



**University of
Zurich**^{UZH}

**Zurich Open Repository and
Archive**

University of Zurich
University Library
Strickhofstrasse 39
CH-8057 Zurich
www.zora.uzh.ch

Year: 2022

The GA4GH Phenopacket schema defines a computable representation of clinical data

Jacobsen, Julius O B ; Baudis, Michael ; Baynam, Gareth S ; Beckmann, Jacques S ; Beltran, Sergi ; Buske, Orion J ; et al

DOI: <https://doi.org/10.1038/s41587-022-01357-4>

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-220996>

Journal Article

Accepted Version



The following work is licensed under a Creative Commons: Attribution 4.0 International (CC BY 4.0) License.

Originally published at:

Jacobsen, Julius O B ; Baudis, Michael ; Baynam, Gareth S ; Beckmann, Jacques S ; Beltran, Sergi ; Buske, Orion J ; et al (2022). The GA4GH Phenopacket schema defines a computable representation of clinical data. *Nature Biotechnology*, 40(6):817-820.

DOI: <https://doi.org/10.1038/s41587-022-01357-4>

The GA4GH Phenopacket schema defines a computable representation of clinical data

Julius O. B. Jacobsen^{1,*}, Michael Baudis^{2,3}, Gareth S. Baynam^{4,5,6}, Jacques S. Beckmann⁷, Sergi Beltran^{8,9,10}, Orion J. Buske¹¹, Tiffany J. Callahan¹², Christopher G. Chute¹³, Mélanie Courtot^{14,15}, Daniel Danis¹⁶, Olivier Elemento¹⁷, Andrea Essenwanger¹⁸, Robert R. Freimuth¹⁹, Michael A. Gargano¹⁶, Tudor Groza²⁰, Ada Hamosh²¹, Nomi L. Harris²², Rajaram Kaliyaperumal²³, Kevin C. Kent Lloyd^{24,25}, Aly Khalifa¹⁹, Peter M. Krawitz²⁶, Sebastian Köhler²⁷, Brian J. Laraway¹², Heikki Lehväslaiho²⁸, Leslie Matalonga⁸, Julie A. McMurry¹², Alejandro Metke-Jimenez²⁹, Christopher J. Mungall²², Monica C. Munoz-Torres¹², Soichi Ogishima³⁰, Anastasios Papakonstantinou⁸, Davide Piscia⁸, Nikolas Pontikos^{31,32}, Núria Queralt-Rosinach²³, Marco Roos²³, Julian Sass¹⁸, Paul N. Schofield^{33,34,35}, Dominik Seelow^{36,37}, Anastasios Siapos³⁸, Damian Smedley¹, Lindsay D. Smith^{15,39}, Robin Steinhaus^{36,37}, Jagadish Chandrabose Sundaramurthi¹⁶, Emilia M. Swietlik^{40,41,42}, Sylvia Thun¹⁸, Nicole A. Vasilevsky⁴³, Alex H. Wagner^{44,45}, Jeremy L. Warner⁴⁶, Claus Weiland⁴⁷, Melissa A. Haendel^{12,*} & Peter N. Robinson^{16,48,*}

¹Queen Mary University of London, William Harvey Research Institute, London EC1M 6BQ, UK

²University of Zurich, Department of Molecular Life Sciences, Zürich 8057, Switzerland

³Swiss Institute of Bioinformatics, Computational Oncogenomics Group, Zürich 8057, CH

⁴King Edward Memorial Hospital, Western Australian Register of Developmental Anomalies and Genetic Services of WA, Perth 6008, AU

⁵University of Western Australia, Faculty of Health and Medical Sciences, Division of Paediatrics, Perth 6008, AU

⁶Telethon Kids Institute, Genetic and Rare Diseases, Perth 6008, AU

⁷University of Lausanne, Faculty of Biology and Medicine, Lausanne CH-1015, Switzerland

⁸CNAG-CRG, Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Bioinformatics Unit, Barcelona 8028, ES

⁹Universitat Pompeu Fabra (UPF), Barcelona 8005, ES

¹⁰Departament de Genètica, Microbiologia i Estadística, Facultat de Biologia, Universitat de Barcelona (UB), Barcelona 8028, ES

¹¹PhenoTips, Toronto M5G 1L5, CA

¹²University of Colorado Anschutz Medical Campus, Center for Health AI, Aurora 80045, CO, USA

