October 2021

# IHCC Cohort Atlas Workshop

Thomas Keane, EMBL-EBI

# Agenda for today

| Time (BST) | Subject | Presenter |
|---|---|---|
| 16:00 - 16:10 | IHCC project and data team overview (10') | Thomas Keane (EMBL-EBI) |
| 16:10 - 16:20 | Atlas overview (10') | Brandon Chan (OICR) |
| 16:20 - 16:35 | Atlas Data harmonization (15') | Melanie Courtot (EMBL-EBI) |
| 16:50 - 18:00 | Live demo (including processing steps) | Carles Garcia, Isuru Liyanage (EMBL-EBI) |

# Data and Infrastructure Team

*"Deliver <u>interoperable cross cohort infrastructure</u> to enable population scale biomolecular <u>data to be accessible across international borders</u> accelerating research and improving the health of individuals resident across continents."*

# Challenges

Challenge 1: Federated cohort data discovery

Challenge 2: Harmonized cohort level metadata

Challenge 3: Cohort access and authorization

Challenge 4: Federated analysis interoperability for research

# Challenges

Challenge 1: Federated cohort data discovery

Challenge 2: Harmonized cohort level metadata

Challenge 3: Cohort access and authorization

Challenge 4: Federated analysis interoperability for research

# Global Standards



Cohort interoperability standards are emerging

Global Alliance for Genomics and Health (GA4GH)
- 8 workstreams (e.g. Discovery, Data Use, Clinical and phenotypes, Ethics, Security)
- Clinical metadata standards (e.g. HL7/FIHR, OMOP etc)

IHCC Data team and GA4GH
- Foundation for all IHCC products to interop with other cohorts
- e.g. IHCC cohort atlas cross queries with other aggregate resources

# Global Standards

Cohort interoperability standards are emerging

Global Alliance for Genomics and Health (GA4GH)
- 8 workstreams (e.g. Discovery, Data Use, Clinical and phenotypes, Ethics, Security)
- Clinical metadata standards (e.g. HL7/FIHR, OMOP etc)

IHCC Data team and GA4GH
- Foundation for all IHCC products to interop with other cohorts
- e.g. IHCC cohort atlas cross queries with other aggregate resources



Global Alliance
for Genomics & Health
Collaborate. Innovate. Accelerate.

Data Use Ontology

Phenopackets
Open and Computable Bioinformation

Beacon Network

Tool Registry Service (TRS)

Tool

Workflow

Workflow Execution Service (WES)

Workflow Execution Engine

# Progress to date

IHCC Cohort Atlas
- Enables researchers to search across IHCC cohort metadata
- 13 cohorts, >100 variables harmonised

Davos Alzheimer's Collaborative
- Example of disease specific expansion of the atlas

Future vision for atlas
- Cornerstone of future IHCC data science platform

**International HundredK+ Cohorts Consortium** IHC

Cohort Name

- Africa Health Research Institute (AHRI) Population Cohort — 1
- Genomics England / 100,000 Genomes Project — 1
- Golestan Cohort Study — 1
- 2 More

Countries

- South Africa — 2
- England — 1
- Iran — 1
- 1 More

Biospecimens
Any | Yes 5 | No 0

Environmental Data
Any | Yes 2 | No 3

Genomic Data
Any | Yes 3 | No 2

← Use the filter panel on the left to customize your cohort search.

Cohorts by Country

Biosample Types
Saliva
Blood
biosample type
Urine

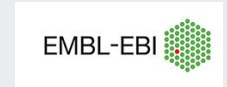Africa He
Golestan
Korean G
Other Co

Showing 1 - 5 of 5 cohorts

| Cohort Name | Countries | Current Enrollment | Genomic Data | Environmental Data | Biospecimen Data | Clinical Data | Data Sharing Potential | PI L |
|---|---|---|---|---|---|---|---|---|
| Africa Health Res... | South Africa | 130000 | ✓ | ✗ | ✓ | ✓ | ✓ | Will |
| Genomics Engla... | England | 100000 | ✓ | ✗ | ✓ | ✓ | ✓ | Mar |
| Golestan Cohort ... | Iran | 50000 | | | | | | Reza Chr Pao Pau Fari Aras |
| | | | ✗ | ✓ | ✓ | ✓ | ✓ | |
| Korean Genome ... | South Korea | 235000 | ✓ | ✓ | ✓ | ✓ | ✓ | Hyu |
| SAPRIN (South Af... | South Africa | 350000 | ✗ | ✗ | ✓ | ✓ | ✓ | Kob |

