

# The Impact of Genomic Variation on Function (IGVF) Consortium

Authors: IGVF Consortium (See below for detailed author list).

Correspondence to: [igvf-markerpapercorrespondence@gowustl.onmicrosoft.com](mailto:igvf-markerpapercorrespondence@gowustl.onmicrosoft.com)

## Abstract

Our genomes influence nearly every aspect of human biology from molecular and cellular functions to phenotypes in health and disease. Human genetics studies have now associated hundreds of thousands of differences in our DNA sequence (“genomic variation”) with disease risk and other phenotypes, many of which could reveal novel mechanisms of human biology and uncover the basis of genetic predispositions to diseases, thereby guiding the development of new diagnostics and therapeutics. Yet, understanding how genomic variation alters genome function to influence phenotype has proven challenging. To unlock these insights, we need a systematic and comprehensive catalog of genome function and the molecular and cellular effects of genomic variants. Toward this goal, the Impact of Genomic Variation on Function (IGVF) Consortium will combine approaches in single-cell mapping, genomic perturbations, and predictive modeling to investigate the relationships among genomic variation, genome function, and phenotypes. Through systematic comparisons and benchmarking of experimental and computational methods, we aim to create maps across hundreds of cell types and states describing how coding variants alter protein activity, how noncoding variants change the regulation of gene expression, and how both coding and noncoding variants may connect through gene regulatory and protein interaction networks. These experimental data, computational predictions, and accompanying standards and pipelines will be integrated into an open resource that will catalyze community efforts to explore genome function and the impact of genetic variation on human biology and disease across populations.

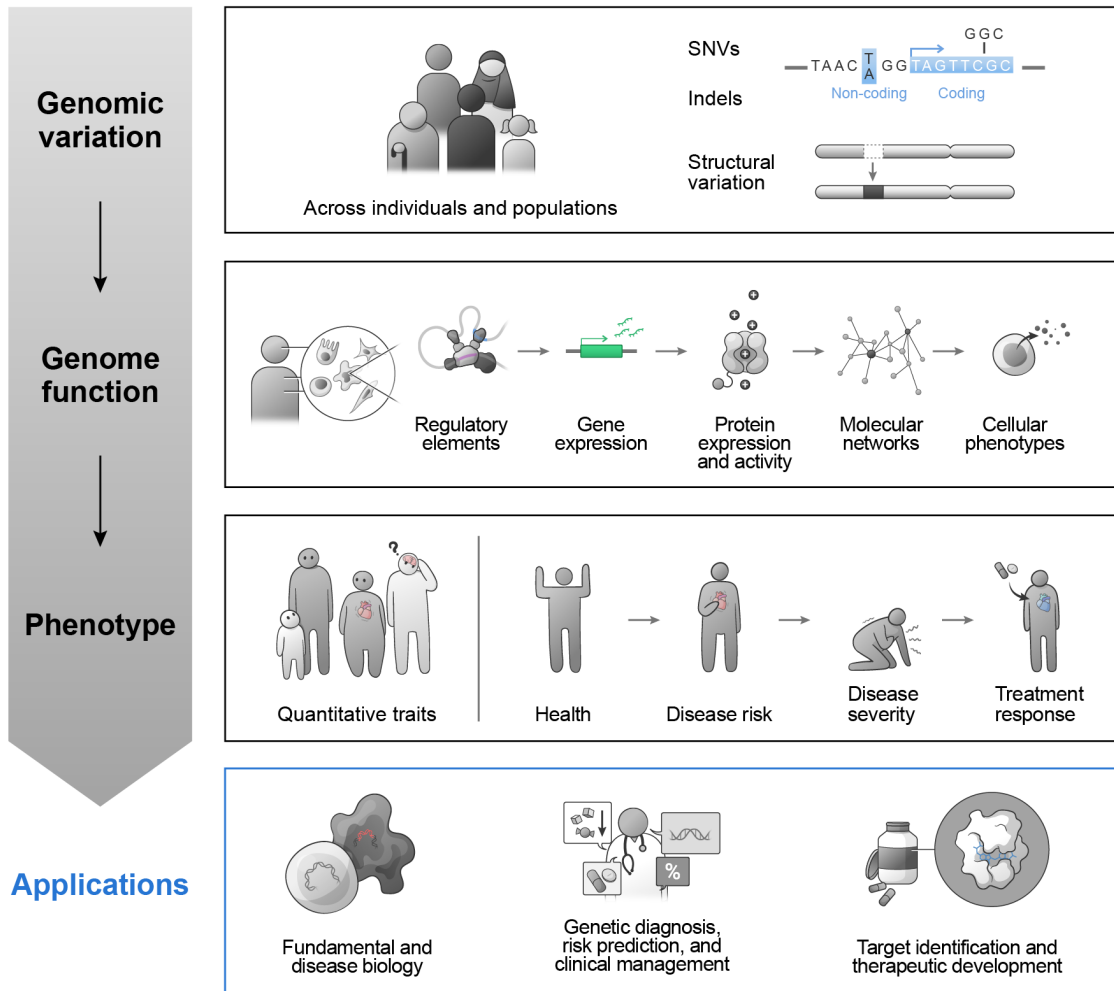
## Introduction

Since the initial sequencing of the human genome, genetic studies have been immensely productive in identifying genomic variants and associating those variants with phenotypes<sup>1-3</sup>. Exome and genome sequencing studies have already observed hundreds of millions of genomic variants, including single-nucleotide variants (SNVs), small insertions and deletions (indels), and larger structural variants (**Fig. 1**)<sup>4,5</sup>. Comparisons within families, case-control cohorts, and population-scale biobanks have now identified hundreds of thousands of associations between such variants and phenotypes in both health and disease<sup>6-12</sup>.

The next challenge is to understand how genomic variation affects molecular and cellular processes (“genome function”) to influence organismal phenotype (**Fig. 1**). At a molecular level, genomic variation can impact the expression, activity, or localization of genes and proteins. Altered gene expression or protein activity can, in turn, impact the activity of other genes and proteins via networks of physical or functional interactions. Changes in molecular networks can then influence the behavior of cells and tissues, and in doing so can influence organismal phenotypes. We note that we use the term “genome function” to refer to these molecular and cellular processes encoded by the genome, and note that this does not necessarily imply “function” in terms of organismal or evolutionary selection.<sup>13,14</sup>

Previous and ongoing efforts have produced breakthroughs in mapping various aspects of genome function, including locating and annotating millions of noncoding regulatory elements in the human genome<sup>15,16</sup>; mapping associations between genomic variants and effects on gene or protein expression across dozens of human tissues<sup>17,18</sup>; profiling hundreds of cell types and states through single-cell measurements of gene expression<sup>19,20</sup>; applying saturation

mutagenesis to analyze coding variants in selected disease genes<sup>21–23</sup>; and characterizing how genes and proteins interact genetically or physically in molecular networks<sup>24–26</sup>. These and other studies have also demonstrated how mapping the impacts of genomic variation on genome function can reveal molecular mechanisms in human biology and disease, guide genetic diagnosis and clinical management, and facilitate the development of novel therapeutics (**Fig. 1**, reviewed in<sup>1,27,28</sup>). In instances when disease mechanisms have been revealed, there are often accompanying advances in understanding basic biology with far-reaching benefits beyond the specific disease of study.



**Figure 1. Genomic variation influences genome function and phenotype.**

Yet, connecting genomic variants to functions and phenotypes continues to prove challenging, and numerous obstacles have blocked rapid progress. The sheer number of genomic variants, both those we have observed already and those we might observe in the future, is immense, and we lack any perturbation-based data for most variants. Due to linkage disequilibrium, most genetic associations with common diseases contain many candidate variants, and the variant(s) that causally affect disease risk are unknown. The vast space of possible molecular and cellular effects has been challenging to study systematically: for example, coding variants might affect protein function via effects on stability, localization, or

interactions with proteins; noncoding variants might affect gene expression through effects on transcription factor binding, chromatin state, and regulatory interactions; and genes and proteins might impact cellular processes through diverse mechanisms involving gene regulatory networks, signaling pathways, protein complexes, and other interactions. Genomic variants, elements, and networks can also have highly cell-type or context-dependent activities, yielding additional complexity given the large number of cell types in the human body. Finally, while previous efforts have largely focused on individual layers of genome function, such as studying coding variants or annotating noncoding regulatory elements, understanding the impact of genomic variation on phenotypes and disease may require a more holistic, integrative understanding of genome function that connects molecular to cellular to physiological processes. Due to these and other challenges, the molecular mechanisms underlying many genetic associations for common diseases remain to be established<sup>2,29</sup>, and genetic diagnosis for rare diseases continues to be hindered by the preponderance of variants of uncertain significance (VUS)<sup>7,30</sup>. New coordinated research activities will be needed to address the scale and scope of these challenges and thereby unlock the vast unrealized potential for understanding human biology and for improving human health<sup>31,32</sup>.

With these challenges in mind, the National Human Genome Research Institute (NHGRI) launched the IGVF Consortium in 2021, with the goal of developing a systematic understanding of the effects of genomic variation on genome function and how these effects shape phenotypes. The Consortium consists of >120 laboratories collaborating on five key activities: (i) Mapping Centers, to analyze regulatory element and gene activity at single-cell resolution across hundreds of cell types; (ii) Functional Characterization Centers, to systematically characterize the molecular and cellular effects of introducing variants or perturbing elements and genes; (iii) Predictive Modeling Projects, to develop and apply computational approaches to comprehensively model the impact of genomic variation on genome function and guide experimental design; (iv) Regulatory Network Projects, to advance network-level understanding of the influence of genetic variation and genome function on cellular and organismal phenotypes; and (v) Data and Administrative Coordinating Centers, to lead development of resources and infrastructure to share IGVF data, standards, and pipelines with the scientific community. IGVF membership and activities are expanding further via Affiliate Membership, a process by which any researcher or research project can apply to join IGVF to drive its vision and execution. Through these activities, the IGVF Consortium aims to generate an extensive resource of experimental data, standardized protocols, and computational tools integrated into a catalog that can be broadly deployed for exploring genome function and the impact of genetic variation on human biology and diseases in diverse populations. Below we describe the goals, strategies, and anticipated deliverables of IGVF (**Box 1**).

#### **Box 1: IGVF goals and approaches**

- Characterize the impact of genomic variants, regulatory elements, and genes on molecular and cellular phenotypes — by analyzing millions of naturally occurring or designed genomic perturbations across dozens of cellular models.
- Identify where and when regulatory elements and genes are active with resolution for individual cell types and states — by applying single-cell mapping technologies across hundreds of biological samples including cellular models, tissues, and environmental contexts.
- Predict the consequences of genomic variation on genome function and phenotype for previously unstudied variants and/or cellular contexts — by developing predictive

computational models that can generalize across contexts and establishing benchmarking pipelines to evaluate and calibrate their accuracy.

- Study diverse cellular and disease systems, types of genomic variation, and aspects of genome function — by developing and applying a “map-perturb-predict” framework in which single-cell mapping, genomic perturbations, and predictive modeling are synergistically combined.
- Create an initial map that annotates the predicted effects of every possible single-nucleotide variant in the human genome on key aspects of genome function — by integrating models for how coding variants might alter protein function, how noncoding variants might affect gene expression, and how noncoding and coding variants might connect within molecular networks.
- Advance our understanding of the impact of genomic variation on disease — by exploring how best to apply IGVF resources to inform genetic diagnosis and to identify biological mechanisms of disease risk.
- Ensure that these advances are applicable to and inclusive of people of diverse sexes, ancestries, and populations — by studying individuals with different genetic backgrounds, assaying and predicting effects of variants observed in diverse populations, and studying diseases disproportionately affecting disadvantaged or under-represented populations.
- Catalyze research by others toward the long-term goal of understanding the impact of genomic variation — by partnering with the broader research community and developing resources and infrastructure to share IGVF data, methods, standards, and pipelines.

## Connecting genomic variation to effects on genome function and phenotype via a map-perturb-predict framework

To create a comprehensive catalog of the effects of genomic variation, IGVF has developed a strategy that integrates three complementary components (**Fig. 2**). One component will be to quantify the activity of regulatory elements and the expression of genes via single-cell mapping. Another will conduct systematic perturbations of variants, regulatory elements, and genes. A third will seek to generalize results to new, unstudied genomic variants and cellular contexts via predictive modeling. By integrating these three components in a map-perturb-predict framework, we aim to achieve substantial synergy across the consortium.

### *Mapping the activities of genes and regulatory elements at single-cell resolution*

Identifying noncoding regulatory elements and genes and mapping their activities across cell types and states is foundational for understanding where and when genomic variation might impact genome function. Due to technological limitations, many previous efforts have lacked this level of resolution. Recent advances in single-cell technologies now enable the generation of comprehensive maps of chromatin state and gene expression in nearly any cell type in the body<sup>19,20</sup>, and computational analysis of these datasets can help to locate candidate regulatory elements, correlate element and gene activities, identify transcription factor (TF) binding regions and footprints, and reveal molecular pathways<sup>33–35</sup>. We will collect single-cell data across

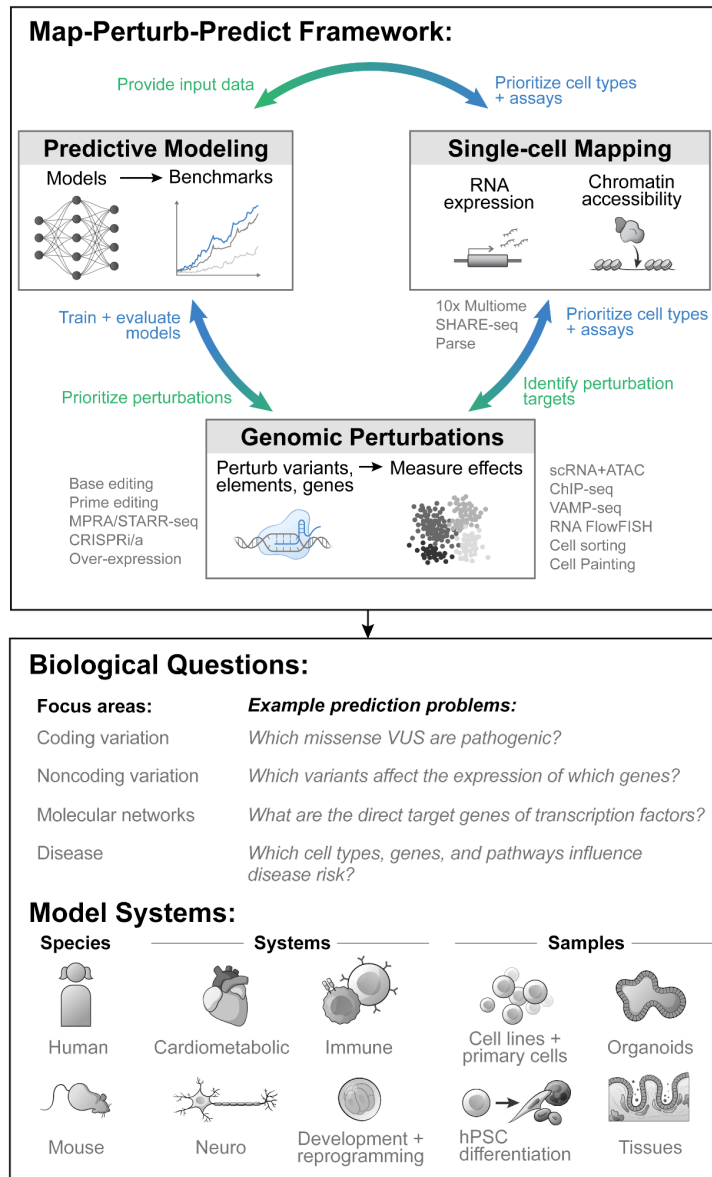
hundreds of cell types and states (see below for biological systems and contexts). We will apply primarily single-nucleus (sn)ATAC-seq and snRNA-seq, including in multiomic formats, and explore new single-cell approaches for TF binding, histone modifications, chromatin interactions, and clonal tracing. These data will provide a foundation for interpreting the effects of functional characterization experiments and building cell-type-specific maps of variant effects.

