

**Reads Meet Rotamers:**

**Structural Biology in the Age of Next Generation Sequencing**

Anurag Sethi\*, Declan Clarke\*, Jieming Chen, Sushant Kumar, Timur Galeev, Lynne Regan,  
and Mark Gerstein<sup>†</sup>

\*Equal Contribution by authors

<sup>†</sup> Corresponding Author: [mark@gersteinlab.org](mailto:mark@gersteinlab.org)

## Abstract:

Structure has historically been interrelated with sequence, usually in the framework of comparing sequences across species sharing a common structural fold. However, the nature of information within the sequence and structure databases is evolving, changing the type of comparisons possible. In particular, we now have a vast amount of personal genome sequences from human populations and a larger fraction of new structures contain interacting proteins within large complexes. Consequently, we have to recast our conception of sequence conservation (considering more selection within the human population) and its relation to structure (now focusing more on interacting surfaces rather than folds). We cover this changing mindset here. [\[\[ANS2MG: Should we add a sentence on networks here?\]\]](#)

## Highlights:

- Increasing amounts of sequencing data adds new dimension to study structural biology.
- Next generation sequencing facilitates analysis of genetic variations.
- Understanding evolutionary constraints acting on proteins remain elusive.
- Essential to integrate sequences, structures, and interaction networks information to rationalize the phenotypic impact of these variations.

## Introduction:

The amount of genomic information is growing at an astonishing pace due to rapid improvements in next-generation sequencing (NGS) technology (Figure 1A) [1]. The nature of biological information stored within biological databases is undergoing a transformation (Figure 1). Before the completion of the human genome project in 2003, we had a large amount of genomic sequence information from different species and structural data in the databases. Due to the technological advances in next-generation sequencing, the amount of human sequence information has grown at an unprecedented pace. Meanwhile, even though the number of protein structures in the PDB database [\cite{10592235}](#) has also increased, the pace of identifying new folds has slowed down indicating that few new folds remain undiscovered. However, a large number of novel domain-domain interactions are detected in the newly deposited structures indicating that the complexity of the structures in the PDB database continues to grow (Figure 1). This trend illustrates an increasing emphasis among structural biologists to treat biomolecules not as individual folds but rather as complex molecular machines that interact and regulate each other as they function within the cellular

Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto

Anurag Sethi 10/24/2015 10:03 PM

Deleted: . . . [1]

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Font color: Auto

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Font color: Auto, Condensed by 0.05 pt

Anurag Sethi 10/24/2015 10:03 PM

Moved (insertion) [1]

Anurag Sethi 10/24/2015 10:03 PM

Deleted: Essential goals of these efforts include the realization of personalized medicine by identifying pathological disease-associated variants

environment. Together, these trends suggest that the stage is set to integrate sequence and structural information to rationalize the effect of variants on protein function.

The identification and characterization of pathological disease-associated variants is an essential goal of genomic sequencing efforts [2,3]. A large number of medically-relevant mutations occur within proteins, some of which are available through databases such as the Online Database of Mendelian Inheritance in Man (OMIM) [4], the Human Gene Mutation Database (HGMD) [5], Humsavar [6], and ClinVar [7]. It is essential to utilize structural information for rationalizing the evolutionary pressure acting on these proteins as well as for developing drugs to combat the effects of disease-causing variants. However, it remains challenging to annotate the physical effects of these mutations on proteins due to the heterogeneous nature of functional constraints acting on a protein family. A protein-coding variant may cause local perturbations, or global changes in structure, or it may have a substantial impact on the protein-protein interaction (PPI) network, and each type of change adds a different layer of functional constraints on the protein. Such analyses are further complicated by the fact that we currently have incomplete knowledge of these constraints, and also by the fact that specific combinations of individually benign variants may cause disease.

While structural data provides an invaluable guide for rationalizing disease-associated variants, we also expect the growing genomic information to be a valuable resource for structural biologists. In particular, as the amount of genomic data continues to grow, we envision a future in which biologists will utilize genetic variation within human population(s) to help interpret their structural data [8]. Population genetic analysis within human proteins has already been used to identify novel species-specific functional constraints within a protein family [9]. In addition, a number of fundamental insights about biological pathways can be garnered by analyzing newly discovered loci associated with a disease [10].

In this review article, we initially explain how genomic information is used to identify pathological disease associated variants as well as variants that are harmful to protein function even within healthy individuals. We later describe how structural information is utilized to understand the harmful effects of different variants. Finally, we discuss how it is necessary to integrate sequence and structural data with a holistic system or network perspective before predicting phenotypic effects of the variants.

▪  
**Classical Sequence Comparison:**

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** to rationalize

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** assortment

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** , incomplete knowledge of these constraints, and how individual variants can be benign but disease-causing in specific combinations

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** functional constraints on the protein. Moreover

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** -

Anurag Sethi 10/24/2015 10:03 PM

**Moved up [1]:** The nature of biological information stored within biological databases is undergoing a transformation (Figure 1).