# 5 Year Implementation



1. IHCC Atlas Content
- Expansion of cohort content
- Expansion of the harmonised metadata model
- Standardisation of core cohort descriptors
- Data model alignment with other cross cohort resources
- New modes of discovery, e.g. data use

# 5 Year Implementation



2. USING THE ATLAS: DATA DISCOVERY

I want to find imaging data I can use in my cancer research

Cohort Atlas

1. BUILDING THE ATLAS: COHORT INTEGRATION

IHCC Registry

DICTIONARY HARMONISATION

Metadata model
GECKO

MAPPING PIPELINE

Cohort metadata

SEMANTIC ALIGNEMENT

Data access request

GA4GH Passport    elixir    Researcher ID

DUO:0000007 disease specific research
MONDO:0004992 cancer

Data Use Ontology

3. FROM THE ATLAS TO DATA ANALYSIS

Data Access Committees

✓ COHORT DATA 1
✓ COHORT DATA 2
✓ COHORT DATA 3
✓ COHORT DATA 4

Trusted research environments

Research and clinical applications

**Disease detection** (the identification of disease), **prediction** (the prediction of risk of disease or therapeutic outcome), **prognostication** (the prediction of oncological outcome), and **response assessment** (the evaluation of change with therapy).

2. Using the Atlas
- Complex interactive queries
- Integration with discovery networks, e.g. Beacon network
- FAIR cohort data access
- Integration of cohort data access processes, e.g. GA4GH researcher Passports

# 5 Year Implementation



3. Federated cohort analysis
- Federated cloud analysis model for cohorts, e.g. AoU, UKB, AD workbench
- Discover, access, analysis
- Atlas integration with cohort platforms
- Exemplar for deployment of GA4GH standards

# Agenda for today

| Time (BST) | Subject | Presenter |
|---|---|---|
| 16:00 - 16:10 | IHCC project and data team overview (10') | Thomas Keane (EMBL-EBI) |
| 16:10 - 16:20 | Atlas overview (10') | Brandon Chan (OICR) |
| 16:20 - 16:35 | Atlas Data harmonization (15') | Melanie Courtot (EMBL-EBI) |
| 16:50 - 18:00 | Live demo (including processing steps) | Carles Garcia, Isuru Liyanage (EMBL-EBI) |

Oct 27, 2021

# IHCC Cohort Atlas Overview

Melanie Courtot / Brandon Chan

# IHCC Cohort Atlas Overview

- The IHCC Cohort Atlas fits directly into the overall IHCC vision and goal to bring large cohorts together to encourage data sharing, improve efficiencies and maximize benefits in addressing scientific questions no party can answer alone.

- The Atlas Browser provides a simple-to-use, easy-to-access UI for researchers to quickly search for a cohort of interest.

- Available cohorts and their metadata are first consolidated and harmonized before being published to the Browser.

- Users can then further explore those cohorts of their own accord.

- **Acts as a gateway or portal for users to access their cohorts of interest and pursue further detailed, specific exploration on each cohort.**

