

Genomic views of distant-acting enhancers

Axel Visel^{1,2}, Edward M. Rubin^{1,2} & Len A. Pennacchio^{1,2}

In contrast to protein-coding sequences, the significance of variation in non-coding DNA in human disease has been minimally explored. A great number of recent genome-wide association studies suggest that non-coding variation is a significant risk factor for common disorders, but the mechanisms by which this variation contributes to disease remain largely obscure. Distant-acting transcriptional enhancers — a major category of functional non-coding DNA — are involved in many developmental and disease-relevant processes. Genome-wide approaches to their discovery and functional characterization are now available and provide a growing knowledge base for the systematic exploration of their role in human biology and disease susceptibility.

Multiple lines of evidence indicate that important functional properties are embedded in the non-coding portion of the human genome, but identifying and defining these features remains a major challenge. An initial estimate of the magnitude of functional non-coding DNA was derived from comparative analysis of the first available mammalian genomes (human and mouse), which indicated that fewer than half of the evolutionary constrained sequences in the human genome encode proteins¹, a prospect that gained further support when additional vertebrate genomes became available for comparative genomic analyses².

The overall impact of these presumably functional non-coding sequences on human biology was initially unclear. A considerable urgency to define their locations and functions came from a growing number of known associations of non-coding sequence variants with common human diseases. Specifically, genome-wide association studies (GWAS) have revealed a large number of disease susceptibility regions that do not overlap protein-coding genes but rather map to non-coding intervals. For example, a 58-kilobase linkage disequilibrium block located at human chromosome 9p21 was shown to be reproducibly associated with an increased risk for coronary artery disease, yet the risk interval lies more than 60 kilobases away from the nearest known protein-coding gene^{3,4}. To estimate the global contribution of variation in non-coding sequences to phenotypic and disease traits, we performed a meta-analysis of ~1,200 single-nucleotide polymorphisms (SNPs) identified as the most significantly associated variants in GWAS published so far (ref. 5, accessed 2 March 2009). Using conservative parameters that tend to overestimate the size of linkage disequilibrium blocks, we found that in 40% of cases (472 of 1,170) no known exons overlap either the linked SNP or its associated haplotype block, suggesting that in more than one-third of cases non-coding sequence variation causally contributes to the traits under investigation.

One possibility that could explain these GWAS hits is that the non-coding intervals contain enhancers, a category of gene regulatory sequence that can act over long distances. A simplified view of the current understanding of the role of enhancers in regulating genes is summarized in Fig. 1. The docking of RNA polymerase II to proximal promoter sequences and transcription initiation are fairly well characterized; by contrast, the mechanisms by which insulator and silencer elements buffer or repress gene regulation, respectively, are less well understood⁶. Transcriptional enhancers are regulatory sequences that can be located upstream of, downstream of or within their target gene and can modulate expression independently of their orientation⁷. In vertebrates, enhancer sequences are thought to comprise densely clustered aggregations of transcription-factor-binding

sites⁸. When appropriate occupancy of transcription-factor-binding sites is achieved, recruitment of transcriptional coactivators and chromatin-remodelling proteins occurs. The resultant protein aggregates are thought to facilitate DNA looping and ultimately promoter-mediated gene activation (see page 199). In-depth studies of individual genes such as *APOE* or *NKX2-5* (reviewed in ref. 9) have shown that many genes are regulated by complex arrays of enhancers, each driving distinct aspects of the messenger RNA expression pattern. These modular properties of mammalian enhancers are also supported by their additive regulatory activities in heterologous recombination experiments¹⁰.

The purely genetic evidence from GWAS does not allow any direct inferences regarding the underlying molecular mechanisms, but a number of in-depth studies of individual loci (see below) suggest that variation in distant-acting enhancer sequences and the resultant changes in their activities can contribute to human disorders. Although we anticipate a variety of other non-coding functional categories such as negative gene regulators or non-coding RNAs to have a role in human disease, in this Review we focus on the role of enhancers and on strategies to define their location and function throughout the genome.

Enhancers in human disease

Beginning with the discovery that an inherited change in the β -globin gene alters one of the coded amino acids and thereby causes sickle-cell anaemia^{11,12}, thousands of mutations in the coding regions of genes have been identified to be responsible for monogenic disorders over the past half century. By contrast, the role of mutations not involving primary gene structural sequences has been minimally explored, largely owing to our inability to recognize relevant non-coding sequences, much less predict their function. The molecular genetic identification of individual enhancers involved in disease has been, in most cases, a painstaking and inefficient endeavour. Nevertheless, a number of successful studies have shown that distant-acting gene enhancers exist in the human genome and that variation in their sequences can contribute to disease. In this section, we discuss three examples in which enhancers were directly shown to play a role in human disease: thalassaemias resulting from deletions or rearrangements of β -globin gene (*HBB*) enhancers, preaxial polydactyly resulting from sonic hedgehog (*SHH*) limb-enhancer point mutations, and susceptibility to Hirschsprung's disease associated with a *RET* proto-oncogene enhancer variant.

The extensive studies of the human globin system and its role in haemoglobinopathies have historically served as a test bed for defining not only the role of coding sequences in disease^{11,12} but also that of non-coding

¹Genomics Division, MS 84-171, Lawrence Berkeley National Laboratory, Berkeley, California 94720, USA. ²US Department of Energy Joint Genome Institute, Walnut Creek, California 94598, USA.

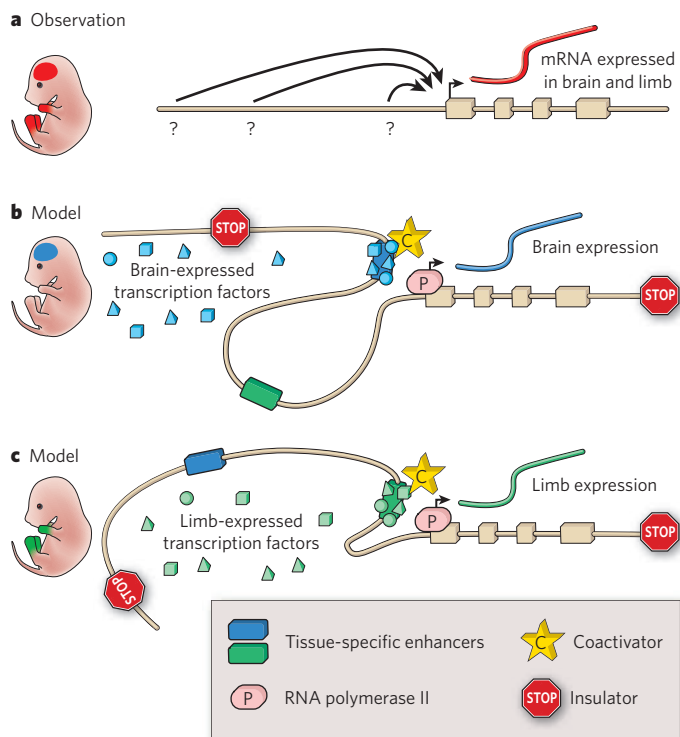


Figure 1 | Overview of gene regulation by distant-acting enhancers.

a, For many genes, the regulatory information embedded in the promoter is insufficient to drive the complex expression pattern observed at the messenger RNA level. For example, a gene could be expressed both in the brain and in the limbs during embryonic development (red), even if the promoter by itself is not active in either of these structures, suggesting that appropriate expression depends on additional sequences that are distant-acting and *cis*-regulatory. However, defining the genomic locations of such regulatory elements (question marks) and their activities in time and space (arrows) is a major challenge. **b, c**, Tissue-specific enhancers are thought to contain combinations of binding sites for different transcription factors. Only when all required transcription factors are present in a tissue does the enhancer become active: it binds to transcriptional coactivators, relocates into physical proximity with the gene promoter (through a looping mechanism) and activates transcription by RNA polymerase II. In any given tissue, only a subset of enhancers is active, as schematically shown in **b** and **c** for the example gene pictured in **a**, whose expression is controlled by two separate enhancers with brain-specific and limb-specific activities. Insulator elements prevent enhancer–promoter interactions and can thus restrict the activity of enhancers to defined chromatin domains. In addition to activation by enhancers, negative regulatory elements (including repressors and silencers) can contribute to transcriptional regulation (not shown).

sequences. The α -thalassaemias and β -thalassaemias are haemoglobinopathies resulting from imbalances in the ratio of α -globin to β -globin chains in red blood cells. The molecular basis of these conditions was initially elucidated in cases in which inactivation or deletion of globin structural genes could be readily identified¹³. However, although gene deletion or sequence changes resulting in a truncated or non-functional gene product explained some thalassaemia cases, for a subset of patients intensive sequencing efforts failed to reveal abnormalities in globin protein-coding sequences. Through extensive long-range mapping and sequencing of DNA from individuals diagnosed with thalassaemia but lacking globin coding mutations, it was eventually discovered that many of these globin chain imbalances were due to deletion or chromosome rearrangements that resulted in the repositioning of distant-acting enhancers required for normal globin gene expression^{14,15}. These early molecular genetic studies revealed a clear role for non-coding regulatory elements as a cause of human disorders through their impact on gene expression. Since then, many such examples of ‘position effects’, defined as changes in the expression of a gene when its location in a chromosome is changed, often by translocation, have been found¹⁶.

In addition to the pathological consequences of the removal or the repositioning of distant-acting enhancers, there are also examples of single-nucleotide changes within enhancer elements as a cause of human disorders. One example of this category of disease-causing non-coding mutation involves the limb-specific long-distance enhancer ZRS (also known as MFCS1) of *SHH* (Fig. 2). This enhancer is located at the extreme distance of approximately 1 megabase from *SHH*, within the intron of a neighbouring gene^{17,18}. Of interest is that, initially, the gene in which the enhancer resides was thought to be relevant for limb development and was therefore named limb region 1 (*LMBR1*)¹⁹. Facilitated by the functional knowledge of the ZRS enhancer from mouse studies, targeted resequencing screens of this enhancer in humans revealed that it is associated with preaxial polydactyly. Approximately a dozen different single-nucleotide variations in this regulatory element have been identified in humans with preaxial polydactyly and segregate with the limb abnormality in families^{18,20}. Studies of the impact of the human ZRS sequence changes have been carried out in transgenic mice, in which the single-nucleotide changes result in ectopic anterior-limb expression during development, consistent with preaxial digit outgrowth²¹. Furthermore, sequence changes in the orthologous enhancers were found in mice, as well as in cats, with preaxial polydactyly^{22,23}, and targeted deletion of the enhancer in mice caused truncation of limbs¹⁷. These studies illustrate the importance of first experimentally identifying distant-acting enhancers in allowing subsequent human genetic studies to explore the potential role of disease-causing mutation in functional non-coding sequences.

Another example of enhancer variation contributing to human disease is provided by the discovery of a common non-coding variant linked to susceptibility to Hirschsprung’s disease. Although multigenic, Hirschsprung’s disease risk is strongly linked to coding mutations in the *RET* proto-oncogene^{24,25}. However, family-based studies have also revealed evidence for Hirschsprung’s disease linked to the *RET* locus in people lacking any accompanying functional *RET* coding mutations. Through the use of multispecies comparisons of orthologous genomic intervals that include and flank *RET*, coupled with *in vitro* and *in vivo* functional studies, an enhancer sequence located in intron 1 of *RET* was identified and found to contain a common variant contributing more than a 20-fold increased risk for Hirschsprung’s disease than rarer alleles in this element^{26,27}. In transgenic mice, this enhancer was shown to be active in the nervous system and digestive tract during embryogenesis in a manner consistent with its putative role in Hirschsprung’s disease²⁷. It is interesting to note that although this enhancer variation is clearly important in disease risk, the variant alone is not sufficient to cause Hirschsprung’s disease, highlighting the complex aetiology of this disorder.

As is evident from these labour-intensive gene-centric studies, enhancers can, in principle, have an important role in disease, but it remains unclear whether these are rare exceptions or whether variation in enhancers contributes to disease on a pervasive scale. Support for the latter comes from a rapidly growing number of examples in which non-coding SNPs linked to disease traits through GWAS were found to affect the expression levels of nearby genes²⁸, suggesting that variation in regulatory sequences may commonly contribute to a wide range of disorders. The results of the recent GWAS, coupled with the role of gene regulation in normal human biology, provide a strong incentive for defining the distant-acting-enhancer architecture of the human genome.

Harnessing evolution

Gene-centric studies have been crucial to defining the general characteristics of gene regulatory regions in specific human disorders, but they have only identified and characterized a limited number of such elements. Systematic large-scale identification of sequences that are likely to be enhancers was first made possible by comparative genomic strategies. These approaches are based on the assumption that the sequences of gene regulatory elements, like those of protein-coding genes, are under negative evolutionary selection, because most changes in functional sequences have deleterious consequences^{29–32}. Thus, it was proposed that statistical measures of evolutionary sequence constraint would provide a way to

identify potential enhancer sequences within the vast amount of non-coding sequence in the human genome. Support for this approach initially came from retrospective comparative genomic analyses of experimentally well-defined enhancers; these analyses revealed that enhancers frequently shared sequence conservation with orthologous regions present in the genomes of other mammals. The observation that DNA conservation identified many of these complex regulatory elements encouraged investigators to move away from blind studies of regions flanking genes of interest towards focusing specifically on non-coding sequences constrained across vertebrate species, culminating in whole-genome studies in which conservation level alone guided experimentation^{32–34}.

