# SISO: A Conceptual Model-based Method for Variant Interpretation Systematization

Mireia **Costa**[1], Alberto **García S.**[1], Ana **Leon**[1] and Oscar **Pastor**[1]

[1]*PROS Research Group, VRAIN Research Institute – Universitat Politècnica de València, Spain*

### Abstract

Variant interpretation is the process by which clinical experts determine if a DNA variant has a significant impact on a patient's health. Current practices in variant interpretation suffer from a lack of traceability and reproducibility due to the chaotic nature of genomic data and the imprecision of existing variant interpretation guidelines. These issues pose substantial barriers to the routine clinical application of variant interpretation. This paper introduces SISO, a conceptual model-based method designed to translate the inherent imprecision of variant interpretation into a concise and well-defined set of steps. The practical utility of the SISO method is demonstrated through a use case involving a variant identified in a patient with suspected familial breast-ovarian cancer syndrome. The SISO method lays the foundations for variant interpretation systematization by guiding decision-making and ensuring reproducibility. As a result, variant interpretation will be a more reliable and consistent process in clinical practice.

### Keywords

Variant Interpretation, Conceptual Modeling, SISO Method

## 1. Introduction

Our DNA sequence is regarded as one of our most distinguishing characteristics as individuals. We humans share more than 99% our DNA sequence. However, these small differences in our DNA contribute to natural human diversity and influence susceptibility to certain diseases or variations in response to typical treatments. These differences are called DNA variants [1]. Given their significance to human health, a key objective of medicine is to understand how DNA variants impact an individual's health, a process known as variant interpretation.

Despite its importance, variant interpretation suffers from several issues that have yet to be resolved. The interpretation process involves weighing data about variants, such as the variant's frequency among the population, whether it has previously been linked to a disease, etc. This data is scattered across thousands of data sources with unique and contradictory content as well as different terminology [2]. This data chaos makes the data used for interpretation difficult to trace and potentially causes the same data to be interpreted differently by different experts

[3]. Additionally, variant interpretation typically follows specialized guidelines that provide a series of recommendations that guide the interpretation by determining whether or not a variant meets specific criteria. However, despite their intentions, these guidelines are often criticized for providing vague definitions that result in subjectivity in their interpretation and inconsistency in their application [4, 5, 6].

These factors lead to a lack of clarity regarding both the criteria each expert uses for interpretation and the evidence they rely on. As a result, variant interpretation suffers from a complete lack of traceability and reproducibility, posing a significant barrier to its application in routine clinical practice [7].

Transforming variant interpretation from an abstract process into a well-defined, repeatable, and reliable procedure is a challenge that requires both breaking down the process into its fundamental elements (i.e., unpacking) and organizing these elements into a coherent, efficient, and standardized framework (i.e., systematization). Unpacking is vital for disentangling the intricate details of variant interpretation by clarifying aspects with implicit or ambiguous definitions and defining a common framework for representation. Conversely, systematization is critical to making the variant interpretation process explicit, guiding decision-making, and ensuring reproducibility.

Previous works have explored the unpacking of variant interpretation through a meta-model called VarClaMM [8], which represents all relevant elements involved in the interpretation process, as detailed in Section 2. Building upon this foundation, the objective of this work is to take the first steps towards variant interpretation systematization. We present SISO, a method conceptually grounded in VarClaMM, that aims to translate the inherent imprecision of the interpretation process into a concise and well-defined set of steps. The primary aim of this research is to present a novel approach that applies the principles of conceptual modeling to the domain of variant interpretation, a non-traditional application area. The proposed SISO method offers a robust framework for addressing the complexities inherent in variant interpretation. Preliminary results from the application of this method in a real use case indicate its potential to enable more standardized, reproducible, and reliable interpretations in clinical practice.

The remaining is organized as follows: Section 2 provides a summary of the VarClaMM model. Section 3 introduces the SISO method. Section 4 demonstrates the utility of the method through a practical use case. Finally, Section 5 concludes the paper.

