# FAIR Genomes F2F @ Utrecht Feb. 6th WP2 & WP3 (Gurnoor, Joeri, et al.)

# WP2/WP3 have been busy

→ FAIR genomes WP2/WP3 zoom video meetings:
  ◆ 2 oct 2019: "Personal and Clinical"
  ◆ 8 nov 2019: "Materials"
  ◆ 22 jan 2020: "Technical"

→ V1.1 of meta model, mostly ontologized

→ 2 posters @ HealthRI 2020

→ A few open issues to be discussed

→ GitHub page updated

FAIR Genomes working document:

https://docs.google.com/spreadsheets/d/1rnLsmE62t15jCwJfx4mCL5USYSeiXNctCA0XPcgprds/edit#gid=914344861

→ **Metadata_version 1.1**

| | |
|---|---|
| Description | Description or example of this attribute |
| Compulsory / Optional | "Compulsory" or "Optional" ... |
| Ontology Term (Meta-data) | Link to ontology and term |
| Ontology Term (Data) | Value type, could be ontology (sub)term or data type like INTEGER, DATE, STRING etc |
| Defined in other projects | Is this attribute present in other projects and if so, how is it represented? |

**Slide format**

# Personal information

| Personal information | | |
|---|---|---|
| Biological sex | optional | biological sex is the quality of a biological organism based on reproductive function or organs |
| Country of residence | Optional | country of residence (environment and epidemology) |
| Ethnicity | Optional | clarify: relates purely to genetic/ethnic |
| Country of Birth | optional | country where patient is born according to official birth certificate |
| Year of birth (if allowed) | optional | year when patient is born according to official birth certificate |
| Patient Status | | [Alive, Dead, Lost in follow-up, Opted-out] |
| Age at death | optional | official age of patient when death occurred |
| Personal ID | compulsory | Anonymous / Non identifiable, |
| Inclusion criterion | optional | an inclusion criterion is a eligibility criterion which defines and states a condition which, if met, makes an entity suitable for a given task or participation in a given experiment. can also have 'age of inclusion' |
| Primary affiliated institute | | |
| Data available in other institutes? | | [clarify why this is needed] |

*Feedback @ HealthRI 2020*

- Biological sex types: Currently we just have male, female and null flavor. There should be more like, "raised as male/female", undetermined, unknown. Not sure if null-flavor can cover both undetermined, unknown
- Biological sex could be renamed as Gender
- Patient can be renamed as Individual
- Ethnicity can be renamed as Population

*No issues (?)*

# Clinical information

| Clinical information | | |
|---|---|---|
| Phenotypic terms | compulsory | phenotypic terms best describing the patient symptoms according to a licenced clinician [we no |
| Unobserved phenotypes | optional | phenotypic terms that were NOT observed |
| Type of phenotypic data | compulsory | e.g. pictures, terms: free text, CSV, HPO terms, CPMS terms , none??? |
| Clinical diagnosis | optional | patient disorder or disorder spectrum as established by licenced clinician |
| Genetic diagnosis (if part of a disease cohort) | optional | genetic diagnosis comprises of the causal variant, HGVS nomenclature, OMIM or Gene symbol |
| Age at diagnosis | optional | patient age when diagnosis was officially established by a licenced clinician |
| Age at last screening (if part of a cohort) | optional | patient age at which a particular screening relevant to cohort selection took place |
| Medication information | optional | any medication taken by the patient, e.g. steroids |
| Dosage | optional | dosage of each medication |
| Family members affected | compulsory | e.g. pedigree. Other affected relatives, or unaffected, but also remarks about consanguinity or RoH assays/ Or: is family information present yes or no and where is it? |
| Family members sequenced | optional | Whether or not and if so, which family members have also been sequencing, helps to classify |
| Procedural history | optional | e.g. liver tumor removed |
| Age of onset | optional | registered and/or self reported patient age at which symptoms for the disorder started to manif |
| First contact with specialised centre [perhaps delete ?? needed?] | optional | Date of first contact with specialised centre with relation to current diagnostic process |

*Issues:*
- *Family information*
- *First contact*

# Material information

| Material information | | | |
|---|---|---|---|
| Sampling TimeStamp | compulsory | Know when the sample was taken;<br>Know when NGS was ordered (in at least cancer it will be relevant to have time/datestamps with respect to disease/treatment/folow-up) to determine if another new NGS is needed for exa | UMLS [C0870078] Sampling |
| Registration TimesStamp | compulsory | the date when sample was entered into the system | |
| Sampling collection protocol | optional | | http://www.ebi.ac.uk/efo/EFO_0005518 |
| Deviations from Sample protocol | optional | any deviations from following the above protocol | NCIT:C25713 |
| Reasons for protocol deviation | optional | why were there deviations from the protocol | |
| Material type | compulsory | Type of material collected, e.g. blood, muscle, bone etc | [C2986062] Material Identifier Type Code |
| Anatomical source | optional | Anatomical source from which this material was derived | UBERON:0001062 |
| Storage conditions | optional | Storage conditions under which this material was kept, but also parafin fixed, fresh frozen, hep | UMLS:C3272596 Storage Condition |
| Expiration date | optional | when is this material allowed to be trown away? or must be thrown away | LOINC:LP173684 |
| Estimated percentage of tumor cells | optional | tumor cell to total cell ratio measurement obtained from this material | UMLS:C4288090, Tumor Cell to Total Cell Ratio Measurement / |
| Amount of input material used | optional | Gives background on how much information can be expected to be extracted from the source | ?? |
| Location of sample (Physical location) | optional | e.g. UMCG department of genetics (OR pURL, PID). do not put your local freezer shelf / box lo | DUO:GAZ_00000448 |
| "is deritative or not"? | optional | sample derived from another sample or not? for example "blood" taken from "tissue", or "aliquot" taken | NCIT:C28355 |