¹³Johns Hopkins University, Schools of Medicine, Public Health, and Nursing, Baltimore 21287, MD, USA

¹⁴European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Hinxton CB10 1SD, UK

¹⁵Ontario Institute for Cancer Research, Adaptive Oncology, Toronto M5G0A3, CA

- ¹⁶The Jackson Laboratory, Genomic Medicine, Farmington 6032, CT, USA
- ¹⁷Weill Cornell Medicine, Caryl and Israel Englander Institute for Precision Medicine, New York 10021, NY, USA
- ¹⁸Berlin Institute of Health at Charité - Universitätsmedizin Berlin, Core Facility Digital Medicine and Interoperability, Berlin 10178, DE
- ¹⁹Mayo Clinic, Department of Artificial Intelligence and Informatics, Rochester 55905, MN, USA
- ²⁰European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Cambridge CB10 1SD, UK
- ²¹Johns Hopkins University, Department of Genetic Medicine, Baltimore 21287, MD, USA
- ²²Lawrence Berkeley National Laboratory, Environmental Genomics and Systems Biology, Berkeley 94720, CA, USA
- ²³Leiden University Medical Center, Human Genetics, Leiden 2333 ZA, NL
- ²⁴UC Davis, Mouse Biology Program, Davis 95618, CA, USA
- ²⁵UC Davis School of Medicine, Department of Surgery, Sacramento 95817, CA, USA
- ²⁶University Hospital Bonn, Bonn, Germany, Institute for Genomic Statistics and Bioinformatics, Bonn 53113, DE
- ²⁷Ada Health GmbH, Berlin 10178, DE
- ²⁸CSC – IT Center for Science, Sensitive Data Services, Espoo FI-02101, FI
- ²⁹CSIRO, The Australian e-Health Research Centre, Herston 4029, AU
- ³⁰Tohoku University, INGEM, Sendai 980-8573, JP
- ³¹University College London, Institute of Ophthalmology, London EC1V 9EL, UK
- ³²Moorfields Eye Hospital, Genetics Service, London EC1V 2PD, UK
- ³³University of Cambridge, Dept of Physiology, Development and Neuroscience, Cambridge CB2 3EG, UK
- ³⁴The Jackson Laboratory, Mammalian Genetics, Bar Harbor ME 04609, ME, USA
- ³⁵The Alan Turing Institute, London NW1 2DB, UK
- ³⁶Bioinformatics and Translational Genetics, Berlin Institute of Health at Charité - Universitätsmedizin Berlin, Berlin 10178, DE
- ³⁷Charité - Universitätsmedizin Berlin, corporate member of Freie Universität Berlin and Humboldt-Universität zu Berlin, Institute of Medical Genetics and Human Genetics, Berlin 13353, DE
- ³⁸Lifebit Biotech Ltd., London W86BD, UK
- ³⁹Global Alliance for Genomics and Health, N/A, Toronto M5G0A3, CA
- ⁴⁰University of Cambridge, Medicine Department, Cambridge CB2 0QQ, UK
- ⁴¹Addenbrooke's Hospital, Respiratory Medicine Department, Cambridge CB2 0QQ, UK
- ⁴²Royal Papworth Hospital, Cambridge Centre for Lung Infection, Cambridge CB2 0AY, UK
- ⁴³University of Colorado, Anschutz Medical Campus, Center for Health AI, Aurora 80045, CO, USA
- ⁴⁴Nationwide Children's Hospital, The Steve and Cindy Rasmussen Institute for Genomic Medicine, Columbus 43215, OH, USA
- ⁴⁵The Ohio State University College of Medicine, Departments of Pediatrics and Biomedical Informatics, Columbus 43215, OH, USA

⁴⁶Vanderbilt University, Departments of Medicine and Biomedical Informatics, Nashville 37235, TN, USA

⁴⁷Senckenberg - Leibniz Institution for Biodiversity and Earth System Research, Data and Modelling Centre, Frankfurt/Main 60325, DE

⁴⁸University of Connecticut, Institute for Systems Genomics, Farmington 6032, CT, USA

