

In the format provided by the authors and unedited.

A harmonized meta-knowledgebase of clinical interpretations of somatic genomic variants in cancer

Alex H. Wagner¹, Brian Walsh², Georgia Mayfield², David Tamborero^{3,4}, Dmitriy Sonkin⁵, Kilannin Krysiak¹, Jordi Deu-Pons^{6,7}, Ryan P. Duren⁸, Jianjiong Gao⁹, Julie McMurry², Sara Patterson¹⁰, Catherine del Vecchio Fitz¹¹, Beth A. Pitel¹², Ozman U. Sezerman¹³, Kyle Ellrott², Jeremy L. Warner¹⁴, Damian T. Rieke¹⁵, Tero Aittokallio^{16,17}, Ethan Cerami¹¹, Deborah I. Ritter^{18,19}, Lynn M. Schriml²⁰, Robert R. Freimuth¹², Melissa Haendel^{2,21}, Gordana Raca^{22,23}, Subha Madhavan²⁴, Michael Baudis²⁵, Jacques S. Beckmann²⁶, Rodrigo Dienstmann²⁷, Debyani Chakravarty⁹, Xuan Shirley Li⁸, Susan Mockus¹⁰, Olivier Elemento²⁸, Nikolaus Schultz⁹, Nuria Lopez-Bigas^{3,6,7}, Mark Lawler²⁹, Jeremy Goecks², Malachi Griffith¹✉, Obi L. Griffith¹✉, Adam A. Margolin² and Variant Interpretation for Cancer Consortium*

¹Washington University School of Medicine, St. Louis, MO, USA. ²Oregon Health and Science University, Portland, OR, USA. ³Pompeu Fabra University, Barcelona, Spain. ⁴Karolinska Institute, Solna, Sweden. ⁵National Cancer Institute, Rockville, MD, USA. ⁶Institute for Research in Biomedicine, Barcelona, Spain. ⁷Catalan Institution for Research and Advanced Studies, Barcelona, Spain. ⁸MolecularMatch, Houston, TX, USA. ⁹Memorial Sloan Kettering Cancer Center, New York, NY, USA. ¹⁰The Jackson Laboratory for Genomic Medicine, Farmington, CT, USA. ¹¹Dana-Farber Cancer Institute, Boston, MA, USA. ¹²Mayo Clinic, Rochester, MN, USA. ¹³Acibadem University, Istanbul, Turkey. ¹⁴Vanderbilt University, Nashville, TN, USA. ¹⁵Charité—Berlin University of Medicine, Berlin, Germany. ¹⁶Institute for Molecular Medicine Finland, Helsinki, Finland. ¹⁷University of Turku, Turku, Finland. ¹⁸Baylor College of Medicine, Houston, TX, USA. ¹⁹Texas Children's Hospital, Houston, TX, USA. ²⁰University of Maryland School of Medicine, Baltimore, MD, USA. ²¹Linus Pauling Institute at Oregon State University, Corvallis, OR, USA. ²²Children's Hospital Los Angeles, Los Angeles, CA, USA. ²³Keck School of Medicine of USC, Los Angeles, CA, USA. ²⁴Georgetown University Medical Center, Washington, DC, USA. ²⁵University of Zurich, Zurich, Switzerland. ²⁶University of Lausanne, Lausanne, Switzerland. ²⁷Vall d'Hebron Institute of Oncology, Barcelona, Spain. ²⁸Weill Cornell Medicine, New York, NY, USA. ²⁹Queen's University Belfast, Belfast, UK. *A list of members and affiliations appears in the Supplementary Note. ✉e-mail: mgriffit@wustl.edu; obigriffith@wustl.edu

Supplementary Note

Differences in somatic and germline variant interpretation

When considered in discovery and translational research endeavors, it is important to determine if a particular variant observed in a gene of interest is *oncogenic* (the variant functionally enables or predisposes towards the development of cancer), as this annotation provides the foundation on which targeted cancer treatment research is based. In contrast, clinical applications are dominated by diagnostic, prognostic, or therapeutic interpretations which in part also depends on underlying variant oncogenicity.

