

Preface

DNA- and RNA-binding proteins play a central role in gene regulation, which includes transcriptional control, alternative splicing, post-translational and transcriptional modifications like methylation and acetylation among other roles. In this way, they control most of the working machinery of the cell in direct or indirect manner. Although more than 60 years ago the structure of DNA was proposed by Watson and Crick, our understanding of how RNA- and DNA-binding proteins interact with the genome and transcriptome remains scarce. One of the most important questions in biology is how a large number of DNA- and RNA-binding proteins find their target, interact and later disassociate. These nucleic acid binding proteins either recognizes the unique structural and chemical signatures of the bases (base readout) which give the specificity or it recognizes a sequence-dependent shape (shape readout).

Methyltransferases are enzymes with diverse folds, which perform methyltransfer to various substrates using mainly S-adenosyl-L-methionine (AdoMet) as a methyl donor. RNA methylation is one of the most crucial post-transcriptional modifications which influences a wide variety of cellular processes like metabolic stabilization of RNA, quality control in protein synthesis, resistance to antibiotics, mRNA reading frame maintenance, splicing, viral nucleoprotein stabilization among others. Specificity in recognition and methylation in ribosomal RNA (rRNA) methyltransferases is very crucial, as rRNA is highly conserved and lack of specificity would influence the stabilization of RNA and thus, will affect the ribosome. In recent years, rRNA modifications which confer resistance to ribosomal antibiotics have also been observed. **The mechanism of recognition to their unique rRNA target site with high selectivity and their evolution still remains an enigma. Thus, the evolution of antibiotic resistance-conferring methyltransferases in pathogenic organisms needs to be investigated from the structural and evolutionary perspective.**

In the last two decades, many global regulators in both eukaryotes and prokaryotes have been discovered, which promiscuously bind to a large number of DNA sequences. In prokaryotes, they are called as ‘Nucleoid-associated proteins’ (NAPs), which influence the transcriptional process and exhibit multi-specificity or promiscuity. They also take part in the formation of many multi-protein complexes. HU and Integration Host Factor (IHF) are NAPs which belong to prokaryotic DNA-binding protein family (DNABII family). HU and IHF play crucial architectural roles in bacterial DNA condensation and additionally play a regulatory role in many cellular processes. Although sharing structural similarity, the DNA binding and bending features of HU and IHF are strikingly different, allowing them to selectively regulate genes from different genomic locations. HU binds to DNA in a sequence promiscuous manner while IHF is moderately sequence specific. **The molecular mechanism of DNA binding multi-specificity (differential specificity with varied binding affinity) of HU/IHF proteins remains unexplored, as little attention has been paid to the determinants at the sequence level.**

Now, the fundamental question which the author attempted to understand is the structural and evolutionary determinants of specificity in DNA- and RNA-binding proteins. The candidate has taken nucleoid-associated protein HU and SPOUT superfamily RNA methyltransferase as model systems. As the very limited number of structural folds makes up the DNA- and RNA-binding proteins, it is intriguing to examine closely related nucleic acid binding domains or folds carrying out specific functions. Also, we observed that some proteins having a particular structural fold (or homologous ancestry) bind to DNA or RNA with high specificity, while its other homolog binds promiscuously. These observations tempted us to find the sequence and structural determinants which guide this phenomenon, not just specific to only a single protein family, but, determinants are of more general nature, where results can possibly be applied to other nucleic acid binding proteins too.

The first part of the thesis reports the **crystal structures of native and AdoMet bound ribosomal RNA Methyltransferase from *Sinorhizobium meliloti* (smMtase)**, by single anomalous dispersion (SAD) phasing on seleno-methionine

substituted crystal, which diffracted to **2.28Å and 2.9 Å** resolutions respectively in space group P2₁2₁2₁. smMtase belong to an rRNA binding SPOUT superfamily protein, which is fused with an RNA binding L30e domain at the N-terminus. We focused our study on these types of proteins among the large superfamily (henceforth termed as SPOUT_{L30}).

