

Musio A, Selicorni A, Focarelli ML et al. (2006) X-linked Cornelia de Lange syndrome owing to *SMC1L1* mutations. *Nature Genet* **38**: 528-530.

Opitz JM (1985) Editorial comment: the Brachmann de Lange syndrome. *Am J Med Genet* **22**: 89-102.

Russel DJ, Nelson BI, Nadya JK et al. (1993) Prenatal diagnosis of congenital diaphragmatic hernia not amenable to prenatal or neonatal repair: Brachmann-de Lange syndrome. *Am J Med Genet* **47**: 1022-1023.

Schoumans J, Wincent J, Barbaro M et al. (2007) Comprehensive mutational analysis of a cohort of Swedish Cornelia de Lange syndrome patients. *Eur J Hum Genet* **15**: 143-149.

Suzuki S, Mukai K, Takei Y et al. (1999) [A autopsy report of Cornelia de Lange syndrome.] *Shounika Rinshou* 52: 1967-1972 (In Japanese.)

Tibboel D, Gaag AVD (1996) Etiological genetic factors in congenital diaphragmatic hernia. *Clin Perinatol* **23**: 689-699.

Tonkin ET, Wang TJ, Lisgo S et al. (2004) *NIPBL*, encoding a homolog of fungal *Scc2*-type sister chromatid cohesion proteins and fly *Nipped-B*, is mutated in Cornelia de Lange syndrome. *Nature Genet* **36**: 636-641.

**Figure legends**

**Fig. 1** Photographs highlighting the patient's symptoms (a and b) including hypertrichosis, short extremities, hypoplasia of the nipple and umbilicus, synophrys of the face, short and upturned nose or anteverted nostrils, long philtrum, thin lip, small mandible, short cervix and single bimanual palmar flexion curve without syndactyly or defects of the fingers.

**Fig. 2** X-ray of the entire body showing the hanging bell-shaped thoracic cage, low pneumatization in the bilateral lungs and a stomach bubble located in the middle thorax.

**Fig. 3** Denaturing high-performance liquid chromatography of the 29th exon of the *NIPBL* gene (upper panel: control, lower panel: patient). Arrow shows the abnormal peak in the translation area (29th exon) of the *NIPBL* gene.

**Table 1. Examination of the umbilical cord and peripheral blood**

	< Cord Blood >	< Patient's Blood > (at 47 minutes after birth)
<b>•Blood gas analysis</b>		
pH	7.276	6.763
pCO <sub>2</sub>	48.9 mmHg	188.0 mmHg
pO <sub>2</sub>	20.5 mmHg	3.2 mmHg
Base excess	-4.3 mmol/l	-16.5 mmol/l
<b>•Blood cell counts</b>		
White blood count	5900 / $\mu$ l	4900 / $\mu$ l
Platelet count	221000 / $\mu$	93000 / $\mu$ l
<b>•Chemistry</b>		
C-reactive protein	<0.05 mg/dl	<0.1 mg/dl
Sodium	140 mmol/l	183 mmol/l
Potassium	4.5 mmol/l	5.3 mmol/l
Calcium	9.5 mg/dl	8.3 mg/dl
Hemoglobin	13.6 g/dl	6.8 g/dl

**Table 2. Findings in the present case**

	<b>This patient's findings</b>
<b>&lt;Cornelia de Lange (1933)&gt;</b>	
<b>1.Mental retardation</b>	?
<b>2.Low birth weight</b>	+
<b>3.Dwarfism</b>	?
<b>4.Microbrachycephaly</b>	+
<b>5.Heavy eyebrows meeting at the midline</b>	+
<b>6.Long eyelashes</b>	+
<b>7.Low ear insertion</b>	-
<b>8.Small hands and feet</b>	+
<b>9.Proximally placed thumb</b>	-
<b>10.Syndactyly of the toes</b>	-
<b>&lt;Beck (1976)&gt;</b>	
<b>1.Low hair line on forehead</b>	+
<b>2.Low hair line on neck</b>	+
<b>3.Long philtrum</b>	+
<b>4.Bushy eyebrows</b>	+
<b>5.Confluent eyebrows</b>	+
<b>6.Thick eyelashes</b>	+
<b>7.Antimongoloid eye slanting</b>	-

<b>8. Anteverted nostrils</b>	<b>+</b>
<b>9. Crescent-shaped mouth</b>	<b>+</b>
<b>10. Thin prolabium</b>	<b>+</b>