Initially, comparisons over extreme evolutionary distances, such as between humans and fish, were deemed most effective for this purpose^{29,31}. Indeed, it was observed through large-scale transgenic mouse and fish studies that many of these non-coding sequences that had been conserved for hundreds of millions of years of evolution were enhancers that drove expression in highly specific anatomical structures during embryonic development. Likewise, so-called ultraconserved non-coding elements, which are blocks of 200 base pairs or more that are perfectly conserved between humans, mice and rats³⁵, were also found to be highly enriched in tissue-specific enhancers, suggesting that the success rate of comparative approaches for enhancer identification depends on scoring criteria, rather than just evolutionary distance³². This idea was further supported by the development of advanced statistical tools designed to quantify evolutionary constraint, from which it became evident that even comparisons between relatively closely related species can be effective predictors of enhancers^{2,36,37}. A large-scale transgenic mouse study that included nearly all non-exonic ultraconserved elements in the human genome revealed that whereas many of them are developmental *in vivo* enhancers, other conserved non-coding sequences that are under similar evolutionary constraint, but are not perfectly conserved between humans and mice, are equally enriched in enhancers³³. These results suggest that ultraconserved elements do not represent a functionally distinct subgroup of conserved non-coding sequences in terms of their enrichment in *in vivo* enhancers but rather that there is a much larger number of non-coding sequences that are under similar evolutionary constraint and are just as enriched in enhancers as are ultraconserved elements.

Independent of the specific algorithms and metrics that were used, most categories of conserved non-coding sequence were found not to be randomly distributed in the genome. Instead, they are located in a highly

biased manner near genes active during development^{2,33–35}, consistent with the observation that a large proportion of these non-coding sequences give robust positive signals in various assays of being tissue-specific *in vivo* enhancers active during development.

Comparative approaches are an effective high-throughput genomic strategy for identifying non-coding sequences that are highly likely to be enhancers, but they have several limitations. First, although conservation is indicative of function, it is not necessarily indicative of enhancer activity, because many other types of non-coding functional element that may have similar conservation signatures are known to exist. Second, even when conservation of non-coding DNA results from enhancer function, conservation cannot predict when and where an enhancer is active in the developing or adult organism. For all identified candidates, experimental studies are needed to decipher the gene-regulatory properties of each element, and these transgenic studies cannot feasibly be scaled to generate truly comprehensive genome-wide data sets.

A perplexing study questioning the importance of extremely conserved enhancers found the lack of an apparent phenotype upon targeted deletion of four independent ultraconserved elements in mice³⁸. General expectations were that non-coding sequences that have been perfectly conserved in mammals for tens of millions of years must be essential and that their deletion should result in severe phenotypes, comparable to those observed upon deletion of the *Shh* limb enhancer and other less well-conserved enhancers^{9,17}. However, mice with deletions of such ultraconserved enhancers were viable, fertile and showed no overt phenotype³⁸. Interpretations of this lack of obvious effect are similar to those of the absence of phenotypes upon deletion of highly conserved protein-coding genes: minor phenotypes may have escaped detection in the assays used; there may have been functional redundancy with other genes or enhancers; or there may have been reductions in fitness that only become apparent over multiple generations or are not easily detected in a controlled laboratory environment. This study highlighted that although extreme non-coding sequence conservation is an effective predictor of the location of enhancers in the genome, the degree of evolutionary constraint is not directly correlated with the severity of anticipated phenotypes.

Sequencing-based enhancer discovery

As a strategy complementary to comparative genomic methods, it has recently become possible to generate genome-wide maps of chromatin marks that can be used to identify the location of enhancers and other

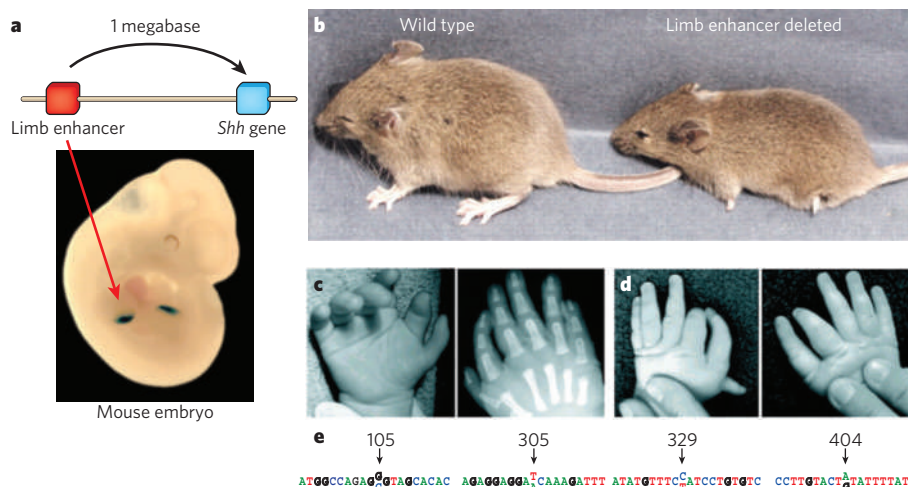


Figure 2 | Consequences of deletion and mutation of the limb enhancer of sonic hedgehog. **a**, The limb enhancer of *Shh* is located approximately 1 megabase away from its target promoter in the intron of a neighbouring gene (*Lmbr1*; exons not shown). In transgenic mouse reporter assays, this non-coding sequence targets gene expression to a posterior region of the developing limb bud (red arrow). (Image reproduced, with permission, from ref. 18.) **b**, Mice with a targeted deletion of this enhancer have severely truncated limbs, which strikingly demonstrates its functional importance

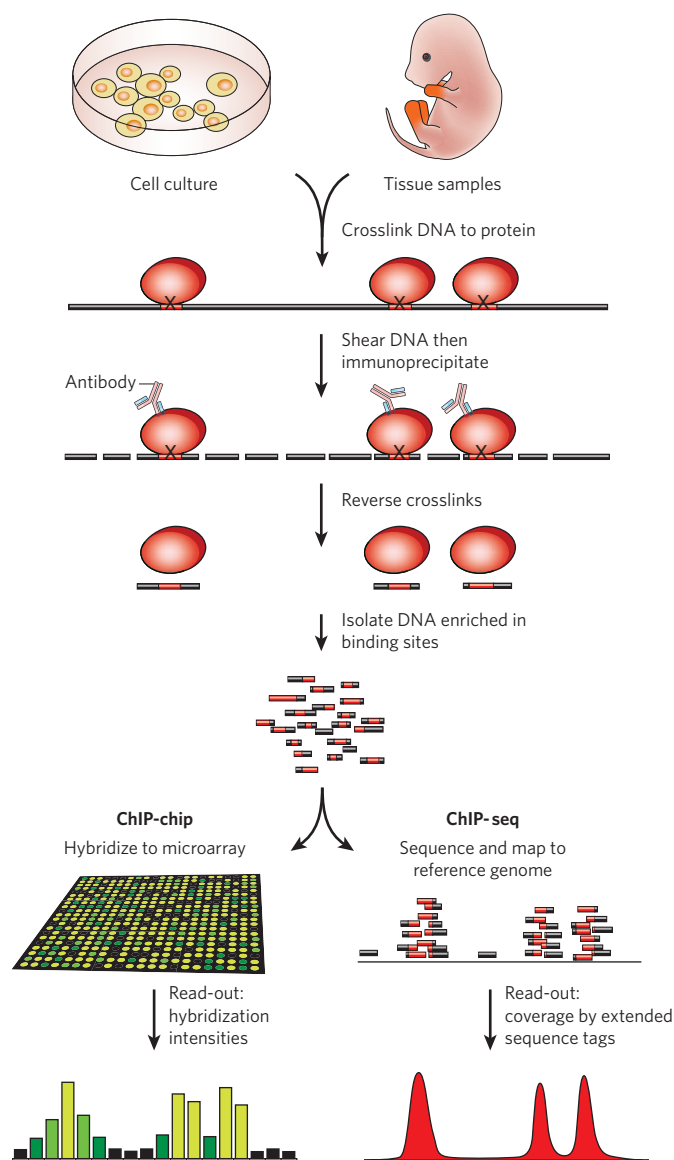
in development. (Reproduced, with permission, from ref. 17.) **c–e**, Point mutations in the orthologous human enhancer sequence result in preaxial polydactyly, emphasizing the potential significance of variation in non-coding functional sequences in both rare and common human disorders: **c** and **d** show the hands of two patients with point mutations in the *SHH* limb enhancer; **e** shows point mutations associated with preaxial polydactyly identified in four unrelated families. (Panels **c** and **d** reproduced, with permission, from ref. 18; panel **e** modified, with permission, from ref. 18.)

Box 1 | Mapping of regulatory elements using ChIP-chip and ChIP-seq

Formaldehyde crosslinking of DNA to proteins that bind to it directly or as part of larger complexes⁷¹, combined with subsequent immunoprecipitation targeting specific DNA-associated proteins (ChIP⁷²), was widely used in the pre-genomic era to study protein–DNA interactions directly in cultured cells or in tissue samples. The top portion of the figure shows a schematic overview of the individual steps involved. They include the molecular fixation of non-covalent protein–DNA interactions, shearing of the crosslinked chromatin, immunoprecipitation with an antibody binding the protein of interest and reversal of crosslinks. In many cases, antibodies that bind to covalently modified proteins are used, for example those that recognize methyl groups at defined amino-acid residues of histones. In the conventional ChIP approach, enrichment of the associated DNA fragments relative to non-immunoprecipitated (‘input’) DNA is quantified for individual proposed binding locations (not shown). This need for quantification at every site of interest initially thwarted the application of ChIP on a genomic scale.

The introduction of DNA microarrays allowed the hybridization-based interrogation of large numbers of potential binding sites in parallel (ChIP-on-chip, or ChIP-chip), thus making it possible to screen entire compact model-organism genomes^{73,74} or large vertebrate genome intervals⁷⁵ in a single experiment (see figure, bottom left). ChIP-chip was used on a massive scale in the Encyclopedia of DNA Elements (ENCODE) pilot project, in which dozens of proteins and protein modifications were initially mapped in a representative 1% portion of the human genome³⁹.

Recently, chromatin immunoprecipitation coupled to massively parallel sequencing (ChIP-seq) has become increasingly used as an alternative to ChIP-chip^{44–47}. The ChIP-seq method is very similar to the experimental set-up of ChIP-chip, except that, in the final step, next-generation sequencing techniques are used to determine the sequence of immunoprecipitated DNA fragments, which are then computationally mapped to the reference genome (see figure, bottom right). Improved sequencing technologies offer the possibility to obtain millions of mappable reads in a single ChIP-seq experiment at moderate cost. The results from ChIP-seq are based on statistical analysis of read counts, which overcomes many of the challenges associated with the quantification and normalization of hybridization signals, and an increasing number of advanced computational ChIP-seq analysis tools are becoming available⁷⁶. ChIP-seq analysis covers by default the entire mappable portion of the reference genome without the need to restrict the analysis to its subregions.



regulatory regions. These genomic approaches have become possible as a result of an improved understanding of the proteins and epigenetic marks found at particular categories of regulatory element, together with concurrently developed technologies that allow traditional chromatin immunoprecipitation (ChIP) techniques to be applied on the scale of whole vertebrate genomes. The initial in-depth studies of 1% of the genome in the Encyclopedia of DNA Elements (ENCODE) pilot project³⁹ were largely based on data sets generated by the ChIP-chip technique (Box 1) and revealed the molecular properties of a variety of regulatory elements.

With respect to enhancer identification, a particularly relevant insight was the identification of specific histone methylation signatures found at enhancers. In contrast to promoters, which are marked by trimethylation of histone H3 at lysine residue 4 (H3K4me3), active enhancers are marked by monomethylation at this position (H3K4me1)⁴⁰. Mapping these marks in the ENCODE regions and, more recently, throughout the entire genome⁴¹ revealed tens of thousands of elements that were predicted to be active enhancers in the examined cell types. Importantly, these predicted enhancers were also frequently associated with the transcriptional coactivators p300 and/or TRAP220 (also known as MED1), raising the possibility that such coactivators might be useful general markers for mapping enhancers. Although it was initially not

clear to what extent the presence of transcriptional coactivators such as p300 is indicative of active rather than inactive enhancers, comparison of DNase I hypersensitivity (a marker of open chromatin structure) in several cell lines throughout the ENCODE regions revealed that the location of cell-line-specific distal DNase-I-hypersensitivity sites correlates with cell-line-specific p300 binding at these sites, providing further support for the possibility that transcriptional coactivators, along with histone modification signatures, may be useful for the mapping of DNA elements with cell-specific and tissue-specific enhancer activities⁴².

Owing to the development of the ChIP-seq technique (Box 1), which has now superseded ChIP-chip as the method of choice for many applications, genome-wide maps for a considerable number of chromatin marks and transcription factors both in humans and mice have become available^{43–55}. These data sets allowed the identification of not only the H3K4me1 and H3K4me3 signatures discussed earlier but also additional chromatin marks present at predicted or validated enhancers, and provided a refined view of their correlation to enhancer activities^{44,51,55}. However, with very few exceptions (see, for example, refs 50 and 54) genome-wide mapping of these and other regulation-associated chromatin marks (Table 1) was done in immortalized cell lines, cultured stem cells or primary cell cultures. Thus, the maps of potentially enhancer-associated marks produced by these studies provided limited insight into

their *in vivo* distribution during embryonic development and in adult organs, most probably concealing the genomic location of enhancers that are inactive in these cells.