#### *Functional characterization of variants, regulatory elements, and genes via genomic perturbations*

Perturbation experiments will be crucial for understanding the causal relationships among variants, regulatory elements, genes, and phenotypes, but until recently have been challenging to apply at sufficient scale. New enabling technologies include high-throughput genetic screens using CRISPR genetic or epigenetic perturbations and over-expression strategies<sup>21,22,36-43</sup>; massively parallel reporter assays (MPRAs) to quantify enhancer and promoter activities of noncoding sequences and their variants<sup>44-50</sup>; and studies of naturally occurring genetic variation to identify and fine-map different types of quantitative trait loci (QTLs)<sup>17,51,52</sup>. IGVF plans to conduct millions of experimental perturbations, including to directly study the effects of naturally occurring or designed DNA variants, and to perturb regulatory elements and genes to build maps of genome function. We will characterize the effects of these perturbations using diverse assays including measurements of chromatin accessibility<sup>53</sup>, gene expression<sup>54-56</sup>, protein expression and activity<sup>24,57-60</sup>, and molecular and cellular phenotypes<sup>61</sup>. These data will enable directly characterizing variants of interest, such as those associated with disease, and provide data to train or evaluate predictive models of variant effects.

#### *Predictive models of genomic variation and genome function*

Genome function is complex, and we cannot expect to experimentally map the effects of all possible variants on all possible activities in all possible cellular contexts. To address this, recent studies have highlighted the possibility of developing powerful predictive models that can make predictions that generalize across contexts — for example, to link genetic variants to effects on TF binding and chromatin accessibility<sup>50,62-65</sup>; identify TF footprints<sup>33,65</sup>; connect regulatory elements to their target genes<sup>65-67</sup>; or identify causal genes and cell types enriched for heritability for complex diseases or traits<sup>68-73</sup>. Equally importantly, successes by CASP<sup>74</sup>, ENCODE<sup>15</sup>, and others<sup>17,75</sup> have illustrated how developing uniform standards, gold-standard datasets, and benchmarking pipelines can catalyze advances throughout the global scientific community by enabling rigorous evaluation of accuracy and direct comparison of alternative strategies. We will leverage new advances in machine learning and large-scale perturbation datasets across cell types and contexts to tackle key prediction problems — including mapping aspects of genome function, interpreting the impact of genomic variation, and guiding the design of future experimental assays such that the data produced will be maximally informative for subsequent predictive modeling. To systematically evaluate and calibrate such models, we will build benchmarking pipelines that compare predictions to perturbation data, including both from IGVF functional characterization experiments and external sources such as QTL, GWAS, and genome sequencing studies. In areas where data collection is already advanced, we will engage the external community by designing prediction challenges with held-out validation datasets produced by IGVF.



**Figure 2. A map-perturb-predict framework to connect genome variation to genome function and phenotype**

*Applying the map-perturb-predict framework to study genomic variation and genome function across cellular and biological systems*

Together, these three activities will form an iterative map-perturb-predict framework (**Fig. 2**) that we will apply to study various aspects of genomic variation, genome function, and phenotype. IGVF projects will investigate single-nucleotide variants, indels, and structural variation, and map the relationships between elements, genes, proteins, and their molecular networks in diverse cellular states and phenotypes (**Fig. 1**). IGVF projects will study a variety of biological systems, including iPSC models (2D and 3D) differentiated into lineages spanning all germ layers; primary cell types relevant to disease areas of interest, including cardiometabolic, immune, neuropsychiatric, and neurodevelopmental diseases; and tissues *in vivo* to inform how

cell-cell interactions and environment alter genome function (**Fig. 2**). The selected models will include dynamic biological processes that will provide insights into how regulatory networks change over time, such as B cell activation and differentiation or fibroblast-to-iPSC reprogramming. While the primary objective of IGVF is to characterize variation and function of the human genome, IGVF studies will also create resources and leverage mouse models for certain studies, such as for *in vivo* CRISPR screens to understand how genes affect cellular phenotypes in a tissue environment, and for comparing the effects of variants, elements, and genes across individuals with different genetic backgrounds. Together, these areas of exploration will yield insights about genomic variation and genome function across diverse areas of biology and help to identify optimal strategies that can be more broadly applied to additional biological systems.

## Map of genome function and variant effects integrating coding variation, noncoding variation, and molecular networks

An integrative deliverable for IGVF will be to generate a preliminary variant-effect map that integrates three key aspects of genome function: gene expression, protein activity, and molecular networks. This draft map would allow querying, for any possible single-nucleotide variant in the genome: Is this variant measured or predicted to (i) impact transcription factor binding, regulatory element activity, and target gene expression in particular cell contexts, for noncoding variants; (ii) impact protein function, for coding variants; and (iii) connect to other genes/proteins via gene regulatory networks and/or protein-interaction networks, for both coding and noncoding variants? We will integrate this map of genome function, along with external data, into a multi-relational knowledge graph<sup>76–78</sup> that can be readily queried by users as part of the IGVF Catalog (see below, **Fig. 3**).

For each of these aspects of genome function, existing computational models have been shown to have some utility in understanding the impact of genomic variation on diseases and traits, but much work is needed to improve the accuracy of these models and to conduct systematic evaluations using more precise and comprehensive gold-standard datasets. To address this, we will establish a pipeline to benchmark all predictions against functional characterization datasets and external human genetics datasets — allowing us to rigorously evaluate and guide interpretation of the draft map. We anticipate that providing genome-wide predictions from the best models, together with a reproducible benchmarking framework, will help launch an iterative and ongoing effort extending beyond IGVF to improve the accuracy of this map over time (**Fig. 3**).

### *Effects of noncoding single-nucleotide variants on regulatory element activity and target gene expression*

In the 99% of our genome that does not encode for proteins, noncoding variants can impact genome function by altering gene expression or regulation. While previous studies have mapped regulatory elements, gene expression patterns, transcription factor binding, and the effects of variants on gene expression in tissues or cells, we still lack models that can make accurate causal inferences about how genomic variation affects gene regulation<sup>79–81</sup>. We will seek to build genome-wide annotations of key components of this *cis*-regulatory code: Which single-nucleotide variants affect transcription factor binding sites, regulatory element activity, and gene expression in *cis*, in which cell types or states, with what magnitude and direction of effect?

To do so, IGVF plans to (i) generate multiomic snRNA+ATAC-seq data to a depth needed to comprehensively identify candidate *cis*-regulatory elements, detect transcription factor footprints<sup>33</sup>, and predict enhancer-gene relationships<sup>34,35,65,67</sup>; (ii) test >1 million noncoding

variants in enhancer activity reporter assays<sup>44,45,49,50,82</sup>; (iii) test >100,000 noncoding variants for effects on expression through fine-mapping of eQTLs or direct CRISPR-based genome editing<sup>17,36–38,40,51</sup>; (iv) measure >100,000 putative regulatory interactions between candidate regulatory elements and nearby genes, for example using dCas9-based epigenome editing paired with single-cell readouts of RNA expression<sup>66,83–86</sup>; (v) and perturb transcription factors to read out effects on gene expression using Perturb-seq<sup>54–56</sup>. The variants and elements studied will include both naturally occurring and designed sequences, which will be critical for building accurate models of the gene regulatory code<sup>87</sup>. Each of these experiments will be conducted in multiple cellular models, so that the data can be used to refine and develop predictive models that can construct maps of noncoding variant effects across many cell types.

### *Effects of single-nucleotide variants in protein-coding genes on function*

For protein-coding sequences, our ability to interpret the functions of genomic variation is based on our knowledge of the genetic code for protein synthesis — which has enabled identifying open reading frames encoding novel proteins, identifying null, frameshift or nonsense variants, and predicting damaging missense variants. However, missense variants and inframe indels remain difficult to interpret, and we still lack a comprehensive understanding of how changes in protein sequence might affect different aspects of protein structure, expression, dynamics, and activity, including the impact on stability, subcellular localization, or interactions with other proteins.

We will improve annotations of how protein-coding missense variants impact protein stability and activity by applying high-throughput technologies<sup>24,57–60</sup> to experimentally characterize the impacts of >200,000 missense variants on various properties of proteins and their phenotypic impacts in cellular models, including protein stability, subcellular localization, cell viability, cell morphology, and protein-protein interactions. These experiments will focus on clinically relevant genes, such as those associated with Mendelian diseases, to provide direct look-up tables for certain genes, and provide data to refine or develop new predictive models to predict the likely impact of any coding variant across the genome.

### *Linking noncoding and coding variants to gene regulatory and protein interaction networks and selected cellular phenotypes*

Upon linking a variant to effects on gene expression or protein activity in *cis*, we will seek to annotate the sets of other genes and proteins linked to the variant in *trans* through molecular networks in a given cell type or state. To a more limited extent, we will explore links to downstream cellular phenotypes. Genes and proteins can work together in many different ways, and it has been challenging to map or infer these sets of functionally related genes and corresponding cellular phenotypes in a comprehensive and cell-type specific fashion.

To construct molecular networks, we will focus on defining (i) gene expression programs, described by sets of genes whose expression levels are correlated across single cells; (ii) gene regulatory networks that infer which transcription factors directly regulate which target genes via particular regulatory sequences; (iii) sets of interacting proteins or protein complexes; and (iv) dynamic changes to these programs/networks across cell state transitions.

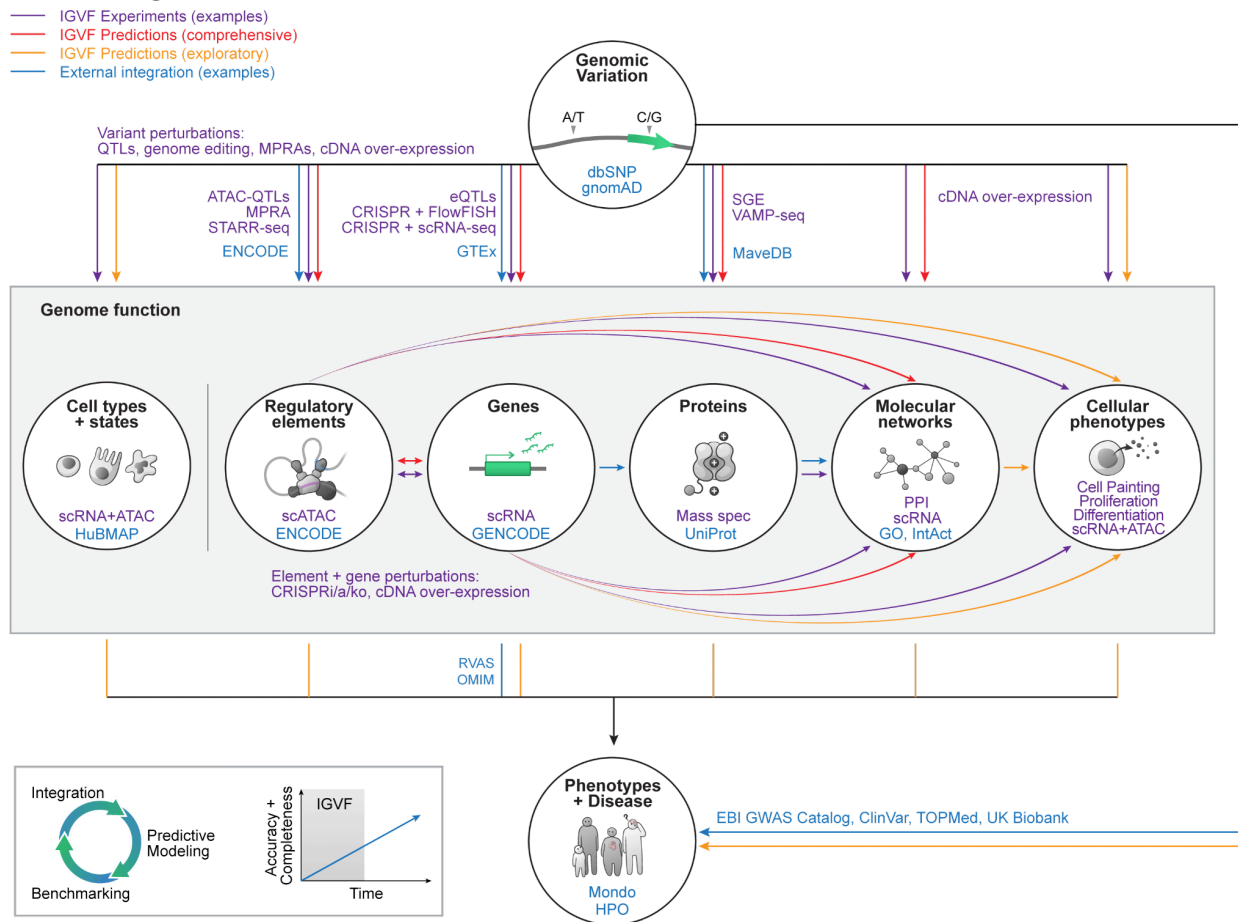
To build these maps, we will collect longitudinal multiomic RNA and ATAC-seq data across dynamic cellular processes including differentiation and reprogramming<sup>88–91</sup>; study how genes and proteins interact in molecular networks, including by mapping protein-protein interactions<sup>24</sup> and conducting large-scale Perturb-seq<sup>54–56</sup>; and assessing how CRISPR-based perturbations or natural genetic variation across individuals affects cellular phenotypes including differentiation, gene expression programs, and cellular states. We will establish benchmarks to evaluate how best to use these data to construct cell-type and state-specific molecular networks and assess the impact of genomic variation on cellular phenotypes.



We anticipate that many aspects of this map of genome function and variant effects will be cell-type specific, with annotations for each of the hundreds of cell types, states, and contexts studied by IGVF. This could be accomplished by developing predictive models that use multiomic snRNA-seq and ATAC-seq as their only cell-type specific input data<sup>34,35,67,92</sup>.

