

# The Progenetix oncogenomic resource in 2021

Qingyao Huang<sup>1,2</sup>, Paula Carrio-Cordo<sup>1,2</sup>, Bo Gao<sup>1,2</sup>, Rahel Paloots<sup>1,2</sup> and Michael Baudis<sup>1,2,\*</sup>

<sup>1</sup>Department of Molecular Life Sciences, University of Zurich, Winterthurerstrasse 190, Zurich 8057, Switzerland

<sup>2</sup>Swiss Institute of Bioinformatics, Winterthurerstrasse 190, Zurich 8057, Switzerland

\*Corresponding author: Tel: +41 44 635 34 86; Email: [michael.baudis@mls.uzh.ch](mailto:michael.baudis@mls.uzh.ch)

Citation details: Huang, Q., Carrio-Cordo, P., Gao, B. *et al.* The Progenetix oncogenomic resource in 2021. *Database* (2021) Vol. 2021: article ID baab043; DOI: <https://doi.org/10.1093/database/baab043>

## Abstract

In cancer, copy number aberrations (CNAs) represent a type of nearly ubiquitous and frequently extensive structural genome variations. To disentangle the molecular mechanisms underlying tumorigenesis as well as identify and characterize molecular subtypes, the comparative and meta-analysis of large genomic variant collections can be of immense importance. Over the last decades, cancer genomic profiling projects have resulted in a large amount of somatic genome variation profiles, however segregated in a multitude of individual studies and datasets. The Progenetix project, initiated in 2001, curates individual cancer CNA profiles and associated metadata from published oncogenomic studies and data repositories with the aim to empower integrative analyses spanning all different cancer biologies. During the last few years, the fields of genomics and cancer research have seen significant advancement in terms of molecular genetics technology, disease concepts, data standard harmonization as well as data availability, in an increasingly structured and systematic manner. For the Progenetix resource, continuous data integration, curation and maintenance have resulted in the most comprehensive representation of cancer genome CNA profiling data with 138 663 (including 115 357 tumor) copy number variation (CNV) profiles. In this article, we report a 4.5-fold increase in sample number since 2013, improvements in data quality, ontology representation with a CNV landscape summary over 51 distinctive National Cancer Institute Thesaurus cancer terms as well as updates in database schemas, and data access including new web front-end and programmatic data access.

Database URL: [progenetix.org](https://progenetix.org)

## Introduction

Copy number aberrations (CNAs) are present in the majority of cancer types and exert functional impact in cancer development (1, 2). As understanding cancer biologies remains one of the main challenges in contemporary medical and life sciences, the number of studies addressing genomic alterations in malignant diseases continues to grow. Progenetix is a publicly accessible cancer genome data resource ([progenetix.org](https://progenetix.org)) that aims to provide a comprehensive representation of genomic variation profiles in cancer, through providing sample-specific CNA profiles and associated metadata as well as services related to data annotation, meta-analysis and visualization. Originally established in 2001 with a focus on data from chromosomal Comparative Genomic Hybridization (CGH) studies (3), the resource has progressively incorporated data from hundreds of publications reporting on genome profiling experiments based on molecular cytogenetics (CGH, genomic arrays) and sequencing (whole-genome or whole-exome sequencing—WGS or WES). Since the last publication dedicated to the Progenetix resource in 2014 (4), changes in content and features of the data repository and its online environment have vastly expanded its scope and utility to the cancer genomics community. For data content, additions include the complete incorporation of the previously separate arrayMap data collection (5, 6) and of datasets from external resources and projects such as The Cancer

Genome Atlas (TCGA; (7, 8)) or cBioPortal (9), as well as the recurrent collection and re-processing of array-based data from National Center for Biotechnology Information (NCBI)'s Gene Expression Omnibus (GEO) or European Molecular Biology Laboratory-European Bioinformatics Institute (EMBL-EBI)'s ArrayExpress (10, 11). Additionally, data content updates have followed the previous methodology of publication-based data extraction where feasible. Beyond the data expansion, a tight integration with projects of the Global Alliance for Genomics and Health (GA4GH (12)) and ELIXIR—such as serving for implementation-driven development of the Beacon application programming interface (API) (13)—has led to an extension of the resource's features as well as adoption and promotion of emerging open data standards.

Here we present the latest updates on data content, structuring, standardization, access and other modifications made to the Progenetix resource.

## Data expansion and new features

### Genomic profiling data

Over the last two decades, thousands of cancer genome studies have used the GEO (14) for deposition of data from array-based experiments. Data from GEO contribute a substantial fraction of the genomic screening data in the Progenetix collection and has again been expanded in both number of

Received 23 February 2021; Revised 16 June 2021; Accepted 30 June 2021

© The Author(s) 2021. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

**Table 1.** Statistics of samples from various data resources

| Data source           | GEO       | ArrayExpress | cBioPortal | TCGA      | Total      |
|-----------------------|-----------|--------------|------------|-----------|------------|
| No. of studies        | 898       | 51           | 38         | 33        | 1939       |
| No. of samples        | 63 568    | 4351         | 19 712     | 22 142    | 138 663    |
| Tumor                 | 52 090    | 3887         | 19 712     | 11 090    | 115 357    |
| Normal                | 11 478    | 464          | 0          | 11 052    | 23 306     |
| Classifications       |           |              |            |           |            |
| ICD-O (Topography)    | 100       | 54           | 88         | 157       | 209        |
| ICD-O (Morphology)    | 246       | 908          | 265        | 140       | 491        |
| NCIt                  | 346       | 148          | 422        | 182       | 788        |
| Collections           |           |              |            |           |            |
| Individuals           | 63 568    | 4351         | 19 712     | 10 995    | 127 549    |
| Biosamples            | 63 568    | 4351         | 19 712     | 22 142    | 138 663    |
| Callsets <sup>a</sup> | 63 568    | 4351         | 19 712     | 22 376    | 138 930    |
| Variants              | 5 514 126 | 118 4170     | 1 778 096  | 2 654 065 | 10 716 093 |

<sup>a</sup>set of variants from one genotyping experiment; ICD-O, International Classification of Diseases for Oncology; NCIt, National Cancer Institute Thesaurus.