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** Before the completion of the human genome project in 2003, we had a large amount of genomic sequence information from different species as well as structural data. Since the technological advances in next-generation sequencing, the amount of human sequence information has grown at a rapid pace. Meanwhile, while structural biologists continue to deposit new structures in the PDB database, the pace of identifying new folds has slowed down indicating that few new folds remain undiscovered. However, the complexity of the structure in the PDB database continues to grow indicating that there is an increasing emphasis among structural biologists to treat biomolecules not as individual folds but rather as complex molecular machines that interact and regulate each other as they function within the cellular environment. Together, these trends suggest that the stage is set to utilize structural information to rationalize the effect of variants on protein function. -

Typically, structural biologists identify functionally constrained regions within a protein family by comparing homologous sequences from different species (Figure 2a) [11,12]. They focus on changes that take place over longer evolutionary timescales by comparing the **reference (or dominant)** sequence within each species rather than focusing on intra-species changes. Nucleotides that do not change across different species are conserved over millions of years and are hence considered to be functionally important. Due to redundancy within the genetic code, some of the changes in the coding regions are silent as they occur without a corresponding change in the protein sequence (synonymous changes). **With rare exceptions**, all synonymous changes and a majority of the nonsynonymous changes are expected to be **neutral or harmful (deleterious) to the protein function**. A small fraction of the nonsynonymous changes can, however, **be** beneficial to the fitness of the species. The **ratio** of nonsynonymous to synonymous variants (dN/dS) is commonly utilized to characterize the selection pressure on the coding regions of the genome (Figure 2) [13]. If the dN/dS ratio for a coding region is less than 1, it indicates that a few of these mutations are harmful or deleterious and that **the protein is** under negative selection. On the other hand, a dN/dS ratio exceeding unity indicates that evolution is promoting a change in the protein sequence and that this protein is under positive selection [9]. Proteins undergoing positive selection may improve the fitness of an organism **to** different environments.

### Introduction to Population Sequencing:

The vast amounts of genomic and exome sequences available **are** providing unique opportunities to characterize genetic variation within the human population. The exome comprises the coding sequences of all protein-coding genes and constitutes approximately 1% of the total genomic sequence [14]. Due to the reduced cost of exome sequencing and better-characterized clinical relevance of variation within the coding regions of the genome, it is more widely used for genetic diagnosis. Variants within an individual's genome are either acquired at birth (germline mutations) or during the person's lifetime (somatic mutations) as a consequence of errors during cell division. While germline mutations are typically present in every cell of the person, somatic mutations only affect certain cells and are typically not passed on to the next generation. There are approximately 74 *de novo* (new) variants that occur during each generation [15]. As only germline mutations are passed on to the next generation, somatic mutations are not under conventional evolutionary selection.

The human genome exhibits extensive variation [16-19]. On average, any individual genome contains 20,000-25,000 coding variants (Table 1), of which 9,000-11,000 are nonsynonymous.

Anurag Sethi 10/24/2015 10:03 PM

Deleted: While

Anurag Sethi 10/24/2015 10:03 PM

Deleted: do exist

Anurag Sethi 10/24/2015 10:03 PM

Deleted: .

Anurag Sethi 10/24/2015 10:03 PM

Deleted: either be harmful (deleterious) or

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Not Highlight

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Not Highlight

Anurag Sethi 10/24/2015 10:03 PM

Deleted: these changes are

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Not Highlight

Anurag Sethi 10/24/2015 10:03 PM

Deleted: in

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Font:Arial, 11 pt

Anurag Sethi 10/24/2015 10:03 PM

Deleted: is

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Font:Arial, 11 pt

The frequency with which a particular variant or allele occurs within a population is used to characterize the evolutionary pressure acting on it as common variants (minor allele frequency greater > 5%) are expected to be benign. However, rare variants (minor allele frequency < 0.5%) are rare either because they are harmful (deleterious) to a protein's function or because the variant has been introduced recently into the population. The ratio of common to rare variants is often used as a proxy to characterize the evolutionary pressure acting on a locus.

Although most of the variants within any particular individual are common, most coding variants manifest as distinct single nucleotide variants (SNVs), each of which occurs very rarely within the human population. About 25-50% of the rare non-synonymous variants within healthy individuals are estimated to be deleterious, suggesting that the human proteome is highly robust to a large number of non-specific perturbations and because most rare deleterious variants are heterozygous implying that the cell also contains a functional copy of the gene [18,19].

Despite the fact that new genomic data is still being produced, about 200,000-500,000 previously unobserved SNVs are still discovered after each personal genome is sequenced, suggesting that there is no saturation observed on human polymorphism data yet [18,19].

Indeed, the number of rare variants continues to grow even after the 1000 Genomes Consortium and Exome Aggregation Consortium data (60,706 individuals) [20] data has become available. As deleterious mutations tend to occur at very low frequencies, we need to continue sequencing a large number of individuals to characterize and catalog these variants and their frequencies within the human population.

As such, we can turn to intra-human comparisons to uncover more human- or domain-specific features (Figure 2). There is, however, an important distinction between interpreting inter- and intra-species conservation due to the huge disparities in the associated evolutionary timescales (Figure 2a-c). While performing such an analysis, one can also align homologous coding regions not only between individuals (Figure 2b), but also within a single human genome (i.e., paralogs), such as proteins originating from the same structural domain family (Figure 2c). In particular, this can be used to elucidate domain-specific features.

Similar to the dN/dS ratio in cross-species comparisons, selective pressure on coding regions can be quantified using fraction of synonymous to nonsynonymous polymorphisms (pN/pS) at any site (Figure 2e). In addition, evolutionary pressure can also be quantified during intra-species comparison using the ratio of rare to common variants at each site as rare variants are under negative selection (Figure 2e). A statistically significant depletion of common variants as

Anurag Sethi 10/24/2015 10:03 PM

Deleted: As deleterious variants are under negative selection, the

Anurag Sethi 10/24/2015 10:03 PM

Deleted: in

Anurag Sethi 10/24/2015 10:03 PM

Deleted: particular

Anurag Sethi 10/24/2015 10:03 PM

Deleted: can be

Anurag Sethi 10/24/2015 10:03 PM

Deleted: it

Anurag Sethi 10/24/2015 10:03 PM

Deleted: (defined as having a minor allele frequency greater than 5%),

Anurag Sethi 10/24/2015 10:03 PM

Deleted: occur

Anurag Sethi 10/24/2015 10:03 PM

Deleted: (defined as having a minor allele frequency less than 0.5%).

