

# Genomic Data Harmonization through the OMOP Standardized Vocabularies

Denys Kaduk<sup>2</sup>, Violetta Komar<sup>2</sup>, Asieh Golozar<sup>6</sup>, Peter Robinson<sup>9</sup>, Alex H. Wagner<sup>10</sup>, Michael Gurley<sup>4</sup>, Seng Chan You<sup>3</sup>, Shin Seojeong<sup>3</sup>, Shaadi Mehr<sup>8</sup>, Andrew Williams<sup>4</sup>, Rimma Belenkar<sup>7</sup>, Meera Patel<sup>7</sup>, Shilpa Ratwani<sup>1</sup>, Ron Miller<sup>1</sup> and Christian Reich<sup>1</sup>,

<sup>1</sup>IQVIA, Plymouth Meeting, PA, USA, <sup>2</sup>Odysseus Data Services Inc., Cambridge, MA, USA, <sup>3</sup>AJOU University, Suwon-si, Gyeonggi-do, South Korea, <sup>4</sup>TUFTS, Medford, MA, USA, <sup>5</sup>Northwestern University, Evanston, IL, USA, <sup>6</sup>Regeneron Pharmaceuticals, Inc, Tarrytown, New York USA, <sup>7</sup>Memorial Sloan Kettering Cancer Center, New York, NY, USA, <sup>8</sup>MMRF Norwalk, Connecticut, USA, <sup>9</sup>The Jackson Laboratory for Genomic Medicine, Hartford, Connecticut Area, USA, <sup>10</sup>Washington University School of Medicine, St. Louis, MO, USA

## Abstract

In a typical observational study, cohorts, exposures, and outcomes can be sufficiently defined through the presence or absence of clinical events encoded by a defined set of concepts. Precision Oncology however is more challenging than that of most other conditions. Accurately integrating EHRs with genomic data for the discovery of clinically actionable variants can generate new insights into disease mechanisms and provide better tools for identification of effective treatments. Currently, the OMOP Standardized Vocabularies lacks a canonical, comprehensive and non-redundant representation of genomic variants. We developed a preliminary simple heuristic for canonization and evaluated six prominent somatic cancer variant knowledgebases to serve as a source of clinically relevant variants to include. We found these to have a low level of overlap and recommend adopting a solution that combines these repositories and incorporates a deduplicated set for the purpose of Standard Concepts.

## Introduction

Clinical research in precision oncology requires concise, standardized and searchable interpretations of detected variants. Many institutions have created collections of clinically relevant variants and formalized interpretations such as disease implication and targeting by pharmaceutical compounds through curation from the biomedical literature. These efforts have resulted in disparate knowledge representation of the genomic information. Currently, there is no comprehensive genomic terminology available in the public domain that would support harmonization of genomic oncology data, which is necessary for standardized analytics in a research network. For OHDSI, representing the various knowledgebases in a standardized manner requires an open, interoperable sharing of variant interpretation and an automated methodology so a comprehensive approach to cancer precision medicine can be developed.

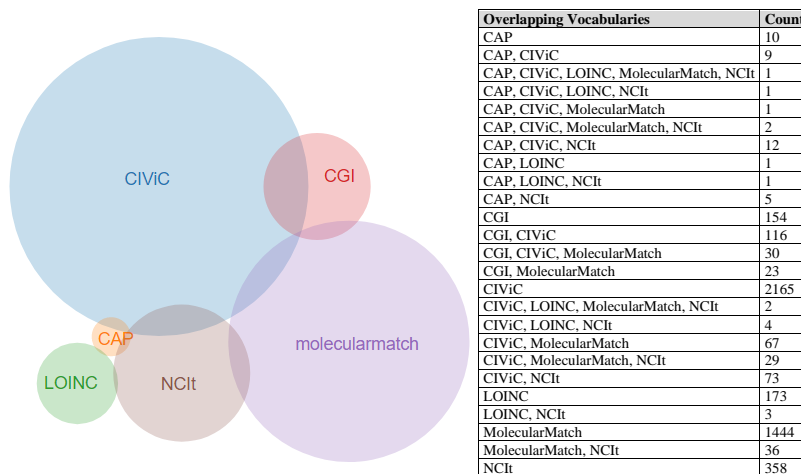
To illustrate the challenges of harmonizing across multiple variant representations, we looked at all interpretations describing the variant AKT1 E17K, a change in the amino acid residue at position 17 in the RAC-alpha serine/threonine-protein kinase protein where glutamic acid has been replaced by lysine, which is implicated in lung and breast cancer. This variant is present in four of the six databases (see below), but each knowledgebase represents this variant differently. The differences are due to representations at the genomic, complementary DNA (cDNA) (which could be in antisense direction) and protein levels, each of which has different reference sequences in different versions from different databases (NCBI, EMBL, LRG):

Dimension	Number of representations
Genomic reference sequences	4
Versions of genomic reference sequences	8
cDNA reference sequences	3
Versions of genomic reference sequences	3
Protein reference sequences	17
Versions of protein reference sequences	17

To address these challenges and to develop interoperability between genomic data and clinical care, a collaboration between the OHDSI Oncology Workgroup and the VICC consortium (1) is seeking to develop a canonical representation of equivalent variants to serve as standard concepts in the OMOP Standardized Vocabularies.

We developed a simple consolidation approach of duplications, insertions, deletions and duplications on the basis of gene symbol, sequence type (g, c, p), reference sequences, versions and locations. As a starting point, we incorporated six such source vocabularies: National Cancer Institute Thesaurus (NCIt), College of American Pathologists Cancer Checklists (CAP), Clinical Interpretation of Variants in Cancer (CIViC), Cancer Genome Interpreter (CGI), MolecularMatch and Logical Observation Identifiers Names and Codes (LOINC) (figure 1).

**Figure 1.** Consolidation of six genomic databases of variants relevant in cancer.



### Conclusion

The limited overlap between cancer variant databases suggests a union of these variants and a heuristic to deduplicate the overlap. The resulting collection can be instantiated as Standard Concepts and evaluated for utility in real-world genomic data research.

### References

1. Wagner, A.H., Walsh, B., Mayfield, G. et al. A harmonized meta-knowledgebase of clinical interpretations of somatic genomic variants in cancer. *Nat Genet* 52, 448–457 (2020)