The research infrastructure IGVF develops to build these maps will set in motion community efforts to expand on this framework by collecting additional datasets, training improved models, generating more accurate maps, and expanding the approach to additional cell types and aspects of genome function. This draft map will also offer immediate opportunities to address questions about the impact of genomic variation and genome function on phenotypes (see next section).

### IGVF Catalog:



**Figure 3.** The IGVF Catalog of genome function and the impact of genomic variation. IGVF will create a catalog linking genomic variation (top) to genome function (middle box) to phenotype (bottom). Purple: Examples of experimental methods applied by IGVF. Red: Relationships where IGVF plans to develop and apply computational models to comprehensively annotate all possible single-nucleotide variants across many cell types. Orange: Relationships where IGVF plans to develop and apply computational methods in a more targeted fashion, for example in the context of certain cellular phenotypes or diseases. Blue: Examples of external resources or ontologies that could interact with and/or be incorporated into this catalog. Abbreviations and citations: dbSNP<sup>93</sup>, gnomAD<sup>4</sup>, ENCODE<sup>94</sup>, GTEx<sup>17</sup>, saturation genome editing (SGE)<sup>22</sup>, Variant Abundance by Massively Parallel sequencing (VAMP-seq)<sup>21</sup>, MaveDB<sup>23</sup>, HuBMAP<sup>19</sup>, GENCODE<sup>95</sup>, UniProt<sup>96</sup>, Gene Ontology (GO)<sup>97</sup>, IntAct Molecular Interaction Database<sup>98</sup>, Mondo Disease Ontology<sup>99</sup>, Human Phenotype Ontology (HPO)<sup>100</sup>, rare variant association studies (RVAS), Online Mendelian Inheritance in Man (OMIM)<sup>101</sup>.

## Exploring the impact of genomic variation and function on disease

The map-perturb-predict framework and the resulting variant-effect maps will provide new resources for the community to study the impact of genomic variation on human diseases and phenotypes, but this goal presents additional challenges.

For many diseases, an individual's risk is likely to be determined by a combination of thousands of independently acting variants<sup>102,103</sup> — including for many diseases presumed to follow Mendelian inheritance patterns, where penetrance and expressivity may include a component of polygenic risk<sup>104</sup>. Molecular networks are highly interconnected — a single variant may influence multiple genes, multiple gene networks, in diverse cell types — making it difficult to determine which of those genes, networks, and cell types are important for disease<sup>1,17,69,70</sup>. Disease susceptibility can involve many different cell types, possibly at specific timepoints, with effects accumulating over decades or in specific environmental contexts<sup>105</sup>. The impact of genomic variation on genome function can also differ across age, sex, populations, and ancestry: expanding human genetic studies across diverse populations has revealed examples where additional disease associations are discovered due to differences in allele frequencies<sup>106</sup>, and some cases in which variant with comparable allele frequencies appear to have different effect sizes on a disease<sup>107–111</sup>.

Toward addressing some of these challenges, we will focus on assessing the impacts of variants in molecular networks and diverse cell contexts and then explore how best to apply this framework to: (i) inform clinical variant interpretation, particularly for rare diseases; (ii) learn about molecular and cellular mechanisms underlying risk for common and rare diseases; and (iii) ensure that lessons about the impact of genomic variation on genome function are applicable across diverse populations. Notably, each of these questions represents a major research area involving many strategies beyond those pursued in IGVF<sup>7–9,28,112–114</sup>, and these exploratory efforts will seek to integrate with other efforts in the field.

### *Interpreting genomic variation to inform genetic diagnosis*

One key use case for variant-effect maps generated by IGVF, particularly for coding variation, will be to inform the clinical interpretation of single-nucleotide VUS in genes with known and suspected links to Mendelian genetic diseases. Indeed, prior work has shown how applying multiplexed assays of variant effect to study individual genes has been translated into powerful evidence for clinical variant interpretation, for example moving 50% of VUS in *BRCA1*, 70% in *TP53*, 74% in *MSH2*, and 90% in *DDX3X* into more definitive pathogenic or benign classifications<sup>58,115,116</sup>. In some cases, updated genetic test results were provided to individuals that had received VUS results when tested for cancer risk, and diagnostic odysseys were ended for families with *DDX3X*-associated neurodevelopmental disease.

To expand this approach to additional diseases, IGVF labs will experimentally measure the effects of thousands of variants in known disease genes, with a particular focus on those where identification of loss-of-function variants is clinically actionable<sup>117,118</sup>. We will then assess the extent to which either these experimental data, or computational predictions of variant impact, enrich for variants previously classified as either pathogenic or benign, and determine whether they can be used to calibrate predictors for clinical applications<sup>119</sup>. These variant-effect maps could ultimately substantially reduce the VUS burden in etiological diagnosis of rare disease<sup>114</sup>. Integration of maps for both coding and noncoding variants could also aid in the development of the next-generation polygenic risk score methodologies for better risk characterization in complex phenotypes<sup>107</sup>.

### *Identifying molecular and cellular mechanisms of disease risk*

Improved variant-effect maps could be transformative for identifying new biological mechanisms that influence genetic risk for disease. In particular, we will seek to understand how best to combine the map-perturb-predict framework and variant-effect maps with human genetic data to nominate variants, genes, cell types, and cellular programs that influence disease risk.

We will study specific diseases and traits, including lipid traits, hematological traits, autoimmune diseases such as systemic lupus erythematosus, cardiometabolic diseases such as coronary artery disease, and neurodegenerative diseases such as Alzheimer's disease. As one example, IGVF investigators are studying variants associated with lipid traits, where GWAS and whole-exome sequencing studies have already identified hundreds of associated noncoding and coding variants, and where certain key genetic pathways involved in lipid handling are already known<sup>11,120-122</sup>. By conducting CRISPR screens to identify variants and regulatory elements that affect lipid uptake in cellular models, testing variant effects on enhancer activity in massively parallel reporter assays, and applying state-of-the-art predictive models, we will evaluate which combinations of experiments and/or predictive models provide the strongest enrichment for disease-associated variation and known causal genes. These combined efforts will help to inform mechanisms of genetic risk for selected diseases, and help to develop strategies to identify causal variants, genes, and pathways for any complex disease.

### *Assessing the impact of variation across populations*

IGVF aims to ensure that insights about the impact of genomic variation and genome function are applicable to and inclusive of people of diverse groups. To do so, we will promote diversity in its functional genomics studies, experimentally study and computationally annotate variants observed in diverse populations, study diseases disproportionately affecting disadvantaged or under-represented populations, and explore the extent to which particular variants might exert the same or different effects due to interactions with genetic background or environment<sup>123-125</sup>.

We will employ a multi-pronged approach encompassing experimental and computational strategies to achieve its goals. In the current design stage, we have incorporated variants, elements, and genes from diverse populations, including those with differential effects on disease. Biological models will include human iPSCs derived from individuals from different ancestries, and genetically diverse mouse lines from the Collaborative Cross<sup>126</sup>. Finally, certain predictive models, such as those linking noncoding variants to chromatin state and gene expression<sup>62,64,65</sup>, can make predictions of variant effects based on measurements obtained in a specific individual, allowing systematic annotation and comparison of variant effects across individuals and groups. Altogether, the data and genome-wide variant-effect maps generated by IGVF will offer insights into variant effects across groups and provide a valuable resource for investigating the effects of variants discovered in diverse populations.

## **Data release and resources**

A major goal of IGVF is to catalyze future research to understand the relationships between genome function, genomic variation, and phenotype, including for biomedical researchers across diverse disciplines and with diverse needs. To do so, we will build the IGVF Data Resource to enable researchers to easily access and apply IGVF datasets, predictions, and methods (<https://igvf.org>).

For researchers who want to explore IGVF data and predictions about genome function and variant effects, we will create the IGVF Catalog. The IGVF Catalog will enable searching for information about specific variants, genomic loci, or genes, and will draw from processed data, analysis products, and computational predictions generated by IGVF as well as external data sources such as dbSNP<sup>127</sup>, FAVOR<sup>128</sup>, gnomAD<sup>4</sup>, GENCODE<sup>95</sup>, topLD<sup>129</sup>, ENCODE<sup>94</sup>, GTEx<sup>17</sup>,

and MaveDB<sup>23</sup> (**Fig. 3**). The IGVF Catalog will be updated several times per year, and all releases will be numbered and archived to maintain reproducibility of studies that depend on its earlier versions. To support users who want to programmatically access IGVF data or analysis products — for example to perform integrative analyses or to develop novel web applications — the IGVF Catalog will also provide a fully featured application programming interface (API) to the underlying knowledge graph.

For researchers who want to access raw or processed data generated by IGVF, we will develop the IGVF Data Portal. The Data Portal will provide web-browser and programmatic access to uniformly processed IGVF datasets, analysis products, and rich metadata, which we anticipate will be useful for users who aim to develop new data analysis methods or predictive models, analyze IGVF data in new ways, or compare their data to IGVF standards. The IGVF Data Portal will follow principles of making data Findable, Accessible, Interoperable, and Reusable (FAIR)<sup>130</sup>. IGVF data and predictions will be made available once they meet pre-defined experimental and computational quality standards. Data will be stored in cloud file buckets to facilitate computing on the data in place and without the need to download data to local servers. Some IGVF data may not have consent for public sharing; such data will be deposited in NHGRI's Analysis, Visualization and Informatics Lab (AnVIL) platform to provide access control in adherence to NIH Policy<sup>131</sup>.

For researchers who want to apply IGVF methods and strategies to additional systems, the Data Portal will also share documentation on IGVF standards, protocols, and best practices for experimental design, data analysis, and predictive modeling. These resources will include computational methods, data formats, and consensus data processing pipelines for key assays and analysis products, such as for single-nucleus RNA-seq and ATAC-seq, CRISPR-based experiments, massively parallel reporter assays, eQTL studies, and others. All data processing code will be released with open-source licenses to enable others to analyze similar data in an identical fashion, and we will strive to make sure that it can be run on compute resources accessible to researchers throughout the global research community.

Finally, for all researchers, we will provide training and support on how to access these IGVF resources. To teach researchers how to find and view IGVF data, we will create instructional streaming videos that we will distribute via the IGVF YouTube channel (<https://www.youtube.com/@igvf>). To teach users how to access data programmatically and use common analytic tools, we will create online notebooks and tutorials demonstrating key uses of the IGVF Portal and the IGVF Catalog. As an additional channel for users to interact with the IGVF Consortium, we will host interactive online seminars and webinars.

Altogether, we expect that these resources will enable a wide range of scientific activities, expanding far beyond the specific studies undertaken by the IGVF Consortium.

## Collaborations and community

Toward advancing our collective efforts to understand genomic variation and genome function — a grand challenge that demands global and interdisciplinary collaboration — IGVF welcomes collaboration with and input from the broader scientific community. Researchers interested in joining IGVF can apply for Affiliate Membership. Affiliate Membership allows investigators to fully participate in working groups and other IGVF collaborations, and thereby help drive the vision, goals, and execution of consortium activities. For more information, visit <https://igvf.org/affiliate-membership/>.

IGVF is actively coordinating with other consortia, including ClinGen<sup>8</sup>, the Genomics Research to Elucidate the Genetics of Rare diseases (GREGoR) consortium, and the Atlas of Variant Effects (AVE) Alliance<sup>132</sup>. These collaborations will facilitate the open exchange and interoperability of genomic data and resources, for example to use common variant naming schema, genome and transcriptome builds, and experimental and analysis pipelines.

Similarly, IGVF data and analysis products will benefit from close interactions with efforts to characterize human genomic variation and assemblies, such as the Human Pangenome Reference Consortium (HPRC)<sup>133</sup>; with efforts to catalog disease-associated variation across ancestries, including All of Us<sup>134</sup>, TOPMed<sup>10</sup>, and other biobanks; with efforts to map the activities of variants, regulatory elements, and genes at single-cell resolution, such as the Human Cell Atlas<sup>20</sup> and HuBMAP<sup>19</sup>; and with efforts to compare and evaluate strategies for interpreting genetic variation associated with disease, such as the International Common Disease Alliance<sup>28</sup>. Strong collaborative ties among such efforts and IGVF will propel scientific advances that shape how both basic and clinical research are performed.

## Outlook and Perspectives

With the rapid expansion of human genetics studies linking variation to disease, the interpretation of the impact of genomic variation on function is currently a rate-limiting step for delivering on the promise of precision medicine. The IGVF Consortium will deploy a coordinated strategy for accelerating progress, including generating large-scale data resources, predictive models, and initial variant-effect maps that will reveal new insights into how genomic variation impacts function and phenotype. The tools, data resource, and strategy developed by IGVF — including new experimental assays, design strategies, predictive models, computational methods, data sharing standards, and more — will provide a foundation to facilitate future efforts. We will prioritize open data and resource sharing, inclusion, and outreach so that all members of the research community can participate and benefit.

While ambitious, IGVF activities do have limitations in scope, and many challenges lie ahead. Genomic technologies, both experimental and computational, are developing rapidly, and balancing the implementation of the newest scalable tools with continuing standards to ensure data interoperability will require attention. While data generation technologies have increased throughput exponentially over the last 15 years, the amount of data needed to build accurate models of genome function is unknown, and fully realizing the goal of mapping the impact of genomic variation on function will require additional advances in both experimental and computational methods. In particular, the development of computational methods to predict synergistic interactions among variants, environments, and time spans of effects that can occur over decades are open problems. We will initially focus on specific biological systems and cellular models according to its members' expertise, but full exploration of the many cell types and disease areas relevant to human biology will require community efforts. IGVF aims for systematic analysis of certain aspects of genome function — gene regulation, protein function, and molecular networks. Additional work is required to explore other important layers of genome function, including effects on nuclear organization and chromatin compartmentalization; RNA splicing, transport, and translation; and impacts on cellular phenotypes, cell-cell interactions, and communication. For all of these challenges, the framework developed by the IGVF Consortium to develop and benchmark methods, refine best practices and standards, and share data and methods will drive scientific discoveries in human health and disease for years to come.

