

Implementing Data Models for the Global Alliance for Genomics and Health

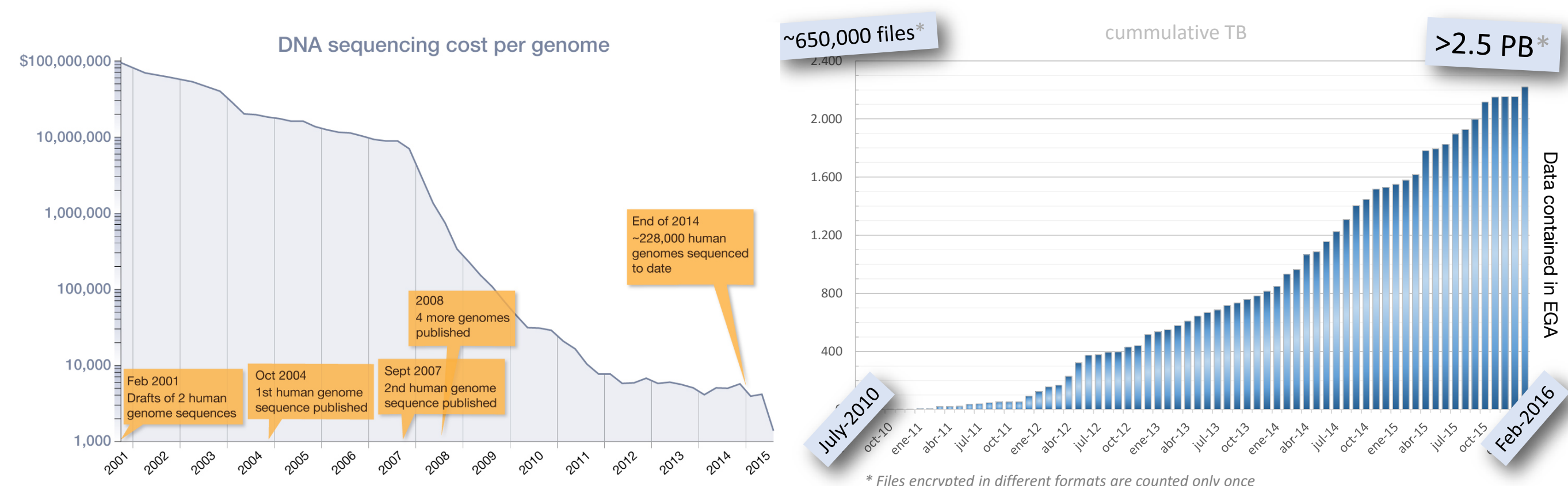


ELIXIR All Hands 2017, 21-23 March, Rome, Italy

Overview

- The advent of sequencing age has enriched our understanding of human malignancies.
- As the cost comes down, the amount of data expands exponentially.
- Large scale comparative study of genome variations is crucial for biomedical research.
- However, data resources are scattered behind firewalls.

