

4-1-2022

## **An Intersection of Computational Biology and Functional Genomics to identify Transcriptional Gene Enhancers and Their Role in Cancer**

Naureen Aslam Khattak  
*Indiana State University*

Rusty Allen Gonser  
*Indiana State University*

Follow this and additional works at: <https://scholars.indianastate.edu/bakerman>

---

### **Recommended Citation**

Aslam Khattak, Naureen and Gonser, Rusty Allen, "An Intersection of Computational Biology and Functional Genomics to identify Transcriptional Gene Enhancers and Their Role in Cancer" (2022). *Bakerman Student Research Awards*. 32.  
<https://scholars.indianastate.edu/bakerman/32>

This Article is brought to you for free and open access by the Cunningham Memorial Library at Sycamore Scholars. It has been accepted for inclusion in Bakerman Student Research Awards by an authorized administrator of Sycamore Scholars. For more information, please contact [dana.swinford@indstate.edu](mailto:dana.swinford@indstate.edu).

1 **An Intersection of Computational Biology and Functional Genomics to identify**  
2 **Transcriptional Gene Enhancers and Their Role in Cancer**

3 Naureen Aslam Khattak<sup>1,2</sup>, Rusty Allen Gonser<sup>1,2</sup>

4 <sup>1</sup>Department of Biology, Indiana State University, Terre Haute, Indiana

5 <sup>2</sup>The Center of Genomic Advocacy (TCGA), Indiana State University, Terre Haute, Indiana

6 [naslam@sycamores.indstate.edu](mailto:naslam@sycamores.indstate.edu), [rusty.gonser@indstate.edu](mailto:rusty.gonser@indstate.edu)

7 **Abstract**

8 Despite the critical role of gene regulation in cell development and differentiation, the major  
9 challenge remains to identify the *cis-regulatory modules* (CRMs). Mainly, these *CRMs* include  
10 enhancers, promoters and insulators that governs the spatiotemporal gene regulation. The gene  
11 regulatory networks are highly dependent on their *CRMs* and mostly consist of DNA motifs and  
12 epigenetic landmarks. The recent advancements in high-throughput sequencing techniques and  
13 comparative genomics analysis accelerate the discovery of enhancers, however the major  
14 obstacles are to identify the genome-wide location of these CRMs, their dynamic nature of  
15 interactions, and *cis/trans* location which could be hundred to thousands base pairs away from  
16 the target gene location. The goal of this literature review is to provide an insight into the *CRMs*  
17 specifically enhancers, how they modulate gene expression, mutations that converts normal cell  
18 into a disease-state such as cancer. Also, this embedded review article is focused on the use of  
19 computational strategies coupled with the biochemical assays to predict functional gene  
20 enhancers. The computational strategies such as window clustering, probabilistic modeling,  
21 phylogenetic footprinting and discriminative modeling are briefly discussed to scan and locate  
22 the putative gene enhancers. Besides these, biochemical techniques such as ChIP-seq, DNA  
23 footprinting, and deletion mapping are briefly reviewed in *Drosophila* to predict functional gene  
24 enhancers and dissecting gene regulatory networks. In addition, this review article may help

25 bench scientists to incorporate bioinformatics tools with biochemical techniques to scan, locate  
26 and verify gene enhancer regions within a cell. With best of our knowledge, this is a first-time  
27 effort to combine insilico, in vitro and in vivo techniques to explore the connections between  
28 CRMs and gene regulation.

### 29 **Keywords**

30 Gene Enhancers, Cis regulatory modules, ChIP-seq, Deletion Mapping, Cancer, Gene  
31 Expression. Cancer, Computational Biology.

### 32 **Introduction**

33 An interesting question in developmental biology is deciphering how multiple cues are integrated  
34 to determine where and under what conditions a specific gene is expressed. Despite all cells in an  
35 organism having the same genetic makeup, only a subset of these genes is expressed in each type  
36 of cell, thereby providing each cell type with a unique identity. The difference in gene expression  
37 is derived by regulatory regions of DNA called enhancers, which bind with specific proteins  
38 called transcription factors (TFs) to regulate gene expression. The challenge is still on identifying  
39 the multiple signals, which coordinate and communicate with each other to drive the expression  
40 of a particular gene. The integration of multiple signals from different Transcription Factors  
41 (TFs) mediates the spatial and temporal gene expression in a tightly controlled environment.  
42 Thus, a specific gene expressed in a specific tissue at a specific time is highly dependent on TFs.  
43 Any mutation these TFs or their associated proteins results in an up and down regulation of gene  
44 expression which leads to multiple disorders such as different type of cancer [1]

45 The TFs and their binding sites are essential for a normal gene activity and controls the rate of  
46 gene transcription. Mainly, enhancer (*CRMs* modules) consists of single or multiple binding sites  
47 for a variety of TFs that modulate a gene expression either in a direct or indirect way. For direct  
48 transcriptional control, an activator (TFs) binds to an enhancer region without any additional

49 support to initiate gene transcription. However, the indirect transcriptional control acts through  
50 cofactors/additional TFs to regulate a gene activity [2]

### 51 **Gene Enhancers**

52 The enhancer consists of sequence-specific DNA binding sites known as Transcription Factor  
53 Binding Sites (TFBS) along with some other signals necessary to regulate a gene expression. The  
54 average size of an enhancer is a few hundred to a few thousands base pairs (bp) long. These  
55 enhancers make a loop-like structure and recruit TFs to regulate a gene expression as depicted in  
56 **figure 1**. Generally, enhancers are divided into two broad categories known as short-range  
57 (proximal) and long-range (distal) enhancers. Super-enhancers are also found in the mammalian  
58 genome where multiple enhancers are present with an array of TFs bound to these enhancers [3].

### 59 **Mechanism of Action**

60 Enhancers acts in a cooperative manner or in a stand-alone mode. Single or multiple TFs can  
61 bind to an enhancer region. If a TF binds to an enhancer and serve as a docking site for another  
62 TFs to activate or repress gene expression, it reflects cooperative or indirect mode of action,  
63 whereas the direct mode of action (stand-alone) shows a TF directly binds to the target site and  
64 act as an activator/repressor without any additional support as shown in **figure 2** [2]. In addition  
65 to direct and indirect mode of action, it also depends upon the (i) binding occupancy of each TF,  
66 and (ii) their orientations with respect to Transcription Start Site (TSS) and enhancer category  
67 (proximal vs. distal enhancer) in driving gene expression (**figure 3**)

68 **Transcription Factors (TFs):** TFs are prerequisite for an efficient gene transcription. This binds  
69 on enhancers regions and drives gene expression. Most of the transcription factors reported so far  
70 consist of DNA binding domain and activation domain. A brief detail of each TF is below

71

### 72 **Helix-loop-Helix motif (HLH)**

73 HLH TFs consist of basic amino acids that contact with DNA and neighboring regions to  
74 mediate dimer formation. Dimer is formed due to flexibility of loop which allows folding  
75 and packing against the other helix (**figure 4**). These TFs play a role in cell development and  
76 differentiation. MyoD1 is an example of HLH TF that binds to E2A protein [3]

### 77 **Helix-turn-Helix motif (HTH)**

78 The HTH binding motif consists of a pair of  $\alpha$  helices separated by a light turn. The second  $\alpha$   
79 helix  
80 lies in the major groove of DNA where it contacts with DNA bases, whereas the first  $\alpha$  helix  
81 make  
82 contact with DNA backbone as shown in **figure 5** [4]

### 83 **Zinc Finger domain**

84 Zinc finger domains are mainly responsible for inducing growth and differentiation. These TFs  
85 were first identified in the *Xenopus* model organism. The residues cysteine and histidine  
86 coordinate with zinc ions and form a zinc finger-like projection. Zinc fingers consist of an  $\alpha$   
87 helix and  $\beta$  sheet held together by zinc ions (**figure 6**). Typically, a finger motif has the  
88 following sequence: Cys-X<sub>2</sub> or <sub>4</sub>-Cys-X<sub>12</sub>-His-X<sub>3-5</sub>-His. Several proteins having zinc finger  
89 binding sites have been identified such as TF SP1[4]

### 90 **Leucine zipper motif**

91 The conserved sequence of Leucine zipper motif was first discovered in eukaryotic proteins and  
92 has a critical function in cell differentiation and development. Leucine zipper contains 4-5  
93 leucine located seven residues apart found in the basic amino acid region. These two regions  
94 spread over 60-80 residues, together constitute bZIP domain. The basic region is held together by  
95 dimerization of adjacent zipper regions, when hydrophobic faces of two zippers interact in

96 parallel orientation, the leucine zipper part stabilizes the protein dimer as shown in **figure 7**  
 97 [4]. Different types of TFs are summarized in **table 1**.

