# Developments in computational cis-regulatory module prediction

**Jeffrey A. Thompson** [1]**, Clare Bates Congdon** [1,*]

[1]*Department of Computer Science, University of Southern Maine, Portland, ME, USA*

Correspondence*:
Clare Bates Congdon
Department of Computer Science, University of Southern Maine, 96 Falmouth Street, Portland, ME, 04104-9300, USA, congdon@usm.maine.edu

## ABSTRACT

Dozens of approaches have been created to predict cis-regulatory modules (CRMs). Two separate reviews, in 2009 and 2010 compared many of these tools and identified the most promising approaches, as well as the need for improvement. Since those reviews were published, a number of new methods have been developed that offer significant advantages over what came before. Additionally, breakthroughs in high throughput biological experimental techniques have occurred that may mitigate or complement the need for computational prediction. In this work we examine the developments in CRM prediction since that time and look at the capabilities of the latest generation of tools and techniques.

Keywords: **cis-regulatory module, CRM, gene regulation, epigenetic, motif, computational prediction, gene regulatory network, GRN**

## 1 INTRODUCTION

Multicellular organisms exhibit complex patterns of differential gene expression to enable the varied roles of their various tissues and organs. Still other patterns of expression are critical to the initial development of the organism and the organization of its body plan. Disruption of these patterns can occur through mutation and exposure to toxins, leading to cancer and other diseases. Therefore, much research has been devoted to understanding how these key differential expression patterns are maintained.

A cell has many opportunities to regulate a gene, but a key element is the set of cis-regulatory modules (CRMs) that control the gene's transcription. A cis-regulatory module (CRM) is a region of noncoding DNA that groups together specific transcription factor binding sites (TFBSs) that together affect gene expression (see Box 1). This combinatory organization allows a relatively small number of transcription factors to participate in the complex differential expression patterns of a larger number of genes. Genes may be regulated by multiple CRMs that control their expression in various contexts, making knowledge of CRMs an important part of understanding gene regulatory networks and disease.

In the past, biological elucidation of regulatory function was a slow and arduous process. Because of this, computational tools were developed to focus biological experiments on high confidence areas for investigation. Recent developments in high-throughput techniques have made volumes of experimental data available, such as DNase-Seq, FAIRE-Seq, and ChIP-Seq, indicating various regulatory pathways, such as chromatin accessibility, histone modification, and transcription factor binding (see Box 1). Some have suggested that such data may supplant the need for computational inference **Hardison and Taylor**

30  (2012). Nevertheless, these experimental data cannot yet tell the whole story. They are limited to a
31  particular cell type, at a particular time, under particular conditions. Publically available data have been
32  generated for a limited number of genomes and cell types and are not yet available for most regulatory
33  proteins, while the cost of generating such comprehensive epigenetic data will likely remain out of reach
34  for many labs for some time to come. Additionally, these high-throughput assays are not capable of
35  conclusively establishing regulatory function. Therefore, computational tools remain an important part of
36  learning about CRMs, perhaps in conjunction with this new high-throughput data.

37     In this review we will briefly discuss the development of the computational prediction of CRMs, as
38  well as biological methods. The most recent review of CRM prediction techniques was in 2012 **Hardison**
39  **and Taylor** (2012). However, the focus of that review was on biological methods of predicting regulatory
40  regions and its discussion of computational methods was focused primarily of general differences between
41  approaches, rather than specific tools that are publically available to researchers. Furthermore, there
42  have been a number of important developments in prediction methods that were either published after
43  that review, or were not covered in it. The emphasis of this review will be on developments in CRM
44  prediction that are available to the public. Earlier reviews provide a comprehensive look at older methods
45  **Elnitski et al.** (2006), **Van Loo and Marynen** (2009), **Su et al.** (2010). Although useful, one issue often
46  overlooked in earlier reviews is that the various methods are often designed with a particular use in mind,
47  but the experimental design used to compare methods tends to favor those that solve just one aspect of the
48  problem. Our focus will be to provide information on the problems particular methods were designed to
49  solve to better enable researchers to pick the tool appropriate for the work in hand.

---

**Box 1 — Pathways of Pre-Transcriptional Regulation**

What follows is a description of some of the major pathways of pre-transcriptional regulation. Transcription is regulated in many ways, and this list cannot be considered comprehensive.

*Transcription factor binding:* transcription factors are proteins that bind to DNA and regulate the transcription of a gene. The regulation can happen through a variety of means, such as recruiting poylmerase II, signaling that transcription should begin, or interfering in some way with other transcription factors.

*Histone modification:* DNA is organized into nucleosomes, which consist of DNA wrapped around histone protein complexes. Tails of the histone proteins extend from the nucleosome and these tails can be modified by the addition of members of acetyl or methyl chemical groups (see Fig. 1). These modifications are still under investigation but have been observed to play a role in gene regulation.

*Chromatin accessibility:* Chromatin is the DNA in combination with the proteins it is bound to. It has a highly organized state in order control the vast length of DNA managed in the small nucleus of the cell. Active regions of chromatin often need to be less compact then inactive regions, in order to be more accessible to ligands (such as transcription factors). Therefore, chromatin accessibility can indicate regions more likely to function in gene regulation.

*DNA Methylation:* DNA methylation involves the attchment of molecules from the methyl functional group to DNA nucleotides. This leads to the compaction of the genome at methylated regions and the repression of gene transcription. Methylation also plays other regulatory roles in the genome, such as the repression of transposons **Espada and Esteller** (2010), **Robertson and Wolffe** (2000).

50

## 2   CIS-REGULATORY MODULES

51  A CRM is a region of DNA that can be bound by a limited set of transcription factors under specific
52  conditions to regulate the expression of a gene. The concept of CRMs is important, in part because

---

transcription factor binding sites are short (approximately 6–20bp) and degenerate (a transcription factor can bind to a variety of sites). Because of this, a transcription factor could theoretically bind in many locations across the genome. It is in combination with other factors and under the correct condition that transcription factors bind to CRMs and affect gene expression. By identifying CRMs, the specificity of predicted transcription factor binding sites can be increased and the combinatorial control of a gene better understood. Therefore, identification of CRMs is critical to revealing the regulatory network that determines when a gene is active.

## 2.1  TYPES OF CRMS

There are a number of different types of CRMs, which are classified by their regulatory effect:

- *Promoters*. Eukaryotic protein-coding gene promoters are a region of noncoding DNA that enables RNA polymerase II to initiate transcription of the gene. The region is usually located close to the transcription start site (TSS) **Taher et al.** (2013). The promoter is sometimes divided into two areas: the core promoter and the proximal promoter **Taher et al.** (2013), **Butler and Kadonaga** (2002). Both are bound by transcription factors, but the core promoter is the minimal region required for transcription. Core promoter studies have shown that most genes have multiple promoters driving alternative transcription **Sandelin et al.** (2007).

- *Enhancers*. Enhancers are regions of noncoding DNA that control the conditions under which a gene is transcribed. They target a promoter and increase the probability of transcription **Walters et al.** (1995). They can be located distally from the gene and some function regardless of orientation **Blackwood and Kadonaga** (1998), **Laimins et al.** (1984). However, in other cases orientation matters **Hozumi et al.** (2013). Enhancers are structurally similar to proximal promoters **Maston et al.** (2006).

- *Silencers*. Silencers also control the conditions under which a gene is transcribed but, contrary to enhancers, they act to repress the transcription of a gene **Ogbourne and Antalis** (1998). Many silencers are similar to enhancers in their relative independence of location and orientation, but others are not. There are a number of different types of silencers **Ogbourne and Antalis** (1998), some of which interfere with transcription factors, the chromatin structure, or the formation of the transcription initiation complex with polymerase II. Also, there is evidence that an enhancer may switch to a silencing function **Jing et al.** (2008).

- *Insulators*. Insulators act to stop regulatory effects from one region affecting another region. This might involve blocking an enhancer, so that it does not affect a particular promoter, or blocking the spread of heterochromatin (compact, silenced chromatin) **Ghirlando et al.** (2012).

## 2.2  CHARACTERISTICS OF CRMS

Numerous studies have revealed properties of CRMs that differentiate them from surrounding sequence. These characteristics are exploited by CRM prediction methods to identify regions of DNA sequence associated with regulatory function.