The author also has conducted a **phylogenetic study, which revealed 11 major clades, out of which we focused our present study in understanding the sequence conservation and variations of 5 (A-E) clades**, for which structural, biochemical and functional data is available. These proteins share homology to antibiotic resistance-conferring methyltransferases. **The availability of experimentally determined structures of native and AdoMet bound smMtase along with an analysis of other homologous crystal structures has enabled a critical examination of factors influencing RNA binding specificity. Also, the thesis reports for the first time an evolutionary and structural inter-connectivity of the three conserved motifs (I-III) in SPOUT superfamily, which is responsible for AdoMet binding and catalysis.** The results highlight that both the location of conserved positive and negatively charged residues influence the RNA binding, specificity, and affinity. The conservation of these residues could be at superfamily, family or at clade level, and the position of these charged residues at specific sites, alters their salt-bridge geometry, which ultimately fixes the conformation of RNA-binding residues, thus defining a particular binding site specific to its cognate RNA. **The study conducted by the author reveals that the dynamics of salt-bridge and other directional interactions like hydrogen bonding and aromatic interactions essentially determines the specificity of SPOUT_{L30}.**

The second part of the thesis reports evolutionary, structural and functional studies on nucleoid-associated proteins HU and IHF. To understand the sequence determinants, which influence the degree of DNA binding specificity, we undertook a phylogenetic study in conjunction with analysis of three-dimensional structures. **The phylogenetic analysis revealed three major clades, belonging to HU, IHF α , and IHF β like proteins with reference to *E. coli*. The author observed statistically significant**

amino acid compositional bias in the DNA binding sites of HU and IHF clade proteins. The author proposes that the molecular mechanisms giving rise to specificity or multi-specificity depend on a combination effect of the amino acid composition of the binding site, its flexibility, ionic and steric constraints. In continuation of this part of the thesis, the candidate examined **the role of protein interacting interface of HU-IHF family proteins, understanding its evolutionary history and utilizing it in designing inhibitors for *Mycobacterium tuberculosis* HU (MtbHU).** The present results give a model example of an evolutionary study of a protein interface of nucleoid-associated protein, which is used to understand the interface and computationally design inhibitors targeting it.

The author was a part of the study (Bhowmick et al. 2014, *Nature communications*) which has determined the crystal structure of *Mycobacterium tuberculosis* HU, inhibited it using stilbene derivatives (SD1 and SD4) which curtailed the Mtb cell growth. In the present thesis, the candidate observed **from microarray analysis that the SD1 stimulon consists of genes involved majorly in lipid biosynthesis pathway, ribosomal genes which affect the overall translation, aerobic respiration pathways, antigenic membrane proteins involved in pathogenicity.** Nearly half of the genes in affected by SD1 are essential in nature, thus could explain the curtailing of cellular growth. The whole study provides a system inspired view of probing as well, inhibiting global regulator HU using novel chemical molecules.

The overall structure of the thesis is further explained through a chapter wise description below:

Chapter 1 | Introduction

The chapter starts with an introduction of the evolution of specificity and promiscuity in substrate recognition by various proteins and the definition of the problem of the thesis. DNA or RNA binding proteins interacts with its substrate/s at its cognate binding pocket, whose nature decides the strength and specificity of the interaction. In both the cases, the nature and the flexibility of the binding site along with other static

and dynamic factors influence the specificity of the protein towards its substrate. Now, the question arises that why do we need specific enzymes or transcription factors if the more promiscuous or multi-specific ones can perform the job? Also, the contrary argument arises to challenge the need for promiscuous ones when specific enzymes or transcription factors can do their job optimally. Broad specificity also provides global regulatory control and response to various stress conditions. These proteins can switch its binding preferences to a certain set of substrates based on the needs of the organism. Highly specific DNA- and RNA-binding proteins cannot perform such global regulatory role. The need for specificity in such macromolecular interaction is to provide greater efficiency, finer regulation, or prevention of deleterious interactions.

The chapter gives a review to promiscuity in enzymes, membrane channel, DNA and RNA binding, immune systems and signaling proteins. It also discusses about RNA methyltransferases which perform one of the most crucial post-transcriptional modifications which influences a wide variety of cellular processes like metabolic stabilization of RNA, quality control in protein synthesis, resistance to antibiotics, mRNA reading frame maintenance, splicing, viral nucleoprotein stabilization and many other processes methylating different RNA species like tRNA, rRNA, snRNA, mRNA, tmRNA, snRNA, snoRNA, miRNA, and viral RNA. The focus of this section is mainly towards class IV SPOUT family methyltransferases, which confers antibiotic resistance in bacteria through modification of nucleotides in 23S ribosomal RNA. Next, the chapter gives a brief account of nucleoid-associated proteins, focusing on HU and IHF family proteins. The chapter discusses binding preferences of HU and IHF toward different DNA substrates, diverse functions of HU like proteins in cellular and extracellular milieu, their DNA protection and association with low complexity regions and its importance as a drug target.