Anurag Sethi 10/24/2015 10:03 PM

Deleted:

Anurag Sethi 10/24/2015 10:03 PM

Deleted: not yet a

Anurag Sethi 10/24/2015 10:03 PM

Deleted: in data

Anurag Sethi 10/24/2015 10:03 PM

Deleted:

Anurag Sethi 10/24/2015 10:03 PM

Deleted:

Anurag Sethi 10/24/2015 10:03 PM

Deleted: .

Anurag Sethi 10/24/2015 10:03 PM

Deleted: .

Anurag Sethi 10/24/2015 10:03 PM

Deleted: .

Anurag Sethi 10/24/2015 10:03 PM

Deleted: (Figure 2b).

Anurag Sethi 10/24/2015 10:03 PM

Deleted: 2

Anurag Sethi 10/24/2015 10:03 PM

Deleted: .

compared to rare variants implies that the site is under higher selective pressure. Furthermore, genomic variants that are increasing in frequency within a human population (positive selection) may help identify a novel gain-of-function event (such as a new protein-protein interaction).

Some of these domain-specific events may be beneficial to the species. Comparative genetics/genomics studies have already uncovered a growing list of genes that might have experienced positive selection during the evolution of human and/or primates [9]. These genes offer valuable insights into understanding the biological processes specific to humans, as well as the evolutionary forces that gave rise to them. It is also important to note that some variants occur in a correlated fashion within the population and these variants are said to be under linkage disequilibrium (LD). Note also that LD is statistically easier to observe for common variants than for rare ones.

### **Deleterious Effects of Variations on Protein Function:**

The patterns of conservation displayed by proteins are the product of a vast array of constraints active throughout its evolutionary history. In this regard, to understand the physical effects that cause a variant to be harmful, we need to consider the multitude of underlying constraints acting on the protein family. Such constraints are often intrinsic to the structure itself: they may include the need to maintain the integrity of functional hinge regions or interior packing geometry or the ability to regulate a protein through post-translational modifications at specific sites. They may also entail that residues at an interaction interface remain topologically compatible with those in the corresponding interface of an interaction partner. We can utilize the structural information in the PDB database to assess the effect of mutations on a protein's stability as nonsynonymous changes that occur within the core of the protein or variants that disrupt the secondary structure of the protein could reduce its stability. Several computational tools based on sequence conservation (inter-species or intra-species) and/or several structural features (the physicochemical characteristics of the amino acid change, solvent accessibility, secondary structure, active site annotations, and protein-protein interfaces) were developed to predict the deleterious effect of sequence variations on a protein's function [22-25]. Disease-associated mutations are highly enriched for residues in the interior of proteins (22% of all mutations in HGMD and OMIM), and active sites of proteins [16-19].

In terms of applying such a catalogue of rules as a means of understanding human disease-associated variants, the fibroblast growth factor receptor provides a case-in-point, several variants in which have been linked to craniofacial defects (Figure 3). The evolutionary constraints listed here provide sensible rationales for how many of these disease-associated

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Not Highlight

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Not Highlight

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Normal1

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** There is one additional confounding factor to consider while identifying disease-associated variants. Genes associated with a disease are identified by detecting deleterious variants that are affecting genes within diseased individuals more often than in healthy populations. This might be misleading, however, because the variants associated with this gene might be correlated with other unanalyzed variants in the genome. Hence, all variants (including the variants within a gene) statistically associated with a disease might not be causative and additional analysis may be required to identify the real disease-causing mutations. We need to annotate the effect of individual variants, however, before we can predict the outcome of a large number of variants.

Each protein has several evolutionary evolutionary constraints imposed upon it based on its biological function. The effect of a deleterious variant can only be understood when all these functional constraints acting on a protein are known and can be considered. The fibroblast growth factor receptor provides a case-in-point (Figure 3). This protein has been shown to host well-documented disease-causing variants that manifest in craniofacial defects in humans. However, several of the disease variants have no clear mechanism of pathogenicity in that they do not fall in any of the protein regions known to be sensitive to amino acid changes. Certainly, a sequence ... [2]

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Font:Not Bold

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** Each protein has several evolutionary constraints imposed upon it based on its biological function. The effect of a deleterious variant can only be understood when all these functiona ... [3]

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Not Highlight

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Font:Arial, 11 pt

Anurag Sethi 10/24/2015 10:03 PM

Deleted: -

... [4]

variants may impart deleterious effects. Importantly, these constraints may act in synergistic ways rather than through isolated mechanisms [cite{23364837, 22153503}]. However, the mechanisms for several other disease-associated variants fail to map to this catalogue, thereby underscoring the need to more comprehensively document sources of constraint. This more comprehensive documentation needs to transcend the native structure itself by including folding pathways, allosteric regulation, and the functional roles of disordered regions or conformational transitions. Such mutations that affect the thermodynamic stability of different allosteric states of a protein [26] are typically ignored while predicting the deleteriousness of a putative variant.