**Table 2.** Data growth by cancer loci

| Cancer loci                                     | No.in 2014 | No.in 2021 |
|---|------------|------------|
| Hematopoietic and reticuloendothelial systems   | 5269       | 18 482     |
| Lymph nodes                                     | 2345       | 5988       |
| Breast  | 2271       | 15 790     |
| Cerebellum                                      | 1439       | 3465       |
| Brain, NOS                                      | 1342       | 6608       |
| Cerebrum  | 1201       | 1712       |
| Liver   | 1180       | 3237       |
| Stomach   | 1155       | 3176       |
| Skin  | 1073       | 3343       |
| Connective, subcutaneous and other soft tissues | 1058       | 2526       |
| Kidney  | 1018       | 3617       |
| Colon   | 1001       | 5182       |
| Ovary   | 733        | 3963       |
| Prostate gland                                  | 735        | 4485       |
| Lung and bronchus                               | 699        | 10 321     |
| Nervous system, NOS                             | 667        | 926        |
| Urinary bladder                                 | 587        | 1961       |
| Cervix uteri                                    | 529        | 1331       |
| Peripheral nerves incl. autonomous              | 523        | 1479       |
| Esophagus                                       | 454        | 1890       |
| Pancreas  | 426        | 1620       |
| Thyroid gland                                   | 404        | 1260       |
| Heart, mediastinum and pleura                   | 383        | 771        |
| Bones, joints and articular cartilage           | 350        | 1205       |
| Spleen  | 278        | 636        |
| Other   | 4522       | 16 268     |
| Total   | 31 642     | 115 359    |

samples and represented platforms. Additionally, we systematically included suitable data from three more resources: ArrayExpress (15), cBioPortal (16) and TCGA(17) project. As in the previous database updates, we have also included data directly derived from publication supplements and from collaborative projects. Table 1 shows statistics of samples within the major sources. Table 2 reports the overall data growth and sample counts stratified by cancer loci since the last update (4).

The ‘ArrayExpress Archive of Functional Genomics Data’, hosted by EMBL-EBI, stores functional genomics data submitted by research groups and projects. In this update, we have incorporated the cancer-related genomic profiles which do not have corresponding GEO entries using our analysis pipeline.

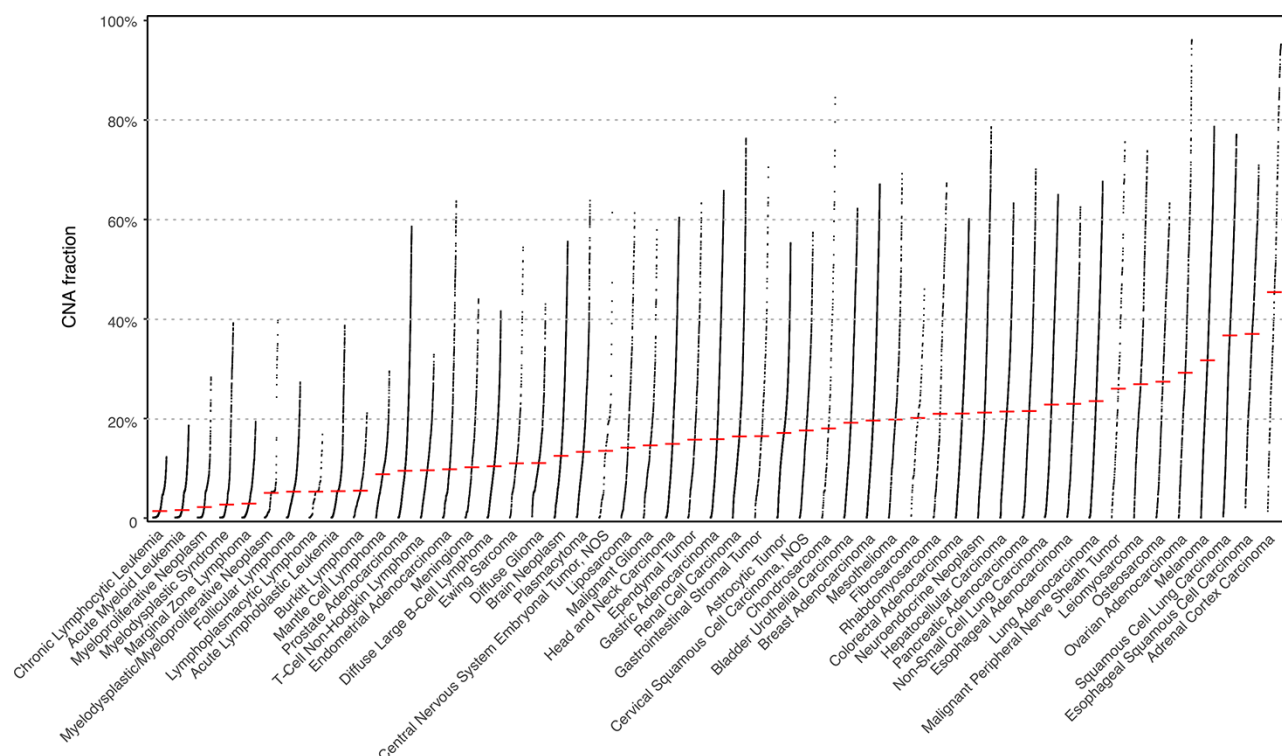
Overall, data from ArrayExpress added 3887 samples from 44 projects, which resolve to 143 distinct cancer types according to the National Cancer Institute Thesaurus (NCIt). Similar to the GEO data acquisition procedure, we have used a combination of text mining methods and expert curation for annotation of technical metadata and biomedical parameter.

The ‘cBioPortal for Cancer Genomics’ is an open-access resource for cancer genomics data, representing different types of molecular screening data from 19 712 samples, derived from 38 studies and mappable to 422 NCIt cancer types. The largest part of genomic data is based on WES analyses from the Memorial Sloan Kettering-Integrated Mutation Profiling of Actionable Cancer Targets or MSK-TARGET (18) pipeline, with CNA data accessed directly as segment files in genome version hg19/Genome Reference Consortium Human Build 37. Data were converted into Genome Reference Consortium Human Build 38 (GRCh38) with the ‘segment-liftover’ tool (19), and oncology classifications as well as relevant clinical data were incorporated into our database.

TCGA project provides a set of multiomics data with extensive structured metadata annotation for a large collection of cancer types, currently through NCBI’s Genomic Data Commons Data Portal (<https://portal.gdc.cancer.gov>). In this update, we incorporated its copy number variation (CNV) profiling data as well as transformed the relevant clinical information into our system (Figure 1).