98 **Table 1: Transcription Factors Domains and their Function**

TFs Domain	Role	Function containing domain	Gene ID (NCBI)	OMIM	Chr location
Homeobox	DNA binding	Numerous <i>Drosophila</i> homeotic genes related genes in other organisms such as <i>Cad, Abd-A/B</i>	CG1759 CG10325 CG11648	600297 142951 142956	2L 3R 3R
Cysteine-histidine zinc finger	DNA binding	<i>TFIIIA, Kruppel, Spl</i>	2971 9314	600860 602253	13 9
Cysteine-cysteine zinc finger	DNA binding	Steroid-thyroid hormone receptor family	*NA	*NA	*NA
Leucine Zipper	Protein Dimerization	<i>C/EBP, c-fos, c-jun, GCN4, c-myc</i>	2353 3725	164810 165160	14 1
Helix-loop-Helix	Protein Dimerization	<i>c-myc, Drosophila daughterless, MyoD, E12, E47</i>	4609 4654	190080 159970	8 11
Proline-rich region	Gene Activation	<i>Yeast GCN4, GAL4, steroid-thyroid</i>	V XVI	856709 855828	--
Glutamine-rich region	Gene Activation	<i>SP1</i>	6667	189906	12
Amphipathic acidic alpha-helix	Gene Activation	<i>CTF/NF1</i>	4763	613113	17

99 \*NA = Not available

## 100 **Enhancers and Their Associated Diseases**

101 Based on the critical importance of enhancers in gene regulation (activation or inhibition of gene  
 102 expression), it is not surprising that any change in the enhancer itself or its associated factors can

103 result in disease. For a detailed review, please see the Jaret M. Karnuta<sup>1</sup> and Peter C. Scacheri,  
104 2018 [7].

### 105 **Enhancers and their Role in Cancer**

106 Disease can appear if a TF is up or down regulated (ectopic activation of gene expression) or  
107 becomes active at the wrong time or in the wrong place. The mis-regulation of gene expression  
108 plays a major role in the development of certain types of cancers. For instance, over-expression  
109 of proto-oncogenes results in 'cancer-causing' oncogenes, particularly the gene responsible for  
110 cell growth. A few examples are the growth factors and their receptors (*erbA*, *fos*, *myb*, and *myc*)  
111 that encode TFs necessary for growth turned into oncogene if there is any mutation associated  
112 with these factors or their genes. Thus, the conversion of these protooncogenes into oncogenes,  
113 which can occur either by mutation (over-expression or under-expression) corresponds to a  
114 difference in gene regulation pattern which results in cancer [8,9]. Table 2 briefly describes TFs,  
115 mutations, and cancer types. A detailed consists of mutation in enhancers regions and their effect  
116 or organism type is provided in the supplementary data files (**Table 3**). The mutations include  
117 insertion/deletion, translocation, inversion, duplications, and point mutations with phenotype  
118 defect, NCBI gene ID, chromosome location and OMIM record are provided below.

119 **Table 2: Different type of cancers and their association with enhancer malfunctioning**

<b>Cancer types</b>	<b>Mutations</b>	<b>Gene ID (NCBI)</b>	<b>Chr</b>	<b>MIM record</b>	<b>Ref.</b>
Breast, prostate	<i>FOXA1</i>	3169	14	602294	10
Lung, AML	<i>RAD21</i>	5885	8	606462	11
Burkitt's lymphoma	<i>IGH/MYC</i> translocation	4609	8	190080	12

Breast, lung	<i>GATA3</i>	2625	10	131320	13
Transitional cell carcinoma	<i>NIPBL</i>	25836	5	608667	14
Urothelial, bladder, breast, head, and neck	<i>CTCF</i>	10664	16	604167	15
B-cell lymphoma, lung	<i>EZH2</i>	2146	7	601573	16
Bladder, glioblastoma, lung, urothelial	<i>STAG2</i>	10735	X	300826	17
Bladder, AML, lung	<i>SMC3</i>	9126	10	606062	18
Bladder, lung, urothelial, and breast	<i>MLL2/MLL3/ MLL4</i>	8085	12	602113	19

120 \*Chr: Chromosome

121 \*OMIM: Online Mendelian Inheritance in Man (database for human disease)

122 An example is *fos* and *jun* TFs. These are normally synthesized transiently in response to growth  
123 promoting signals and act to activate the genes encoding specific proteins required for cellular  
124 growth. If for any reason these proteins are continuously expressed (over-expression), they act to  
125 promote the continuous growth in the absence of growth factors and are capable of transforming  
126 normal cells into cancer cells. In contrast to above examples, in some cases the TFs failed to  
127 regulate correctly (at correct time and place) results in an inappropriate gene activity. This  
128 indicates that, as with other cellular processes, gene expression is subject to complex regulatory  
129 mechanisms, the failure of which can be as devastating as the failure of the basic process

### 130 **Computational Strategies for Enhancer Prediction**

131 Computational search algorithms are widely used to identify the enhancers regions or the  
132 hotspots areas where an enhancer or their associated TFs may be existing. These search



133 algorithms either used experimental data to get a fine-tuned matching for *CRMs* or based on  
134 mathematical or statistical models to get a reasonable prediction. The most accurate strategy for  
135 predicting *CRMs* varies and depends upon the question of interest. For instance, among all  
136 different types of prediction tools, which predictive tool stands high for identifying the *CRMs*?  
137 Are they time- efficient? Are there any false positive and negative results? Did they use  
138 experimental data or based on putative annotation? What is the limitation of the algorithms in  
139 terms of predicting short vs. long enhancers or in other words what is the input size limitation of  
140 the tool? For details, please refer to Su, J., Teichmann, S. A., & Down, T. A. [23]. The *insilico*  
141 methods used for enhancers predictions are shown in **figure 8** are roughly classified into four  
142 categories briefly described below.

- 143 1. **Window clustering** involves significant clustering of high densities of binding sites  
144 within a sequence window.
- 145 2. **Probabilistic modelling** consists of identifying sequences that resemble a statistical  
146 model of a binding site cluster more than a model of background DNA.
- 147 3. **Phylogenetic footprinting** searches for high density regions of binding sites conserved  
148 between closely related species.
- 149 4. **Discriminative modelling** seeks to identify set of signals on regulatory regions that can  
150 maximize the differences between regulatory regions and non-regulatory regions. Many  
151 methods are hybrids of two or more strategies.

152 **Window Clustering:** The literal meaning of the word “clustering” is to “group together” based  
153 on similar properties. The same approach is utilized to predict the *CRMs* in the genome based on  
154 statistical analyses. The method uses high density TFBS and groups them together based on  
155 statistical observation. These significant clusters are then the hotspot for finding the TFs. The (i)

156 MSCAN [24] (ii) MCAST [25] and (iii) CisPlusFinder [26] tools are based on clustering  
157 methods. The input data for these tools consist of motif library against single genome and  
158 multiple sequences alignment respectively as shown in **figure 9**.

159 **Probabilistic Modelling:** This computational approach utilizes the Hidden Markov Model  
160 (HMM) to generate a set of *CRMs* sequences which are based on a combination of a set of  
161 TFBS. Common tools include ClusterBuster [27], Stubb, StubbMS [26], MorphMS [28],  
162 CisModule [28], and MultiModule [29]. The difference between these tools is mentioned below

- 163 1) The ClusterBuster, Stubb, and CisModule stands on Multiple Sequences  
164 Alignment (MSA)
- 165 2) StubbMS and Morphs MS are diverse in their first step of execution. StubbMS  
166 provides fixed alignment by using Lagan [23]. On the other hand, MorphMS sums  
167 up (using probability) all the possible alignments based on their binding sites.
- 168 3) CisModule and MultiModule predict *CRMs* in a single step. The CisModule  
169 follows Bayesian inference to find the binding sites and location of the *CRMs*,  
170 however, the Multimodule uses the same strategy but adds the comparative  
171 genomic information to complete the analyses.

172 Some of the tools

173 present in window clustering (CisPlusFinder) and Probabilistic Modelling (StubbMS, MorphMS  
174 and MultiModule) also follow the Phylogenetic fingerprinting methods based on Multiple  
175 sequence alignment approach [23]

176 **Discriminative Modelling:** In this method for instance HexDiff [26], the input consists of a  
 177 hexamer or set of nucleotides (6-mer) with high frequency and differentiate between *CRMs* and  
 178 non-*CRMs*.

179 **Phylogenetic Footprinting:** This method uses phastCons score [27] which serves as  
 180 independent control and takes sequence conservation as an input. The phastCons score considers  
 181 the evolutionary distance between species which are followed by the Hidden *Markov Model*.  
 182 Different bioinformatics tools are now available which provide information about regulatory  
 183 elements as well as TFs binding sites and target genes of regulatory elements. A brief overview  
 184 of different bioinformatics software and databases are available in **table 4**. A brief list of  
 185 databases used for identifying Transcription Factors (TFs) and their binding sites are provided in  
 186 **table 5**.