### Networks as a Framework for Understanding Deleterious Variants:

While structural and sequence information are invaluable in providing a rationale for the deleterious effects of certain disease-causing and rare variations, it is often difficult to interpret the phenotypic effects of an individual variant without considering the broader cellular context. As proteins are extensively involved in protein-DNA interactions (gene regulatory network), protein-RNA interactions (post-transcriptional regulation), and protein-protein interactions (PPI) within the cellular milieu, variants that disrupt these interactions could potentially affect the viability of the cell. We refer the reader to comprehensive essays on the phenotypic effect of noncoding variation [27,28], and focus instead on deleterious effects of variants on the protein-protein interaction (PPI) network here.

Various experimental and computational approaches have been applied to characterize the PPI network in several model organisms and human beings [29,30] and these networks have been invaluable in interpreting the role of evolutionary constraints on a protein family [cite{}]. In the PPI network, a node represents a protein, while an edge represents an interaction between the two proteins connected by the edge. Proteins that are highly interconnected in PPI networks (hubs) are under strong negative selection while proteins under positive selection in humans tend to occur at the periphery of the network [31]. Proteins that are more central in an integrated "multinet" formed by integrating biological networks from different context (PPI, metabolic, post-translational modification, gene regulatory network, etc.) are under negative selection within human populations [32]. In agreement with this, perturbations to hub proteins are more likely to be associated with diseases than non-hub proteins [33]. The PPI networks are organized in a modular fashion as proteins associated with the same function are more likely to interact with one another [34] and proteins associated with similar diseases tend to occur within the same module [33]. The system properties of the network have also been useful in interpreting how the human proteome is robust even in the presence of a large number of deleterious variants within

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** they are present in.

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** here

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** that

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** may have

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** .

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** constraints

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** at the periphery of the network are

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** pooling

healthy individuals. Most deleterious variants observed in healthy individuals occur in peripheral regions of the interactome. Such limited effects may result as a consequence of compensatory mutations or functional redundancy [35]. On the other hand, cancer-associated somatic deleterious variations occur in the internal regions of the interactome and tend to have larger structural consequences on the PPI network.

The interactome provides a convenient platform to measure the impact of a deleterious variant on the cell. As shown in Figure 4, a deleterious variant can either remove a protein (such a node effect would naturally also result in the removal of all the associated edges) from the PPI network by making a protein nonfunctional or it could lead to the loss of just one or more of its interactions (edgetic effects). Mutations at a PPI interface can have drastic effects on the biomolecular binding constant and several sequence and structure-based methods have been proposed to identify these interaction hotspots [36,37]. Even though we have incomplete information on the structures of protein complexes (Figure 1), it has been predicted that about 12% of all the HGMD and OMIM mutations occur at a PPI interaction [38] while approximately 28% of experimentally-tested HGMD missense mutations affect one or more interactions, thus underscoring the importance of these interactions for annotating rare variants and disease-associated mutations [39].

In an effort to bridge the information gained from individual structures with network properties of the interactome, Kim, et al., [40] combined the experimentally determined interactome with structural information from the iPFam database to form the structural interaction network (SIN) and were able to obtain a higher-resolution understanding of the selection constraints on the hubs. Using structural information, the hubs were classified into different groups based on the number of distinct interfaces utilized for biomolecular complex formation and they showed that the number of distinct interfaces is a better proxy for evolutionary pressure acting on the hub rather than the number of edges in the PPI network. Consistent with this interpretation, hub proteins in the PPI network contain a higher fraction of disease-causing mutations on their solvent exposed surface, as compared to non-hub proteins suggesting that a larger fraction of a hub's disease-associated mutations could affect its interactions [40].

Hub proteins interact with a large number of partners and tend to be more flexible and conformationally heterogenous than non-hub proteins [41]. Furthermore, the number of distinct interfaces in hub proteins is correlated with degrees of conformational heterogeneity [41]. To the

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** , as a deleterious variant would have a larger effect on the structure of the PPI network if it occurs on a hub. As shown in Fig. 4, A

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** While the discovery of structural folds has saturated, the discovery of new domain-domain interactions continues to grow (Figure 1).

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** ,

Anurag Sethi 10/24/2015 10:03 PM

**Formatted:** Not Highlight

Anurag Sethi 10/24/2015 10:03 PM

**Formatted:** Not Highlight



extent that variants may enable or disable certain conformational states from being visited, such mutations could potentially affect protein complex formation and signaling pathways, and this has not yet been examined very closely. As deleterious mutations that affect hubs in networks tend to have a larger effect on the structures, they would also cause large changes in the PPI network. Proteins can utilize different interfaces for different (sets of) interactions, so multiple mutations on the same protein can be associated with drastically different diseases depending on the afflicted interface. Such mutations would have different edgetic effects on the protein's interaction network - by breaking or weakening one of its interactions while the rest of its interactions remain intact - and a large proportion of HGMD and OMIM mutations are predicted to have edgetic effects on the PPI network [39,42].

It should also be noted that the hubs in PPI networks also tend to contain higher degrees of disordered regions, and these regions typically become well-ordered upon ligand or protein binding [49,50]. Disease-associated mutations are enriched within disordered regions of the protein as they could affect post-translational modifications and/or protein-protein interaction sites \cite{24830552, 22080206}. The assessment of a mutation on the activity of an intrinsically disordered protein is even more challenging because it would be dependent upon the effects of these mutations upon the unfolded ensemble or the structure gained in the presence of its interaction partner. Due to their inherent flexibility, the unfolded ensembles of disordered proteins are especially difficult to characterize using either experimental or computational techniques [51,52], making variant annotation in the context of disordered proteins an uphill task. However, the phenotypic effect of mutations on the functional viability of a disordered protein is important because a number of proteins also change their interaction partners in a tissue-specific manner based upon the dominant isoform of the protein in that tissue \cite{22749400}. Recent evidence suggests that many mutations occurring on these alternatively-spliced disordered motifs may drive cancer [53]. Therefore it is important to understand the phenotypic effects of sequence variations in the disordered regions.

