In Silico Identification of Prioritized Interacting Domains in Primary Immunodeficiency Disease Causing Genes

Suresh Kumar Ramadoss and Sujatha Mohan

Abstract—Primary immunodeficiency diseases (PIDs) are complex and intrinsic genetic disorders that leads to immune dysfunction. We have developed "Resource of Asian Primary Immunodeficiency Diseases", an open access database on PIDs. In the current study, we propose a heuristic approach of PID gene mutation data analysis based on the functional domain interactions. Through this approach, a list of functionally significant domains that disrupts PID genes' protein-protein interactions(PPI) associated with disease mutation have been identified. Moreover, the domains to be associated with immune diseases and function based on observed PID gene mutations have also been prioritized for further molecular characterization of disease pathogenesis.

Index Terms—Domain-Domain Interactio, Disease Mutation, DIMA ,HitPredict.

I. INTRODUCTION

Primary immunodeficiency diseases (PIDs) are genetic disorders resulting in abnormalities in the development and maintenance of immune system. Patients with these intrinsic defects have common and overlapping manifestations which pose daunting task to clinicians in providing definitive diagnosis based on determined sequence variations with observed phenotype. Towards this end, we have launched an open access PID database designated "Resource of Asian Primary Immunodeficiency Diseases (RAPID)" [1], a webbased compendium of molecular alterations and gene expression at the mRNA and protein levels of all PID genes reported from PID patients in the public literature. The database also includes other pertinent information about protein-protein interactions, mouse studies and microarray gene expression profiles in various organs and cells of the immune system. Along with this, RAPID also contains DNA sequencing protocols for selected PID genes. RAPID can be accessed at http://rapid.rcai.riken.jp.

At present, RAPID comprises a total of 207 PID genes out of which 196 genes are reported with over 4600 unique disease-causing mutation data obtained from more than 1500 citations.

Several studies indicate the importance of interacting domains and its disease association [2-7]. Also, there are proven results that explains well about the mutation(s) in

Email: sujatha@rcai.riken.jp.

domain-domain interaction sites disrupt the protein binding and its subsequent function [8-12]. It is apparent that protein sequence-specific functional domains play an important role in various key biological events including protein-protein interactions, post-translational modifications (PTMs) and so on. The diverse combination of domains gives the range of protein functions. The identification and understanding of protein domain henceforth provides an insight to its function. We understand that any mutations occurring in such functional domains would cause varying effects in the subsequent biological events. Moreover, intra-molecular and inter-molecular interactions of proteins take place through binding to the specific sites defined in 3-D domains. These types of interactions are frequently occurring in any signal transduction network [13,14]. Also, the availability of disease-specific mutation data at the sequence level should aid in characterizing these domains. With this collective knowledge, we propose a heuristic integrated in silico approach to analyze PID gene-specific mutations observed in the spanning region of the functional domains. With the successful implementation of prescribed methods, we intend to identify and prioritize the domains involved in immune function based on PID gene mutation data. This kind of systematic study on the reported mutations occurring in PID gene interacting domains will shed light on the domain function and genotype-phenotype correlation of PIDs thereby further enhance our current knowledge of PID pathogenesis.

II. METHODS

Data sources

A. Identification and Characterization Of PID Gene Interacting Domains

The complex descriptive biological data mining is more suitable using BioMarts' MartView system [15]. The Pfam database [16] is having one of the largest known collection of protein family domains. In order to obtain defined domain spanning region for PID genes and its interacting partners along with the defined gene ontology (GO) terms, we have used MartView with Pfam as latter being selected as a primary source signature database. Apart from this, MartView also hosted other necessary information about protein annotation and sequence features classification from InterPro [17]. Furthermore, to introduce stringency in the generated list of InterPro domains, we have considered only active site, binding site, conserved site, domain, family and post-translational modifications (PTMs) defined as 'InterPro type'. The compiled results of domain spanning regions for PID gene and its interactors are derived from BioMart by selecting Uniprot/Swiss-Prot [18] as source protein database

Manuscript received May 20, 2011; revised June 25, 2011.

Suresh Kumar Ramadoss. Author is with the Research Unit for Immunoinformatics, Research Center for Allergy and Immunology (RCAI), RIKEN Yokohama Institute, Japan.