**+, Present; ?, Not detected due to early death.**

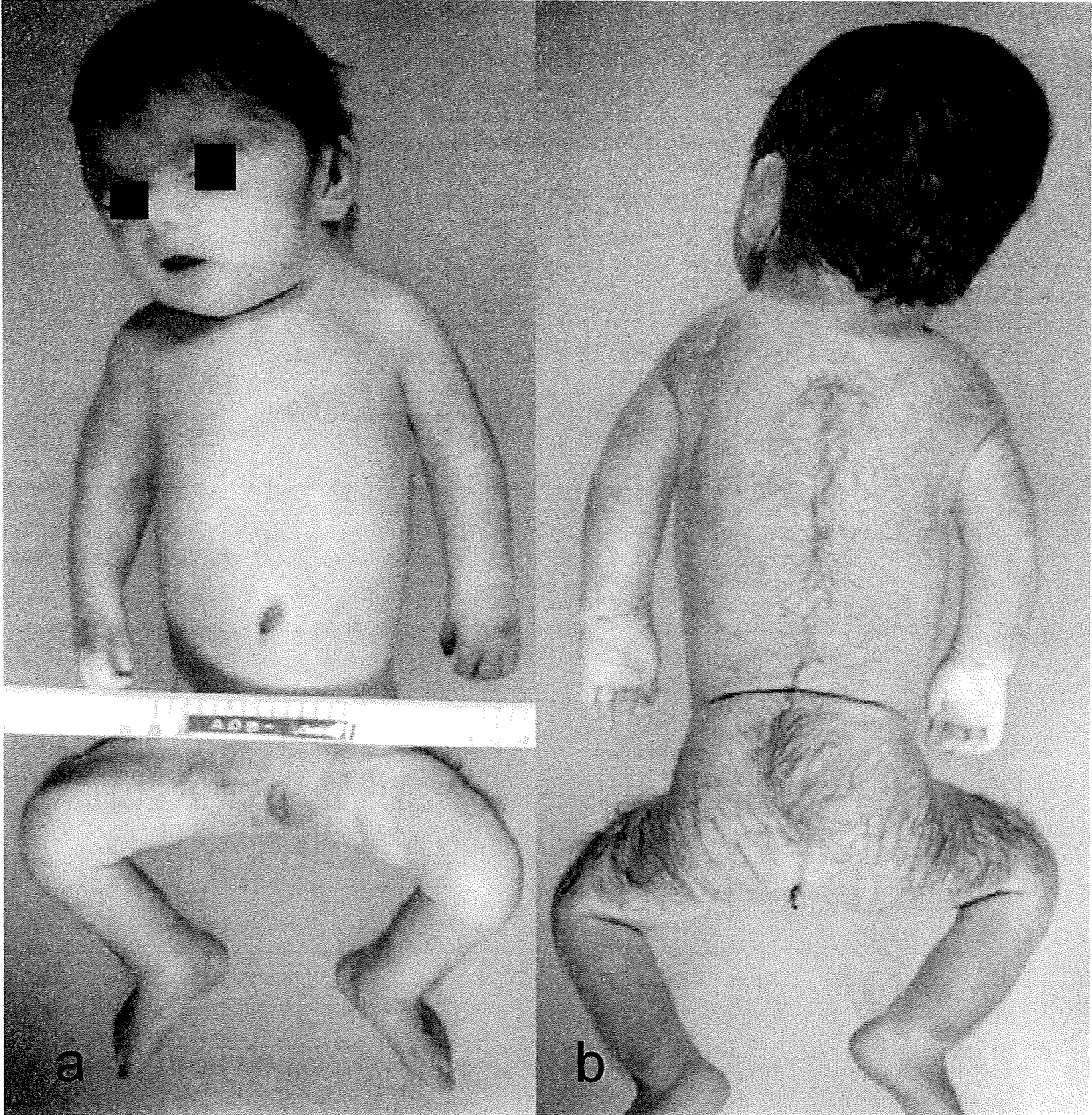


Fig 1



Fig 2

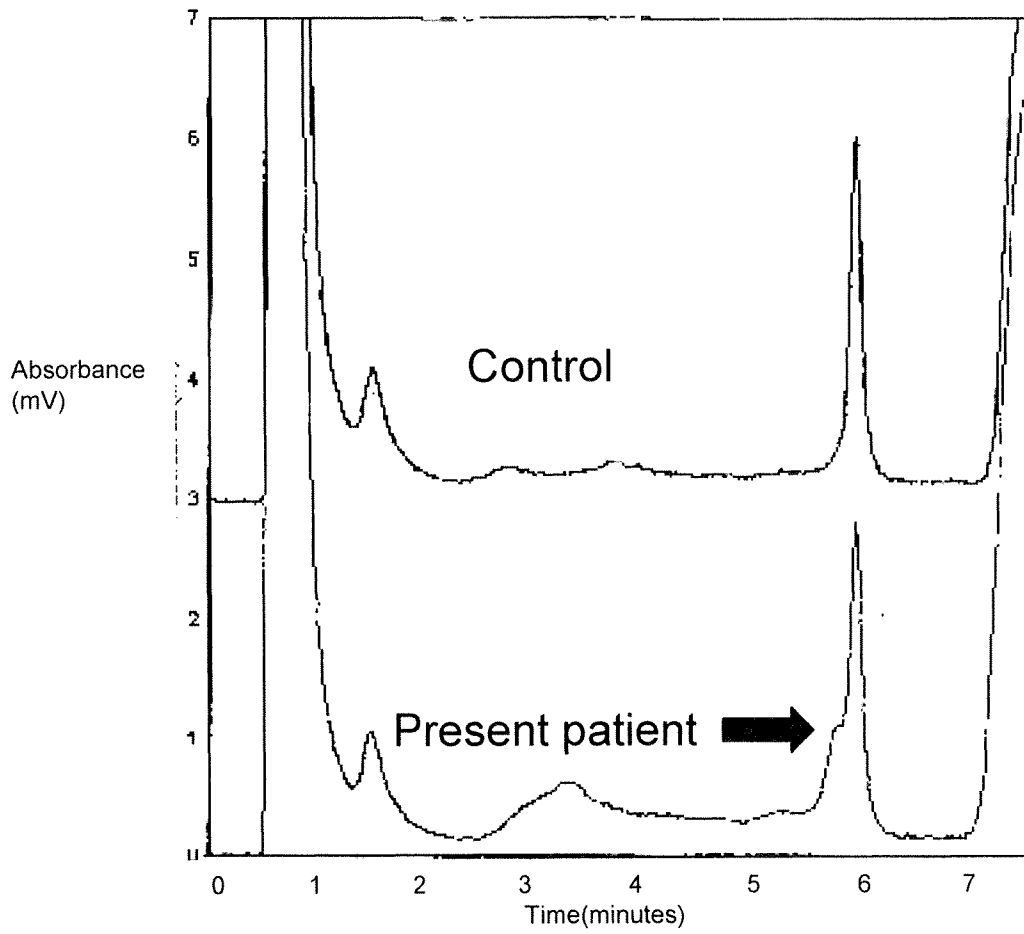


Fig 3



## Planning the Human Variome Project: The Spain Report\*

Jim Kaput,<sup>1†</sup> Richard G.H. Cotton,<sup>2,3†</sup> Lauren Hardman,<sup>2†</sup> Michael Watson,<sup>4†</sup> Aida I. Al Aqeel,<sup>5</sup> Jumana Y. Al-Aama,<sup>6</sup> Fahd Al-Mulla,<sup>7</sup> Santos Alonso,<sup>8</sup> Stefan Aretz,<sup>9</sup> Arleen D. Auerbach,<sup>10</sup> Bharati Bapat,<sup>11</sup> Inge T. Bernstein,<sup>12</sup> Jong Bhak,<sup>13</sup> Stacey L. Bleoo,<sup>14</sup> Helmut Blöcker,<sup>15</sup> Steven E. Brenner,<sup>16</sup> John Burn,<sup>17</sup> Mariona Bustamante,<sup>18</sup> Rita Calzone,<sup>19</sup> Anne Cambon-Thomsen,<sup>20</sup> Michele Cargill,<sup>21</sup> Paola Carrera,<sup>22</sup> Lawrence Cavedon,<sup>23</sup> Yoon Shin Cho,<sup>24</sup> Yeun-Jun Chung,<sup>25</sup> Mireille Claustres,<sup>26</sup> Garry Cutting,<sup>27</sup> Raymond Dalglish,<sup>28</sup> Johan T. den Dunnen,<sup>29</sup> Carlos Díaz,<sup>30</sup> Steven Dobrowolski,<sup>31</sup> M. Rosário N. dos Santos,<sup>32</sup> Rosemary Ekong,<sup>33</sup> Simon B. Flanagan,<sup>34</sup> Paul Fliceck,<sup>35</sup> Yoichi Furukawa,<sup>36</sup> Maurizio Genuardi,<sup>37</sup> Ho Ghang,<sup>13</sup> Maria V. Golubenko,<sup>38</sup> Marc S. Greenblatt,<sup>39</sup> Ada Hamosh,<sup>40</sup> John M. Hancock,<sup>41</sup> Ross Hardison,<sup>42</sup> Terence M. Harrison,<sup>43</sup> Robert Hoffmann,<sup>44</sup> Rania Horaitis,<sup>2</sup> Heather J. Howard,<sup>2</sup> Carol Isaacson Barash,<sup>45</sup> Neskuts Izagirre,<sup>46</sup> Jongsun Jung,<sup>24</sup> Toshio Kojima,<sup>47</sup> Sandrine Laradi,<sup>48</sup> Yeon-Su Lee,<sup>49</sup> Jong-Young Lee,<sup>24</sup> Vera L. Gil-da-Silva-Lopes,<sup>50</sup> Finlay A. Macrae,<sup>51</sup> Donna Maglott,<sup>52</sup> Makia J. Marafie,<sup>53</sup> Steven G.E. Marsh,<sup>54</sup> Yoichi Matsubara,<sup>55</sup> Ludwine M. Messiaen,<sup>56</sup> Gabriela Möslein,<sup>57</sup> Mihai G. Netea,<sup>58</sup> Melissa L. Norton,<sup>59</sup> Peter J. Oefner,<sup>60</sup> William S. Oetting,<sup>61</sup> James C. O'Leary,<sup>62</sup> Ana Maria Oller de Ramirez,<sup>63</sup> Mark H. Paalman,<sup>64</sup> Jillian Parboosingh,<sup>65</sup> George P. Patrinos,<sup>66</sup> Giuditta Perozzi,<sup>67</sup> Ian R. Phillips,<sup>68</sup> Sue Povey,<sup>64</sup> Suyash Prasad,<sup>69</sup> Ming Qi,<sup>70</sup> David J. Quin,<sup>71</sup> Rajkumar S. Ramesar,<sup>72</sup> C. Sue Richards,<sup>73</sup> Judith Savige,<sup>74</sup> Dagmar G. Scheible,<sup>75</sup> Rodney J. Scott,<sup>76</sup> Daniela Seminara,<sup>77</sup> Elizabeth A. Shephard,<sup>78</sup> Rolf H. Sijmons,<sup>79</sup> Timothy D. Smith,<sup>2</sup> María-Jesús Sobrido,<sup>80</sup> Toshihiro Tanaka,<sup>81</sup> Sean V. Tavtigian,<sup>82</sup> Graham R. Taylor,<sup>83</sup> Jon Teague,<sup>84</sup> Thoralf Töpel,<sup>85</sup> Mollie Ullman-Cullere,<sup>86</sup> Joji Utsunomiya,<sup>47</sup> Henk J. van Kranen,<sup>87</sup> Mauno Vihinen,<sup>88</sup> Elizabeth Webb,<sup>2</sup> Thomas K. Weber,<sup>89</sup> Meredith Yeager,<sup>90</sup> Young I. Yeom,<sup>91</sup> Seon-Hee Yim,<sup>92</sup> and Hyang-Sook Yoo,<sup>93</sup> on behalf of contributors to the Human Variome Project Planning Meeting

<sup>1</sup>Division of Personalised Nutrition and Medicine, FDA/National Center for Toxicological Research, Jefferson, Arkansas

<sup>2</sup>Genomic Disorders Research Centre, Melbourne, Australia

<sup>3</sup>Faculty of Medicine, Dentistry and Health Sciences, University of Melbourne, Parkville, Australia

<sup>4</sup>American College of Medical Genetics, Bethesda, Maryland

Received 2 December 2008; accepted revised manuscript 22 December 2008.

Published online 20 March 2009 in Wiley InterScience (www.interscience.wiley.com). DOI 10.1002/humu.20972

**ABSTRACT:** The remarkable progress in characterizing the human genome sequence, exemplified by the Human Genome Project and the HapMap Consortium, has led to the perception that knowledge and the tools (e.g., microarrays) are sufficient for many if not most biomedical research efforts. A large amount of data from diverse studies proves this perception inaccurate at best, and at worst, an impediment for further efforts to characterize the variation in the human genome. Because variation in genotype and environment are the fundamental basis to understand phenotypic variability and heritability at the population level, identifying the range of human genetic variation is crucial to the development of personalized nutrition and medicine. The Human Variome Project (HVP; <http://www.humanvariomeproject.org/>) was proposed initially to systematically collect mutations that cause human disease

and create a cyber infrastructure to link locus specific databases (LSDB). We report here the discussions and recommendations from the 2008 HVP planning meeting held in San Feliu de Guixols, Spain, in May 2008.

Hum Mutat 30, 496–510, 2009. © 2009 Wiley-Liss, Inc.

**KEY WORDS:** variome; genome; mutation; database; genetic disease

## Introduction

The completion of the consensus sequence of the human genome [Lander et al., 2001; Venter et al., 2001] ushered in the “postgenomic era” of science—that is, experiments could be designed using the reference sequence of the genome without need for additional sequencing efforts. Subsequent publication of the human haplotype map, an analysis of sequence diversity in 270 individuals from four ancestral populations [Frazer et al., 2007; International HapMap Consortium, 2003, 2004], provided knowledge for building reagents for further genetic analyses. The knowledge and sequence information provided the resources to analyze the genetic contribution to virtually all measurable phenotypes. These efforts and the resulting databases complemented the long-standing efforts by geneticists to locate, identify, and characterize mutations that cause monogenic and polygenic diseases in humans, an effort begun by McKusick and colleagues

\*Affiliations for all authors are listed in the Appendix. Members of working groups, their affiliations, institutions, and email addresses can be found at <http://www.humanvariomeproject.org>.

<sup>†</sup>The first four authors contributed equally to the preparation of this manuscript.

Additional Supporting Information may be found on the online version of this article.

<sup>†</sup>Correspondence to: Jim Kaput, Division of Personalized Nutrition and Medicine, FDA/National Center for Toxicological Research, 3900 NCTR Road, Jefferson, AR 72079. E-mail: James.kaput@fda.hhs.gov

in the 1950s (reviewed in [McKusick, 2006]) and now cataloged in the Online Mendelian Inheritance in Man (<http://www.ncbi.nlm.nih.gov/sites/entrez?db=omim>) [Amberger et al., 2009] and the Human Gene Mutation Database (<http://www.hgmd.cf.ac.uk>) [Stenson et al., 2008]. NCBI's dbGaP (Genotype and Phenotype database: <http://www.ncbi.nlm.nih.gov/gap>) [Mailman et al., 2007], HuGENet (Human Genome Epidemiology Network: <http://www.cdc.gov/genomics/hugenet/>), EBI's EGA (European Genotype Archive: <http://www.ebi.ac.uk/ega/>), FINDbase (Frequency of INherited Disorders database: <http://www.findbase.org>) [van Baal et al., 2007] and GAD (the Genetic Association Database: <http://geneticassociationdb.nih.gov/>) [Becker et al., 2004], are repositories for data from population studies associating genetic variation with phenotypes. Most of these databases are study oriented, and analyze existing polymorphisms rather than focusing on the discovery of new genetic variants.