The development of procedures for oncogenicity classification in somatic interpretation guidelines is an unmet need in the oncology domain. ClinGen, ACMG, AMP, ASCO, VICC and CAP are collaboratively developing such guidelines to enable consistent and comprehensive assessment of somatic variant oncogenicity.

Consensus and recommendations for the elements of an Interpretation

To provide readily searchable, standardized interpretations across knowledgebases, we evaluated the structure of cancer variant interpretations across the core dataset (**Figure 1**). Our first challenge was to develop a consensus for the minimum required data elements that constitute a cancer variant interpretation. These minimal elements include a gene identifier, variant name, cancer subtype (tumor type and organ), clinical implication (diagnostic, prognostic, therapeutic, or predisposing biomarker), provenance of supporting evidence (e.g., PubMed identifier), and curation source. In addition, we recommended ascribing a tiered level of support for the evidence contributing to the interpretation. Each VICC knowledgebase (**Supplementary Table 1**) provided cancer variant interpretations as structured data meeting these requirements.

Difference in curation strategy

The differing curation and data modeling strategies of these knowledgebases serve as contributing factors for their dramatically different variant content. CGI uses invited expert curators to build out its knowledgebase of therapeutic biomarkers. Similarly, JAX-CKB, MolecularMatch, OncoKB, and PMKB create interpretations from in-house expert panels, but the provenance of these interpretations are each credited to the entire panel instead of an individual curator. The CIViC knowledgebase provides an interface for anyone to curate cancer variant interpretations from the literature, and the curated content is then reviewed by expert editors. The need for only 1 (CGI) or 2 (CIViC) curators to generate interpretations may contribute to the greater breadth of variants in these resources compared to others. PMKB has far fewer variants than the others as a result of the way interpretations are modeled; many of the interpretations in PMKB apply to broad variant representations (described above as categorical variants) which may be used for multiple interpretations. In fact, there are more interpretations in PMKB than in OncoKB, but OncoKB has more specific variants associated with their interpretations. Additional differences exist between the focus of these knowledgebases that are not accounted for in our analysis. For instance, JAX-CKB and MolecularMatch each also aggregate clinical trial data, and CGI and OncoKB both maintain large “oncogenic” annotation lists, but none of these additional data are included in the analyses presented in this manuscript.

Grouping of disease terms to top-level disease concepts

The aggregated knowledge across the core dataset describes 357 distinct disease concepts from the Disease Ontology (DO)⁴⁷ across 12,497 interpretations (**Supplementary Table 4**). These diseases range from highly specific (e.g. *DOID:0080164 - myeloid and lymphoid neoplasms with eosinophilia and abnormalities of PDGFRA, PDGFRB, and FGFR1*) to generalized (e.g. *DOID:162 - cancer*). To compare the variant interpretations to disease type, we used the expert-curated “TopNodeCancerSlim” DO mapping⁵⁰ that describes 58 common, top-level disease terms (TopNode terms) across several major datasets, including The Cancer Genome Atlas (TCGA), International Cancer Genome Consortium (ICGC), and COSMIC.^{3,51,52}