187 **Table 4: Bioinformatics software and databases for predicting *cis-regulatory modules* in**  
 188 **genome**

Tools	Principle	Input	Website
MSCAN	Window Clustering	Single Genome	<a href="http://www.cisreg.ca/cgi-bin/mscan/">http://www.cisreg.ca/cgi-bin/mscan/</a> MSCAN
MCAST			<a href="http://alternate.meme-suite.org">http://alternate.meme-suite.org</a>
CisPlusFinder			<a href="http://jakob.genetik.uni-koeln.de/bioinformatik/people/nora/nora.html">http://jakob.genetik.uni-koeln.de/bioinformatik/people/nora/nora.html</a>
ClusterBuster	Phylogenetic footprinting	MSA	<a href="http://zlab.bu.edu/cluster-buster/">http://zlab.bu.edu/cluster-buster/</a>
Stubb			<a href="http://stubb.rockefeller.edu/">http://stubb.rockefeller.edu/</a>
StubbMS			<a href="http://stubb.rockefeller.edu/">http://stubb.rockefeller.edu/</a>
MorphMS			<a href="http://veda.cs.uiuc.edu/Morphalign/">http://veda.cs.uiuc.edu/Morphalign/</a>

			<a href="#">supplement/</a>
CisModule		Motif Library	<a href="http://www.stat.ucla.edu/~zhou/">http://www.stat.ucla.edu/~zhou/</a> <a href="#">CisModule/</a>
MultiModule			<a href="http://www.stat.ucla.edu/~zhou/">http://www.stat.ucla.edu/~zhou/</a> <a href="#">MultiModule/index.html</a>
EEL	Window Clustering		<a href="http://www.cs.helsinki.fi/u/kpalin/EEL/">http://www.cs.helsinki.fi/u/kpalin/EEL/</a>
RP	Discriminative Modeling	CRM annotations	<a href="http://www.bx.psu.edu/projects/rp/">http://www.bx.psu.edu/projects/rp/</a>
HexDiff	Discriminative Modeling		<a href="http://www.ics.uci.edu/~bobc/hexdiff.html">http://www.ics.uci.edu/~bobc/hexdiff.html</a>
PhylCRM	Phylogenetic Footprinting	Single Genome	<a href="http://the_brain.bwh.harvard.edu/PhylCRM/">http://the_brain.bwh.harvard.edu/PhylCRM/</a>
EMMA	Phylogenetic Footprinting/ Probabilistic Modeling	Motif Library	<a href="https://www.bioinformatics.nl/cgi-bin/emboss/emma">https://www.bioinformatics.nl/cgi-bin/emboss/emma</a>

189 \*\*For more details on these methods, please refer to Su, J., Teichmann, S. A., & Down, T. A.  
190 (2010).

191 With the availability of the microarray expression data analysis, several tools are published to  
192 predict CRMs in tissue and stage specific manner for example LRA [28], Cluster Scan [29]  
193 Composite Module Analyst [30], Module Miner [31]. The methods which predict TFBS based on  
194 user-defined dataset includes but not limited to Module Scanner [64], TargetExplorer [32], and  
195 CisModScan [33]. These computer-based tools are unable to find out novel CRM instead of that  
196 they look for the binding sites within a defined sequence.

**Table 5: Databases for identifying the Transcription Factors (TFs) and their binding sites**

Database	Acronym	Principle	Model Organism	Website
Transcription factor prediction database	DBD	Database of predicted TFs in completely sequenced genomes. Superfamily, Pfams and Hidden Markov model libraries-based prediction	<i>B. Subtilis</i> <i>C. Elegans</i> , <i>D. melanogaster</i> <i>E. coli</i> , <i>H. Sapiens</i> , <i>M. Musculus</i> <i>S. cerevisiae</i> .	<a href="http://www.transcriptionfactor.org">www.transcriptionfactor.org</a>
Transcription factor 2 DNA	TF2DNA	Predict TFs binding motifs using experimental and theoretical data source	<i>E. coli</i> , <i>C. Elegans</i> , <i>D. melanogaster</i> , <i>M. Musculus</i> , <i>S. cerevisiae</i> , <i>H. Sapiens</i> ,	<a href="http://www.fiserlab.org/tf2dna_db/">http://www.fiserlab.org/tf2dna_db/</a>
JASPAR	JASPAR CORE	A curated, non-redundant set of profiles, derived from published and experimentally defined TF binding sites for eukaryotes uses position weight matrices (PWM)	Eukaryotes <i>Vertebrata</i> , <i>Nematoda</i> , <i>Insect</i> , <i>Plantae</i> , <i>Fungi</i> , <i>Urochordata</i>	<a href="http://jaspar.genereg.net">http://jaspar.genereg.net</a>
Gene Transcription Regulation Database	GTRD	Database of TFs identified by ChIP-seq experiments for human and mouse.	<i>M. Musculus</i> , <i>H. Sapiens</i> .	<a href="http://gtrd.biouml.org">http://gtrd.biouml.org</a>
TRANSFAC	TRANSFAC	a database on TFs and their DNA binding sites	Eukaryotic transcription factors	<a href="http://genexplain.com/transfac/#section0">http://genexplain.com/transfac/#section0</a>
AnimalTFDB	AnimalTFDB	Annotations from the NCBI Entrez Gene and Ensembl	Animal Genomes	<a href="http://bioinfo.life.hust.edu">http://bioinfo.life.hust.edu</a>

		databases, including basic information, gene phenotypes, homologous genes, and Gene Ontology (GO) Classification of transcription cofactors; (iii) TF binding sites information; (iv) the GWAS phenotype related information of human TFs.		<a href="http://u.cn/AnimalTFDB/">u.cn/AnimalTFDB/</a>
--	--	--	--	--

198

### 199 **Biochemical Techniques for Enhancer Prediction using Transgenic Animal Models**

200 As mentioned earlier, the TFs are specific proteins that binds to the enhancer region and regulate  
 201 a gene activity. The DNA footprinting, Deletion mapping and Chromatin Immunoprecipitation  
 202 (ChIP-seq) are frequently used techniques in laboratory to identify the location and binding  
 203 occupancy of TFs on the enhancer regions in genome as shown in **figure 12**. Although these  
 204 techniques stand on the same basic principle of identifying TFBS, steps used in these methods  
 205 are entirely different and are comparable. The overall principle and steps used in these strategies  
 206 are given below:

#### 207 **Deletion Mapping: TFBS Identification Techniques**

208 This technique takes the advantage of deleting various parts the promoter region of a gene and  
 209 measured the transcription activity. In this method, mutant gene (depending upon the deletion  
 210 results in either (i) increase in the transcription activity (ii) decrease in transcription activity or  
 211 (iii) in some cases no /little effect on the gene regulation. In general, a few nucleotide deletions  
 212 have little or no effect on the gene transcription, however if the deletion hit the regions which is  
 213 important for the binding of the TFs, then it will show subsequent decrease in the transcriptional  
 214 activity. On the other hand, if the deletion occurred in the region which is responsible for the

215 repression of the gene activity, then it might result in increased gene transcription because the  
216 repressor binding site is no longer available to inhibit the gene expression as shown in figure 13.

217 The different scenarios are listed below

218 i. If the deletion occur/falls in the regions which is important for TF binding (prevents  
219 binding of TFs) to activate transcription the level will decrease, the transcriptions  
220 (38%)

221 ii. Deletion of regions also might increase transcription in a case that TF binds to a  
222 region that inhibit transcriptions, but deletion of that inhibiting region will increase  
223 transcription (114%)

224 iii. Deletion of other regions might have no or little effect [35]

### 225 **DNA footprinting: TFBS Identification Techniques**

226 DNA footprinting techniques benefit from the action of nuclease enzymes. The nuclease is the  
227 class of enzymes that degrade DNA by breaking down the DNA phosphodiester bond. Based on  
228 this fact, if the DNA sequence is treated with nuclease such as DNase I, the free DNA will be  
229 digested easily whereas the DNA that is bound to a protein, also known as protected DNA, will  
230 remain intact. After digestion process, the bound protein is removed and the DNA sequences  
231 specific for the TFs are identified (**figure 14**) [35]

### 232 **Chromatin Immunoprecipitation by Sequencing (ChIP-seq)**

233 Chromatin Immunoprecipitation is most widely used technique to identify the genome-wide  
234 location of TFs. The overall steps involve in this technique are (i) Isolation of cultured cells or a  
235 particular tissue (ii) Cross-linking of DNA (iii) DNA sonication (iv) Immunoprecipitation using  
236 antibody and (v) identification of DNA-protein bound sequences either by sequencing (ChIP-  
237 seq) or DNA hybridization (ChIP-chip). The critical steps are the cross-linking of DNA within  
238 the cell. Mostly, formaldehyde is used to cross-link the TFs to the DNA sites at which they are

239 bound in the living cells. After that, DNA is sonicated into different fragments. The DNA  
240 fragmentation is entirely depending upon the experiments needs, sequencing system, and the  
241 basic purpose of doing the ChIP-seq. For instance, on average 200-300 base pair (bp) is  
242 preferable size for DNA fragmentation which will then proceed for immunoprecipitation.  
243 However, in some cases, it might be around 400-500 bp.