A large amount of other mutation or gene variation data, however, is likely to exist on servers in laboratories scattered throughout the world. Each of these databases may contain valuable data for other studies and for the medical practitioner. The Human Variome Project (HVP; <http://www.humanvariomeproject.org/>) was previously proposed to systematically collect mutations that cause human disease [Cotton et al., 2007a,b; Ring et al., 2006] and create a cyber infrastructure to link locus specific databases (LSDBs). Local experts would curate individual LSDBs, but each would have similar architecture, ontologies, and data elements, allowing for interoperability. Links to national and international databases such as at the National Center for Biotechnology Information (NCBI: <http://www.ncbi.nlm.nih.gov/>) and the European Bioinformatics Institute (EBI: <http://www.ebi.ac.uk/>) would consolidate the knowledge of the curation done by local experts. We report here the discussions and recommendations from the 2008 HVP planning meeting held in San Feliu de Guixols, Spain, in May 2008, to further the development of the HVP.

The theoretical rationale for resequencing genes from individuals in diverse populations is that the existing databases have focused primarily on Europeans and their descendants, and are therefore a relatively narrow subdivision of the entire range of human genetic diversity. Published data supporting a concerted resequencing effort for monogenic and complex diseases come from independent and unlinked studies:

- The molecular basis of 3,770 diseases caused by mutations in 2,239 genes was known as of September 2008 (<http://www.ncbi.nlm.nih.gov/Omim/mimstats.html>). Many monogenic diseases are caused by different mutations in one gene, and all monogenic diseases are known to have variable age of onset, severity, and outcome (e.g., [McKusick, 2007; Ropers, 2007]). Differences in monogenic disease phenotype may be caused by variations in location of mutations, by modifier genes that interact with the disease causing allele [McKusick, 2007], and by gene-environment interactions [Ordovas and Corella, 2006]. Characterizing causative mutations in familial and sporadic cases in diverse populations is warranted for a full understanding of each disease.
- The molecular basis of over 3,700 other phenotypes are either suspected to be Mendelian disorders or are unknown (<http://www.ncbi.nlm.nih.gov/Omim/mimstats.html>).
- Simple monogenic traits may be caused by different gene variants. The ability to hydrolyze lactose as an adult, which is called lactase persistence, occurs in ~30% of the world's population [Lomer et al., 2008]. Expression of the lactase gene postweaning has been associated with a C/T variant at position -13910 from the start of the lactase (*LCT*) gene in Finnish families and 236 individuals from four populations (Germany, Italy, South Korea, and Finland) [Enattah et al., 2002]. However, lactase persistence is associated with a different variant (G/C at -14010 from *LCT*) in Kenyans, Tanzanians, and Sudanese [Tishkoff et al., 2007]. Other populations with higher percentages of individuals with lactose tolerance have not been analyzed [Montgomery et al., 2007]. Variation in amounts of lactose required to induce intestinal bloating and diarrhea, severity, and age of onset are observed in reference populations (lactose intolerant) and in populations where the lactase persistence variants are more common [Lomer et al., 2008; Montgomery et al., 2007].
- Over 290 studies associating polymorphisms in methylene tetrahydrofolate reductase (*MTHFR*) with various disease or physiological conditions have been published (<http://www.cdc.gov/genomics/search.htm> → *MTHFR*). The most studied variants are c.677C>T (p.A222V) and c.1298A>C (p.E429A). Marini et al. [2008] recently resequenced 564 individuals of diverse genetic ancestry (Coriell Institute panels: <http://cct.coriell.org/Sections/BrowseCatalog/Populations.aspx?PgId=4>) and discovered 14 nonsynonymous changes including 11 alleles with frequencies <1% along with the common alleles p.A222V, p.E429A, and p.R594Q [Marini et al., 2008]. Increased levels of folate restored *MTHFR* activity to the normal range in four of the five variants. The sequence heterogeneity and remediation of enzyme activity by folate supports a greater emphasis on the ~600 cofactor-dependent enzymes in the human proteome. Because many cofactors are derived from diet, such studies may identify individuals who require higher concentrations of vitamins for optimal health. Analyses of populations using HapMap data or their derivative reagents also provide justification for the need to resequence genes in diverse populations.
- Published HapMap data analyzed by a novel algorithm identified chromosomal regions with a high *F<sub>st</sub>* (Fixation index, a measure of population differentiation) between three ancestral populations (European, Chinese, and African) [Myles et al., 2008b]. These regions encoded genes involved in carbohydrate metabolism, skeletal development, and pigmentation. Such allele frequency differences may explain, for example, the differential effect in incidence of obesity and type 2 diabetes between Europeans and Pima Indians who consume similar Western diets [Schulz et al., 2006].
- Twenty-five single nucleotide polymorphisms (SNPs) linked to six chronic diseases in genome-wide association studies (GWAS) were analyzed in ~1,000 individuals from 53 populations [Myles et al., 2008a]. Several risk alleles were absent from some populations and several were present at 100% frequency, indicating that the allele may contribute uniquely to disease in the European population. Other polymorphisms in these genes or in other genes within the non-European populations are likely to contribute to disease incidence and severity.
- Allele frequencies of 873 tag SNPs in 82 candidate genes involved in hypertension, type 2 diabetes, obesity, dyslipidemia, or metabolic syndrome varied in 54 populations [Hancock et al., 2008]. Some of these genes were linked to adaptation to cold climates and others likely to selection by other environmental factors. Exposure to infectious agents such as malaria

[Kwiatkowski, 2005], high altitude [Beall, 2007] and the food environment [Perry et al., 2007] have also been shown to select for certain genotypes.

- Culture is also known to influence allele frequencies among populations. A gradient of genotypes exists from northwest to southeast in Europe [Price et al., 2008]. At a fine mapping level, allele frequencies mirror geography, and by inference, national culture, within Europe [Lao et al., 2008; Novembre et al., 2008].

Variation at the genomic and gene levels demonstrates that existing data and reagents will not be sufficient to identify genes involved in maintaining health or those that contribute to the incidence and severity of disease. The newly initiated 1000 Genomes project (<http://www.1000genomes.org>), which is being organized by an international consortium, will employ genome-wide resequencing and targeted coding region sequencing in a total of approximately 1,500 individuals from three human populations: Europeans, Africans, and Asians [Lang, 2008; Siva, 2008]. Each of these populations will be represented by a number of subpopulations consisting of approximately 100 individuals including Yoruba in Ibadan, Nigeria; Japanese in Tokyo; Chinese in Beijing; Utah residents with ancestry from northern and western Europe; Toscani in Italy; and other populations to be determined. The goal of this international effort is to characterize alleles with frequencies of approximately 1% genome-wide and less than 1% in coding regions. The phenotype of individuals sequenced in the project will not be analyzed. Hence, these data will be used for improving selection of reagents and designs for GWAS, and will not be directly focused on identification of disease causing genes. Further information is available from the project website referenced above.

## The Human Variome Project

The Human Variome Project differs from these other efforts in developing and fostering an international effort to systematically identify genes, their mutations, and their variants associated with phenotypic variability and indications of human disease. The HVP is an international effort linking clinical, medical, and research

laboratories for developing knowledge housed within linked databases. This knowledge would be accessible to the research and medical communities to improve research strategies and clinical medical practice. The key objectives of the project are described in Box 1. An example of the need for the HVP as applied to neurological disorders has recently been published [Cotton et al., 2008].

The HVP Planning Meeting conducted concurrent meetings that discussed (1) classifying genetic variation from unlinked clinical medicine or research laboratories, (2) capturing data from diagnostic and service laboratories, (3) assessment of pathogenicity, (4) data transfer, (5) data integration access, (6) funding and governance, (7) emerging countries' initiative and involvement, (8) ethics: existing and anticipated concerns, (9) attribution and publication, and (10) pilot projects (Table 1). Reports for classifying and capturing genetic variation from laboratories (i.e., committee reports 1 and 2) have been combined for sake of brevity. Detailed reports for these sessions are provided in the online Supporting Information for this article. Below are synopses of the main outcomes and recommendations.

### Classifying Genetic Variation From Unlinked Clinical Medicine, Research, or Service Laboratories

The cyber infrastructure for biological data is extensive [Stein, 2008], but still not fully integrated or developed. The HVP is relying upon these databases for data element definition, storage, management, retrieval, and nomenclature. For example, the NCBI provides a gene-centric index for mutation nomenclature, Human Genome Organization (HUGO; <http://www.hugo-international.org/>) has a naming scheme for genes, and the Human Genome Variation Society (HGVS; <http://www.hgvs.org/>) [Cotton and Horaitis, 2000] provides for naming of mutations. The Cancer Bioinformatics Grid (caBIG - <https://cabig.nci.nih.gov/>) [Fenstermacher et al., 2005] exemplifies the need for interoperability, common languages, data standards, and sharing. In addition to the disease-specific (e.g., caBIG) and international databases, LSDBs in individual laboratories and institutes exist but are not easily linked to the rest of the bioinformatic community.

Although the data management infrastructure continues to expand, curated genetic data are scattered: no coordinated effort

#### Box 1: Key objectives of the HVP

1. Capture and archive all human gene variation associated with human disease in a central location with mirror sites in other countries. Data governance will ensure security and integrity through the use of auditing and security technologies but, nevertheless, allow searching across all genes using a common interface.
2. Provide a standardized system of gene variation nomenclature, reference sequences, and support systems that will enable diagnostic laboratories to use and contribute to total human variation knowledge.
3. Establish systems that ensure adequate curation of human variation knowledge from gene-specific (locus-specific), country-specific, or disease-specific database perspective to improve accuracy, reduce errors, and develop a comprehensive data set comprising all human genes.
4. Facilitate the development of software to collect and exchange human variation data in a federation of gene-specific (locus-specific), country-specific, disease-specific, and general databases.
5. Establish a structured and tiered mechanism that clinicians can use to determine the health outcomes associated with genetic variation. This will work as a dialogue between those who use human variation data and those who provide them. Clinicians will be encouraged to provide data and will have open access to complete variation data.
6. Create a support system for research laboratories that provides for the collection of genotypic and phenotypic data together using the defined reference sequence in a free, unrestricted and open access system and create a simple mechanism for logging discoveries.
7. Develop ethical standards to ensure open access to all human variation data that are to be used for global public good and address the needs of "indigenous" communities under threat of dilution in emerging countries.
8. Provide support to developing countries to build capacity and to fully participate in the collection, analysis and sharing of genetic variation information.
9. Establish a communication and education program to collect and spread knowledge related to human variation knowledge to all countries of the world.
10. Continue to carry out research within the opportunities presented by investigation of human genetic variation and to present these findings to users of this information for the benefit of all.

exists to harness and harmonize these efforts, data, and knowledge. Gene mutation and variation data are generated from and used by diagnostic, epidemiological experiments, research laboratories, and clinicians, each of which has different missions, ability, or willingness to curate information, and resources (Table 2).