Mapping GENIE diseases to Disease Ontology

GENIE samples are annotated with a diverse array of Oncotree ontology (oncotree.mskcc.org) disease codes, with 81% (539 / 667) of Oncotree diseases represented in the dataset. Over 53% (286 / 539) of the Oncotree diseases from GENIE do not link to DO through cross-references, of which 34% (96 / 286) do not have any cross-references (**Supplementary Table 9**). This lack of cross-references among GENIE diseases is significantly higher than the 20% (133 / 667) of all Oncotree terms lacking cross-references (OR = 2.0, $p = 1.1e-5$; Fisher's exact test, two-sided), suggesting that terms used to describe individual patient cancers (e.g. *Well-Differentiated Neuroendocrine Tumor of the Rectum*) are less likely to map to other knowledgebases than high-level parent terms (e.g. *colorectal cancer*). Despite this, 65% of GENIE patients had a disease term map to DO, indicating that the common cancers among this cohort are more likely to be cross-referenced adequately for mapping. Further evaluation confirmed a significant enrichment of more frequently observed disease terms among the terms that mapped to DO, compared to those that did not ($U = 31094.0$, $p = 4.8e-3$; Mann-Whitney U test, two-sided; see **Online Methods**).

In addition to the above limitations, matching patient variants to diseases with our strategy is highly dependent upon ontology structure. For example, CGI exhibits a 54% reduction of exact matches due to the mapping of CGI terms to DO. To illustrate this problem, patients with DOID:3008, *invasive ductal carcinoma*, account for 11% of the unmatched diseases (or 6% of the overall reduction). 82% of these patients have a variant matching a CGI interpretation for DOID:3458, *breast adenocarcinoma*. This is a sibling term to *breast lobular carcinoma* and *breast ductal carcinoma*. As a result, GENIE *invasive ductal carcinoma* patients (*invasive ductal carcinoma* is the sole descendant of *breast ductal carcinoma*) do not match, as our match strategy requires ancestral relationships between concepts (see **Extended Data Figure 4** for details).

Improvements from harmonization

We observed large gains in overlapping terms between resources across variants, diseases, and drugs. Importantly, this harmonization allowed us to search relationships between patient and interpretation disease terms, improving precise matching between patient and interpretations of clinical significance. We noted that there was little need to harmonize gene identifiers, as each knowledgebase had independently selected HGNC gene symbols as a reference, enabling easy and direct comparison of genes. This underscores both the utility of standardized data, as well as the need for adoption of similar standards (such as those described in this work) to drive direct comparison of variant interpretations. In our analysis of the variants and diseases of the harmonized interpretations, we observed that frequent top-level cancer terms mirror cancers with high incidence and mortality. We also noted that a large percentage of these interpretations described a relatively small number of gene-disease relationships.

Future goals

Additionally, we are building inference tools to automatically identify the concepts users are querying in real time. We also will be expanding our effort to harmonize and present interpretations of various non-coding variants, structural variants beyond gene pairs, and aggregate markers like microsatellite instability status. A prioritized long-term goal is the development of standards and techniques for interpretations of combined germline and somatic variations. Similarly, we are building guidelines and methods to enable automated consensus recommendations. Finally, we are seeking out additional knowledgebases of clinical interpretations of variants to harmonize and share with the broader cancer genomics community, and building an API specification which they may use to incorporate their own interpretations.

50. Wu, T.-J. et al. Generating a focused view of disease ontology cancer terms for pan-cancer data integration and analysis. *Database* **2015**, bav032 (2015).

51. Cancer Genome Atlas Research Network. Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature* **455**, 1061–1068 (2008).

52. International Cancer Genome Consortium et al. International network of cancer genome projects. *Nature* **464**, 993–998 (2010).

Participating Knowledgebase Agreement

Dear VICC leaders,

I would like to take this opportunity to express my commitment to participating in the Variant Interpretation for Cancer Consortium (VICC). As you know, we are one of several institutions engaged in the challenge of curating knowledge to annotate cancer genome mutations associated with evidence of pathogenicity or linked to relevant treatment options. Specifically, we have developed the [INSTITUTION RESOURCE (INSTITUTION NAME)]. While this resource serves the specific needs of our own institution, we recognize that there is **clear value in sharing knowledge of cancer-variant-treatment associations**. Such sharing will increase confidence where interpretations overlap, fill gaps, reduce redundancy, and leverage disparate domain expertise. To this end, we support your plans to coordinate global efforts for curation and help develop a community resource for cross-knowledgebase queries under the auspices of the Global Alliance for Genomics and Health (GA4GH). As a VICC participant we agree with the following data sharing principles, developed through community discussion at GA4GH meetings and calls.