Cost comes down & Data goes up



Data held in silos & unshared

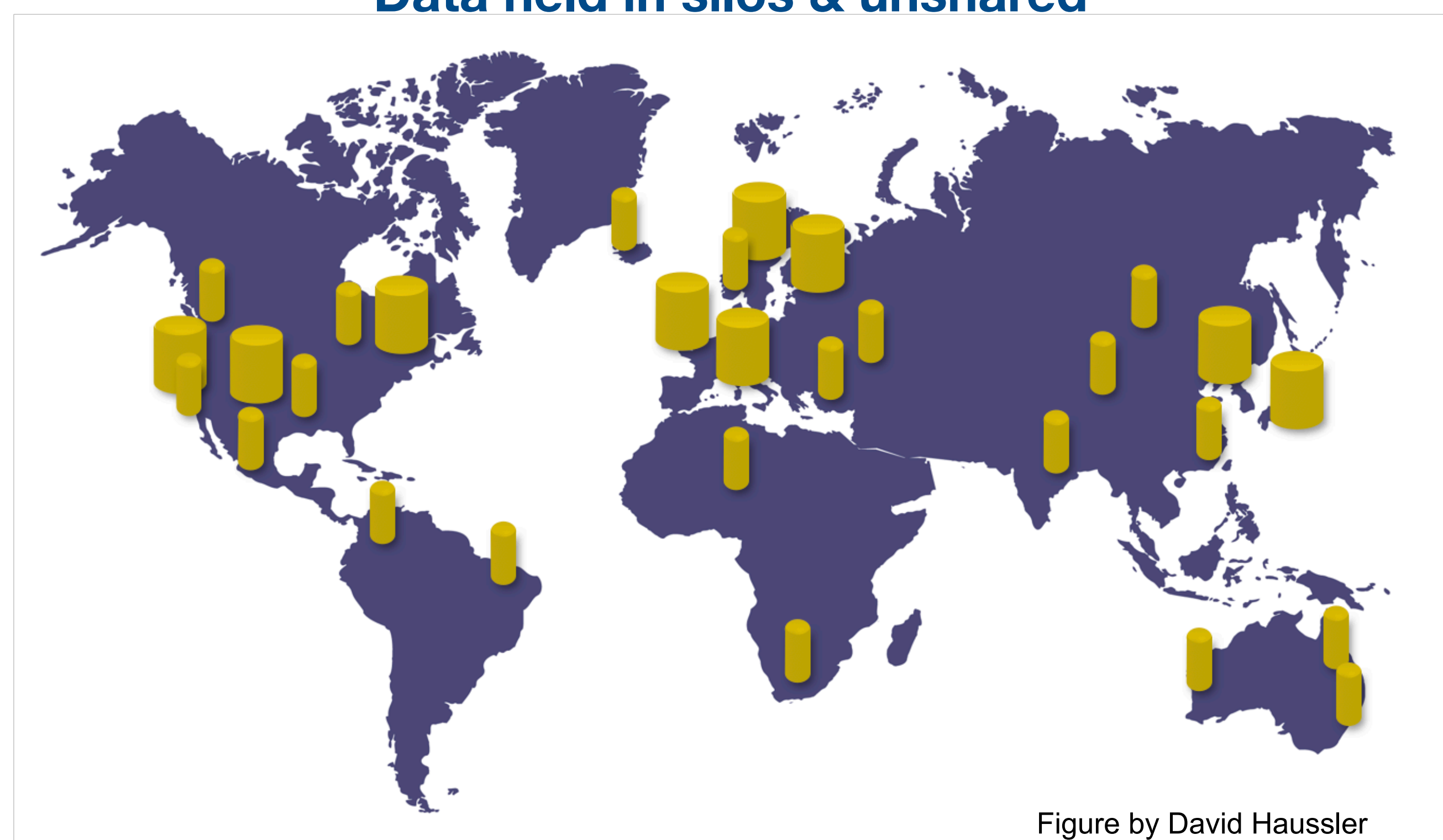


Figure by David Haussler

The arraymap.org Cancer Genome Resource arrayMap

- The arrayMap resource has been established as curated oncogenomic resource, focussing on genomic arrays and copy number aberration (CNA) profiles.
- The underlying data is being extracted from NCBI's Gene Expression Omnibus (GEO), EBI's ArrayExpress, and, importantly, through targeted mining of publication data.
- It is for cancer related genome data and clinical use, such as the diagnostic validations as well as target evaluation for personalized therapeutic approaches.