- *Clustering*. CRMs are composed of clusters of transcription factor binding sites. Although individual sites can occur at random, this is significantly less likely to happen for clusters of binding sites **Wagner** (1999). Therefore, clusters of transcription factor binding sites may indicate a CRM. Clustering can be divided into two basic types:

  - *Homotypical clusters*. The clusters are composed of a number of repeats of the same basic binding site. For a while this was thought to possibly be a general CRM feature, but it turns out to be true of only a subtype of CRMs **Li et al.** (2007).
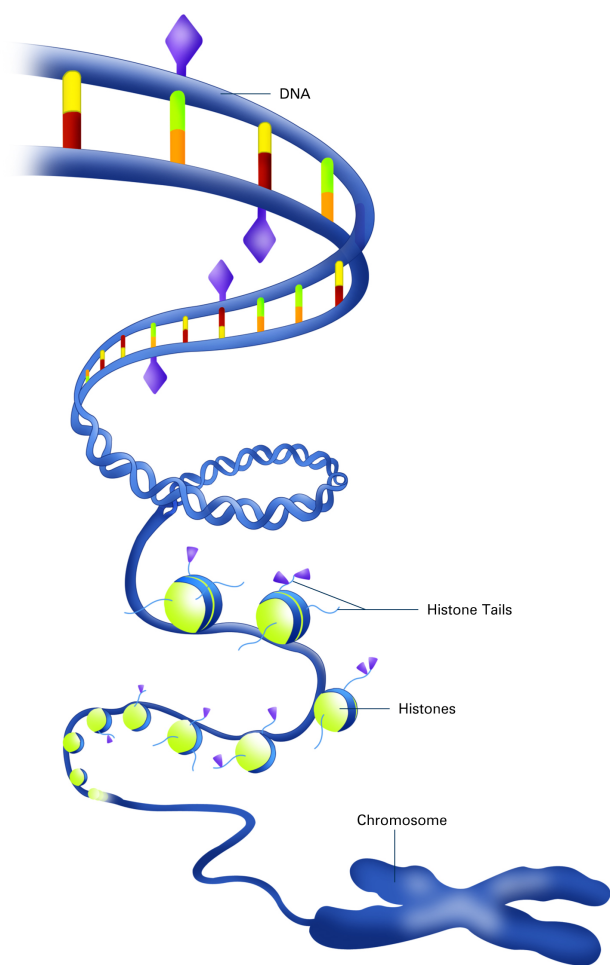
- *Heterotypical clusters*. The clusters contain a binding sites for a number of different transcription factors **Roth et al.** (2007).

- *Conservation*. CRMs have been shown to be more highly conserved than flanking regions **Kwon et al.** (2011), **Li et al.** (2007). Additionally, TFBSs within CRMs can exhibit an even higher level of constraint **Kwon et al.** (2011). Because these regions are functionally important, purifying selection has the effect of conserving them. This is not to say that all CRMs will be conserved across diverse species. It is likely that many lineage specific CRMs exist, possibly contributing to species diversity. The frequency of lineage specific CRMs is not known.

- *GC-content*. Some CRMs have been shown to have elevated GC-content compared to other noncoding sequence **Li et al.** (2007), **Saxonov et al.** (2006).

- *Motif Synergy*. CRMs are composed of specific groups of transcription factor binding sites. Transcription factors bound to these sites exhibit synergistic effects (the degree of regulation is greater than the contribution of a particular factor). These transcription factors may need to be certain distances from or orientations to each other. Therefore, the relationship of specific binding sites can be used in recognizing CRMs.

- *Size*. The largest database of biologically validated CRMs is RedFly **Gallo et al.** (2011), which contains hundreds of validated functional CRMs in *Drosophila*, which have been reduced to minimal size required to regulate the gene. An examination of these CRMs showed the length of CRMs to be variable (most were between 100 and 2100 bp in length with an average of 760bp) **Li et al.** (2007). The size characteristics of CRMs in other genomes may be different.

- *Epigenetic Properties*. A number of biochemical marks are associated with active CRMs. These include chromatin accessibility, certain histone modifications, and methylation patterns **Calo and Wysocka** (2013), **Bannister and Kouzarides** (2011), **Consortium et al.** (2011), **Cheng et al.** (2011). Of course, these markers of CRMs are not at the sequence level, where most past research has been focused. There have recently been concerted efforts to map all known epigenetic marks to a few genomes.

The characteristics of CRMs listed above are those that have generally been found to be useful in in distinguishing CRMs from background sequences. There may be others.

## 3   COMPUTATIONAL PREDICTION OF CRMS

As early as 1985, a computational tool had been written to detect eukaryotic gene promoters **Claverie and Sauvaget** (1985). This early tool searched for manually defined patterns of sequences thought specific to heatshock and glucocorticoid regulated promoters. Development of additional techniques for promoter identification followed, such as one that analyzed the density of known transcription factor binding patterns in promoters vs. non-promoters and built a promoter recognition profile **Prestridge** (1995) that could be used to recognize similar promoters and another **Chen et al.** (1997), which analyzed the density of 5-10bp strings in known promoters in order to build a profile.

Following those tools, a number of approaches were developed to identify clusters of position weight matrices (PWMs). PWMs capture the information from multiple binding sites for a transcription factor in an attempt to describe the variety of sites a transcription factor might bind to (see Table 1). A PWM shows how often a given base was seen at a particular position in the binding sites for a transcription factor. PWMs were first used to search for transcription factor binding sites in 1996 **Fickett** (1996). However, they were soon used to search for CRMs **Wasserman and Fickett** (1998), often by searching for clusters of particular PWMs that might regulate under a specific context. Since that time, various methods of searching for clusters of PWMs has been the predominant method of predicting CRMs computationally **Elnitski et al.** (2006), **Van Loo and Marynen** (2009), **Su et al.** (2010).

**Figure 1.** DNA is organized into nucleosomes (regions of DNA wrapped around histone proteins). Modifications to the tails of these histones, such as the addition of members of the methyl or acetyl groups, can affect the regulation of genes (shown as triangles). Certain marks are known to be associated with cis-regulatory modules. Methyl groups might also be attached directly to the DNA, repressing transcription of a gene (shown as diamonds). Image courtesy of NIGMS, provided by Crabtree & Company.

**Table 1.** Position Weight Matrix for Arnt: The log-likelihood of the occurence of each base at each position in the binding site is shown.

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|------|------|------|------|------|------|
| A | -0.31 | 1.88 | -4.70 | -4.70 | -4.70 | -4.70 |
| C | 1.64 | -4.70 | 1.96 | -4.70 | -4.70 | -4.70 |
| G | -4.70 | -2.12 | -4.70 | 1.96 | -4.70 | 1.96 |
| T | -4.70 | -4.70 | -4.70 | -4.70 | 1.96 | -4.70 |

138    Even the first method to use PWMs anticipated the need for more information and recommended using
139 conservation analysis when the data were available **Wasserman and Fickett** (1998). As noted in Section
140 2, many CRMs are conserved by purifying selection. Therefore, it is possible to use conservation to help
141 identify such CRMs. For a while, it was difficult to find sufficient orthologous data to use, and most of

142  the early tools did not use conservation (e.g. **Sinha et al.** (2003), **Frith** (2003)). Nevertheless, it was soon
143  incorporated as a filter in a number of tools (e.g. **Sinha et al.** (2004), **Sinha and He** (2007)).

144      Conservation can also be used as the main feature in CRM detection (rather than as an additional filter).
145  Methods of this sort were also developed around the same time **Kolbe et al.** (2004). Neither clustering
146  nor conservation can provide information about the context of when and where the gene a CRM regulates
147  is active, unless the binding sites predicted are for transcription factors known to be active in a particular
148  context. They simply identify regions that are more likely to regulate the gene and possibly suggest the
149  other genes involved in that regulation.

150      By 2010, several dozen different methods had been developed for CRM prediction. One independent
151  assessment using benchmark data **Su et al.** (2010) found that the Regulatory Potential method **Kolbe**
152  **et al.** (2004) performed the best on the human genome, although this method does not identify binding
153  sites within the CRM. The method compares two-way alignments of human and mouse DNA and classifies
154  alignments based on whether or not short fixed sequences in the alignments are more typical of regulatory
155  or neutral regions based on statistical models. For *Drosophila*, the most successful method was MorphMS
156  **Sinha and He** (2007), which combines clustering of PWMs and conservation to make predictions, with
157  the conservation measured as conserved binding sites, rather than the more typical alignments. Despite
158  the application of increasingly sophisticated techniques, no one method had emerged that could reliably
159  predict functional CRMs from DNA sequence data alone **Su et al.** (2010). Most methods that depend
160  primarily on detecting clusters of binding sites using PWMs rely on prior knowledge of processes that the
161  CRMs being searched for might be involved in so that suitable PWMs can be selected. This limits them to
162  detecting CRMs controlling the gene's known pathways. Furthermore, CRMs do not always exhibit tight
163  clustering of binding sites. Methods that are capable of selecting the PWMs that might regulate a gene
164  from a larger library had not yet seen much success **Van Loo and Marynen** (2009). Neither had *de novo*
165  approaches that did not rely on PWMs **Su et al.** (2010).