Consortia authors

Myles Axton¹, Lawrence Babb², Cornelius F. Boerkoel³, Bimal P. Chaudhari^{4,5}, Hui-Lin Chin^{6,7}
Michel Dumontier⁸, David P. Hansen⁹, Harry Hochheiser¹⁰, Veronica A. Kinsler^{11,12}, Hanns
Lochmüller^{13,14,15} Alexander R. Mankovich¹⁶, Gary I. Saunders¹⁷, Panagiotis I. Sergouniotis¹⁸,
Rachel Thompson¹³ & Andreas Zankl^{19,20,21}

¹Wiley, Inc, Research, Hoboken 7030, NJ, USA

²Broad Institute of MIT and Harvard, Cambridge 2142, MA, USA

³University of British Columbia, Medical Genetics, Vancouver V6H3N1, CA

⁴Nationwide Children's Hospital, The Steve and Cindy Rasmussen Institute for Genomic Medicine, Divisions of Neonatology, Genetics and Genomic Medicine, Columbus 43215, OH, USA

⁵The Ohio State University College of Medicine, Department of Pediatrics, Columbus 43210, OH, USA

⁶Khoo Teck Puat-National University Children's Medical Institute, National University Hospital, Department of Paediatrics, Singapore 119074, Singapore

⁷Women's Hospital of British Columbia, Provincial Medical Genetics, Vancouver V6H3N1, CA

⁸Maastricht University, Institute of Data Science, Maastricht 6229 EN, NL

⁹CSIRO, Australian e-Health Research Centre, Brisbane 4027, AU

¹⁰University of Pittsburgh, Biomedical Informatics, Pittsburgh 15206, PA, USA

¹¹Great Ormond St Hospital for Children, Paediatric Dermatology, London WC1N 3JH, UK

¹²Francis Crick Institute, Mosaicism and Precision Medicine Laboratory, London NW1 1AT, UK

¹³Children's Hospital of Eastern Ontario Research Institute, Molecular Biomedicine, Ottawa K1H 8L1, CA

¹⁴University of Ottawa, Brain and Mind Research Institute, Department of Cellular and Molecular Medicine, Ottawa K1H 8M5, CA

¹⁵The Ottawa Hospital, Neuromuscular Centre, Ottawa K1Y 4E9, CA

¹⁶Philips Research North America, Precision Diagnosis & Image-Guided Therapy, Cambridge 2141, MA, USA

¹⁷ELIXIR, ELIXIR Hub, Cambridge CB10 1SD, UK

¹⁸University of Manchester, Division of Evolution, Infection and Genomics, Manchester M13 9PT, UK

¹⁹The University of Sydney, Faculty of Medicine and Health, Sydney 2006, AU

Correspondence should be addressed to J.J, M. & P.R: (e-mail j.jacobsen@qmul.ac.uk; melissa@tislabs.org; peter.robinson@jax.org)

To the editor. Despite great strides in the development and wide acceptance of standards for exchanging structured information about genomic variants, progress in standards for computational phenotype analysis for translational genomics has lagged behind. Phenotypic features (signs, symptoms, laboratory and imaging findings, results of physiological tests, etc.) are of high clinical importance, yet exchanging them in conjunction with genomic variation is often overlooked or even neglected. In the clinical domain, substantial work has been dedicated to the development of computational phenotypes.¹ Traditionally, these approaches have largely relied on rule-based methods and large sources of clinical data to identify cohorts of patients with or without a specific disease.²⁻⁵ However, they were not developed to enable deep phenotyping of abnormalities, to facilitate computational analysis of interpatient phenotypic similarity, or to support computational decision support. To address this, the Global Alliance for Genomics and Health⁶ (GA4GH) has developed the Phenopacket schema, which supports exchange of computable longitudinal case-level phenotypic information for diagnosis of and research on all types of disease, including Mendelian and complex genetic diseases, cancer, and infectious diseases (Fig 1). The Phenopacket software is available at <https://github.com/phenopackets/>.

The ‘PhenotypicFeature’ is the central element of the Phenopacket schema. A ‘PhenotypicFeature’ can be used to describe any phenotypic characteristic (often, but not necessarily, clinical abnormalities), including signs and symptoms, laboratory findings, histopathology findings, imaging, and electrophysiological results, along with modifier and qualifier concepts. Each phenotypic feature is described using an ontology term. Although the Phenopacket schema does not mandate which ontology to use, it provides recommendations,

such as the Human Phenotype Ontology⁷ (HPO) for rare diseases and the National Cancer Institute Thesaurus (NCIT) for transmission of information about a cancer specimen (e.g., pathological staging or more detailed information about histology or tumor markers).⁸ Within the schema, it is possible to indicate whether an abnormality was excluded during the diagnostic process (e.g., whether a morphological cardiac defect was excluded by echocardiography) or to use other optional HPO terms to denote the severity, frequency (e.g., number of occurrences of seizures per week), laterality (e.g., unilateral), or other pattern of a phenotypic feature in the patient being described. Finally, the onset (and if applicable the resolution) of specific features can be indicated.