In a recent ChIP-seq study targeted at the prediction of enhancers that are active in a particular tissue during embryonic development, the transcriptional coactivator p300 was mapped in chromatin directly derived from embryonic mouse tissues, including the forebrain, the midbrain and the limb buds⁵⁶. Overall, several thousand p300 peaks were identified from these three tissues, with the vast majority of genome regions only being significantly enriched in one of the three tissues and located in non-coding regions distal from known promoters. Transgenic mouse experiments with almost 100 of these sequences revealed that they are developmental enhancers in almost all cases. More importantly, the tissue-specific occupancy by p300 as identified by ChIP-seq could in most cases also accurately predict the *in vivo* patterns of expression driven by these enhancers, providing an important advantage over comparative genomic methods for enhancer identification. The study also showed global enrichment in tissue-specific p300 peaks near genes that are expressed in the same tissue, again consistent with the proposed function of these genomic regions as active transcriptional enhancers.

These experimentally predicted genome-wide sets of *in vivo* enhancers also made it possible to address the controversial issue of the extent to which evolutionary conservation is a hallmark of *in vivo* enhancers⁵⁷. Several studies have shown that highly conserved non-coding elements are enriched in developmental *in vivo* enhancers^{32–34}. However, some observations have challenged such a generalized correlation between sequence conservation and enhancer activity: experimental analysis of individual loci suggested that a large proportion of enhancers cannot be detected by comparative genomics⁵⁸; the molecular marks of a surprisingly large proportion of sequences in the ENCODE regions suggested that regulatory functions are not, or are only weakly, conserved⁵⁹; and histone methylation present at orthologous loci in humans and mice did not correlate with overall increased levels of sequence conservation⁵⁹. In contrast to these findings, approximately 90% of the tissue-specific p300 peaks identified by ChIP-seq in developing mouse tissues overlapped regions that are under detectable evolutionary constraint⁵⁶. There may be variation in the degree of evolutionary constraint of enhancers that are active in different types of cell or developing tissue, but these data suggest that developmental enhancers that can be identified through p300 binding are commonly evolutionarily constrained.

Although preliminary, the selected studies reviewed here highlight the clear potential of mapping various chromatin marks for identifying and predicting the activity of transcriptional enhancers on a genome-wide scale. The continued progress in throughput increase and the cost reductions of next-generation sequencing technologies offer an increasingly powerful genome-wide means of identifying specific DNA–protein interactions. We anticipate that high-resolution genome-wide *in vivo* maps of chromatin marks will become available for comprehensive series of developing and adult tissues in normal states, as well as diseased states, providing multilayered *in vivo* annotations of the non-coding portion of our genome. It is important to realize that, despite this expected progress, we will continue to need parallel *in vitro* and *in vivo* biological studies to understand the functions associated with chromatin marks and to study conclusively the mechanisms by which sequence variation in distant-acting enhancers contributes to disease.

Defining the targets

The methods described here have considerably improved our ability to identify enhancers and their associated activity patterns on a genomic scale, but a remaining important challenge is to determine the relationships between enhancers and genes. Comparing ChIP-chip or ChIP-seq data with transcriptome data from microarrays or RNA-seq⁶⁰ can provide highly suggestive clues to the identity of the target gene of a given enhancer in a given tissue, but such comparisons do not provide the direct evidence for enhancer–promoter interactions that would be desirable in mapping tissue-specific regulatory networks on a genomic scale.

Early circumstantial evidence suggested that long-distance regulation of genes by enhancers occurs through the formation of physical chromatin loops, but it only became possible to study such interactions systematically through the introduction of the chromosome conformation capture (3C) assay and its derivative technologies⁶¹. Similar to ChIP, the 3C approach relies on formaldehyde crosslinking to capture DNA–DNA interactions directly in intact cells or cell nuclei. Previously suggested pairs of interacting sites are subsequently tested and validated one by one through the quantification of crosslinking events. In one of many examples demonstrating the utility of 3C in the analysis of distant-acting vertebrate enhancers, this technique was recently used⁶² to study chromatin interactions at the *Shh* locus, whose role in limb development was discussed in detail earlier. Using the 3C technique, it was demonstrated that the limb-specific long-range enhancer located in an intron of the *Lmbr1* gene directly interacts with increased frequency with the *Shh* promoter in limb buds but not in other tissues tested, providing important mechanistic support for its proposed role in *Shh* gene regulation in limb development. As an alternative approach to 3C, RNA tagging and recovery of associated proteins (RNA TRAP) can also be used to establish physical proximity between distal non-coding sequences and actively transcribed genes; this was first demonstrated in the mouse β -globin gene locus⁶³.

This work and other gene-centric studies (for more examples, see refs 64 and 65) were critical in shaping our understanding of enhancer–promoter interactions. However, they have the fundamental limitation that only one or very few previously proposed interactions between specific loci can be assayed per experiment. This limitation was partly overcome through the use of microarrays to analyse entire 3C libraries (called chromosome conformation capture-on-chip⁶⁶ and circular chromosome conformation capture⁶⁷, both known as 4C). By applying this approach to fetal liver and brain, it was demonstrated that the β -globin gene locus control region (LCR) makes reproducible tissue-specific contacts with other loci predominantly located on the same chromosome but in some cases dozens of megabases away from the LCR⁶⁶. Of possible relevance to the adoption of this approach for enhancer discovery is that reproducible interactions with other chromosome regions were also observed in the brain, where the LCR is thought to be inactive.

The 4C approaches are a significant improvement, but they still preclude the generation of truly genome-wide interaction networks because each experiment only reveals the genome-wide interactions of a single site of interest. This problem is partly alleviated by the chromosome conformation capture carbon copy (5C) method⁶⁸, in which a complex 3C library generated through multiplexed PCR is analysed by large-scale sequencing to generate a comprehensive ‘many-to-many’ interaction map of DNA–DNA interactions. However, owing to the need for specific

Table 1 | Selected major categories of non-coding functional element

Category	Function	Selected associated chromatin marks*
Promoter	Region that is located immediately upstream of a protein-coding gene, and binds to RNA polymerase II; where transcription is initiated	RNA polymerase II ⁴⁴ , H3K4me3 (ref. 40) (active promoters)
Enhancer	Region that activates transcription, often in a temporally and spatially restricted manner, by acting on a promoter. Enhancers can be located far from target promoters and are orientation independent	p300 (refs 40, 56), H3K4me1 (ref. 40)
Insulator	Separates active from inactive chromatin domains and interferes with enhancer activity when placed between an enhancer and promoter	CTCF ^{44,53}
Repressor/silencer	Negative regulators of gene expression	REST ⁴⁵ , SUZ12 (refs 69, 70)

*Many additional chromatin marks were found to correlate with one or several of these categories of regulatory element. Detailed descriptions of these markers and their respective binding characteristics at different types of regulatory sequence element can be found in refs 40, 41, 44, 51 and 55.

primers for each possible interacting fragment and the sequencing depth required for analysis of the resultant libraries, the application of 5C has so far been restricted to the in-depth analysis of single loci or chromosome regions.

As an alternative genome-wide approach, antibody-based methods might be used to restrict the analysis space in which DNA–DNA interactions are studied to a size that can be affordably analysed using currently available sequencing technologies. One possibility is to couple a chromatin–interaction paired-end tag (ChIA–PET) sequencing strategy to a ChIP step that enriches for chromatin fragments bound to a specific transcription factor or other chromatin mark of interest⁶⁴. Although the technical feasibility of this approach remains to be demonstrated, it has remarkable potential for enhancer discovery. This is because its application to general enhancer-associated marks such as p300 or histone methylation^{40,56} might identify, in a single step, enhancers active in a tissue of interest, as well as their respective target genes.

Perspective

Genetic and medical resequencing studies have been advanced by knowledge about the structure of protein-coding genes and a detailed understanding of the relationship between mRNA sequences and the primary structures of the proteins they encode. Through such studies, disease links have been established for a sizeable proportion of the ~20,000 protein-coding genes in the human genome. By contrast, a very limited number of changes in gene regulatory sequences have so far been linked to human disease. Consequently, an important motivation for functionally annotating the non-coding portion of the human genome and the *cis*-regulatory elements that it contains is to assess the relationship between variations in non-coding sequences and human disease. In the absence of genome-wide catalogues of functionally annotated regulatory elements, how these elements impact on human biology, as well as disease, will remain an untested hypothesis.

Despite advances in relevant technologies, functionally characterizing the distant-acting-enhancer architecture of the human genome in its entirety will be an enormous undertaking, owing to the great number of data points needed, which include dozens of tissues and cell types, as well as developmental states and possibly disease states.

A further challenge will be to link distant-acting enhancers to the genes they regulate. Linking enhancers to their cognate gene will allow the further assignment of these functional sequences to their basic ‘gene’ unit of heredity, for collective resequencing analysis.

Although we have focused on distant-acting enhancers here, there are other categories of functional element in the non-coding portion of the genome (for example insulators, negative regulators, promoters and non-coding RNAs), and they will also be crucial targets for large-scale identification and characterization. It is expected that technologies similar to those described here for enhancers will make it possible to explore their roles in human biology and disease. ■