## 2. Background: The VarClaMM Meta-model

The VarClaMM meta-model (Figure 1) represents the most important elements of the variant interpretation process: i) the variant itself (depicted in green), ii) the constructs that conform to the variant interpretation guidelines (depicted in pink), iii) the results of evaluating a guideline over a variant to obtain its interpretation (depicted in orange), and iv) the variant related data required to perform the interpretation (depicted in blue).

Below, we provide a high-level description of VarClaMM, focusing on the most important elements for understanding the logic behind the SISO method. A more in depth description can be found in [8].
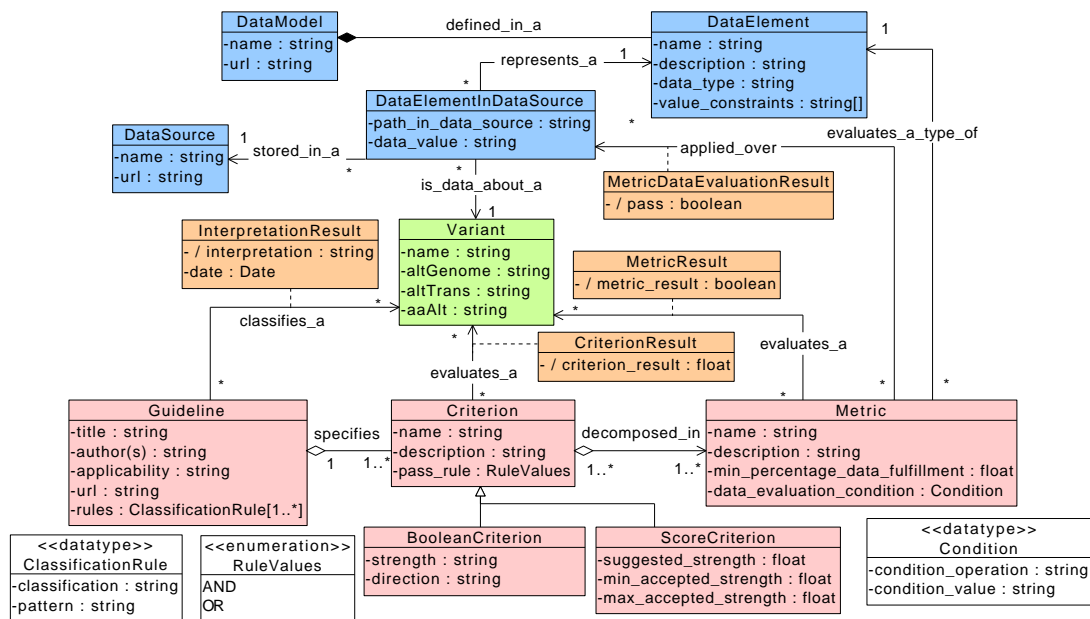
**Figure 1:** VarClaMM meta-model.

## Guideline

Each variant interpretation guideline is represented in the GUIDELINE class, which is characterized by its *applicability*, which refers to the specific disease, gene, or variant type context in which the guideline is applicable. An example is the ACMG-AMP 2015 guidelines [9], which only apply to mendelian diseases.

## Criteria and Metrics

Each GUIDELINE evaluates different CRITERION to obtain the classification of a variant. For instance the ACMG-AMP 2015 guidelines defines 28 criteria, one of which is the PVS1 CRITERION. This criterion evaluates if *the variant is null and if it is in a gene where null variants are to cause disease.* Guidelines consider criteria of one of two types: BOOLEANCRITERION, which evaluate to true or false, and SCORECRITERION, whose evaluation returns a numerical value.

Both types of criteria define specific conditions whose fulfillment determines whether the criterion is met. In this model, we call these conditions METRICS. Consider the example of the PVS1 criterion mentioned above. In its definition, this criterion establishes two well-differentiated, independent conditions: i) the variant must be null, and ii) the gene affected must cause disease through null variants. Our model represents these two conditions as metrics associated with the PVS1 criterion.

**Variant related data**

To determine whether the condition established by a Metric is fulfilled, we need to evaluate data about the variant under study. Each kind of data to be evaluated (e.g., an allele frequency) is represented in the DataElement class. VarClaMM forces the different DataElements to be structured according to a DataModel in order to have shared and standard definition of the data used.