Sujatha Mohan. Author (To whom correspondence should be addressed) is with the Research Unit for Immunoinformatics, Research Center for Allergy and Immunology (RCAI), RIKEN Yokohama Institute, Japan. Phone:+81-45-503-7034;Fax:+81-45-503-9694.

with its respective UniprotKB accession IDs as input data.

B. Protein-Protein Interaction (PPI)

For our analysis, the list of binary interaction data directly observed through experimental studies as reported in HitPredict [19] have been generated. HitPredict is a comprehensive resource of high confidence protein-protein interactions. To generate human interaction data sets, we have sorted and consolidated a list of non-redundant interaction pairs of known PID genes along with the binary interaction data as available in HPRD recent release [20].

C. Domain-Domain Interaction (DDI)

Domain Interaction Map (DIMA) is a comprehensive resource of functional and physical interactions among conserved protein domains [21]. In general, the domaindomain interactions and structurally known interactions have also been integrated from iPfam [22] and 3DID [23] databases in DIMA. We used this resource for mapping all available domain-domain interaction pairs representing PID genes and its interactors. Also, we obtained available DDI data for Homo sapiens with the default parameters as given in DIMA tool. In order to maintain reliable and accuracy of selected data sets, we have considered only experimentally proven DDI data for the final domain interaction analysis.

III. IMPLEMENTATION AND RESULTS

There are several factors involved in causing a disease phenotype by a particular mutation such as loss of protein stability and, functional protein interactions thereby, damaging essential inter and intra molecular interactions. To overcome huge computational task of analyzing ever increasing number of genes and its mutations found across full length sequence, we have implemented our approach using available open web resources for the analysis of protein interactions and domain-domain interactions.

For this analysis, we have obtained PID gene-specific mutation data from RAPID as on December 2010. To understand and dissect-out disease causing mutations observed in the defined functional domain spanning region of PID gene that disrupt the intensity of interaction, we mapped the PID gene mutations and its interacting domains. The workflow of our approach is depicted in Fig. 1. To accomplish this task, Python scripts have been written to extract and analyze the data from various open-accessible resources. First, we obtained the Pfam domain spanning region along with its domain function from MartView by submitting the Uniprot IDs of all PID genes and its interactors as well as for those missed entries, necessary details have been collected from available published literature. Next, we processed PPI data and obtained 2603 non-redundant interaction pairs for all available PID genes.

This is followed by generation of DDI data for the interaction pairs by mapping HitPredict and DIMA. In this step, if both the Pfam IDs of interacting domains in DIMA matches with both PID gene and its interactors' Pfam IDs respectively, then those interaction pairs were confirmed as true interacting partners. In other words, the interaction pairs were excluded if either interacting partner or its respective domains were not matched with each other. This stringent criterion has been implemented throughout the analysis in order to control and eliminate ambiguous entries thereby presenting only the best possible results. Among PID gene interaction pairs, we could identify only 341 pairs having the associated DDI evidences as reported from DIMA. Subsequently, we scanned the RAPID mutation data which are mapped to the protein-coding regions and obtained 105 PID genes. The mutation frequency has been calculated using the following formula:



Fig 1. Overview of in silico approach to prioritize interacting domain associated with high frequency of PID gene-specific mutations. Domain spanning region for PID gene and its interactors are retrieved from BioMart. PPI and DDI data sets are referred from HitPredict & HPRD and DIMA resources respectively. Functional domain annotation and prioritization are performed using UniProt and RAPID to gather PID gene-specific mutations with more than 80% frequency along with integration of available PPI and DDI data sets.

Moreover, we also observed that varied frequency of mutations has been occurred in the interacting domains of individual PID genes. But, we considered only the genes having the mutation frequency of 80% and above to prioritize its respective domains. Using this approach, we could filter 39 PID genes having frequency of 80% and above for the observed mutation in its interacting domains as well as a total of 33 prioritized functional domains associated with it. The overall statistics of interacting domain analysis and the list of prioritized domains are shown in Tables I and II respectively. The compiled results are given in a supplementary file and it can be accessed at http://rapid.rcai.riken.jp/RAPID/mut/domain results.xls.