1. Mouse Genome Sequencing Consortium. Initial sequencing and comparative analysis of the mouse genome. *Nature* **420**, 520–562 (2002).
2. Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* **15**, 1034–1050 (2005).
3. Helgadottir, A. *et al.* A common variant on chromosome 9p21 affects the risk of myocardial infarction. *Science* **316**, 1491–1493 (2007).
4. McPherson, R. *et al.* A common allele on chromosome 9 associated with coronary heart disease. *Science* **316**, 1488–1491 (2007).
5. Hindorf, L. A., Junkins, H. A., Mehta, J. P. & Manolio, T. A. A catalog of published genome-wide association studies. *OPG: Catalog Published Genome-Wide Assoc. Studies* <<http://www.genome.gov/gwastudies>> (2009).
6. Maston, G. A., Evans, S. K. & Green, M. R. Transcriptional regulatory elements in the human genome. *Annu. Rev. Genomics Hum. Genet.* **7**, 29–59 (2006).
This paper is a comprehensive overview of functional classes of gene regulatory sequence, including many disease-relevant examples identified through gene-centric studies.
7. Banerji, J., Rusconi, S. & Schaffner, W. Expression of a β -globin gene is enhanced by remote SV40 DNA sequences. *Cell* **27**, 299–308 (1981).
8. Panne, D. The enhanceosome. *Curr. Opin. Struct. Biol.* **18**, 236–242 (2008).
9. Visel, A., Bristow, J. & Pennacchio, L. A. Enhancer identification through comparative genomics. *Semin. Cell Dev. Biol.* **18**, 140–152 (2007).
10. Visel, A. *et al.* Functional autonomy of distant-acting human enhancers. *Genomics* **93**, 509–513 (2009).
11. Ingram, V. M. Gene mutations in human haemoglobin: the chemical difference between normal and sickle cell haemoglobin. *Nature* **180**, 326–328 (1957).
12. Pauling, L. *et al.* Sickle cell anemia, a molecular disease. *Science* **110**, 543–548 (1949).
13. Kan, Y. W. *et al.* Deletion of α -globin genes in haemoglobin-H disease demonstrates multiple α -globin structural loci. *Nature* **255**, 255–256 (1975).
14. Kioussis, D., Vanin, E., deLange, T., Flavell, R. A. & Grosfeld, F. G. β -Globin gene inactivation by DNA translocation in $\gamma\beta$ -thalassaemia. *Nature* **306**, 662–666 (1983).
15. Semenza, G. L. *et al.* The silent carrier allele: β thalassaemia without a mutation in the β -globin gene or its immediate flanking regions. *Cell* **39**, 123–128 (1984).
16. Kleinjan, D. A. & Lettice, L. A. Long-range gene control and genetic disease. *Adv. Genet.* **61**, 339–388 (2008).
17. Sagai, T., Hosoya, M., Mizushima, Y., Tamura, M. & Shiroishi, T. Elimination of a long-range *cis*-regulatory module causes complete loss of limb-specific *Shh* expression and truncation of the mouse limb. *Development* **132**, 797–803 (2005).
This paper shows that deletion of the distant-acting limb enhancer of the *Shh* gene in mice causes severe limb truncation, providing a model example of the requirement for enhancers in mammalian development.
18. Lettice, L. A. *et al.* A long-range *Shh* enhancer regulates expression in the developing limb and fin and is associated with preaxial polydactyly. *Hum. Mol. Genet.* **12**, 1725–1735 (2003).
19. Clark, R. M., Marker, P. C. & Kingsley, D. M. A novel candidate gene for mouse and human preaxial polydactyly with altered expression in limbs of *Hemimelic extra-toes* mutant mice. *Genomics* **67**, 19–27 (2000).
20. Furniss, D. *et al.* A variant in the sonic hedgehog regulatory sequence (ZRS) is associated with triphalangeal thumb and deregulates expression in the developing limb. *Hum. Mol. Genet.* **17**, 2417–2423 (2008).
21. Masuya, H. *et al.* A series of ENU-induced single-base substitutions in a long-range *cis*-element altering Sonic hedgehog expression in the developing mouse limb bud. *Genomics* **89**, 207–214 (2007).
22. Lettice, L. A., Hill, A. E., Devenney, P. S. & Hill, R. E. Point mutations in a distant sonic hedgehog *cis*-regulator generate a variable regulatory output responsible for preaxial polydactyly. *Hum. Mol. Genet.* **17**, 978–985 (2008).
23. Lettice, L. A. *et al.* Disruption of a long-range *cis*-acting regulator for *Shh* causes preaxial polydactyly. *Proc. Natl Acad. Sci. USA* **99**, 7548–7553 (2002).
24. Bolk, S. *et al.* A human model for multigenic inheritance: phenotypic expression in Hirschsprung disease requires both the *RET* gene and a new 9q31 locus. *Proc. Natl Acad. Sci. USA* **97**, 268–273 (2000).
25. Gabriel, S. B. *et al.* Segregation at three loci explains familial and population risk in Hirschsprung disease. *Nature Genet.* **31**, 89–93 (2002).
26. Emison, E. S. *et al.* A common sex-dependent mutation in a *RET* enhancer underlies Hirschsprung disease risk. *Nature* **434**, 857–863 (2005).
27. Grice, E. A., Rochelle, E. S., Green, E. D., Chakravarti, A. & McCallion, A. S. Evaluation of the *RET* regulatory landscape reveals the biological relevance of a HSCR-implicated enhancer. *Hum. Mol. Genet.* **14**, 3837–3845 (2005).
28. Cookson, W., Liang, L., Abecasis, G., Moffatt, M. & Lathrop, M. Mapping complex disease traits with global gene expression. *Nature Rev. Genet.* **10**, 184–194 (2009).
29. Aparicio, S. *et al.* Detecting conserved regulatory elements with the model genome of the Japanese puffer fish, *Fugu rubripes*. *Proc. Natl Acad. Sci. USA* **92**, 1684–1688 (1995).
30. Loots, G. G. *et al.* Identification of a coordinate regulator of interleukins 4, 13 and 5 by cross-species sequence comparisons. *Science* **288**, 136–140 (2000).
31. Nobrega, M. A., Ovcharenko, I., Afzal, V. & Rubin, E. M. Scanning human gene deserts for long-range enhancers. *Science* **302**, 413 (2003).
32. Pennacchio, L. A. *et al.* In vivo enhancer analysis of human conserved non-coding sequences. *Nature* **444**, 499–502 (2006).
33. Visel, A. *et al.* Ultraconservation identifies a small subset of extremely constrained developmental enhancers. *Nature Genet.* **40**, 158–160 (2008).
34. Woolfe, A. *et al.* Highly conserved non-coding sequences are associated with vertebrate development. *PLoS Biol.* **3**, e7 (2005).
35. Bejerano, G. *et al.* Ultraconserved elements in the human genome. *Science* **304**, 1321–1325 (2004).
36. Prabhakar, S. *et al.* Close sequence comparisons are sufficient to identify human *cis*-regulatory elements. *Genome Res.* **16**, 855–863 (2006).
37. Cooper, G. M. *et al.* Distribution and intensity of constraint in mammalian genomic sequence. *Genome Res.* **15**, 901–913 (2005).
38. Ahituv, N. *et al.* Deletion of ultraconserved elements yields viable mice. *PLoS Biol.* **5**, e234 (2007).
This paper shows that deletion of several ultraconserved non-coding sequences in mice may not result in obvious phenotypes, demonstrating that even extreme evolutionary constraint does not necessarily indicate that a non-coding sequence is required for viability.
39. The ENCODE Project Consortium. Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature* **447**, 799–816 (2007).
40. Heintzman, N. D. *et al.* Distinct and predictive chromatin signatures of transcriptional promoters and enhancers in the human genome. *Nature Genet.* **39**, 311–318 (2007).
This paper identifies a histone H3K4 differential methylation signature that distinguishes promoters from enhancers, providing a chromatin-based tool for genome-wide enhancer prediction.
41. Heintzman, N. D. *et al.* Histone modifications at human enhancers reflect global cell-type-specific gene expression. *Nature* **459**, 108–112 (2009).
42. Xi, H. *et al.* Identification and characterization of cell type-specific and ubiquitous chromatin regulatory structures in the human genome. *PLoS Genet.* **3**, e136 (2007).
43. Wei, C. L. *et al.* A global map of p53 transcription-factor binding sites in the human genome. *Cell* **124**, 207–219 (2006).
This paper describes mapping of protein–DNA interactions by ChIP coupled with conventional capillary-based sequencing of concatenated paired-end tags (ChIP–PET), a conceptual predecessor of the ChIP–seq approach.
44. Barski, A. *et al.* High-resolution profiling of histone methylations in the human genome. *Cell* **129**, 823–837 (2007).
45. Johnson, D. S., Mortazavi, A., Myers, R. M. & Wold, B. Genome-wide mapping of *in vivo* protein–DNA interactions. *Science* **316**, 1497–1502 (2007).

46. Robertson, G. *et al.* Genome-wide profiles of STAT1 DNA association using chromatin immunoprecipitation and massively parallel sequencing. *Nature Methods* **4**, 651–657 (2007).
This paper is one of several independently published early ChIP-seq studies validating the method for genome-wide mapping of transcription-factor-binding sites.
47. Mikkelsen, T. S. *et al.* Genome-wide maps of chromatin state in pluripotent and lineage-committed cells. *Nature* **448**, 553–560 (2007).
This paper is one of several independently published early ChIP-seq studies providing some of the first genome-wide data sets of several histone modifications in different mouse cell types and examining their correlation with functional genome features.
48. Zhao, X. D. *et al.* Whole-genome mapping of histone H3 Lys4 and 27 trimethylations reveals distinct genomic compartments in human embryonic stem cells. *Cell Stem Cell* **1**, 286–298 (2007).
49. Chen, X. *et al.* Integration of external signaling pathways with the core transcriptional network in embryonic stem cells. *Cell* **133**, 1106–1117 (2008).
50. Wederell, E. D. *et al.* Global analysis of *in vivo* Foxa2-binding sites in mouse adult liver using massively parallel sequencing. *Nucleic Acids Res.* **36**, 4549–4564 (2008).
51. Robertson, A. G. *et al.* Genome-wide relationship between histone H3 lysine 4 mono- and tri-methylation and transcription factor binding. *Genome Res.* **18**, 1906–1917 (2008).
52. Ku, M. *et al.* Genomewide analysis of PRC1 and PRC2 occupancy identifies two classes of bivalent domains. *PLoS Genet.* **4**, e1000242 (2008).
53. Cuddapah, S. *et al.* Global analysis of the insulator binding protein CTCF in chromatin barrier regions reveals demarcation of active and repressive domains. *Genome Res.* **19**, 24–32 (2009).
54. Gao, N. *et al.* Dynamic regulation of *Pdx1* enhancers by Foxa1 and Foxa2 is essential for pancreas development. *Genes Dev.* **22**, 3435–3448 (2008).
55. Wang, Z. *et al.* Combinatorial patterns of histone acetylations and methylations in the human genome. *Nature Genet.* **40**, 897–903 (2008).
56. Visel, A. *et al.* ChIP-seq accurately predicts tissue-specific activity of enhancers. *Nature* **457**, 854–858 (2009).
57. Cooper, G. M. & Brown, C. D. Qualifying the relationship between sequence conservation and molecular function. *Genome Res.* **18**, 201–205 (2008).
58. McGaughey, D. M. *et al.* Metrics of sequence constraint overlook regulatory sequences in an exhaustive analysis at *phox2b*. *Genome Res.* **18**, 252–260 (2008).
59. Bernstein, B. E. *et al.* Genomic maps and comparative analysis of histone modifications in human and mouse. *Cell* **120**, 169–181 (2005).
60. Wang, Z., Gerstein, M. & Snyder, M. RNA-Seq: a revolutionary tool for transcriptomics. *Nature Rev. Genet.* **10**, 57–63 (2009).
61. Dekker, J., Rippe, K., Dekker, M. & Kleckner, N. Capturing chromosome conformation. *Science* **295**, 1306–1311 (2002).
62. Amano, T. *et al.* Chromosomal dynamics at the *Shh* locus: limb bud-specific differential regulation of competence and active transcription. *Dev. Cell* **16**, 47–57 (2009).
63. Carter, D., Chakalova, L., Osborne, C. S., Dai, Y. F. & Fraser, P. Long-range chromatin regulatory interactions *in vivo*. *Nature Genet.* **32**, 623–626 (2002).
64. Fullwood, M. J., Wei, C. L., Liu, E. T. & Ruan, Y. Next-generation DNA sequencing of paired-end tags (PET) for transcriptome and genome analyses. *Genome Res.* **19**, 521–532 (2009).
65. Miele, A. & Dekker, J. Long-range chromosomal interactions and gene regulation. *Mol. Biosyst.* **4**, 1046–1057 (2008).
66. Simonis, M. *et al.* Nuclear organization of active and inactive chromatin domains uncovered by chromosome conformation capture-on-chip (4C). *Nature Genet.* **38**, 1348–1354 (2006).
67. Zhao, Z. *et al.* Circular chromosome conformation capture (4C) uncovers extensive networks of epigenetically regulated intra- and interchromosomal interactions. *Nature Genet.* **38**, 1341–1347 (2006).
68. Dostie, J. *et al.* Chromosome conformation capture carbon copy (5C): a massively parallel solution for mapping interactions between genomic elements. *Genome Res.* **16**, 1299–1309 (2006).
69. Lee, T. I. *et al.* Control of developmental regulators by Polycomb in human embryonic stem cells. *Cell* **125**, 301–313 (2006).
70. Squazzo, S. L. *et al.* Suz12 binds to silenced regions of the genome in a cell-type-specific manner. *Genome Res.* **16**, 890–900 (2006).
71. Van Lente, F., Jackson, J. F. & Weintraub, H. Identification of specific crosslinked histones after treatment of chromatin with formaldehyde. *Cell* **5**, 45–50 (1975).
72. Solomon, M. J., Larsen, P. L. & Varshavsky, A. Mapping protein–DNA interactions *in vivo* with formaldehyde: evidence that histone H4 is retained on a highly transcribed gene. *Cell* **53**, 937–947 (1988).
73. Ren, B. *et al.* Genome-wide location and function of DNA binding proteins. *Science* **290**, 2306–2309 (2000).
74. Iyer, V. R. *et al.* Genomic binding sites of the yeast cell-cycle transcription factors SBF and MBF. *Nature* **409**, 533–538 (2001).
75. Horak, C. E. *et al.* GATA-1 binding sites mapped in the β -globin locus by using mammalian chip-chip analysis. *Proc. Natl Acad. Sci. USA* **99**, 2924–2929 (2002).
76. Barski, A. & Zhao, K. Genomic location analysis by ChIP-seq. *J. Cell. Biochem.* **107**, 11–18 (2009).

Acknowledgements We thank M. Blow, S. Deutsch and A. Sczyrba for help with computational analysis of GWAS data and C. Attanasio for comments. L.A.P. and E.M.R. were supported by the Berkeley Program for Genomic Applications (funded by the US National Heart, Lung, and Blood Institute), and the Director, Office of Science, Office of Basic Energy Sciences, US Department of Energy, under contract number DE-AC02-05CH11231. L.A.P. was also supported by the US National Human Genome Research Institute.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Correspondence should be addressed to L.A.P. (lapennacchio@lbl.gov).



Review

Enhancer identification through comparative genomics

Axel Visel^b, James Bristow^{a,b}, Len A. Pennacchio^{a,b,*}

^a U.S. Department of Energy Joint Genome Institute, Walnut Creek, CA 94598, USA

^b Genomics Division, MS 84-171, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA

Abstract

With the availability of genomic sequence from numerous vertebrates, a paradigm shift has occurred in the identification of distant-acting gene regulatory elements. In contrast to traditional gene-centric studies in which investigators randomly scanned genomic fragments that flank genes of interest in functional assays, the modern approach begins electronically with publicly available comparative sequence datasets that provide investigators with prioritized lists of putative functional sequences based on their evolutionary conservation. However, although a large number of tools and resources are now available, application of comparative genomic approaches remains far from trivial. In particular, it requires users to dynamically consider the species and methods for comparison depending on the specific biological question under investigation. While there is currently no single general rule to this end, it is clear that when applied appropriately, comparative genomic approaches exponentially increase our power in generating biological hypotheses for subsequent experimental testing. It is anticipated that cardiac-related genes and the identification of their distant-acting transcriptional enhancers are particularly poised to benefit from these modern capabilities.

Published by Elsevier Ltd.

Keywords: *Cis*-regulatory; Comparative genomics; Enhancer; Review; Transgenic

Contents

1. Introduction	00
2. Role of non-coding sequences in development and human disease	00
2.1. Modularity of transcriptional regulation by enhancers	00
2.2. Spatiotemporal precision of developmental enhancers	00
2.3. Enhancers are required for vertebrate development	00
2.4. Enhancers contribute to human disease	00
2.5. Challenges	00
3. Enhancer identification by comparative genomic strategies	00
3.1. Pre-genome-scale comparative approaches	00
3.2. Using genomic data in comparative approaches	00
3.2.1. Deep comparisons: human–fish	00
3.2.2. Extreme conservation within mammals	00
3.2.3. Comparison of close species: phylogenetic shadowing	00
4. Tools and resources for comparative genomics	00
4.1. Identification of candidate regions at a genomic scale	00
4.1.1. Aligning genome sequences	00
4.1.2. Scoring conservation in aligned genome sequences	00
4.2. Experimental validation of <i>cis</i> -regulatory elements	00
4.3. Enhancer browser: large-scale data set of in vivo-validated enhancers	00

* Corresponding author at: Genomics Division, One Cyclotron Road, MS 84-171, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA.
Tel.: +1 510 486 7498; fax: +1 510 486 4229.
E-mail address: LAPennacchio@lbl.gov (L.A. Pennacchio).