# 4  HIGH-THROUGHPUT BIOLOGICAL SCREENING

166  Most of the approaches developed prior to 2010 made CRM predictions based on the DNA sequence.
167  However, there are numerous epigenetic markers associated with gene regulation **Bernstein et al.** (2007),
168  which can be tested with an array of high-throughput techniques. Until relatively recently, the time
169  and expense associated with assays for these marks limited their use. However, breakthroughs in high-
170  throughput techniques have made large amounts of epigenetic data available. Therefore, interest has grown
171  in using these marks to indicate regulatory regions biologically. A few large efforts have started to map
172  regulatory markers genome-wide. The ENCODE consortium has used many of these assays across the
173  human genome **Dunham et al.** (2012) and these data are publically available through the UCSC Genome
174  Browser **Rosenbloom et al.** (2013), **Kent et al.** (2002). Similar projects exist for *Drosophila melanogaster*
175  and *Caenorhabditis elegans* **Celniker et al.** (2009), and for *Mus musculus* **Stamatoyannopoulos et al.**
176  (2012).

177      Here, we will briefly describe some of the major assays used for high-throughput analysis of epigenetic
178  data. Although there are other methods available, these will serve to illustrate the types of data that are
179  available.

## 4.1  CHIP-SEQ

180  Chromatin immunoprecipitation followed by sequencing (ChIP-seq) identifies proteins bound to DNA
181  **Kim and Ren** (2006), **Barski et al.** (2007), **Ren et al.** (2000). This can be directed at identifying
182  transcription factors bound to chromatin, histone modifications (that mark regulatory regions), RNA
183  polymerase II binding, and more. One limitation of this technique is its requirement for an antibody
184  that binds to the protein targeted by the assay. Antibodies are not available for all regulatory proteins.

## 4.2 DNASE-SEQ

185 DNase hypersensitive site sequencing (DNase-seq) identifies regions where the DNase I enzyme digests
186 chromatin more readily **Crawford et al.** (2006) (which are then sequenced with high-throughput
187 techniques). These regions of open chromatin more accessible to regulatory proteins. Active regulatory
188 regions, bound by transcription factors, form DHSs **Cockerill** (2011), making them a useful tool in
189 identifying candidate regulatory regions.

## 4.3 FAIRE-SEQ

190 Formaldehyde assisted isolation of regulatory elements followed by sequencing (FAIRE-seq) **Giresi**
191 **and Lieb** (2009), **Giresi et al.** (2007) is another technique for indetifying regions of accessible
192 chromatin. Formaldehyde cross-links more efficiently to nucleosomes, allowing nucleosome-free DNA
193 to be identified. These regions are associated with regulatory activity.

## 4.4 RRBS

194 Reduced representation bisulfite sequecing (RRBS), is a method that can be used to identify methylated
195 cytosine. Methylated DNA is often repressed, reducing transcription. The chromatin is treated with
196 sodium bisulfite, converting unmethylated cytosine to uracil **Frommer et al.** (1992). The remaining
197 cytosines are methylated. The reduced representation refers to the fact that the genome is fragmented and
198 selected for regions with high CpG dinucleotide content **Dunham et al.** (2012), **Meissner et al.** (2005)
199 where methylation represses transcription. The fragments can then be sequenced using high-throughput
200 techniques **Meissner et al.** (2008).

## 4.5 5C

201 Chromosome Conformation Capture Carbon Copy (5C) is used to identify interactions between regions
202 of chromatin **Dostie et al.** (2006). For example, it can help identify the interaction of enhancers with
203 promoters **Sanyal et al.** (2012). It is essentially a high-throughput version of Chromosome Conformation
204 Capture (3C), a method that uses formaldehyde to crosslink chromatin and then detects the crosslinked
205 regions **Dekker et al.** (2002).

206     These biological assays have brought researchers much closer to goal of quickly and affordably being
207 able to test genomes directly for regions that control gene activity. Nevertheless, that goal remains
208 unrealized. There is not yet a test that can simply and reliably identify enhancers and other CRMs, as
209 well as the transcription factor binding sites within them. Rather, these assays reveal the characteristics of
210 the chromatin, some of which are associated with regulatory regions, but they are pieces of a larger puzzle.
211 Although they undoubtedly open new avenues for investigation, they do not represent the entire regulatory
212 code. However, used in conjunction with, or as part of, computational methods, they are likely to further
213 our understanding of these complex processes. Of course, at the time being, these data are available for
214 an extremely limited number of genomes. Researchers working with model organisms outside the scope
215 of the few large projects mentioned at the beginning of this section seldom have the resources to run all of
216 these assays genome-wide. Additionally, there have been developments in the high-throughput validation
217 of candidate regulatory regions. One method, developed for *Strongylocentrotus purpuratus*, allows the
218 simultaneous validation, and quantification of effect, of up to 130 candidate CRMs **Nam and Davidson**
219 (2012). Such techniques should soon make it easier to work with computational methods and validate
220 their predictions.

# 5  DEVELOPMENTS IN CRM PREDICTION

Prior to the last major computational CRM prediction reviews **Su et al.** (2010), **Van Loo and Marynen** (2009), most approaches made CRM predictions based on clustering of TFBSs, conservation of noncoding DNA, or both, using a variety of techniques. Most of these approaches made predictions based on the DNA sequence alone. Nearly all new methods have moved beyond these techniques to include other sources of information, instead of, or in addition to, these CRM characteristics.

 CRM prediction can mean different things. In reality, most methods attempt to solve some particular aspect of the problem. There are significant differences between some tools in this regard. In fact, this presents an opportunity for the researcher, since different tools can be used to uncover various possible aspects of regulation. Here, we will provide information relevant to making an informed decision about when a particular tool might be useful. Given the variety of approaches, a meaningful experiment comparing the predictive value of each tool would be difficult to design. Other attempts to do so have arbitrarily ignored some extant methods because they did not fit the experimental design **Su et al.** (2010), **Klepper et al.** (2008). Here, we will focus on describing the design, capability, and usability of each system in order to facilitate an informed decision of the appropriate tool. It is important to note that there are inherent trade-offs in any approach to prediction and that it is unlikely that any one approach is the best choice in all scenarios. This was found to be the case with the older methods **Su et al.** (2010), and as we will show, it remains true with the newer ones as well.

## 5.1  I-CISTARGET

**Location:** http://med.kuleuven.be/lcb/i-cisTarget
**Type:** web page
**Input:** a list of gene IDs or ChIP peak locations
**Species:** *D. melanogaster*

 i-cisTarget **Herrmann et al.** (2012) predicts CRMs that may regulate an input set of genes or be related to an input set of ChIP peak locations. It exploits nearly all of the CRM characteristics (Section 2.2) as part of its search. Currently it works only for the *D. melanogaster*. The algorithm consists of the following key steps:

1. *Partition the Noncoding Genome.* The algorithm starts by partitioning all the noncoding DNA. This partition is performed on the basis of conservation by using PhastCons conservation scores **Siepel et al.** (2005). The subsets of the partition are centered on regions of conservation.

2. *Detecting Homotypical Clusters.* ClusterBuster **Frith** (2003) is run repeatedly on the partition using one PWM at a time from a library of over 6000 PWMs. It is also run on orthologous regions from other *Drosophila* genomes identified by the liftOver tool in the UCSC Genome Browser **Rosenbloom et al.** (2008). In this way, the partition is ranked according to homotypical clustering and conservation.

3. *In Vitro Events.* The regions of the partition are also scored for in vitro events (iVEs). These are from ChIP-seq or ChIP-chip experiments that assay histone modifications or the binding of transcription factors.

4. *Identify Candidate Regulatory Regions and Filter Input Locations.* Candidate regulatory regions are defined for each gene. The user can select one of the definitions, such as one that includes the 5kb upstream, the 5' UTR, and the first intron. The input set of ChIP locations (if any) is filtered to include only those that overlap at least 40% (user settable) or more of a candidate regulatory region or ChIP peak region (iVEs from step 3).

5. *Calculate Enrichment.* Given a set of co-expressed genes, or a set of ChIP peak locations, i-cisTarget calculates top-ranked regions of the partition that are enriched for the input. Any or all of the calculated ranks can be used for the enrichment analysis.

265    6. *Predict Enhancers.* The regions of the partition that were most enriched are then predicted as
266       enhancers. These candidate enhancers can be further scanned for homotypic or heterotypic clustering
267       of binding sites. Since clustering of binding sites is used, the locations of binding sites are also
268       predicted.

