

Metadata harmonization at scale:

FAIRification of Genomic Annotations WG

Sveinung Gundersen ^{1*}, Anna Bernasconi ^{2*}, Nathan Sheffield ^{3*}, Adam Wright ^{4*}

1ELIXIR Norway / University of Oslo (UiO), Norway

2Politecnico di Milano, Italy

3University of Virginia, Virginia, USA

4Ontario Institute for Cancer Research, Canada

WG Co-chairs

Experiment Matrix

Showing 16574 results

Assay

Assay type

Assay title

Target category

TF ChIP-seq	443	2182	933	415	399	225	305	385	275	121	2	67	1
Histone ChIP-seq	59	188	56	120	2	15	5	54	65	3			
DNase-seq	7	83	20	15	16	7	25	40	8			1	
total RNA-seq	8	39	32	23	11	8	23	59	15	5	2	2	
Mint ChIP-seq	42	91	15	3	20	8	7				1	1	
polyA plus RNA-seq													
ATAC-seq													
microRNA-seq													
scRNA-seq													
snATAC-seq													

Immense resources have gone into the generation of data that relates genomic positions to functional and structural aspects, in large consortia such as ENCODE, ICGC, FANTOM and FAANG, as well as in smaller research projects. This is complemented by a current push to sequence millions of species in biodiversity projects – all of which needs characterisation of genes and other genomic features.

2003 → 2021

2008 → 2018

2011 → 2017

2005 → 2014

2014 → 2020

2015 →

2016 →

2018 →

2019 →

Genome browser meme by Jedidiah Carlson

Functional genomics data are typically provided as tracks for visualisation in genome browsers, or made downloadable for non-visual analysis. The use of reference genomes as one-dimensional coordinate systems constitutes a powerful unified model for data analysis. Interestingly, the condensed nature of coordinate-based data is ideally suited for data-driven approaches!

Pangenomes challenge the concept of one-dimensional reference genomes. While a consensus has not yet formed on how to annotate pangenomes, the value of reusing the abundance of existing functional genomics data in this new paradigm is clearly evident. Moreover, data-driven discovery by applying AI to functional annotations still holds considerable untapped potential. To these ends, FAIR (Findable, Accessible, Interoperable, Reusable) metadata remains key.

Flatten hierarchical JSON

Normalise relational tables

Map fields across rel. tables

Data cleanup

Map/convert ontology terms

Build hierarchical JSON

Validate w/JSON Schema

Harmonised JSON documents

Harmonisation and FAIRification of metadata from different sources is a time-consuming and error-prone process, often implemented in *ad hoc* scripts that are difficult to maintain and scale. While each source exhibits particularities in schema, API and content, many operations are common. It should therefore be possible to streamline the process of developing maintainable metadata transformation pipelines.

Functional genomics data are often provided through dedicated data portals, such as the ENCODE data portal (left). However, the metadata are provided according to distinct models and interfaces. Additionally, data from smaller research projects are often available only with limited metadata.

Planned FAIRification infrastructure & WG deliverables

FGA-WG Guidelines for enabling scalable and maintainable metadata transformation pipelines

FGA-WG Strategy for providing persistent identifiers (PIDs) Possibly: content-derived

FGA-WG Strategy for persistent and public deposition of harmonised metadata

FGA-WG Minimal Metadata Schema

FGA-WG Registry

Public repositories

Academic journals

Research data reuse loop

Research tool integration + Direct access of API

FGA-WG Uniform Search API

Third-party services harvest and share public metadata

QR code

https://tinyurl.com/rda-genomics

Join our working group in RDA:

FAIRification of Genomic Annotations WG (FGA-WG)

We aim to produce a set of recommendations that together will define a community-oriented infrastructure to make it easier to discover and reuse genomic annotations in a range of contexts.

WG meetings currently every 1st and 3rd Tuesday per month, at 2pm UTC.

New Global Support Services for RDA Working Groups

Relevant roles:

– data producer

– tool developer

– biocurator

– domain expert

– analytical end user

We do not want to reinvent the wheel.

If you have knowledge to share, please come see us!