## References

1. Claussnitzer, M. *et al.* A brief history of human disease genetics. *Nature* **577**, 179–189 (2020).
2. Loos, R. J. F. 15 years of genome-wide association studies and no signs of slowing down. *Nat. Commun.* **11**, 5900 (2020).
3. Green, E. D. *et al.* Strategic vision for improving human health at The Forefront of Genomics. *Nature* **586**, 683–692 (2020).
4. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
5. Collins, R. L. *et al.* A structural variation reference for medical and population genetics. *Nature* **581**, 444–451 (2020).
6. Sollis, E. *et al.* The NHGRI-EBI GWAS Catalog: knowledgebase and deposition resource. *Nucleic Acids Res.* **51**, D977–D985 (2023).
7. Landrum, M. J. *et al.* ClinVar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res.* **42**, D980–5 (2014).
8. Rehm, H. L. *et al.* ClinGen—the clinical genome resource. *N. Engl. J. Med.* **372**, 2235–2242 (2015).
9. Zhou, W. *et al.* Global Biobank Meta-analysis Initiative: Powering genetic discovery across human disease. *Cell Genom* **2**, 100192 (2022).
10. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature* **590**, 290–299 (2021).
11. Backman, J. D. *et al.* Exome sequencing and analysis of 454,787 UK Biobank participants. *Nature* **599**, 628–634 (2021).
12. Karczewski, K. J. *et al.* Systematic single-variant and gene-based association testing of thousands of phenotypes in 394,841 UK Biobank exomes. *Cell Genomics* **2**, 100168 (2022).
13. Doolittle, W. F., Brunet, T. D. P., Linquist, S. & Gregory, T. R. Distinguishing between ‘function’ and ‘effect’ in genome biology. *Genome Biol. Evol.* **6**, 1234–1237 (2014).
14. Kellis, M. *et al.* Defining functional DNA elements in the human genome. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 6131–6138 (2014).
15. ENCODE Project Consortium *et al.* Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature* **583**, 699–710 (2020).
16. Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–330 (2015).
17. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science* **369**, 1318–1330 (2020).
18. Võsa, U. *et al.* Large-scale cis- and trans-eQTL analyses identify thousands of genetic loci and polygenic scores that regulate blood gene expression. *Nat. Genet.* **53**, 1300–1310 (2021).
19. HuBMAP Consortium. The human body at cellular resolution: the NIH Human Biomolecular Atlas Program. *Nature* **574**, 187–192 (2019).
20. Regev, A. *et al.* The Human Cell Atlas. *eLife* **6**, (2017).

21. Matreyek, K. A. *et al.* Multiplex assessment of protein variant abundance by massively parallel sequencing. *Nat. Genet.* **50**, 874–882 (2018).
22. Findlay, G. M. *et al.* Accurate classification of BRCA1 variants with saturation genome editing. *Nature* **562**, 217–222 (2018).
23. Esposito, D. *et al.* MaveDB: an open-source platform to distribute and interpret data from multiplexed assays of variant effect. *Genome Biology* **20**, 223 (2019).
24. Luck, K. *et al.* A reference map of the human binary protein interactome. *Nature* **580**, 402–408 (2020).
25. Szklarczyk, D. *et al.* The STRING database in 2021: customizable protein-protein networks, and functional characterization of user-uploaded gene/measurement sets. *Nucleic Acids Res.* **49**, D605–D612 (2021).
26. Pacini, C. *et al.* Integrated cross-study datasets of genetic dependencies in cancer. *Nat. Commun.* **12**, 1661 (2021).
27. Visscher, P. M. *et al.* 10 Years of GWAS Discovery: Biology, Function, and Translation. *Am. J. Hum. Genet.* **101**, 5–22 (2017).
28. International Common Disease Alliance. ICDA Recommendations and White Paper. Available online: <https://drive.google.com/file/d/16SVJ5lbneN9hB9E03PZMhpescAN527HO/view>.
29. Abdellaoui, A., Yengo, L., Verweij, K. J. H. & Visscher, P. M. 15 years of GWAS discovery: Realizing the promise. *Am. J. Hum. Genet.* **110**, 179–194 (2023).
30. Rehm, H. L. & Fowler, D. M. Keeping up with the genomes: scaling genomic variant interpretation. *Genome Med.* **12**, 5 (2019).
31. Lappalainen, T. & MacArthur, D. G. From variant to function in human disease genetics. *Science* **373**, 1464–1468 (2021).
32. Findlay, G. M. Linking genome variants to disease: scalable approaches to test the functional impact of human mutations. *Hum. Mol. Genet.* **30**, R187–R197 (2021).
33. Hu, Y. *et al.* Single-cell multi-scale footprinting reveals the modular organization of DNA regulatory elements. *bioRxiv* (2023) doi:10.1101/2023.03.28.533945.
34. Granja, J. M. *et al.* ArchR is a scalable software package for integrative single-cell chromatin accessibility analysis. *Nat. Genet.* **53**, 403–411 (2021).
35. Kartha, V. K. *et al.* Functional inference of gene regulation using single-cell multi-omics. *Cell Genom* **2**, (2022).
36. Cuella-Martin, R. *et al.* Functional interrogation of DNA damage response variants with base editing screens. *Cell* **184**, 1081–1097.e19 (2021).
37. Morris, J. A. *et al.* Discovery of target genes and pathways at GWAS loci by pooled single-cell CRISPR screens. *Science* **380**, eadh7699 (2023).
38. Martin-Rufino, J. D. *et al.* Massively parallel base editing to map variant effects in human hematopoiesis. *Cell* **186**, 2456–2474.e24 (2023).
39. Anzalone, A. V., Koblan, L. W. & Liu, D. R. Genome editing with CRISPR-Cas nucleases, base editors, transposases and prime editors. *Nat. Biotechnol.* **38**, 824–844 (2020).

40. Hanna, R. E. *et al.* Massively parallel assessment of human variants with base editor screens. *Cell* **184**, 1064–1080.e20 (2021).
41. Klann, T. S. *et al.* CRISPR–Cas9 epigenome editing enables high-throughput screening for functional regulatory elements in the human genome. *Nature Biotechnology* **35**, 561–568 (2017).
42. Fulco, C. P. *et al.* Systematic mapping of functional enhancer–promoter connections with CRISPR interference. *Science* **354**, 769–773 (2016).
43. Canver, M. C. *et al.* BCL11A enhancer dissection by Cas9-mediated in situ saturating mutagenesis. *Nature* **527**, 192–197 (2015).
44. Arnold, C. D. *et al.* Genome-wide quantitative enhancer activity maps identified by STARR-seq. *Science* **339**, 1074–1077 (2013).
45. Melnikov, A. *et al.* Systematic dissection and optimization of inducible enhancers in human cells using a massively parallel reporter assay. *Nat. Biotechnol.* **30**, 271–277 (2012).
46. Bergman, D. T. *et al.* Compatibility rules of human enhancer and promoter sequences. *Nature* **607**, 176–184 (2022).
47. Vockley, C. M. *et al.* Massively parallel quantification of the regulatory effects of noncoding genetic variation in a human cohort. *Genome Res.* **25**, 1206–1214 (2015).
48. Klein, J. C. *et al.* A systematic evaluation of the design and context dependencies of massively parallel reporter assays. *Nat. Methods* **17**, 1083–1091 (2020).
49. Patwardhan, R. P. *et al.* High-resolution analysis of DNA regulatory elements by synthetic saturation mutagenesis. *Nat. Biotechnol.* **27**, 1173–1175 (2009).
50. Agarwal, V. *et al.* Massively parallel characterization of transcriptional regulatory elements in three diverse human cell types. *bioRxiv* (2023) doi:10.1101/2023.03.05.531189.
51. Wang, Q. S. *et al.* Leveraging supervised learning for functionally informed fine-mapping of cis-eQTLs identifies an additional 20,913 putative causal eQTLs. *Nat. Commun.* **12**, 3394 (2021).
52. Xu, Y. *et al.* An atlas of genetic scores to predict multi-omic traits. *Nature* **616**, 123–131 (2023).
53. Gate, R. E. *et al.* Genetic determinants of co-accessible chromatin regions in activated T cells across humans. *Nat. Genet.* **50**, 1140–1150 (2018).
54. Adamson, B. *et al.* A Multiplexed Single-Cell CRISPR Screening Platform Enables Systematic Dissection of the Unfolded Protein Response. *Cell* **167**, 1867–1882.e21 (2016).
55. Dixit, A. *et al.* Perturb-Seq: Dissecting Molecular Circuits with Scalable Single-Cell RNA Profiling of Pooled Genetic Screens. *Cell* **167**, 1853–1866.e17 (2016).
56. Replogle, J. M. *et al.* Mapping information-rich genotype-phenotype landscapes with genome-scale Perturb-seq. *Cell* (2022) doi:10.1016/j.cell.2022.05.013.
57. Sahni, N. *et al.* Widespread macromolecular interaction perturbations in human genetic disorders. *Cell* **161**, 647–660 (2015).
58. Fayer, S. *et al.* Closing the gap: Systematic integration of multiplexed functional data resolves variants of uncertain significance in BRCA1, TP53, and PTEN. *Am. J. Hum. Genet.* **108**, 2248–2258 (2021).
59. Starita, L. M. *et al.* Massively Parallel Functional Analysis of BRCA1 RING Domain Variants.



- Genetics* **200**, 413–422 (2015).
60. Sun, S. *et al.* An extended set of yeast-based functional assays accurately identifies human disease mutations. *Genome Res.* **26**, 670–680 (2016).
  61. Bray, M.-A. *et al.* Cell Painting, a high-content image-based assay for morphological profiling using multiplexed fluorescent dyes. *Nat. Protoc.* **11**, 1757–1774 (2016).
  62. Avsec, Ž. *et al.* Base-resolution models of transcription-factor binding reveal soft motif syntax. *Nat. Genet.* **53**, 354–366 (2021).
  63. Beer, M. A. Predicting enhancer activity and variant impact using gkm-SVM. *Hum. Mutat.* **38**, 1251–1258 (2017).
  64. Chen, K. M., Wong, A. K., Troyanskaya, O. G. & Zhou, J. A sequence-based global map of regulatory activity for deciphering human genetics. *Nat. Genet.* **54**, 940–949 (2022).
  65. Avsec, Ž. *et al.* Effective gene expression prediction from sequence by integrating long-range interactions. *Nat. Methods* **18**, 1196–1203 (2021).
  66. Fulco, C. P. *et al.* Activity-by-contact model of enhancer-promoter regulation from thousands of CRISPR perturbations. *Nat. Genet.* **51**, 1664–1669 (2019).
  67. Sakaue, S. *et al.* Tissue-specific enhancer-gene maps from multimodal single-cell data identify causal disease alleles. *bioRxiv* (2022) doi:10.1101/2022.10.27.22281574.
  68. Weeks, E. M. *et al.* Leveraging polygenic enrichments of gene features to predict genes underlying complex traits and diseases. *Nat. Genet.* (2023) doi:10.1038/s41588-023-01443-6..
  69. Nasser, J. *et al.* Genome-wide enhancer maps link risk variants to disease genes. *Nature* **593**, 238–243 (2021).
  70. Schnitzler, G. R. *et al.* Mapping the convergence of genes for coronary artery disease onto endothelial cell programs. *bioRxiv* 2022.11.01.514606 (2022) doi:10.1101/2022.11.01.514606.
  71. Finucane, H. K. *et al.* Heritability enrichment of specifically expressed genes identifies disease-relevant tissues and cell types. *Nat. Genet.* **50**, 621–629 (2018).
  72. Forgetta, V. *et al.* An effector index to predict target genes at GWAS loci. *Hum. Genet.* **141**, 1431–1447 (2022).
  73. Ghousaini, M. *et al.* Open Targets Genetics: systematic identification of trait-associated genes using large-scale genetics and functional genomics. *Nucleic Acids Res.* **49**, D1311–D1320 (2021).
  74. Kryshtafovych, A., Schwede, T., Topf, M., Fidelis, K. & Moult, J. Critical assessment of methods of protein structure prediction (CASP)-Round XIV. *Proteins* **89**, 1607–1617 (2021).
  75. The Critical Assessment of Genome Interpretation Consortium. CAGI, the Critical Assessment of Genome Interpretation, establishes progress and prospects for computational genetic variant interpretation methods. *arXiv [q-bio.GN]* (2022).
  76. Hogan, A. *et al.* Knowledge Graphs. *arXiv [cs.AI]* (2020).
  77. Feng, F. *et al.* GenomicKB: a knowledge graph for the human genome. *Nucleic Acids Res.* **51**, D950–D956 (2023).
  78. Chandak, P., Huang, K. & Zitnik, M. Building a knowledge graph to enable precision medicine. *Sci Data* **10**, 67 (2023).

79. Karollus, A., Mauermeier, T. & Gagneur, J. Current sequence-based models capture gene expression determinants in promoters but mostly ignore distal enhancers. *Genome Biol.* **24**, 56 (2023).
80. Ambrosini, G. *et al.* Insights gained from a comprehensive all-against-all transcription factor binding motif benchmarking study. *Genome Biol.* **21**, 114 (2020).
81. Moore, J. E., Pratt, H. E., Purcaro, M. J. & Weng, Z. A curated benchmark of enhancer-gene interactions for evaluating enhancer-target gene prediction methods. *Genome Biol.* **21**, 17 (2020).
82. Inoue, F. *et al.* A systematic comparison reveals substantial differences in chromosomal versus episomal encoding of enhancer activity. *Genome Res.* **27**, 38–52 (2017).
83. Xie, S., Duan, J., Li, B., Zhou, P. & Hon, G. C. Multiplexed Engineering and Analysis of Combinatorial Enhancer Activity in Single Cells. *Mol. Cell* **66**, 285–299.e5 (2017).
84. Gasperini, M. *et al.* A Genome-wide Framework for Mapping Gene Regulation via Cellular Genetic Screens. *Cell* **176**, 1516 (2019).
85. Reilly, S. K. *et al.* Direct characterization of cis-regulatory elements and functional dissection of complex genetic associations using HCR–FlowFISH. *Nat. Genet.* **53**, 1166–1176 (2021).
86. Schraivogel, D. *et al.* Targeted Perturb-seq enables genome-scale genetic screens in single cells. *Nat. Methods* **17**, 629–635 (2020).
87. White, M. A. Understanding how cis-regulatory function is encoded in DNA sequence using massively parallel reporter assays and designed sequences. *Genomics* **106**, 165–170 (2015).
88. Ma, S. *et al.* Chromatin Potential Identified by Shared Single-Cell Profiling of RNA and Chromatin. *Cell* **183**, 1103–1116.e20 (2020).
89. McGinnis, C. S. *et al.* MULTI-seq: sample multiplexing for single-cell RNA sequencing using lipid-tagged indices. *Nat. Methods* **16**, 619–626 (2019).
90. Daniel, B. *et al.* Divergent clonal differentiation trajectories of T cell exhaustion. *Nat. Immunol.* **23**, 1614–1627 (2022).
91. Rebboah, E. *et al.* Mapping and modeling the genomic basis of differential RNA isoform expression at single-cell resolution with LR-Split-seq. *Genome Biol.* **22**, 286 (2021).
92. Pratapa, A., Jalihal, A. P., Law, J. N., Bharadwaj, A. & Murali, T. M. Benchmarking algorithms for gene regulatory network inference from single-cell transcriptomic data. *Nat. Methods* **17**, 147–154 (2020).
93. Sherry, S. T. *et al.* dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res.* **29**, 308–311 (2001).
94. The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
95. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res.* **47**, D766–D773 (2019).
96. UniProt: the universal protein knowledgebase in 2023. *Nucleic Acids Res.* **51**, D523–D531 (2023).
97. Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat. Genet.* **25**, 25–29 (2000).
98. Del Toro, N. *et al.* The IntAct database: efficient access to fine-grained molecular interaction data.