4.3.1. Experimental data	00
4.3.2. Computational data set	00
5. Conclusions and perspectives	00
Acknowledgements	00
References	00

1. Introduction

One of the most intriguing features of biology is the identical DNA content across all cells within an organism and yet the ability of this genetic information to dictate the enormous cellular diversity within the body. Rather, cell type complexity arises predominantly from vast temporal and spatial differences in gene expression during development. The principal mechanism underlying this gene expression diversity across cell types is dynamic gene regulation induced by a variety of interacting transcription factors which are also encoded by our genome and subject to tight regulation [1–3]. Transcription factors recognize specific target sequences located within gene promoters and/or more distant acting *cis*-regulatory regions, and function to either enhance or repress a given gene's cellular expression. Through this highly orchestrated process, higher organisms have been able to evolve beyond the limitations of unicellularity to create complex forms and functions, including the development of the cardiovascular system.

Insights into this complexity are beginning to emerge for the human genome with the availability of a complete genomic sequence template [4,5]. This starting point has led to the identification of the ~25,000 genes in the human genome, albeit work remains to be done in deciphering all of their functions. Gene identification was greatly facilitated by having access to protein sequence databases and “expressed sequence tags” where computational algorithms for gene identification could subsequently be built based upon knowledge gained from these experimental datasets. In contrast, the availability of the human genome sequence alone provided no additional clues as to the precise locations of distant-acting gene enhancers. Challenges included the large non-coding search space in the human genome (~98% of 3×10^9 bp), the small size and degenerate nature of transcription factor binding sites, and most importantly the lack of experimental training sets for computational methods to identify such sequences in a global manner. The recent determination of additional genome sequences from other vertebrates has proven to be powerful at identifying the location of candidate distant-acting *cis*-regulatory elements based on their evolutionary conservation across appropriately distanced species.

In this review, we describe the use of comparative genomics as an increasingly powerful strategy for sequence-based enhancer identification. In particular, we provide an overview of selected computational tools and resources that are useful for the identification of enhancers involved in development and/or specific gene function. We end by highlighting the challenges arising from the identification of large numbers of putative enhancers through comparative genomics and the need to develop high throughput functional assays to determine their spatiotemporal *in vivo* activity at a genomic scale.

2. Role of non-coding sequences in development and human disease

Traditionally, most studies of the genetic networks underlying vertebrate development have focused on the proteins that are involved, since they are – compared to regulatory sequences – generally easier to identify and more readily accessible to a variety of experimental methods. However, these proteins are generally limited to functional activity only in tissues where they are expressed, thereby stressing the importance of understanding the intricacies of gene regulation to comprehend regulatory networks in their entirety. In this section, we provide a brief overview of insights gained from gene-centric in-depth studies. While the list of examples described here is by no means exhaustive, it illustrates some of the major properties and characteristics of distant-acting *cis*-regulatory elements and exemplifies their important role in vertebrate development and human disease.

2.1. Modularity of transcriptional regulation by enhancers

A characteristic feature of enhancers is the modular mode by which they regulate gene expression. One of many insightful examples for these properties can be obtained by examination of the human apolipoprotein E (*APOE*) locus. At least six distinct sequence elements flanking this gene control different aspects of *APOE* expression. Namely, the enhancement of kidney expression has been ascribed to the promoter [6], while elements located downstream of the gene include two liver-specific enhancers [7,8], a skin enhancer [6,9], two multiple tissue enhancers directing gene expression to adipocytes, macrophages and brain astrocytes [9,10], and a distal brain-specific enhancer [11]. It is worth noting that each of these discrete elements are on the order of several hundred basepairs in length and are scattered across 42 kb. A second example where the modularity of transcriptional regulation has been experimentally studied in great detail is the cardiac homeobox gene *Nkx2-5* (*Csx*). This gene is required for heart development [12] and series of deletions and transgenic reporter experiments were used to dissect both its proximal and distal regulatory regions [13–18]. These studies revealed that at least five distinct elements target *Nkx2-5* gene expression to specific sub-regions of the developing heart as well as to non-cardiac tissues and it has been suggested that this regulatory complexity played a important role in the evolution of the multi-chambered mammalian heart [19]. Thus, modular transcriptional regulation appears to be a common mechanism of complex gene regulation and a number of gene-centric studies beyond the selected examples of *APOE* and *Nkx2-5* have further supported the concept that the complex expression patterns of genes across tissues regularly arise from the combined activity of multiple elements.

2.2. Spatiotemporal precision of developmental enhancers

Another remarkable feature of enhancers is the high spatiotemporal precision with which they regulate gene expression. One example of the tight restriction of the timing and tissue-specificity of enhancer activity during embryonic development is the *Hoxd11* locus. Deletion of a single *Hoxd11* regulatory element in mice delays expression of both *Hoxd10* and *Hoxd11* during somitogenesis, but at later stages normal expression of *Hoxd10* and *Hoxd11* is restored [20]. It is hypothesized that this partial gene expression rescue is mediated by complementary regulatory elements present in this region. Since only a subset of anatomical regions lack *Hoxd11* expression temporally, this gene regulatory deletion results in vertebral patterning and specification defects but of lesser severity than complete *Hoxd11* gene knockouts.

The *Hoxd11* locus thus demonstrates how a single enhancer regulates a relatively subtle, yet functionally important spatiotemporal sub-aspect of the expression pattern of a key developmental gene. The general picture emerging from this and other similar gene-centric studies is that the high spatiotemporal precision of single enhancers – in combination with their modular mode of action – has allowed complex gene expression patterns to evolve. This is particularly the case for many developmentally important genes, whose expression patterns appear to be frequently the result of the orchestrated activity of several different enhancers with distinct spatiotemporal activity patterns. Importantly, these single elements tend to be more restricted in their tissue specificity than the mRNA expression patterns to which they contribute, providing researchers with reagents for tissue-specific targeting of gene expression.

2.3. Enhancers are required for vertebrate development

Like mutations in the protein-coding portion of genes, deletions or mutations of regulatory elements can result in developmental defects, such as in the *Hoxd11* locus (see Section 2.2). Another example from the Hox gene family is the 200 bp “early enhancer” (EE) of the *Hoxc8* gene. Deletion of this enhancer results in delayed expression of the Hoxc8 protein and in skeletal defects that recapitulate aspects of the *Hoxc8*^{-/-} phenotype [21], demonstrating that this regulatory element is required for normal embryogenesis. As a third example, deletion of three brain-specific enhancers of *Otx2* [22,23] revealed that they are required for maintaining normal expression levels of *Otx2* in the developing brain. While deletion of these enhancers did not result in obvious phenotypes, compound heterozygous embryos in which one *Otx2* allele was null and the other allele was an *Otx2* enhancer deletion displayed defects in brain development. These results support that while each of these elements is not absolutely required for viability, they play an important role in embryonic development through their coordinated and quantitative effects on gene expression.

Of note, defects resulting from deletion or mutation of regulatory elements are usually restricted to the tissue in which they drive expression. This property can be exploited to study gene functions that are otherwise difficult to assess experimentally.

For example, the role of *Hand2* in craniofacial development cannot be studied by targeted deletion of the gene itself because *Hand2*^{-/-} embryos die from cardiac abnormalities before the differentiation of craniofacial features. However, deletion of a branchial arch-specific *Hand2* enhancer in mice results in craniofacial defects including cleft palate and mandibular hypoplasia, demonstrating a role both for this enhancer and the *Hand2* gene in craniofacial development [24]. These studies allowed for the dissection of the regulatory architecture of this locus through the separate assessment of the roles of this gene in cardiac and craniofacial development. Another important possibility arising from the identification of tissue-specific enhancers is the possibility to use them to drive the expression of Cre recombinase. Such constructs can be used to generate tissue-specific knockouts by introducing flanking LoxP sites to the gene of interest [25]. For example, the conditional Cre/Lox-mediated deletion of *Mef2c* using a myocardial-specific enhancer has been used to examine the role of *Mef2c* beyond developmental stages at which mice with a complete deletion of *Mef2c* die from cardiovascular defects [26]. Thus, even in cases where the deletion of an enhancer is insufficient to abolish gene expression in a particular tissue, the enhancer can be used to study the function of the respective gene in a tissue-specific manner.

Indeed, many enhancers do not cause an overt phenotype beyond changes in expression levels of the target gene when experimentally deleted in mice. Examples include tissue- or cell type-specific enhancers for *Engrailed2* [27], *Fgf4* [28], *Gata1* [29] or *MyoD* [30]. An obvious explanation for the frequent absence of phenotypes in enhancer deletion experiments is that often only one aspect of a complex endogenous mRNA expression pattern is affected, while expression of the gene in other tissues or at other stages is maintained. This higher spatiotemporal restriction is therefore expected to result in generally milder effects than deletion of entire genes. A second explanation is functional redundancy, which might be more common among regulatory elements than it is among protein-coding genes. While being sufficient to drive expression in reporter assays, many enhancers could be dispensable for normal development and physiology because their function is complemented by other regulatory elements with similar tissue specificity. Such redundancy of regulatory elements has, for instance, been directly shown for the TCR- γ locus, where a deletion of two enhancers results in severe reduction in γ - δ -thymocytes, whereas single deletion of either element did not cause a major immunological phenotype [31]. Functional redundancy does not imply that these enhancers are functionally less important and that their deletion does not reduce reproductive fitness. Rather it indicates that many enhancers are involved in fine-tuning gene expression. These findings also raise the possibility that functional redundancies are a factor in the comparative studies described below, since they might result in reduced evolutionary conservation of such elements.

2.4. Enhancers contribute to human disease

As a result of our limited knowledge about the location of most enhancers in the genome, the contribution of distant acting

mutations to human disease has so far not been explored on a large scale. One of the few known examples is the limb-specific ZRS long-distance enhancer of *Sonic hedgehog* (*SHH*). This element is located at the extreme distance of one megabase from the gene it regulates, residing in the intron of a neighboring gene. Genetic lesions affecting this element cause polydactyly both in human individuals and in mutant mouse strains, demonstrating the crucial role of enhancers during mammalian development [32]. Elimination of the conserved intronic region in which this enhancer is embedded results in severe limb truncations in mice, strongly supporting human disease studies [33]. Even point mutations in this regulatory element cause human preaxial polydactyly [34], offering an explanation why many enhancers are highly constrained and therefore often conserved across long evolutionary distances. While hundreds of regulatory mutations contributing to human disease have been reported [35], most of them affect promoter regions whose precise location is known for many human genes. It is expected that with growing numbers of identified human enhancers it will become possible to target systematic screens increasingly for regulatory mutations in this distant-acting class of gene regulatory elements.

2.5. Challenges

The selected examples above highlight the important role of enhancers in development and disease. However, it must be emphasized that the vast majority of distant-acting regulatory sequences in the mammalian genome has so far not been experimentally characterized either in vitro or in vivo and their overall contribution to human disease remains unclear. Two major challenges have rendered large-scale studies of developmental enhancers difficult. First, the absence of suitable prediction methods continues to present a major obstacle for identifying the location of these elements, especially for those that act over long distances. Second, the limited number of known developmental enhancers has largely prevented prediction by computational analysis because no suitable training sets of enhancers characterized by standardized experimental methods have been available. In consequence, our understanding of the sequence features involved in enhancer function remains limited to gene-centric studies and single elements. In the next sections, we will describe recent efforts to tackle both of these problems. Namely, recently developed methods and computational tools for comparative genomics have significantly improved our ability to identify the location of putative enhancers in the human genome and provide a starting point for large-scale experimental characterization of enhancers.

3. Enhancer identification by comparative genomic strategies

Cross-species sequence comparisons were shown to be an efficient approach to identify putative functional regions in non-coding DNA even before whole genome sequences of humans and other vertebrates became available. Many variations on this theme have been presented, including variation of the species being compared and different comparison methods, yet they

all rely on the same basic principle that functionally relevant sequences are under purifying selection, whereas non-functional regions are subject to genetic drift and become increasingly different between species with increasing phylogenetic distance. As a result, functional sequences generally stand out as more “conserved” than non-functional sequences when genomic sequences of different species are compared. Sequence conservation between different species can thus be used to identify putative functional regions, and many of these will be *cis*-regulatory elements.

3.1. Pre-genome-scale comparative approaches

Bottom-up approaches provided the early foundation for the utility of cross-species comparisons for the identification of *cis*-regulatory elements in the genomic sequence of a gene of interest (for early examples, see references [36,37]). In the absence of publicly available whole-genome sequence data and specialized computational tools for these purposes, this strategy usually included cloning and sequencing of orthologous non-coding sequences from two or more organisms, manual alignment and identification of conserved regions at the nucleotide level, often focusing on transcription factor binding sites. In reference to experimentally exploring these sequences through DNase footprinting, such approaches became known as “phylogenetic footprinting”.

Such gene-centric studies provided an important proof of principle, but the hypothesis that sequence conservation is a universal predictor of non-coding regulatory sequences was difficult to verify conclusively in the absence of sequence data for genome-wide comparisons. Thus, the prospect of genome-wide comparative identification of *cis*-regulatory regions was early recognized as an important motivation to sequence the genomes of the mouse and other vertebrates in addition to the human genome [38,39].