- We will commit to sharing at least a minimal set of data elements for cancer variant interpretations including: gene symbol, variant name, cancer subtype (tumor type and organ), clinical implication (drug sensitivity, drug resistance, adverse response, diagnostic, or prognostic), source (e.g., PubMed identifier) and curation group.
- We agree that to avoid patient data privacy concerns, the project will focus on only clinical interpretations of variants derived from published findings (literature, conference proceedings, and clinical trial records), not individual patient/variant-level observations. Thus, there should be no possibility of linking variants to individuals.
- We agree to share [SPECIFY: all OR a significant proportion] of our interpretations (with at least the minimal required data elements) accumulated by our ongoing curation efforts.
 - This content will be released under a permissive license (free and non-exclusive for at least research use).
 - Software developed as part of this data sharing initiative will be released in public repositories (e.g., github) with open source licenses.
 - Public APIs will be developed to facilitate access to our data for use by the VICC.
 - Wherever possible, data sharing will be facilitated by use of the existing GA4GH schemas, APIs and demonstration implementations.
 - Interpretations made available by our institution will also be made available as cross-knowledgebase bulk downloads.

I look forward to working together with the Variant Interpretation for Cancer Consortium to further our common goal of improving genomics-guided precision medicine for cancer patients.

Sincerely,

[INSTITUTION REPRESENTATIVE]

VICC Members

Tero Aittokallio^{1,2,3,4}; Lawrence Babb⁵; Michael Baudis⁶; Jacques S Beckmann⁷; Anas Belouali⁸; Andrew Biankin⁹; Misha Bouzinier¹⁰; Steven E Brenner¹¹; Alberto Cambrosio¹²; Jonah Campbell¹²; Ethan Cerami¹³; Debyani Chakravarty¹⁴; David Chang⁹; Brad Chapman¹⁵; Thomas Conway¹⁶; Christopher L Corless¹⁷; Melanie Courtot¹⁸; Robert Currie¹⁹; Catherine del Vecchio Fitz¹³; Jordi Deu-Pons^{20,21}; Rodrigo Dienstmann²²; Kenneth Doig¹⁶; Ryan P Duren²³; Daniel Durkin²⁴; Korneel Duyvesteyn²⁵; Olivier Elemento²⁶; Jonathan Ellis²⁷; Kyle Elrott¹⁷; Kenneth W Eng²⁶; Aidan Flynn^{28,29}; Robert R Freimuth³⁰; JianJiong Gao¹⁴; Martina Gasull²¹; Moritz Gerstung¹⁸; William B Glen³¹; Jeremy Goecks¹⁷; Santiago Gonzalez²¹; Sara Gosline³²; Malachi Griffith³³; Obi L Griffith³³; Melissa Haendel^{17,34}; Maximilian Haeussler¹⁹; David Heckerman³⁵; Oliver Hofmann³⁶; Peter Horak³⁷; Sarah Hunt¹⁸; Ivan Jelas³⁸; Vaidehi Jobanputra^{39,40}; Rachel Karchin⁴¹; Ian King⁴²; Liang K Koh⁴³; Ana Krepischi⁴⁴; Kilannin Krysiak³³; Mario Lamping³⁸; Melissa Landrum⁴⁵; Mark Lawler⁴⁶; Jennifer Lee⁴⁷; Jonas Leichsenring⁴⁸; Michele L Lenoue-Newton⁴⁹; Paul Leo²⁷; Huei S Leong¹⁶; Xuan S Li²³; Xuelu Liu¹³; Nuria Lopez-Bigas^{20,21,50}; Chris Love¹⁶; Ravi Madduri⁵¹; Subha Madhavan⁸; Sameer Malhotra²⁶; Adam Margolin¹⁷; David L Masica⁴¹; Georgia Mayfield¹⁷; Matthew D McCoy⁸; Clay McLeod⁵²; Christine Micheel⁴⁹; Susan Mockus²⁴; Victoria Muenzer³⁸; Christopher J Mungall⁵³; Rishi Nag¹⁸; Kevin Osborn¹⁹; Ravi Pandya⁵⁴; Nicole Park⁵⁵; Sara Patterson²⁴; Michael Piechotta⁵⁶; Beth A Pitel³⁰; Gordana Raca^{57,58}; Erin Ramos⁵⁹; Shruti Rao⁸; Gunnar Ratsch^{60,61}; Iker Reyes²¹; Damian T Rieke^{38,62}; Deborah I Ritter^{63,64}; Peter Rogan⁶⁵; Jeffrey Rosenfeld⁶⁶; Sameek Roychowdhury⁶⁷; Gabe Rudy⁶⁸; Gina Rueter³⁸; Chris Sander^{5,13}; Andrea Sboner²⁶; Lynn M Schriml⁶⁹; Nikolaus Schultz¹⁴; Ozman U Sezerman⁷⁰; Sandra Siesing⁴⁸; Lillian Siu⁵⁵; Heidi J Sofia⁵⁹; Dmitriy Sonkin⁷¹; Vipin Sreedharan⁶¹; Albrecht Stenzinger⁴⁸; Andrew I Su⁷²; David Tamborero^{50,73}; Bin T Teh⁴³; Nora C Toussaint⁶¹; Eliezer Van Allen^{5,13,74,75}; Nicole A Vasilevsky¹⁷; Etienne Vignola-Gagne¹²; Ioannis Vlachos⁷⁴; Andra Waagmeester⁷⁶; Alex H Wagner³³; Brian Walsh¹⁷; Jeremy L Warner⁴⁹; Joachim Weischenfeldt²⁸; Trish Whetzel¹⁸; Julia Wilson⁷⁷; Chunlei Wu⁷²; Andrew Yates¹⁸; Andrey Zaparyi⁷⁸