Chapter 2 | Crystallization, Single anomalous dispersion (SAD) phasing and refinement of ribosomal RNA Methyltransferase from *Sinorhizobium meliloti*

Chapter 2 presents the crystallization, data collection and structure refinement of an rRNA methyltransferase from *Sinorhizobium meliloti* (smMtase) in native and AdoMet bound form. Cloning and expression were performed at New York Structural

Genomics Research Consortium (NYSGRC). The protein crystallizes only in the presence of a specific combination of divalent metal ions. In general, metals are known to play an important role in oligomerization and crystallization of proteins, and a recently described crystallization strategy utilized surface mutations that support the formation of metal-mediated crystal contacts.

Initial screening for crystallization condition(s) was performed with Crystal Screen HT and Index Screen HT from Hampton Research, USA, using the sitting drop vapor diffusion method in 96-well Intelliplates (Art Robbins). Sitting drops with a 1:1 mixture of 0.8 μ l each of the protein solution (20 mg/ml in 20mM HEPES pH7.5, 150 mM NaCl, 10% vol/vol Glycerol, 0.1% wt/vol Sodium Azide, 0.5mM TCEP) and reservoir solution were equilibrated against 40 μ l of reservoir solution. Small rod-like crystals were observed in a condition containing 5 mM Cobalt(II) chloride hexahydrate (CoCl_2), 5 mM Nickel(II) chloride hexahydrate (NiCl_2), 5 mM Cadmium chloride hydrate (CdCl_2), 5 mM Magnesium chloride hexahydrate (MgCl_2), 0.1 M HEPES pH 7.5, 12% wt/vol PEG 3,350. This condition was refined to obtain better crystals. Sitting drops of 2 μ l (1:1 protein and reservoir solution) in a 24-well plate (Hampton Research, USA) were set up with varying concentrations of PEG (12%, 9% and 6% wt/vol) and a combination of the above four different metal chlorides: 0.1 M HEPES buffer, pH7.5, with varying concentrations of PEG drops were set up with (i) only CoCl_2 , (ii) only NiCl_2 , (iii) only CdCl_2 (iv) CoCl_2 and NiCl_2 (v) CoCl_2 and CdCl_2 (vi) NiCl_2 and CdCl_2 (vii) CoCl_2 , NiCl_2 and CdCl_2 and (viii) CoCl_2 , NiCl_2 , CdCl_2 and MgCl_2 . After two days, a number of small crystals were observed only in the condition containing all the metals. Since the MgCl_2 was missing in all the other conditions except for the condition viii, the only condition, which gave crystals, 5 mM MgCl_2 was added to all the other drops. After three more days, crystals larger than those from the condition with all four metals were observed in the condition with NiCl_2 , CdCl_2 and MgCl_2 with 9% wt/vol PEG 3350 as precipitant. This gave the hint that Mg^{2+} was perhaps essential for crystallization, hence, we repeated crystallization attempts by first incubating the protein with 100mM MgCl_2 for one hour prior to crystallization trials with conditions containing 0.1 M HEPES pH 7.5, 12% wt/vol PEG 3,350 and

one of the five divalent metal ions (Ni^{2+} , Cd^{2+} , Zn^{2+} , Cu^{2+} or Co^{2+}). Subsequent experiments confirmed that crystals can be grown in conditions containing Mg^{2+} in combination with any of the individual metals Ni^{2+} , Cd^{2+} , Zn^{2+} , Cu^{2+} or Co^{2+} , suggesting that at least two metals are required for crystallization. It was also observed that the condition containing Co^{2+} or Ni^{2+} in combination with Mg^{2+} yielded the best crystals. **The protein crystallized in the space-group $\text{P2}_1\text{2}_1\text{2}_1$ and diffracts to a maximum resolution of 3.1 Å initially which after further improvement of crystals diffracted to 2.28 Å for the native protein. AdoMet bound smMtase crystal diffracted to 2.9 Å** (see Table 1).