TABLE I. STATISTICS OF PID GENE-SPECIFIC MUTATION-BASED INTERACTING DOMAIN ANALYSIS

Analyzed data set	Total number of analyzed data set	
PID gene interaction pairs	2603	
Domain-domain interaction pairs obtained from mapping PPI and DDI data sets	341	
PID genes' domain-domain interaction pairs with reported mutations	199	
PID genes with domain-domain interaction data	105	
PID genes with observed mutation frequency (≥ 80%) in the functional domain	39	
Prioritized Pfam domains having mutation frequency (> 80%)	33	

IV. DISCUSSION

Protein-protein interactions occurs through three levels of contacts viz., domain binding to other domain, domain binding to short protein motif, or motif binding to another motif [24,25]. It is quite obvious that experimental studies on the identification and characterization of functional domains involved in the disease-causing mutations are not easily doable because of laborious, time-consuming and expensive procedures. On the other hand, selection and implementation of integrated *in silico* approach for such studies employing various bioinformatics tools face its own challenges and limitations.

To bridge this gap, in the present work, we considered PID genes and its interactors binding through its respective domains wherever the experimentally confirmed data have become available.

Our results suggest that the prioritized domains have a significant role in immune function. For *in silico* evaluation of our analyzed results, we screened a few selected high frequency mutation entries for thorough description of PID gene-specific functional domain annotation as mentioned below:

Lectin C domain in C-type lectin domain family 7, member A (CLEC7A) is mainly responsible for recognizing pathogens and immune regulation [26]. Serine proteases including trypsin domain is particularly involved in innate immune response [27]. Also, ELANE gene encodes Elastase 2, neutrophil protein containing trypsin domain interacts with LPA having reported mutation frequency of 93.65% (as per the data shown in the supplementary file -'PID Gene Domain Analysis' sheet) that causes characteristic neutropenia in the diagnosed patients. The collagen domain in human Surfactant protein A plays a major role in innate immune defense mechanism [28]. These studies clearly illustrated the potential role of prioritized domains (Table II) in human immune system. The list of domains highlighted in bold letters (in Table II) are having at least one of the defined UniProt sequence feature description viz., phosphorylation site, sites of disulfide bond, active site, binding site and other regions of interest. Thus, these annotated sequence features serve as added evidences for its functional importance.

Earlier our collaborative lab published a mutation evaluation tool [29] using SIFT program [30]. It is a webbased integrated bioinformatics tool to identify and analyze novel and known human gene mutations causing genetic diseases. This study clearly demonstrates a new dimension towards identification and description of PID causing mutations in its functional domains in human disease pathogenesis.

In a similar way, our integrated approach should provide necessary insight into functionally defined domains based on the frequency of observed mutations. Such mutation based domain analysis in the PID genes should help further to demonstrate genotype-phenotype correlation wherever functional studies are not available in the published literature. Although, there are many factors involved in deducing the effect of a particular mutation, we consider the conserved functional domain as our primary aspect in analyzing the mutation data. Based on this collective domain interaction data, it has been observed that about 56% (in

average) of total mutations are found in the interacting domains of PID genes. However, for individual PID gene, the frequency of mutations in interacting domains varies from 2.17% to 100%. The obtained data clearly demonstrates that more functional annotation of domains in disease causing genes can be identified through such in silico analysis. To check this hypothesis, we used domain prediction using context (DPUC) [31] and DomFOLD [32] tools. These tools provide computationally predicted domains for given protein sequences. We found that there are more regions predicted as domains apart from known domains, for the given PID gene amino acid sequences having mutations other than known DDI region (data not shown). A note of caution is that successful implementation of such studies should require appropriate selection of predicted domain-domain interaction data along with other supporting evidences including functional annotation. Overall, these domain regions are considered to be functionally involved in protein interactions and have biological significance in any immune signaling pathways. PID causing through such genetic defects must involve a disruptive interaction in the functional pathway leading to expression of specific phenotypic traits. These observations clearly suggest us more integrated computational methods are to be implemented for in-depth analysis of protein domain prediction along with other available supportive experimental evidences that should provide global picture on the effect of functional domain mutations in the disease pathogenesis. Also, these domains should provide further clue for screening the PID candidate genes, as identified by our group using the Support Vector Machine (SVM) learning algorithm [33], having similar mode of DDI pairs involved in a disease-specific pathway.

It is also apparent that disrupted interaction in a signaling pathway leads to a phenotype trait [34-36]. To delve further, we will incorporate and integrate these results in our ongoing project in the development and construction of PID specific immune signaling pathways.