Other key elements of the schema are ‘Measurement’, which is used to capture quantitative (i.e., numerical), ordinal (e.g., absent/present), or categorical measurements; ‘Biosample,’ a description of biological material obtained from the individual represented in the Phenopacket and used for phenotypic, genotypic, or other -omics analysis; and ‘MedicalAction,’ which includes a hierarchical representation of medical actions, including medications, procedures, and other actions taken for clinical management. The ‘Treatment’ element is a subelement of ‘MedicalAction’ and represents administration of a pharmaceutical agent, broadly defined as prescription and over-the-counter medicines, vaccines, and other therapeutic agents, such as monoclonal antibodies or chimeric antigen receptor (CAR)-T-cell-therapy.

The ‘Interpretation’ element specifies interpretations of genomic findings. This element leverages complementary resources developed by the GA4GH Genomic Knowledge Standards Work Stream: the Variation Representation Specification (VRS) and VRS Added Tools for Interoperable Loquacious Exchange (VRSATILE).⁶ Further information on this and other elements is available in the online documentation (<https://phenopacket-schema.readthedocs.io/>).

The Phenopacket schema was designed to support several use cases. Many of these use cases have been successfully implemented and tested in the community, particularly in the field of rare disease diagnostics and biobanking, whereas others, such as electronic health record integration, are in the process of being implemented (Supplementary Table 1).

The Phenopacket schema (version 2.0) was formally reviewed and approved as a GA4GH standard⁶ in 2021. It is designed to be interoperable with other relevant standards, including the traditional PED (pedigree format) file as well as the GA4GH pedigree standard, the GA4GH Beacon,⁹ and the GA4GH Variation Representation Specification. The GA4GH has committed to coordinate its activities and future roadmaps with those of other standards development organizations, including the International Organization for Standardization (ISO) Technical Subcommittee for Genomics Informatics (ISO/TC215/SC1) and HL7 Clinical Genomics (CG). Consequently, a Fast Interoperable Healthcare Resources (FHIR) implementation guide for Phenopacket interoperability is being developed and the Phenopacket schema is in the process of ISO certification (Supplementary Table 2).

The variant call format (VCF) standard for storing genotyping data allowed a wide range of research groups to write software for analyzing such data.¹⁰ The GA4GH Phenopacket schema aspires to be similarly transformative in the landscape of genome analysis using phenotype data. Multiple providers of phenotypic data include patients and clinicians, via a variety of mechanisms, including clinical notes and electronic health records, interfaces such as FHIR, app-based entry, and mobile devices. The Phenopacket schema acts as a common model that can capture data from many sources with a unified software representation and in turn can be used by multiple receivers of the phenotypic information, including journals, databases, registries, and clinical laboratories. Phenopackets can support diverse users and use cases, including patient-matchmaking services, diagnostics, and cohort identification. Software has become an essential resource for genomic medicine. We anticipate that the Phenopacket schema will encourage the development of a collection of software for the analysis of genomic data in the context of clinical information that will accelerate innovation and discovery. Genomic data will become ever more important in translational research and clinical care in the coming years and decades. The Phenopacket schema represents a standard for capturing clinical data and integrating it with genomic data that will help to obtain the maximal utility of this data for understanding disease and developing precision medicine approaches to therapy.