Table 1

| Parameters | Native smMtase | smMtase bound to AdoMet |
|---|-----------------------------------|-----------------------------------|
| Diffraction source | BM14 ESRF | Home source |
| Wavelength (Å) | 0.954 | 1.542 |
| Space group | $\text{P2}_1\text{2}_1\text{2}_1$ | $\text{P2}_1\text{2}_1\text{2}_1$ |
| <i>a, b, c</i> (Å) | 65.82, 84.30, 112.04 | 64.85, 82.89, 108.80 |
| α, β, γ (°) | 90, 90, 90 | 90, 90, 90 |
| Resolution range (Å) | 50-2.27 (2.31-2.27) | 50-2.95 (3.0-2.95) |
| Completeness (%) | 98.0 (98.4) | 100 (100) |
| Redundancy | 7.9 (7.7) | 6.7 (6.3) |
| $\langle I/\sigma(I) \rangle$ | 46.98 (2.86) | 28.51 (2.10) |
| CC1/2 # | 0.919 | 0.86 |
| No. of molecules in asymmetric unit | 2 | 2 |

Ab initio phasing was attempted using the SHELXC/D/E pipeline with HKL2MAP GUI. smMtase was built extending on the poly-alanine trace generated by SHELXE and iterative model building with COOT. The initial diffraction data (at 3.1 Å), resulted in an initial model with clear density for most of the residues of the C-terminal SPOUT domain. However, we were not able to unambiguously define all the residues, particularly the N-terminal L30 domain from this data, as this domain was relatively disordered and the resolution of the data was also low. The AutoBuild Wizard used XTRIAGE, RESOLVE, and Phenix refine was used to assess the quality of the data, for density modification and automatic building and refinement respectively. The resolution of the data was further improved by looking for an appropriate Cryo solution. A crystal cryo-cooled with 10% DMSO diffracted to

resolution $\sim 2.7\text{\AA}$. Although this data was helpful in improving the model, it was not sufficient to completely define all the residues of the L30 domain. Crystals were grown with Mg^{2+} and with only one transition metal Co^{2+} or Ni^{2+} gave better quality crystals and also diffracted to $\sim 2.3\text{\AA}$, where mother liquor was supplemented with 10% DMSO and 0.1M LiSo_4 in the cryo-solution. The model was improved by iterative cycles of model building with COOT and refinement with Refmac5. The final model contained most of the residue except for the residues 1-15, 71-77 and 285-288 of chain A and residues 1-15, 71-77, 260-262 and 284-288 of chain B. **The final round of refinement was carried out with restrained refinement with TLS parameters which resulted R-factors of $R_{\text{work}} = 20.6\%$ and $R_{\text{free}} = 25.9\%$ (resolution range 67.3 to 2.28 \AA) for the native smMtase structure.**

The AdoMet soaked crystals diffracted to 2.9 \AA . Since we did not observe considerable variation between the cell-dimensions of native and AdoMet soaked crystals, the apo structure (protein atoms only) was refined against the data collected from the AdoMet soaked crystals. The structure was subjected 20 cycles of rigid body refinement followed by ten cycles of restrained refinement in refma5. **The final structure was refined with $R_{\text{work}} = 19.7\%$ and $R_{\text{free}} = 26.7\%$** (for in the resolution range 65.9 to 2.9 \AA) for AdoMet bound smMtase. The final refinement statistics for native and AdoMet bound smMtase structure is presented in Table 2

Table 2

| Property | Native Mtase | AdoMet bound Mtase |
|---|--|---|
| Resolution (\AA) | 67.36- 2.28 | 65.94- 2.90 |
| PDB ID | 5KZK | 5L0Z |
| R_{work} | 0.206 | 0.191 |
| R_{free} | 0.259 | 0.265 |
| Fo-Fc correlation | 0.93 | 0.93 |
| Average B, all atoms (\AA^2) | 39.0 | 41.0 |
| Ramachandran Outlier | Analyzed:513/521(98%) Favored:496 (94.1%) Allowed:17 (5.9%) Outlier:0 | Analyzed: 517/523 (99%) Favored: 492 (91%) Allowed: 22 (9%) Outlier: 0 |

Chapter 2 also gives details of phasing and refinement of the native and AdoMet bound structures.