269    The first three of these steps are pre-calculated and do not need to be run each time.

270    i-cisTarget is available through a web interface that is user-friendly. The minimum input is a list of gene
271    IDs for co-expressed genes or locations for ChIP peaks. The user can select from a number of preset
272    regulatory regions and the overlap fraction. The enrichment score threshold lets the user set the stringency
273    of recovery required for enrichment. The ROC threshold sets the fraction of regions that are considered
274    for being "top-ranked," when checking the enrichment of top-ranked regions. Locations of known CRMs
275    can be provided as well.

276    The output is a list of features that were enriched in the input, along with scores and links to the locations
277    the features were found. i-cisTarget is the only method we are aware of to combine PWM scanning,
278    clustering, conservation, and epigenetic data to predict enhancers. Nevertheless, there are some notable
279    limitations. i-cisTarget is currently available only for *D. melanogaster*. Also, the transcription factor
280    binding site clustering technique is applicable only to homotypical clustering - but these characterize only
281    a subset of CRMs (Section 2.2), although the regions from the partition that are predicted as enhancers
282    can be further scanned for heterotypic clustering. Finally, the feature labels in the output are not very
283    enlightening. The nomenclatures are based on the data source. Each has its own abbreviations and there
284    is no reference list, or crosslinked information to explain their meaning (e.g. BDTNP-da_2_050307).

285    A distinguishing feature of this method is the ability to extract CRMs on the basis of enriched features
286    in co-expressed genes. The results for i-cisTarget are organized by feature. Candidate targets of a given
287    feature are listed by clicking on a link. However, a number of features can be selected and common targets
288    found.

## 5.2   CRMMINER

289    **Location:** http://www.biomedcentral.com/1471-2105/13/25/additional
290    **Type:** downloadable (Linux)
291    **Input:** DNA test and control sequences in FASTA format and a library of PWMs
292    **Species:** any with PWMs available
293

294    CrmMiner **Girgis and Ovcharenko** (2012) searches noncoding DNA near coexpressed genes to identify
295    CRMs that may be responsible for their co-regulation. Unlike the best performing systems in prior reviews
296    **Su et al.** (2010), **Van Loo and Marynen** (2009), it does not use conservation. Searching for clusters of
297    PWMs can present certain problems. Often, the user must specify a limited set of PWMs to scan for.
298    That is, they must understand something of how the gene is regulated in advance. In other cases, a larger
299    library of PWMs has been allowed, but in the past this has resulted in a large number of false positives.
300    CrmMiner moves beyond homotypic or heterotypic clustering of binding sites on the DNA sequence,
301    using the concept of the "regulatory signature". The signature is composed of pairs of motifs that act
302    synergistically. CrmMiner learns this signature from the input, but validated CRMs are not required for it
303    to work.

304    Input to CrmMiner consists of two sets of sequences (mixed and control), and a library of PWMs.
305    The mixed and control sequences are further subdivided into three groups (training, validation, and test).
306    The mixed should contain sequences suspected of containing regulatory regions, mixed with background
307    sequences. The control set contains sequences unlikely to contain CRMs. The PWMs must be from the
308    TRANSFAC database **Matys et al.** (2003), or in the same format. CrmMiner then proceeds through the
309    following steps:

1. *Scan for PWMs.* CrmMiner scans all the input sequences for matches against the library of PWMs. This is performed using the motif scanner MAST **Bailey and Gribskov** (1998).

2. *Identify Enriched Pairs of Motifs.* CrmMiner identifies pairs of motifs from the scan that are near one another, occur multiple times in the mixed data, and are enriched in the mixed data compared to the control data.

3. *Identify Enriched Sequences.* CrmMiner identifies sequences from the mixed set that are enriched (compared to the control) when using the motif pairs from the last step to select them. It then creates a list of motif pairs in these sequences that represent a regulatory signature for the co-expressed genes.

4. *Training.* CrmMiner trains a Bayesian classifier on the mixed sequences to find the scores for enrichment that will maximally separate the mixed and control sequences using the regulatory signature.

5. *Validation.* CrmMiner uses the regulatory signature to identify CRMs in the validation set for both the mixed and control data. It keeps trying different parameters until it performs well on both the training and validation data.

6. *Testing.* CrmMiner runs on the test data to identify more CRMs.

Although it does not directly use epigenetic data, CrmMiner is able to make predictions that use it indirectly. That is because of the way it utilizing training data. Rather than being dependent on validated CRMs, CrmMiner expects input data to be mixed. A researcher can input sequences that were identified by biochemical markers as being interesting (see Section 4). Therefore, it can learn associations of motifs that act synergistically within tissue specific enhancers to construct a context specific regulatory signature. Because this approach is not dependent on particular signatures, it is quite flexible (the sequences could be identified by a variety of assays). Also, unlike many past approaches, the user does not need to select PWMs thought active in a particular biological context.

The output is a list of genomic locations and scores for candidate CRMs, as well as the pairs of motifs in the discovered regulatory signature.

Nevertheless, CrmMiner does have a couple of significant limitations. First, it depends on TRANSFAC PWMs. TRANSFAC requires licensing fees that may be an obstacle for some users. Still, it is possible to convert PWMs from other databases (such as JASPAR **Mathelier et al.** (2014)), to the TRANSFAC format **Thomas-Chollier et al.** (2008). Second, installation of CrmMiner is fairly technical, even for those comfortable with running tools on the command line. It requires the user to download other dependencies, compile, and install them, find the location of installed dependencies, and edit configuration files. Furthermore, CrmMiner does not have a way of selecting candidate CRMs out of unbroken sequence. The user must supply short sequences to be tested as candidate CRMs.

### 5.3 MATRIXCATCH

**Location:** http://gnaweb.helmholtz-hzi.de/cgi-bin/MCatch/MatrixCatch.pl
**Type:** webpage and downloadable (Windows and Linux)
**Input:** DNA sequences in plain text, FASTA, or EMBL format
**Species:** any with PWMs available

MatrixCatch **Deyneko et al.** (2013) does not predict CRMs exactly, it predicts composite elements (CEs) in one or more sequences. These are pairs of transcription factor binding sites that are predicted to act synergistically, in a manner similar to CrmMiner. Here, however, the pairs of motifs are not used to rank CRM predictions, they are used on their own as what might be thought of as mini-CRMs.

Similar to CrmMiner, MatrixCatch does not utilize epigenetic or conservation data to make predictions, but like most recent CRM prediction methods, it utilizes information beyond the DNA sequence. In MatrixCatch, this information comes from a database of known transcription factor interactions called

355 TransCompel **Kel-Margoulis et al.** (2002). It is also possible for users to upload their own interaction
356 data.

357 MatrixCatch uses a library of PWMs to create a model of CEs, based on the known interactions in
358 TransCompel. These models are then used to scan the input sequences for similar CEs. The output
359 provides a set of transcription factor pairs, their locations, orientations, and distance between each element
360 of the pairs. It also provides a graphical display of motif locations along the input sequence.

361 MatrixCatch is very simple to use, and the output is easy to interpret. This cannot be said of all CRM
362 prediction techniques. To our knowledge, it is the only available method for predictions based on validated
363 transcription factor interactions. A possible limitation of this approach is that MatrixCatch does not predict
364 full CRMs, but it should be kept in mind that there is no method that can accurately define the size of
365 a CRM. MatrixCatch at least provides additional information about interactions between transcription
366 factors within a CRM that is not available with most other methods.

## 5.4 CORECLUST

367 **Availability:** http://sourceforge.net/projects/coreclust/
368 **Type:** downloadable (any with Java)
369 **Input:** text file with orthologous (putative) regulatory regions or regions near known co-regulated genes,
370 text files with PWMs, and FASTA file with sequences to search for co-regulated CRMs
371 **Species:** any
372

373 CORECLUST **Nikulova et al.** (2012) attempts to learn the regulatory code controlling a particular gene
374 and builds a model that can be used to predict coregulated genes, genome-wide. Like many of the previous
375 era of CRM predictions methods, it relies on a Hidden Markov Model (HMMs). This is a statistical model
376 of the probability of states following each other in a sequence **Eddy** (2004). Therefore, it is readily
377 applicable to questions surrounding DNA sequences that may transition from regulatory to background to
378 coding sequence.