Acknowledgements

The authors gratefully acknowledge insight and feedback from Marian H. Adly, Pier Luigi Buttigieg, Nour Gazzaz, Janine Lewis, Manuel Posada de la Paz and Maria Taboada. This work was supported by 7RM1HG010860-02 (NHGRI). Additional funding was as follows. PNR was supported by NLM contract #75N97019P00280, NIH NHGRI RM1HG010860, NIH OD R24OD011883, NIH NICHD 1R01HD103805-01. HH was supported by NIH OD R24OD011883. GIS was supported by ELIXIR, the research infrastructure for life-science data. CGC was supported by NIH NCATS U24TR002306. KCL was supported by NIH OD 5UM1OD023221. MB was supported by BioMedIT Network project of Swiss Institute of Bioinformatics (SIB) and Swiss Personalized Health Network (SPHN). AHW was supported by NIH NHGRI K99HG010157, NIH NHGRI R00HG010157. CJM, MAH, MCM-T, JAM, DD were supported by NIH NHGRI RM1HG010860, NIH OD R24OD011883. AM-J was supported by Australian Genomics. Australian Genomics is supported by the National Health and Medical Research Council (GNT1113531). DS, JOBJ were supported by NIH NHGRI RM1HG010860, NIH OD R24OD011883, NIH NICHD 1R01HD103805-01. MD was supported by NIH NHGRI U54HG004028, NIH NHGRI 5U01HG008473-03, NIH NCATS OT2TR003434-01S1U54HG008033-01. GSB was supported by Roy Hill Community Foundation, Angela Wright Bennett Foundation, McCusker Charitable Foundation, Borlaug Foundation, Stan Perron Charitable Foundation. LB was supported by NIH NHGRI U41HG006834 (Clinical Genome Resource). MC was supported by EMBL-EBI Core Funds and Wellcome Trust GA4GH award number 201535/Z/16/Z. AH was supported by NIH NHGRI 1U41HG006627, NIH NHGRI 1U54HG006542, NIH NHGRI 1RM1HG010860. PNS was supported by The Alan Turing Trust. NLH was supported by NIH NHGRI RM1HG010860, NIH OD R24OD011883, U.S. Department of Energy Contract DE-AC02-05CH11231. NP was supported by Moorfields Eye Charity. NQ-R was supported by EU Horizon 2020 research and innovation programme grant agreement 825575 (EJP-RD). OE was supported by NIH grants UL1TR002384, R01CA194547, P01CA214274 LLS SCOR grants 180078-01, 7021-20, Starr Cancer Consortium Grant I11-0027. HL was supported by CIHR Foundation Grant on Precision Health for Neuromuscular Diseases FDN-167281. RT was supported by CIHR postdoctoral fellowship award MFE-171275. LDS was supported by Genome Canada and NIH NHGRI

U24HG011025. SO was supported by AMED. DP, LM, AP, SB, MR, RK were supported by EU Horizon 2020 research and innovation programme grant agreements 779257 (Solve-RD) and 825575 (EJP-RD). RRF was supported by NLM contract #75N97019P00280.

Competing interests

SK is an employee of Ada Health GmbH. NP is a director of Phenopolis Ltd. OE is supported by Janssen, Johnson and Johnson, Volastra Therapeutics, AstraZeneca and Eli Lilly research grants. He is scientific advisor and equity holder in Freenome, Owkin, Volastra Therapeutics and One Three Biotech. ARM is an employee of Philips Research North America. OJB is an employee of PhenoTips. MA is an editor employed by Wiley. AS is an employee of Lifebit Biotech Ltd.

1. Richesson, R. & Smerek, M. Electronic health records-based phenotyping. *Rethinking clinical trials: A living textbook of pragmatic clinical trials* **2016**, (2014).
2. Hripcsak, G. & Albers, D. J. Next-generation phenotyping of electronic health records. *J. Am. Med. Inform. Assoc.* **20**, 117–121 (2013).
3. Shivade, C. *et al.* A review of approaches to identifying patient phenotype cohorts using electronic health records. *J. Am. Med. Inform. Assoc.* **21**, 221–230 (2014).
4. Wei, W.-Q. & Denny, J. C. Extracting research-quality phenotypes from electronic health records to support precision medicine. *Genome Med.* **7**, 41 (2015).
5. Richesson, R. L., Sun, J., Pathak, J., Kho, A. N. & Denny, J. C. Clinical phenotyping in selected national networks: demonstrating the need for high-throughput, portable, and computational methods. *Artif. Intell. Med.* **71**, 57–61 (2016).
6. Rehm, H. L. *et al.* GA4GH: International policies and standards for data sharing across genomic research and healthcare. *Cell Genom* **1**, (2021).
7. Köhler, S. *et al.* The Human Phenotype Ontology in 2021. *Nucleic Acids Res.* **49**, D1207–D1217 (2021).
8. Sioutos, N. *et al.* NCI Thesaurus: a semantic model integrating cancer-related clinical and molecular information. *J. Biomed. Inform.* **40**, 30–43 (2007).
9. Fiume, M. *et al.* Federated discovery and sharing of genomic data using Beacons. *Nat. Biotechnol.* **37**, 220–224 (2019).
10. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).