We propose to extend the same approach in identifying the right set of data for domain-motifs and motif-motif interactions including orthologous species' data. This will certainly broaden our research scope for a set of characterized short motifs with its own generic definitions. It is now quite obvious that structural information is also critical and necessary for thorough understanding and dissecting out selected domain level interactions. Therefore, we intend to further study mutation-specific domain-domain interactions based on the collected structural interaction data as given in iPfam.

V. CONCLUSIONS

With successful implementation of our integrated approach, it would assist in deriving genotype-phenotype correlation thereby improving phenotype-based genetic analysis of PID genes. Moreover, this kind of mutation analysis should augment well with the understanding of domain specific interaction and its impact on the disease pathogenesis. Eventually, it should also facilitate clinicians in confirming early PID diagnosis and proper therapeutic interventions.

Domain name	Pfam ID	Molecular function class	Reported PID gene	Distribution of mutations (average mutation frequency in %)
7tm_1	PF00001	G-protein coupled receptor protein signaling pathway	FPR1	100
ABC_membrane	PF00664	ATPase activity, coupled to transmembrane movement of substances	TAP2	100
Actin	PF00022	Protein binding	ACTB	100
ApoL	PF05461	Lipoprotein metabolic process	APOL1	100
C1q	PF00386	Complement activation	C1QA	100
Collagen	PF01391	Connective tissue formation	MBL2	100
CUB	PF00431	Complement activation; Cell signaling; Tissue repair	MASP2	100
Cytochrom_B558a	PF05038	Heme binding	СҮВА	100
Death	PF00531	Signal transduction	FADD	100
DUF1650	PF07856	Mediation of CRAC channel activity	ORAI1	100
FAT	PF02259	Protein binding	PRKDC	100
Fibrinogen_C	PF00147	Signal transduction; Receptor binding	FCN3	100
IL6Ra-bind	PF09240	Cytokine binding	CSF2RA	100
Interfer-bind	PF09294	Ligand binding	IL10RB	100
Lectin_C	PF00059	Carbohydrate-binding activity	CLEC7A	100
MACPF	PF01823	Transmembrane channel formation	C8A	100
Peptidase_C14	PF00656	Cysteine-type endopeptidase activity	CASP10; CASP8	100
PX	PF00787	Cell communication; Phosphoinositide binding	NCF4	100
RAG2	PF03089	DNA recombination; DNA binding	RAG2	100
Sushi	PF00084	Complement control protein	CFHR1	100
Tetraspannin	PF00335	Signal transduction	CD81	100
UPAR_LY6	PF00021	Membrane attack complex inhibition factor	CD59	100
V_ATPase_I	PF01496	Proton-transporting two-sector ATPase complex, proton-transporting domain	TCIRG1	100
V-set	PF07686	Cell-surface receptor	CD8A; CD79B	100
WD40	PF00400	Signal transduction; Transcription regulation	CORO1A	100
A_deaminase	PF00962	Deaminase activity; Purine ribonucleoside monophosphate biosynthetic process	ADA	98.08
Trypsin	PF00089	Serine-type endopeptidase activity; Proteolysis	CFD; ELANE	96.83
MFS_1	PF07690	Transmembrane transport	SLC37A4; SLC46A1	96.66
Ras	PF00071	Guanine nucleotide exchange factors interaction	NRAS; RAC2; RAB27A	94.44
An_peroxidase	PF03098	Heme binding; Peroxidase activity	МРО	90.91
PNP_UDP_1	PF01048	catalytic activity; nucleoside metabolic process	NP	90.48
Serpin	PF00079	Serine-type endopeptidase inhibitor activity	SERPING1	83.77
Cobalamin_bind	PF01122	Cobalamin binding	TCN2	80

TABLE II. List of PID Gene Domains Reported with Disease-Causing Mutations* from RAPID

*PID gene domains having 80% and above percentage of mutations are considered. Note that bold entries have mutations in functional sites as described in Uniprot's sequence features

VI. FUTURE PERSPECTIVES

Future efforts will aim towards integrated data analyses, thereby providing more comprehensive picture of PID biology, i.e. a prerequisite for identification of any potential diagnostic and prognostic markers along with improved patients' therapeutic modalities.