- Nucleic Acids Res.* **50**, D648–D653 (2022).
99. Vasilevsky, N. A. *et al.* Mondo: Unifying diseases for the world, by the world. *bioRxiv* (2022) doi:10.1101/2022.04.13.22273750.
  100. Köhler, S. *et al.* The Human Phenotype Ontology in 2021. *Nucleic Acids Res.* **49**, D1207–D1217 (2021).
  101. Amberger, J. S., Bocchini, C. A., Scott, A. F. & Hamosh, A. OMIM.org: leveraging knowledge across phenotype-gene relationships. *Nucleic Acids Res.* **47**, D1038–D1043 (2019).
  102. Zhang, Y., Qi, G., Park, J.-H. & Chatterjee, N. Estimation of complex effect-size distributions using summary-level statistics from genome-wide association studies across 32 complex traits. *Nat. Genet.* **50**, 1318–1326 (2018).
  103. O'Connor, L. J. The distribution of common-variant effect sizes. *Nat. Genet.* **53**, 1243–1249 (2021).
  104. Lewis, C. M. & Vassos, E. Polygenic risk scores: from research tools to clinical instruments. *Genome Med.* **12**, 44 (2020).
  105. Hekselman, I. & Yeger-Lotem, E. Mechanisms of tissue and cell-type specificity in heritable traits and diseases. *Nat. Rev. Genet.* **21**, 137–150 (2020).
  106. Uffelmann, E. *et al.* Genome-wide association studies. *Nature Reviews Methods Primers* **1**, 1–21 (2021).
  107. Weissbrod, O. *et al.* Leveraging fine-mapping and multipopulation training data to improve cross-population polygenic risk scores. *Nat. Genet.* **54**, 450–458 (2022).
  108. Heid, I. M. *et al.* Meta-analysis identifies 13 new loci associated with waist-hip ratio and reveals sexual dimorphism in the genetic basis of fat distribution. *Nat. Genet.* **42**, 949–960 (2010).
  109. Goossens, G. H., Jocken, J. W. E. & Blaak, E. E. Sexual dimorphism in cardiometabolic health: the role of adipose tissue, muscle and liver. *Nat. Rev. Endocrinol.* **17**, 47–66 (2021).
  110. Rajabli, F. *et al.* Ancestral origin of ApoE  $\epsilon$ 4 Alzheimer disease risk in Puerto Rican and African American populations. *PLoS Genet.* **14**, e1007791 (2018).
  111. Blue, E. E., Horimoto, A. R. V. R., Mukherjee, S., Wijsman, E. M. & Thornton, T. A. Local ancestry at APOE modifies Alzheimer's disease risk in Caribbean Hispanics. *Alzheimers. Dement.* **15**, 1524–1532 (2019).
  112. Baxter, S. M. *et al.* Centers for Mendelian Genomics: A decade of facilitating gene discovery. *Genet. Med.* **24**, 784–797 (2022).
  113. Costanzo, M. C. *et al.* The Type 2 Diabetes Knowledge Portal: An open access genetic resource dedicated to type 2 diabetes and related traits. *Cell Metab.* **35**, 695–710.e6 (2023).
  114. Richards, S. *et al.* Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet. Med.* **17**, 405–424 (2015).
  115. Scott, A. *et al.* Saturation-scale functional evidence supports clinical variant interpretation in Lynch syndrome. *Genome Biol.* **23**, 266 (2022).
  116. Radford, E. J. *et al.* Saturation genome editing of DDX3X clarifies pathogenicity of germline and somatic variation. *medRxiv* 2022.06.10.22276179 (2022) doi:10.1101/2022.06.10.22276179.

117. Wojcik, M. H. *et al.* Beyond the exome: what's next in diagnostic testing for Mendelian conditions. *ArXiv* (2023) doi:10.1002/ajmg.a.63053.
118. Miller, D. T. *et al.* ACMG SF v3.1 list for reporting of secondary findings in clinical exome and genome sequencing: A policy statement of the American College of Medical Genetics and Genomics (ACMG). *Genet. Med.* **24**, 1407–1414 (2022).
119. Pejaver, V. *et al.* Calibration of computational tools for missense variant pathogenicity classification and ClinGen recommendations for PP3/BP4 criteria. *Am. J. Hum. Genet.* **109**, 2163–2177 (2022).
120. Graham, S. E. *et al.* The power of genetic diversity in genome-wide association studies of lipids. *Nature* **600**, 675–679 (2021).
121. Musunuru, K. & Kathiresan, S. Genetics of Common, Complex Coronary Artery Disease. *Cell* **177**, 132–145 (2019).
122. Hamilton, M. C. *et al.* Systematic elucidation of genetic mechanisms underlying cholesterol uptake. *Cell Genom* **3**, 100304 (2023).
123. Martin, A. R. *et al.* Clinical use of current polygenic risk scores may exacerbate health disparities. *Nat. Genet.* **51**, 584–591 (2019).
124. Shi, H. *et al.* Population-specific causal disease effect sizes in functionally important regions impacted by selection. *Nat. Commun.* **12**, 1098 (2021).
125. Hou, K. *et al.* Causal effects on complex traits are similar for common variants across segments of different continental ancestries within admixed individuals. *Nat. Genet.* **55**, 549–558 (2023).
126. Threadgill, D. W., Miller, D. R., Churchill, G. A. & de Villena, F. P.-M. The collaborative cross: a recombinant inbred mouse population for the systems genetic era. *ILAR J.* **52**, 24–31 (2011).
127. Sherry, S. T., Ward, M. & Sirotkin, K. dbSNP-database for single nucleotide polymorphisms and other classes of minor genetic variation. *Genome Res.* **9**, 677–679 (1999).
128. Zhou, H. *et al.* FAVOR: functional annotation of variants online resource and annotator for variation across the human genome. *Nucleic Acids Res.* **51**, D1300–D1311 (2023).
129. Huang, L. *et al.* TOP-LD: A tool to explore linkage disequilibrium with TOPMed whole-genome sequence data. *Am. J. Hum. Genet.* **109**, 1175–1181 (2022).
130. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data* **3**, 160018 (2016).
131. Schatz, M. C. *et al.* Inverting the model of genomics data sharing with the NHGRI Genomic Data Science Analysis, Visualization, and Informatics Lab-space. *Cell Genom* **2**, (2022).
132. Fowler, D. M. *et al.* An Atlas of Variant Effects to understand the genome at nucleotide resolution. *Genome Biol.* **24**, 147 (2023).
133. Wang, T. *et al.* The Human Pangenome Project: a global resource to map genomic diversity. *Nature* **604**, 437–446 (2022).
134. All of Us Research Program Investigators *et al.* The 'All of Us' Research Program. *N. Engl. J. Med.* **381**, 668–676 (2019).

## Author List

**Writing group (ordered by contribution):** Jesse M. Engreitz, Heather A. Lawson, Harinder Singh, Lea M. Starita, Gary C. Hon, Hannah Carter, Nidhi Sahni, Timothy E. Reddy, Xihong Lin, Yun Li, Nikhil V. Munshi, Maria H. Chahrour, Alan P. Boyle, Benjamin C. Hitz, Ali Mortazavi, Mark Craven, Karen L. Mohlke, Luca Pinello, Ting Wang

**Steering Committee Co-Chairs (alphabetical by last name):** Anshul Kundaje, Karen L. Mohlke, Feng Yue

**Working Group and Focus Group Co-Chairs (alphabetical by last name):**

**Catalog:** Michael I. Love, Lea M. Starita, Feng Yue

**Characterization:** Gary C. Hon, Martin Kircher, Timothy E. Reddy

**Computational Analysis, Modeling, and Prediction:** Xihong Lin, Jian Ma, Predrag Radivojac

**Project Design:** Brunilda Balliu, Jesse M. Engreitz, Nidhi Sahni

**Mapping:** Nina P. Farrell, Brian A. Williams

**Networks:** Hannah Carter, Danwei Huangfu

**Standards and Pipelines:** Anshul Kundaje, Luca Pinello

**Cardiometabolic:** Nikhil V. Munshi, Chong Y. Park, Thomas Quertermous

**Cellular Programs and Networks:** Hannah Carter, Jishnu Das

**Coding Variants:** Michael A. Calderwood, Douglas M. Fowler, Predrag Radivojac, Lea M. Starita, Marc Vidal

**CRISPR:** Lucas Ferreira, Luca Pinello

**Defining and Systematizing Function:** Mark Craven, Sean D. Mooney, Vikas Pejaver

**Enumerating Variants:** Benjamin C. Hitz, Jingjing Zhao

**Evolution:** Steven Gazal, Evan Koch, Steven K. Reilly, Shamil Sunyaev

**Imaging:** Anne E. Carpenter,

**Immune:** Jason D. Buenrostro, Christina S. Leslie, Rachel E. Savage

**Impact on Diverse Populations:** Stefanija Giric, Yun Li

**iPSC:** Chongyuan Luo, Kathrin Plath

**MPRA:** Alejandro Barrera, Michael I. Love, Max Schubach

**Noncoding Variants:** Jesse M. Engreitz, Jill E. Moore, Nidhi Sahni

**Neuro:** Nadav Ahituv, Maria H. Chahrour

**Phenotypic Impact and Function:** Kushal Dey, Xihong Lin

**QTL/Statgen:** Brunilda Balliu, Ingileif Hallgrimsdottir, Kyle Gaulton, Saori Sakaue

**Single Cell:** Sina Boeshaghi, Anshul Kundaje, Eugenio Mattei, Ali Mortazavi, Surag Nair, Lior Pachter, Austin Wang

**Characterization Awards (contact PI, MPIs (alphabetical by last name), other members (alphabetical by last name)):**

**UM1HG011966:** Jay Shendure<sup>1,5,171,172,173</sup>, Nadav Ahituv<sup>2</sup>, Martin Kircher<sup>3,4</sup>, Vikram Agarwal<sup>1,174</sup>, Andrew Blair<sup>2</sup>, Theofilos Chalkiadakis<sup>4</sup>, Florence M. Chardon<sup>1</sup>, Pyaree M Dash<sup>4</sup>, Chengyu Deng<sup>2</sup>, Nobuhiko Hamazaki<sup>1</sup>, Pia Keukeleire<sup>3</sup>, Connor Kubo<sup>1</sup>, Jean-Benoît Lalanne<sup>1</sup>, Thorben Maass<sup>3</sup>, Beth Martin<sup>1</sup>, Troy McDiarmid<sup>1</sup>, Mai Nobuhara<sup>2</sup>, Nicholas F Page<sup>2</sup>, Sam Regalado<sup>1</sup>, Max Schubach<sup>4</sup>, Jasmine Sims<sup>2</sup>, Aki Ushiki<sup>2</sup>, Jingjing Zhao<sup>2</sup>

**UM1HG011969:** Lea M. Starita<sup>1,5</sup>, Douglas M. Fowler<sup>1,5</sup>, Sabrina M. Best<sup>1</sup>, Gabe Boyle<sup>1</sup>, Nathan Camp<sup>6</sup>, Silvia Casadei<sup>1</sup>, Estelle Y. Da<sup>7</sup>, Moez Dawood<sup>5,8</sup>, Samantha C. Dawson<sup>6</sup>, Shawn Fayer<sup>1</sup>, Audrey Hamm<sup>1</sup>, Richard G. James<sup>6</sup>, Gail P. Jarvik<sup>1</sup>, Abbye E. McEwen<sup>1,5,9</sup>, Nick Moore<sup>7</sup>, Lara A. Muffley<sup>1</sup>, Sriram Pendyala<sup>1</sup>, Nicholas A. Popp<sup>1</sup>, Mason Post<sup>1</sup>, Alan F. Rubin<sup>7</sup>, Jay Shendure<sup>1,5,171,172,173</sup>, Nahum T. Smith<sup>1</sup>, Alan F. Rubin<sup>1</sup>, Jeremy Stone<sup>5</sup>, Malvika Tejura<sup>1</sup>, Ziyu R. Wang<sup>1</sup>, Melinda K. Wheelock<sup>1</sup>, Ivan Woo<sup>1</sup>, Brendan D. Zapp<sup>1</sup>

**UM1HG011972:** Jesse M. Engreitz<sup>10,11,12,61</sup>, Thomas Quertermous<sup>13</sup>, Dulguun Amgalan<sup>10,11</sup>, Aradhana Aradhana<sup>10</sup>, Sophia M. Arana<sup>10</sup>, Michael C. Bassik<sup>10</sup>, Julia R. Bauman<sup>10</sup>, Asmita Bhattacharya<sup>10</sup>,

Xiangmeng Shawn Cai<sup>10,11,22</sup>, Ziwei Chen<sup>14</sup>, Stephanie Conley<sup>10,11</sup>, Salil Deshpande<sup>15</sup>, Benjamin R. Doughty<sup>10</sup>, Peter P. Du<sup>10</sup>, Casey Gifford<sup>10,11,16,17</sup>, William J. Greenleaf<sup>10,161</sup>, Andreas R. Gschwind<sup>10</sup>, Katherine Guo<sup>10,11</sup>, Sarasa Isobe<sup>18</sup>, Evelyn Jagoda<sup>12,13</sup>, Nimit Jain<sup>10</sup>, Hank Jones<sup>10,11</sup>, Helen Y. Kang<sup>10,11</sup>, Samuel H. Kim<sup>162</sup>, YeEun Kim<sup>163</sup>, Sandy Klemm<sup>10</sup>, Anshul Kundaje<sup>10,14</sup>, Soumya Kundu<sup>14</sup>, Mauro Lago-Docampo<sup>18</sup>, Yannick C. Lee-Yow<sup>10,11</sup>, Roni Levin-Konigsberg<sup>10</sup>, Daniel Y. Li<sup>13</sup>, Dominik Lindenhofer<sup>19</sup>, X. Rosa Ma<sup>10,11</sup>, Georgi K. Marinov<sup>10</sup>, Gabriella E. Martyn<sup>10,11</sup>, Eyal Metzl-Raz<sup>10</sup>, Joao P. Monteiro<sup>13</sup>, Michael T. Montgomery<sup>10,11</sup>, Kristy S. Mualim<sup>11,20,164</sup>, Chad Munger<sup>10,11</sup>, Glen Munson<sup>12</sup>, Tri C. Nguyen<sup>10,11</sup>, Trieu Nguyen<sup>13</sup>, Brian T. Palmisano<sup>13</sup>, Anusri Pampari<sup>14</sup>, Chong Y. Park<sup>13</sup>, Marlene Rabinovitch<sup>18</sup>, Markus Ramste<sup>13</sup>, Judhajeet Ray<sup>12</sup>, Kevin R. Roy<sup>10,21</sup>, Oriane M. Rubio<sup>11</sup>, Julia M. Schaepe<sup>22</sup>, Gavin Schnitzler<sup>13</sup>, Jacob Schreiber<sup>10</sup>, Disha Sharma<sup>13</sup>, Maya U. Sheth<sup>10,11,22</sup>, Huitong Shi<sup>13</sup>, Vasundhara Singh<sup>12</sup>, Riya Sinha<sup>23</sup>, Lars M. Steinmetz<sup>10,19,21</sup>, Jason Tan<sup>10,11,14</sup>, Anthony Tan<sup>10,11</sup>, Josh Tycko<sup>10</sup>, Raeline C. Valbuena<sup>10</sup>, Valeh Valiollah Pour Amiri<sup>10</sup>, Mariëlle J.F.M. van Kooten<sup>10</sup>, Alun Vaughan-Jackson<sup>10</sup>, Anthony Venida<sup>10</sup>, Chad S. Weldy<sup>13</sup>, David Yao<sup>10</sup>, Tony Zeng<sup>10,11</sup>, Ronghao Zhou<sup>10,11</sup>