3.2. Using genomic data in comparative approaches

Even before sufficient sequence data for whole-genome comparisons became available, the merits of comparative approaches for enhancer identification were confirmed in studies that involved the sequencing of large genomic intervals. For example, Götting et al. [40] sequenced a 320 kb interval of the stem cell leukemia (SCL) locus in human, mouse and chicken to identify regulatory candidate regions. A subset of these regions corresponded to known regulatory elements and functional testing of previously uncharacterized conservation peaks led to the discovery of a new neural enhancer in the SCL locus. In another study, Loots et al. [41] identified multiple non-coding elements regulating the human interleukin-4, -5, and -13 genes by sequencing and aligning one megabase of human chromosome 5 and the orthologous mouse genome region. These results lent further support to the notion that conservation of non-coding sequences can be used to predict functional regions including regulatory elements in genomic sequence data.

The publication of the mouse and the pufferfish genomes in 2002 marked the kick-off for genome-wide comparative

approaches since they allowed for the first time systematic large-scale comparisons of the human with non-human vertebrate genomes [42,43]. Comparative analysis of the human and mouse genomes was particularly productive because their size is similar, 90% of these genomes are organized in syntenic blocks in which the respective order of genes is maintained, and in an initial analysis 40% of the two genomes were found to be alignable at the nucleotide level. Interestingly, while only ~1.5% of the human and mouse genome encode proteins, ~5% of these mammalian genomes were estimated to be under purifying selection, suggesting that much more than protein-encoding functions are constrained within our genome [43]. However, a multitude of functions can potentially be embedded into non-protein-coding DNA, including activating and repressing regulatory binding sites, known and unknown functional RNA types, and structural chromatin features. Most of these cannot be reliably predicted by existing computational methods; therefore, the functional relevance of constrained non-coding regions remained initially obscure.

Subsequent functional testing of such conserved regions revealed, however, that one of the predominant functions of constrained non-coding DNA seems in fact to be the tissue-specific spatial and temporal regulation of gene expression. One of the likely reasons for this is the large size of many enhancer sequences, conserved over hundreds of basepairs, which makes it possible to identify them through whole genome comparisons. In what follows, we provide an overview of comparative strategies that have so far been successfully used to find such *cis*-regulatory elements (for a more detailed discussion of general considerations regarding comparisons over different evolutionary distances, including the advantages and limitations of distant and close comparisons, see reference [44]).

3.2.1. Deep comparisons: human–fish

In the pre-genomic era, studies focusing on single genes suggested that distant evolutionary comparison could be useful to identify regulatory regions involved in core aspects of vertebrate development. For example over 10 years ago, Aparicio et al. [45] used comparisons between mouse and pufferfish (*Takifugu rubripes*) to identify functional regulatory elements in the *Hoxb4* locus based on non-coding conservation. These and other results demonstrated that deep comparisons are an efficient tool for enhancer prediction, but genome-wide application was not possible at the time since none of these vertebrate genome sequences were available.

A more recent study systematically exploited the remarkable potential of such distant vertebrate sequence comparisons to identify gene enhancers at the scale of larger genomic intervals [46]. In this work, the gene-sparse regions surrounding the human *DACH* locus were scanned for sequences that are not only highly conserved among mammals, but also had considerable sequence conservation in *Xenopus* as well as in pufferfish. Using an *in vivo* enhancer assay, these extremely conserved regions were found to be highly enriched for enhancers that drive tissue-specific gene transcription during embryogenesis. In fact, many of the conserved elements that are currently being tested in a large-scale transgenic *in vivo* screen in our

laboratory (see Section 4.3) were identified using human–fish conservation.

There are, however, several important limitations to distant comparative approaches. First, their high specificity is accompanied by moderate sensitivity. Depending on the alignment method, the comparative strategy, and the stringency of the applied filters, previously reported numbers of conserved non-coding elements identified by human–fish comparisons vary between 1400 [47] and 5700 [48]. Compared to estimates of the total number of protein-coding genes in the human genome [49], this is up to an order of magnitude lower, suggesting that many regulatory regions are missed by such distant comparisons. Second, to aggravate this problem, many elements with such extremely deep conservation occur in clusters around genes implicated in transcriptional regulation and development (*trans-dev* genes). For example, 85% of the 1400 human–fish CNSs described by Woolfe et al. [47] are found in clusters of five or more elements. In total, only 165 distinct clusters were identified and 93% of these clusters are associated with *trans-dev* genes. In contrast, the majority of genes with other functions are not associated with any deeply conserved elements, despite modular regulation of gene expression in time and space. Third, extremely distant comparisons are expected to identify predominantly regulatory elements that are involved in molecular, developmental or physiological mechanisms that exist in both species under consideration, thereby explaining why they are anciently conserved. Human–fish comparisons would therefore be of limited utility for studies of enhancers that are involved in mammalian-specific developmental processes. As an example, we performed comparative analysis retrospectively on a subset of heart-specific *cis*-regulatory sequences originally identified through functional studies. These elements drive gene expression in the anterior heart field, a transient developmental structure, and heart regions derived from it [50]. The vast majority lacked conservation outside of mammals, which may be partially due to differences in heart development between mammals and non-mammalian vertebrates (Fig. 1B).

3.2.2. Extreme conservation within mammals

If conventional comparative criteria such as 70% identity over at least 100 bp are used, human–rodent comparisons are of limited use for identification of enhancer elements. This is due to the fact that these two species share a relatively short divergence time since their last common ancestor which results in their high overall similarity even in non-functional genome regions. This results in the identification of an excess of elements as illustrated by the observation that ~40% of the human and mouse genome are alignable, yet only ~5% of the human genome are estimated to be under purifying selection [43]. In consequence, using human–mouse comparisons with relatively relaxed percent identity parameters for enhancer prediction is very sensitive, but results in a false-positive rate that is too high to be useful for most applications [58,59].

While an obvious solution is to seek more distant species for human genome comparison, this problem can be partially overcome by using more stringent conservation criteria in human–rodent comparisons alone. Human–rodent “ultracon-

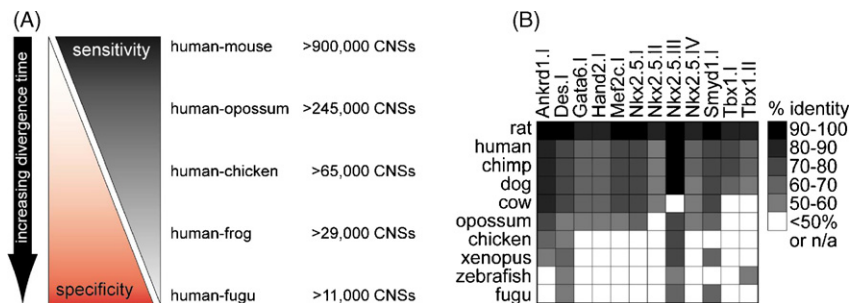


Fig. 1. Trade offs in comparative genomics of non-coding DNA based on different phylogenetic distances. (A) With simple definitions of CNSs, conservation depth can be used to calibrate specificity vs. sensitivity in comparative enhancer prediction. Closer sequence comparisons such as human–mouse provide a significant amount of non-coding conservation which provides strong sensitivity to identify known putative function, but at the cost of poor specificities. In contrast, human–fish comparison yields relatively little non-coding conservation and hence poor sensitivity to identify putative function, but with strong specificities for those conserved elements it does identify. (B) Known heart enhancers lack deep sequence conservation. In this illustrative example, retrospective comparative analysis of 12 known heart-specific *cis*-regulatory elements in 11 vertebrate genomes reveals limited sensitivity of deep comparisons for detecting mammalian heart-specific enhancers (% identity refers to mouse as the base genome). Most of these elements are only minimally conserved beyond mammals and would have been missed by human–fish comparisons. These data indicate that biological context is an important factor for comparative-based approaches, though on occasion heart enhancers are anciently conserved to fish. For detailed description and experimental characterization of these elements, see references [13,15,18,51–57].

served” elements are one such class of extremely conserved human–rodent sequences and are defined as sequences of 200 bp or more that are 100% identical between human, mouse and rat [60]. Thus, these sequences are at the extreme end of the conserved human–mouse continuum which is exemplified by there only being approximately 250 of such elements that do not overlap with protein-coding sequences in our genome. The function of these elements has not been exhaustively explored, but studies of single ultraconserved elements [46,61] as well as their genomic localization in clusters near key developmental genes [62] suggest that many of them may be long-range modulators of gene transcription.

While ultraconserved elements are highly likely to be enhancers or other functional elements, their value for large-scale prediction of enhancers is limited because they represent only a relatively small subset of the functionally conserved sequences in the human genome. Their low total number indicates a poor sensitivity, suggesting that many or most functional elements will be missed if ultraconservation alone is used to screen a genomic interval of interest. Moreover, because of the extreme conservation criteria of ultraconserved elements, most of them coincide with regions that are also conserved between human and fish. However, it has recently been suggested that statistically more rigorous methods than the original concept of ultraconservation might provide a way to extract larger populations with ultra-like constraints from human–rodent comparisons, increasing the sensitivity while maintaining the specificity associated with ultraconserved elements [48] (see Section 4.1.2). Computational tools to exploit this concept are becoming increasingly available [48,63,64].

3.2.3. Comparison of close species: phylogenetic shadowing

For studying regulatory elements related to aspects of biology that are specific to humans or primates, but do not exist in more distant species such as rodents, distant comparisons will only be useful in cases where previously existing regulatory features have assumed a new function in the primate

lineage. However, distant comparisons will miss elements that have evolved more recently and are possibly specific to the primate phylogenetic branch. On the other hand, comparison with other primates does not yield useful results when conventional sequence comparison is performed due to the relatively short period since the last common ancestor in the primate branch, e.g. ~25 million years for humans and Old World monkeys [65]. This is exemplified to a severe degree in comparisons of human and chimpanzee, which separated from their common ancestor ~7 million years ago. Between these two genomes ~99% of all nucleotides are conserved [66], rendering conventional comparative approaches useless because virtually all regions of the genome appear highly similar. This problem can be overcome using a “phylogenetic shadowing” approach [67]. In this method, the sequences of multiple, evolutionary close species such as humans, apes and monkeys are aligned. This depth of several species provides the nucleotide diversity that would otherwise be achieved through more distant pair-wise comparisons such as human–mouse. Moreover, this approach incorporates a molecular phylogenetic model to consider the phylogenetic relationships among the different species that are compared such that changes that occurred in a closely related species are given more power than those in a more distantly related species. Phylogenetic shadowing requires aligned sequences from multiple closely related species and has therefore so far only been used in the context of studies focusing on particular loci of interest [67,68]. However, this method will likely become increasingly used for the identification of regulatory elements as more and more closely related genomes become available [69].

4. Tools and resources for comparative genomics

A number of tools are available to identify conserved non-coding elements in genome sequences. In this section, we will provide an overview of computational approaches and web-based resources to interrogate and browse the human genome for such elements and retrieve their sequences for experimental studies. We also discuss approaches for experimental

Table 1
Selected interactive genome browsing tools for the identification of vertebrate CNSs

Identification of conserved elements	Available at	URL	Based on alignment	Display/download
Percent identity plot (PiP) [71,72]	Vista Genome Browser [73]	http://pipeline.lbl.gov	SLAGAN (pair-wise, glocal ^a) [74]	Percent identity curves; display and download of elements with adjustable threshold identity percentage
	Dcode ECR Browser [75]	http://ecrbrowser.dcode.org	BLASTZ (pair-wise, local) [76]	Percent identity plots or curves; display and download of elements with adjustable threshold identity percentage
PhastCons [64]	UCSC Genome Browser [77,78]	http://genome.ucsc.edu	MULTIZ (multiple, local) [79]	UCSC Genome Browser “Most Conserved” track; download of elements with adjustable constraint threshold
Gumby [48]	Vista Genome Browser [73]	http://pipeline.lbl.gov	SLAGAN (pair-wise, glocal ^a) [74]	“RankVista” <i>p</i> -value bar plots; display and download of elements with adjustable threshold <i>p</i> -value
	Vista Enhancer Browser	http://enhancer.lbl.gov	MLAGAN (multiple, global) [70]	Browsable list of human–mouse–rat CNSs; direct link to developmental enhancer assay results where available

^a “Glocal”, global alignments allowing local rearrangements.

characterization of developmental enhancers and describe the Vista Enhancer Browser as a public database of experimentally validated enhancers. Relevant web addresses and references describing each of the listed resources are provided in Table 1.

4.1. Identification of candidate regions at a genomic scale

Identification of conserved elements by comparison of genomes from different species is generally a two-step process. First, homologous regions of two or more different genomes are aligned at the nucleotide level, so that for each nucleotide position in the reference genome a best fit with the nucleotide at the respective position in the other genome(s) is determined. Second, based on this alignment, the different genomes are compared at the nucleotide level and statistical methods are used to identify regions where the sequence is more constrained (i.e. similar between the different organisms) than what would be expected for neutrally evolving DNA.