244 The rule of thumb is to have a decent DNA fragment neither too short, nor too long. In both  
245 cases, it will affect the results. The shorter DNA fragment can lose the TFBS, and longer DNA  
246 fragment might give false positive results. A proactive approach is to check the DNA fragments  
247 on gel to get a clear idea of the size. For the immunoprecipitation, controls (positive and  
248 negative) are crucial to get a successful result along with at least 2-3 biological replicates. After  
249 that, based on the selective techniques i.e., ChIP-seq or ChIP-chip, further steps are carried out.  
250 The former technique is followed by sequencing method (next generation sequencing) and later  
251 apply DNA hybridization concept (DNA labeling with fluorophore or radiolabeled isotopes) to  
252 identify the DNA sequences (figure 16) The sequencing data will be analyzed to find the peaks.  
253 The peaks are the area with probability of TFs bound on the genome. For that, raw sequencing  
254 reads are mapped on the reference genome (FASTQ file format; raw sequence data) and then the  
255 reads are aligned to the genome by using alignment Softwares such as Bowtie2. The next step is  
256 the peak detection (MACS software) and visualization via computational tools such as Integrated  
257 Genome Viewer (IGV) visualization. The selection of alignment software and visualization is  
258 entirely depended upon the user (**figure 17**) [35]

## 259 **Conclusion and Future Directions**

260 Recent advances accelerate the discovery of the *CRMs*, but many questions remain unanswered.  
261 The understanding of the molecular mechanism that governs direct and indirect interactions of  
262 TFs to dissect the genetic role of TFs is still under investigation. Although computational



263 methods are widely employed to predict potential candidates to save time and resources, the use  
264 of experimentally verified data is still under way to get maximum confidence. Apart from that,  
265 the bench-side scientist is paying more attention to produce high-quality, unbiased, reproducible  
266 datasets which can increase the success of training datasets used for *in silico* CRMs prediction.  
267 However, to gain a more complete picture of the role of these factors both in normal cellular  
268 function and in disease processes, their function to promote or inhibit gene transcription, the  
269 consequences of mutations in these CRMs regions, and more specifically the functionality of  
270 distal enhancers need to be examined.

271 The future directions may be emphasized on predicting CRMs by using a combinatorial  
272 methodology of conserved intra and inter-specific motifs (evolutionary signatures), epigenetic  
273 marks, landscape of histone modifications, novel and new experimental techniques, more  
274 specific and sensitive bioinformatics tools to detangle the e mystery behind these CRMs and their  
275 function in shaping gene expression level.

## 276 **Abbreviation**

277 **ChIP-seq:** Chromatin Immunoprecipitation sequencing

278 **CRMs:** Cis-Regulatory Modules

279 **TFs:** Transcription Factors

280 **TFBs:** Transcription Factors Binding Sites

## 281 **Declarations**

## 282 **Ethics approval and consent to participate**

283 Not applicable

## 284 **Consent for publication**

285 All authors are agreed to publish this review article

## 286 **Availability of data and materials**

287 Not applicable: Data sharing is not applicable to this article as no datasets were generated or  
288 analyzed during the current study.

### 289 **Competing interests**

290 The authors have no conflict of interest to declare.

### 291 **Funding**

292 We would like to thank The Center of Genomic Advocacy at Indiana State University and  
293 Sigma-Xi grant to support this project.

### 294 **Authors' contributions**

295 NAK and RAG conceived the ideas. RAG designed methodology. NAK led the writing of the  
296 manuscript. All authors contributed critically to the drafts and gave final approval for  
297 publication.

### 298 **Acknowledgements**

299 We greatly acknowledge The Center of Genomic Advocacy at Indiana State University to  
300 provide necessary support to complete this manuscript.

### 301 **References**

- 302 1) Visel A, Rubin EM, Pennacchio LA. Genomic views of distant-acting  
303 enhancers. *Nature*. 2009; 461:199–205.
- 304 2) Pennacchio LA, Bickmore W, Dean A, Nobrega MA, Bejerano G. Enhancers: five  
305 essential questions. *Nat Rev Genet*. 2013;14(4):288–295. doi:10.1038/nrg3458
- 306 3) David S. Latchman. Transcription factors: an overview. *Current Status Review. Int. J.*  
307 *Exp. Path.* (1993) 74, 417-422
- 308 4) Pabo, C. O., & Sauer, R. T. (1992). Transcription factors: structural families and  
309 principles of DNA recognition. *Annual review of biochemistry*, 61, 1053–1095.  
310 <https://doi.org/10.1146/annurev.bi.61.070192.005201>

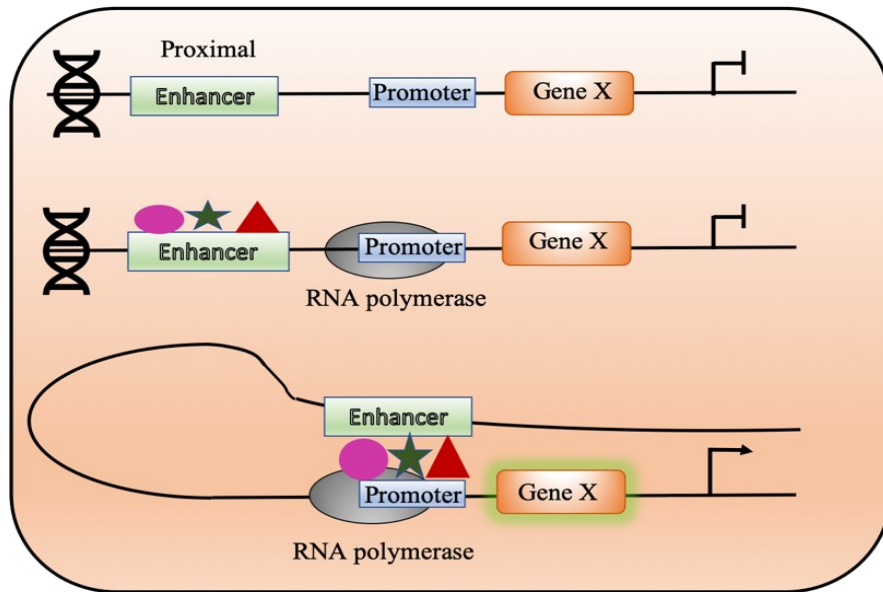
- 311 5) Radovick, S., Nations, M., Du, Y., Berg, L. A., Weintraub, B. D., & Wondisford, F. E.  
312 (1992). A mutation in the POU-homeodomain of Pit-1 responsible for combined pituitary  
313 hormone deficiency. *Science (New York, N.Y.)*, 257(5073), 1115–1118.  
314 <https://doi.org/10.1126/science.257.5073.1115>
- 315 6) Crossley, M., & Brownlee, G. G. (1990). Disruption of a C/EBP binding site in the factor  
316 IX promoter is associated with haemophilia B. *Nature*, 345(6274), 444–446.  
317 <https://doi.org/10.1038/345444a0>
- 318 7) Karnuta, J. M., & Scacheri, P. C. (2018). Enhancers: bridging the gap between gene  
319 control and human disease. *Human molecular genetics*, 27(R2), R219–R227.  
320 <https://doi.org/10.1093/hmg/ddy167>
- 321 8) Radovick, S., Nations, M., Du, Y., Berg, L. A., Weintraub, B. D., & Wondisford, F. E.  
322 (1992). A mutation in the POU-homeodomain of Pit-1 responsible for combined pituitary  
323 hormone deficiency. *Science (New York, N.Y.)*, 257(5073), 1115–1118.  
324 <https://doi.org/10.1126/science.257.5073.1115>
- 325 9) Crossley, M., & Brownlee, G. G. (1990). Disruption of a C/EBP binding site in the factor  
326 IX promoter is associated with haemophilia B. *Nature*, 345(6274), 444–446.  
327 <https://doi.org/10.1038/345444a0>
- 328 10) Karnuta, J. M., & Scacheri, P. C. (2018). Enhancers: bridging the gap between gene  
329 control and human disease. *Human molecular genetics*, 27(R2), R219–R227.  
330 <https://doi.org/10.1093/hmg/ddy167>
- 331 11) Cox, P. M., & Goding, C. R. (1991). Transcription and cancer. *British journal of*  
332 *cancer*, 63(5), 651–662. <https://doi.org/10.1038/bjc.1991.151>
- 333 12) Kandoth, C., McLellan, M. D., Vandin, F., Ye, K., Niu, B., Lu, C., Xie, M., Zhang, Q.,  
334 McMichael, J. F., Wyczalkowski, M. A., Leiserson, M., Miller, C. A., Welch, J. S.,

- 335 Walter, M. J., Wendl, M. C., Ley, T. J., Wilson, R. K., Raphael, B. J., & Ding, L. (2013).  
336 Mutational landscape and significance across 12 major cancer types. *Nature*, *502*(7471),  
337 333–339. <https://doi.org/10.1038/nature12634>
- 338 13) Adams, E. J., Karthaus, W. R., Hoover, E., Liu, D., Gruet, A., Zhang, Z., Cho, H.,  
339 DiLoreto, R., Chhangawala, S., Liu, Y., Watson, P. A., Davicioni, E., Sboner, A.,  
340 Barbieri, C. E., Bose, R., Leslie, C. S., & Sawyers, C. L. (2019). FOXA1 mutations alter  
341 pioneering activity, differentiation and prostate cancer phenotypes. *Nature*, *571*(7765),  
342 408–412. <https://doi.org/10.1038/s41586-019-1318-9>
- 343 14) Cheng H, Zhang N, Pati D. Cohesin subunit RAD21: From biology to disease. *Gene*.  
344 2020 Oct 20; 758:144966. doi: 10.1016/j.gene.2020.144966. Epub 2020 Jul 17. PMID:  
345 32687945.
- 346 15) Nishijima, A., Noguchi, Y., Narukawa, K., Takano, H., & Oshikawa, G. (2019). [*Rinsho*  
347 *ketsueki*] *The Japanese journal of clinical hematology*, *60*(10), 1425–1430.  
348 <https://doi.org/10.11406/rinketsu.60.142>
- 349 16) Takaku, M., Grimm, S. A., & Wade, P. A. (2015). GATA3 in Breast Cancer: Tumor  
350 Suppressor or Oncogene?. *Gene expression*, *16*(4), 163–168.  
351 <https://doi.org/10.3727/105221615X14399878166113>
- 352 17) Gao, D., Zhu, B., Cao, X., Zhang, M., & Wang, X. (2019). Roles of NIPBL in  
353 maintenance of genome stability. *Seminars in cell & developmental biology*, *90*, 181–  
354 186. <https://doi.org/10.1016/j.semcdb.2018.08.005>
- 355 18) Debaugny, R. E., & Skok, J. A. (2020). CTCF and CTCFL in cancer. *Current opinion in*  
356 *genetics & development*, *61*, 44–52. <https://doi.org/10.1016/j.gde.2020.02.021>

