



HVP/DA/001-02/EN

2014-05-08

WG05

GDSDBAC

DRAFT FOR APPROVAL

GENE/DISEASE SPECIFIC VARIANT DATABASE QUALITY PARAMETERS

Notice

This document is not an HVP Standard or Guideline. This document is a Draft that has been distributed for approval by the Gene/Disease Specific Database Advisory Council. This document is subject to change without notice. **USE AT YOUR OWN RISK!** Because this is an unapproved draft, this document must not be utilized for any conformance/compliance purposes.

Authors

Mauno Vihinen, Lund University, Lund, Sweden,
Christophe Beroud, INSERM UMR_S910, Marseille, France,
Andrew Devereau, National Genetics Reference Laboratories, Manchester, UK,
John Hancock, University of Cambridge, Cambridge, UK,
Peter Taschner, Leiden University Medical Center, Leiden, The Netherlands

Editor

Timothy D. Smith, Human Variome Project International Ltd, Australia

Published by:

Human Variome Project International Limited
Level 5, 234 Queensberry Street, The University of Melbourne VIC 3010, AUSTRALIA



Copyright © 2014 by the listed members of the Human Variome Project Working Group WG05: Variant Database Quality Assessment and licensed under the terms of the Creative Commons Attribution-ShareAlike 4.0 International License.

Published 2014-05-08. Printed in Australia.

PDF: ISBN <<PDF_ISBN>>
Print: ISBN <<PRINT_ISBN>>

1 **Contents**

2 **Foreword**.....5

3 This Document5

4 **Important Notice**6

5 **Introduction**7

6 **1 Scope**.....7

7 **2 Quality Assessment Principles**7

8 **3 Quality Assessment Criteria**.....7

9 3.1 Data quality.....8

10 3.2 Technical quality.....8

11 3.3 Accessibility.....8

12 3.4 Timeliness9

13 **4 Assessing Quality**.....9

14 **5 Relevant references**9

15

1

1 **Foreword**

2 The Human Variome Project is an international consortium of researchers, policy makers and healthcare
3 professionals committed to the free and open collection, curation, interpretation and sharing of genomic
4 knowledge.

5 The Human Variome Project Consortium envisions a world where the availability of and access to genetic
6 variation information is not an impediment to diagnosis and treatment; where the burden of genetic
7 disease on the human population is significantly decreased; and where the sharing of genetic variation
8 information is standard clinical practice.

9 To facilitate worldwide and interoperable sharing of genomic knowledge, the Human Variome Project
10 Consortium produces Standards and Guidelines. HVP Standards are those systems, procedures and
11 technologies that the Human Variome Project Consortium has determined shall be used by the
12 community. These carry more weight than the less prescriptive HVP Guidelines, which cover those
13 systems, procedures and technologies that the Human Variome Project Consortium has determined would
14 be beneficial for the community to adopt.

15 HVP Standards and Guidelines are central to supporting the work of the Human Variome Project
16 Consortium and cover a wide range of fields and disciplines, from ethics to nomenclature, data transfer
17 protocols to collection protocols for clinical data. They can be thought of as both technical manuals and
18 scientific documents, and while the impact of HVP Standards and Guidelines differ, they are both
19 generated in a similar fashion.

20 HVP Standards and Guidelines make the collection, curation and sharing of information more efficient
21 and reliable by establishing consistent protocols that can be universally understood. They facilitate
22 interconnection of and interoperability between different systems.

23 HVP Standards and Guidelines represent a consensus of the Human Variome Project Consortium, each
24 member of which has had the opportunity to participate in the development and review of each standard
25 and guideline. In addition, as every effort is made to include all interests in the activity, HVP Standards
26 and Guidelines can be considered to be representative of all interests concerned within the scope of each
27 Standard or Guideline.

28 The Human Variome Project defines consensus as significant agreement between all affected parties
29 covered by the scope of the standard or guideline. Consensus requires that all views and objections be
30 considered, and that a concerted effort be made toward their resolution.

31 More information on the Human Variome Project is available at the Project's website
32 (<http://www.humanvariomeproject.org/>). Procedures for the development of HVP Standards and
33 Guidelines can be found in *PD06-2011: Standards Development Process*, available at
34 <http://short.variome.org/PD06-2011>.

35 **This Document**

36 This document has been prepared the HVP Working Group: WG05: Variant Database Quality
37 Assessment. The Gene/Disease Database Advisory Council acted as Sponsoring Council.

38 An Exposure Draft of this Document was released for comment to the Human Variome Project
39 Consortium on 2014-02-13. The consultation period ended 2014-04-14.

1 A Draft for Approval was submitted to the Gene/Disease Specific Database Advisory Council 2014-05-
2 08.

3 **Important Notice**

4 HVP Standards and Guidelines are not intended to replace or substitute for any applicable legislation or
5 regulation in any jurisdiction, or any institutional policy or funding agreement that a genetic variation
6 information resource is operating under. Implementers of HVP Standards and Guidelines are responsible
7 for determining and complying with all appropriate ethical and cultural protection practices and all
8 applicable laws, regulations, policies and agreements.

9

1 Introduction

Thousands of variation databases have been developed during the last decades. They have widely varying contents, resources, as well as quality. Quality is becoming ever more important as the use of these databases by the community increases and users start to combine and compare information in the different resources.

Computer science and IT communities have discussed and developed criteria for database platform/system quality, however, there are no widely accepted systematic criteria and evaluation systems. Each domain has developed its own. Closest to this field comes BioDBcore, which has developed core attributes of biological databases. However, it is not a quality scheme.