# IHCC Cohort Atlas Browser

https://ihccglobal.org/

# Search & Filter Controls

# Query Bar



Query bar shows the actual full query that you have constructed

# Query Results

# Cohort Summary Visualizations

# Link Out to Cohorts for Further Exploration

https://www.overture.bio/

# Data harmonisation

Mélanie Courtot
EMBL-EBI

# Building a common framework

1. **Data models** to represent both access conditions and cohort data
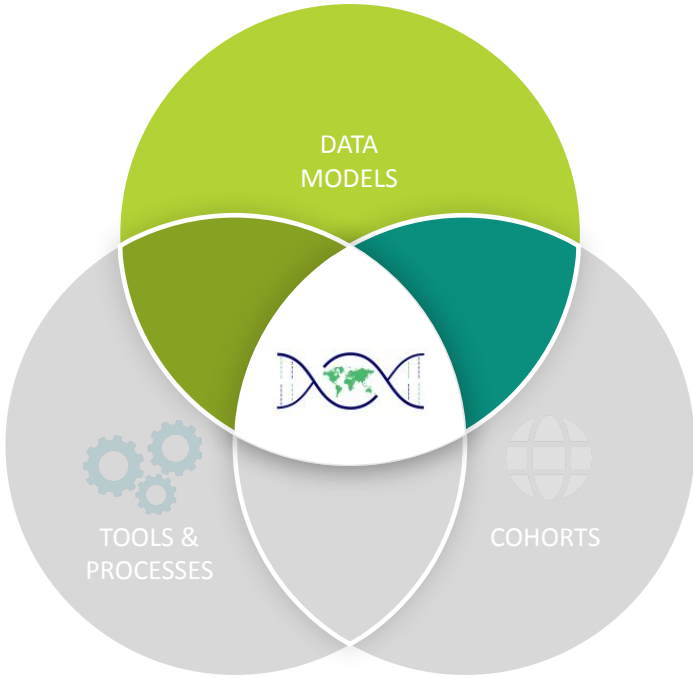2. **Tools** and processes for implementations
3. Deployment over **clinical cohorts**

# Genomics Cohort Knowledge Ontology (GECKO)

- **Commonly used attributes to describe cohort metadata**
- **"Medication", "sample type", "genomics datatypes"...**

https://www.ebi.ac.uk/ols/ontologies/gecko
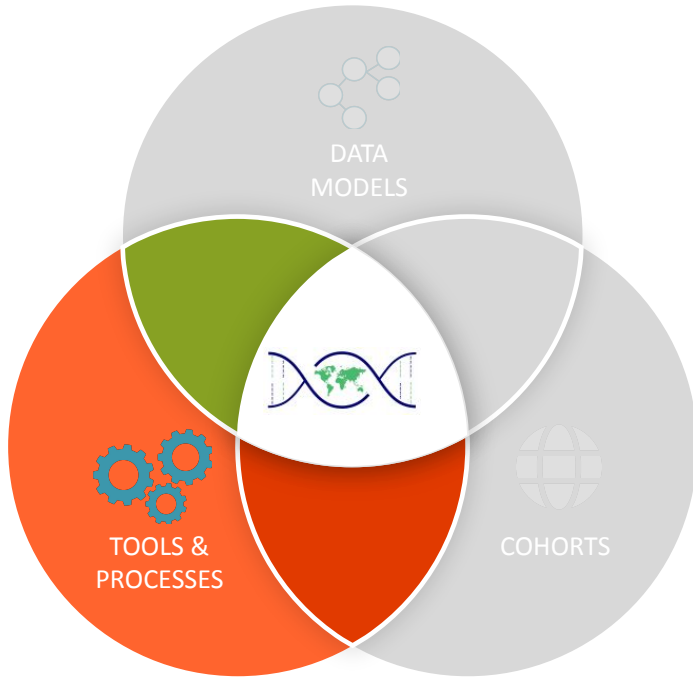https://github.com/IHCC-cohorts/GECKO

Fiona
Brinkman

Registry and mapping

# IHCC cohort registry

- **Human readability of cohort dictionaries**
- **Version and change detection for update**
- **Built on the EMBL-EBI Ontology Lookup Service platform**



https://registry.ihccglobal.app/index

# IHCC cohort mappings



- **Stores mapping between GECKO and cohort terms**
- **Accessible through APIs**
- **Parameter to bridge between mappings If A ⇔ B and B ⇔ C then can infer A⇔ C**

https://mapping.ihccglobal.app

# Automated mapping pipeline for cohort owners



IHCC cohort registry

Applying these techniques to clinical cohorts...