- 357 19) Lue, J. K., & Amengual, J. E. (2018). Emerging EZH2 Inhibitors and Their Application  
358 in Lymphoma. *Current hematologic malignancy reports*, 13(5), 369–382. [https://doi.org/  
359 10.1007/s11899-018-0466-6](https://doi.org/10.1007/s11899-018-0466-6)
- 360 20) Aquila, L., Ohm, J., & Woloszynska-Read, A. (2018). The role of STAG2 in bladder  
361 cancer. *Pharmacological research*, 131, 143–149.  
362 <https://doi.org/10.1016/j.phrs.2018.02.025>
- 363 21) Hill, V. K., Kim, J. S., & Waldman, T. (2016). Cohesin mutations in human  
364 cancer. *Biochimica et biophysica acta*, 1866(1), 1–11.  
365 <https://doi.org/10.1016/j.bbcan.2016.05.002>
- 366 22) Dhar, S. S., Zhao, D., Lin, T., Gu, B., Pal, K., Wu, S. J., Alam, H., Lv, J., Yun, K.,  
367 Gopalakrishnan, V., Flores, E. R., Northcott, P. A., Rajaram, V., Li, W., Shilatifard, A.,  
368 Sillitoe, R. V., Chen, K., & Lee, M. G. (2018). MLL4 Is Required to Maintain Broad  
369 H3K4me3 Peaks and Super-Enhancers at Tumor Suppressor Genes. *Molecular  
370 cell*, 70(5), 825–841.e6. <https://doi.org/10.1016/j.molcel.2018.04.028>
- 371 23) Su, J., Teichmann, S. A., & Down, T. A. (2010). Assessing computational methods of  
372 cis-regulatory module prediction. *PLoS computational biology*, 6(12), e1001020.  
373 <https://doi.org/10.1371/journal.pcbi.1001020>
- 374 24) Alkema, W. B., Johansson, O., Lagergren, J., & Wasserman, W. W. (2004). MSCAN:  
375 identification of functional clusters of transcription factor binding sites. *Nucleic acids  
376 research*, 32(Web Server issue), W195–W198. <https://doi.org/10.1093/nar/gkh387>
- 377 25) Grant, C. E., Johnson, J., Bailey, T. L., & Noble, W. S. (2016). MCAST: scanning for  
378 cis-regulatory motif clusters. *Bioinformatics (Oxford, England)*, 32(8), 1217–1219.  
379 <https://doi.org/10.1093/bioinformatics/btv750>

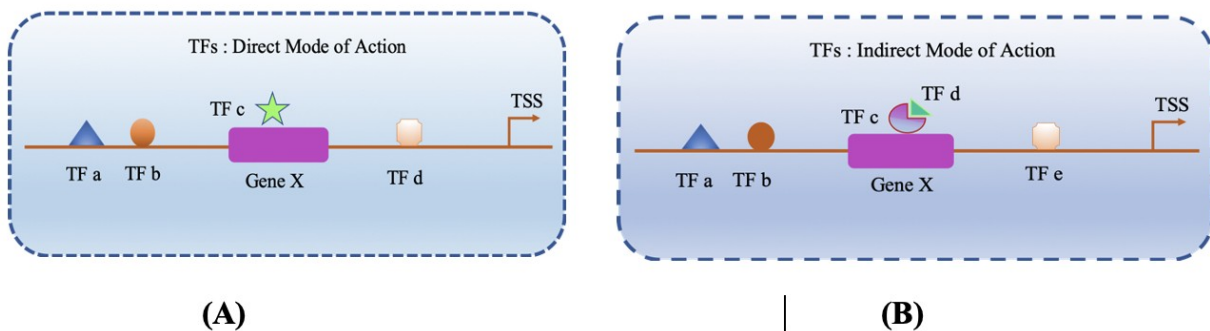
- 380 26) Sinha, S., van Nimwegen, E., & Siggia, E. D. (2003). A probabilistic method to detect  
381 regulatory modules. *Bioinformatics (Oxford, England)*, *19 Suppl 1*, i292–i301.  
382 <https://doi.org/10.1093/bioinformatics/btg1040>
- 383 27) Frith, M. C., Li, M. C., & Weng, Z. (2003). Cluster-Buster: Finding dense clusters of  
384 motifs in DNA sequences. *Nucleic acids research*, *31(13)*, 3666–3668.  
385 <https://doi.org/10.1093/nar/gkg540>
- 386 28) Zhou, Q., & Wong, W. H. (2004). CisModule: de novo discovery of cis-regulatory  
387 modules by hierarchical mixture modeling. *Proceedings of the National Academy of*  
388 *Sciences of the United States of America*, *101(33)*, 12114–12119. [https://doi.org/10.1073/](https://doi.org/10.1073/pnas.0402858101)  
389 [pnas.0402858101](https://doi.org/10.1073/pnas.0402858101)
- 390 29) Zhou, Qing; Wong, Wing Hung (2007). Coupling hidden Markov models for the  
391 discovery of Cis -regulatory modules in multiple species. *Ann. Appl. Stat.* 1, no. 1, 36--  
392 65. doi:10.1214/07-AOAS103. <https://projecteuclid.org/euclid.aoas/1183143728>
- 393 30) Palin, K., Taipale, J. & Ukkonen, E. (2006). Locating potential enhancer elements by  
394 comparative genomics using the EEL software. *Nat Protoc* 1, 368–374  
395 <https://doi.org/10.1038/nprot.2006.56>
- 396 31) King, D. C., Taylor, J., Elnitski, L., Chiaromonte, F., Miller, W., & Hardison, R. C.  
397 (2005). Evaluation of regulatory potential and conservation scores for detecting cis-  
398 regulatory modules in aligned mammalian genome sequences. *Genome research*, *15(8)*,  
399 1051–1060. <https://doi.org/10.1101/gr.3642605>
- 400 32) Chan, B.Y., Kibler, D. (2005). Using hexamers to predict cis-regulatory motifs in  
401 *Drosophila*. *BMC Bioinformatics* 6, 262. <https://doi.org/10.1186/1471-2105-6-262>

- 402 33) Warner, J. B., Philippakis, A. A., Jaeger, S. A., He, F. S., Lin, J., & Bulyk, M. L. (2008).  
403 Systematic identification of mammalian regulatory motifs' target genes and  
404 functions. *Nature methods*, 5(4), 347–353. <https://doi.org/10.1038/nmeth.1188>
- 405 34) Alberts B, Johnson A, Lewis J, et al. *Molecular Biology of the Cell*. 4th edition. New  
406 York: Garland Science; 2002. Available from:  
407 <https://www.ncbi.nlm.nih.gov/books/NBK21054/>
- 408 35) Haseloff J, Siemering KR (2005) The uses of green fluorescent protein in plants. In:  
409 Green Fluorescent Protein. Wiley, pp 259–284. doi:[10.1002/0471739499.ch12](https://doi.org/10.1002/0471739499.ch12)
- 410 36) Jefferson RA, Kavanagh TA, Bevan MW. GUS fusions: beta-glucuronidase as a  
411 sensitive and versatile gene fusion marker in higher plants. *EMBO J*. 1987;6(13):3901–  
412 3907.
- 413 37) Goldsbrough AP, Tong Y, Yoder JI, Tong Y (1996) *Lc* as a non-destructive visual  
414 reporter and transposition excision marker gene for tomato. *Plant J* 9:927–933
- 415 38) Bischof J, Maeda RK, Hediger M, Karch F, Basler K. An optimized transgenesis system  
416 for *Drosophila* using germ-line-specific phiC31 integrases. *Proc Natl Acad Sci U S A*.  
417 2007;104(9):3312–3317. doi:[10.1073/pnas.0611511104](https://doi.org/10.1073/pnas.0611511104)
- 418 39) A.H. Brand, N. Perrimon. Targeted gene expression as a means of altering cell fates and  
419 generating dominant phenotypes. *Development* 1993, 118: 401-415
- 420 40) Koen J. T. Venken, Hugo J. Bellen. Transgenesis upgrades for *Drosophila melanogaster*  
421 *Development* 2007 134: 3571-3584; doi: [10.1242/dev.00568](https://doi.org/10.1242/dev.00568)



422  
423

424 **Figure 1:** The 3C loop-formation of enhancer to activate/repress a gene. In first case, the  
425 proximal enhancer is present near to the gene promoter site but there is no TFs attach to the  
426 enhancer. In second case, three different TFs (circle, star, and triangle) binds to the enhancer  
427 region but unable to activate the gene expression. In this scenario, the TFs are most likely the  
428 repressor proteins which inhibit the gene expression. In third case, the TFs binds and results in a  
429 loop-like formation which enables enhancer region to come into the proximity of the promoter  
430 and RNA polymerase to activate the gene expression.