4.1.1. Aligning genome sequences

For the alignment step, a range of whole genome methods has been developed and several relevant programs are listed in Table 1. These generally fall into two categories: local and global alignment approaches. Local methods compare relatively short intervals of genomic sequences with each other and return the best match between two genomes for each sub-region. However, because they do not take into account the region surrounding these matches, they can result in false hits, e.g. returning a paralogous sequence instead of the true ortholog. In contrast, global methods align entire syntenic regions and are less prone to return false-positive matches, but fail to recognize homologous regions that have been locally rearranged by translocations or inversions. Finally, “glocal” alignment [70] is a global alignment strategy that allows for local rearrangements, thereby eliminating some of the problems associated with local-only or global-only alignments.

While all three types of alignments have been successfully used for comparative identification of functional elements, it is important to keep in mind that they will often return slightly different results for a particular genome region of interest. Thus, trial and error approaches are appropriate to maximize the likelihood of biological discovery.

4.1.2. Scoring conservation in aligned genome sequences

For defining highly conserved elements in aligned genomes, there is also a range of computational tools available. We focus here on a small subset of such tools that is of particular relevance for the identification of candidate enhancer sequences in the human genome by biomedical investigators (Fig. 2). The most straightforward way to identify highly constrained elements in genome alignments are pair-wise percent identity plots. When using local alignment methods such as BLASTZ [76], the length and percent identity of each aligned segment can be directly converted into a sequence plot [71] (Fig. 2A). Alternatively, for two globally aligned sequences, a sliding window of user-defined size (e.g. 100 bp) is moved along the alignment and returns for each nucleotide position the percentage of identity within the window [72] (Fig. 2B). Conserved non-coding sequences (CNSs) are in both cases defined by a user-specified threshold, e.g. as regions exceeding 70% identity over at least 80 bp.

Percent identity plots have been widely used because the concept is simple and readily implemented, but they have several important limitations. For example, they do not allow direct multi-species comparisons, but rather multiple species can be indirectly considered by aligning the pair-wise alignments to the same reference genome. Moreover, they do not take into account the evolutionary distance between the species that are being compared. When using the same threshold (e.g. 70% identity, ≥ 100 bp), the choice of the species being compared can be used to roughly calibrate sensitivity versus specificity (Fig. 1A). For instance, CNSs identified by comparison of distant species such as human–fish are highly enriched in functional enhancers [46]. However, the relatively small number of such elements

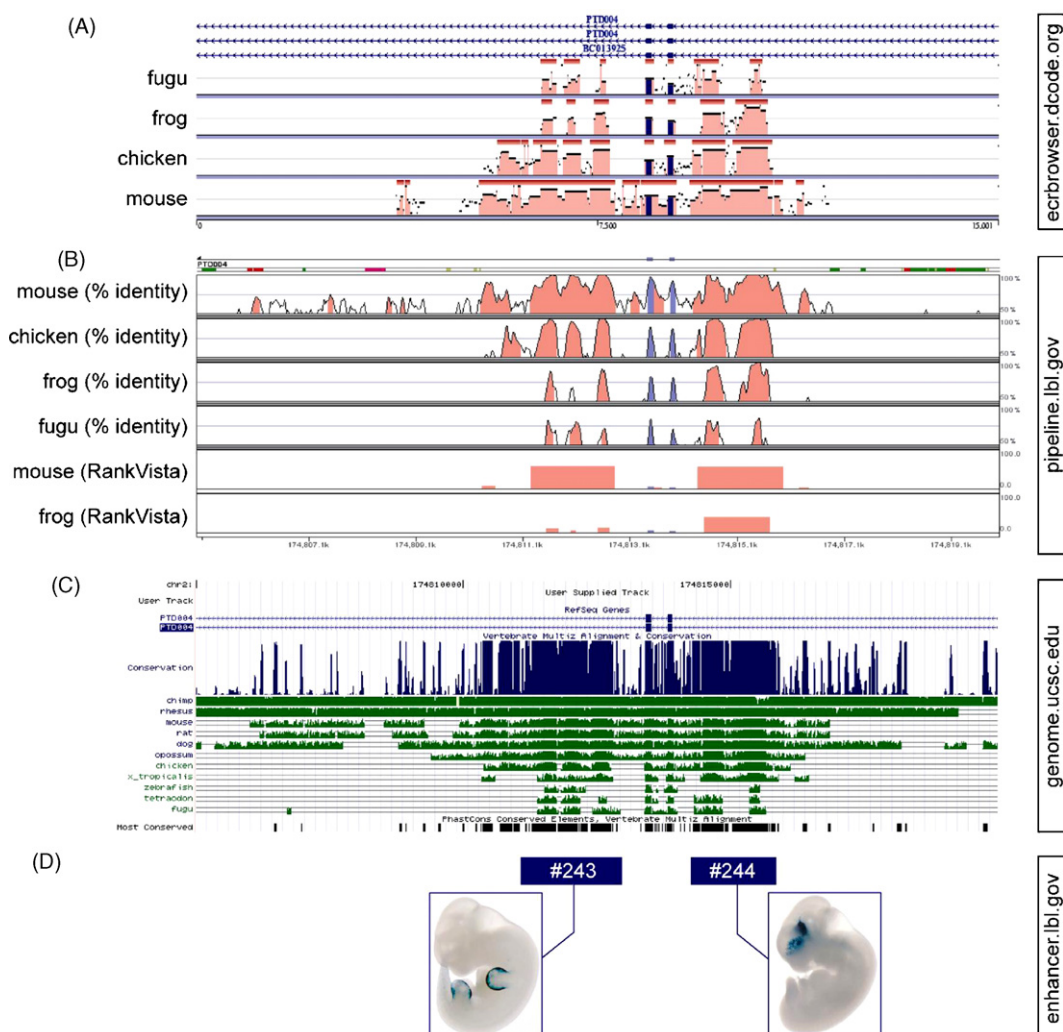


Fig. 2. Sequence display of the same human genome region by various tools for comparative analysis. A 15 kb region comprising two exons of the GTP-binding protein *PTD004* is shown (chr2: 174,805,000–174,820,000; hg17). (A) Percent identity plots as displayed in the Dcode ECR Browser. (B) Percent identity tracks and RankVista tracks in the Vista Genome Browser. RankVista tracks are based on *p*-values of conserved elements determined by the Gumbly algorithm. (C) Conservation and PhastCons ("Most Conserved") tracks in the UCSC Genome Browser. (D) Experimental results for two CNSs in the Vista Enhancer Browser. See Table 1 for relevant references.

detected by this strategy indicates that it fails to capture many functional sequences (see Section 3.2.1). In contrast, comparison of close species such as human–mouse identifies hundreds of thousands of elements and is thus more sensitive, but suffers from a high false-positive rate when such elements are tested for their tissue-specific enhancer activity in functional assays [58]. The problem of low specificity in percent-identity types of comparisons between close species can be partially alleviated by using more stringent threshold parameters. For example, human–mouse–rat "ultra"-conservation of 100% for ≥ 200 bp [60] is similarly successful for enhancer identification as deep human–fish conservation (AV, LAP, unpublished observations), but is even less sensitive by an order of magnitude (see Section 3.2.2).

Recently a new generation of advanced, mathematically and statistically rigorous tools have become available that allow direct multi-species (*n*-way) comparisons while also considering phylogenetic branch length and local neutral background substi-

tution rates [48,64]. Importantly, these methods do not require a single pre-specified evolutionary distance [64] (Fig. 2C) and provide high specificity even in pair-wise comparisons of relatively close species such as human and mouse [48] (Fig. 2B). Moreover, they use statistical tests to assign quantitative scores to elements, allowing a user to rank all elements within a given genomic interval according to the significance of their constraint. We have started to explore the relative value of these different comparative methods for prediction of tissue-specific enhancers by testing elements predicted by different methods in a transgenic reporter assay (see below), where we find that these more advanced comparative tools are indeed superior to simple percent identity plots in their ability to predict functional enhancers.

In order to browse the human or other vertebrate genomes for the presence of elements identified using the different methods described above, a variety of public resources is available online. We provide a list of such sites in Table 1, limiting our

selection to those resources that provide pre-aligned sequences and elements identified by the methods described above.

4.2. Experimental validation of *cis*-regulatory elements

An array of experimental approaches is available to assess the potential for putative regulatory elements to influence the expression of genes. These include *in vitro* methods for determination of consensus binding sites of specific transcription factors, evaluation of potential accessibility of putative transcription factor binding sites (TFBSs) by DNase I hypersensitivity assays, electrophoretic mobility shift assays, and chromatin immunoprecipitation assays to determine the binding sites of a specific transcription factor within the genome. While this field has experienced considerable progress in the past, all of these methods, even when used in combination, are generally insufficient to successfully predict the location of a particular enhancer element or its tissue-specificity in an animal, prompting the need to validate and characterize putative enhancers in suitable *in vivo* assays.

Methods for *in vivo* testing of enhancer activities have been described for several vertebrate model organisms, including zebrafish and *Xenopus* [40,47]. In this article we will, however, focus on experimental approaches employing the mouse for determining the *in vivo* activity of candidate human enhancer sequences. Due to their shared phylogeny as mammals, the mouse is a suitable model for many aspects of human development, physiology, and disease. Importantly, mice are among the mammalian model organisms for which transgenic techniques have been available for many years, enabling the easy and efficient introduction of reporter constructs into the genome.

In order to study the *in vivo* properties of human enhancers, and in particular their ability to drive tissue-specific expression during embryonic development, we have recently set up a pipeline for testing of putative enhancers in transgenic mice (Fig. 3). We identify candidate elements by comparative criteria, such as human–fish comparison [46,80] or “ultra”-conservation between humans and rodents [60,61] (Fig. 2D). Then we assess the potential of such candidate regulatory regions experimen-

tally in a transgenic mouse enhancer assay [81,82]. Candidate regions are PCR-amplified from human genomic DNA and cloned into a reporter vector in which they are fused to a minimal heat shock protein 68 promoter and a beta-galactosidase reporter gene. On its own, this vector does not drive beta-galactosidase gene expression in mammalian embryonic tissues [81,82], but when fused to a DNA fragment with gene enhancer properties, spatial and temporal patterns of expression can be robustly and reproducibly characterized. This construct is injected into one of the two pronuclei of fertilized mouse oocytes, where it integrates into the genomic DNA at a random position, usually in multiple copies. The oocytes are then implanted into pseudo-pregnant females; embryos are harvested at embryonic day 11.5 and stained for β -galactosidase activity using X-Gal as a chromogenic substrate.

We chose this particular stage of development for analysis for several reasons. (1) Many human–fugu and ultra-conserved elements reside near genes that are expressed in early development [60,62]. (2) Whole embryo staining at this time-point enables the global identification of enhancer expression features without bias for particular tissues. (3) This is a key time-point during organogenesis at which most structures are present. Our preliminary studies of ~150 human–fugu elements indicate that this time-point is able to catch enhancer activities for >40% of the fragments tested, in contrast to moderately conserved human–rodent fragments where less than 5% of fragments behave as enhancers at this time-point [58]. Due to position effects that can alter *in vivo* enhancer characteristics as a result of the transgene integration site, we generate >5 independent transgenic animals per injection and require that at least 3 of these independent founders for each construct show reproducible spatial expression characteristics before assigning a conserved element an associated regulatory activity.

Compared with the generation of traditional BAC or YAC transgenic lines, use of this transient transgenic method results in a dramatically increased throughput that allows us to currently test 500 elements per year. This assay has previously been used in numerous gene-centric studies, where its reproducibility and high spatiotemporal resolution has provided valuable insights

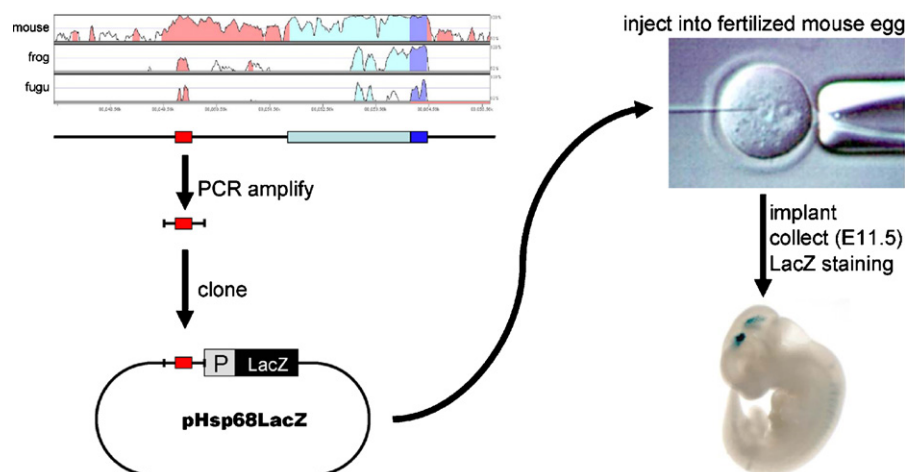


Fig. 3. Experimental design. Identification (example alignment displayed as Vista track), cloning and transgenic testing of candidate enhancer sequences.