# Next steps

**Pipeline can be reused**



https://www.davosalzheimerscollaborative.org
Video demo: https://vimeo.com/505253841

# Next steps

**Pipeline can be reused**

**Model can be extended**



https://www.davosalzheimerscollaborative.org
Video demo: https://vimeo.com/505253841

Morris
Swertz

https://directory.bbmri-eric.eu/

# IHCC cohort atlas

Cohort presentation and display

Reference to external cohort sites

Intuitive filtering by cohort metadata & data dictionary attributes

https://ihccglobal.org



Christina Yung

Philip Awadalla

https://ihcc-cohorts.github.io/DataDictionaryMapping.html

| | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| 1 | **IHCC Data Harmonization** | | | | | | |
| 2 | Please fill the Metadata and Terminology sheets to the best of your ability. | | | | | | |
| 3 | Email to | ihcc-browser@googlegroups.com | | | | | |
| 4 | | | | | | | |
| 5 | **Metadata** | | | | | | |
| 6 | Cohort ID | Required | The short identifier for the cohort. Must be one word, uppercase, and unique among the IHCC cohorts | | | | |
| 7 | Cohort Name | Required | The title of the cohort. | | | | |
| 8 | Description | Required | The description of the cohort, in one or two sentences. | | | | |
| 9 | License | Required | The link to the license of the cohort data dictionary. | | | | |
| 10 | Rights | Optional | An optional text description of the usage rights for the cohort's data dictionary. | | | | |
| 11 | Website | Optional | A link to the cohort website. | | | | |
| 12 | PI/Lead | Required | The cohort PI/Lead (comma-separated for more than one lead). | | | | |
| 13 | Countries | Required | Country(ies) in which the cohort participants are located (comma-separated). | | | | |
| 14 | Data Sharing | Required | TRUE/FALSE value for if the cohort has IRB approval for data sharing. | | | | |
| 15 | Enrollment Data | Required* | Enrollment years and numbers for the cohort. Enrollment start year and total enrollment values are required, other values are optional. | | | | |
| 16 | Available Datatypes | Required | TRUE/FALSE values for the types of data collected by the cohort. | | | | |
| 17 | | | | | | | |
| 18 | **Terminology** | | One row for each data dictionary term. | | | | |
| 19 | Label | Required | The label of your term. | | | | |
| 20 | Parent Term | Optional | The parent term or category: this is optional, but if filled, the value must be the label of another term in this sheet. | | | | |
| 21 | Definition | Optional | The definition of your term. | | | | |
| 22 | Internal ID | Optional | If your term has an internal id, like a database identifier, you can add it here. | | | | |
| 23 | Comment | Optional | Other comments for the IHCC data harmonization team. | | | | |
| 24 | | | | | | | |
| 25 | | | | | | | |

◄ ► | **Instructions** | Metadata | Terminology | + |

| | A | B | C |
|---|---|---|---|
| 1 | **Identifiers** | | |
| 2 | **Cohort ID** | | |
| 3 | **Cohort Name** | | |
| 4 | | | |
| 5 | **Metadata** | | |
| 6 | **Description** | | |
| 7 | **License** | | |
| 8 | **Rights** | | |
| 9 | **Website** | | |
| 10 | **PI/Lead** | | |
| 11 | **Countries** | | |
| 12 | **Data Sharing** | | |
| 13 | | | |
| 14 | **Enrollment** | | |
| 15 | **Enrollment Start Year** | | |
| 16 | **Enrollment End Year** | | |
| 17 | **Total Enrollment** | | |
| 18 | **Target Enrollment** | | |
| 19 | | | |
| 20 | **Available Datatypes** | Please select TRUE or FALSE | |
| 21 | **Biospecimens** | | |
| 22 | **Environmental Data** | | |
| 23 | **Genomic Data** | | |
| 24 | **Phenotypic/Clinical Data** | | |
| 25 | | | |
| 26 | | | |
| 27 | | | |
| 28 | | | |
| 29 | | | |