ACKNOWLEDGMENT

The authors thank the research scientists at Institute of Bioinformatics, India for their collaboration in developing RAPID. Also, we thank all PID physicians involved in the PID Japan project as well as RAPID – PID experts and Dr. Ashwini Patel, Human Genome Center, Institute of Medical Science, University of Tokyo for their valuable inputs and suggestions.

REFERENCES

- Keerthikumar, S., R. Raju, *et al.* (2009^a). "RAPID: Resource of Asian Primary Immunodeficiency Diseases." *Nucleic Acids Res* 37(Database issue): D863-7
- [2] Argentaro, A., J. C. Yang, *et al.* (2007). "Structural consequences of disease-causing mutations in the ATRX-DNMT3-DNMT3L (ADD) domain of the chromatin-associated protein ATRX." *Proc Natl Acad Sci U S A* 104(29): 11939-44
- [3] Pippal, J. B., Y. Yao, et al. (2009). "Structural and functional characterization of the interdomain interaction in the mineralocorticoid receptor." *Mol Endocrinol* 23(9): 1360-70
- [4] Perreau, V. M., S. Orchard, *et al.* "A domain level interaction network of amyloid precursor protein and Abeta of Alzheimer's disease." *Proteomics* 10(12): 2377-95
- [5] Limviphuvadh, V., S. Tanaka, et al. (2007). "The commonality of protein interaction networks determined in neurodegenerative disorders (NDDs)." *Bioinformatics* 23(16): 2129-38
- [6] George, R. A., J. Y. Liu, et al. (2006). "Analysis of protein sequence and interaction data for candidate disease gene prediction." *Nucleic* Acids Res 34(19): e130
- [7] Gandhi, P. N., S. G. Chen, A. L. Wilson-Delfosse (2009). "Leucinerich repeat kinase 2 (LRRK2): a key player in the pathogenesis of Parkinson's disease." *J Neurosci Res* 87(6): 1283-95
- [8] Jothi, R., P. F. Cherukuri, *et al.* (2006). "Co-evolutionary analysis of domains in interacting proteins reveals insights into domain-domain interactions mediating protein-protein interactions." *J Mol Biol* 362(4): 861-75.
- [9] Kelly, L., R. Karchin, *et al.* (2007). "Protein interactions and disease phenotypes in the ABC transporter superfamily." *Pac Symp Biocomput*: 51-63.
- [10] Kobayakawa, T., S. Yamada, et al. (2009). "Single nucleotide polymorphism that accompanies a missense mutation (Gln488His) impedes the dimerization of Hsp90." Protein J 28(1): 24-8.
- [11] Tateishi, H., M. Yano, et al. (2009). "Defective domain-domain interactions within the ryanodine receptor as a critical cause of diastolic Ca2+ leak in failing hearts." Cardiovasc Res 81(3): 536-45.
- [12] Pang, E. and K. Lin "Yeast protein-protein interaction binding sites: prediction from the motif-motif, motif-domain and domain-domain levels." *Mol Biosyst* 6(11): 2164-73.
- [13] Stein, A., R. A. Pache, *et al.* (2009). "Dynamic interactions of proteins in complex networks: a more structured view." *FEBS J* 276(19): 5390-405.
- [14] Shimizu, K. and H. Toh (2009). "Interaction between intrinsically disordered proteins frequently occurs in a human protein-protein interaction network." *J Mol Biol* 392(5): 1253-65.
- [15] Haider, S., B. Ballester, et al. (2009). "BioMart Central Portalunified access to biological data." Nucleic Acids Res 37(Web Server issue): W23-7
- [16] Finn, R. D., J. Mistry, et al. (2010). "The Pfam protein families database." Nucleic Acids Res 38(Database issue): D211-22
- [17] Hunter, S., R. Apweiler, et al. (2009). "InterPro: the integrative protein signature database." Nucleic Acids Res 37(Database issue): D211-5