A Cancer Genome Resource with 60,000+ aCGH arrays

BRAIN TUMOURS	5593 samples ↗	62977 genomic array profiles
BREAST CANCER	8329 samples ↗	914 experimental series
COLORECTAL CANCER	3157 samples ↗	267 array platforms
PROSTATE CANCER	991 samples ↗	245 ICD-O cancer entities
STOMACH CANCER	1062 samples ↗	

Visualization of Cancer Genome Profiling

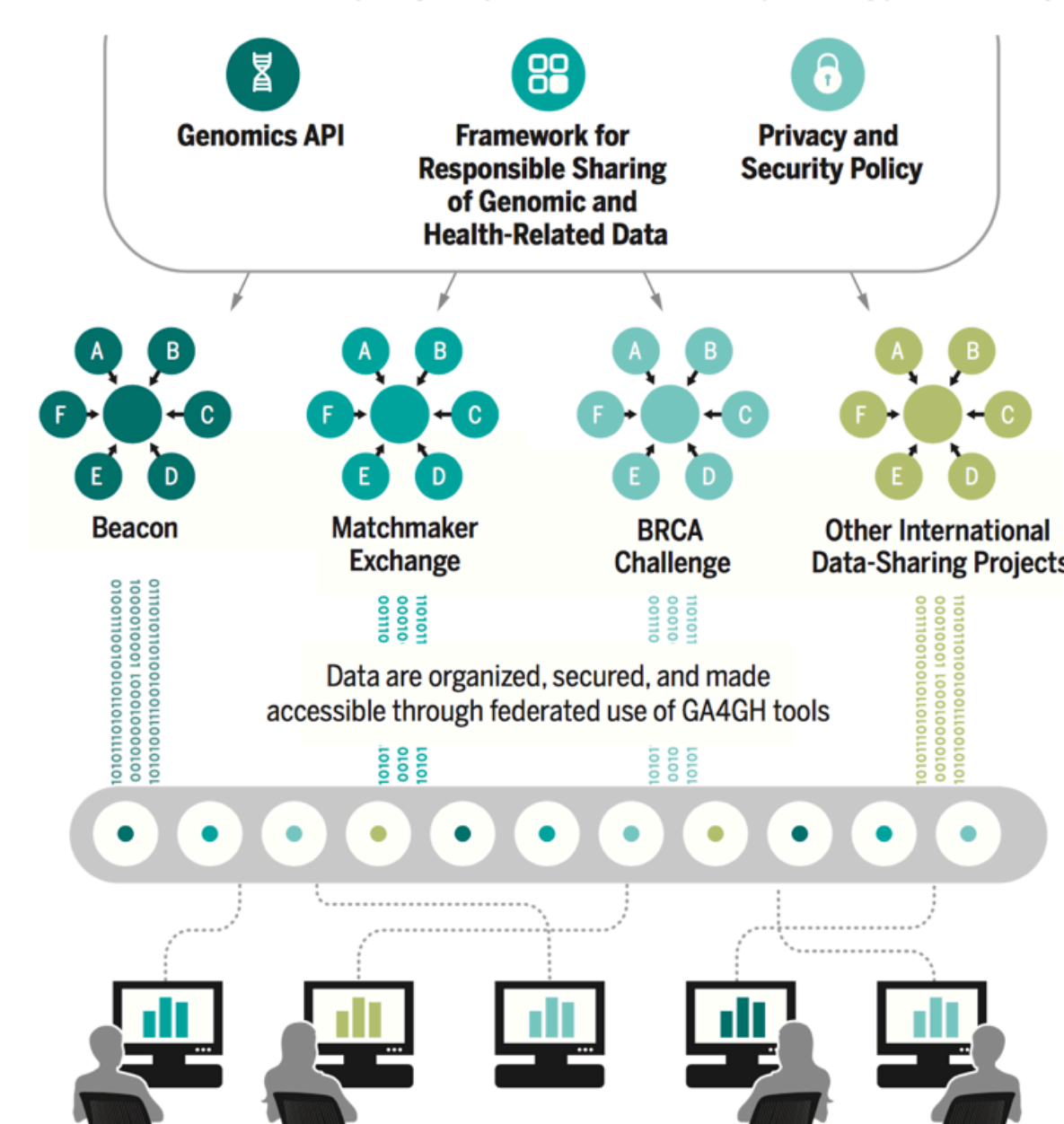


Global Alliance for Genomic and Health

GA4GH was founded by leading scientists in biology, medicine, computational research, data security as well as law and ethics. The aim is:

- to develop standards for the representation and exchange of genome data and supporting information,
- to promote the implementation of legal and ethics frameworks and procedures related with the use of this data for research purposes.