Instructions | Metadata | Terminology | +

|   | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | **Label** | **Parent Term** | **Definition** | **Internal ID** | **Comment** |
| 2 |   |   |   |   |   |
| 3 |   |   |   |   |   |
| 4 |   |   |   |   |   |
| 5 |   |   |   |   |   |
| 6 |   |   |   |   |   |
| 7 |   |   |   |   |   |
| 8 |   |   |   |   |   |
| 9 |   |   |   |   |   |
| 10 |   |   |   |   |   |
| 11 |   |   |   |   |   |
| 12 |   |   |   |   |   |
| 13 |   |   |   |   |   |
| 14 |   |   |   |   |   |
| 15 |   |   |   |   |   |
| 16 |   |   |   |   |   |
| 17 |   |   |   |   |   |
| 18 |   |   |   |   |   |
| 19 |   |   |   |   |   |
| 20 |   |   |   |   |   |
| 21 |   |   |   |   |   |
| 22 |   |   |   |   |   |
| 23 |   |   |   |   |   |
| 24 |   |   |   |   |   |
| 25 |   |   |   |   |   |
| 26 |   |   |   |   |   |

Instructions | Metadata | Terminology | +

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | Label | Parent Term | Definition | Internal ID | Comment |
| 2 | Nationality | | The country of origin of the individual (the country of which the person is/was a citizen by birth) | | |
| 3 | Surname | | Family or last name | | |
| 4 | First Name | | The first given name of the individual | | |
| 5 | Sex | | Male (1), female(2) or unknown(0) | | |
| 6 | Birth Date | | Date of Birth | | |
| 7 | Employment Type | | Employed, Unemployed, Part time, Retired | | |
| 8 | Birth Weight | | Birth weight in grams | | |
| 9 | Blood cell count | | Complete blood cell count | | |
| 10 | Glucose | | Fasting glucose level | | |
| 11 | | | | | |
| 12 | | | | | |
| 13 | | | | | |
| 14 | | | | | |
| 15 | | | | | |
| 16 | | | | | |
| 17 | | | | | |
| 18 | | | | | |
| 19 | | | | | |
| 20 | | | | | |
| 21 | | | | | |
| 22 | | | | | |
| 23 | | | | | |