VICC Member Affiliations

¹Institute for Molecular Medicine Finland, Helsinki, Finland; ²Institute for Cancer Research, Oslo, Norway; ³University of Oslo, Oslo, Norway; ⁴University of Turku, Turku, Finland; ⁵Broad Institute of MIT and Harvard, Cambridge, MA, USA; ⁶University of Zurich, Zurich, Switzerland; ⁷University of Lausanne, Lausanne, Switzerland; ⁸Georgetown University Medical Center, Washington D.C., USA; ⁹University of Glasgow, Glasgow, Scotland; ¹⁰InterSystems Corporation, Cambridge, MA, USA; ¹¹UC Berkeley, Berkeley, CA, USA; ¹²McGill University, Montreal, Québec, Canada; ¹³Dana-Farber Cancer Institute, Boston, MA, USA; ¹⁴Memorial Sloan Kettering Cancer Center, New York, NY, USA; ¹⁵Harvard T.H. Chan School of Public Health, Boston, MA, USA; ¹⁶Peter MacCallum Cancer Centre, Melbourne, Australia; ¹⁷Oregon Health and Science University, Portland, OR, USA; ¹⁸European Molecular Biology Laboratory - European Bioinformatics Institute, Cambridge, UK; ¹⁹UC Santa Cruz, Santa Cruz, CA, USA; ²⁰Catalan Institution for Research and Advanced Studies, Barcelona, Spain; ²¹Institute for Research in Biomedicine (IRB Barcelona), The Barcelona Institute of Science and Technology, Barcelona, Spain; ²²Vall d'Hebron Institute of Oncology, Barcelona, Spain; ²³MolecularMatch, Houston, TX, USA; ²⁴The Jackson Laboratory for Genomic Medicine, Farmington, CT, USA; ²⁵Hartwig Medical Foundation, Amsterdam, Netherlands; ²⁶Weill Cornell Medicine, New York, NY, USA; ²⁷Queensland University of Technology, Woolloongabba, Australia; ²⁸University of Copenhagen, Copenhagen, Denmark; ²⁹Rigshospitalet, Copenhagen, Denmark; ³⁰Mayo Clinic, Rochester, MN, USA; ³¹Medical University of South Carolina, Charleston, SC, USA; ³²Sage Bionetworks, Seattle, WA, USA; ³³Washington University School of Medicine, St. Louis, MO, USA; ³⁴Linus Pauling Institute at Oregon State University, Corvallis, OR, USA; ³⁵Human Longevity Inc, San Diego, CA, USA; ³⁶University of Melbourne, Melbourne, Australia; ³⁷National Center for Tumor Diseases Heidelberg, Heidelberg, Germany; ³⁸Charité – Berlin University of Medicine, Berlin, Germany; ³⁹New York Genome Center, New York, NY, USA; ⁴⁰Columbia University Irving Medical Center, New York, NY, USA; ⁴¹Johns Hopkins University, Baltimore, MD, USA; ⁴²University of Toronto, Toronto, Ontario, Canada; ⁴³National Cancer Centre Singapore, Singapore, Singapore; ⁴⁴National Institute of Science and