The value of a DataElement for a given variant comes from external DataSources. This concrete value is represented in the DataElementInDataSource class. For example, the concept of allele frequency (e.g. a DataElement) takes the value 1.647e-05 in ExAC DataSource and 1.36e-06 in the GnomAD DataSource for the variant FANCI:c.669+1G>T.

**Evaluation results**

The Guideline, Criterion, and Metric constructs are evaluated over a specific Variant to obtain its interpretation. The interpretation of a variant based on a particular Guideline is represented in the InterpretationResult class. This interpretation is calculated by applying a set of *rules* over each CriterionResult. Each CriterionResult is calculated by applying the *pass_rule* (which takes the value *AND* when all the metrics must be fulfilled and *OR* otherwise) over the results of each criterion's Metric.

Finally, each MetricResult is calculated by evaluating all the DataElementInDataSource for the variant corresponding to the DataElement evaluated by the Metric. The model also represents the results of evaluating each individual DataElementInDataSource in the MetricDataEvaluationResult class.

## 3. SISO Method

The SISO method builds on the unpacking achieved with VarClaMM and establishes the foundation for systematizing variant interpretation by transforming it into a series of concrete and well-defined steps. Our method guides the use of the VarClaMM meta-model to define an interpretation framework that can be consistently applied to any given variant. Consequently, SISO supports users in applying variant interpretation guidelines, ensuring that accurate and reproducible results are obtained. This is accomplished through four stages, which are detailed below.

**1. (S)elect an interpretation guideline.**

In this step, the expert must select the most suitable guideline for the interpretation. Here it is crucial to consider the context in which the interpretation is being performed, as interpretation guidelines have specific applicabilities (see Section 2). There are three key aspects to consider when selecting a guideline:

- Variant type: Certain guidelines are tailored to specific types of variants. For instance, the ACMG-ClinGen 2019 guidelines address *copy number variants*, while the ACMG 2020 guidelines focus on variants located in *mitochondrial DNA*.
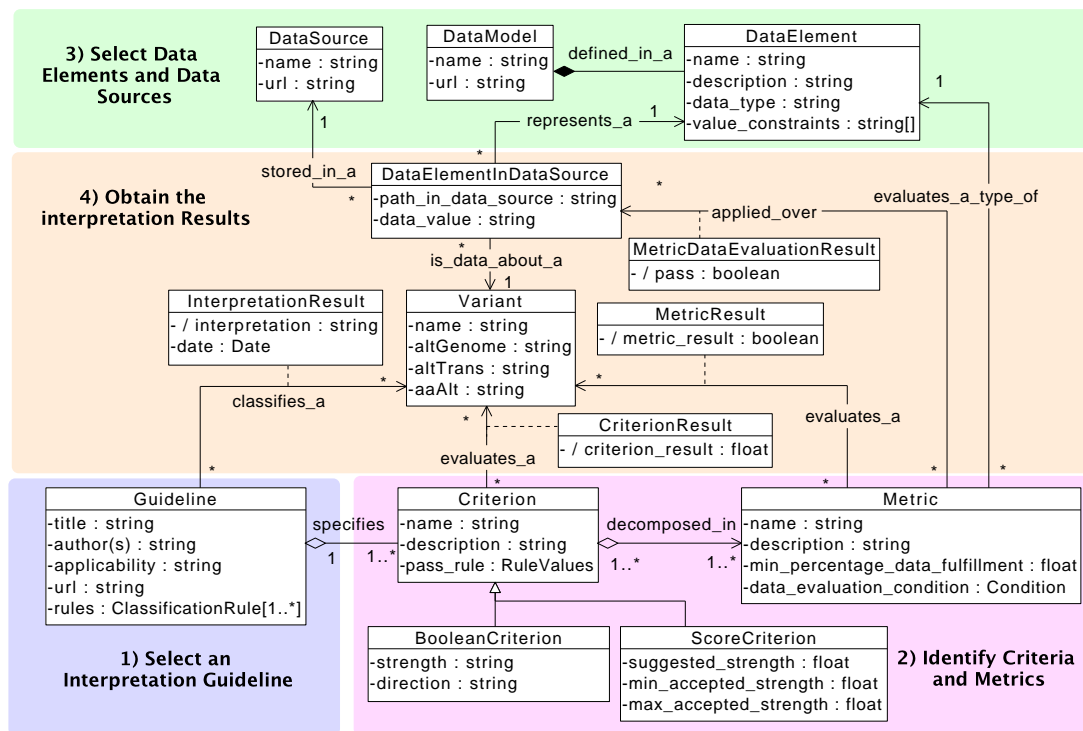
**Figure 2:** Schematic representation of the SISO method over the components of the VarClaMM meta-model.