379 Essentially, CORECLUST trains a model of CRMs controlling a particular gene or set of genes and
380 uses the model to search for CRMs that are controlled in a similar manner. The model it builds may
381 represent part of the regulatory code for genes that need to be active in a particular context. This model is
382 based not only on the composition of particular binding sites in a CRM, it considers the arrangement
383 and spacing of binding sites to attempt to capture the interactions between transcription factors. In
384 considering conservation, CORECLUST does not perform a classic alignment, but instead uses the
385 PWMs to find the similarity in arrangment of binding sites between orthologous sequences. In these
386 ways, CORECLUST follows from and extends the methods of the most successful CRM predictors from
387 the previous generation.

388 CORECLUST relies heavily on the idea of motif interactions and allows for different distributions to be
389 used in modelling binding site spacing. It builds a model of CRM structure considering motif composition,
390 frequency, arrangment, and distribution. An interesting result of this approach is a detailed description of
391 CRM structure in the training regions, which in itself may be useful. A limitation of this approach is that
392 CORECLUST is not able to use a library of PWMs to build a CRM model. The user must have some idea
393 of which factors regulate the training sequences beforehand.

394 CORECLUST is a downloadable program with a command line interface, which may be difficult for
395 some users. It is relatively easy to install and comes with examples of how to run it. The input is not overly
396 difficult to format (a failure of many command line programs). A number of optional parameters can be
397 set to fine tune control of the program.

### 5.5 IMOGENE

**Location** http://mobyle.pasteur.fr/cgi-bin/portal.py#forms::imogene and https://github.com/hrouault/Imogene
**Type:** webpage and downloadable (Linux)
**Input:** a set of genomic coordinates of validated CRMs
**Species:** Eutherian and Drosophilae

Imogene **Rouault et al.** (2014) is not a general CRM predictor. It predicts CRMs genome-wide that are found to be similar to input CRMs. This allows a researcher to take a small number of experimentally validated CRMs and use them to find putative regulatory regions that may be active under the same conditions as the input. In addition, Imogene performs *de novo* motif prediction at the same time that CRMs are predicted, generating PWMs that are based on locations within the input. By obviating the need for a database of PWMs, such as those provided by TRANSFAC **Matys et al.** (2003) or JASPAR **Mathelier et al.** (2014), Imogene provides the additional advantage of identifying PWMs that are particularly apt to the species or context under consideration. The PWMs generated by Imogene can then be compared to matrices from PWM databases to provide clues about which specific transcription factor may be regulating the predicted CRMs.

Imogene is provided as both a web interface and downloadable program. The web interface is limited to CRM identification in Drosophilae and Eutherians. The input is given as a set of genomic coordinates, but these must be lifted over to the fruit fly or mouse genomes used in the paper (*D. melanogaster* release 5 and mm9 respectively).

Imogene is based on a three step process:

1. *Expand training data.* The training data is expanded with orthologous sequences based on alignments with related genomes.
2. *Identify PWMs.* The training sequences are used to identify a set of motifs, of user specified width, that have a significantly different distribution in the training sequences than in background sequences. These motifs are scored based on their distribution and their conservation.
3. *Predict CRMs.* The PWMs identified in the previous step are now used to scan the genome for intergenic regions with a similar motif content to the training CRMs. These newly identified CRMs are considered to function under similar condition to the training CRMs.

A procedure is given by the authors for using Imogene to predict the type of a CRM. That is, given a genome region that is thought to be a CRM, what genomic context does it function in? The details of this approach are beyond the scope of this review.

### 5.6 CHROMHMM

**Location:** http://compbio.mit.edu/ChromHMM/
**Type:** downloadable (any with Java)
**Input:** a set of BED files with aligned reads of chromatin modification marks
**Species:** any with epigenetic data available

ChromHMM is not a dedicated CRM prediction tool **Ernst and Kellis** (2012). It is a tool for analyzing the state of chromatin by finding characteristic combinations and arrangements of chromatin modification marks. This is an unsupervised learning approach that builds models of chromatin state based on patterns of these marks. Nevertheless, this is a powerful tool for discovering CRMs that regulate a gene under particular conditions.

439     Chromatin state is an important part of gene regulation and is closely connected with cis-regulatory
440 modules. Certain chromatin modification marks are associated with promoters and enhancers, but
441 individual marks do not act as switches that activate a module for regulation. It is only by examining the
442 combination and arrangement of these modifications, active enhancers and promoters can be recognized.

443     ChromHMM, as the name implies, is based on a hidden Markov model. A model is built for the
444 probability of moving from one chromatin state to another as one proceeds through a sequence. The
445 number of states is specified by the user. The sequence, in this case, is build from the combination of
446 chromatin modification marks present at a location. Therefore, although ChromHMM cannot directly
447 predict CRMs, it can indicate the probability of a particular region being in a particular chromatin state.
448 The researcher can then identify states that seem probable for association with enhancers and promoters.

449     ChromHMM is a downloadable program that must be run on the command line. For users comfortable
450 with running command line programs, it presents no major obstacles.


## 6   DISCUSSION

451 The development of high-throughput assays for epigenetic marks related to gene regulation is a powerful
452 new source of data for revealing the complex code behind differential gene expression. Nevertheless, these
453 data have not lessened the utility of computational approaches to CRM prediction. In their 2012 review,
454 Hardison and Taylor proposed a multi-faceted approach **Hardison and Taylor** (2012). They suggested
455 starting with mapping epigenetic features to a region of interest. Necessarily this means a context or
456 contexts of interest as well, since biological assays are only relevant to a particular cell type and condition.
457 They then suggested applying conservation and binding site analysis to regions predicted to be interesting
458 by epigentic features. However, determining the precise combination of marks associated with an active
459 CRM is not necessarily straightforward. As we discussed earlier (Section 4), these data do not, in and of
460 themselves, indicate cis-regulatory modules. We have yet to determine the full range of epigenetic changes
461 that regulate genes, or to understand the full implication of those we do know about. There are dozens of
462 known histone marks alone **Tan et al.** (2011), although some of these marks may be redundant **Zentner**
463 **and Henikoff** (2013). Furthermore, epigenetic marks are not a binary proposition, in which a mark is
464 either observed or not under a particular condition. The chromatin state is dynamic and most assays must
465 be interpreted for reliability or strength of signal. Therefore, computational methods of analyzing the
466 subtle interactions of epigenetic marks are a useful tool in trying to understand the epigenetic code.

467     Given the suggestions of Hardison and Taylor, and the past success in CRM prediction using clusters
468 of known PWMs and conservation, we expected that more recent approaches would build on previous
469 methods, incorporating epigenetic data to improve predictions. In one case, this is exactly what we found.
470 i-cisTarget incorporates (Section 5.1) clustering of known TFBSs (both homotypical and heterotypical),
471 conservation of genomic segments, and ChIP-seq data to make predictions. However, none of the other
472 publically available, recently developed systems took this approach. In fact, the most common feature
473 utilized by the methods we reviewed is motif synergy. Albeit in different ways, motif synergy is used by
474 CrmMiner, MatrixCatch, and CORECLUST. Another point raised by Hardison and Taylor is the need for
475 *de novo* motif prediction, in addition to scanning for PWMs. The full complement of transcription factors
476 is not known for any genome, but perhaps more significantly, the full range of binding sites for those
477 transcription factors is even less well understood. Therefore, methods of *de novo* motif prediction, such as
478 that used by Imogene, might useful in learning unknown regulatory pathways. On the more extreme end
479 of using epigenetic data in predictions lies ChromHMM, which learns chromatin states (including active
480 enhancers and promoters) without reference to DNA sequence.

481     Interestingly, all of the tools in this review have unique capabilities. We believe it is important to consider
482 the strengths of a method for the questions being pursued.

483     i-cisTarget is currently only available for *D. melanogastar*, but for those working with this organism it
484 should be an excellent choice. The number of features considered by i-cisTarget is more comprehensive

than any other tool, allowing researchers to pick out high confidence regions for further research. The presentation of results by feature allows features to be combined to find targets that are common to them. Furthermore, although initial scanning is only for regions of homotypic clustering, any regions associated with a particular feature can be further scanned for heterotypic clustering of binding sites. The incorporation of epigenetic data offers a significant advantage to i-cisTarget over the previous generation of tools, especially since the results are not limited to only those correlated with epigenetic marks. The biggest limitation of this approach may be the reliance on homotypic clusters of binding sites during the initial stages of the algorithm, although it should be borne in mind that any algorithm that uses conservation may miss lineage-specific CRMs.

CrmMiner's ability to self-tune its parameters through repeated training and validation cycles is a powerful technique. Furthermore, a publically available tool capable of identifying regions of the genome with a characteristic pattern of synergistic motif interactions is an exciting development that brings us one step closer to cracking the regulatory code. It is unfortunate that CrmMiner is somewhat challenging to install and configure. It is also left to the user to segment the input into parts that may represent regions of regulation. This seems like an important step that should not be left to chance. i-cisTarget answers this challenge by centering its segmentation based on regions of conservation, which is at least a reasonable motivation, given that CRMs frequently exhibit higher levels of conservation.