Technology in Oncogenomics, São Paulo, Brazil; ⁴⁵National Center for Biotechnology Information, Bethesda, MD, USA; ⁴⁶Queen's University Belfast, Belfast, UK; ⁴⁷Frederick National Laboratory for Cancer Research, Rockville, MD, USA; ⁴⁸Heidelberg University Hospital, Heidelberg, Germany; ⁴⁹Vanderbilt University, Nashville, TN, USA; ⁵⁰Pompeu Fabra University, Barcelona, Spain; ⁵¹University of Chicago, Chicago, IL, USA; ⁵²St. Jude Children's Research Hospital, Memphis, TN, USA; ⁵³Lawrence Berkeley National Laboratory, Berkeley, CA, USA; ⁵⁴Microsoft, Redmond, WA, USA; ⁵⁵University Health Network, Toronto, Ontario, Canada; ⁵⁶Humboldt University of Berlin, Berlin, Germany; ⁵⁷Children's Hospital Los Angeles, Los Angeles, CA, USA; ⁵⁸Keck School of Medicine of USC, Los Angeles, CA, USA; ⁵⁹National Human Genome Research Institute, Bethesda, MD, USA; ⁶⁰Swiss Institute of Bioinformatics, Lausanne, Switzerland; ⁶¹Swiss Federal Institute of Technology in Zurich, Zurich, Switzerland; ⁶²Berlin Institute of Health, Berlin, Germany; ⁶³Baylor College of Medicine, Houston, TX, USA; ⁶⁴Texas Children's Hospital, Houston, TX, USA; ⁶⁵University of Western Ontario, London, Ontario, Canada; ⁶⁶Rutgers University, New Brunswick, NJ, USA; ⁶⁷Ohio State University, Columbus, OH, USA; ⁶⁸Golden Helix, Bozeman, MT, USA; ⁶⁹University of Maryland School of Medicine, Baltimore, MD, USA; ⁷⁰Acibadem University, Istanbul, Turkey; ⁷¹National Cancer Institute, Rockville, MD, USA; ⁷²Scripps Research, La Jolla, CA, USA; ⁷³Karolinska Institute, Stockholm, Sweden; ⁷⁴Harvard Medical School, Boston, MA, USA; ⁷⁵Brigham and Women's Hospital, Boston, MA, USA; ⁷⁶Micelio, Antwerp, Belgium; ⁷⁷Wellcome Sanger Institute, Hinxton, UK; ⁷⁸Dell EMC, Hopkinton, MA, USA

Supplementary Equations

Calculations for evaluating non-harmonized aggregate content

Non-harmonized element counts:

Given the set of interpretation elements:

$$E = \{genes, variants, diseases, drugs\}$$

and the set of knowledgebase resources:

$$R = \{cgi, civic, jax, molecularmatch, oncokb, pmkb\}$$

For each element $e \in E$, and each resource $r \in R$, we created a unique set of *non-harmonized* element values observed in that resource, $S_{e,r}$.