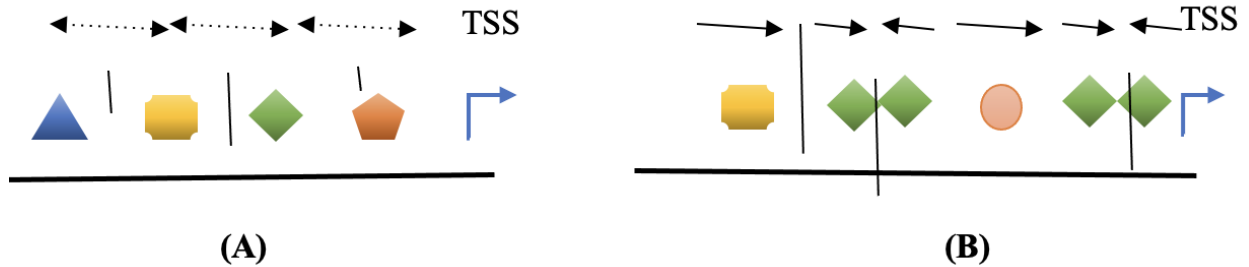


431  
432

433 **Figure 2:** A simple schematic of the direct and indirect mode of action of a TF. (A) The direct  
434 mode requires only a single TF to bind and regulate gene expression. (B) On the other hand, the

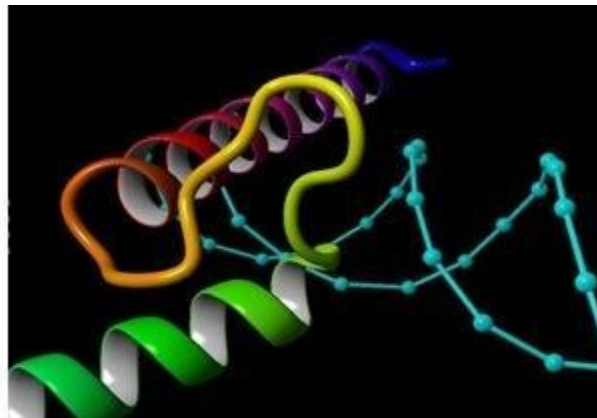


434 indirect mode may have multiple TFs or a single TF with multiple binding sites with other co-  
435 factors or TFs to regulate gene expression.

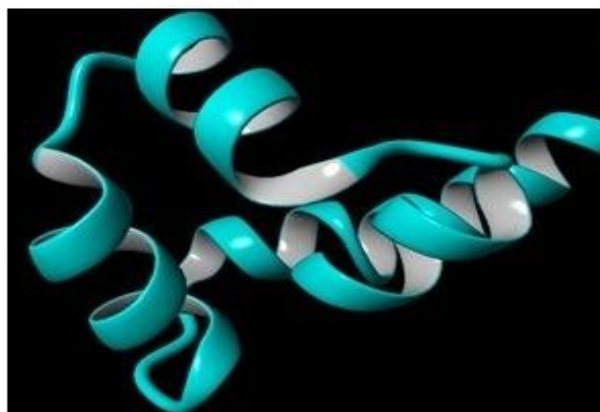


436  
437 **Figure 3:** (A) shows the distance between each TFs whereas (B) focuses on the orientation of  
438 TFs. The vertical bar represents the binding affinity of each of the TF. The smaller vertical bar  
439 shows less binding affinity and vice versa.

440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468

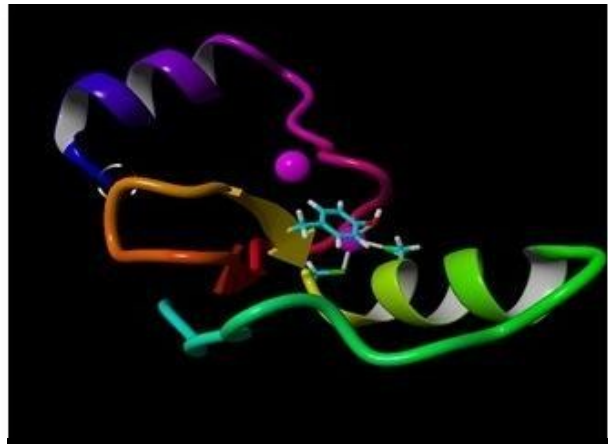


**Figure 4:** Helix-loop-helix

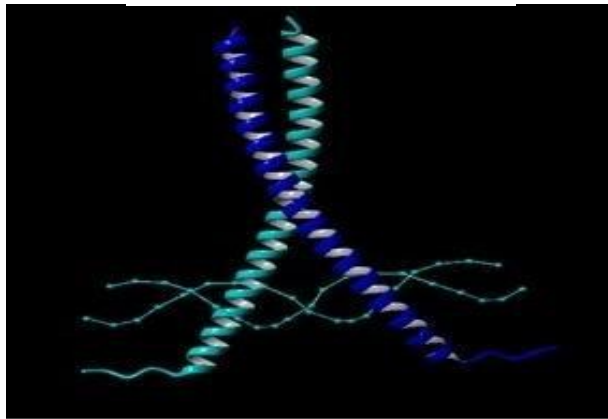


**Figure 5:** Helix turn Helix

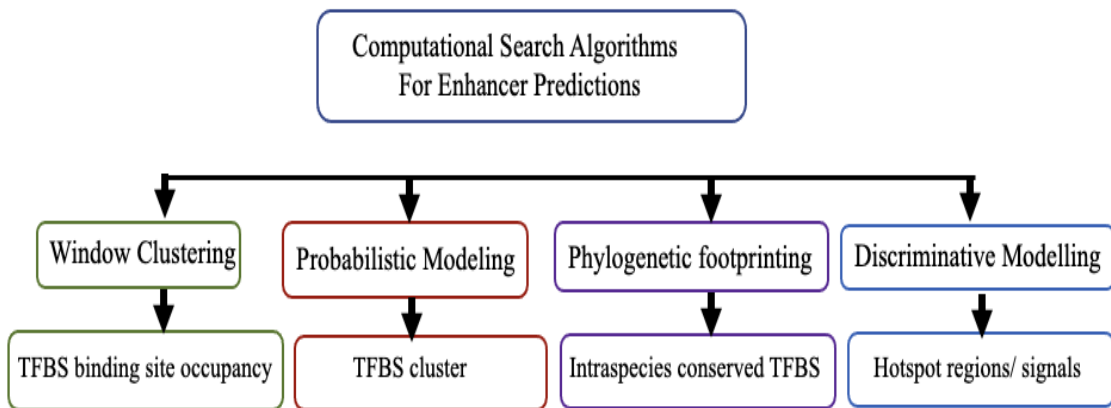
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512