- Gene location: Some guidelines consider unique aspects of variants within specific genes. Examples include the ClinGen Guidelines for the FBN1 gene and the CanVIG guidelines for the CHEK2 gene.
- Disease relevance: Certain diseases have distinct characteristics that must be taken into account when determining if a variant is pathogenic. Examples of such guidelines are the ACGS guidelines for rare diseases, the Ambry Genetics guidelines for autosomal dominant and X-linked diseases, and the ACMG-AMP 2015 guidelines for Mendelian diseases.

It could be the case that no existing guideline covers the specific needs of the interpretation context. In that case, the method supports the creation of a new guideline. In both cases, the GUIDELINE must be completely characterized following the VarClaMM meta-model. The selection of the guideline will completely condition the results of the interpretation and, consequently, is vital for ensuring the quality and reliability of the interpretation results.

## 2. (I)dentify the criteria and define the metrics of the guideline.

Once the interpretation guideline has been selected, the next step involves identifying the criteria and the metrics needed for its evaluation. This would be straightforward if the guidelines provided precise definitions. However, as stated in Section 1, this is often not the case.

For instance, consider the PVS1 criterion of the ACMG-AMP 2015 guidelines mentioned above. For this criterion, the metrics are clear: i) the variant must be null, and ii) the gene affected must cause disease through null variants. Now, let's examine the BS1 criterion of the same guidelines, which assesses whether *the variant's allele frequency is greater than expected for that specific disease*. In this case, the definition makes the frequency cut-off entirely dependent on the expert performing the interpretation [5]. As a result, one expert might consider a cut-off of 0.5%, while another, stricter expert might set a cut-off of 1%. Consequently, even though both experts claim to apply the same criterion, they are not truly evaluating the same thing.

On the other hand, there are cases where certain criteria from the selected guideline cannot be applied due to a lack of resources. For instance, the PS3 criterion of the ACMG-AMP 2015 guidelines evaluates whether there are *in vitro or in vivo functional studies supporting damaging effects*. However, only 36% of clinical experts are able to obtain this kind of evidence, as functional studies require significant monetary and time investments [10]. Consequently, claiming that a certain guideline is used not always implies that all its criteria are applied.

This highlights that merely stating the guideline used is insufficient to fully understand how a variant has been interpreted. To properly define the interpretation process, it is fundamental to specify which criteria of the guideline are actually applied, how each criterion is evaluated through the definition of metrics, and how these metrics are combined to obtain a criterion's result. This can be achieved thanks to the composition relationships between GUIDELINE and CRITERION and CRITERION and METRIC defined in the VarClaMM meta-model.

## 3. (S)elect the required data elements and data sources.

Once the criteria and their respective metrics are defined, the next step is to identify the specific data (e.g., the DATAELEMENT) that each identified metric needs to evaluate. For example, within the PVS1 criterion, the metric *the variant must be null* evaluates the consequence of the variant, whereas the metric *the gene affected must cause disease through null variants* evaluates the gene's disease mechanism. In the context of the BS1 criterion, the selected metric will evaluate the variant's frequency. According to VarClaMM, the representation of these DATAELEMENTs is dependent on the chosen DATAMODEL.

With the DATAELEMENTs identified, the next step is to determine which data sources provide the required data. The primary challenge here is that this knowledge is scattered across thousands of data sources. In a recent study, Costa et al. [11] performed a comparative analysis of several genomic data sources and concluded that none of them provide complete data about variants and that, in some cases, the data they provide is not concordant.