Because clinical laboratories are not required or encouraged to deposit genotype or phenotype data into publicly available databases, data sharing ranges from complete to none. Efforts to encourage and develop ongoing data collection have begun and range from commercial enterprises to funded grant programs (Table 3). Each of these programs is developing unique solutions for the barriers of time, cost, concerns of patient confidentiality, institutional review board requirements, maintenance of quality assurance, and difficulties in obtaining clinical information from referring centers. These challenges reinforce the concept that clinical laboratories should not be expected to develop and curate public databases. However, clinical laboratories should be expected to contribute data. Developing a standard open software suite such as the Leiden Open Variation Database (LOVD; <http://www.lovd.nl>) [Fokkema et al., 2005] for these initiatives will allow existing tools, for example, the Universal Mutation Database (UMD; <http://www.umd.be/>) [Beroud et al., 2000] software, to query across the cyberspace of LSDBs to retrieve and analyze data (Table 3). Relying upon a common database design, language, and interoperability will enforce quality standards across clinical and research laboratories. Nevertheless, clinical and research laboratories may have processes and quality measures that would require a “data warning” for select entries or datasets. Some of the issues and requirements to initiate adoption of these ideas are described in Box 2. Once these individual LSDBs are developed and curated locally, ethically appropriate data elements can be deposited in national or international databases (NCBI or EBI).

**Table 1. Organization of Committee Reports<sup>a</sup>**

Committees Meeting in Spain	Box <sup>b</sup>
Capturing and classifying genetic variation	
Unlinked clinical and research laboratories	2
Diagnostic and service laboratories	2
Assessment of pathogenicity	
Data transfer	4, 5
Data Integration access	4, 5
Emerging countries	6
Ethics	7
Funding and governance	8, 9
Attribution and Publication	
Pilot Projects, including InSiGHT	10, 11

Box 1 describes Key Objectives of the HVP. Box 3 outlines Standards, Validation, Quantification, Transparency for data.

<sup>a</sup>See text for details.

<sup>b</sup>Outlines of each committee report are available in the online Supporting Information.

**Table 2. Laboratories Involved in Mutation/Variation Discovery and Use**

Laboratory	Mission	Quality Assurance <sup>a</sup>	Curation <sup>b</sup>	Database	Other
Diagnostic or Reference	Service, no research	In place	None	Local, if any	Cost pressure
Epidemiology or Public Health	Associate existing variations to phenotype	Provider dependent	Population	Local and global	No new variant data generation
Research	Generate new knowledge	None	Individual	Local	Continuity risk

<sup>a</sup>Quality assurance is required in diagnostic or reference laboratory. QA in epidemiology studies would depend on the service provider.

<sup>b</sup>Curation would be at the population level in epidemiological studies and the individual level in research laboratories. Not all research laboratories would curate variants or mutants.

## Pathogenicity and Clinical Utility

Understanding the consequences of genetic variation depends upon the simultaneous collection and documentation of phenotypic data for each variant (e.g., [Cotton et al., 2007a; Crawford and Nickerson, 2005; Kaput, 2008; Kaput et al., 2005; Ring et al., 2006; Taylor et al., 2001]). The correlation between genome and phenotype (pathogenicity) is the basis for the clinical benefit. The two broad principles for assessing pathogenicity or phenotype linked to a genetic variant are that (1) multiple data elements must be integrated and (2) data elements and the integration process must have standards, validation, quantification, and transparency (Box 3).

The omics sciences are now capable of generating large but disparate (e.g., genomic vs. metabolomic) datasets that may be used in research but also clinical applications. Although assessing pathogenicity will be an ongoing, iterative process, several specific recommendations are warranted.

### Genetic and genomic data

Gene marker analysis is an important step in the clinical diagnosis of pediatric and adult genetic disorders. The issues associated with clinical genetic testing are well recognized for inherited cancer syndromes, where missense variants represent 10–30% of test results (e.g., [Eisinger, 2008; Metcalfe et al., 2008; Stoffel et al., 2008]). Many of these variants are classified as having an uncertain effect unless strong genetic epidemiologic and/or functional evidence exists associating them with the syndrome. The Breast Cancer Information Consortium (BIC; <http://research.nhgri.nih.gov/bic/>) classifies *BRCA1* and *BRCA2* variants as pathogenic only if the probability of pathogenicity, (usually based on statistical genetic approaches) is, definitely pathogenic >0.99; likely pathogenic 0.95–0.99; uncertain 0.05–0.949; likely not pathogenic/little clinical significance (LCS) 0.001–0.049; neutral or LCS <0.001 [Plon et al., 2008]. For other genes, there are no standards. For each type of data, old and new, the principles of standards, validation, quantification, and transparency apply. Many of the expected gene variants have not been identified or characterized in world populations. Hence, complete resequencing of a gene proven or suspected to be involved in monogenic and polygenic diseases will be required to determine causal linkages between genes and phenotype.

### Standardizing existing clinical phenotype and pathology

A fundamental problem with assessing phenotypes is the diversity of the underlying molecular pathways that cause disease, and as a consequence the heterogeneity in clinical manifestations, age of onset, severity, complications, and age of death. Other groups [Kaput et al., 2005; Kathiresan et al., 2008;

**Table 3. Examples of Clinical–Research Translational Initiatives with Database Suites**

Program	Name	Agency	Mission	Method	Website
ARUP	Clinical laboratory	University of Utah (owner)	Galactosemia	Locus specific	<a href="http://www.aruplab.com/">http://www.aruplab.com/</a>
CETT	Collaboration, education, test translation	NIH Office of Rare Diseases	Collaboration between clinical laboratory, research, and advocate group	Clinical laboratories collect phenotypic data from clinicians through online forms prior to genetic testing. Genotypic and phenotypic data to be placed in a database	<a href="http://www.cettprogram.org/">http://www.cettprogram.org/</a>
DMuDB	Diagnostic Mutation Database	Certus Technology Ltd, UK National Genetics Reference Laboratory	Mutation database for diagnostic genetics services	~5,000 entries including 2,900 BRCA1 and BRCA2 entries from 12 laboratories (851 unique variants, 436 of which are not in BIC), HNPCC (827 entries from 9 labs), and the following unique variants: 177 MLH1, 178 MSH2, 43 MSH6, 116 APC, 70 Lowe Syndrome, 38 Sotos Syndrome, 150 Cystic Fibrosis, 96 Neurofibromatosis, and 101 X-Linked Retinitis Pigmentosa	<a href="http://secure.dmu-db.net/">http://secure.dmu-db.net/</a>
eyeGENE	National Ophthalmic Disease Genotyping Network	NIH National Eye Institute	The discovery of the genetic causes of ocular diseases; accurate diagnostic genotyping to patients with inherited eye diseases; patient information; develop populations for study; develop resources; improve phenotype descriptions	Patient meeting minimal clinical criteria sign research and a clinical testing consent forms, are provided counseling. Clinician enters the clinical information into the online data fields that populate the eyeGENE database, while the clinical laboratory enters the genotype information. The database is not public and is password protected for privacy	<a href="http://www.nei.nih.gov/resources/eyegene.asp">http://www.nei.nih.gov/resources/eyegene.asp</a>
Idbases	Immunodeficiency Database	Tampere University Hospital	Establish database for every immunodeficiency or provide links to those maintained elsewhere.	Maintain 122 public IDbases with 3 under construction (2008). These databases contain data for 5359 patients, 27 Immunodeficiency mutation databases are maintained by others	<a href="http://bioinf.uta.fi/base_root/">http://bioinf.uta.fi/base_root/</a>
NCCRCG	National Coordinating Center for the Genetics and Newborn Screening Regional Collaborative Groups	American College of Medical Genetics	Develop standardized laboratory and clinical language for electronic medical records that will be utilized in newborn screening data collection	Genotype and phenotype information, with long-term follow-up of newborn screening patients, will be collected on the 54 conditions with over 150 genes	<a href="http://www.nccrcg.org/">http://www.nccrcg.org/</a>
LOVD 2.0	Leiden Open Variation Database LAMP = Linux, Apache, MySQL, PHP	Leiden University Medical Centre	Software and database to establish, manage and curate a Web-based gene specific DNA variant database (LSDB)	Sequence variant data and patient data stored separately—variations in several genes can be linked to one patient. The pathogenicity of each variant can be set by submitter and curator, allowing one to easily separate disease causing from nonpathogenic variants.	<a href="http://www.lovd.nl">http://www.lovd.nl</a>
UMD	Universal Mutation Database <sup>h</sup>	INSERM—Montpellier	Creation of relational databases that run on both Macintosh and PC platforms and can create dynamic HTML pages for direct on-line queries via the Web.	The UMD software uses the 4th Dimension <sup>h</sup> package from 4D. The UMD central tool can query multiple Locus Specific DataBases (LSDB) developed with the Universal Mutation Database <sup>h</sup> software	<a href="http://www.umd.be/">http://www.umd.be/</a>

Makinen et al., 2008; Rosenzweig et al., 2002; Wong, 2006; Zaninotto et al., 2007] have proposed using disease as a classifier (e.g., type 2 diabetes), but rely on quantitative measures of phenotype (e.g., fasting glucose, fasting insulin) as a means to reduce subjective assignments of disease [Tracy, 2008]. Because the HVP seeks to collect data from laboratories and clinics, phenotype templates are needed to define ranges of (1) minimum sets of clinical data, (2) range of subset data, and (3) maximum datasets. Such hierarchical template structures will allow scientists in all countries to participate in data and sample collection. Developing the LOVD/UMD tools (Table 3, and below) for the HVP will also require a means to validate data and

data quality prior to implementation across laboratories. Clinical and pathology data standards must be developed by experts in each genetic disorder for interpreting the effects of genetic variation.

