

Advancing the Clinical Usefulness of Next-Generation Sequencing: Addressing the Need for Genomic Data Standards

Ira M. Lubin, PhD, FACMG and Edward R. Lockhart, PhD on behalf of the Clinical Grade Variant File Workgroup*
Centers for Disease Control and Prevention, Atlanta, GA

Members of the Clinical Grade Variant File Workgroup: Nazreen Anzi¹, Larry Babbs², Dennis Ballinger³, Himani Bhatt⁴, Deanna M. Church⁵, Shaun Conroy⁶, Alden A. Dima⁷, Karan Ellouk⁸, Timothy Fennell⁹, J. Bradley Holmes¹⁰, Fiona Hyland¹¹, Lisa Kilman¹², Melissa Landrum¹³, Edward R. Lockhart¹⁴, Ira M. Lubin¹⁵, Donna Maglott¹⁶, Elliott Marquis¹⁷, Gibor Matth¹⁸, Anu Nakaaj¹⁹, John Pfeiffer²⁰, Mala Ramanah²¹, Heidi Rehm²², Somak Roy²³, Marc L. Sallit²⁴, Chris Saunders²⁵, Stephen Sherry²⁶, Zviana Tozak²⁷, Rebecca Trudy²⁸, Mollie Milman-Culler²⁹, Karl Woelberding³⁰, Elizabeth Worshay³¹, Alexander Watt Zaranek³², and Justin Zook³³.
¹Phoenix Children's Hospital, Phoenix, Arizona; ²Laboratory for Molecular Medicine, Partners Healthcare Personalized Medicine, Boston, MA; ³Compass Genomics, Mountain View, California; ⁴Center for Devices and Radiological Health, Food and Drug Administration, Silver Spring, MD; ⁵Pharmacia, Menlo Park, California; ⁶Genomics Division, National Institute of Standards and Technology, Gaithersburg, MD; ⁷Department of Human Genetics, University of Utah School of Medicine, Salt Lake City, UT; ⁸Novartis Institute, Boston, MA; ⁹National Center for Biotechnology Information, National Institutes of Health, Bethesda, MD; ¹⁰Life Technologies, Carlsbad, CA; ¹¹Division of Laboratory Programs, Standards, and Services, Centers for Disease Control and Prevention, Atlanta, GA; ¹²Genome Cambridge, Bedford, UK; ¹³Department of Pathology and Immunology, Washington University School of Medicine, St. Louis, MO; ¹⁴Department of Pathology, Brigham and Women's Hospital and Harvard Medical School, Cambridge, MA; ¹⁵Department of Pathology, University of Pittsburgh Medical Center, Pittsburgh, PA; ¹⁶illumina, San Diego, CA; ¹⁷Clinical and Translational Informatics, Dana-Farber Cancer Institute, Boston, MA; ¹⁸Department of Pathology, University of Utah and the Institute for Clinical and Experimental Pathology, Associated Regional and University Pathologists (ARUP) Laboratories, Salt Lake City, Utah; ¹⁹Department of Pediatrics, Medical College of Wisconsin; ²⁰Personal Genome Project, Harvard Medical School, Boston, MA

Introduction

Systems interoperability that supports genomic applications does not exist because of the absence of a framework for data exchange. To address this challenge, the Centers for Disease Control and Prevention, working with federal partners (National Center for Biotechnology Information, National Institute of Standards and Technology, and the Food and Drug Administration) the HL7 Clinical Genomics Workgroup, and others, established the Clinical Grade Variant File Workgroup to address this shortcoming. This workgroup is proposing a "clinical-grade" variant template that is designed to describe the variant set generated during NGS testing prior to clinical assessment and the clinically relevant findings to support clinical applications. The workgroup identified principles and made recommendations for a constrained dataset that promotes consistency and minimizes ambiguity in the representation of genomic data.