MatrixCatch deserves high marks for ease of use. It is by far the easiest tool to use of those reviewed. Although it does not predict entire CRMs, such a claim would be farfetched for most tools to make, since there are no clear features marking the beginning and end of CRMs. The visualizations provided by MatrixCatch make it easy to see regions of clustering and the interactions among predicted binding sites are based on validated interactions, which is an interesting approach.

CORECLUST has a couple of innovations that are interesting. Similar to a number of the methods we reviewed, CORECLUST considers not just the clustering of motifs but also their relative arrangement, in order to capture synergistic effects. Although it uses conservation as a filter, it considers conservation only of binding sites and their arrangement, rather than performing an alignment, which has been shown to be advantageous **Su et al.** (2010). CORECLUST's biggest limitation is that it is not able to consider a giant library of PWMs in the way that MatrixCatch, CrmMiner, and especially i-cisTarget are able to. Instead, the user must select PWMs thought to be relevant to the regulation of the gene at hand. However, the model it builds of the training data may be useful on its own, to understand the interactions regulating a known CRM, as well as to use the model for identifying similar CRMs. Furthermore, CORECLUST is able to search genome-wide for similar CRMs.

Imogene is the only method reviewed that can perform *de novo* motif prediction at the same time it searches for CRMs. It takes an interesting approach to using conservation, by expanding the training data with orthologous regions from related genomes. This method may be particularly useful for identifying the regulation of genes involved in poorly studied pathways or organisms.

ChromHMM is the only tool other than i-cisTarget that we reviewed to use epigenetic data as part of its search. It is not limited to a particular genome like i-cisTarget (although there are few genomes available that are mapped with epigenetic data – see Section 4). Its biggest limitation is that it does not directly identify candidate CRMs. It identifies chromatin states, some of which will be associated with CRMs, such as enhancers. Predictions using ChromHMM that are linked to particular CRMs are available through the UCSC genome browser.

There is not yet a publically available tool that fully implements the suggestions of Hardison and Taylor. That is, there is no one tool that helps researchers to analyze cis-regulatory modules by starting with epigenetic data and allowing other sources of information to be brought in to aid that analysis, including conservation and clustering of PWMs. However, as Hardison and Taylor note, epigenetic data are not yet available for most organisms. Furthermore, there are limitations in the kind of data that are available. For example, transcription factor binding assayed by ChIP-seq is only useful for transcription factors that antibodies exist for. Therefore, we believe there is utility in having a range of approaches that are able to predict CRMs based on different kinds of data.

## 6.1 RECOMMENDATIONS

535 For those working with *Drosophila* who want to identify CRMs or other regulatory regions associated
536 with a set of genes or epigenetic marks, we recommend starting with i-cisTarget. The web-based
537 interface makes it easy to use and the option of including epigenetic data in the analysis provides a clear
538 advantage over the previous generation of tools and should provide high confidence candidates for further
539 experimentation.

540 For those who would like to identify CRMs that are similar to those controlling genes known to be
541 co-regulated (and for which the PWMs are known), we recommend CORECLUST as a starting point.
542 Although similar to a number of previous methods, CORECLUST provides two innovations that have
543 been shown to significantly improve its performance: 1) it considers motif synergy, 2) it considers
544 conservation of binding sites. For this more limited case, CORECLUST showed impressive performance
545 during validation.

546 If the set of transcription factors regulating a set of genes is unknown, we recommend starting with either
547 CrmMiner or MatrixCatch. There are advantages to either choice. CrmMiner's sophisticated learning
548 algorithm helps it to distinguish CRMs from background sequences with a high degree of precision.
549 Therefore, if the highest confidence candidates are desired, CrmMiner is a good choice (although one
550 should expect the sensitivity to be lower). MatrixCatch also shows improved performance compared to
551 previous methods, and its web-based interface is significantly easier to use than CrmMiner's command
552 line interface.

553 Following any of these methods, we recommend running Imogene if the species of interest is a Eutherian
554 or *Drosophilae*. It has a simple, web-based interface and perform *de novo* motif prediction in conjunction
555 with CRM prediction. Although it uses a statistical approach to find motifs with a significantly different
556 distribution in the training sequences compared to background sequence, it uses orthologous sequences to
557 expand the training data. Based on the training, Imogene builds PWMs from the input and uses them to
558 search the test data. This is a significant advantage to Imogene, since it can describe unknown transcription
559 factor binding sites.

560 Finally, if one did not start with i-cisTarget, but there are epigenetic data available for the species being
561 studied, we recommend the use of ChromHMM. This will indicate significant combinations of epigenetic
562 marks that indicate various chromatin states in relation to the data. It will help focus the study on the
563 regions of the most interest. Although one limitation of epigenetic data is that it is only valid for a
564 particular cell type, this is also a significant advantage, since it can focus attention on cell types that
565 are of interest for a particular line of work.

## 6.2 OTHER TOOLS

566 In this review, we have focused on methods that are publically available to researchers. That is, tools that
567 maintain a website or location where the tool can be downloaded without prior permission of the authors.
568 Other methods do exist. We have developed GAMMI **Gagne and Congdon** (2012) and GAMI-CRM
569 **Thompson and Congdon** (2014), which both use genetic algorithms as part of a heuristic approach
570 to identifying CRMs. GAMMI takes a library of motifs that were identified by another algorithm and
571 identifies sets that exhibit conserved clustering. GAMI-CRM identifies clusters of conserved sites, which
572 are identified *de novo* and also facilitates the use of epigenetic information. ChroModule **Won et al.** (2013)
573 predicts CRMs based on models of the continuous histone modification data, as opposed to the discrete
574 peaks used by ChromHMM. It showed improved performance compared to ChromHMM in its validation.
575 CRFEM **Gan et al.** (2014) is another approach that identifies CRMs and binding sites *de novo* by finding
576 clusters of overrepresented motifs and scoring them using other features, such as epigenetic data.

577 There are also a couple of tools that incorporate other methods at part of a larger regulatory analysis
578 tool, such as MotifLab **Klepper and Drabløs** (2013) for general use and **Kazemian et al.** (2011) Genome

579 Surveyor for *Drosophila*. These tools allow the user to run various analyses using the integrated methods
580 and to visualize the results.


## 7 CONCLUSIONS

581 Despite the wealth of epigenetic data that are available for some genomes (human, fruit fly, and mouse),
582 most computational methods are not yet making use of it. Those that are use only a subset of the available
583 data. This leaves open the possibility of far more sophisticated methods that predict CRMs active in
584 a particular context, elucidate the the gene regulatory network, and more accurately identify the genes
585 activated or repressed by particular CRMs.

586 Epigenetic data carries both significant advanatages and disadvantages. In the future, we hope to see
587 more tools that integrate these data with other methods of prediction, to take full advantages of the
588 strengths of each. Ideally, flexibility in how the data are used will be maintained so that researchers
589 can choose the characteristics that are the most important to them in their work. Biology, perhaps more
590 than most sciences, is full of exceptions to the "rules" we discover. In order to discover these exceptions,
591 we need to know both what passed the filters we create and what did not.


## DISCLOSURE/CONFLICT-OF-INTEREST STATEMENT

592 The authors declare that the research was conducted in the absence of any commercial or financial
593 relationships that could be construed as a potential conflict of interest.


## AUTHOR CONTRIBUTIONS

594 Jeffrey A. Thompson and Clare Bates Congdon developed the concept for the structure and content of this
595 manuscript. Jeffrey A. Thompson researched and wrote the initial draft. Clare Bates Congdon critically
596 revised the manuscript. Both authors reviewed and approved the final version of the manuscript.