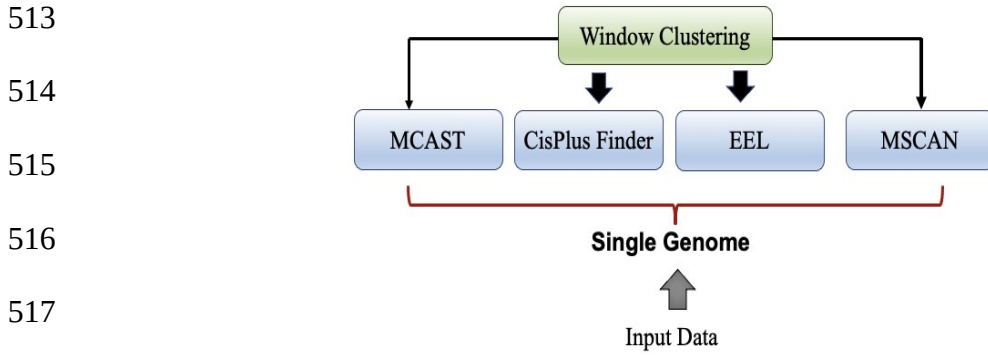
**Figure 6:** Zinc finger



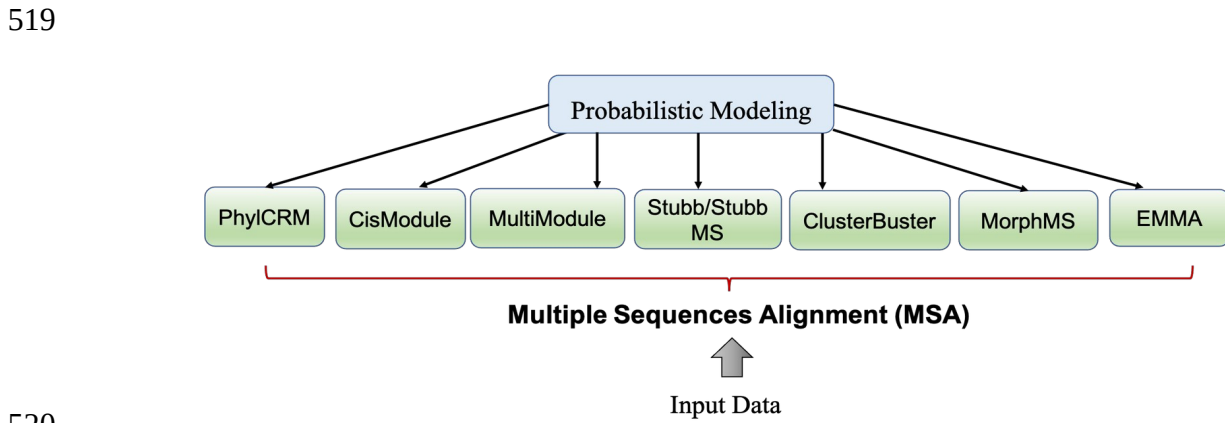
**Figure 7:** Leucine zipper motif



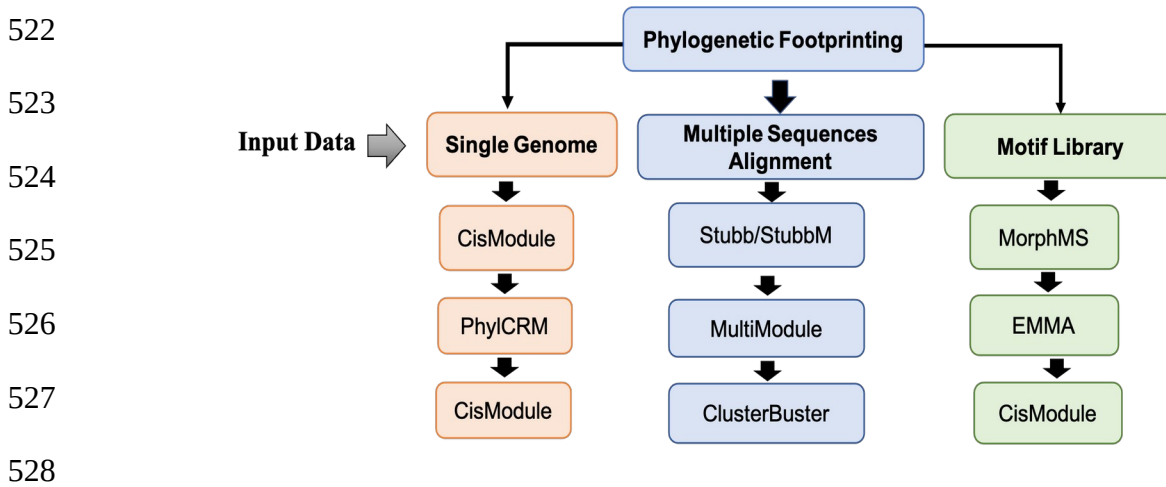
**Figure 8:** The common strategies used in computational search algorithm for identifying the *CRMs* regions in genome.



518 **Figure 9:** Window Clustering for Enhancer



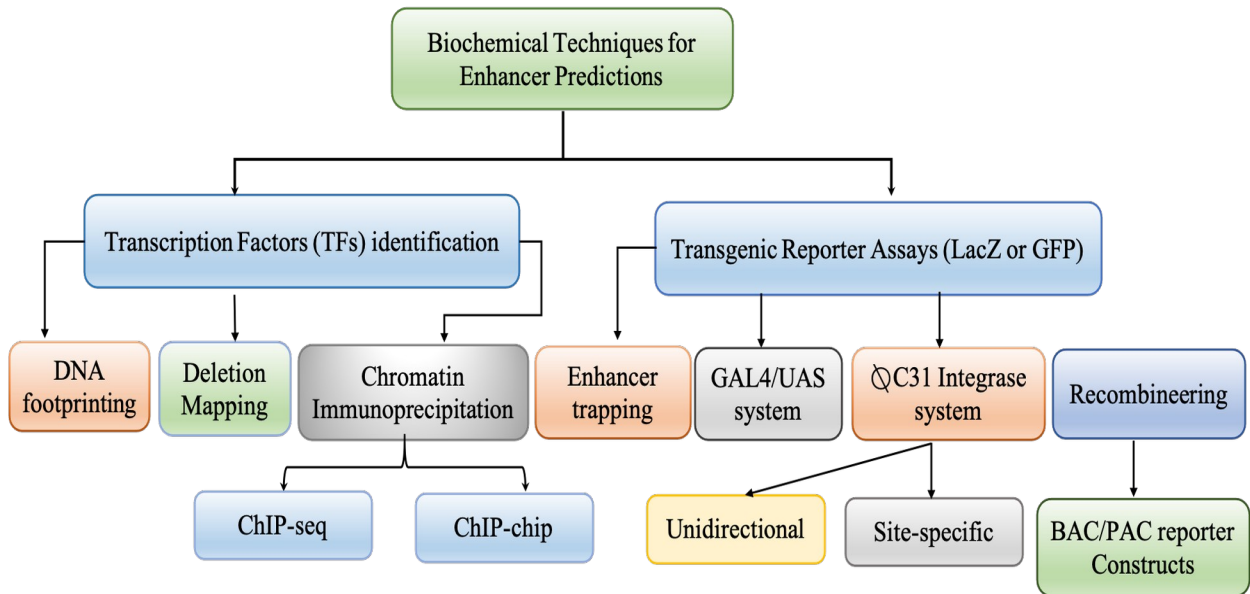
520 **Figure 10:** Different Softwares under the category of Probabilistic Modeling approach



529 **Figure 11:** Phylogenetic footprinting approach with the input data type and Software for

530 enhancer predictions

531

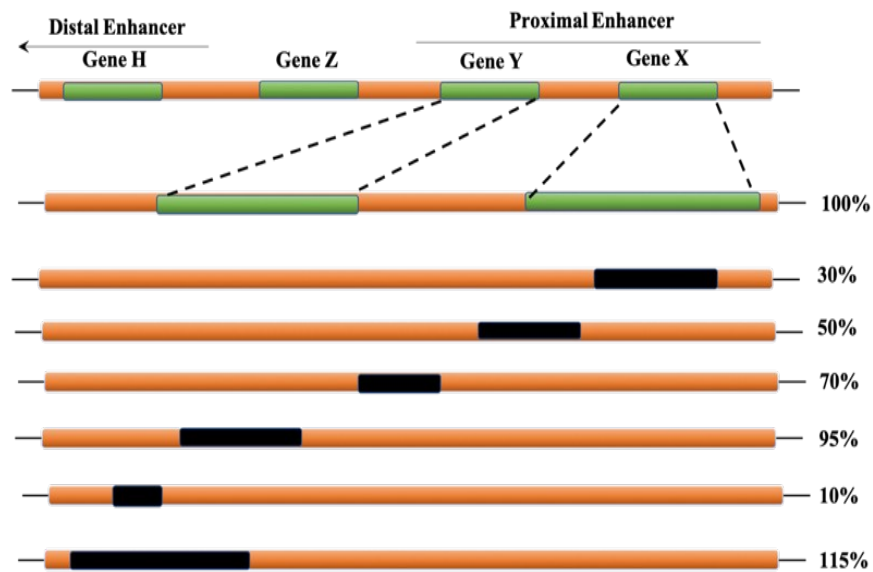


532

533 **Figure 12:** Biochemical methods to identify the Transcription Factor Binding sites and Enhancer