For example, consider the PM2 criterion of the ACMG-AMP 2015 guidelines, which states that *the variant must be absent from control populations*. In practice, this criterion is evaluated by a metric that checks whether the variant has an allele frequency of 0 in population databases. Two of the most important population databases are ExAC [12] and GnomAD [13]. For the specific case of MITF:c.1A>G, the variant is absent in ExAC, thus meeting the PM2 criterion. However, in the GnomAD database, this variant has a frequency of 7e-07, meaning the PM2 criterion would not be met. This example illustrates the discrepancies that can arise when using different DATASOURCEs, underscoring the importance of carefully selecting and cross-referencing genomic data sources to ensure traceability of the interpretation results, as even

with the same criteria and metrics the results may differ depending on the DataSources used.

**4. (O)btain the classification results.**

The initial three steps of the SISO method have enabled us to develop a comprehensive framework for variant interpretation. This framework thoroughly details the guidelines, criteria, and metrics to be evaluated, the necessary data for evaluation, and the sources from which this data can be obtained. Such detailed documentation will significantly enhance the traceability and reproducibility of the interpretation process.

The final step is to "execute" this framework to interpret a given variant. To do this, we first need to obtain the specific value of each DataElement for the variant in question (e.g., the DataElementInDataSource). Next, we calculate all the MetricDataEvaluationResults. From these results, we can determine the outcomes of each metric and criterion, ultimately leading to the final ClassificationResult of the analyzed variant.

## 4. Use case

Here we present a practical application of the SISO method to illustrate its utility in real-world scenarios. Specifically, we employed the SISO method to guide the interpretation of the c.191G>A variant in the BRCA1 gene. This variant has been identified in a Chilean patient suspected of having a familial breast-ovarian cancer syndrome. Below, a description of how each stage has been performed is provided.

**1. Select an interpretation guideline.**

As detailed in Section 3, selecting the optimal interpretation guideline requires considering the variant type, the gene, and the disease under study. The variant under study is of type deletion, which means it removes a portion of the gene where is located, in this case, the BRCA1 gene. To our knowledge, there is no interpretation guidelines focusing exclusively on deletion variants. However, the ClinGen institution has developed a guideline that focus on BRCA1 variants associated with familial breast-ovarian cancer syndrome [14]. Consequently, we selected this guideline as the most appropriate for our use case.

**2. Identify the criteria and define the metrics of the guideline.**

The selected guideline comprises 40 BooleanCriterion that must be evaluated to obtain the InterpretationResult of the variant. Detailing with all 40 criteria is beyond the scope of this paper. Thus, we focus on a single criterion (i.e., PM2). It is important to underscore that while we focus on this single criterion, the method has been applied in its entirety to calculate the InterpretationResult.

The PM2 criterion evaluates if *the variant is absent from controls in an outbred population, from gnomAD v2.1 (non-cancer, exome only subset) and gnomAD v3.1 (non-cancer). Region around the variant must have an average read depth > 25.* From the description of this criterion, two metrics can be identified: i) the variant must be absent from the population, and ii) the average read depth must be greater than 25. In this guideline, even though its not common practice, the
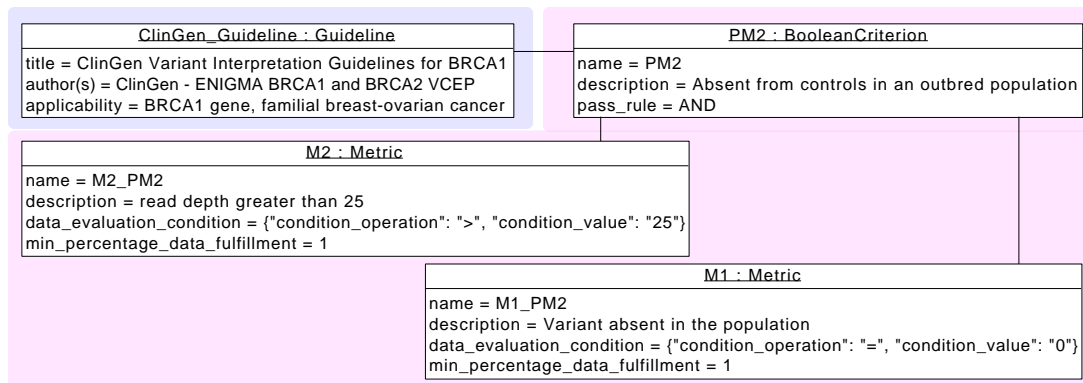
**Figure 3:** ClinGen guideline and PM2 criterion and metrics instantiation.