## Data processing update

Genomic profiling data in Progenetix originates from a large number of studies, which are based on different molecular-cytogenetics- and sequencing-based technologies. In order to maximize qualitative homogeneity of the final CNA calls, we prefer to download source files with the least amount of pre-processing and apply our in-house data processing pipeline from the arrayMap project (5). Currently, our analysis workflow handles the raw-data-based processing for 13 Affymetrix single nucleotide polymorphism (SNP) array platforms, including nine genome-wide arrays—10K (GPL2641), 50K (Hind240 and Xba240; GPL2004 and GPL2005), 250K (Nsp and Sty; GPL3718 and GPL3720), Genome-wide SNP (5.0 and 6.0; GPL6894 and GPL6801, respectively), CytoScan (750K and HD; GPL18637 and GPL16131) arrays (GPL-prefixed platform coding in brackets according



**Figure 1.** The currently available CNA data points in Progenetix and TCGA Progenetix database contain 115 357 cancer samples with 92 307 mapped to the 51 defined critical nodes in NCI ontology tree and 23 050 samples not mapped to the tree (black), whereas TCGA repository contains 11 090 samples with 9103 samples mapped and 1987 samples not mapped to the tree (black). Colors of the stacked bar plot (left) match the branch colors on NCI ontology tree (right).

to GEO standard)—as well as the four cancer-specific ‘Oncoscan’ arrays - GPL18602, GPL13270, GPL15793 and GPL21558 (accessible through GitHub repository *baudis-group/a.m.\_process*). Our current model treats the most prevalent copy number as the baseline and derives the relative copy number gain and loss per sample based on the assumption that the relative gene dosage imbalance exerts pathophysiological effects in cancer biology.

### Allele-specific copy number variation

For the subset of SNP-array-based experiments—where the status of both alleles can be evaluated separately—we have analyzed allele-specific copy number data (ASCN) and incorporated 35 897 loss of heterozygosity (LOH) profiles into the database. ASCN potentiates new analysis on the same samples. First, probe-wise it gives an overview of germline variant landscape, as used in determining the ancestry background. Second, it allows detection of LOH events, including copy-number-neutral event (CN-LOH), which e.g. can be commonly observed in hematological malignancies due to a selective process for duplication of minor disease-prone germline alleles (20, 21). Lastly, it acts as a second reference for CNA to combat the variability caused by known wave artifacts from array technologies (22). For all SNP arrays, we have implemented a pipeline to determine probe-wise B-allele frequency (BAF) of SNP probes and perform subsequent segmentation (23, 24). Subsequently, we use ASCN to assess ancestry provenance of the samples (25) and store the LOH regions of the samples in our genomic variants database.

## Metadata updates

### NCIt ontology mapping

Since its establishment, Progenetix has made use of the ‘International Classification of Diseases in Oncology’, 3rd Edition (ICD-O-3) (26) for cancer sample classification. While the combination of the ICD-O Morphology and Topography coding systems depicts diagnostic entities with high specificity, the current ICD-O is limited in its representation of hierarchical concepts and does not easily translate to modern ontologies. In comparison, NCIt (access through <http://bioportal.bioontology.org/ontologies/NCIT>) is a dynamically developed hierarchical ontology, which empowers layered data aggregation and transfer between classification systems and resources. However, due to the comparatively recent development and ongoing expansions, NCIt terms are rarely used in primary sample annotations. In the recent Progenetix update, we performed a data-driven generation of ICD-O—NCIt mappings and added the derived NCIt codes to all (existing and new) samples (mapping available through GitHub repository ‘progenetix/ICDOntologies’; manuscript in preparation), to take advantage of NCIt’s hierarchical structure for data retrieval, analysis and exchange (Figure 4B).

### Data summary based on the NCIt hierarchy tree

All cancer samples in Progenetix have been annotated with an NCIt code, resulting in currently 788 distinct NCIt terms. However, as the definition of increasingly specific NCIt terms outruns their incorporation into the hierarchical tree, so far 98 of these terms are not represented in the tree hierarchy. For

better illustration, we define 51 prominent nodes under which we summarize and visualize the data collection (see Supplementary material for the selection procedure). This brings about additional 324 (60 in TCGA) terms not mappable to the selected nodes, resulting in 23 050 (1987 for TCGA) samples excluded from the summary tree counts (black bar in left panel of Figure 1). For terms with multiple occurrences in the tree, we define the preferred path to the selected node by prioritizing morphology-based separation. The sample collection in Progenetix compared to TCGA is summarized with reference to the NCIt coding system (Figure 1; Supplementary Table S1).