534 regions in genome

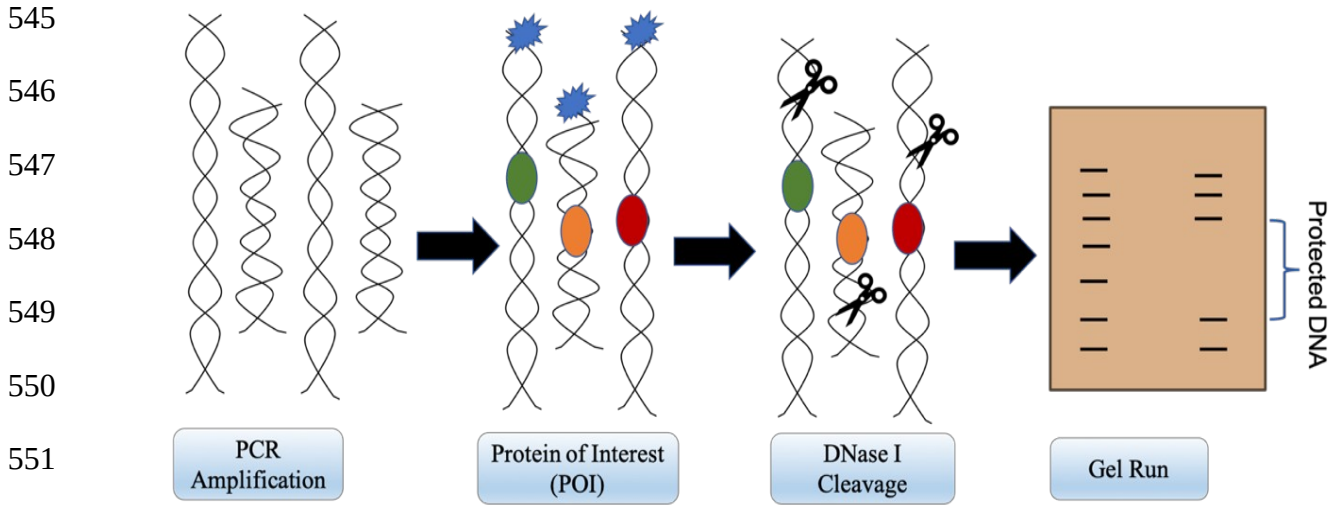
535



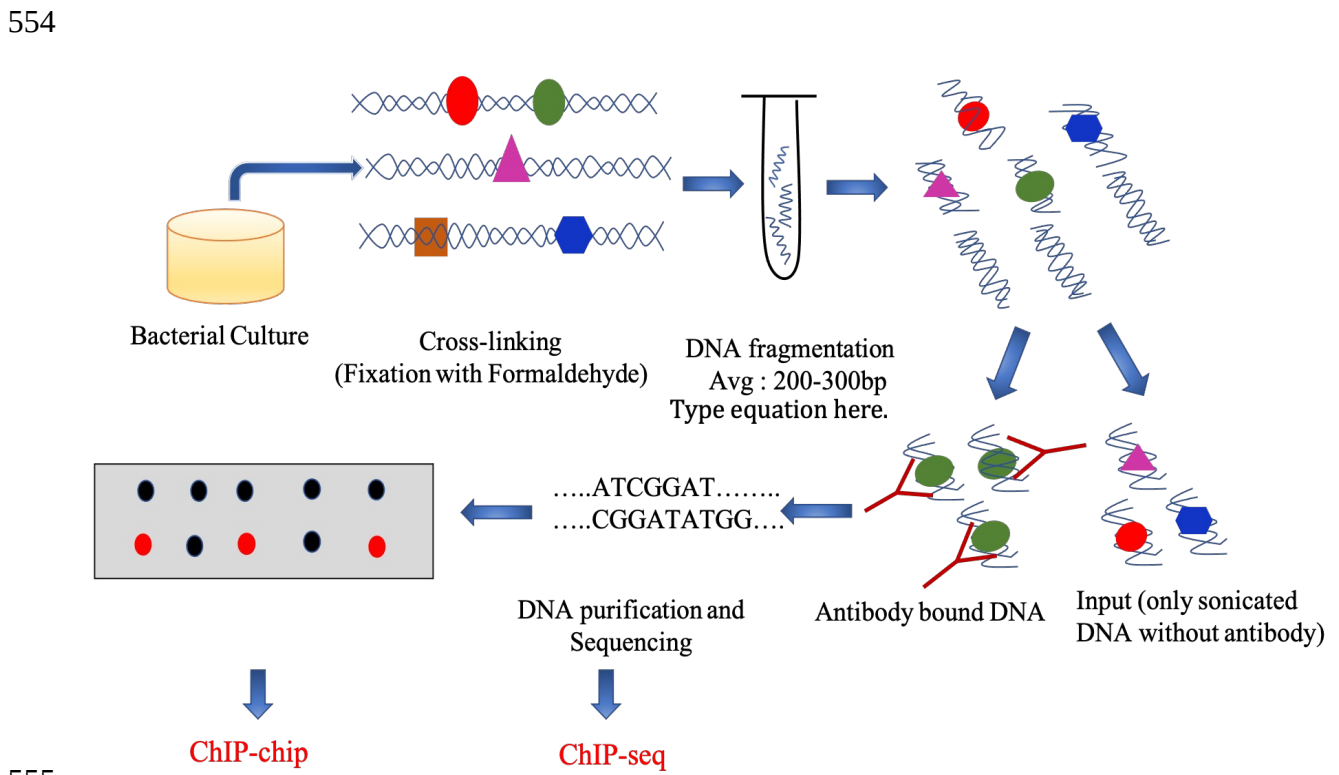
542

543 **Figure 13:** The deletion mapping technique apply to remove regions in gene of interest to see the

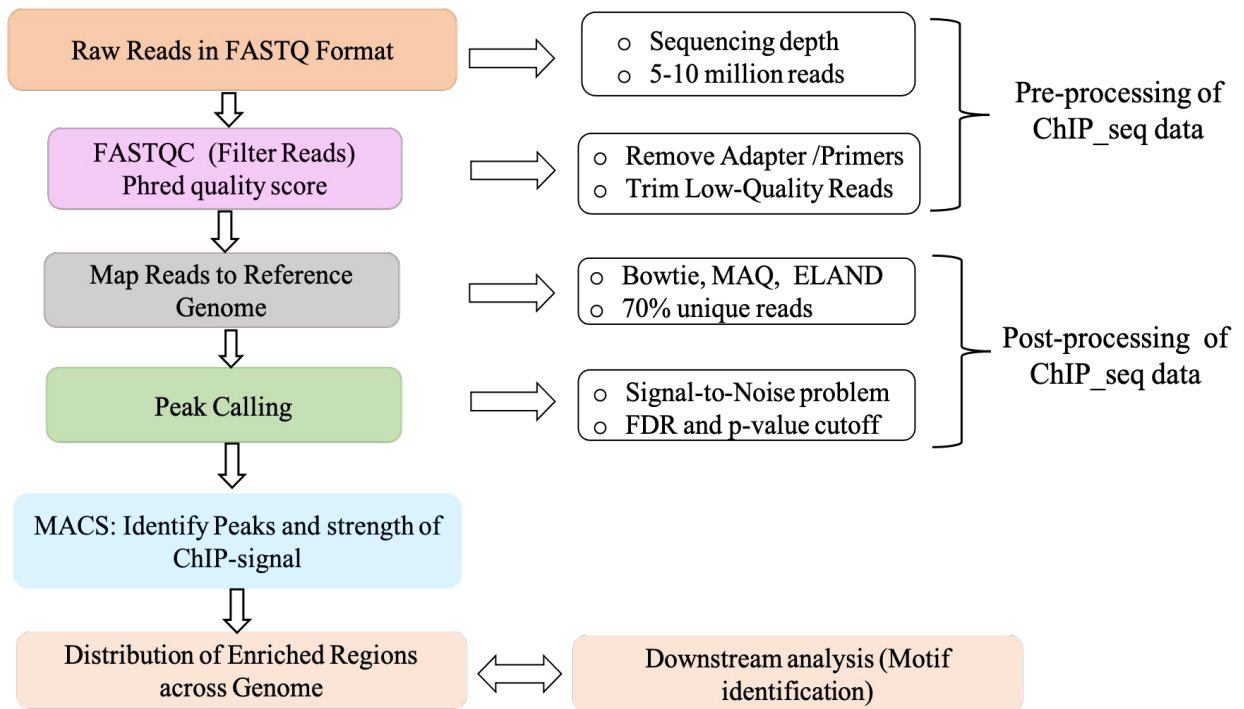
544 effect of these deleted regions on gene expression



553 **Figure 14:** The DNA footprinting strategy to identify DNS-Protein binding sites.



556 **Figure 16:** The ChIP -sequence procedure for identifying the TFs binding sites on enhancer.



557

558

**Figure 17:** The post ChIP-seq data analysis pipeline

559

560

561

562

563

564

565

566

567

568

569

570

571

573 Table 3: Mutations in the Enhancer Regions and their Effect on Organism Phenotype.

Mutations	Phenotype defect	Enhancer defect	Gene ID (NCBI)	Chr	MIM record
<b>Insertion/ Deletion</b>	X-linked deafness type 3 (DFN3)	Multiple deletions 900 kb from <i>POU3F4</i>	5456	X	300039
	Split-hand-split food malformation (SHFM)	7q21.3 deletion affecting enhancer sequences within <i>DYNCH3</i>	1780	7	603772
	Autosomal dominant adult-onset demyelinating leukodystrophy	Deletion eliminating <i>TAD</i> , allowing for enhancer adoption of <i>LMNB1</i>	395342 4001	24 and 5	150340
	Preaxial polydactyly	13 bp insertion in the zone of polarizing activity regulatory sequence (ZRS) affecting sonic hedgehog (SHH) expression.	64327	7	605522
	Aniridia Involves 11p13,	Involves 11p13, downstream of <i>PAX6</i>	5080	11	607108
	Pierre Robin sequence	1 Mb away from <i>SOX9</i> . Abrogates binding of <i>MSX1</i> in vitro studies.	6662	17	608160
<b>Translocation</b>	Split-hand syndrome	t (2; 7) (p25.1; q22), separates limb enhancers in <i>DYNCH3</i> from <i>DLX5/6</i> .	1780	7	603772
	<b>Inversion</b>	Limb syndactyly	Enhancer adoption by SHH induced by a 7q inversion.	6469	7
Hand-foot-genital syndrome		Hand-foot-genital syndrome Chromosome 7 inversion causing a HOXA13 enhancer	3209	7	142959

		delocalization.			
<b>Duplication</b>	Disorders of sex development (DSD)	16p13.3 duplication of <i>GNG13</i> and <i>SOX8</i> enhancers. 600 kb upstream of <i>SOX9</i>	51764 30812 6662	13, 8 and 9	607298 605923 608160
	Keratolytic Winter Erythema	Duplication of enhancer upstream of <i>CTSB</i> .	1508	8	116810
	Haas-type polysyndactyly and Laurin-Sandrow syndrome	Microduplications in SHH limb enhancer ZRS	6469 64327	7	600725 605522
<b>Point Mutations</b>	Holoprosencephaly	460 kb upstream of SHH resulting in loss of SHH brain enhancer-2 activity.	6469	7	600725
	Van der Woude syndrome	Mutation in <i>IRF6</i> enhancer, abrogating p63 and E47 binding.	3664	1	607199
	Preaxial polydactyly	Various point mutations in the ZRS enhancer (e.g., 295 T>C)	64327	7	605522
	Hirschsprung disease	Common non-coding variant within an enhancer like sequence in <i>RET</i> intron 1.	5979	10	164761