**UM1HG011989:** Marc Vidal<sup>24,25</sup>, Michael A. Calderwood<sup>24,25,26</sup>, Anne E. Carpenter<sup>27</sup>, Beth A. Cimini<sup>27</sup>, Georges Coppin<sup>24,25,26,28</sup>, Atina G. Coté<sup>29,30,31</sup>, Marzieh Haghghi<sup>27</sup>, Tong Hao<sup>24,25,26</sup>, David E. Hill<sup>24,25,26</sup>, Jessica Lacoste<sup>29,30</sup>, Florent Laval<sup>24,25,26,28,184</sup>, Chloe Reno<sup>29,30</sup>, Frederick P. Roth<sup>29,30,31,32</sup>, Shantanu Singh<sup>27</sup>, Kerstin Spirohn-Fitzgerald<sup>24,25</sup>, Mikko Taipale<sup>29,30</sup>, Tanisha Teelucksingh<sup>29</sup>, Maxime Tixhon<sup>24,25,26,185</sup>, Anupama Yadav<sup>24,25,26</sup>, Zhipeng Yang<sup>24,25,26</sup>

**UM1HG011996:** Gary C. Hon<sup>33,34,35</sup>, W. Lee Kraus<sup>33,34</sup>, Nikhil V. Munshi<sup>36,37</sup>, Daniel A. Armendariz<sup>33</sup>, Maria H. Chahrour<sup>38,211,212,213,214</sup>, Ashley E. Dederich<sup>39</sup>, Lauretta El Hayek<sup>38</sup>, Sean C. Goetsch<sup>36</sup>, Kiran Kaur<sup>38</sup>, Hyung Bum Kim<sup>33</sup>, Melissa K. McCoy<sup>39</sup>, Mpathi Z. Nzima<sup>33</sup>, Carlos A. Pinzón-Arteaga<sup>40</sup>, Bruce A. Posner<sup>39</sup>, Daniel A. Schmitz<sup>37</sup>, Sushama Sivakumar<sup>36,37</sup>, Anjana Sundarraj<sup>33</sup>, Lei Wang<sup>33</sup>, Yihan Wang<sup>33</sup>, Jun Wu<sup>37</sup>, Lin Xu<sup>40,41</sup>, Jian Xu<sup>42</sup>, Leqian Yu<sup>37</sup>, Yanfeng Zhang<sup>40</sup>, Huan Zhai<sup>33</sup>, Qinbo Zhou<sup>40</sup>

**UM1HG012003:** Hyejung Won<sup>43,44</sup>, Michael I. Love<sup>43,204</sup>, Karen L. Mohlke<sup>43</sup>, Jessica L. Bell<sup>43,44</sup>, K. Elaine Broadaway<sup>43</sup>, Katherine N. Degner<sup>43,44</sup>, Amy S. Etheridge<sup>43</sup>, Stefanija Giric<sup>43</sup>, Beverly H. Koller<sup>43</sup>, Yun Li<sup>43,204</sup>, Won Mah<sup>43,44</sup>, Wancen Mu<sup>204</sup>, Kimberly D. Ritola<sup>44,205</sup>, Jonathan D. Rosen<sup>43</sup>, Sarah A. Schoenrock<sup>43,44</sup>, Rachel A. Sharp<sup>43,44</sup>

**UM1HG012010:** Luca Pinello<sup>45,61</sup>, Daniel Bauer<sup>47,48</sup>, Guillaume Lettre<sup>49,50</sup>, Richard Sherwood<sup>51</sup>, Basheer Becerra<sup>47,48</sup>, Logan J. Blaine<sup>45,52</sup>, Lucas Ferreira<sup>52,53</sup>, Matthew J. Francoeur<sup>51</sup>, Ellie N. Gibbs<sup>51</sup>, Nahye Kim<sup>45,54,217</sup>, Emily M. King<sup>45,54,217,218</sup>, Benjamin P. Kleinstiver<sup>45,54,217</sup>, Estelle Lecluze<sup>49</sup>, Zhijian Li<sup>45,46</sup>, Zain M. Patel<sup>45,46</sup>, Quang Vinh Phan<sup>51</sup>, Jayoung Ryu<sup>45,52</sup>, Marlena L. Starr<sup>53</sup>, Ting Wu<sup>48,53</sup>

**UM1HG012053:** Charles A. Gersbach<sup>55,56</sup>, Gregory E. Crawford<sup>56,57</sup>, Timothy E. Reddy<sup>58</sup>, Andrew S. Allen<sup>58</sup>, William H. Majoros<sup>58</sup>, Nahid Iglesias<sup>55,56</sup>, Alejandro Barrera<sup>56,58</sup>, Ruhi Rai<sup>56</sup>, Revathy Venukuttan<sup>56</sup>, Boxun Li<sup>55,56</sup>, Taylor Anglen<sup>56,59</sup>, Lexi R. Bounds<sup>55,56</sup>, Marisa C. Hamilton<sup>56</sup>, Siyan Liu<sup>56</sup>, Sean R. McCutcheon<sup>55,56</sup>, Christian D. McRoberts Amador<sup>56,60</sup>, Samuel J. Reisman<sup>56,59</sup>, Maria A. ter Weele<sup>55,56</sup>, Josephine C. Bodle<sup>55,56</sup>, Helen L. Streff<sup>55,56</sup>, Keith Siklenka<sup>58</sup>, Kari Strouse<sup>58</sup>

**Mapping Awards (contact PI, MPis (alphabetical by last name), other members (alphabetical by last name)):**

**UM1HG011986:** Jason D. Buenrostro<sup>61,62</sup>, Bradley E. Bernstein<sup>61,63</sup>, Juliana Babu<sup>61,62</sup>, Guillermo Barreto Corona<sup>61</sup>, Kevin Dong<sup>61</sup>, Fabiana M. Duarte<sup>61,62</sup>, Neva C. Durand<sup>61</sup>, Charles B. Epstein<sup>61</sup>, Kaili Fan<sup>61,62,98</sup>, Nina P. Farrell<sup>61</sup>, Elizabeth Gaskell<sup>61</sup>, Amelia W. Hall<sup>61</sup>, Alexandra M. Ham<sup>61</sup>, Mei K. Knudson<sup>61</sup>, Eugenio Mattei<sup>61</sup>, Rachel E. Savage<sup>61,62</sup>, Noam Shores<sup>61</sup>, Siddarth Wekhande<sup>61</sup>, Cassandra M. White<sup>61</sup>, Wang Xi<sup>61,62</sup>

**UM1HG012076:** Ansuman T. Satpathy<sup>64,65,66</sup>, M. Ryan Corces<sup>73,74,187</sup>, Serena H. Chang<sup>73,74,187</sup>, Iris M. Chin<sup>73,74,187</sup>, James M. Gardner<sup>75,76</sup>, Zachary A. Gardell<sup>73,74,187</sup>, Jacob C. Gutierrez<sup>64,66</sup>, Alia W. Johnson<sup>73,74,187</sup>, Lucas Kampman<sup>73,74,187</sup>, Maya Kasowski<sup>64,72</sup>, Caleb A. Lareau<sup>64,65,66</sup>, Vincent Liu<sup>64,66</sup>, Leif S. Ludwig<sup>67,68</sup>, Christopher S. McGinnis<sup>64,65,66</sup>, Shreya Menon<sup>73,74,187</sup>, Adam W. Turner<sup>73,74,187</sup>, Chun J. Ye<sup>69,70,71</sup>, Yajie Yin<sup>64,66</sup>, Wenxi Zhang<sup>64</sup>

**UM1HG012077:** Ali Mortazavi<sup>188,189</sup>, Barbara J. Wold<sup>190,191</sup>, Sina Boeshaghi<sup>190</sup>, Maria Carilli<sup>192</sup>, Dayeon Cheong<sup>188</sup>, Ghassan Filibam<sup>188</sup>, Kim Green<sup>188,193</sup>, Ingileif Hallgrimsdottir<sup>190</sup>, Shimako Kawauchi<sup>189</sup>, Charlene Kim<sup>190</sup>, Heidi Liang<sup>189</sup>, Rebekah Loving<sup>190</sup>, Laura Luebbert<sup>190</sup>, Grant MacGregor<sup>188</sup>, Angel G Merchan<sup>190</sup>, Lior Pachter<sup>190,194</sup>, Elisabeth Rebboah<sup>188</sup>, Fairlie Reese<sup>188,189</sup>, Narges Rezaie<sup>188,189</sup>, Jasmine Sakr<sup>189,195</sup>, Delaney Sullivan<sup>190</sup>, Nikki Swarna<sup>192</sup>, Diane Trout<sup>190</sup>, Sean Upchurch<sup>190</sup>, Ryan Weber<sup>188</sup>, Brian A. Williams<sup>190</sup>

**Predictive Modeling Awards (contact PI, MPis (alphabetical by last name), other members (alphabetical by last name)):**

**U01HG011952:** Alan P. Boyle<sup>196,197</sup>, Christopher P. Castro<sup>196</sup>, Elysia Chou<sup>196</sup>, Fan Feng<sup>196</sup>, Andre Guerra<sup>198</sup>, Yuanhao Huang<sup>196</sup>, Linghua Jiang<sup>196</sup>, Jie Liu<sup>196</sup>, Ryan E. Mills<sup>196,197</sup>, Weizhou Qian<sup>196</sup>, Tingting Qin<sup>196</sup>, Maureen A. Sartor<sup>196,198</sup>, Rintsen N. Sherpa<sup>196</sup>, Jinhao Wang<sup>196</sup>, Yiqun Wang<sup>196</sup>, Joshua D. Welch<sup>196</sup>, Zhenhao Zhang<sup>196</sup>, Nanxiang Zhao<sup>196</sup>

**U01HG011967:** Andrew S. Allen<sup>58</sup>, Sayan Mukherjee<sup>77,78,79</sup>, C. David Page<sup>58</sup>, Shannon Clarke<sup>58</sup>, Richard W. Doty<sup>58</sup>, Yuncheng Duan<sup>80</sup>, Raluca Gordan<sup>58,79</sup>, Kuei-Yueh Ko<sup>58</sup>, Shengyu Li<sup>58</sup>, Boyao Li<sup>58</sup>, William H. Majoros<sup>58</sup>, Timothy E. Reddy<sup>58</sup>, Alexander Thomson<sup>58</sup>

**U01HG012009:** Soumya Raychaudhuri<sup>51,82</sup>, Alkes Price<sup>83,84,82</sup>, Shamil Sunyaev<sup>51,52</sup>, Thahmina A. Ali<sup>81</sup>, Kushal K. Dey<sup>81,83</sup>, Arun Durvasula<sup>83,85</sup>, Manolis Kellis<sup>46,220</sup>, Evan Koch<sup>52</sup>, Saori Sakaue<sup>51,82</sup>

**U01HG012022:** Predrag Radivojac<sup>86</sup>, Lilia M. Iakoucheva<sup>87</sup>, Tulika Kakati<sup>87</sup>, Sean D. Mooney<sup>88</sup>, Yile Chen<sup>88</sup>, Mariam Benazouz<sup>88</sup>, Vikas Pejaver<sup>89,90</sup>, Shantanu Jain<sup>86,215</sup>, Daniel Zeiberg<sup>86</sup>, M. Clara De Paolis Kaluza<sup>86</sup>, Michelle Velyunskiy<sup>86</sup>

**U01HG012039:** Mark Craven<sup>91</sup>, Audrey Gasch<sup>92</sup>, Kunling Huang<sup>93</sup>, Yiyang Jin<sup>91</sup>, Qiongshi Lu<sup>91</sup>, Jiacheng Miao<sup>91</sup>, Michael Ohtake<sup>94</sup>, Eduardo Scopel<sup>92</sup>, Robert D. Steiner<sup>95,96,97</sup>, Yuriy Sverchkov<sup>91</sup>

**U01HG012064:** Zhiping Weng<sup>98</sup>, Manuel Garber<sup>98</sup>, Xihong Lin<sup>84,100</sup>, Yu Fu<sup>98</sup>, Natalie Haas<sup>98</sup>, Xihao Li<sup>43,84,204</sup>, Nishigandha Phalke<sup>98</sup>, Shuo C. Shan<sup>98</sup>, Nicole Shedd<sup>98</sup>, Eric Van Buren<sup>84</sup>, Tianxiang Yu<sup>98</sup>, Yi Zhang<sup>101</sup>, Hufeng Zhou<sup>84</sup>

**U01HG012064:** Anshul Kundaje<sup>10,14</sup>, Alexis Battle<sup>102,103,104,105</sup>, Ziwei Chen<sup>14</sup>, Salil Deshpande<sup>15</sup>, Jesse M. Engreitz<sup>10,11,12,61</sup>, Livnat Jerby<sup>10</sup>, Eran Kotler<sup>10</sup>, Soumya Kundu<sup>10,14</sup>, Andrew R. Marderstein<sup>64</sup>, Georgi K. Marinov<sup>10</sup>, Stephen B. Montgomery<sup>10,64,106</sup>, Surag Nair<sup>14</sup>, AkshatKumar Nigam<sup>10,14</sup>, Evin M. Padhi<sup>64</sup>, Anusri Pampari<sup>14</sup>, Aman Patel<sup>14</sup>, Jonathan Pritchard<sup>10</sup>, Ivy Raine<sup>10</sup>, Vivekanandan Ramalingam<sup>10</sup>, Kameron Rodrigues<sup>64</sup>, Jacob M. Schreiber<sup>10</sup>, Arpita Singhal<sup>14</sup>, Riya Sinha<sup>15</sup>, Valeh Valiollah Pour Amiri<sup>10</sup>, Austin T. Wang<sup>14</sup>

**Network Projects (contact PI, MPis (alphabetical by last name), other members (alphabetical by last name)):**

**U01HG012041:** Harinder Singh<sup>107</sup>, Jishnu Das<sup>107</sup>, Nidhi Sahni<sup>108,109,110</sup>, Marisa Abundis<sup>111</sup>, Deepa Bisht<sup>112</sup>, Trirupa Chakraborty<sup>107</sup>, Jingyu Fan<sup>107</sup>, David R. Hall<sup>107</sup>, Zarifeh H. Rarani<sup>107</sup>, Abhinav Jain<sup>112</sup>, Babita Kaundal<sup>112</sup>, Swapnil Keshari<sup>107</sup>, Daniel McGrail<sup>113,114</sup>, Nicholas A. Pease<sup>107</sup>, Vivian F. Yi<sup>107</sup>, S. Stephen Yi<sup>115,116</sup>