Ultimately, we want to develop an integrative framework to understand the effects of deleterious variants on the phenotype of the cell. However, a mutation typically displays tissue-specific phenotypic effects, hence an understanding of functional constraints on a protein should also incorporate tissue-specific information. While the gene regulatory network is being mapped out in a developmental time point and cell type-dependent fashion by several international consortia [43,44] the PPI network is largely treated in a static fashion. Recent work has tried to integrate

Anurag Sethi 10/24/2015 10:03 PM  
Moved (insertion) [2]

Anurag Sethi 10/24/2015 10:03 PM  
Moved (insertion) [3]

proteome and gene expression profiles with PPI networks to create tissue-specific networks [cite{24550720, 23399932}](#)[45]. However, these studies typically neglect the protein isoform even though the protein's interactions are dependent on its isoform [46,47]. A structural study on the effect of sequence variations on isoform-dependent PPI complexes has not been performed and would improve the prediction of phenotypic effects due to missense mutations. However, it is likely that the high costs in resources associated with studying isoform-specific assays in various cell types have impeded these types of studies. We anticipate that isoform-specific protein-protein interaction network annotation will become easier and more accessible in the near future, which will present new opportunities to better annotate such networks.

### Conclusions:

The exponential growth in genomic data has demonstrated that a surprisingly large amount of genomic variation is present within the human population, and this data has also helped identify a vast number of rare variants and disease-associated variants. Though the motivation of developing methods to annotate the effects of variants that cause human disease are clear, it remains challenging to do so as it requires bridging disparate sources of information together to understand the functional constraints on a protein family. It is essential to utilize structural information to rationalize the effect of variants. The network properties of the protein in addition to sequence and structural information regarding the nonsynonymous amino acid changes need to be considered within a single framework before predicting the phenotypic impact of an amino acid change.

### Acknowledgements:

#### Acknowledgments:

We acknowledge support from NIH and the AL Williams Professorship funds. [DC acknowledges the support of the NIH Predoctoral Program in Biophysics \(T32 GM008283-24\).](#)

Anurag Sethi 10/24/2015 10:03 PM

**Moved up [2]:** The assessment of a mutation on the activity of an intrinsically disordered protein is even more challenging because it would be dependent upon the effects of these mutations upon the unfolded ensemble or the structure gained in the presence of its interaction partner.

Anurag Sethi 10/24/2015 10:03 PM

**Moved up [3]:** Recent evidence suggests that many mutations occurring on these alternatively-spliced disordered motifs may drive cancer [53]. Therefore it is important to understand the phenotypic effects of sequence variations in the disordered regions. -

Anurag Sethi 10/24/2015 10:03 PM

**Deleted: Effect of Mutations on Disordered Regions:** .

... [5]

Anurag Sethi 10/24/2015 10:03 PM

**Deleted:** Due to their flexibility, the unfolded ensembles of disordered proteins are especially difficult to characterize using either experimental and computational techniques [51,52], making variant annotation in the context of disordered proteins an uphill task. However, the phenotypic effect of mutations on the functional viability of a disordered protein is important because a number of proteins also change their interaction partners in a tissue-specific manner based upon the dominant isoform of the protein in that tissue.

## Figure Captions:

**Figure 1:** The pace of novel fold discovery has begun to saturate, and while the volume of X-Ray crystal structures and structurally-resolved protein-protein interactions has continued to grow. However, the pace with which personal genomic sequencing databases are growing is considerably greater than the pace at which structure databases are growing.

**Figure 2: Evolutionary conservation in different contexts.** Evolutionary conservation can be inferred via sequence comparison in different contexts. (A) The examination of sequence conservation in orthologous sequences across multiple species looks at a longer evolutionary timescale. (B) The examination of the enrichment of rare variants (or depletion of common variants) in the same genomic element across multiple individuals within a single species or population looks at a shorter evolutionary timescale. Here, the red diamonds denote variants that are rare in a single human population (found in only one or a small number of individuals) and the blue diamonds denote variants that are commonly found in multiple individuals in the population. (C) The examination of sequence conservation in similar protein domain sequences within a single genome can reveal species- and domain-specific conservation that might be important to the structure or function of the domain family. (D) To illustrate (C), we use ankyrin protein domains as an example. We translate the DNA sequence of each ankyrin domain into its amino acid sequence. In order to relate the positions of the linear sequence of an ankyrin repeat domain to their structural locations, we then specifically paint each of the six ankyrin domains found in the structure of the human Notch 1 ankyrin domain (PDB ID: 1YYH) similar to the sequence profile in (E). (E) The top plot in this panel is the sequence profile of an ankyrin repeat domain with 30 amino acids, colored by position left to right, from green to yellow, corresponding to the coloring of the motifs of the human Notch 1 PDB structure in (D). In the sequence profile, the height of the amino acid letters connotes the degree of conservation of a particular residue at a specific location along the ankyrin repeat; the degree of conservation is computed using relative entropy in bits of information. To examine evolutionary conservation in more detail, the sequence profile can be further analysed with genomic variant profiles. For example, for each of the position along the ankyrin motif, the second plot shows the absolute numbers of variants binned into four categories: cyan bars show the number of variants that are common (c) and synonymous (s); blue bars for variants that are common and non-synonymous (ns); pink bars, rare (r) and synonymous; red bars, rare and non-synonymous. Subsequently, we can derive log ratios from these numbers to demonstrate an enrichment (or depletion) of categories of variants, in order to gain further

Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto

Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto

Anurag Sethi 10/24/2015 10:03 PM  
Deleted: of exome

Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto

Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto

Anurag Sethi 10/24/2015 10:03 PM  
Deleted: (A) Examining sequence

Anurag Sethi 10/24/2015 10:03 PM  
Deleted: of homologous

Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font:Bold, Italic

Anurag Sethi 10/24/2015 10:03 PM  
Deleted: Examining

Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font:Bold, Italic

Anurag Sethi 10/24/2015 10:03 PM  
Deleted: (C) Examining sequence conservation of similar

Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font:Bold, Italic

Anurag Sethi 10/24/2015 10:03 PM  
Deleted: , (D) reveals

Anurag Sethi 10/24/2015 10:03 PM  
Deleted: Genomic variant profile across the domain can be further analysed, e.g. by comparing the number of non-synonymous (ns) relative to synonymous (s) variants and comparing the number of rare (r) relative to common (c) variants. Here,

Anurag Sethi 10/24/2015 10:03 PM  
Deleted: found in the human genome

Anurag Sethi 10/24/2015 10:03 PM  
Deleted: The sequence profile of an ankyrin repeat motif is painted green to yellow, corresponding to the structure of

Anurag Sethi 10/24/2015 10:03 PM  
Deleted: of

Anurag Sethi 10/24/2015 10:03 PM  
Deleted: repeats

Anurag Sethi 10/24/2015 10:03 PM  
Deleted: ).

biological insights. Here, the third subplot displays a general enrichment of rare variants relative to common variants across the entire motif, suggesting a uniform evolutionary importance of the ankyrin domain in the human population. However, the fourth subplot exhibits a depletion of nonsynonymous variants relative to the synonymous variants at more conserved motif positions (in the sequence profile), hinting at only a subset of positions being of particular functional importance to the ankyrin domain family.

**Figure 3:** (A) The fibroblast growth factor receptor is shown in complex with FGF2 (PDB 1IIL), along with the loci of HGMD variants (orange spheres). (B) Various structural annotations (i.e., a "catalogue of constraints") are shown in sequence space. Hinge residues are taken from HingeMaster [54], buried residues are identified using NACCESS [55], protein-protein interaction residues are defined to be those within 4.5 Angstroms of the co-crystallized growth factor, and post-translational modification sites are taken from UniProt. HGMD loci shown as holo circles coincide with the catalogue of constraints, and may thus likely be rationalized in light of such constraints. However, a large number of HGMD loci (shown in filled orange circles) fail to overlap with these annotations, highlighting the need to consider alternative sources of constraint.

**Figure 4:** Various mechanisms of SNP-induced disruption in protein-protein interaction networks. A SNP that destabilizes a hub protein can ablate all associated interactions (A). SNPs disrupting different interfaces of the hub may interfere with interactions active in different tissues (B, C). Blue (hub protein), Yellow (nodes expressed in tissue1), Green (nodes expressed in tissue2), Turquoise (node expressed in tissue3). Mutation in cystathionine  $\beta$ -synthase (CBS) leads to metabolic disease called Homocystanuria. Among many HGMD SNPs impacting this protein, experimental evidence [cite{19888216}] suggest that I278T mutation leads to destabilization of CBS, which further disrupts of all three important interactions involving this protein and this is equivalent to removing a node from the PPI network. Mutation in EFHC1 gene, which has been implicated in epilepsy, presents a good example of edgetic effect [cite{25910212}]. This mutation perturbs interaction of EFHC1 with ZBED1 and TCF4. While the perturbed interaction between EFHC1 and ZBED1 interfere with cell proliferation [cite{17220279}], on the other hand disturbance in EFHC1 and TCF4 interaction influence the neuronal differentiation process [cite{17878293}].

- Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto
- Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Normal1
- Anurag Sethi 10/24/2015 10:03 PM  
Deleted: -
- Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto
- Anurag Sethi 10/24/2015 10:03 PM  
Deleted: (
- Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto
- Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto
- Anurag Sethi 10/24/2015 10:03 PM  
Deleted:
- Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto
- Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto
- Anurag Sethi 10/24/2015 10:03 PM  
Deleted: Only the sequence of the resolved growth factor receptor (chain E in PDB 1IIL) is shown.
- Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto
- Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto
- Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto
- Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto
- Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto

Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Normal1

## References and recommended reading

1. [Stephens ZD, Lee SY, Faghri F, Campbell RH, Zhai C, Efron MJ, Iyer R, Schatz MC, Sinha S, Robinson GE: \*\*Big Data: Astronomical or Genomical?\*\* \*PLoS Biol\* 2015, \*\*13\*\*:e1002195.](#)
2. Offit K: **Personalized medicine: new genomics, old lessons.** *Hum Genet* 2011, **130**:3-14.
3. Chin L, Andersen JN, Futreal PA: **Cancer genomics: from discovery science to personalized medicine.** *Nat Med* 2011, **17**:297-303.
4. Hamosh A, Scott AF, Amberger JS, Bocchini CA, McKusick VA: **Online Mendelian Inheritance in Man (OMIM), a knowledgebase of human genes and genetic disorders.** *Nucleic Acids Res* 2005, **33**:D514-517.
5. Stenson PD, Mort M, Ball EV, Howells K, Phillips AD, Thomas NS, Cooper DN: **The Human Gene Mutation Database: 2008 update.** *Genome Med* 2009, **1**:13.
6. **The Universal Protein Resource (UniProt) in 2010.** *Nucleic Acids Res* 2010, **38**:D142-148.
7. Landrum MJ, Lee JM, Riley GR, Jang W, Rubinstein WS, Church DM, Maglott DR: **ClinVar: public archive of relationships among sequence variation and human phenotype.** *Nucleic Acids Res* 2014, **42**:D980-985.
8. Sulkowska JI, Morcos F, Weigt M, Hwa T, Onuchic JN: **Genomics-aided structure prediction.** *Proc Natl Acad Sci U S A* 2012, **109**:10340-10345.
9. Voight BF, Kudaravalli S, Wen X, Pritchard JK: **A map of recent positive selection in the human genome.** *PLoS Biol* 2006, **4**:e72.
10. Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorf LA, Hunter DJ, McCarthy MI, Ramos EM, Cardon LR, Chakravarti A, et al.: **Finding the missing heritability of complex diseases.** *Nature* 2009, **461**:747-753.
11. Chothia C, Lesk AM: **The relation between the divergence of sequence and structure in proteins.** *Embo j* 1986, **5**:823-826.
12. Durbin R, et al.: *Biological Sequence Analysis*: Cambridge University Press; 1998.
13. Kryazhimskiy S, Plotkin JB: **The population genetics of dN/dS.** *PLoS Genet* 2008, **4**:e1000304.
14. Ng SB, Turner EH, Robertson PD, Flygare SD, Bigham AW, Lee C, Shaffer T, Wong M, Bhattacharjee A, Eichler EE, et al.: **Targeted capture and massively parallel sequencing of 12 human exomes.** *Nature* 2009, **461**:272-276.
15. Veltman JA, Brunner HG: **De novo mutations in human genetic disease.** *Nat Rev Genet* 2012, **13**:565-575.
16. Abecasis GR, Altshuler D, Auton A, Brooks LD, Durbin RM, Gibbs RA, Hurles ME, McVean GA: **A map of human genome variation from population-scale sequencing.** *Nature* 2010, **467**:1061-1073.
17. Tennessen JA, Bigham AW, O'Connor TD, Fu W, Kenny EE, Gravel S, McGee S, Do R, Liu X, Jun G, et al.: **Evolution and functional impact of rare coding variation from deep sequencing of human exomes.** *Science* 2012, **337**:64-69.

Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto

Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Red

Anurag Sethi 10/24/2015 10:03 PM  
Formatted: Font color: Auto

18. Abecasis GR, Auton A, Brooks LD, DePristo MA, Durbin RM, Handsaker RE, Kang HM, Marth GT, McVean GA: **An integrated map of genetic variation from 1,092 human genomes.** *Nature* 2012, **491**:56-65.
19. Khurana E, Fu Y, Colonna V, Mu XJ, Kang HM, Lappalainen T, Sboner A, Lochovsky L, Chen J, Harmanci A, et al.: **Integrative annotation of variants from 1092 humans: application to cancer genomics.** *Science* 2013, **342**:1235587.
20. Exome Aggregation Consortium (ExAC) on World Wide Web URL: <http://exac.broadinstitute.org>
21. Wang Z, Moutl J: **SNPs, protein structure, and disease.** *Hum Mutat* 2001, **17**:263-270.
22. Kumar P, Henikoff S, Ng PC: **Predicting the effects of coding non-synonymous variants on protein function using the SIFT algorithm.** *Nat Protoc* 2009, **4**:1073-1081.
23. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, Kondrashov AS, Sunyaev SR: **A method and server for predicting damaging missense mutations.** *Nat Methods* 2010, **7**:248-249.
24. Bromberg Y, Rost B: **SNAP: predict effect of non-synonymous polymorphisms on function.** *Nucleic Acids Res* 2007, **35**:3823-3835.
25. Li B, Krishnan VG, Mort ME, Xin F, Kamati KK, Cooper DN, Mooney SD, Radivojac P: **Automated inference of molecular mechanisms of disease from amino acid substitutions.** *Bioinformatics* 2009, **25**:2744-2750.
26. Perica T, Kondo Y, Tiwari SP, McLaughlin SH, Kemplen KR, Zhang X, Steward A, Reuter N, Clarke J, Teichmann SA: **Evolution of oligomeric state through allosteric pathways that mimic ligand binding.** *Science* 2014, **346**:1254346.
27. Ward LD, Kellis M: **Interpreting noncoding genetic variation in complex traits and human disease.** *Nat Biotechnol* 2012, **30**:1095-1106.
28. Albert FW, Kruglyak L: **The role of regulatory variation in complex traits and disease.** *Nat Rev Genet* 2015, **16**:197-212.
29. Rual JF, Venkatesan K, Hao T, Hirozane-Kishikawa T, Dricot A, Li N, Berriz GF, Gibbons FD, Dreze M, Ayivi-Guedehoussou N, et al.: **Towards a proteome-scale map of the human protein-protein interaction network.** *Nature* 2005, **437**:1173-1178.
30. Rolland T, Tasan M, Charlotheaux B, Pevzner SJ, Zhong Q, Sahni N, Yi S, Lemmens I, Fontanillo C, Mosca R, et al.: **A proteome-scale map of the human interactome network.** *Cell* 2014, **159**:1212-1226.
31. Kim PM, Korbelt JO, Gerstein MB: **Positive selection at the protein network periphery: evaluation in terms of structural constraints and cellular context.** *Proc Natl Acad Sci U S A* 2007, **104**:20274-20279.
32. Khurana E, Fu Y, Chen J, Gerstein M: **Interpretation of genomic variants using a unified biological network approach.** *PLoS Comput Biol* 2013, **9**:e1002886.
33. Goh KI, Cusick ME, Valle D, Childs B, Vidal M, Barabasi AL: **The human disease network.** *Proc Natl Acad Sci U S A* 2007, **104**:8685-8690.
34. Sharan R, Ulitsky I, Shamir R: **Network-based prediction of protein function.** *Mol Syst Biol* 2007, **3**:88.
35. Garcia-Alonso L, Jimenez-Almazan J, Carbonell-Caballero J, Vela-Boza A, Santoyo-Lopez J, Antinolo G, Dopazo J: **The role of the interactome in the maintenance of deleterious variability in human populations.** *Mol Syst Biol* 2014, **10**:752.
36. Ofra Y, Rost B: **Protein-protein interaction hotspots carved into sequences.** *PLoS Comput Biol* 2007, **3**:e119.
37. Aytuna AS, Gursoy A, Keskin O: **Prediction of protein-protein interactions by combining structure and sequence conservation in protein interfaces.** *Bioinformatics* 2005, **21**:2850-2855.
38. Gao M, Zhou H, Skolnick J: **Insights into Disease-Associated Mutations in the Human Proteome through Protein Structural Analysis.** *Structure* 2015, **23**:1362-1369.