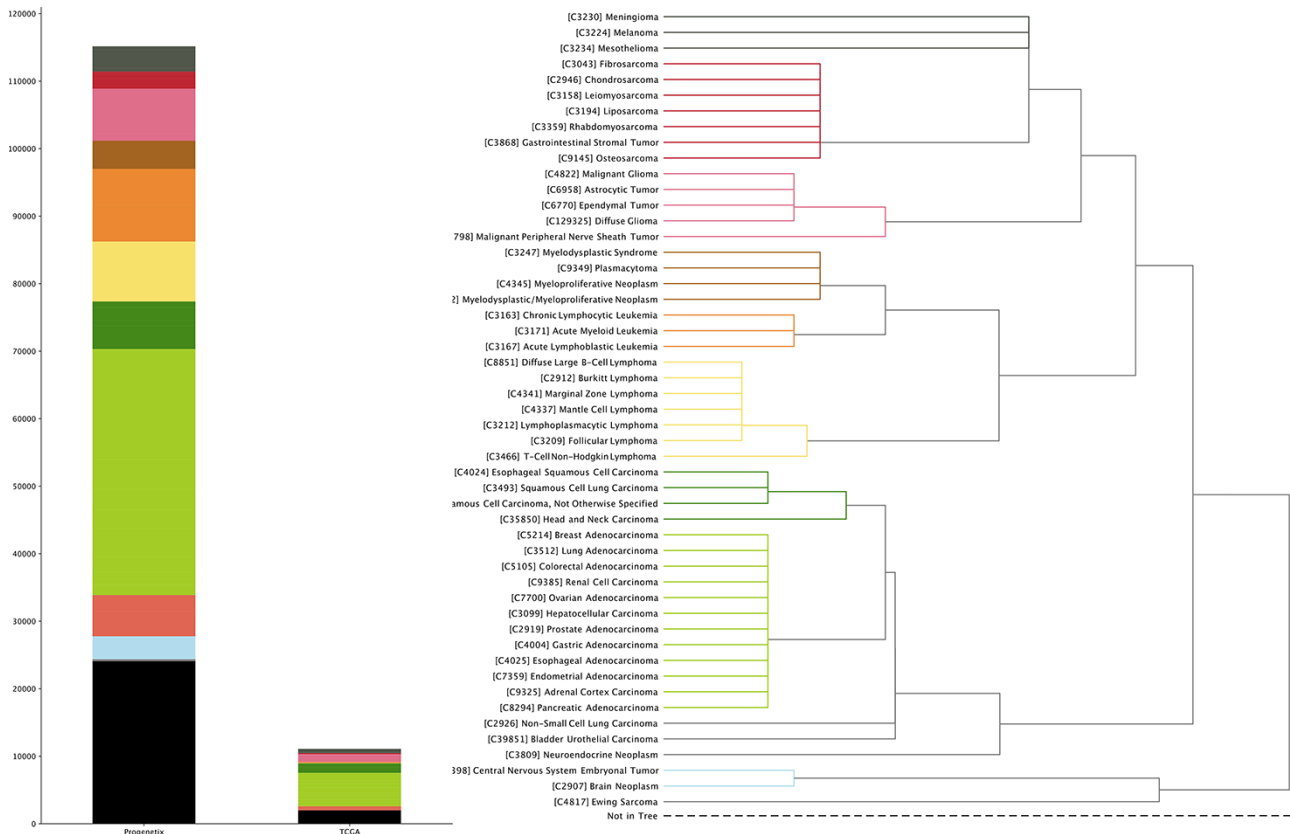
### CNV data content by cancer type

With cancer genomes grouped in the 51 NCIt nodes, we assessed their differences in the CNV landscape. The fraction of genome with a copy number alteration (CNV fraction) varies widely among the cancer types with a global median of 0.121 (Figure 2; Supplementary Figure S1). Among the most studied cancer types, breast carcinoma shows a consistent CNV profile as an earlier analysis with frequent chr1q, 8q, 16p, 17q, 20 gain and 8p, 16q, 17p, 18, 22q loss (27); the CNV patterns in cervical (chr3 gain) and colorectal (chr7, 8q, 13, and 20q gain and 8p, 17p, and 18 loss) carcinoma also correspond with previous observation (28), similar to T-cell non-Hodgkin lymphoma (29), myelodysplastic syndrome (30) and a number of malignant

epithelial tumors (31). In addition, we also present the genome-wide LOH profile in the evaluated NCIt nodes clustered by their LOH landscape (average LOH profiles of 42 out of 51 with at least 20 samples are shown in Supplementary Figure S2; (32)). LOH profile of a cancer genome complements its CNV profile with the information of allelic loss. Here we highlight a few prominent patterns, which have been previously reported: chr3p and 9 in esophageal squamous cell carcinoma (33, 34); chr18q in colorectal carcinoma (35); and chr13q, 16q and 17p in hepatocellular carcinoma (36).

### Uberon anatomy ontology

While the ICD-O topography system provides organ- and substructure-specific mapping rooted in traditional clinical and diagnostic aspects of a ‘tumor entity’, ‘UBERON’ is a cross-species anatomical structural ontology system closely aligned with developmental processes (37). Its relationship structure allows integrative queries linking multiple databases (e.g. Gene Ontology (38) and Protein Ontology (39)) and description logic query within the same organism (linking related organs) and between model animals and humans. In this resource update, we have mapped all existing ICD-O T codes to ‘UBERON’ terms and additionally provided those as part of the ‘Monarch’ initiative (40), with our latest mapping table (made available through a GitHub repository ‘progenetix/icdot2uberon’).



**Figure 2.** The genomic CNV fraction across 51 NCIt umbrella nodes. Each dot represents one sample's CNV fraction range from 0 to 1 and the red horizontal line indicates median CNV of the respective cancer type. Each cancer type contains between 104 and 11 804 CNV profiles (median 904; See Supplementary Table S1).



## Provenance by geography

As part of the curated metadata provided in the sample representation, we have included geographic point coordinates for each individual sample. As this information is often missing from individual sample annotations, we have previously applied a mapping procedure to assign the samples' approximate geographic origins (41, 42). For samples with the submitter's contact available from repository entries, a default point location in the corresponding city was used—otherwise that of the corresponding author of the associated publication was used. Associated publications were also explored for more detailed descriptions of sample origin. Point coordinates for each city were obtained using the external geographic database GeoNames ([www.geonames.org](http://www.geonames.org)), as detailed previously.

## Provenance by ancestry group

While providing a good approximation for the geographic origin of cancer profiling data, which can e.g. be useful for epistemic validation and decision processes, the geographic location of the studies provides limited specificity regarding individual sample provenance, especially when assessing correlations between genomic variants and ancestral population background. Beyond the scope of high-penetrance variants like mutations in the BRCA1/2 (43, 44) or RetinoBlastoma (RB) genes (45) in cancer predisposition, other studies have asserted an influence of genetic background on tumor development (46–49). Previously we have developed a method for deriving ancestry groups from unmasked germline variants in cancer genomes, based on reference populations studied in the 1000 Genomes Project (25). For samples in Progenetix with accessible SNP data, population groups were assigned based on the reference categories mapped to Human Ancestry Ontology terms (Supplementary Table S2). Where available, the respective data are now represented under the 'populations provenance' schema for the corresponding biosample entries.

## Updated data access modalities

Since the last release, we have adopted the GA4GH data schema standards and migrated to Phenopackets (50)-formatted response delivery with modified data access points in the user interface. Information about API methods are provided through the documentation pages (<https://info.progenetix.org/categories/API>).

