animal agriculture.

- Allows users to list ontologies of importance, flag ontologies that need improvement or are missing.

- Users submit improvements that are forwarded to underlying ontologies for update.



https://data.faang.org/ontology

EMBL-EBI

# Embargos and third party restrictions

- Necessary in many contexts, but must recognise the dampening of accelerated research if they limit open data reuse.

- FAANG encourages prepublication data archiving under Fort Lauderdale agreements to facilitate data reuse.

- Recently updated its data sharing statement to make data reuse expectations clearer.

- Still lots needed as a community to develop data openness and clarity of labelling.

EMBL-EBI

# Potential challenge of Multiple references to data reuse

- Potential need to reanalyse past data and for researchers to agree and switch to new reference in coordination.

- Ensuring data is comparable and does not go out of date.

- Future advances of graph representations and scalability of genome browsers are needed.

# Some key Gaps for FAANG and wider communities to address

- **Standardised** ontologies across animals and crops, and **improving existing ontologies** from model organism and medical focus.

- Standardisation of **minimal metadata standards** between larger projects.

- Making data reuse conditions clearer, clearer labelling in molecular archives and data portals. **Machine readable third party constraints**.

- Need for **cloud based preconfigured analyses** to ease entry point and further standardise G2P analysis.

- Managing multiple references and graph genomics.

- These are some of the key points of focus for FAANG, AG2PI and Agbiodata.

EMBL-EBI

# Concluding Remarks

- FAANG promotes reuse through datasets having rich validated metadata, mandatory protocols, prepublication data sharing, standardised assays, standardised and documented analyses and an intuitive data portal.

- FAANG has produced functional maps of key animal species and continues to broaden its scope (FAANG to fork), whilst maintaining its core values of reusability and openness.

- Keen to address the data torrent grand challenges, in coordination with wider communities, to further promote data reuse, openness and standardisation.
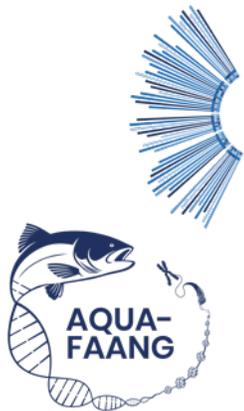


Photo: Peter Harrison

Google 'From FAANG to fork'

EMBL-EBI

# Acknowledgments

peter@ebi.ac.uk

@peterwharrison

**EMBL-EBI**

- Alexey Sokolov (FAANG DCC Project Leader)
- Akshatha Nayak (FAANG DCC Bioinformatician)
- Koosum Roochun (FAANG DCC Bioinformatician)
- Raheela Aslam (DCC Web developer)

EuroFAANG

Horizon 2020
European Union Funding
for Research & Innovation

AQUA-FAANG

BovReg
Understanding cattle genomes

GENE-SWitCH

- H2020 EuroFAANG collaborators in particular form:
  - AQUA-FAANG
  - BovReg
  - GENE-SWitCH
- Members of FAANG steering committee and working groups
- Fantastic FAANG community for their continued data submission, feedback and support

e!Ensembl

wellcome

https://data.faang.org/

EMBL-EBI