criterion also specifies the data sources to be used: the non-cancer exome dataset of gnomAD 2.1 and the non-cancer dataset of gnomAD 3.1. Therefore, in this definition it is crucial to distinguish between the definition of the metrics and the specific data sources in which this condition must be evaluated.

Figure 3 illustrates the instantiation of the ClinGen guideline, specifically highlighting the PM2 criterion and its associated metrics, as defined by the VarClaMM meta-model. It is important to note that the evaluation of the variant's absence in the population will be determined by evaluating whether the variant's frequency is equal to zero.

### 3. Select the required data elements and data sources.

The first task in this step is to determine which DATAELEMENT each metric will evaluate. As specified in Section 3, the selection of DATAELEMENT depends on the DATAMODEL chosen as the framework for data representation. For our purposes, we have selected the VarDaM DATAMODEL, which has been specifically designed for variant interpretation data representation [15].

The first metric evaluates the frequency of a variant in a population, while the second metric assesses read_depth[1]. Both concepts are represented as attributes of a class called ALLELE-FREQUENCY in the VarDaM model. Figure 4 illustrates the instantiation of each DATAELEMENT within the context of VarDaM.

Once we have identified the required DATAELEMENTS, the next step is to determine the DATASOURCES from which to obtain them. In this example, the task is straightforward because the criterion definition specifies the necessary data sources: the non-cancer exome dataset of gnomAD 2.1 and the non-cancer dataset of gnomAD 3.1.

### 4. Obtain the classification results.

The previous steps have allowed us to define the interpretation framework. Now, we need to execute it for the selected variant. To do so, first, we need to collect the values of the selected

---

[1]The number of times a specific variant in the DNA is read during the sequencing process.
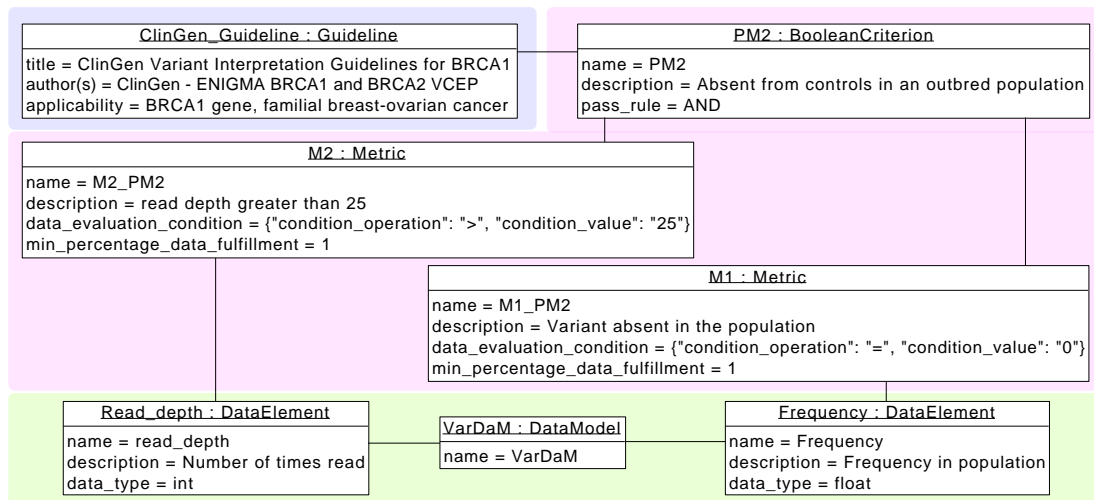
**Figure 4:** Data elements for the PM2 criterion's metrics.