The Clinical Grade Variant Template

Section 1 – General Information	Section 2 – Clinically Relevant Findings	Section 3 – Sequence Dataset (prior to clinical assessment)
<ol style="list-style-type: none"> 1. Patient information^a 2. Indication for Test^b 3. Test Name^c 4. Specimen^d 5. Method^e 6. Referring physician and facility 7. Laboratory performing test 8. Date of Initial Test 9. Date of Initial Report 10. Summary of re-analysis^f 	<ol style="list-style-type: none"> 1. Gene/Variant/Haplotype/Amino acid change(s)^g 2. Classification with supporting references^h 3. Quality metricsⁱ 4. Limitations^j 5. Re-analyzed data^k 	<p>(The Variant Call Format (VCF) is the recommended specification for storing sequence data)</p> <pre>##fileformat=VCFv4.1 ##1=chr1:100000000-100000000 ##2=chr1:100000000-100000000 ##3=chr1:100000000-100000000 ##4=chr1:100000000-100000000 ##5=chr1:100000000-100000000 ##6=chr1:100000000-100000000 ##7=chr1:100000000-100000000 ##8=chr1:100000000-100000000 ##9=chr1:100000000-100000000 ##10=chr1:100000000-100000000 ##11=chr1:100000000-100000000 ##12=chr1:100000000-100000000 ##13=chr1:100000000-100000000 ##14=chr1:100000000-100000000 ##15=chr1:100000000-100000000 ##16=chr1:100000000-100000000 ##17=chr1:100000000-100000000 ##18=chr1:100000000-100000000 ##19=chr1:100000000-100000000 ##20=chr1:100000000-100000000 ##21=chr1:100000000-100000000 ##22=chr1:100000000-100000000 ##23=chr1:100000000-100000000 ##24=chr1:100000000-100000000 ##25=chr1:100000000-100000000 ##26=chr1:100000000-100000000 ##27=chr1:100000000-100000000 ##28=chr1:100000000-100000000 ##29=chr1:100000000-100000000 ##30=chr1:100000000-100000000 ##31=chr1:100000000-100000000 ##32=chr1:100000000-100000000 ##33=chr1:100000000-100000000 ##34=chr1:100000000-100000000 ##35=chr1:100000000-100000000 ##36=chr1:100000000-100000000 ##37=chr1:100000000-100000000 ##38=chr1:100000000-100000000 ##39=chr1:100000000-100000000 ##40=chr1:100000000-100000000 ##41=chr1:100000000-100000000 ##42=chr1:100000000-100000000 ##43=chr1:100000000-100000000 ##44=chr1:100000000-100000000 ##45=chr1:100000000-100000000 ##46=chr1:100000000-100000000 ##47=chr1:100000000-100000000 ##48=chr1:100000000-100000000 ##49=chr1:100000000-100000000 ##50=chr1:100000000-100000000 ##51=chr1:100000000-100000000 ##52=chr1:100000000-100000000 ##53=chr1:100000000-100000000 ##54=chr1:100000000-100000000 ##55=chr1:100000000-100000000 ##56=chr1:100000000-100000000 ##57=chr1:100000000-100000000 ##58=chr1:100000000-100000000 ##59=chr1:100000000-100000000 ##60=chr1:100000000-100000000 ##61=chr1:100000000-100000000 ##62=chr1:100000000-100000000 ##63=chr1:100000000-100000000 ##64=chr1:100000000-100000000 ##65=chr1:100000000-100000000 ##66=chr1:100000000-100000000 ##67=chr1:100000000-100000000 ##68=chr1:100000000-100000000 ##69=chr1:100000000-100000000 ##70=chr1:100000000-100000000 ##71=chr1:100000000-100000000 ##72=chr1:100000000-100000000 ##73=chr1:100000000-100000000 ##74=chr1:100000000-100000000 ##75=chr1:100000000-100000000 ##76=chr1:100000000-100000000 ##77=chr1:100000000-100000000 ##78=chr1:100000000-100000000 ##79=chr1:100000000-100000000 ##80=chr1:100000000-100000000 ##81=chr1:100000000-100000000 ##82=chr1:100000000-100000000 ##83=chr1:100000000-100000000 ##84=chr1:100000000-100000000 ##85=chr1:100000000-100000000 ##86=chr1:100000000-100000000 ##87=chr1:100000000-100000000 ##88=chr1:100000000-100000000 ##89=chr1:100000000-100000000 ##90=chr1:100000000-100000000 ##91=chr1:100000000-100000000 ##92=chr1:100000000-100000000 ##93=chr1:100000000-100000000 ##94=chr1:100000000-100000000 ##95=chr1:100000000-100000000 ##96=chr1:100000000-100000000 ##97=chr1:100000000-100000000 ##98=chr1:100000000-100000000 ##99=chr1:100000000-100000000 ##100=chr1:100000000-100000000</pre>

Figure 1. This template is proposed to support systems interoperability by defining and constraining a set of data fields described in three sections:

Section 1 – General Information	Section 2 – Clinically Relevant Findings	Section 3 – Sequence Dataset (prior to clinical assessment)
<ul style="list-style-type: none"> • Patient information (name, date of birth, sex) • Indication for testing (with associated ICD code) • Test name (Standard, common, LOINC) • Specimen (Type /site of origin, germline / somatic) • Method (Platform, Software, etc., references) • Summary of Re-analysis, when applicable 	<ul style="list-style-type: none"> • Variant/gene and/or haplotype, amino acid changes: (HGNC/HGVS descriptions + common names) • Pathogenicity classification and supporting data • Quality metrics (Primarily of clinical relevance) • Limitations (e.g. evidence of data) • when new knowledge or different indication of testing presents 	<ul style="list-style-type: none"> • Chromosome number or reference to an alternate assembly derived from a GRC reference assembly, when applicable • Position of nucleotide within a chromosome mapped against a GRC versioned reference assembly (or other standard) • Identifiers of the represented variant (e.g., dbSNP) • Reference base of the genome assembly • Alternate base; non-reference allele • Quality score for assertion made in the ALT field • Filter_status; lists whether the site/position has passed filters • Additional genomic information (i.e. genotype, base call quality, read depth, etc.)

Data Standards are Applied During Testing

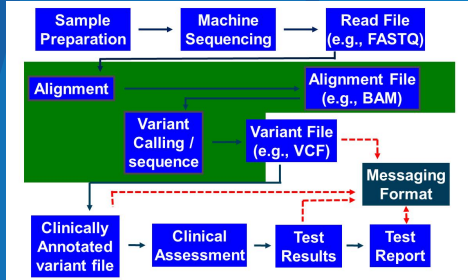


Figure 2. Sequence alignment and variant calling are steps during the informatics analysis that require attention to assure consistent position assignments and descriptions of sequences generated during NGS analysis. The same applies to clinical annotation and assessment but these topics were not addressed by this workgroup.

Clinical Applications Supported by Genomic Data

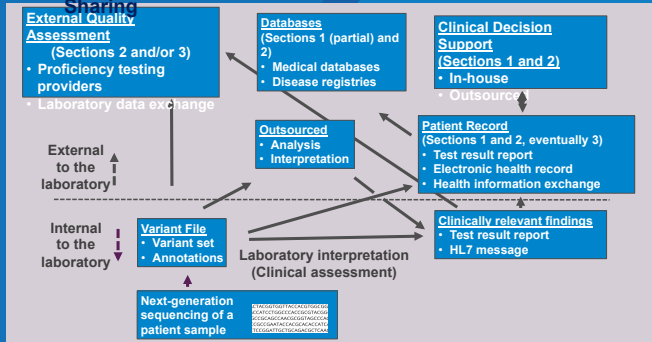


Figure 3. Clinical applications requiring the sharing of genomic data. Sections 1, 2, and/or 3 of the clinical-grade variant template can be combined and shared as applicable to specific clinical applications. This selected combination is proposed to inform the content of an HL7 message.

Use Case Example: Where Information Sharing is

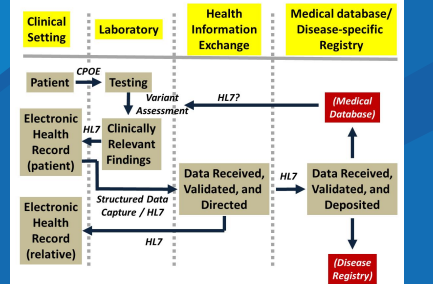


Figure 4. Use case example. The clinical-grade variant template is proposed to inform the development of the HL7 message. Structured data capture protocols are useful for normalizing data among varied practice settings. This example provides a model for communicating genomic data to a relative's medical setting at a distant location from the index patient, deposition of genomic data into a medical database or disease registry, and extraction of genomic data from a medical database to inform a clinical assessment.

Summary

1. Standards for genomic data representation are required to minimize ambiguity in the description of sequencing findings generated from clinical next-generation sequencing.
2. The clinical-grade variant template identifies a constrained genomic dataset that describes the sequence generated from NGS before and after analysis for clinical relevance.
3. The application of standard conventions requires that certain laboratory methods be adopted, such as alignment against a versioned reference assembly.
4. The clinical-grade variant template was developed to inform approaches for genomic data integration into health IT systems using standard messaging conventions (e.g., HL7). These outcomes will support systems interoperability that is essential for accurate and reliable use of genomic data in diverse healthcare settings.

For additional information, please contact:
Ira Lubin, PhD, FACMG at ilubin@cdc.gov

www.cdc.gov | Contact CDC at: 1-800-CDC-INFO or www.cdc.gov/info

The findings and conclusions in this report are those of the authors and do not necessarily represent the official position of the Centers for Disease Control and Prevention.

Center for Surveillance, Epidemiology, and Laboratory Services
Division of Laboratory Programs, Standards, and Services