Instructions | Metadata | Terminology | +

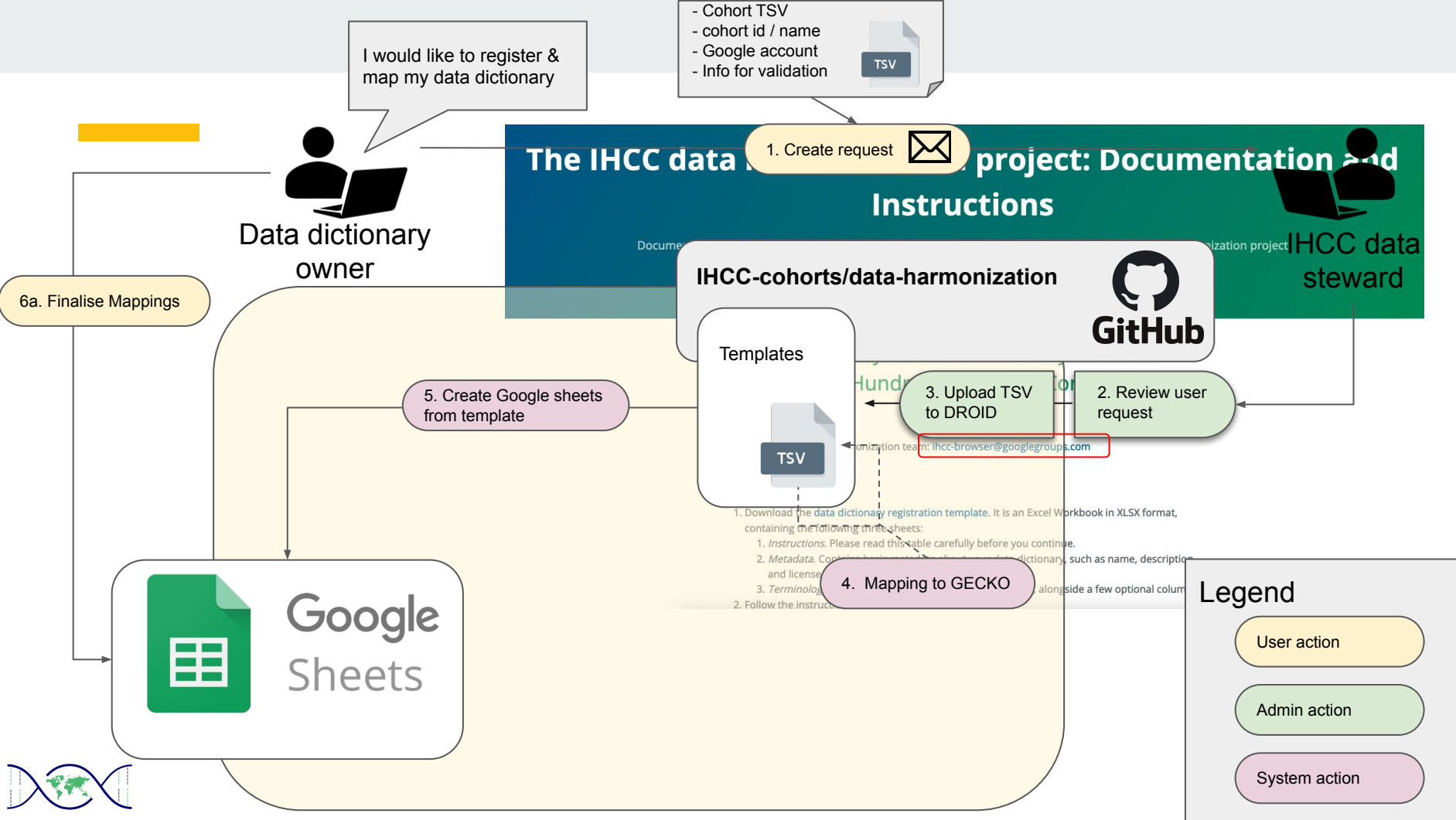| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | Label | Parent Term | Definition | Internal ID | Comment |
| 2 | Nationality | | The country of origin of the individual (the country of which the person is/was a citizen by birth) | | |
| 3 | Surname | | Family or last name | | |
| 4 | First Name | | The first given name of the individual | | |
| 5 | Sex | | Male (1), female(2) or unknown(0) | | |
| 6 | Birth Date | | Date of Birth | | |
| 7 | Employment Type | | Employed, Unemployed, Part time, Retired | | |
| 8 | Birth Weight | | Birth weight in grams | | |
| 9 | Blood cell count | | Complete blood cell count | | |
| 10 | Glucose | | Fasting glucose level | | |
| 11 | | | | | |
| 12 | | | | | |
| 13 | | | | | |
| 14 | | | | | |
| 15 | | | | | |
| 16 | | | | | |
| 17 | | | | | |
| 18 | | | | | |
| 19 | | | | | |
| 20 | | | | | |
| 21 | | | | | |
| 22 | | | | | |
| 23 | | | | | |

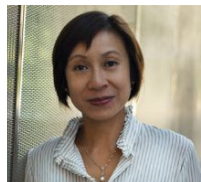Instructions | Metadata | Terminology | +

# **Acknowledgements**



**Thomas Keane**

**Philip Awadalla**

Christina Yung

Rosi Bajari

Giselle Kerry

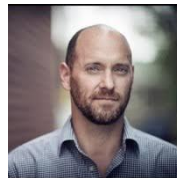Melanie Courtot

Eric Plummer

Minh Ha

Brandon Chan

Carles Garcia

Isuru Liyanage

Dan Brake

Chris Lunt

**KNOCEAN**

James Overton

Rebecca Jackson

Nicolas Matentzoglu

elixir

CINECA

OICR
Ontario Institute
for Cancer Research

EMBL-EBI

NIH National Institutes of Health
*Turning Discovery Into Health*

Global Genomic
Medicine Collaborative

This project has received funding from the European Union's Horizon 2020
Research and Innovation Programme under grant agreement No. 825775.