#### *Standards linking in vitro functional studies with clinical results*

If a cellular function can be established that appears to correlate with the clinical syndrome, then in vitro assays could be used to classify whether a variant retains or loses function. Standards for performing and interpreting the assays are crucial if these methods are to be accepted as a mechanism for classifying variants clinically

### Box 2: Developing interoperable clinical and research locus specific databases

1. Security, consent, and ethical approval
2. Develop standards and guidelines for diagnostic labs to address database contribution
3. Develop quality standards for genotype and interpretation
4. Develop recommendations for formats for submission of data
5. Provide models for collecting clinical information
6. Provide recommendation for funding for database efforts, including bioinformatics
7. Scalability: can small databases readily handle the output from next gen sequencers?
8. Laboratory pipeline from instruments to database (LIMS: Database interface).

### Box 3: Standards, validation, quantification, and transparency

1. **Standards**—How data (including clinical, demographic, pathologic, genetic, lifestyle and environment, computational, and laboratory) are collected, recorded, and utilized should be standardized so that results are as unambiguous as possible to users.
2. **Validation**—Before conclusions are drawn, new types of data should be compared with validated sets of variants that are known to be (or are highly likely to be) pathogenic and neutral, respectively.
3. **Quantification**—Uncertainty exists around all of the data types described above. That uncertainty should be quantified as best as possible. When properly quantified, results from all data can be combined to produce a final probability that a variant is pathogenic.
4. **Transparency**—Uncertainty can best be addressed when the source and quality of all data are known. Understanding genetic variation will be an ongoing, iterative process. Proper reporting of a conclusion regarding pathogenicity will require disclosure of all data sources, their quality, the methods of combining them, and the date of the conclusion.

[Couch et al., 2008]. In cancer genetics, the correlation of mismatch repair defects with Lynch syndrome (Hereditary Nonpolyposis Colorectal Cancer) is probably the most well-established example [Ou et al., 2007]; even so, the principle that multiple data elements must be integrated to achieve classification should be respected.

#### Computational studies

Multiple studies in recent years have confirmed the value of comparative sequence analysis in helping to predict whether a missense variant is pathogenic or not (reviewed in [Tavtigian et al., 2008]). However, the issues of standards, validation, and transparency also apply to computational methods. Most importantly, the quality of a multiple sequence alignment is critical to their accuracy [Ahola et al., 2006, 2008]. The choice of ortholog sequences that is used, the quality of cDNA or genomic data, and the methods used to construct the alignment are all important features.

#### Computational studies—predictive algorithms

Some methods have already been validated on curated data sets of variants, establishing their Negative Predictive Value (NPV, the proportion of predictions of “neutral” that are actually neutral) and Positive Predictive Value (PPV, the proportion of predictions of “pathogenic” that are actually pathogenic) [Chan et al., 2007; Chao et al., 2008]. However, several of the more commonly used algorithms have been updated. Algorithms exist for coding region variants and predictions of altered splice sites [Nalla and Rogan, 2005; Spurdle et al., 2008]. New methods, including both rule-based and machine-learning approaches, are being developed [Tavtigian et al., 2008], and in the future, algorithms to assess other noncoding functions of DNA are anticipated.

#### Data Transfer and Databasing—Gene and Locus Specific Databases

Historically, gene variation data were first collected for specific gene(s) causing a Mendelian disorder or a change in the phenotype.

These listings usually were initiated and driven by the interests of an expert using the collection for research, clinical or diagnostic applications. Currently there are over 700 such LSDBs (<http://www.hgvs.org/dblist/glsdb.html>) [Horaitis et al., 2007], mostly Web accessible. Complete collection and expert curation of gene sequence variants and their coupling to phenotypic consequences (if any), will be essential for proper future healthcare and research. Data Transfer and Databasing plans are outlined in Box 4.

#### Integrating Data and Providing Access

The breadth and depth of information available about human variation are rapidly expanding as new technologies (e.g., omics and imaging) analyze health and pathogenicities. Uncomplicated methods of access are needed for multiple user communities with differing expertise in genetics, clinical medicine, nutrition, physiology, and probably the public. Information generated by the HVP will have many of those dimensions, ranging from how variants are identified, the type of variant, the physiological parameters associated with the variant, and where and how records are maintained and accessed. For example, information for rare variants, SNPs, microsatellites, small insertions/deletions will range from reports of one variant in one individual, such as might be gathered from manual data entry and maintained in medical records, to large-scale screening of thousands of genomes in defined populations. How the data are represented will be a challenge because information will come from published literature or text-rich resources such as OMIM or GeneReviews (<http://www.geneclinics.org/profiles/all.html>) [Pagon et al., 2002], include explicit records in LSDBs and/or genome-wide resources such as dbSNP or genome browsers. The scope of the data, which includes locations of variants on reference sequences to phenotypes in humans and model organisms, increases the dimensionality. Hence, the scale of information for human genetic variation and linked phenotype can range from a single text document to petabytes of raw data derived from sequencing thousands of genomes to a high level of coverage and accuracy. Model databases capable of accessing known variants will be developed from pilot projects and provide a resource for clinicians, patient advocates, and the public (i.e., education).

#### Box 4: Data transfer and databasing standards for LSDB

1. Develop a program to recruit curators of LSDBs who have yet to join the HVP.
2. LSDBs will be standardized and use freely available and compatible software packages such as Universal Media Disc (UMD) or Leiden Open Source Variation Database (LOVD).
3. Develop standard formats for published data by disseminating rules for LSDB submission to journals. A new standard of publication should call for as much phenotypic data as possible. Similarly, funding agencies should require submission of all gene variant data resulting from their support.
4. LSDBs should develop and follow quality standards that are measureable and transparent. The HGVS should develop a quality label for LSDBs following all standards.
5. To ensure permanent access, LSDBs should share a copy of their data to central repositories (e.g., HGVS, NCBI, EBI, Gen2Phen). In addition, to improve easy access, LSDBs should consider automated data exchanges in real time with continuously funded national databases (NCBI, EBI).
6. LSDBs can include data from expression or other in vitro functional studies, but these data should be clearly marked as in vitro.
7. New LSDBs are needed for every gene related to a Mendelian disorder. Lost or abandoned LSDBs should be revived. HGVS and the Gen2Phen consortium could play an active role here, soliciting new curators as well as finding financial support.
8. Duplication of effort should be avoided, but if it occurs, LSDBs with common genes and variants should share data and entries.
9. An internationally acceptable database policy statement, covering ethical and privacy issues for gene variant and phenotype data collection and sharing, is urgently needed.

#### Accessing known variants

Given the dimensionality issues and the challenges they represent, the HVP will help develop an infrastructure to identify variants relative to a reference standard and allow facile linking of data with appropriate tools. Notification schemes would be developed to indicate missing data for each variant. The HVP is identifying a progression of doable tasks with milestones for each. For example, the HVP will link with one such effort, the on-going collaboration among GEN2PHEN (<http://www.gen2phen.org>), EBI, and NCBI to develop an international set of standard Locus Specific Genomic Sequence reference (LRGS/RefSeqGene: [http://www.gen2phen.org/docs/LRG\\_Specification\\_Summary\\_version\\_9.pdf](http://www.gen2phen.org/docs/LRG_Specification_Summary_version_9.pdf)). With that foundation in place, active LSDBs can report their variants in LRG/RefSeqGene coordinates to centralized databases (EBI/NCBI) to be accessioned to dbSNP (<http://www.ncbi.nlm.nih.gov/projects/SNP/>) [Sherry et al., 2001]. The scope of information to be centralized requires further debate, but could include items described in Box 5. The HVP will ensure that other standards in addition to the reference sequence are developed and/or used, such as human gene nomenclature (HUGO Gene Nomenclature Committee (HGNC); <http://www.genenames.org>) [Bruford et al., 2008], variant nomenclature (Human Genome Variation Society; <http://www.hgvs.org/mutnomen>) [den Dunnen and Antonarakis, 2000], variant accessions (dbSNP), and names for diseases or diagnostic tests in medical records (Logical Observation Identifier Names and Codes, LOINC) (e.g., [McDonald et al., 2003]). Reference gene-specific standards to support LSDBs and genetic testing groups (see <http://www.gen2phen.org> and <http://www.ncbi.nlm.nih.gov/RefSeq/RSG>) are under development. In addition, assigning accessions in dbSNP to known variants, either directly to dbSNP or Ensembl (<http://www.ensembl.org>), or with PhenCode (<http://phencode.bx.psu.edu>) [Giardine et al., 2007] serving as a data collection center are being implemented.

#### Public education

With the understanding that large data sets are only as useful as their ease of access, the HVP can also foster portal sites to help direct users to tools and resources of interest, and identify areas requiring additional development. Tutorials comparing and contrasting different resources might also be commissioned. Topics that must be addressed are how to access information in the published literature, the effect of variation on transcription, the significance of conserved noncoding regions on phenotypic

#### Box 5: Data Elements for LSDB

1. HGVS name and historical name(s) of the variant
2. Data submitter (LSDB or individual researcher)
3. URL back to data submitter per variant
4. PubMed UIDs associated with the variant
5. Number of observations of a variant
6. MIM number of the disease(s) or phenotype(s) of interest
7. MIM number of the allele of interest
8. Clinical interpretation
9. Clinical or research groups testing for this variant

variation, identification of mRNA splice sites, descriptions on the structure and maintenance of LSDBs, and strategies for interpreting the impact of single (e.g., SNPs) and multiple (e.g., haplotypes) variants in one allele or in multiple alleles on the phenotype of interest. The HVP intends to make greater efforts to educate the public about the project's importance and the benefits of research participation.