## Data standards

In many genomic repositories, databases are structured around experimental outcomes (e.g. variants from a DNA sequencing experiments as collections of VCF files). Recent attempts in evaluating sensible meta-schemas for the representation of genomic variants and related biological or technical metadata, especially with respect to empowering data federation over flexible, networked resources, have led to a set of emerging meta-models and data schemas (51). The data storage and representation models for the Progenetix resource have been designed to comply with concepts developed by the previous GA4GH Data Working Group (12, 52) and subsequent GA4GH work streams, documented e.g. by the 'SchemaBlocks' initiative (<http://schemablocks.org>). One of

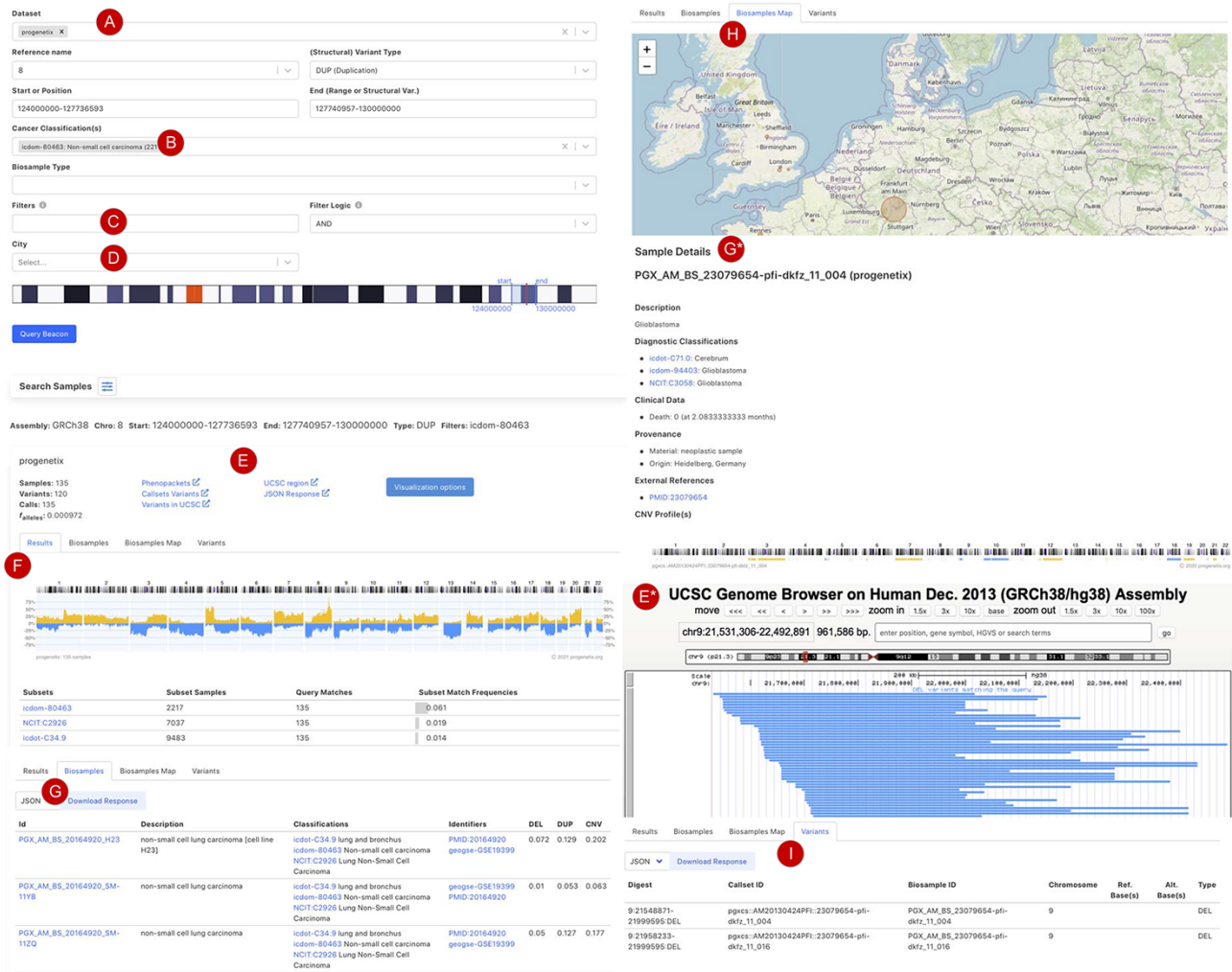
the core concepts is the 'individual—biosample(s) - variants' meta-model, which is applicable to cancer-related analyses with potentially multiple samples representing different stages in the course of disease as well as the underlying genomic background. This hierarchical model provides a solid representation and connection between the physical source of the data and the logical genotyping information and adapts to various scenarios for data aggregation and analysis.

## User interface

The completely re-designed user interface provides flexibility and versatility in query parameters and types and optimized the response delivery. Technically, the query interface for retrieval of sample specific data is built on top of a forward-looking implementation of the GA4GH Beacon API (13) with features from the upcoming version 2 of this standard.

Figure 3 shows the current web interface to perform a CNA query with start and end position range with filter options for cancer type, tissue location, morphology, cell line or geographic location. The top panel of the result page shows a summary with the number of matched samples, variants, calls and the frequency of alleles containing the CNA (Figure 3E). The 'Phenopackets' link returns a json document of biosamples with the phenopacket-formatted response. The 'UCSC region' links externally to a University of California Santa Cruz (UCSC) browser track providing an overview of the genomic elements which map to the region of the observed variants. Also, customized visualization is enabled in the linked page 'visualization options', e.g. for selected chromosomal regions and grouping by subsets or studies. The lower panel is organized in four sections: (i) the 'Result' tab (Figure 3F) shows the genome-wide CNA by the percentage of samples with yellow (+) as CN gain and blue (–) as CN loss. Below the CNA plot is a table showing the list of subsets as defined by ICD-O-3 and NCI Ontology terms sorted by frequency of matched samples within that subset. (ii) the 'Biosamples' tab (Figure 3G) shows information of matched biosamples, i.e. description, classifications and external identifiers. The table can be downloaded in json or csv format. The further detail of the biosample can be accessed by clicking the biosample id. (iii) The 'Biosamples Map' tab (Figure 3H) shows a world map with the matched geological locations highlighted. (iv) the 'Variants' tab (Figure 3I) shows the variant 'digest' (concatenated format with chromosome, start and end position, and type of the CNA) and its corresponding biosample and callset. Likewise, the table can be downloaded in json or csv format.