Chapter 3 | Structural insights into rRNA recognition based on the crystal structures of native and AdoMet bound rRNA Mtase from *Sinorhizobium meliloti*

Chapter 3 describes the crystal structures of native and AdoMet bound smMtase and details the evolutionary and structural analysis of smMtase which is aimed to provide structure based rationalizations of RNA binding specificity and factors that determine catalysis. A comprehensive structural analysis of native and AdoMet bound Mtase was performed to understand the knot region, mechanism of catalysis, AdoMet binding, intermolecular interactions which stabilize the dimer interface. Interactions between the two subunits (L30e and SPOUT), which determines the flexibility of N-terminal L30e domain. **The AdoMet binding site is formed by Thr212, Gly235, Gln238, Ile255 and other conserved hydrophobic residue constituting the conserved motifs (I-III).** Three conserved sequence patterns were identified which characterizes SPOUT methyltransferase superfamily: motif I ([PAVH]-X-N-X-G-X₃-R) responsible for catalysis, motif II (h-[VLIM]-h-G-X-E-X₂-G-h), and motif III (I-P-X₆-S-L-N-h), where 'h' is hydrophobic and X is any residue is involved in AdoMet binding. Our evolutionary analysis found 11 major clades belonging to SPOUT family proteins fused with the L30 domain, majorly involved in rRNA methyltransfer activity. **The analysis by the author shows that the conserved motifs (I-III) are interconnected with hydrogen bonding and salt-bridge interaction and form the AdoMet binding and catalytic site.**

In smMtase, Arg107, Arg108, Pro111, in the L30 domain of the catalytic monomer are specific to clade A, while Lys53 is conserved in SPOUT superfamily, forms the RNA recognition/binding residues. In the SPOUT domain of the catalytic monomer, Arg153, and Arg179, both conserved across SPOUT superfamily, form the RNA recognition site. We also observed that most, conserved positively charged residues are constrained by salt-bridges, to make specific interaction with RNA backbone. **We also found Arg33, Lys38, Lys39, Arg41 and Arg107 contributes to the RNA recognition in the L30 domain of the non-catalytic monomer,** with only Arg41 being conserved, while others are specific to the clade. Arg141, Arg143, and Arg254 (all clade-specific residues) forms the potential

RNA-binding residues in the non-catalytic SPOUT domain. **Our study hints that the salt-bridge is crucial for RNA interaction, and loss of negatively charged positions can also influence RNA binding and methyltransfer activity.** Our study supports the model in which the initial binding is contributed by the C-terminal domain, which is further adjusted by the N-terminal L30e domain for efficient catalysis. So, it appears as if there is a two- step verification process in the methyltransfer of these enzymes. From these results, we infer that not only the catalytic monomer but the non-catalytic monomer also plays a crucial role in RNA recognition and structural rearrangements of the RNA. The analyses presented in chapter3 are in agreement with other biochemical and structural studies on SPOUT_{L30} proteins.

Chapter 4 | Structural and evolutionary analyses reveal determinants of DNA binding specificities of nucleoid-associated proteins HU and IHF

Chapter 4 details the evolutionary and structural analyses which are aimed towards understanding the differences in DNA binding specificity among nucleoid-associated protein paralogs HU and IHF. Specificity and differential binding affinity of HU/IHF proteins towards their target binding sites play a crucial role in their regulatory dynamics. Decades of biochemical and genomic studies have been carried out for HU and IHF like proteins. **Yet, questions related to their DNA binding specificity, and differential ability to bend DNA thus affecting the binding site length remained unanswered. In addition, the problem has not been investigated from an evolutionary perspective. Our phylogenetic analysis revealed three major clades belonging to HU, IHF α and IHF β like proteins with reference to *E.coli*. We carried out a comparative analysis of three-dimensional structures of HU/IHF proteins to gain insight into the structural basis of clade division. The present study revealed three major features which contribute to differential DNA binding specificity of HU/IHF proteins, I) conformational restriction of DNA binding residues due to salt-bridge formation II) the enrichment of alanine in the DNA binding site increasing conformational space of flexible side chains in its vicinity and III) nature of DNA binding residue (Arg to Lys bias in different clades) which**

interacts differentially with DNA bases. Differences in the dimer stabilization strategies between HU and IHF were also observed. Our analysis reveals a comprehensive evolutionary picture, which rationalizes the origin of multi-specificity of HU/IHF proteins using sequence and structure-based determinants, which could also be applied to understand differences in binding specificities of other nucleic acid binding proteins.