A federated data ecosystem. To share genomic data globally, this approach furthers medical research without requiring compatible data sets or compromising patient identity.



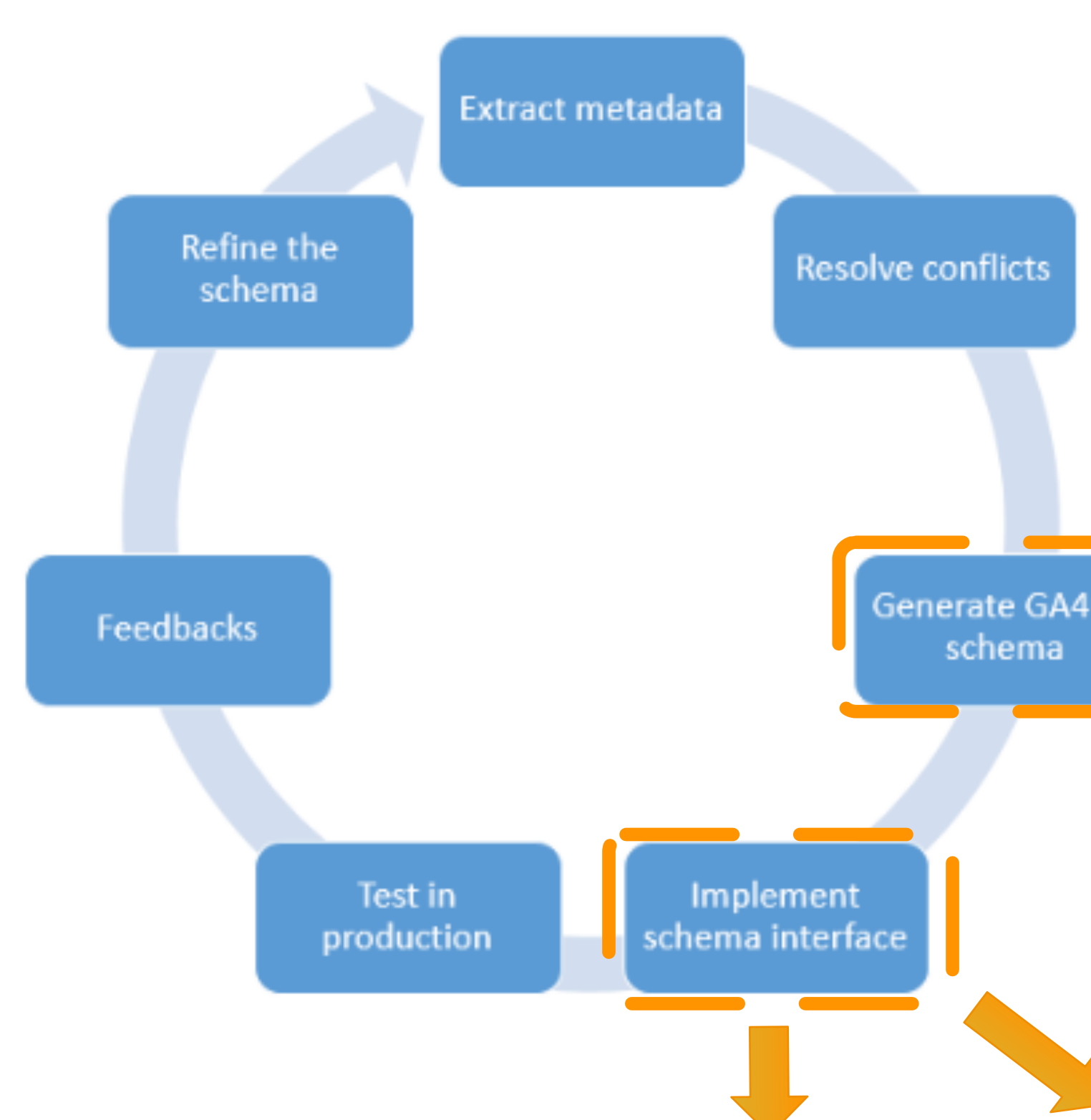
GA4GH task teams

REFERENCE GENOMES TASK TEAM	G2P TASK TEAM
RNASEQ TASK TEAM	VARIANT ANNOTATION TASK TEAM
META DATA TASK TEAM	FILEFORMATS TASK TEAM
CONTAINERS AND WORKFLOWS TASK TEAM	BENCHMARKING TASK TEAM
BEACON TASK TEAM	MATCHMAKER TASK TEAM
READS TASK TEAM	CANCER GENE TRUST
DIRECTORY AND STREAMING API	VICC

Schema Development

- The arrayMap to GA4GH development pioneers the definition of data formats and software implementations for genomic and associated metadata.
- We are developing modern data schemas to facilitate annotation and mining of biomedical attributes as well as provenance of physical or procedural objects related to genomic data.
- Since the first prototype version of the GA4GH schema in 2014, considerable progress has been made in the schema development and its integration with ontologies.