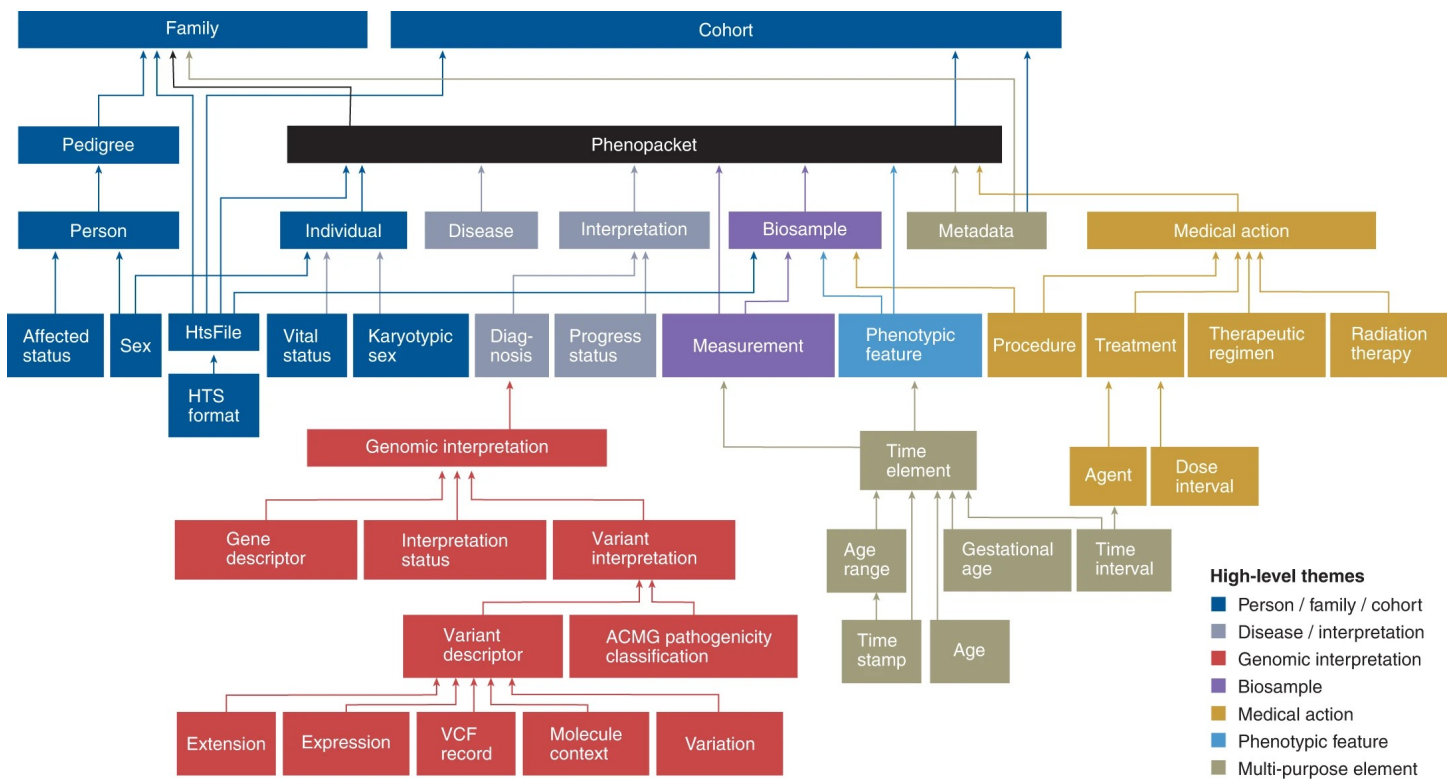


Figure 1. Phenopacket schema overview. The GA4GH Phenopacket schema consists of several optional elements, each of which contains information about a certain topic, such as phenotype, variant or pedigree. An element can contain other elements, which allows a hierarchical representation of data. For instance, Phenopacket contains elements of type Individual, PhenotypicFeature, Biosample, and so on. Individual elements can therefore be regarded as building blocks that are combined to create larger structures. Colors represent the major themes of elements within the schema.