#### Attribution and Publication

A major barrier to the development of comprehensive analyses of human genetic variation, and therefore LSDBs, is the reward system for clinicians and academic researchers. Nonacademic clinicians typically receive little or no credit for contributing to the scientific literature, and academic clinicians and basic researchers cannot easily persuade journals to publish the 50th variant of a gene that has an observable effect on phenotype. These cultures may be difficult to change overnight, but specific steps could be instituted immediately to promote the submission of genotype-phenotype data to LSDBs and to "reward" contributions to team projects. Database entries could be a mandatory quality control standard for clinical laboratories and clinicians. For researchers, a publication or Web-based system establishing microattribution and community annotation of mutations (e.g., <http://www.wiki.genes.org>) and cited data will enable measurable contributions to the scientific knowledge base [Axton, 2008; Hoffmann, 2008]. Similarly, database journals would also serve this task by providing a forum for publishing gene variation data that would be eventually deposited in the PubMed literature database. One such journal, the open access *Human Genomics and Proteomics* (HGP; <http://www.sage-hindawi.com/journals/hgp>), which is affiliated with FINDbase, will focus on studies characterizing causative

mutation and/or biomarker frequency spectra. Accepted contributions including datasets will be linked in FINDbase and deposited in PubMed [Patrinos and Petricoin, 2009]. Journals, tenure, and promotion committees, and funding agencies would be encouraged to cite these contributions and citations of LSDB and international national databases [Patrinos and Brookes, 2005]. The HVP recommends that researchers cite these attributions and citations in their *curricula vitae* to foster the transition of the academic culture. The same trend and recommendations for development of coherent tools are valid for the recognition of contribution to the setting, use and sharing of any bioresource such as biobanks [Cambon-Thomsen, 2003; Kauffmann and Cambon-Thomsen, 2008] and international efforts like P3G are being developed in the same spirit [Knoppers et al., 2008].

### **Developing/Emerging Countries—Ensuring a Worldwide Collection**

Although ~90% of known SNPs are shared between Asians, Europeans, and Africans, 80% of private SNPs are found within Asian and Africans [Hinds et al., 2005; Jorde and Wooding, 2004]. The recent sequencing of the Watson [Wheeler et al., 2008], Venter [Levy et al., 2007], Kriek (<http://www.sciencedaily.com/releases/2008/05/080526155300.htm>), West African [Bentley et al., 2008], and Han Chinese [Wang et al., 2008] genomes, along with gene-specific resequencing efforts (see Introduction), suggest that a large number of SNPs and other sequence variation exists in the human population. Estimates from African genetic diversity and the Pan Asian SNP initiative suggest that 80–90% of human genomic variation resides in the world's emerging countries. Any formal attempt to identify the extent of genomic variation must include geographical regions that have not been included in haplotype mapping projects. Although the Population Reference Sample (POPRES) will address some of these missing populations [Nelson et al., 2008], this effort is designed as a mapping project, not one focused on functional polymorphisms or mutations. Hence, the main focus of the HVP effort is the inclusion and analyses of clinical samples from diverse ethnic groups.

The distinct advantage of some ethnic populations is the opportunity to study genetic diseases due to consanguinity, large family size, and potential founder effects (e.g., [Bittles, 2001, 2002; Saadallah and Rashed, 2007]). Emerging nations will be regarded as major contributors to the VARIOME project. However, biomedical research has not been the focus of resource poor countries even though such activities are likely to produce economic and health benefits for all [Daar et al., 2002; Singer and Daar, 2001]. Education of healthcare providers, the public, and government officials is needed for demonstrating the universal nature of the HVP, the need to include populations in developing countries, and the benefits from cooperating in biomedical research [Bhan et al., 2007; Cohen et al., 2008; Seguin et al., 2008; Tindana et al., 2007]). Certain populations may mistrust research involving genetic analyses or fear that results can be misused to support discrimination or worse (<http://www.eubios.info/ASIAE/BIAE201.htm>): Malay-Muslims, Chinese, and Indians in Singapore expressed anxiety about breach of confidentiality, the misuse of their DNA for cloning, and possibilities of being diagnosed with disease [Wong et al., 2004]. Community-based participatory research collaborations may provide forums for addressing cultural and ethical concerns of biomedical research [McCabe-Sellers et al., 2008].

Analyzing the extent of human genomic variation creates an ideal opportunity for the developed and the developing nations of the world to forge meaningful partnerships and to work together in an unprecedented way, initially to identify variation causing

disease, and then to understand how general variation contributes to human phenotypic diversity. By ensuring that all nations and ethnic groups have an equal and fair opportunity to share data and technology, we will provide evidence-based information that all populations can benefit from a global society health network. The primary objectives for including populations in emerging countries are described in Box 6.

The HVP appreciates the genomic sovereignty/equality for all countries to be involved in the Human Variome Project and acknowledges the value of “human capital” within all populations. Real and tangible benefits of the HVP to improve health will be generated for participating populations; the voluntary participation of the greatest number of countries would ensure a more general applicability, and it is hoped that many countries will decide to participate.

### **Progress in Developing Ethical Guidelines for LSDBs: Principles to Practice to Implementation**

Ethical issues remain of vital concern to the Human Variome Project. Participating researchers are committed to adhering to the highest ethical principles governing research, data sharing, and ultimately enabling this new knowledge to benefit all of humanity as much as possible. Ethical guidelines specifically for LSDBs were previously published [Cotton et al., 2005].

LSDBs may contain a large amount of phenotypic data. Most LSDBs post a considerable amount of data on public Web sites and increasingly this information may be accessible through genome browsers. Although the best intention of the HVP is to ensure that participants are acknowledged as a group, without any risk of identification, a specific challenge occurs in the case of rare mutations associated with distinctive clinical features. Because epistatic and environment interactions (reviewed in [Kaput, 2008]) alter age of onset, severity, complications and outcomes for monogenic and polygenic phenotypes, it may be necessary to analyze entire genomes for personalized healthcare. Such polygenic analyses generate data that could be used for reidentifying individual patients [Craig et al., 2008].

Other ethical concerns may be minimized by improving communication about the project and its goals through multiple channels such as print and broadcast media, local community outreach, and internet sources such as the HGVS Web site. The HVP will develop an ethics review committee with a subcommittee focused on issues related to LSDB for (1) providing counsel when dilemmas arise, (2) overseeing guidelines, (3) identifying best practices, (4) determining how best to ensure privacy in all situations, (5) formulating how to handle data for which explicit consent does not exist or is not possible to achieve, and (6) developing a consent form that is consistent for all LSDBs but which can be adapted to the requirements of individual countries. Such consent would contain, for example, a recontact clause. The specific recommendation and open questions are outlined in Box 7.

The HVP, through its Ethics Working Group, is committed to (1) soliciting, collecting and analyzing consent forms in order to develop a model consent form that can provide greater consistency across all LSDBs, (2) seeking the input of relevant clinical genetic societies for comment, and (3) using that input to develop ethical standards for LSDBs.

### **Funding Mechanisms and Governance**

Funding for collecting data of mutations causing single gene disorders has traditionally been difficult due to the extreme



### Box 6: Emerging populations: objectives for inclusion

1. Develop procedures to include neglected, underresourced clinicians/researchers in emerging countries, by creating networks among those of common interest.
2. Recognize the need and challenges for recruiting participants. Considerable resources are needed to ensure an adequate consenting process as well as data collection, both of which can be cumbersome.
3. Ensure that the databases and software suite of tools that are being used or developed can be utilized by those in an undercapacitated environment, and that the software pipelines transfer genotype data easily between clinic and laboratory (e.g., Locus Specific Databases [LSDB]).
4. Researchers and their academic institutions will develop an infrastructure in anticipation of the Coordinating Centre working with international organizations (UN, UNICEF, WHO, OECD) to create a framework to facilitate interactions with national governments and regional organizations such as NEPAD (New Programme for Africa's Development) downward through national governments, their Departments of Health, Education, Science and Technology, among others.
5. To ensure that ethical, legal, religious, and social issues of emerging countries are considered in all portals of the HVP. The HVP will enlist the regulatory/policy arm of WHO/UN to acknowledge and deal with issues pertinent to Economic, Educational, Ethical, Legal and Social Issues.
6. To foster communication, interaction, and research networks between the developing and the developed countries.
7. To keep abreast of new analytical methods that may impact or influence the ability to maintain confidentiality of data.

fragmentation of the field even though mutations affect 60% of all individuals in a lifetime [Baird et al., 1988]. The funding possibilities are more likely if the international HVP is treated as a concerted effort. Given the limitations of existing knowledge (see Introduction and [Patrinos and Brookes, 2005]), this initiative will benefit research in many fields and impact prevention and clinical care of disease. Specific focus areas for developing funding streams are described in Box 8.

#### Governance

Those dedicated to assisting themselves and others in their clinics by collecting mutation/variants causing inherited disease have in the past acted in isolation; reinventing wheels wastes funds and time. The HGVS, formerly HUGO-MDI, was formed to alleviate this problem and accelerate the collection and management of information on mutation causing disease(s). The HVP was named and initiated to define the aims of this activity (Melbourne in 2006) by an extremely high-profile group of experts in all types of genetic variation analyses. The Genomic Disorders Research Center (Melbourne, Australia; <http://www.genomic.unimelb.edu.au>) was voted as coordinating office, to continue its function of facilitating independent projects started at the inception of the center in 1996. In association with world experts, Deloitte (<http://www.deloitte.com>) developed a business plan, which was approved by the HVP Planning Group. The HVP Planning Meeting was a designated activity of this plan and other sections of the plan will follow (e.g., board function to oversee the coordinating office's function and support it). The business plan calls for a broadly defined community of interested stakeholders to develop the HVP (Box 9).

Open questions that must be resolved in future meetings are (1) the extent of data sharing between patient records and research databases, (2) appropriate descriptions of data elements, and (3) data ownership and confidentiality.