Figure 4 shows the additional functional interfaces and services provided by the Progenetix project. Users can search for publications or studies by publication title, author names or the geographic location of the research center. Then, navigation extends to the summary of publications with the number of samples catalogued by technology and availability in database as well as options to visualize the associated samples (Figure 4A). Users can also access samples from the NCI hierarchical tree or other classification systems (e.g. ICD-O and UBERON) to select a subset of cancer types for summary statistics and visualization (Figure 4B). Alternatively, users can also upload their own data for single or multiple samples to visualize genome-wide CNA (Figure 4C). In addition, a list of studies and cohorts can be selected in the navigation menu, including arrayMap (probe-specific arrays from



**Figure 3.** Beacon-style query using fuzzy ranges to identify biosamples with variants matching the CNA range This example queries for a continuous, focal duplication covering the complete MYC gene's coding region with  $\leq 6$  Mb in size. A: Filter for dataset; B: filter for cancer classification (NCIt and ICD-O-3 ontology terms available); C: additional filter, e.g. CelloSaurus; D: additional filter for geographic location; E: external link to UCSC browser to view the alignment of matched variants; F: cancer type classification sorted by frequency of the matched biosamples present in the subset; G: list of matched biosamples with description, statistics and reference. More detailed biosample information can be viewed through 'id' link to the sample detail page; H: matched variants with reference to biosamples can be downloaded in json or csv format.

published studies (5)), diffuse intrinsic pontine glioma cohort (53) and the 'cancer signature' cohort (54). All the functionalities and provided services are detailed in the documentation pages at [info.progenetix.org](http://info.progenetix.org), which invite request submission through the GitHub 'issues' tracker.

Other improvements

Genome version update

All samples have been updated to GRCh38. The process has been completed in a step-wise manner. Preferably, for samples with available probe-specific array data, either GRCh38 mapped platform data files were used for re-processing of the original files or alternatively, a lift-over of the probe data and subsequent re-segmentation was performed. For those cases where only called CNA data had been collected, we applied our recently published 'segment-liftover' tool (19) for the efficient re-mapping of continuous segments. Overall, more than 99.99% of probes and more than 99% of segments could be recovered successfully.

Cell line collection

Cancer cell lines are important models for understanding the molecular mechanisms of malignant diseases and have a prominent role in pharmacological screening procedures. Besides the primary tumor data, the Progenetix data collection also includes genomic profiling experiments using *in vitro* models. Recently, we introduced a systematic update of cell line annotations based on 'CelloSaurus', a comprehensive knowledge resource on cell line data with extensive annotations and mappings to a variety of classifications and ontologies (55). We meticulously assigned CelloSaurus ids for the cancer cell line samples as well as the ICD-O morphology and topography codes based on the NCIt term annotated by CelloSaurus. At this time, Progenetix includes a total of 5764 samples corresponding to 2162 different cancer cell lines, representing 259 different cancer types (NCIt). While so far we provide the option to search for cell lines by applying a 'cellosaurus' filter either in the web interface (e.g. 'cellosaurus: CVCL\_0030' for 'HeLa' cell line samples) or in the API query, work on a dedicated cell line data access tool is underway.

## A Progenetix Publication Collection

The current page lists articles describing whole genome screening (WGS, WES, aCGH, cCGH) experiments in cancer, registered in the Progenetix publication collection. For each publication the table indicates the numbers of samples analysed with a given technology and if sample profiles are available in Progenetix and/or arraymap (array source files).

Please [contact us](#) to alert us about additional articles you are aware of.

Search

City

Range (km)

| Publications (37) |   | Samples |      |     |     |     |    |
|-------------------|---|---------|------|-----|-----|-----|----|
| id                | Publication   | cCGH    | aCGH | WES | WGS | pgx | am |
| PMID:29739401     | Kiz\ndig P, Giesen C, Jackson H, Bodenmiller B, Papassotiropoulos et al. (2018): Limited utility of tissue micro-arrays in detecting intra-tumoral heterogeneity in stem cell characteristics and tumor ...<br>J Transl Med 16(1), 2018   | 0       | 372  | 0   | 0   | 0   | 0  |
| PMID:29556019     | Pillonel V, Juskevicius D, Ng CKY, Bodmer A, Zetti et al. (2018): High-throughput sequencing of nodal marginal zone lymphomas identifies recurrent BRAF mutations.<br>Leukemia, 2018  | 0       | 22   | 8   | 0   | 0   | 0  |
| PMID:27658049     | Riba J, Renz N, Niemöller C, Bleul S, Pfeifer D, Stosch JM, Metzeler KH, Hackanson B, Lübbert M, Duyster J, Koltay P, Zengerle R, Claus R, Zimmermann S, Becker H. (2016): Molecular Genetic Characterization of Individual Cancer Cells Isolated via Single-Cell Printing. ...<br>PloS one   | 0       | 3    | 0   | 0   | 3   | 0  |
| PMID:27491809     | Fusco N, Geyer FC, De Filippo MR, Martelotto LG, Ng CK, Piscuoglio S, Guerini-Rocco E, Schultheis AM, Fuhrmann L, Wang L, Jungbluth AA, Burke KA, Lim RS, Vincent-Salomon A, Bamba M, Moritani S, Badve SS, Ichihara S, Ellis IO, Reis-Filho JS, Weigelt B. (2016): Genetic events in the progression of adenoid cystic carcinoma of the breast to high-grade ...<br>Modern pathology : an official journal of the United States and Canadian Academy of Pathology, Inc | 0       | 5    | 0   | 0   | 5   | 0  |
| PMID:26632267     | Kovac M, Blattmann C, Ribi S, Smida J, Mueller NS, Engert F, Castro-Giner F, Weischenfeldt J, Kovacova M, Krieg A, Andreou D, Tunn PU, Dürr HR, Rechl H, Schaser KD, Melcher I, Burdach S, Kulozik A, Specht K, Heinemann K, Fulda S, Bielack S, Jundt G, Tomlinson I, Korb J, Nathrath M, Baumhoer D. (2015): Exome sequencing of osteosarcoma reveals mutation signatures reminiscent of BRCA deficiency. ...<br>Nature communications                                | 0       | 113  | 0   | 0   | 113 | 0  |

## B Cancer Types

Cancer Classification:  Dataset:

Filter cancer:  Collapse all Expand 1 level

☒ Chronic Lymphocytic Leukemia (3535)

- NCIT:C3262: Neoplasm (111840 samples)
  - NCIT:C3263: Neoplasm by Site (106563 samples)
  - NCIT:C4741: Neoplasm by Morphology (106398 samples)
  - NCIT:C27134: Hematopoietic and Lymphoid Cell Neoplasm (24961 samples)
    - NCIT:C3161: Leukemia (11451 samples)
      - NCIT:C3172: Myeloid Leukemia (5635 samples)
      - NCIT:C3483: Chronic Leukemia (3622 samples)
      - ☒ NCIT:C3163: Chronic Lymphocytic Leukemia (3535 samples)

## C Data visualization Upload

Drag and drop some files here, or click to select files.

File format

Data has to be submitted as tab-delimited .tsv segment files. An example file is being provided [here](#).