## ACKNOWLEDGEMENT

## REFERENCES

601 Bailey, T. L. and Gribskov, M. (1998), Combining evidence using p-values: application to sequence
602    homology searches., *Bioinformatics*, 14, 1, 48–54
603 Bannister, A. J. and Kouzarides, T. (2011), Regulation of chromatin by histone modifications, *Cell
604    research*, 21, 3, 381–395
605 Barski, A., Cuddapah, S., Cui, K., Roh, T.-Y., Schones, D. E., Wang, Z., et al. (2007), High-resolution
606    profiling of histone methylations in the human genome, *Cell*, 129, 4, 823–837
607 Bernstein, B. E., Meissner, A., and Lander, E. S. (2007), The mammalian epigenome, *Cell*, 128, 4,
608    669–681

609  Blackwood, E. M. and Kadonaga, J. T. (1998), Going the distance: a current view of enhancer action,
610      *Science*, 281, 5373, 60–63
611  Butler, J. E. and Kadonaga, J. T. (2002), The rna polymerase ii core promoter: a key component in the
612      regulation of gene expression, *Genes & development*, 16, 20, 2583–2592
613  Calo, E. and Wysocka, J. (2013), Modification of enhancer chromatin: what, how, and why?, *Molecular
614      cell*, 49, 5, 825–837
615  Celniker, S. E., Dillon, L. A. L., Gerstein, M. B., Gunsalus, K. C., Henikoff, S., Karpen, G. H., et al.
616      (2009), Unlocking the secrets of the genome, *Nature*, 459, 7249, 927–930
617  Chen, Q. K., Hertz, G. Z., and Stormo, G. D. (1997), PromFD 1.0: a computer program that predicts
618      eukaryotic pol II promoters using strings and IMD matrices, *Computer applications in the biosciences:
619      CABIOS*, 13, 1, 29–35
620  Cheng, C., Shou, C., Yip, K. Y., and Gerstein, M. B. (2011), Genome-wide analysis of chromatin features
621      identifies histone modification sensitive and insensitive yeast transcription factors, *Genome Biol*, 12,
622      11, R111
623  Claverie, J.-M. and Sauvaget, I. (1985), Assessing the biological significance of primary structure
624      consensus patterns using sequence databanks. i. heat- shock and glucocorticoid control elements in
625      eukaryotic promoters., *Comput Appli Biosci*, 1, 2, 95–104
626  Cockerill, P. N. (2011), Structure and function of active chromatin and DNase i hypersensitive sites:
627      Active chromatin and DNase i hypersensitive sites, *FEBS Journal*, 278, 13, 2182–2210
628  Consortium, E. P. et al. (2011), A user's guide to the encyclopedia of dna elements (encode), *PLoS
629      biology*, 9, 4, e1001046
630  Crawford, G. E., Holt, I. E., Whittle, J., Webb, B. D., Tai, D., Davis, S., et al. (2006), Genome-wide
631      mapping of dnase hypersensitive sites using massively parallel signature sequencing (mpss), *Genome
632      research*, 16, 1, 123–131
633  Dekker, J., Rippe, K., Dekker, M., and Kleckner, N. (2002), Capturing chromosome conformation,
634      *science*, 295, 5558, 1306–1311
635  Deyneko, I. V., Kel, A. E., Kel-Margoulis, O. V., Deineko, E. V., Wingender, E., and Weiss, S. (2013),
636      Matrixcatch-a novel tool for the recognition of composite regulatory elements in promoters, *BMC
637      bioinformatics*, 14, 1, 1–10
638  Dostie, J., Richmond, T. A., Arnaout, R. A., Selzer, R. R., Lee, W. L., Honan, T. A., et al. (2006),
639      Chromosome conformation capture carbon copy (5c): a massively parallel solution for mapping
640      interactions between genomic elements, *Genome research*, 16, 10, 1299–1309
641  Dunham, I., Kundaje, A., Aldred, S. F., Collins, P. J., Davis, C. A., Doyle, F., et al. (2012), An integrated
642      encyclopedia of DNA elements in the human genome, *Nature*, 489, 7414, 57–74
643  Eddy, S. R. (2004), What is a hidden markov model?, *Nature biotechnology*, 22, 10, 1315–1316
644  Elnitski, L., Jin, V. X., Farnham, P. J., and Jones, S. J. (2006), Locating mammalian transcription factor
645      binding sites: A survey of computational and experimental techniques, *Genome Research*, 16, 12,
646      1455–1464
647  Ernst, J. and Kellis, M. (2012), Chromhmm: automating chromatin-state discovery and characterization,
648      *Nature methods*, 9, 3, 215–216
649  Espada, J. and Esteller, M. (2010), Dna methylation and the functional organization of the nuclear
650      compartment, *Seminars in cell & developmental biology*, 21, 2, 238–246
651  Fickett, J. W. (1996), Quantitative discrimination of MEF2 sites., *Molecular and cellular biology*, 16, 1,
652      437–441
653  Frith, M. C. (2003), Cluster-buster: finding dense clusters of motifs in DNA sequences, *Nucleic Acids
654      Research*, 31, 13, 3666–3668
655  Frommer, M., McDonald, L. E., Millar, D. S., Collis, C. M., Watt, F., Grigg, G. W., et al. (1992), A
656      genomic sequencing protocol that yields a positive display of 5-methylcytosine residues in individual
657      DNA strands., *Proceedings of the National Academy of Sciences*, 89, 5, 1827–1831
658  Gagne, D. J. and Congdon, C. B. (2012), Preliminary results for GAMMI: Genetic algorithms for motif-
659      module inference, in X. Li, ed., Proceedings of the 2012 IEEE Congress on Evolutionary Computation
660      (Brisbane, Australia), 1309–1316

Gallo, S. M., Gerrard, D. T., Miner, D., Simich, M., Des Soye, B., Bergman, C. M., et al. (2011), Redfly v3. 0: toward a comprehensive database of transcriptional regulatory elements in drosophila, *Nucleic acids research*, 39, suppl 1, D118–D123

Gan, Y., Guan, J., Zhou, S., and Zhang, W. (2014), Identifying cis-regulatory elements and modules using conditional random fields, *IEEE/ACM Transactions on computational biology and bioinformatics*, 11, 1

Ghirlando, R., Giles, K., Gowher, H., Xiao, T., Xu, Z., Yao, H., et al. (2012), Chromatin domains, insulators, and the regulation of gene expression, *Biochimica et Biophysica Acta (BBA)-Gene Regulatory Mechanisms*, 1819, 7, 644–651

Giresi, P. G., Kim, J., McDaniell, R. M., Iyer, V. R., and Lieb, J. D. (2007), FAIRE (formaldehyde-assisted isolation of regulatory elements) isolates active regulatory elements from human chromatin, *Genome Research*, 17, 6, 877–885

Giresi, P. G. and Lieb, J. D. (2009), Isolation of active regulatory elements from eukaryotic chromatin using FAIRE (formaldehyde assisted isolation of regulatory elements), *Methods*, 48, 3, 233–239

Girgis, H. Z. and Ovcharenko, I. (2012), Predicting tissue specific cis-regulatory modules in the human genome using pairs of co-occurring motifs, *BMC bioinformatics*, 13, 1, 25

Hardison, R. C. and Taylor, J. (2012), Genomic approaches towards finding cis-regulatory modules in animals, *Nature Reviews Genetics*, 13, 7, 469–483

Herrmann, C., Van de Sande, B., Potier, D., and Aerts, S. (2012), i-cistarget: an integrative genomics method for the prediction of regulatory features and cis-regulatory modules, *Nucleic acids research*, 40, 15, e114–e114

Hozumi, A., Yoshida, R., Horie, T., Sakuma, T., Yamamoto, T., and Sasakura, Y. (2013), Enhancer activity sensitive to the orientation of the gene it regulates in the chordategenome, *Developmental Biology*, 375, 1, 79–91, doi:10.1016/j.ydbio.2012.12.012

Jing, H., Vakoc, C. R., Ying, L., Mandat, S., Wang, H., Zheng, X., et al. (2008), Exchange of gata factors mediates transitions in looped chromatin organization at a developmentally regulated gene locus, *Molecular cell*, 29, 2, 232–242

Kazemian, M., Brodsky, M. H., and Sinha, S. (2011), Genome surveyor 2.0: cis-regulatory analysis in drosophila, *Nucleic acids research*, 39, suppl 2, W79–W85

Kel-Margoulis, O. V., Kel, A. E., Reuter, I., Deineko, I. V., and Wingender, E. (2002), Transcompel®: a database on composite regulatory elements in eukaryotic genes, *Nucleic acids research*, 30, 1, 332–334

Kent, W. J., Sugnet, C. W., Furey, T. S., Roskin, K. M., Pringle, T. H., Zahler, A. M., et al. (2002), The human genome browser at UCSC, *Genome Research*, 12, 6, 996–1006

Kim, T. H. and Ren, B. (2006), Genome-wide analysis of protein-DNA interactions, *Annual Review of Genomics and Human Genetics*, 7, 1, 81–102

Klepper, K. and Drabløs, F. (2013), Motiflab: a tools and data integration workbench for motif discovery and regulatory sequence analysis, *BMC bioinformatics*, 14, 1, 9