#### Pilot Projects

The HVP has established a partnership with the International Society for Gastrointestinal Hereditary Tumors (InSiGHT; <http://www.insight-group.org>) to collect and classify a large set of missense variants associated with hereditary colorectal cancer. This effort will develop HVP's prototype system for interpreting variants observed in clinical genetic testing (Box 10). InSiGHT consists of a multidisciplinary scientists focused particularly on

the Mendelian disorders predisposing to colorectal cancer (Familial Adenomatous Polyposis, Lynch Syndrome, and MUTYH related polyposes). This effort is an ongoing project but also a model or pilot for the HVP. InSiGHT has conducted several multidisciplinary studies of Hereditary Non Polyposis Colon Cancer (HNPCC) patients that (1) require the development of a disease-specific model for integrating databases across laboratories, (2) establish standards for data consistency for phenotypes (which include graphic pedigrees), (3) address confidentiality, and (4) develop a template for consent. Clinicians from multiple countries and regions are contributing and committing to the development of these systems. Some of these issues cross disciplinary boundaries and are being addressed by other committees of the HVP. The InSiGHT consortium's roadmap includes providing access to clinicians, a vital resource that will serve patients in the immediate future, and as a model for other genes and phenotypes. Among the first efforts involved uploading large datasets of mismatch repair variants generated by national consortia and laboratories into the InSiGHT MMR, an LSDB that uses the LOVD platform. Data transfer into a mirrored central database (e.g., NCBI or EBI) is also planned with an initial reciprocal agreement with the Health Data Integration project at the Australian CSIRO Centre (<http://www.ict.csiro.au/HAIL/Abstracts/2004/UmaSrinivasan.htm>). Other pilot projects are described in Box 11.

#### Discussion

The vision of the HVP, to catalog and access all information related to human disease variation, is ambitious. One can conceptualize the challenge as a multidimensional, fluid matrix, with all ~20,000 genes as column headings and rows of potentially thousands of variants as descriptors. In addition, third and further dimensions would annotate other biological parameters, for example, clinical and/or metabolic phenotype, microarray expression, proteomics, protein interactions, nutrient intakes, physical activity, and other functional phenotypic and epidemiological data. Separate dimensions that must be linked to these variants are the main effect of gene-environment interactions (e.g., [Lim et al., 2007]). These data elements relate to the cells in the initial two dimensional matrix because each may affect the genetic expression of the mutation or gene variant. These dimensions are domains of knowledge that must be integrated for understanding biological processes.

A predetermined bioinformatics structure to accommodate this matrix with forced fields for data entry is notionally appealing but

### Box 7: Ethics subgroup recommendations and remaining issues

1. Develop a Common Ethical Framework that governs all aspects of the HVP. This guideline will ensure strict privacy protection for all participants and provide guidance for optimizing informed consent. It will strive to ensure wide access and information sharing, global benefit sharing, while preserving individual country norms, guidelines, and legal requirements. It will establish a clear basis on which to protect individuals, families, and communities from harm.
2. Curation and all communications about work involved in a particular LSDB must be transparent and clear as to the nature, purpose and limits of that particular LSDB.
3. Different types of data must be handled in ethically appropriate manners, including prospective, published, and archived data. All data must, furthermore, be handled according to the conditions under which consent was originally obtained, and privacy strictly maintained (see footnote).
4. Ideally, donors should be informed when their data is being transmitted to a LSDB and to the Web. However, difficulties in achieving this goal arise. Is consent required prior to transmission? What is the ethically appropriate way to handle this requirement when consent is not possible (no recontact clause, impossible to recontact, or an individual has died. If consent for transmission was not part of the original consent, must recontact occur to obtain consent? If consent is not obtained is it unethical to transmit even anonymized data, or particular types of anonymized data such as that governing smaller more easily identified populations/individuals?).
5. What constitutes appropriate consent for prospective data collection?
6. How should data that are not tied to consent be handled? That is, in the absence of explicit approval or refusal can data be used and if so, should it to be deidentified or anonymized, or not used?
7. Is access to deceased individuals' information ethically appropriate? Is familial consent required if individuals did not authorize access prior to their death?
8. In fostering standardization, there is a need to consistently define "deidentified" and "anonymized" and ensure consistent use across the LSDBs.
9. Create an Ethics Review Committee that can ensure that the HVP adheres to the utmost ethical principles as well as serve as counsel in the event of issues in need of resolution. For example, the Committee could determine whether rare mutations ought to have more stringent protections than common mutations.
10. Ensure strict protection of privacy in cases of rare mutations in rare disease or unique combinations of clinical features, where identification of individuals is at risk.
11. Ensure against harms to vulnerable individuals or groups of individuals, including groups defined on the basis of data findings.
12. Disclosure of additional data, even for clinical reasons to medical professionals, cannot occur unless prior authorization given by donor in original consent.
13. Define conditions for removal from database.
14. Virtual or actual medical relationships with an enquirer are unwise and are discouraged (e.g., a company or person who performs the genetic test should not provide advice on the use of the results).
15. Determine the appropriate limits to access within LSDBs to prevent individual identification. (ex. Should parts of the LSDB be password protected?)
16. Consider ethical issues arising from transferring data from LSDB to Web browsers; refer to the Ethics Committee for consensus building.

### Box 8: Development of justifications for funding

1. Focus on scientific questions and hypotheses that can be addressed using the variome database. For example, finding out what proportion of cases in each single gene disorder have unexpected phenotype(s) due to presumably modifier genes will require substantially reduced effort when collection is complete, up to date and easily accessible, and this information will be available at all geographical locations. Provide a flow of data from patient via clinician/laboratory and curator to central browsers/databases will be automated and seamless.
2. Both gene specific and disease specific programs can be developed when applying for funding (i.e., the InSiGHT colon cancer project).
3. Develop plans for grant applications for all venues. For examples, applications to international, national, state or province or department, country health departments; research agencies (e.g., NIH, NCI, MRC, NHMRC, EC, OECD); charitable organizations, foundations, philanthropic trusts and individuals; companies interested in data management and sequence data acquisition. Many of these efforts can be pursued in parallel. Successful applications will provide leverage for additional ones.
4. Demonstrate feasibility and productivity of approaches before application (pilot projects). The most notable at present is the InSiGHT collaboration and the "Adopt a gene" approach.
5. Lobby for general support of the HVP as well as for individual grants.
6. Proceed with systematical implementation of the Adopt a Gene approach, which has now been initiated, for funding of curators.

practically impossible. The reality is that LSDBs, which capture the core information in any one of these domains of information, are developed by experts and curated with invaluable skill and experience. To force any change on these individual efforts would risk inestimable loss of activity by the curators and threaten data of individuals in populations. The challenge then is to integrate the existing and developing information within existing databases and public resources into a system based on this matrix of domains—the vision of the HVP.

The task of developing "super searching" software to interrogate the global information and relate it across the HVP matrix in a user-friendly fashion for enquiry represents a bioinformatics challenge already embraced by the numerous national projects (e.g., [Stein,

2008]). This challenge can be met with resources applied to software development or existing applications that allow searches to locate all information across all domains of the international data matrix. Hansen's SRS (sequence retrieval system) approach (<http://e-hrc.net/hdi/>), or the novel Genome Commons Navigator [Brenner, 2007], supported the Berkeley Computational Biology Center (<http://ccb.berkeley.edu/ccb/index.html>), both target this concept. The Navigator also seeks to provide algorithms for potential interpretation of pathogenicity.

The approaches presented at the HVP planning meeting contribute to this goal of integrating biological information. The HVP efforts are also consistent with the newly emerging initiative to develop standards for scientific disciplines and research

### Box 9: HVP governing plan to be developed with stakeholders

1. Examine overall governance strategy outside the facilitating activities of the coordinating office.
2. Examine the potential use of the term Human Variome Programme to be applied to all variation activities outside the initial remit of the Human Variome Project as initiated in Melbourne. For example, the utility of the datasets for genome wide association studies, HapMap projects, diagnostics for personalized nutrition and medicine.
3. Evaluate the need for a term or body such as the Human Variome Organisation (HUVVO) and examine whether Human Proteome Organization (HUPO: <http://www.hupo.org/>) is a satisfactory model both for the "HVPprogramme" and HVPProject."
4. Define HVP projects and HVP collaborating projects. The first example of a collaboration project will be with the InSiGHT consortium.
5. Develop draft criteria for new HVP initiatives.
6. Develop a HVP international consortium after piloting a project in Australia.
7. Develop a Governance/Strategic working group from the Planning Group and others.
8. Encourage standards in the various facets of the HVP and harmonization if this is not possible.
9. Encourage the development of miniature devices for complete genome sequencing.

### Box 10: The InSiGHT committee description and plans for the prototype HVP pilot

The problem of interpreting the pathogenicity of missense variants is most commonly encountered in genetic testing for cancer predisposition syndromes. Genetic testing is commonly performed for Hereditary Nonpolyposis Colorectal Cancer (HNPCC, Lynch Syndrome), which results from inherited pathogenic mutations in at least four genes responsible for DNA Mismatch Repair (MMR). Between 15–30% of MMR variants detected during genetic testing for Lynch syndrome are missense, and most are rare and not interpretable regarding pathogenicity [Peltomaki and Vasen, 1997]. InSiGHT has established committees addressing phenotype, curation, virtual histology, and funding, which provide advice and support for the Interpretation Committee comprised of experts ranging from clinicians to basic scientists. Its two short-term goals are to (1) produce a paper discussing standards for data types, classification, and integration; and (2) convene to classify as many variants as possible from various international MMR gene mutation databases. The InSiGHT Interpretation Committee is undertaking the following actions:

1. Consolidation of multiple MMR gene databases
2. Assessment quality standards for phenotypic data, including tumor phenotype
3. Establishment of standards for clinical, epidemiologic, and pathologic data
4. Formation of standards for in vitro assays
5. Standardization and quality assurance methods for computational data and methods
6. Validation of methods for integrating these multiple data sources
7. Using databases to organize data and allowing access to the biomedical community

### Box 11: Other pilot and relevant projects

1. Kuwait (<http://www.al-mulla.org/>) established a Molecular Genetics Diagnostic Service Division, integrated within a research laboratory, which is focused on delivering state-of-the-art diagnostic service for the Kuwaiti population. These services provide the research laboratory with a unique opportunity to meet probands with various genetic diseases.
2. The Danish National HNPCC system has had 12 months' experience through the Web-based, but closed Danish Health Data Network ([www.sundhed.dk](http://www.sundhed.dk)). This model is relevant for all genetic diseases. A lesson to date here is that systems are successful where the output improves the user's activity and efficiency.
3. Systems to support public and personal health in nutrition are being developed through two international nutrigenomics consortia for type 2 diabetes and micronutrient genomics (Web sites within <http://www.nugo.org>, free registration). These efforts introduce a nutritional domain in the data matrix that relate variants of the relevant genes to nutrition profiles.
4. The Rare Metabolic Disease database (RAMEDIS <http://www.ramedis.de>) [Topel et al., 2006] serves the needs of specialist metabolic centers, including their clinical management, and provides access for an understanding of normal biological processes often first informed by the study of these rare disorders. These too could be linked into the broader HVP plan. The interface of microarray expression data into the HVP plan offers yet another very challenging data dense domain, which will need careful interpretation with respect to its utility in clinical medicine.
5. The UK has developed the central MutDB some time ago to collect all UK variation data and its effect from participating laboratories (<http://mutdb.org/>) [Dantzer et al., 2005].

strategies: the MIBBI Project (minimal reporting guidelines for biological and biomedical investigations) [Taylor, 2007; Taylor et al., 2008]. The challenge to catalog and access this vast body of information relating to human biology and behavior is immense, but the HVP is leading this endeavor through international collaborations and harmonized protocols. The development of this network of LSDBs and the knowledge they generate and maintain will be beneficial not only for the genetic research community, but also for researchers in nutrition, toxicology, teratology, physiology—virtually all biological research arenas, but perhaps most importantly for the translation of basic research for improving personal and public health. The future is indeed exciting.