**Figure 4.** Demonstration of further functionality pages: A. Publication search; B. NCIT hierarchical tree navigation A: Cancer-genomics-associated publications are recorded with number of samples stratified by technology used. The publications can be filtered by keywords; B: Part of the sample subsets contained in Progenetix under the hierarchical NCIT classification tree. It allows for selection of sample subsets at different levels; C: User can upload custom segment files for data visualization.



## Conclusion

The Progenetix resource provides an extensive collection of oncogenomic data with a focus on individual genome-wide CNA profiles and the use of modern ontologies and data schemas to render curated biological and technical meta-data, as well as thorough references to external repositories and annotation resources. Through aggregation of data from thousands of individual research studies as well as several consortium-derived collections, to our knowledge Progenetix database currently constitutes the largest public, freely accessible resource for pre-computed CNA profiles and associated phenotypic information and additional metadata dedicated to cancer studies. While the application of uniform genomic data formats and a benchmarked data processing pipeline minimizes biases from separate studies, the forward-looking implementation of emerging ontology standards facilitates the integrative and comparative analysis across a vast range of cancer types. The tight integration with GA4GH product development and standardization processes guarantees the compatibility with emerging data federation approaches and the widest re-utilization of the resource's data. For the future, besides the continuous maintenance and expansion of the existing data types, we will work toward enhancing clinical and diagnostic annotation, expanding cross-database references and the types of genomic variant data as well as active data sharing and integration through networked services and platforms.

## Supplementary data

Supplementary data are available at *Database* Online.

## Acknowledgements

We would like to thank Amos Bairoch for support with the cell line annotations. Improvements in data annotation concepts were highly influenced through the GA4GH community.

## Funding

The Progenetix database does not receive dedicated funding support. Work on the Beacon API has been supported through the the European life-sciences Infrastructure (ELIXIR) Beacon 2019–2021 implementation study and under the BioMedIT Network project of Swiss Institute of Bioinformatics (SIB) and Swiss Personalized Health Network (SPHN). Bo Gao has been recipient of a fellowship by the China Scholarship Council (CSC).

## References

- Hanahan,D. and Weinberg,R.A. (2011) Hallmarks of cancer: the next generation. *Cell*, **144**, 646–674.
- Albertson,D.G., Collins,C., McCormick,F. *et al.* (2003) Chromosome aberrations in solid tumors. *Nat. Genet.*, **34**, 369–376.
- Baudis,M. and Cleary,M.L. (2001) Progenetix. net: an online repository for molecular cytogenetic aberration data. *Bioinformatics*, **17**, 1228–1229.
- Cai,H., Kumar,N., Ai,N. *et al.* (2014) Progenetix: 12 years of oncogenomic data curation. *Nucleic Acids Res.*, **42**, D1055–D1062.
- Cai,H., Kumar,N. and Baudis,M. (2012) ArrayMap: a reference resource for genomic copy number imbalances in human malignancies. *PLoS One*, **7**, e36944.
- Cai,H., Kumar,N. and Bagheri,H.C. *et al.* (2014) Chromothripsis-like patterns are recurring but heterogeneously distributed features in a survey of 22,347 cancer genome screens. *BMC Genomics*, **15**, 82.
- Cancer Genome Atlas Research Network, Weinstein,J.N., Collisson,E.A. *et al.* (2013) The cancer genome atlas pan-cancer analysis project. *Nat. Genet.*, **45**, 10, 1113–1120.
- National Cancer Institute. (2013) *The Cancer Genome Atlas Program*.
- Gao,J., Aksoy,B.A., Dogrusoz,U. *et al.* (2013) Integrative analysis of complex cancer genomics and clinical profiles using the cBioPortal. *Sci. Signal.*, **6**, pl1.
- National Center for Biotechnology Information (NCBI). (2002) *Gene Expression Omnibus*.
- The European Bioinformatics Institute (EMBL-EBI). (2003) *Array-Express*.
- Global Alliance for Genomics and Health. (2016) GENOMICS. A federated ecosystem for sharing genomic, clinical data. *Science*, **352**, 1278–1280.
- Fiume,M., Cupak,M., Keenan,S. *et al.* (2019) Federated discovery and sharing of genomic data using Beacons. *Nat Biotechnol*, **37**, 220–224.
- Edgar,R., Domrachev,M. and Lash,A.E. (2002) Gene expression omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Res.*, **30**, 207–210.
- Athar,A., Füllgrabe,A., George,N. *et al.* (2019) ArrayExpress update—from bulk to single-cell expression data. *Nucleic Acids Res.*, **47**, D711–D715.
- Cerami,E., Gao,J., Dogrusoz,U. *et al.* (2012) The cBio cancer genomics portal: an open platform for exploring multidimensional cancer genomics data. *Cancer Discov.*, **2**, 401–404.
- Cancer Genome Atlas Research Network. (2008) Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature*, **455**, 1061–1068.
- Cheng,D.T., Mitchell,T.N., Zehir,A. *et al.* (2015) Memorial sloan kettering-integrated mutation profiling of actionable cancer targets (MSK-IMPACT): a hybridization capture-based next-generation sequencing clinical assay for solid tumor molecular oncology. *J. Mol. Diagn.*, **17**, 251–264.
- Gao,B., Huang,Q. and Baudis,M. (2018) segment\_liftover : a Python tool to convert segments between genome assemblies [version 1; referees: awaiting peer review]. *F1000Research*, **7**, 319.
- O'Keefe,C., McDevitt,M.A. and Maciejewski,J.P. (2010) Copy neutral loss of heterozygosity: a novel chromosomal lesion in myeloid malignancies. *Blood*, **115**, 2731–2739.
- Mullighan,C.G., Goorha,S., Radtke,I. *et al.* (2007) Genome-wide analysis of genetic alterations in acute lymphoblastic leukaemia. *Nature*, **446**, 758–764.
- Ai,N., Cai,H., Solovan,C. *et al.* (2016) CNARA: reliability assessment for genomic copy number profiles. *BMC Genomics.*, **17**, 799.
- Ortiz-Estevéz,M., Bengtsson,H. and Rubio,A. (2010) ACNE: a summarization method to estimate allele-specific copy numbers for Affymetrix SNP arrays. *Bioinformatics*, **26**, 1827–1833.
- Olshen,A.B., Venkatraman,E.S., Lucito,R. *et al.* (2004) Circular binary segmentation for the analysis of array-based DNA copy number data. *Biostatistics*, **5**, 557–572.
- Huang,Q. and Baudis,M. (2020) Enabling population assignment from cancer genomes with SNP2pop. *Sci. Rep.*, **10**, 1–9.
- World Health Organization and others (2013). *International Classification of Diseases for Oncology (ICD-O) 1st revision 3rd edition*.
- Cai,H., Gupta,S., Rath,P. *et al.* (2015) ArrayMap 2014: an updated cancer genome resource. *Nucleic Acids Res.*, **43**, D825–D830.