*Issues:*
- ***Sample location*** *vs* ***data location*** *in technical*

| Technical Information | | | |
|---|---|---|---|
| Sequencing date | Yes (compulsory) | Date when NGS was performed | GENEPIO:0000069 |
| Sample prep kit | Yes (optional) | e.g. Agilent QXT or Agilent XT | GENEPIO_0000081 |
| Sequencing platform | Yes (compulsory) | e.g. Illumina NextSeq500, Nanopore Gridlon, PacBio, Sanger, IonTorrent | GENEPIO_0000071 |
| Sequencing data type | Yes (compulsory) | e.g. Whole exome sequencing (WES), whole genome sequencing (WGS) | NCIT:C18881 |
| PCR-free yes/no | Yes (optional) | e.g. WGS may be done PCR-free | NCIT:C17003 |
| Sequencing average read depth | Yes (compulsory) | e.g. 100x, 30x, 42x | NCIT:C155320 |
| Enrichment panel used | Yes (compulsory) | e.g. Agilent SureSelect v7, custom hotspot panel, none | NCIT:C154307 |
| UMIs present yes/no | Yes (optional) | Does the sequencing technique use Unique Molecular Identifiers? | EFO:0010199, UMI barcode |
| Read length | Yes (compulsory) | e.g. PE 150 bp | NCIT:C153362 |
| Insert size | Yes (optional) | e.g. 350 bp, 200 bp | NA |
| Location of data | Yes (compulsory) | e.g. UMCG department of genetics | DUO:GAZ_00000448 |
| Type of data stored | Yes (compulsory) | e.g. FASTQ, BAM, CRAM, VCF (v 4.0, 4,1, 4.2, 4.3) | EDAM:FORMAT_1915 |
| Algorithms used | Yes (optional) | e.g. BWA MEM, GATK Haplotype caller: link to protocol | NCIT:C16275 |
| Bioinformatic protocols used | Yes (optional) | The bioinformatic protocol, workflow or SOP that was followed | EDAM:DATA_2531 |
| Special parameters used | Yes (optional) | e.g. alternative gap lenght (BWA MEM -w ..) | NCIT:C44175 |
| Follows international WGS guidelines | Yes (optional) | Does the DNA sequencing analysis follow existing international guidelines, and if so, which one? | NCIT:C17564 |

*Issues:*

- *Insert size ontology /* Insert size should be **Observed insert size**
- ***Sample location** vs **data location** in technical*

# For **any empty** value, please elaborate

https://www.hl7.org/fhir/v3/NullFlavor/cs.html

## *some examples:*

| NAV | temporarily unavailable | Information is not available at this time but it is expected that it will be available later. |
|---|---|---|
| ASKU | asked but unknown | Information was sought but not found (e.g., patient was asked but didn't know) |
| NASK | not asked | This information has not been sought (e.g., patient was not asked) |
| NA | not applicable | Known to have no proper value (e.g., last menstrual period for a male). |
| MSK | masked | There is information on this item available but it has not been provided by the sender due to security, privacy or other reasons. There may be an alternate mechanism for gaining access to this information. |
| TRC | trace | The content is greater than zero, but too small to be quantified. |

# In progress: connecting to...

- → 1+ MG
- → SolveRD
- → EJP-RD
- → X-omics
- → GA4GH
- → CINECA
- → Phenopackets
- → Biosamples

- → RIVM WG?
- → Illumina WG?
- → BioSchemas (!)
- → WP5: 'FAIR variants' ?

*see*

https://raw.githubusercontent.com/LUMC-BioSemantics/ERN-common-data-elements/master/images/complete_data_model_v2.0.png

https://github.com/Xomics/GenomicsMeta-data

https://phenopackets-schema.readthedocs.io/en/latest/building-blocks.html

https://submission.ebi.ac.uk/api/docs/guide_getting_started.html

*How to publish FAIR genomes guidelines for standardizing meta-data?*

➔ Meta-data release
  ◆ Github [https://github.com/fairgenomes/information] and Github wiki
➔ Data-model
  ◆ Reviewed by stakeholders
    ● (clinicians, diagnostics, pathologist)
  ◆ Technical review
  ◆ Community review

# Future work

→ Discuss further in FAIR genomes & with other experts
  ◆ Finalize attributes, descriptions, values, link to other projects / standards
  ◆ Integrate with WP4

→ Demonstrator projects
  ◆ Try out the FAIRgenomes schema on real genomics datasets

→ From meta-data to a 'real' data model ?
  ◆ Semantic data model (RDF or JSON-LD)
  ◆ Rules based validator (RDF Shacl or JSON-LD markups)

→ Apply to daily practice, use in care & research
  ◆ (VKGL, VKGN, KMBP……)

→ Publish to share our insights & FAIRify the rest of the world

# Thanks