The criteria described in computer science and information technology quality papers are more suitable to assess database platforms (e.g. DMuDB, LOVD, MUTbase, UMD etc.) than for variant database content. For some platforms, there might be little difference due to limited customization and flexibility. Even in the same installation using the latest version of LOVD, one gene variant database might be of much better quality than another depending on the curator(s).

1 Scope

This document details the quality assessment parameters that might be used to evaluate the quality of genetic variation databases and stimulate database curators to make improvements where they are needed.

How these quality assessment criteria might be implemented in a quality accreditation scheme, including the criteria that should be used to assess compliance with each parameter is outside the scope of this document.

2 Quality Assessment Principles

The guiding principle used throughout this document is simplicity while maintaining a good overview of the quality of the database and its content. The goal is to capture the major issues pertinent to quality balanced with a simple system for database users to comprehend and database administrators and curators to apply.

3 Quality Assessment Criteria

Quality evaluation for genetic variation databases should consider four specific areas. The basic principle is to have a simple system, which still can provide a good overview of the quality. There could be overall ranking in the style of hotels with stars, which could then be broken down to individual scores. The main evaluation areas suggested are:

- a) Data quality
- b) Technical quality
- c) Accessibility
- d) Timeliness

These four quality evaluation criteria can be divided into more specific items:

1 **3.1 Data quality**

- 2 a) Database scope and purpose
- 3 b) Contents (patient data, if so, how much, fulfillment of minimal requirements¹)
- 4 c) Completeness of data
- 5 d) Coverage of variations
- 6 e) Accuracy, error rate
- 7 f) Consistency (language, spelling, reference sequences) - type of curation
- 8 g) Integration to other resources
- 9 h) Use of standards (gene names, reference sequences, variation nomenclature, data models,
- 10 ontologies ...)
- 11 i) Date stamps
- 12 j) Authority, curatorial team competence
- 13 k) Contact details
- 14 l) Use of references
- 15 m) Data collection, sources
- 16 n) Definition of pathogenicity used
- 17 o) What kind of data (e.g. NGS data, SNVs)
- 18 p) Range of numerical values (min to max)
- 19 q) Correct units
- 20 r) No data lost when inputted
- 21 s) Consents, privacy
- 22 t) Ethical issues
- 23 u) Public/nonpublic data

24 **3.2 Technical quality**

- 25 a) Database management system, suitability
- 26 b) Speed of access
- 27 c) Quality control measures implemented
- 28 d) Use of automatic steps (e.g. HGVS names with Mutalyzer)
- 29 e) How corrections made
- 30 f) Reliability (uptime)
- 31 g) Version history
- 32 h) Functional links
- 33 i) Use on different browsers
- 34 j) Data security (operating system, backups, firewall)

35 **3.3 Accessibility**

- 36 a) Design
- 37 b) Readability
- 38 c) User friendly, logical interface
- 39 d) Navigation
- 40 e) Language, correctness
- 41 f) Ease of use
- 42 g) Consistency on the site

¹ HVP Standards and Guidelines for minimum requirements for molecular and clinical data are currently under development.

- 1 h) Interactivity
- 2 i) Available freely, is registration required
- 3 j) How to contact
- 4 k) Help, support, tutorials
- 5 l) Documentation (purpose, scope, motivation, copyright, disclaimer, database policy, data items,
6 annotation guidelines)
- 7 m) Searchability (search engine/possibilities)
- 8 n) Downloadability of data/search results
- 9 o) Format(s) of output (for further analyses)
- 10 p) Graphics (use, clarity etc)
- 11 q) Contacts to community, support groups etc
- 12 r) Other modes of access, esp. web services access

13 **3.4 Timeliness**

- 14 a) Update frequency - including whether it is currently being maintained
- 15 b) Currency of updates (how new data included)
- 16 c) Versioning policy (are old versions available?)

17 **4 Assessing Quality**

18 Although consideration of how a database quality assessment program might be implemented is outside
19 the scope of this document, the following issues may be of use when such a program is developed.

20 Technical information required to perform a quality audit of a database might be provided by the database
21 administrator. This kind of self-assessment or would save time for the quality evaluators. For popular
22 platforms (e.g. LOVD) scripts could be developed to automatically collect a number of facts.

23 **5 Relevant references**

- 24 Celli, J, R Dagleish, M Vihinen, PE Taschner, and JT den Dunnen. "Curating gene variant databases
25 (LSDBs): Toward a universal standard." *Human Mutation* 33 (2012): 291-297.
- 26 Cotton, RGH, et al. "Recommendations for locus-specific databases and their curation." *Human Mutation*
27 29 (2008): 2-5.
- 28 den Dunnen, JT, et al. "Sharing data between LSDBs and central repositories." *Human Mutation* 30
29 (2009): 493-495.
- 30 Gaudet, P, et al. "Towards BioDBcore: a community-defined information specification for biological
31 databases." *Nucleic Acids Research* 39 (2011): D7-10.
- 32 Smedley, D, et al. "Finding and sharing: new approaches to registries of databases and services for the
33 biomedical sciences." *Database (Oxford)*, 2010: doi:10.1093/database/baq1014.
- 34 Vihinen, M, JT den Dunnen, R Dagleish, and RGH Cotton. "Guidelines for establishing locus specific
35 databases." *Human Mutation* 33 (2012): 298-305.

36