Chapter 5 | Phylogenetic Studies and Inhibitor Design targeting Protein interacting Interface of Nucleoid-Associated Protein HU

Chapter 5 details the evolutionary and structural analyses of the protein interacting interface, at the alpha helical region, which was for MtbHU as a case study to inhibit HU and Topoisomerase I interaction,. In continuation of this part of the thesis, the candidate examined the role of the protein binding interface of HU-IHF family proteins, understanding its evolutionary history and utilizing it in designing inhibitors for *Mycobacterium tuberculosis* HU (MtbHU). **In an earlier study, by our group and with contributions from the author, the crystal structure of MtbHU was determined, and was inhibited using stilbene derivatives which inhibited the Mtb cell growth. It motivated us to understand the evolutionary and structural characteristics of the HU protein binding (HU_{pb}) interface, which has not being investigated previously for HU or for any other NAPs. Results from maximum likelihood phylogenetic analysis and ancestral reconstruction points that the residue position 16 (w.r.t. MtbHU) could be a major factor of formation of pocket at HU_{pb} interface.** The conservation and importance of other residue positions in the HU_{pb} interface are discussed in this chapter. If it is polar (as in the case of MtbHU), or charged residue, a pocket is formed as in contrast to the apolar residues which shield it. **The docking results aimed at targeting the MtbHU_{pb} interface discovered compounds like maltotetraose, valrubicin, iodixanol, enalkiren, indinavir, carfilzomib, oxytetracycline, quinalizarin, raltitrexed, epigallocatechin and their analogs exhibit high scoring binding at the interface,** which was further screened for better scoring compounds.

Our present result discussed in the chapter gives a model example of an evolutionary study of a protein interface of nucleoid-associated protein, which is used to understand the interface and computationally design inhibitors targeting it. This strategy could be in general useful for designing inhibitors all types of protein-protein interfaces, where evolutionary studies can guide to direct which interfaces are difficult (flatter interface) or which are easier (interface with pocket, which can act as an anchor for the core of inhibitory compound) to target.

Chapter 6 | *Mycobacterium tuberculosis* HU inhibitor affects crucial genes involved in Protein synthesis, Aerobic respiration and Lipid metabolism

Chapter 6 details analysis of the transcriptional profile of Mtb cells treated with MtbHU inhibitor SD1 and can provide us a glimpse to the regulatory network of MtbHU. A study in the present thesis can serve as a model example of using chemical tools (inhibitors) to unravel the regulatory network of a nucleoid-associated protein, which acts as a global regulator, controlling many transcriptional units. Understanding the gene regulation of various NAPs under different biological conditions has become important for gaining insight into the global regulatory network of the organism. **In a previous study, we inhibited HU, which is essential in *Mycobacterium tuberculosis* with Stilbene like compounds (SD1 and SD4) which disrupted its DNA binding ability.** In this present study, the author analyzed the transcription profile of response to treatment of Mtb cells by SD1 which is a MtbHU inhibitor, to identify those genes/operons which are differentially expressed. We combined our analysis with Gene Ontology (GO); Pathway-based enrichment and comparative analysis of various gene transcription profiles from tuberculosis database (TBDB), along with ChIP-qPCR validation. **A large number of genes involved in mycolic acid biosynthesis, energy metabolism, and protein translation were observed to be down-regulated, which affects various crucial metabolic and information processing pathways.** The result also shows some overlaps with HU regulons from *E. coli* and *Salmonella enterica*, while its control over genes/pathways like ribosomal proteins and mycolic acid biosynthesis pathways were found new in *M. tuberculosis*. The results presented in chapter 6 suggest HU may coordinate the expression of genes involved in virulence

which can contribute the organism's pathogenicity. **Our results thus provide a glimpse of the HU regulon of *M. tuberculosis* investigated using its inhibitor as a chemical tool, which highlights the role HU as a global regulator.**

Chapter 7 | Summary and Future directions

Chapter 7 summarizes the general question asked in the thesis, its important findings, and outlines the future directions of the work