Anurag Sethi 10/24/2015 10:03 PM

Formatted: Font color: Auto

39. Sahni N, Yi S, Taipale M, Fuxman Bass JI, Coulombe-Huntington J, Yang F, Peng J, Weile J, Karras GI, Wang Y, et al.: **Widespread macromolecular interaction perturbations in human genetic disorders.** *Cell* 2015, **161**:647-660.
40. Kim PM, Lu LJ, Xia Y, Gerstein MB: **Relating three-dimensional structures to protein networks provides evolutionary insights.** *Science* 2006, **314**:1938-1941.
41. Bhardwaj N, Abyzov A, Clarke D, Shou C, Gerstein MB: **Integration of protein motions with molecular networks reveals different mechanisms for permanent and transient interactions.** *Protein Sci* 2011, **20**:1745-1754.
42. Wang X, Wei X, Thijssen B, Das J, Lipkin SM, Yu H: **Three-dimensional reconstruction of protein networks provides insight into human genetic disease.** *Nat Biotechnol* 2012, **30**:159-164.
43. **An integrated encyclopedia of DNA elements in the human genome.** *Nature* 2012, **489**:57-74.
44. Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, Heravi-Moussavi A, Kheradpour P, Zhang Z, Wang J, Ziller MJ, et al.: **Integrative analysis of 111 reference human epigenomes.** *Nature* 2015, **518**:317-330.
45. Magger O, Waldman YY, Ruppin E, Sharan R: **Enhancing the prioritization of disease-causing genes through tissue specific protein interaction networks.** *PLoS Comput Biol* 2012, **8**:e1002690.
46. Ellis JD, Barrios-Rodiles M, Colak R, Irimia M, Kim T, Calarco JA, Wang X, Pan Q, O'Hanlon D, Kim PM, et al.: **Tissue-specific alternative splicing remodels protein-protein interaction networks.** *Mol Cell* 2012, **46**:884-892.
47. Buljan M, Chalancon G, Eustermann S, Wagner GP, Fuxreiter M, Bateman A, Babu MM: **Tissue-specific splicing of disordered segments that embed binding motifs rewires protein interaction networks.** *Mol Cell* 2012, **46**:871-883.
48. Dunker AK, Lawson JD, Brown CJ, Williams RM, Romero P, Oh JS, Oldfield CJ, Campen AM, Ratliff CM, Hipps KW, et al.: **Intrinsically disordered protein.** *J Mol Graph Model* 2001, **19**:26-59.
49. Kim PM, Sboner A, Xia Y, Gerstein M: **The role of disorder in interaction networks: a structural analysis.** *Mol Syst Biol* 2008, **4**:179.
50. Oldfield CJ, Dunker AK: **Intrinsically disordered proteins and intrinsically disordered protein regions.** *Annu Rev Biochem* 2014, **83**:553-584.
51. Eliezer D: **Biophysical characterization of intrinsically disordered proteins.** *Curr Opin Struct Biol* 2009, **19**:23-30.
52. Sethi A, Tian J, Vu DM, Gnanakaran S: **Identification of minimally interacting modules in an intrinsically disordered protein.** *Biophys J* 2012, **103**:748-757.
53. Colak R, Kim T, Michaut M, Sun M, Irimia M, Bellay J, Myers CL, Blencowe BJ, Kim PM: **Distinct types of disorder in the human proteome: functional implications for alternative splicing.** *PLoS Comput Biol* 2013, **9**:e1003030.
54. Flores SC, Keating KS, Painter J, Morcos F, Nguyen K, Merritt EA, Kuhn LA, Gerstein MB: **HingeMaster: normal mode hinge prediction approach and integration of complementary predictors.** *Proteins* 2008, **73**:299-319.
55. Hubbard S, Thornton J: **NACCESS, Computer Program.** Edited by: Department of Biochemistry Molecular Biology, University College London; 1993.