**U01HG012047:** Hao Wu<sup>117</sup>, Sreeram Kannan<sup>118</sup>, Hongjun Song<sup>119</sup>, Jingli Cai<sup>120</sup>, Ziyue Gao<sup>117</sup>, Ronni Kurzion<sup>119</sup>, Julia I. Leu<sup>117</sup>, Fan Li<sup>117</sup>, Dongming Liang<sup>117</sup>, Guo-li Ming<sup>119</sup>, Kiran Musunuru<sup>120</sup>, Qi Qiu<sup>117</sup>, Junwei Shi<sup>121</sup>, Yijing Su<sup>119</sup>, Sarah Tishkoff<sup>117</sup>, Ning Xie<sup>117</sup>, Qian Yang<sup>119</sup>, Wenli Yang<sup>120</sup>, Hongjie Zhang<sup>117</sup>, Zhijian Zhang<sup>119</sup>

**U01HG012051:** Danwei Huangfu<sup>122,123</sup>, Michael A. Beer<sup>124</sup>, Ronald Cutler<sup>125</sup>, Rachel A. Glenn<sup>122,123,126</sup>, Renhe Luo<sup>122,123</sup>, Jin Woo Oh<sup>124</sup>, Milad Razavi-Mohseni<sup>124</sup>, Dustin Shigaki<sup>124</sup>, Simone Sidoli<sup>125</sup>, Thomas Vierbuchen<sup>122,123</sup>, Jieli Yan<sup>122,123</sup>, Yunxiao Yang<sup>124</sup>

**U01HG012059:** Maike Sander<sup>127</sup>, Hannah Carter<sup>128</sup>, Kyle J. Gaulton<sup>127</sup>, Bing Ren<sup>129,130</sup>, Weronika Bartosik<sup>129</sup>, Hannah S. Indralingam<sup>129</sup>, Adam Klie<sup>131</sup>, Hannah Mummey<sup>131</sup>, Mei-Lin Okino<sup>132</sup>, Gaowei Wang<sup>127</sup>, Nathan R. Zemke<sup>129</sup>, Kai Zhang<sup>129</sup>, Han Zhu<sup>127</sup>

**U01HG012079:** Chongyuan Luo<sup>133</sup>, Kathrin Plath<sup>134</sup>, Noah Zaitlen<sup>135</sup>, Brunilda Balliu<sup>136,137,138</sup>, Jason Ernst<sup>134,137</sup>, Justin Langerman<sup>134</sup>, Terence Li<sup>133</sup>, Yu Sun<sup>134</sup>

**U01HG012103:** Christina S. Leslie<sup>199</sup>, Alexander Y. Rudensky<sup>200,201</sup>, Preethi K. Periyakoil<sup>199</sup>, Vianne R. Gao<sup>199</sup>, Melanie H. Smith<sup>202</sup>, Norman M. Thomas<sup>199</sup>, Laura T. Donlin<sup>202,203</sup>, Amit Lakhnanpal<sup>202</sup>, Kaden M. Southard<sup>199</sup>, Rico C. Ardy<sup>199</sup>

**Data and Administrative Coordinating Center Awards (contact PI, MPIs (alphabetical by last name), other members (alphabetical by last name)):**

**U24HG012103:** J. Michael Cherry<sup>10</sup>, Mark B. Gerstein<sup>166,167,168,169,170</sup>, Kalina Andreeva<sup>10</sup>, Pedro R. Assis<sup>10</sup>, Beatrice Borsari<sup>166,167</sup>, Eric Douglass<sup>10</sup>, Shengcheng Dong<sup>10</sup>, Idan Gabdank<sup>10</sup>, Keenan Graham<sup>10</sup>, Benjamin C. Hitz<sup>10</sup>, Otto Jolanki<sup>10</sup>, Jennifer Jou<sup>10</sup>, Meenakshi S. Kagda<sup>10</sup>, Jin-Wook Lee<sup>10</sup>, Mingjie Li<sup>10</sup>, Khine Lin<sup>10</sup>, Stuart R. Miyasato<sup>10</sup>, Joel Rozowsky<sup>166,167</sup>, Corinn Small<sup>10</sup>, Emma Spragins<sup>10</sup>, Forrest Y. Tanaka<sup>10</sup>, Ian M. Whaling<sup>10</sup>, Ingrid A. Youngworth<sup>10</sup>, Cricket A. Sloan<sup>10</sup>

**U24HG012103:** Ting Wang<sup>139,140</sup>, Feng Yue<sup>175,176</sup>, Eddie Belter<sup>140</sup>, Xintong Chen<sup>175</sup>, Rex L. Chisholm<sup>178</sup>, Sarah Cody<sup>140</sup>, Patricia Dickson<sup>180</sup>, Changxu Fan<sup>139</sup>, Lucinda Fulton<sup>140</sup>, Heather A. Lawson<sup>139</sup>, Daofeng Li<sup>139</sup>, Tina Lindsay<sup>140</sup>, Yu Luan<sup>175</sup>, Yuan Luo<sup>179</sup>, Huijue Lyu<sup>175</sup>, Xiaowen Ma<sup>139</sup>, Jian Ma<sup>165</sup>, Juan Macias-Velasco<sup>139</sup>, Karen H. Miga<sup>186</sup>, Kara Quaid<sup>139</sup>, Nathan Stitzel<sup>181</sup>, Barbara E. Stranger<sup>177</sup>, Chad Tomlinson<sup>140</sup>, Juan Wang<sup>175</sup>, Wenjin Zhang<sup>139</sup>, Bo Zhang<sup>182</sup>, Guoyan Zhao<sup>139,183,216</sup>, Xiaoyu Zhuo<sup>139</sup>

**IGVF Affiliate Member Projects (Contact PIs, other members (alphabetical by last name)):**

Kristen Brennand<sup>219</sup>

Alberto Ciccia<sup>210</sup>, Samuel B. Hayward<sup>210</sup>, Jen-Wei Huang<sup>210</sup>, Giuseppe Leuzzi<sup>210</sup>, Angelo Tagliatela<sup>210</sup>, Tanay Thakar<sup>210</sup>, Alina Vaitsiankova<sup>210</sup>

Kushal K. Dey<sup>46,141</sup>, Thahmina A. Ali<sup>141</sup>

Steven Gazal<sup>142,143,144</sup>, Artem Kim<sup>142</sup>

H. Leighton Grimes<sup>209</sup>, Nathan Salomonis<sup>209</sup>

Rajat Gupta<sup>13</sup>, Shi Fang<sup>13</sup>, Vivian Lee-Kim<sup>13</sup>

Matthias Heinig<sup>145,146,147</sup>, Corinna Losert<sup>145,146</sup>

Thouis R. Jones<sup>12</sup>, Elisa Donnard<sup>12</sup>, Maddie Murphy<sup>12</sup>, Elizabeth Roberts<sup>12</sup>, Susie Song<sup>12</sup>

Jill E. Moore<sup>98</sup>

Sara Mostafavi<sup>221,222</sup>, Alexander Sasse<sup>221</sup>, Anna Spiro<sup>221</sup>

Len A. Pennacchio<sup>148,151</sup>, Momoe Kato<sup>148</sup>, Michael Kosicki<sup>148</sup>, Brandon Mannion<sup>148</sup>, Neil Slaven<sup>148</sup>

Axel Visel<sup>148,151</sup>

Katherine S. Pollard<sup>152,153,154</sup>, Siron Drusinsky<sup>152,153</sup>, Sean Whalen<sup>152</sup>

John Ray<sup>1,172,208</sup>, Ingrid A. Harten<sup>172</sup>, Ching-Huang Ho<sup>172</sup>

Steven K. Reilly<sup>223</sup>

Neville E. Sanjana<sup>149,150</sup>, Christina Caragine<sup>149,150</sup>, John A. Morris<sup>149,150</sup>

Davide Seruggia<sup>155,156</sup>, Ana Patricia Kutschat<sup>155,156</sup>, Sandra Wittibschlager<sup>155,156</sup>

Han Xu<sup>108</sup>, Rongjie Fu<sup>108</sup>, Wei He<sup>108</sup>, Liang Zhang<sup>108</sup>

S. Stephen Yi<sup>157,158</sup>, Daniel Osorio<sup>157,158</sup>

**NHGRI Program Management (alphabetical by last name):** Zo Bly<sup>159</sup>, Stephanie Callouri<sup>160,206</sup>, Daniel A. Gilchrist<sup>160</sup>, Carolyn M. Hutter<sup>160</sup>, Stephanie A. Morris<sup>160</sup>, Michael J. Pazin<sup>160</sup>, Ella K. Samer<sup>160,207</sup>

**Affiliations:**

1. Department of Genome Sciences, University of Washington, Seattle, WA, USA
2. Department of Bioengineering and Therapeutic Sciences, Institute for Human Genetics, University of California San Francisco, San Francisco, California, USA



3. Institute of Human Genetics, University Medical Center Schleswig-Holstein, University of Lübeck, 23562 Lübeck, Germany
4. Exploratory Diagnostic Sciences, Berlin Institute of Health at Charité-Universitätsmedizin Berlin, 10117 Berlin, Germany
5. Brotman Baty Institute for Precision Medicine, Seattle, WA., USA
6. Center of immunotherapy and Immunity, Seattle Children's Research Institute, Seattle, WA, USA
7. Bioinformatics Division, WEHI, Parkville, VIC, Australia
8. Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX, USA
9. Department of Laboratory Medicine and Pathology, University of Washington, Seattle, WA, USA
10. Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA
11. Basic Science and Engineering Initiative, Stanford Children's Health, Betty Irene Moore Children's Heart Center, Stanford, CA, USA
12. The Novo Nordisk Foundation Center for Genomic Mechanisms of Disease, Broad Institute of MIT and Harvard, Cambridge, MA, USA
13. Division of Cardiovascular Medicine, School of Medicine, Stanford University
14. Department of Computer Science, Stanford University School of Medicine, Stanford, CA, USA
15. Institute for Computational and Mathematical Engineering, Stanford University, Stanford, CA, USA
16. Department of Pediatrics, Stanford University School of Medicine, Stanford, CA, USA
17. Institute for Stem Cell Biology and Regenerative Medicine, Stanford University School of Medicine, Stanford, CA, USA
18. Division of Pediatric Cardiology and Cardiovascular Institute, Stanford University School of Medicine, Stanford University
19. European Molecular Biology Laboratory (EMBL), Genome Biology Unit, Heidelberg, Germany
20. Department of Biology, Stanford University, Stanford, CA, USA
21. Stanford Genome Technology Center, Palo Alto, CA., USA
22. Department of Bioengineering, Stanford University School of Engineering, Stanford, CA, USA
23. Department of Biomedical Informatics, Stanford University School of Medicine, Stanford, CA, USA
24. Center for Cancer Systems Biology (CCSB), Dana-Farber Cancer Institute, Boston, MA, USA
25. Department of Genetics, Blavatnik Institute, Harvard Medical School, Boston, MA, USA
26. Department of Cancer Biology, Dana-Farber Cancer Institute, Boston, MA, USA
27. Imaging Platform, Broad Institute of Harvard and MIT, Cambridge, Massachusetts
28. Laboratory of Viral Interactomes, GIGA Institute, University of Liège, Liège, Belgium
29. Donnelly Centre for Cellular and Biomolecular Research (CCBR), University of Toronto, Toronto, Ontario, Canada
30. Department of Molecular Genetics, University of Toronto, Toronto, Ontario, Canada.
31. Lunenfeld-Tanenbaum Research Institute (LTRI), Sinai Health System, Toronto, Ontario, Canada
32. Department of Computer Science, University of Toronto, Toronto, Ontario, Canada
33. Cecil H. and Ida Green Center for Reproductive Biology Sciences, University of Texas Southwestern Medical Center, Dallas, TX, USA
34. Department of Obstetrics and Gynecology, University of Texas Southwestern Medical Center, Dallas, TX, USA
35. Department of Bioinformatics, University of Texas Southwestern Medical Center, Dallas, TX, US
36. Department of Internal Medicine, Division of Cardiology, University of Texas Southwestern Medical Center, Dallas, TX, USA
37. Department of Molecular Biology, University of Texas Southwestern Medical Center, Dallas, TX, USA
38. Eugene McDermott Center for Human Growth and Development, University of Texas Southwestern Medical Center, Dallas, TX, USA
39. Department of Biochemistry, University of Texas Southwestern Medical Center, TX, USA
40. Quantitative Biomedical Research Center, Peter O'Donnell Jr. School of Public Health, University of Texas Southwestern Medical Center, Dallas, TX, U.S.A
41. Department of Pediatrics, Division of Hematology/Oncology, University of Texas Southwestern Medical Center, Dallas, TX, U.S.A

42. Children's Medical Center Research Institute, University of Texas Southwestern Medical Center, TX, USA
43. Department of Genetics, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA
44. Neuroscience Center, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA
45. Department of Pathology, Harvard Medical School, Boston, MA, USA
46. Broad Institute of MIT and Harvard, Boston, MA, USA
47. Division of Hematology/Oncology, Boston Children's Hospital, Boston, MA, USA
48. Department of Pediatrics, Harvard Medical School, Boston, MA, USA
49. Montreal Heart Institute, Montreal, Quebec, H1T 1C8, Canada
50. Département de Médecine, Université de Montréal, Montréal, Quebec, H3T 1J4, Canada
51. Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA 02115
52. Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA
53. Division of Hematology/Oncology, Boston Children's Hospital, Boston, MA, USA
54. Center for Genomic Medicine and Department of Pathology, Massachusetts General Hospital, Boston, MA, USA
55. Department of Biomedical Engineering, Duke University, Durham, NC, USA
56. Center for Advanced Genomic Technologies, Duke University, Durham, NC
57. Department of Pediatrics, Duke University, Durham, NC, USA
58. Department of Biostatistics and Bioinformatics, Duke University Medical Center, Durham, NC
59. Department of Cell Biology, Duke University Medical Center, Durham, NC, USA
60. Department of Pharmacology and Cancer Biology, Duke University Medical Center, Durham, NC, USA
61. Gene Regulation Observatory, The Broad Institute of MIT and Harvard, Cambridge, MA, USA
62. Department of Stem Cell and Regenerative Biology, Harvard University, Cambridge, MA, USA
63. Department of Cancer Biology, Dana-Farber Cancer Institute, Boston, MA, USA
64. Department of Pathology, Stanford University, Stanford, CA, USA
65. Parker Institute for Cancer Immunotherapy, San Francisco, CA, United States
66. Gladstone-UCSF Institute of Genomic Immunology, San Francisco, CA, 94158, USA.
67. Berlin Institute of Health at Charité – Universitätsmedizin Berlin, Berlin, Germany
68. Max-Delbrück-Center for Molecular Medicine in the Helmholtz Association (MDC), Berlin Institute for Medical Systems Biology (BIMSB), Berlin, Germany
69. Institute for Human Genetics, Department of Medicine, Division of Rheumatology, University of California, San Francisco, CA, United States
70. Parker Institute for Cancer Immunotherapy, San Francisco, CA, United States
71. Chan Zuckerberg Biohub, San Francisco, CA, United States
72. Sean N Parker Center for Allergy and Asthma Research, Stanford University, Stanford, CA, USA
73. Gladstone Institute of Neurological Disease, San Francisco, CA, USA
74. Department of Neurology, University of California San Francisco, San Francisco, CA, USA
75. Department of Surgery, University of California San Francisco, San Francisco, CA, USA
76. Diabetes Center, University of California San Francisco, San Francisco, CA, USA
77. Department of Statistical Science, Duke University, Durham, NC, USA
78. Department of Mathematics, Duke University, Durham, NC, USA
79. Department of Computer Science, Duke University, Durham, NC, USA
80. Department of Biology, Duke University, Durham NC, USA
81. Computational and Systems Biology Program, Memorial Sloan Kettering Cancer Center, New York NY, USA
82. Department of Medical and Population Genetics, Broad Institute, Cambridge, MA, USA
83. Department of Epidemiology, Harvard T.H.Chan School of Public Health, Boston, MA, USA
84. Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA
85. Department of Genetics, Harvard Medical School, Boston, MA, USA
86. Khoury College of Computer Sciences, Northeastern University, Boston, MA 02115, USA
87. Department of Psychiatry, University of California San Diego, La Jolla, CA 92093, USA
88. Department of Biomedical Informatics and Medical Education, University of Washington, Seattle, WA 98195, USA
89. Institute for Genomic Health, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA

90. Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA
91. Department of Biostatistics and Medical Informatics, University of Wisconsin, Madison, WI, USA
92. Department of Genetics, University of Wisconsin, Madison, WI, USA
93. Department of Statistics, University of Wisconsin, Madison, WI, USA
94. Department of Computer Sciences, University of Wisconsin, Madison, WI, USA
95. Department of Pediatrics, University of Wisconsin, Madison, WI, USA
96. Prevention Genetics Inc., Marshfield, WI, USA
97. Marshfield Clinic, Marshfield, WI, USA
98. Program in Bioinformatics and Integrative Biology, UMass Chan Medical School, Worcester, MA, USA
99. Department of Data Science, Dana-Farber Cancer Institute, Boston, MA
100. Department of Statistics, Harvard University, Cambridge, MA
101. Department of Data Science, Dana-Farber Cancer Institute, Boston, MA
102. Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD, USA
103. Malone Center for Engineering in Healthcare, Johns Hopkins University, Baltimore, MD, USA
104. Department of Computer Science, Johns Hopkins University, Baltimore, MD, USA
105. Department of Genetic Medicine, Johns Hopkins University, Baltimore, MD, USA
106. Department of Biomedical Data Science, Stanford University, Stanford, CA, USA
107. Departments of Immunology and Computational and Systems Biology, University of Pittsburgh, Pittsburgh, PA
108. Department of Epigenetics and Molecular Carcinogenesis, University of Texas MD Anderson Cancer Center, Houston, TX, USA
109. Department of Epigenetics and Molecular Carcinogenesis, University of Texas MD Anderson Cancer Center, Houston, TX, USA
110. Department of Bioinformatics and Computational Biology, The University of Texas MD Anderson Cancer Center, Houston, TX, USA
111. Departments of Immunology and Computational and Systems Biology, University of Pittsburgh, Pittsburgh, PA, USA
112. Department of Epigenetics and Molecular Carcinogenesis, University of Texas MD Anderson Cancer Center, Houston, TX, USA
113. Center for Immunotherapy and Precision Immuno-Oncology, Cleveland Clinic, Cleveland, OH, USA
114. Center for Immunotherapy and Precision Immuno-Oncology, Cleveland Clinic, Cleveland, OH, USA
115. Livestrong Cancer Institutes, Department of Oncology, and Department of Biomedical Engineering, The University of Texas at Austin, Austin, TX, USA
116. Interdisciplinary Life Sciences Graduate Programs (ILSGP), and Oden Institute for Computational Engineering and Sciences (ICES), The University of Texas at Austin, Austin, TX, USA
117. Department of Genetics, University of Pennsylvania, Philadelphia, PA., USA
118. Department of Electrical & Computer Engineering, University of Washington, Seattle, WA., USA
119. Department of Neuroscience, University of Pennsylvania, Philadelphia, PA., USA
120. Department of Medicine, University of Pennsylvania, Philadelphia, PA., USA
121. Department of Cancer Biology, University of Pennsylvania, Philadelphia, PA., USA
122. Developmental Biology Program, Sloan Kettering Institute, New York, NY, USA
123. Center for Stem Cell Biology, Sloan Kettering Institute for Cancer Research, New York, NY 10065, USA
124. Department of Biomedical Engineering and McKusick-Nathans Department of Genetic Medicine, Johns Hopkins University; Baltimore, MD 21218, USA
125. Department of Biochemistry, Albert Einstein College of Medicine, Bronx, NY 10461, USA
126. Weill Cornell Graduate School of Medical Sciences, Weill Cornell Medicine, 1300 York Avenue, New York, NY 10065, USA
127. Department of Pediatrics, University of California, San Diego, USA
128. Department of Medicine, University of California, San Diego, USA
129. Department of Cellular and Molecular Medicine, University of California, San Diego, CA, USA
130. Center for Epigenomics, University of California, San Diego

131. Bioinformatics and Systems Biology Program, University of California, San Diego, CA USA
132. Biomedical Sciences Program, University of California, San Diego, CA USA
133. Department of Human Genetics, University of California Los Angeles, Los Angeles, CA USA
134. Department of Biological Chemistry, David Geffen School of Medicine, University of California Los Angeles, Los Angeles, CA, USA
135. Department of Neurology, University of California Los Angeles, Los Angeles, CA, USA
136. Department of Pathology and Laboratory Medicine, University of California Los Angeles, Los Angeles, CA, USA
137. Department of Computational Medicine, University of California Los Angeles, Los Angeles, CA, USA
138. Department of Biostatistics, University of California Los Angeles, Los Angeles, CA, USA
139. Department of Genetics, Washington University, St. Louis, MO., USA
140. McDonnell Genome Institute, Washington University School of Medicine, Saint Louis, MO, USA
141. Sloan Kettering Institute, Memorial Sloan Kettering Cancer Center, New York, NY, USA
142. Center for Genetic Epidemiology, Department of Population and Public Health Sciences, Keck School of Medicine, University of Southern California, CA, USA
143. Department of Quantitative and Computational Biology, University of Southern California, CA, USA
144. Norris Comprehensive Cancer Center, Keck School of Medicine, University of Southern California, CA, USA
145. Institute of Computational Biology, Helmholtz Zentrum Munich, Neuherberg, Germany
146. Department of Computer Science, School of Computation, Information and Technology, Technical University Munich, Munich, Germany
147. Munich Heart Alliance, DZHK (German Center for Cardiovascular Research), Munich, Germany
148. Lawrence Berkeley National Laboratory, Berkeley, CA 94720
149. New York Genome Center, New York, NY USA
150. Department of Biology, New York University, New York, NY USA
151. DOE Joint Genome Institute, Berkeley, CA USA
152. Gladstone Institutes, San Francisco, CA USA
153. University of California, San Francisco, CA USA
154. Chan Zuckerberg Biohub - San Francisco, San Francisco, CA USA
155. St. Anna Children's Cancer Research Institute (CCRI), Vienna, Austria
156. CeMM Research Center for Molecular Medicine of the Austrian Academy of Sciences, Vienna, Austria
157. Livestrong Cancer Institutes, Department of Oncology, and Department of Biomedical Engineering, The University of Texas at Austin, Austin, TX, USA
158. Interdisciplinary Life Sciences Graduate Programs (ILSGP), and Oden Institute for Computational Engineering and Sciences (ICES), The University of Texas at Austin, Austin, TX, USA
159. Division of Genomic Medicine, National Human Genome Research Institute (NHGRI), National Institutes of Health (NIH), Bethesda, MD, USA
160. Division of Genome Sciences, National Human Genome Research Institute (NHGRI), National Institutes of Health (NIH), Bethesda, MD, USA
161. Department of Applied Physics, Stanford University, Stanford, CA, USA 94305
162. Cancer Biology Program, Stanford University School of Medicine, Stanford, CA, USA 94305
163. Immunology Graduate Program and Department of Pathology, Stanford University School of Medicine, Stanford, CA, USA 94305
164. Department of Plant Biology, Carnegie Institution for Science, Stanford, CA 94305, USA
165. Computational Biology Department, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, USA
166. Program in Computational Biology & Bioinformatics, Yale University, New Haven, CT, USA
167. Department of Molecular Biophysics & Biochemistry, Yale University, New Haven, CT, USA
168. Department of Computer Science, Yale University, New Haven, CT, USA
169. Department of Statistics & Data Science, Yale University, New Haven, CT, USA
170. Department of Biomedical Informatics & Data Science, Yale University, New Haven, CT, USA
171. Howard Hughes Medical Institute, Seattle, WA, USA

172. Systems Immunology, Benaroya Research Institute, Seattle, WA, USA
173. Allen Discovery Center for Cell Lineage Tracing, Seattle, WA, USA
174. mRNA Center of Excellence, Sanofi Pasteur Inc., Waltham, MA, USA
175. Department of Biochemistry and Molecular Genetics, Feinberg School of Medicine Northwestern University, Chicago, IL, USA
176. Robert H. Lurie Comprehensive Cancer Center of Northwestern University, Chicago, IL, USA
177. Center for Genetic Medicine, Department of Pharmacology, Northwestern University, Chicago, IL, USA
178. Center for Genetic Medicine and Department of Cell and Developmental Biology, Feinberg School of Medicine, Northwestern University, Chicago, IL, USA
179. Department of Preventive Medicine, Feinberg School of Medicine, Northwestern University
180. Department of Pediatrics, Washington University, St. Louis, MO., USA
181. Department of Medicine, Washington University, St. Louis, MO., USA
182. Department of Developmental Biology, Washington University, St. Louis, MO., USA
183. Department of Pathology and Immunology, Washington University, St. Louis, MO., USA
184. TERRA Teaching and Research Centre, University of Liège, Gembloux, Belgium
185. Computational Biology and Bioinformatics, Université Libre de Bruxelles, Brussels, Belgium
186. UC Santa Cruz Genomics Institute, University of California Santa Cruz, Santa Cruz, CA, USA
187. Gladstone Institute of Data Science and Biotechnology, Gladstone Institutes, San Francisco, CA, USA
188. Department of Developmental and Cell Biology, UC Irvine, Irvine, CA., USA
189. Center for Complex Biological Systems, UC Irvine, Irvine, CA., USA
190. Division of Biology and Biological Engineering, California Institute of Technology, Pasadena, CA USA
191. Richard N. Merkin Institute for Translational Research, California Institute of Technology, Pasadena, CA., USA
192. Division of Chemistry and Chemical Engineering, California Institute of Technology, Pasadena, CA., USA
193. Department of Neurobiology and Behavior, UC Irvine, Irvine, CA., USA
194. Department of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA., USA
195. Department of Pharmaceutical Sciences, UC Irvine, Irvine, CA., USA
196. Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI, USA
197. Department of Human Genetics, University of Michigan, Ann Arbor, MI, USA
198. Department of Biostatistics, School of Public Health, University of Michigan, Ann Arbor, MI, USA
199. Computational and Systems Biology Program, Memorial Sloan Kettering Cancer Center, New York, NY, USA
200. Howard Hughes Medical Institute and Immunology Program at Sloan Kettering Institute, New York, NY, USA
201. Ludwig Center for Cancer Immunotherapy, Memorial Sloan Kettering Cancer Center, New York, NY, USA
202. Division of Rheumatology, Department of Medicine, Hospital for Special Surgery, New York, NY, USA
203. Weill Cornell Medical College and Graduate School, New York, NY, USA
204. Department of Biostatistics, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA
205. Department of Pharmacology, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA
206. Department of Environmental Health Sciences, Mailman School of Public Health Columbia University, New York, NY, USA
207. Masters of Physician Assistant Studies Program, Colorado Mesa University, Grand Junction CO, USA
208. Department of Immunology, University of Washington, Seattle, WA, USA
209. Cincinnati Children's Hospital, Cincinnati OH, USA
210. Department of Genetics and Development, Institute for Cancer Genetics, Herbert Irving Comprehensive Cancer Center, Columbia University Irving Medical Center, New York, NY, USA
211. Department of Neuroscience, University of Texas Southwestern Medical Center, Dallas, TX, USA

212. Department of Psychiatry, University of Texas Southwestern Medical Center, Dallas, TX, USA
213. Center for the Genetics of Host Defense, University of Texas Southwestern Medical Center, Dallas, TX, USA
214. Peter O'Donnell Jr. Brain Institute, University of Texas Southwestern Medical Center, Dallas, TX, USA
215. Altos Labs Inc., 1300 Island Drive, Redwood City, CA, USA
216. Department of Neurology, Washington University, St. Louis, MO., USA
217. Department of Pathology, Massachusetts General Hospital, Boston, MA, USA
218. PhD Program in Biological and Biomedical Sciences, Harvard University, Boston, MA, USA
219. Departments of Psychiatry and Genetics, Division of Molecular Psychiatry, Department of Genetics, Wu Tsai Institute, Yale University School of Medicine, New Haven, CT, USA
220. MIT Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA, USA
221. Paul G. Allen School of Computer Science and Engineering, University of Washington, Seattle, WA, USA
222. Canadian Institute for Advanced Research, Toronto, ON, Canada
223. Department of Genetics, Yale School of Medicine, New Haven, CT, USA

## Author Contributions

Jesse M. Engreitz, Heather A. Lawson, and Harinder Singh co-led the Writing Group. Jesse M. Engreitz, Heather A. Lawson, Harinder Singh, Lea Starita, Gary C. Hon, Hannah Carter, Nidhi Sahni, Timothy E. Reddy, Xihong Lin, Yun Li, Nikhil Munshi, Maria Chahrour, Benjamin Hitz, and Ali Mortazavi wrote initial text based on input from PIs, the Writing Group, and Working Group and Focus Group Co-Chairs. Jesse M. Engreitz, Alan Boyle, and Jayoung Ryu developed figures. All authors contributed to developing the vision and goals of the IGVF Consortium, outlining the project, and editing the manuscript. The role of the NHGRI Program Management in the preparation of this paper was limited to coordination and scientific management of the IGVF Consortium.

## Acknowledgements

This work was supported by the NIH NHGRI IGVF Program (UM1HG011966, UM1HG011969, UM1HG011972, UM1HG011989, UM1HG011996, UM1HG012003, UM1HG012010, UM1HG012053, UM1HG011986, UM1HG012076, UM1HG012077, U01HG011952, U01HG011967, U01HG012009, U01HG012022, U01HG012039, U01HG012064, U01HG012069, U01HG012041, U01HG012047, U01HG012051, U01HG012059, U01HG012079, U01HG012103, U24HG012012, U24HG012070), NIH NCI (R01CA197774), and the Novo Nordisk Foundation (NNF21SA0072102). Figures were illustrated by SciStories LLC. We thank members of the IGVF External Consultants Panel (Guillaume Bourque, Prashant Mali, Judy Cho, Barbara Engelhardt, and Olga Troyanskaya) for critical feedback on the manuscript.