DATAELEMENTS from the specified DATASOURCES (e.g., obtain the DATAELEMENTINDATASOURCE). For the c.191G>A variant, we found a frequency of 0 (i.e., it is absent) and a mean read_depth of 31 in the non-cancer dataset of gnomAD 3.1. Similarly, the variant has an approximate read_depth of 40 and a frequency of 0 in the non-cancer exome dataset of gnomAD 2.1.

Consequently, as the conditions established by the two metrics are met for all DATAELE-MENTSINDATASOURCE, all METRICDATAEVALUATIONRESULT are set to true. Based on this, we can determine the METRICRESULT and the CRITERIONRESULT. The METRICRESULT for both defined metrics is true, as all METRICDATAEVALUATIONRESULT are true. Therefore, since both METRICRESULT are true, the CRITERIONRESULT is also true.

Finally, the INTERPRETATIONRESULT is derived from all the CRITERIONRESULT. After applying all criteria and metrics defined by the ClinGen guideline —not just the one displayed here— the variant is classified as *Pathogenic*. This indicates that the variant is responsible for the familial breast-ovarian cancer syndrome exhibited by the patient. The final model instantiation for this variant is shown in Figure 5.

## 5. Conclusions

In conclusion, the SISO method has set the basis for variant interpretation systematization by providing a structured framework grounded in the VarClaMM meta-model. This method addresses the inherent imprecision and variability in current variant interpretation practices by defining clear steps for carrying out variant interpretation. By applying this method to the c.191G>A variant in the BRCA1 gene, we have demonstrated its practical utility and the ability to achieve consistent and reproducible results. This case study underscores the potential of the SISO method to enhance the accuracy and reliability of variant interpretation in clinical
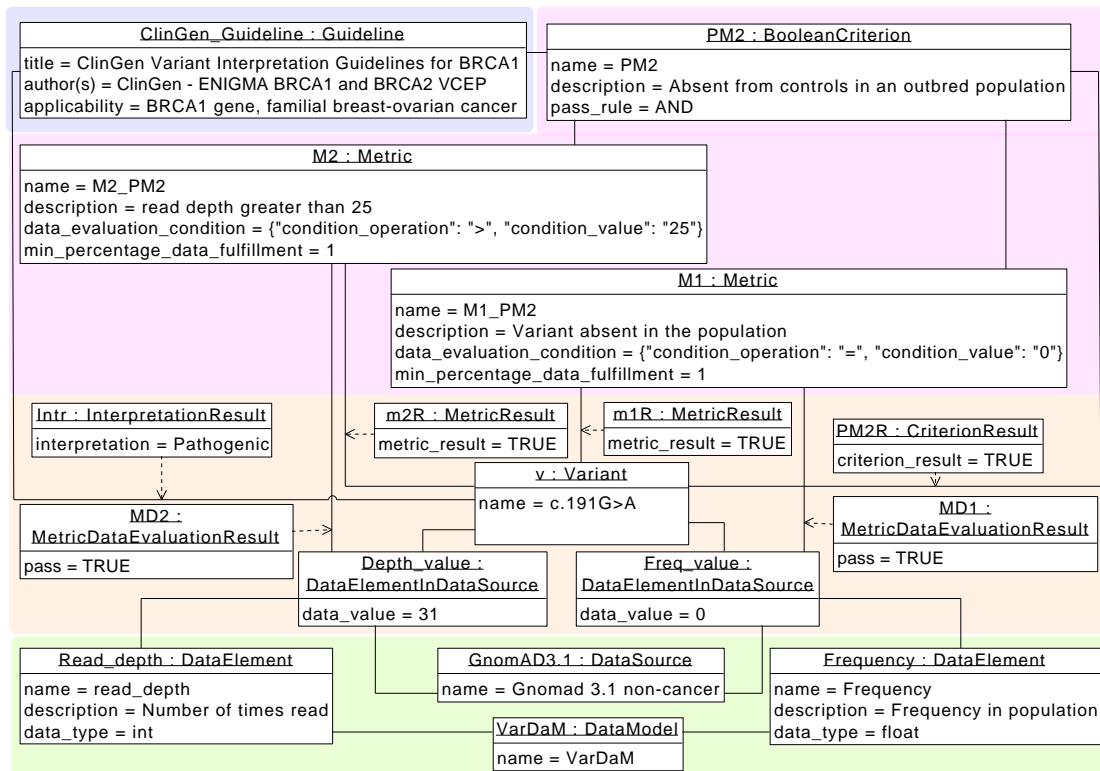
**Figure 5:** Complete instantiation of VarClaMM for the c.191G>A variant. For clarity, only the data from GnomAD 3.1 has been displayed. The connection between the METRICS and DATAELEMENTS has been omitted for the same reason.

settings.