### Disclaimer

This work includes contributions from, and was reviewed by, the FDA. This work has been approved for publication by this agency but it does not necessarily reflect official agency policy.

### Acknowledgments

The authors thank Maria Mendoza and Donna Mendrick of FDA/NCTR for critically reviewing the manuscript. Mike Parker and Helen Firth are acknowledged for their contribution to the work leading up to Box 7. J.M.H. thanks CASIMIR (funded by the European Commission under contract number LSHG-CT-2006-037811) for financial support.

## Appendix

### Author Affiliations

<sup>1</sup>Division of Personalised Nutrition and Medicine, FDA/National Center for Toxicological Research, Jefferson, Arkansas

<sup>2</sup>Genomic Disorders Research Centre, Melbourne, Australia

<sup>3</sup>Faculty of Medicine, Dentistry and Health Sciences, University of Melbourne, Parkville, Australia

<sup>4</sup>American College of Medical Genetics, Bethesda, Maryland

<sup>5</sup>Department of Paediatrics, Riyadh Military Hospital, Kingdom of Saudi Arabia

<sup>6</sup>Princess Al Jawhara Center for Hereditary Disorders, King Abdulaziz University, Saudi Arabia

<sup>7</sup>Molecular Pathology Unit, Kuwait University, Kuwait

<sup>8</sup>Department of Genetics, Physical Anthropology and Animal Physiology, University of the Basque Country, Spain

<sup>9</sup>Institute of Human Genetics, University of Bonn, Germany

<sup>10</sup>Laboratory of Human Genetics and Hematology, The Rockefeller University, New York, New York

<sup>11</sup>Lab Medicine & Pathobiology, Mount Sinai Hospital, University of Toronto, Canada

<sup>12</sup>Surgical gastroenterology, Hvidovre Hospital, Copenhagen, Denmark

<sup>13</sup>Korean Bioinformation Center, KRIBB, ChungNam, South Korea

<sup>14</sup>Medical Genetics, University of Alberta, Edmonton, Canada

<sup>15</sup>Genome Analysis, HZI—Helmholtz Centre for Infection Research, Braunschweig, Germany

<sup>16</sup>University of California, Berkeley, California

<sup>17</sup>Institute of Human Genetics, International Centre for Life, United Kingdom

<sup>18</sup>Centre for Genomic Regulation and Center for Network Biomedical Research on Epidemiology and Public Health, Barcelona, Spain

<sup>19</sup>Genetic Service ASL, Napoli, Italy

<sup>20</sup>Epidemiology and Public Health Analyses, Inserm, Toulouse, France

<sup>21</sup>Human Genetics, Navigenics, Redwood Shores, California

<sup>22</sup>Unit of Genomics for Diagnostics of Human Disease and Laboraf, San Raffaele Scientific Institute, Milano, Italy

<sup>23</sup>Victorian Research Laboratory, NICTA (National ICT Australia), Parkville, Australia

<sup>24</sup>Division of Structural and Functional Genomics, Korean National Institute of Health, Seoul, South Korea

<sup>25</sup>Department of Microbiology and Genomics, Catholic University of Korea, Seoul, South Korea

<sup>26</sup>Université Montpellier 1, Faculté de Médecine and CHU, Laboratoire de Génétique Moléculaire, Montpellier, France

<sup>27</sup>Johns Hopkins University School of Medicine, Institute of Genetic Medicine, Baltimore, Maryland

<sup>28</sup>Department of Genetics, University of Leicester, United Kingdom

<sup>29</sup>Human and Clinical Genetics, Leiden University Medical Center, The Netherlands

<sup>30</sup>European Projects Management and Coordination Office, Fundació IMIM, Barcelona, Spain

<sup>31</sup>Research & Development, Idaho Technology, Inc., Salt Lake City, Utah

<sup>32</sup>Centro de Genética Médica J., Magalhães, INSA, Portugal

<sup>33</sup>Department of Genetics, Evolution and Environment, University College London, United Kingdom

<sup>34</sup>Molecular Genetics, Pathology Queensland, Royal Brisbane and Women's Hospital, Herston, Australia

<sup>35</sup>EMBL—European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, United Kingdom

<sup>36</sup>Division of Clinical Genome Research, Institute of Medical Science, The University of Tokyo, Japan

<sup>37</sup>Clinical Pathophysiology, University of Florence, Italy

<sup>38</sup>Institute of Medical Genetics, Tomsk, Russia

<sup>39</sup>College of Medicine, University of Vermont, Burlington, Vermont

<sup>40</sup>OMIM and Institute of Genetic Medicine, Johns Hopkins University School of Medicine, Baltimore, Maryland

<sup>41</sup>Bioinformatics Group, MRC Harwell, United Kingdom

<sup>42</sup>The Pennsylvania State University, University Park, Pennsylvania

<sup>43</sup>Health Sciences Library and The Victorian Mental Health Library, Royal Melbourne Hospital, Parkville, Australia

<sup>44</sup>Computer Science and Artificial Intelligence Laboratory, MIT, Cambridge, Massachusetts

<sup>45</sup>Genetics, Ethics & Policy Consulting Inc., Boston, Massachusetts

<sup>46</sup>University of the Basque Country, Spain

<sup>47</sup>Advanced Computational Sciences Department, RIKEN, Kanagawa, Japan

<sup>48</sup>Establissement Francais Du sang, Auvergne-Loire, France

<sup>49</sup>Functional Genomics Branch, National Cancer Center, Goyang, South Korea

<sup>50</sup>Department of Medical Genetics, State University of Campinas, Brazil

<sup>51</sup>Department of Colorectal Medicine and Genetics, Royal Melbourne Hospital, Australia

<sup>52</sup>OMIM, NCBI, Baltimore, Maryland

<sup>53</sup>Clinical Genetics, Kuwait Medical Genetics Centre, Kuwait

<sup>54</sup>HLA Informatics Group, Anthony Nolan Research Institute and Department of Haematology, UCL Cancer Institute, London, United Kingdom

<sup>55</sup>Department of Medical Genetics, Tohoku University School of Medicine, Sendai, Japan

<sup>56</sup>Department of Genetics, University of Alabama, Birmingham, Alabama

<sup>57</sup>St. Josephs-Hospital Bochum-Linden, Germany

<sup>58</sup>Department of Medicine, Radboud University Nijmegen Medical Center, The Netherlands

<sup>59</sup>Genome Medicine, London, United Kingdom

<sup>60</sup>Institute of Functional Genomics, University of Regensburg, Germany

<sup>61</sup>Department of Medicine, Genetics, Institute of Human Genetics, University of Minnesota, Minneapolis, Minnesota

<sup>62</sup>Genetic Alliance, Washington, DC

<sup>63</sup>Pediatric Clinic Department, School of Medicine, National University of Córdoba, Argentina

<sup>64</sup>Human Mutation, Hoboken, New Jersey

<sup>65</sup>Medical Genetics, University of Calgary, Alberta, Canada

<sup>66</sup>MGC—Department of Cell Biology and Genetics, Faculty of Medicine and Health Sciences, Erasmus University Medical Center, Rotterdam, The Netherlands and University of Patras, Department of Pharmacy, Patras, Greece

<sup>67</sup>INRAN—National Research Institute on Food & Nutrition, Rome, Italy

<sup>68</sup>School of Biological and Chemical Sciences, Queen Mary, University of London, United Kingdom

<sup>69</sup>Genzyme Therapeutics Ltd, Oxford, United Kingdom

<sup>70</sup>ADINOVO Center for Genetic & Genomic Medicine, The First Affiliated Hospital of Zhejiang University School of Medicine, Hangzhou, P.R. China, and University of Rochester Medical Center, New York

<sup>71</sup>Funding Health Information Policy, Department of Human Services, Melbourne, Australia

<sup>72</sup>Institute of Infectious Diseases and Molecular Medicine, University of Cape Town, South Africa

<sup>73</sup>OHSU DNA Diagnostic Lab, Oregon Health & Science University, Portland, Oregon

<sup>74</sup>Department of Medicine, The University of Melbourne, Epping, Australia

<sup>75</sup>Metabolic Department, Klinik fuer Kinder-und Jugendmedizin, Reutlingen, Germany

<sup>76</sup>Discipline of Medical Genetics, Faculty of Health, University of Newcastle, Australia

<sup>77</sup>Division of Cancer Control and Population Sciences, National Cancer Institute, NIH, Bethesda, Maryland

<sup>78</sup>Department of Structural and Molecular Biology, University College London, United Kingdom

<sup>79</sup>Department of Genetics, University Medical Center Groningen, Groningen, The Netherlands

<sup>80</sup>Fundacion Publica Galega de Medicina Xenomica, Santiago de Compostela, Spain, and Center for Network Biomedical Research on Rare Diseases, Institute of Health Carlos III, Madrid, Spain

<sup>81</sup>Centre for Genome Medicine, RIKEN, Tokyo, Japan

<sup>82</sup>International Agency for Research on Cancer (IARC), Lyon, France

<sup>83</sup>Regional DNA Laboratory, Cancer Research, UK Mutation Detection Facility, Leeds, United Kingdom