28. Ried,T., Hu,Y., Difilippantonio,M.J. *et al.* (2012) The consequences of chromosomal aneuploidy on the transcriptome of cancer cells. *Biochimica et Biophysica Acta (BBA)-Gene Regulatory Mechanisms.*, **1819**, 784–793.
29. da Silva Almeida,A.C., Abate,F., Khiabani,H. *et al.* (2015) The mutational landscape of cutaneous T cell lymphoma and sezary syndrome. *Nat. Genet.*, **47**, 1465–1470.
30. Xu,L., Gu,Z.-H., Li,Y. *et al.* (2014) Genomic landscape of CD34+ hematopoietic cells in myelodysplastic syndrome and gene mutation profiles as prognostic markers. *Proc. Natl. Acad. Sci.*, **111**, 8589–8594.
31. Baudis,M. (2007) Genomic imbalances in 5918 malignant epithelial tumors: an explorative meta-analysis of chromosomal CGH data. *BMC Cancer*, **7**, 226.
32. Cordo,P.C. and Baudis,M. (2021) Copy number variant heterogeneity among cancer types reflects inconsistent concordance with diagnostic classifications. *BioRxiv*.
33. Tarmin,L., Yin,J. and Zhou,X. *et al.* (1994) Frequent loss of heterozygosity on chromosome 9 in adenocarcinoma and squamous cell carcinoma of the esophagus. *Cancer Res*, **54**, 6094–6096.
34. Kuroki,T., Trapasso,F., Yendamuri,S. *et al.* (2003) Allele loss and promoter hypermethylation of VHL, RAR- $\beta$ , RASSF1A, and FHIT tumor suppressor genes on chromosome 3p in esophageal squamous cell carcinoma. *Cancer Res*, **63**, 3724–3728.
35. Armaghany,T., Wilson,J.D., Chu,Q. *et al.* (2012) Genetic alterations in colorectal cancer. *Gastrointestinal Cancer Research: GCR*, **5**, 19.
36. Nishida,N., Fukuda,Y., Kokuryu,H. *et al.* (1992) Accumulation of allelic loss on arms of chromosomes 13q, 16q and 17p in the advanced stages of human hepatocellular carcinoma. *nt. J. Cancer*, **51**, 862–868.
37. Mungall,C.J., Torniai,C., Gkoutos,G.V. *et al.* (2012) Uberon, an integrative multi-species anatomy ontology. *Genome Biol.*, **13**, R5.
38. Gene Ontology Consortium. (2008) The gene ontology project in 2008. *Nucleic Acids Res*, **36**, D440–D444.
39. Natale,D.A., Arighi,C.N., Barker,W.C. *et al.* (2007) Framework for a protein ontology. In:*BMC bioinformatics*, Springer, S1.
40. Mungall,C.J., McMurtry,J.A., Köhler,S. *et al.* (2017) The monarch initiative: an integrative data and analytic platform connecting phenotypes to genotypes across species. *Nucleic Acids Res.*, **45**, D712–D722.
41. Carrio-Cordo,P. and Baudis,M. (2018) Mountains and chasms: surveying the oncogenomic publication landscape. *Oncology*, **1**–12.
42. Carrio-Cordo,P., Acheson,E., Huang,Q. *et al.* (2020) Geographic assessment of cancer genome profiling studies. *Database*, **2020**.
43. Miki,Y., Swensen,J., Shattuck-Eidens,D. *et al.* (1994) A strong candidate for the breast and ovarian cancer susceptibility gene BRCA1. *Science*, **266**, 66–71.
44. Wooster,R., Bignell,G., Lancaster,J. *et al.* (1995) Identification of the breast cancer susceptibility gene BRCA2. *Nature*, **378**, 789–792.
45. Friend,S.H., Bernards,R., Rogelj,S. *et al.* (1986) A human DNA segment with properties of the gene that predisposes to retinoblastoma and osteosarcoma. *Nature*, **323**, 643–646.
46. Amundadottir,L.T., Sulem,P., Gudmundsson,J. *et al.* (2006) A common variant associated with prostate cancer in European and African populations. *Nat. Genet.*, **38**, 652–658.
47. Stacey,S.N., Manolescu,A., Sulem,P. *et al.* (2007) Common variants on chromosomes 2q35 and 16q12 confer susceptibility to estrogen receptor-positive breast cancer. *Nat. Genet.*, **39**, 865–869.
48. Tenesa,A., Farrington,S.M., Prendergast,J.G.D. *et al.* (2008) Genome-wide association scan identifies a colorectal cancer susceptibility locus on 11q23 and replicates risk loci at 8q24 and 18q21. *Nat. Genet.*, **40**, 631–637.
49. Wu,C., Hu,Z., Yu,D. *et al.* (2009) Genetic variants on chromosome 15q25 associated with lung cancer risk in Chinese populations. *Cancer Res.*, **69**, 5065–5072.
50. Jacobsen,J.O.B., Robinson,P.N. and Mungall,C.J. (2019) *Phenopackets Schema*.
51. Wagner,A.H., Babb,L., Alterovitz,G. *et al.* (2021) The GA4GH variation representation specification (VRS): a computational framework for the precise representation and federated identification of molecular variation. *BioRxiv*.
52. Lawler,M., Siu,L.L., Rehm,H.L. *et al.* (2015) Clinical working group of the global alliance for genomics and health, (GA4GH), all the World's a stage: facilitating discovery science and improved cancer care through the global alliance for genomics and health. *Cancer Discov.*, **5**, 1133–1136.
53. Mackay,A., Burford,A., Carvalho,D. *et al.* (2017) Integrated molecular meta-analysis of 1,000 pediatric high-grade and diffuse intrinsic pontine glioma. *Cancer Cell*, **32**, 520–e5.
54. Gao,B. and Baudis,M. (2021) Signatures of discriminative copy number aberrations in 31 cancer subtypes. *BioRxiv*, **12**, 732.
55. Bairoch,A. (2018) The cellosaurus, a cell-line knowledge resource. *J. Biomol. Tech.: JBT*, **29**, 25–38.