Future work will focus on developing technological support for each stage of the SISO method. This advancement will enable the automated application of each step, achieving full systematization of variant interpretation. By automating these processes, we aim to enhance the precision and efficiency of variant interpretation, ultimately contributing to the broader goal of advancing medicine. This will ensure that patients receive the most accurate and effective clinical care tailored to their unique genetic characteristics.

## Acknowledgments

# References

[1] National Cancer Institute, Variant, 2024. URL: https://www.cancer.gov/publications/dictionaries/genetics-dictionary/def/variant.

[2] D. Rigden, et al., The 2023 nucleic acids research database issue and the online molecular biology database collection, Nucleic Acids Research 51 (2023) D1–D8.

[3] A. Furqan, et al., Care in specialized centers and data sharing increase agreement in hypertrophic cardiomyopathy genetic test interpretation, Circulation: Cardiovascular Genetics 10 (2017) e001700.

[4] M. S. Lebo, et al., Data sharing as a national quality improvement program: reporting on BRCA1 and BRCA2 variant-interpretation comparisons through the Canadian Open Genetics Repository (COGR), Genetics in Medicine 20 (2018) 294–302.

[5] Y.-E. Kim, et al., Challenges and considerations in sequence variant interpretation for mendelian disorders, Annals of Laboratory Medicine 39 (2019) 421.

[6] S. M. Harrison, et al., Clinical laboratories collaborate to resolve differences in variant interpretations submitted to ClinVar, Genetics in Medicine 19 (2017) 1096–1104.

[7] N. Agaoglu, et al., Consistency of variant interpretations among bioinformaticians and clinical geneticists in hereditary cancer panels, European Journal of Human Genetics 30 (2022).

[8] M. Costa, et al., A reference meta-model to understand DNA variant interpretation guidelines, in: J. P. A. Almeida, J. Borbinha, G. Guizzardi, S. Link, J. Zdravkovic (Eds.), Conceptual Modeling - 42nd International Conference, ER 2023, Lisbon, Portugal, November 6-9, 2023, Proceedings, volume 14320 of *Lecture Notes in Computer Science*, Springer, 2023, pp. 375–393. doi:10.1007/978-3-031-47262-6\_20.

[9] S. Richards, et al., Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the american college of medical genetics and genomics and the association for molecular pathology, Genetics in medicine 17 (2015) 405–423.

[10] E. Zirkelbach, et al., Managing variant interpretation discrepancies in hereditary cancer: Clinical practice, concerns, and desired resources, Journal of Genetic Counseling 27 (2018) 761–769.

[11] M. Costa, et al., The consequences of data dispersion in genomics: a comparative analysis of data sources for precision medicine, BMC Medical Informatics and Decision Making 23 (2023) 256.

[12] K. Karczewski, et al., The exac browser: Displaying reference data information from over 60 000 exomes, Nucleic acids research 45 (2016). doi:10.1093/nar/gkw971.

[13] K. Karczewski, et al., The mutational constraint spectrum quantified from variation in 141,456 humans, Nature 581 (2020) 434–443. doi:10.1038/s41586-020-2308-7.

[14] ClinGen - ENIGMA BRCA1 and BRCA2 VCEP, Clingen enigma brca1 and brca2 expert panel specifications to the acmg/amp variant interpretation guidelines for brca1 version 1.1.0, 2024. URL: https://cspec.genome.network/cspec/ui/svi/doc/GN092.

[15] M. Costa, et al., Comprehensive representation of variation interpretation data via conceptual modeling, 2023, pp. 25–34. doi:10.1007/978-3-031-47112-4_3.