Klepper, K., Sandve, G. K., Abul, O., Johansen, J., and Drablos, F. (2008), Assessment of composite motif discovery methods, *BMC bioinformatics*, 9, 1, 123

Kolbe, D., Taylor, J., Elnitski, L., Eswara, P., Li, J., Miller, W., et al. (2004), Regulatory potential scores from genome-wide three-way alignments of human, mouse, and rat, *Genome research*, 14, 4, 700–707

Kwon, A. T., Chou, A. Y., Arenillas, D. J., and Wasserman, W. W. (2011), Validation of skeletal muscle cis-regulatory module predictions reveals nucleotide composition bias in functional enhancers, *PLoS computational biology*, 7, 12, e1002256

Laimins, L. A., Gruss, P., Pozzatti, R., and Khoury, G. (1984), Characterization of enhancer elements in the long terminal repeat of moloney murine sarcoma virus., *Journal of virology*, 49, 1, 183–189

Li, L., Zhu, Q., He, X., Sinha, S., and Halfon, M. S. (2007), Large-scale analysis of transcriptional cis-regulatory modules reveals both common features and distinct subclasses, *Genome Biol*, 8, 6, R101

Maston, G. A., Evans, S. K., and Green, M. R. (2006), Transcriptional regulatory elements in the human genome, *Annu. Rev. Genomics Hum. Genet.*, 7, 29–59

Mathelier, A., Zhao, X., Zhang, A. W., Parcy, F., Worsley-Hunt, R., Arenillas, D. J., et al. (2014), Jaspar 2014: an extensively expanded and updated open-access database of transcription factor binding profiles, *Nucleic acids research*, 42, D1, D142–D147

Matys, V., Fricke, E., Geffers, R., Gößling, E., Haubrock, M., Hehl, R., et al. (2003), Transfac®: transcriptional regulation, from patterns to profiles, *Nucleic acids research*, 31, 1, 374–378

Meissner, A., Gnirke, A., Bell, G. W., Ramsahoye, B., Lander, E. S., and Jaenisch, R. (2005), Reduced representation bisulfite sequencing for comparative high-resolution dna methylation analysis, *Nucleic acids research*, 33, 18, 5868–5877

Meissner, A., Mikkelsen, T. S., Gu, H., Wernig, M., Hanna, J., Sivachenko, A., et al. (2008), Genome-scale DNA methylation maps of pluripotent and differentiated cells, *Nature*

Nam, J. and Davidson, E. H. (2012), Barcoded dna-tag reporters for multiplex cis-regulatory analysis, *PloS one*, 7, 4, e35934

Nikulova, A. A., Favorov, A. V., Sutormin, R. A., Makeev, V. J., and Mironov, A. A. (2012), Coreclust: identification of the conserved crm grammar together with prediction of gene regulation, *Nucleic acids research*, 40, 12, e93–e93

Ogbourne, S. and Antalis, T. (1998), Transcriptional control and the role of silencers in transcriptional regulation in eukaryotes, *Biochem. J*, 331, 1–14

Prestridge, D. S. (1995), Predicting pol II promoter sequences using transcription factor binding sites, *Journal of molecular biology*, 249, 5, 923–932

Ren, B., Robert, F., Wyrick, J. J., Aparicio, O., Jennings, E. G., Simon, I., et al. (2000), Genome-wide location and function of DNA binding proteins, *Science*, 290, 5500, 2306–2309

Robertson, K. D. and Wolffe, A. P. (2000), Dna methylation in health and disease, *Nature Reviews Genetics*, 1, 1, 11–19

Rosenbloom, K., Taylor, J., Schaeffer, S., Kent, J., Haussler, D., and Miller, W. (2008), Phylogenomic resources at the ucsc genome browser, in Phylogenomics (Springer), 133–144

Rosenbloom, K. R., Sloan, C. A., Malladi, V. S., Dreszer, T. R., Learned, K., Kirkup, V. M., et al. (2013), ENCODE data in the UCSC genome browser: year 5 update, *Nucleic Acids Research*, 41, D56–D63

Roth, C. L., Mastronardi, C., Lomniczi, A., Wright, H., Cabrera, R., Mungenast, A. E., et al. (2007), Expression of a tumor-related gene network increases in the mammalian hypothalamus at the time of female puberty, *Endocrinology*, 148, 11, 5147–5161

Rouault, H., Santolini, M., Schweisguth, F., and Hakim, V. (2014), Imogene: identification of motifs and cis-regulatory modules underlying gene co-regulation, *Nucleic acids research*, gku209

Sandelin, A., Carninci, P., Lenhard, B., Ponjavic, J., Hayashizaki, Y., and Hume, D. A. (2007), Mammalian rna polymerase ii core promoters: insights from genome-wide studies, *Nature Reviews Genetics*, 8, 6, 424–436

Sanyal, A., Lajoie, B. R., Jain, G., and Dekker, J. (2012), The long-range interaction landscape of gene promoters, *Nature*, 489, 7414, 109–113

Saxonov, S., Berg, P., and Brutlag, D. L. (2006), A genome-wide analysis of cpg dinucleotides in the human genome distinguishes two distinct classes of promoters, *Proceedings of the National Academy of Sciences of the United States of America*, 103, 5, 1412–1417

Siepel, A., Bejerano, G., Pedersen, J. S., Hinrichs, A. S., Hou, M., Rosenbloom, K., et al. (2005), Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes, *Genome research*, 15, 8, 1034–1050

Sinha, S. and He, X. (2007), MORPH: probabilistic alignment combined with hidden markov models of cis-regulatory modules, *PLoS computational biology*, 3, 11, e216

Sinha, S., Schroeder, M. D., Unnerstall, U., Gaul, U., and Siggia, E. D. (2004), Cross-species comparison significantly improves genome-wide prediction of cis-regulatory modules in drosophila, *BMC bioinformatics*, 5, 1, 129

Sinha, S., van Nimwegen, E., and Siggia, E. D. (2003), A probabilistic method to detect regulatory modules, *Bioinformatics*, 19, i292–i301

Stamatoyannopoulos, J. A., Snyder, M., Hardison, R., Ren, B., Gingeras, T., Gilbert, D. M., et al. (2012), An encyclopedia of mouse dna elements (mouse encode), *Genome biology*, 13, 8, 418

Su, J., Teichmann, S. A., and Down, T. A. (2010), Assessing computational methods of cis-regulatory module prediction, *PLoS Comput. Biol.*, 6, 12, e1001020

Taher, L., Smith, R. P., Kim, M. J., Ahituv, N., and Ovcharenko, I. (2013), Sequence signatures extracted from proximal promoters can be used to predict distal enhancers, *Genome biology*, 14, 10, R117

767 Tan, M., Luo, H., Lee, S., Jin, F., Yang, J. S., Montellier, E., et al. (2011), Identification of 67 histone
768 marks and histone lysine crotonylation as a new type of histone modification, *Cell*, 146, 6, 1016–1028
769 Thomas-Chollier, M., Sand, O., Turatsinze, J.-V., Defrance, M., Vervisch, E., Brohée, S., et al. (2008),
770 Rsat: regulatory sequence analysis tools, *Nucleic acids research*, 36, suppl 2, W119–W127
771 Thompson, J. A. and Congdon, C. B. (2014), Gami-crm: Using de novo motif inference to detect cis-
772 regulatory modules, in Proceedings of the 2014 IEEE Congress on Evolutionary Computation (IEEE)
773 Van Loo, P. and Marynen, P. (2009), Computational methods for the detection of cis-regulatory modules,
774 *Briefings in Bioinformatics*, 10, 5, 509–524, doi:10.1093/bib/bbp025
775 Wagner, A. (1999), Genes regulated cooperatively by one or more transcription factors and their
776 identification in whole eukaryotic genomes., *Bioinformatics*, 15, 10, 776–784
777 Walters, M. C., Fiering, S., Eidemiller, J., Magis, W., Groudine, M., and Martin, D. (1995), Enhancers
778 increase the probability but not the level of gene expression, *Proceedings of the National Academy of*
779 *Sciences*, 92, 15, 7125–7129
780 Wasserman, W. W. and Fickett, J. W. (1998), Identification of regulatory regions which confer muscle-
781 specific gene expression, *Journal of molecular biology*, 278, 1, 167–181
782 Won, K.-J., Zhang, X., Wang, T., Ding, B., Raha, D., Snyder, M., et al. (2013), Comparative annotation
783 of functional regions in the human genome using epigenomic data, *Nucleic acids research*, gkt143
784 Zentner, G. E. and Henikoff, S. (2013), Regulation of nucleosome dynamics by histone modifications,
785 *Nature structural & molecular biology*, 20, 3, 259–266

## FIGURES