For **genes** ($S_{genes, r \in R}$), we used HGNC gene symbols, which were provided by each knowledgebase.

Gene symbols were almost universally provided across interpretations, although some interpretations do not have associated genes.

For **variants** ($S_{variants, r \in R}$), we extracted the genomic coordinates (chromosome, start, stop) from each resource and created a unique set of those variants. JAX-CKB and OncoKB do not provide genomic coordinates for variants. When applicable, we split records by the appropriate delimiter to separate out multiple variants. For CGI, we also did minimal HGVS parsing for chr/start/stop when gDNA HGVS strings were provided.

For **diseases** ($S_{diseases, r \in R}$), we extracted the disease term from each knowledgebase and transformed it to lowercase text. PMKB represents diseases as a combination of tissue and tumor type, which we transformed to a compound string joined by a space (e.g., *Tissue: Breast* and *Type: Adenocarcinoma* became *Disease: breast adenocarcinoma*).

For **drugs** ($S_{drugs, r \in R}$), we extracted the drug term from each knowledgebase and transformed it to lowercase text. As many interpretations contain more than one drug, we identified the delimiting character for each resource where multiple drugs are represented as a single string and split the string on the delimiter (e.g., the single string “*dabrafenib + trametinib*” was treated as the two strings “*dabrafenib*” and “*trametinib*”).

We did not perform this analysis for **evidence** levels, as there is no shared meaning behind unharmonized evidence levels across resources (**Table 1**).

The size of the set of unique values $\left| S_{e \in E, r \in R} \right|$ was recorded in **Supplementary Table 3**. For example, $\left| S_{genes, cgi} \right|$ (cell B3 of **Supplementary Table 3**) represents the 183 unique gene symbols observed in the Cancer Genome Interpreter.

Harmonized element counts:

For each element $e \in E$, and each resource $r \in R$, we created a unique set of *harmonized* element values observed in that resource, $S'_{e, r}$. These values were determined via the element harmonization routines specified in the previous **Online methods** sections.

The sizes of the sets of unique values $\left| S'_{e \in E, r \in R} \right|$ were recorded in **Supplementary Table 3**. For example, $\left| S'_{genes, cgi} \right|$ (cell C3 of **Supplementary Table 3**) represents the 182 unique gene symbols derived from the Cancer Genome Interpreter after harmonization.

Summary of gains from harmonization:

To measure the degree to which harmonization improved the consistency of observed terms across knowledgebases, we first calculated the total *Terms Evaluated*. For each element e , this consisted of calculating both the total *Terms Evaluated* from both non-harmonized (t_e) and harmonized (t'_e) valuesets:

$$t_e = \sum_{r \in R} \left| S_{e, r} \right| \text{ and } t'_e = \sum_{r \in R} \left| S'_{e, r} \right|$$

For each element e we also recorded the total *Unique Terms* across the union of terms from all resources, by measuring the cardinality of both the non-harmonized (u_e) and harmonized (u'_e) valuesets:

$$u_e = \left| \bigcup_{r \in R} S_{e, r} \right| \text{ and } u'_e = \left| \bigcup_{r \in R} S'_{e, r} \right|$$

Finally, for each element e we calculate the *Overlap* as the percentage reduction of the total Terms Evaluated to the total Unique Terms, for both the non-harmonized (o_e) and harmonized (o'_e) valuesets:

$$o_e = 1 - \frac{u_e}{t_e} \text{ and } o'_e = 1 - \frac{u'_e}{t'_e}$$