- [18] UniProt consortium. "The Universal Protein Resource (UniProt) in 2010." Nucleic Acids Res 38(Database issue): D142-8
- [19] Patil, A., K. Nakai, *et al.* (2011). "HitPredict: a database of quality assessed protein-protein interactions in nine species." *Nucleic Acids Res.* (in press)
- [20] Keshava Prasad, T. S., R. Goel, et al. (2009). "Human Protein Reference Database--2009 update." Nucleic Acids Res 37(Database issue): D767-72
- [21] Luo, Q., P. Pagel, et al. (2010). "DIMA 3.0: Domain Interaction Map." Nucleic Acids Res. (in press)
- [22] Finn, R. D., M. Marshall, et al. (2005). "iPfam: visualization of protein-protein interactions in PDB at domain and amino acid resolutions." *Bioinformatics* 21(3): 410-2.
- [23] Stein, A., A. Ceol, et al. (2011). "3did: identification and classification of domain-based interactions of known threedimensional structure." *Nucleic Acids Res* 39(Database issue): D718-23.
- [24] Aloy, P. and R. B. Russell (2002). "Interrogating protein interaction networks through structural biology." *Proc Natl Acad Sci U S A* 99(9): 5896-901.
- [25] Aloy, P., H. Ceulemans, et al. (2003). "The relationship between sequence and interaction divergence in proteins." J Mol Biol 332(5): 989-98.
- [26] Cambi, A. and C. G. Figdor (2003). "Dual function of C-type lectinlike receptors in the immune system." *Curr Opin Cell Biol* 15(5): 539-46
- [27] Dale, C. and N. Vergnolle (2008). "Protease signaling to G proteincoupled receptors: implications for inflammation and pain." J Recept Signal Transduct Res 28(1-2): 29-37
- [28] Garcia-Verdugo, I., G. Wang, et al. (2002). "Structural analysis and lipid-binding properties of recombinant human surfactant protein a derived from one or both genes." *Biochemistry* 41(47): 14041-53.
- [29] Hijikata, A., R. Raju, et al. "Mutation@A Glance: an integrative web application for analysing mutations from human genetic diseases." DNA Res 17(3): 197-208
- [30] Ng, P. C. and S. Henikoff (2003). "SIFT: Predicting amino acid changes that affect protein function." *Nucleic Acids Res* 31(13): 3812-4
- [31] Ochoa, A., M. Llinas, et al. (2011). "Using context to improve protein domain identification." BMC Bioinformatics 12: 90.
- [32] Roche, D. B., M. T. Buenavista, et al. (2011). "The IntFOLD server: an integrated web resource for protein fold recognition, 3D model quality assessment, intrinsic disorder prediction, domain prediction and ligand binding site prediction." *Nucleic Acids Res.*
- [33] Keerthikumar, S., S. Bhadra, *et al.* (2009). "Prediction of candidate primary immunodeficiency disease genes using a support vector machine learning approach." *DNA Res* 16(6): 345-51.
 [34] Gong, Z., Y. W. Cho, *et al.* (2009). "Accumulation of Pax2
- [34] Gong, Z., Y. W. Cho, et al. (2009). "Accumulation of Pax2 transactivation domain interaction protein (PTIP) at sites of DNA breaks via RNF8-dependent pathway is required for cell survival after DNA damage." J Biol Chem 284(11): 7284-93
- [35] Jiang, W., R. Sordella, et al. (2005). "An FF domain-dependent protein interaction mediates a signaling pathway for growth factorinduced gene expression." Mol Cell 17(1): 23-35
- [36] Kann, M. G. (2007). "Protein interactions and disease: computational approaches to uncover the etiology of diseases." *Brief Bioinform* 8(5): 333-46.

Ramadoss Suresh Kumar: This author was born on the 3rd November 1979 at Salem, Tamilnadu, India. In the year 2004, he obtained his Master of Technology (M.Tech) degree in Bioinformatics from Sathyabama University, India.

Currently he is working as RESEARCH ASSOCIATE since June 2009 at Research unit for Immunoinformatics, RCAI, RIKEN Yokohama Institute, Japan. He served as SYSTEM PROGRAMMER at kCube Consultancy Services, India during 2006-2009.

Mohan Sujatha: This author was born on the 7th January 1969 at Srirangam, Tiruchirapalli, Tamil Nadu, India. She obtained her Master of Science (M.Sc.) and Master of Philosophy (M.Phil.) degrees in biological sciences from University of Madras, India.In the year 1999, Doctorate of Philosophy (Ph.D.) in Biological Sciences was awarded from University of Madras, India

Since 2007, she has been working as a research unit leader at Research unit for Immunoinformatics, rcai, riken yokohama Institute, Japan.