Iterative Efforts



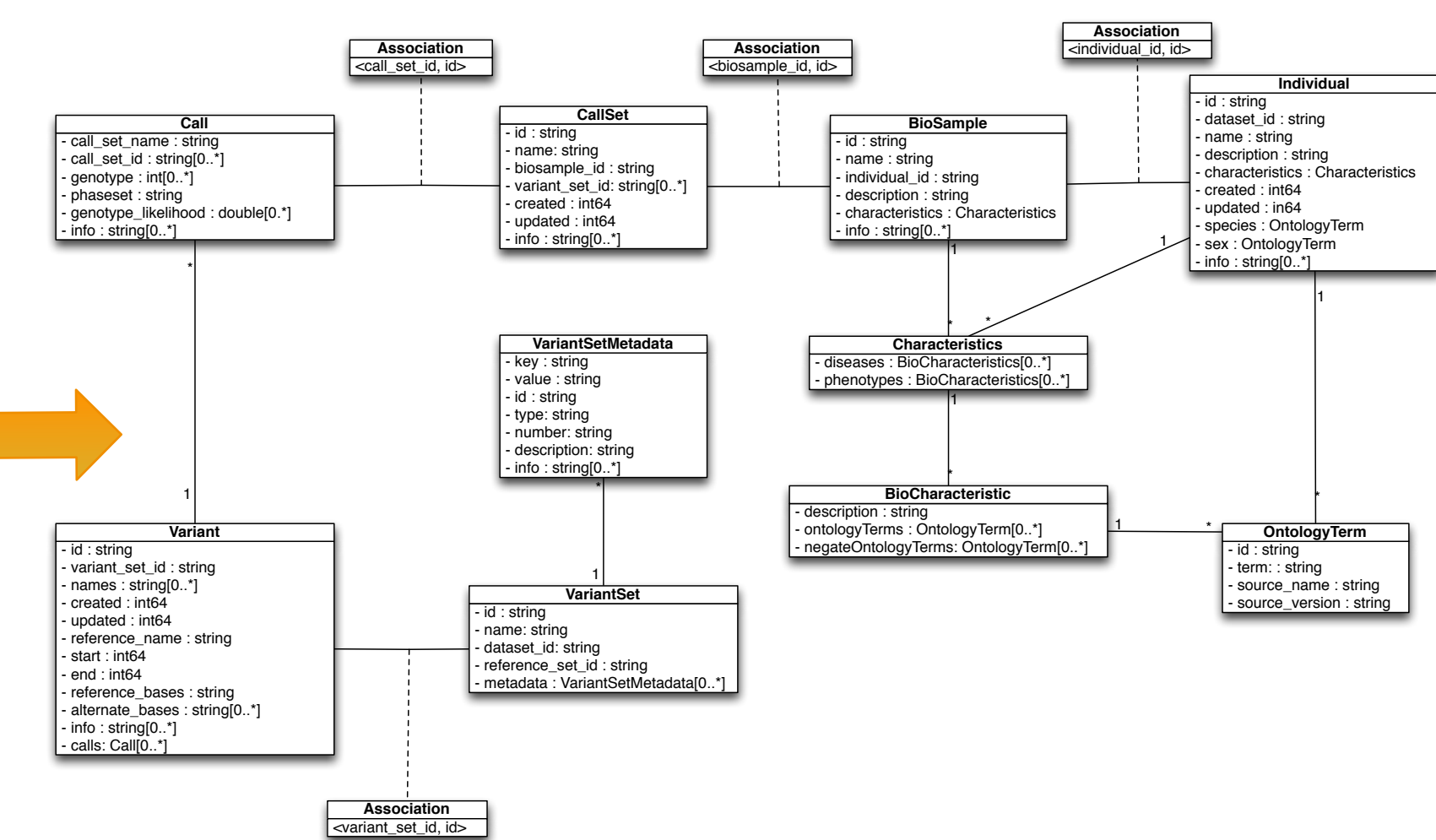
arrayMap-beacon Implementation

Beacon arrayMap
Beacon v0.4 Implementation for arrayMap.

Reference name: 9
Start: 42049214
Length: 26740
Assembly ID: GRCh36
Dataset Ids: (9440/3) 9440/3: Glioblastoma, NOS (2278)
Alternate bases: DEL (Deletion)
Confidence Interval (Start position): 500
Confidence Interval (End position): 500
Match type: Any

Beacon Query Beacon Info

arrayMap Data Translation



arrayMap-ga4gh Implementation @github.com/progenetix/arraymap2ga4gh

```
1 {
2   "url": "https://progenetix.org/arraymap2ga4gh/13658888",
3   "name": "13658888",
4   "individual_id": "PGIND_090322223",
5   "id": "AP_090322223",
6   "characteristics": [
7     {
8       "diseases": [
9         {
10           "ontologyTerms": [
11             {
12               "termLabel": "Chronic Lymphocytic Leukemia",
13               "termId": "NCIT:C3163",
14             },
15             {
16               "termLabel": "B-cell chronic lymphocytic leukemia/small lymphocytic lymphoma",
17               "termId": "SNOMED-98233",
18             },
19             {
20               "termLabel": "B-cell chronic lymphocytic leukemia/small lymphocytic lymphoma",
21               "termId": "SNOMED-98233",
22             },
23             {
24               "termLabel": "hematopoietic and reticuloendothelial systems",
25               "termId": "ICD01:C42",
26             },
27             {
28               "negatedOntologyTerms": [
29               ],
30               "description": "Chronic Lymphocytic Leukemia"
31             }
32           ],
33           "phenotypes": [
34             {
35               "term": "T1",
36               "death": "B",
37               "country": "Sweden",
38               "age": "17-64",
39               "redirection": "null",
40               "followup_months": 60,
41               "age": "17-64",
42               "pubmed_id": "13484635",
43               "sex": "female",
44               "age": "50",
45               "city": "Uppsala"
46             }
47           ],
48           "update": "2017-02-18T17:15:02.388Z",
49           "create": "2017-02-18T17:15:02.388Z"
50         }
51       ]
52     }
53   ]
54 }
```

Contact:

Bo Gao, Michael Baudis
Institute of Molecular Life Sciences,
Swiss Institute of Bioinformatics,
University of Zürich, Switzerland
Email: bo.gao@imls.uzh.ch



Universität
Zürich ^{UZH}



Swiss Institute of
Bioinformatics



Global Alliance
for Genomics & Health



@ELIXIREurope