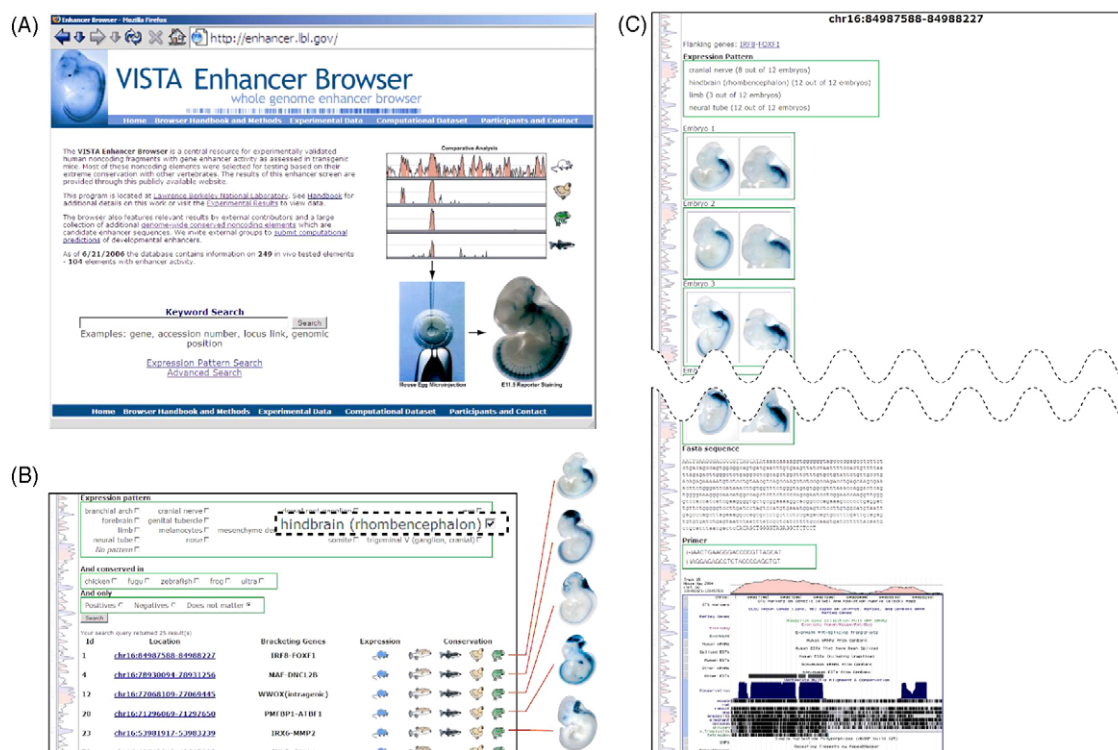


Fig. 4. Retrieving data from the Vista Enhancer Browser. (A) Entry page with basic query function. (B) Advanced search page with query form for experimental data. The results of a search for enhancers with hindbrain expression are shown. Each row in the results table corresponds to one experimental data set. A representative embryo is shown for the first five data sets. (C) Full data set display mode. Top: coordinates of element, neighboring genes, anatomical description of expression patterns and pictures of representative embryos. Note that each embryo is an independent transgenic F0 animal. Overview pictures and magnified views of expression sites are provided; all images can be downloaded at high resolution. Bottom: sequence of element, PCR primers used for cloning and conservation profile linked to UCSC Genome Browser.

into the *in vivo* activities of single elements of interest. This increase in throughput allows application of this method at a genomic scale, without requiring guidance by their neighboring genes.

4.3. Enhancer browser: large-scale data set of *in vivo*-validated enhancers

In order to make the results of our enhancer screen available to the scientific community, we have established a public database, the Vista Enhancer Browser, which is available at <http://enhancer.lbl.gov> (Fig. 4A). This browser houses two principal kinds of data: (1) experimental results from our *in vivo* screen and (2) a large collection of vertebrate non-coding sequences that are evolutionary conserved at varying distances.

4.3.1. Experimental data

The experimental results of our transgenic *in vivo* screen constitute the core data set of the enhancer browser. Each tested fragment has an associated dataset (Fig. 4C) consisting of sequence-related information and the experimental results. Sequence-related information includes the genomic coordinates, names of neighboring genes, PCR primers used to amplify the element from human genomic DNA, and an overview of the conservation in various species. The results of the transgenic enhancer assay are provided both in the form of pictures

of embryos with representative reporter gene activity and in anatomical annotation format. To be considered positive in our assay, an element has to drive reporter gene expression in the same anatomical structure in at least three independent transgenic embryos. Elements in which no such reproducibility is observed, although a sufficient number of transgenic embryos was generated (generally at least five transgenics confirmed by PCR genotyping) are reported as negative and no pictures of the embryos are shown. For positive elements, a selection of representative embryos is displayed. The images for each embryo can be retrieved as high-resolution files and are often supplemented by images at higher magnification or from more informative angles than the standard sagittal overview of the whole-mount specimen.

To enable searches of our data as well as bulk downloads, we annotate the tissue specificity of each positive enhancer identified using a list of anatomical terms that is largely consistent with existing standardized nomenclature [83]. We thus provide the ratio of X-Gal stained-positive embryos versus all transgenic embryos separately for each structure (Fig. 4C). A text-based query function is available on the front page of the enhancer browser. Using this feature, the database can also be searched by genomic coordinates, gene names, accession number and Entrez Gene IDs. An additional comprehensive search tool is available for more advanced queries of the database. This includes searches for enhancers that are specific for a particular

anatomical structure of interest (Fig. 4B) and/or restriction of the search to elements of a user-defined conservation depth (e.g. human–frog or human–fugu).

4.3.2. Computational data set

In addition to the experimental and external data, the enhancer browser also provides a genome-wide computationally generated set of more than 145,000 highly conserved elements for which no experimental data from the transgenic assay is available. These elements were identified using Gumbo/RankVista with globally aligned human–mouse–rat sequences [48]. Only elements with a $p \leq 0.001$ that do not overlap known mRNAs or spliced expressed sequence tags were considered for this data set. All of these elements were then checked for their conservation in chicken, frog, zebrafish and pufferfish to determine the conservation depth which is provided at the website. While we plan to test some subsets of this large collection of highly conserved elements in the future, the major purpose of this collection is to provide users with an easily accessible list of candidate regions for genomic intervals of interest for analysis in complementary computational and experimental approaches. Similar datasets can be obtained from other resources listed in Table 1.

5. Conclusions and perspectives

While gene regulation studies were possible in the pre-genome era, they were exceedingly expensive and time-consuming. Distant enhancers flanking a gene of interest were usually painstakingly identified through historic deletion series in transgenic animals. These experiments occurred sequentially in a largely trial and error fashion until the minimum sequence necessary to drive a given expression pattern was identified. Retrospective comparative analysis reveals that many of these functionally identified fragments strongly overlap with highly conserved regions of the human genome. For example, the distal liver-specific enhancer of *APOE*, a protein that impacts cholesterol metabolism, cardiovascular and Alzheimer's disease, was originally identified through such testing of many overlapping gene fragments in transgenic mice [6,7], but retrospective comparative analysis revealed that simple percent identity plot human–mouse comparisons would have readily identified this hepatic control region [84]. As is the case for numerous regulatory elements, had comparative data been available prior to beginning these experiments, hypotheses based on sequences under evolutionary constraint could have directly guided these studies from their inception.

Today, with this background experience, we are privileged to begin studies with computational sequence analysis followed by functional investigations. Such an approach can occur on a gene-by-gene basis or at a whole genome level of analysis. As a caveat, we should emphasize that comparative-based approaches are not without limitations. Some enhancers will lack conservation or may be missed by current computational tools, as illustrated in this article by the relatively weak conservation of many experimentally identified enhancers involved in heart development (Fig. 1B). While the thought of more vertebrate species genomic

sequences is a daunting data management task, their availability will without doubt further improve our ability to know which species to compare to address which biological question and allow additional flexibility in the choice of organisms used in multi-species analyses.

Importantly, the possibility of deep alignments across a wide range of vertebrate taxa will also increasingly allow us to address the relation between non-coding sequences and phenotypic diversity. One paradigmatic example to this end was the analysis of the aforementioned *Hoxc8* early enhancer in a panel of mammals that suggested that evolution of this enhancer contributed to the differences in axial morphology distinguishing baleen whales from other mammals [85]. While this study in the pre-genomic era relied on targeted sequencing of this regulatory element in a large number of species in the mammalian clade, the ever-growing number of available vertebrate sequences will increasingly allow for similar such studies at genomic scale.

The moderate-scale experimental testing of candidate enhancers through transgenic approaches such as that described here are expected to provide larger training sets for improved computational predictions of what activities conserved sequences are likely to contain. The first level of annotation in this area is occurring on the most highly (human–rodent “ultra”) and deepest (human–fish) conserved elements in the human genome. These classes of conserved non-coding elements are enriched near genes active in early development and this is not universally applicable for all types of known enhancers. Rather, they will serve to demonstrate how one can go from comparative sequence data to their functional testing to using the resulting dataset to computationally predict additional such enhancer elements in the larger human genome. It is anticipated that through such an iterative process we will learn vital clues as to developmental enhancer function and that this knowledge will translate into a deeper understanding of the regulation of both developmental and non-developmental genes in vertebrates.

Acknowledgements

L.A.P. was supported by grant HL066681, Berkeley-PGA, under the Programs for Genomic Applications, funded by National Heart, Lung, & Blood Institute, and HG003988 funded by National Human Genome Research Institute. Research was performed under Department of Energy Contract DE-AC02-05CH11231, University of California, E.O. Lawrence Berkeley National Laboratory. A.V. was supported by an American Heart Association postdoctoral fellowship.

References

- [1] Stathopoulos A, Levine M. Genomic regulatory networks and animal development. *Dev Cell* 2005;9:449–62.
- [2] Levine M, Tjian R. Transcription regulation and animal diversity. *Nature* 2003;424:147–51.
- [3] Davidson EH. Genomic regulatory systems: development and evolution. 1st ed. San Diego: Academic Press; 2001.
- [4] Venter JC, Adams MD, Myers EW, Li PW, Mural RJ, Sutton GG, et al. The sequence of the human genome. *Science* 2001;291:1304–51.

- [5] Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial sequencing and analysis of the human genome. *Nature* 2001;409:860–921.
- [6] Simonet WS, Bucay N, Pitas RE, Lauer SJ, Taylor JM. Multiple tissue-specific elements control the apolipoprotein *e/c-i* gene locus in transgenic mice. *J Biol Chem* 1991;266:8651–4.
- [7] Simonet WS, Bucay N, Lauer SJ, Taylor JM. A far-downstream hepatocyte-specific control region directs expression of the linked human apolipoprotein *e* and *c-i* genes in transgenic mice. *J Biol Chem* 1993;268:8221–9.
- [8] Allan CM, Walker D, Taylor JM. Evolutionary duplication of a hepatic control region in the human apolipoprotein *e* gene locus. Identification of a second region that confers high level and liver-specific expression of the human apolipoprotein *e* gene in transgenic mice. *J Biol Chem* 1995;270:26278–81.
- [9] Grehan S, Tse E, Taylor JM. Two distal downstream enhancers direct expression of the human apolipoprotein *e* gene to astrocytes in the brain. *J Neurosci* 2001;21:812–22.
- [10] Shih SJ, Allan C, Grehan S, Tse E, Moran C, Taylor JM. Duplicated downstream enhancers control expression of the human apolipoprotein *e* gene in macrophages and adipose tissue. *J Biol Chem* 2000;275:31567–72.
- [11] Zheng P, Pennacchio LA, Le Goff W, Rubin EM, Smith JD. Identification of a novel enhancer of brain expression near the *apoe* gene cluster by comparative genomics. *Biochim Biophys Acta* 2004;1676:41–50.
- [12] Lyons I, Parsons LM, Hartley L, Li R, Andrews JE, Robb L, et al. Myogenic and morphogenetic defects in the heart tubes of murine embryos lacking the homeo box gene *nkx2-5*. *Genes Dev* 1995;9:1654–66.
- [13] Chi X, Chatterjee PK, Wilson 3rd W, Zhang SX, Demayo FJ, Schwartz RJ. Complex cardiac *nkx2-5* gene expression activated by noggin-sensitive enhancers followed by chamber-specific modules. *Proc Natl Acad Sci USA* 2005;102:13490–5.
- [14] Tanaka M, Wechsler SB, Lee IW, Yamasaki N, Lawitts JA, Izumo S. Complex modular *cis*-acting elements regulate expression of the cardiac specifying homeobox gene *csx/nkx2.5*. *Development* 1999;126:1439–50.
- [15] Searcy RD, Vincent EB, Liberatore CM, Yutzey KE. A *gata*-dependent *nkx-2.5* regulatory element activates early cardiac gene expression in transgenic mice. *Development* 1998;125:4461–70.
- [16] Reecy JM, Li X, Yamada M, DeMayo FJ, Newman CS, Harvey RP, et al. Identification of upstream regulatory regions in the heart-expressed homeobox gene *nkx2-5*. *Development* 1999;126:839–49.
- [17] Lien CL, McAnally J, Richardson JA, Olson EN. Cardiac-specific activity of an *nkx2-5* enhancer requires an evolutionarily conserved *smad* binding site. *Dev Biol* 2002;244:257–66.
- [18] Lien CL, Wu C, Mercer B, Webb R, Richardson JA, Olson EN. Control of early cardiac-specific transcription of *nkx2-5* by a *gata*-dependent enhancer. *Development* 1999;126:75–84.
- [19] Schwartz RJ, Olson EN. Building the heart piece by piece: modularity of *cis*-elements regulating *nkx2-5* transcription. *Development* 1999;126:4187–92.
- [20] Zakany J, Gerard M, Favier B, Duboule D. Deletion of a *hoxd* enhancer induces transcriptional heterochrony leading to transposition of the sacrum. *EMBO J* 1997;16:4393–402.
- [21] Juan AH, Ruddle FH. Enhancer timing of *hox* gene expression: deletion of the endogenous *hoxc8* early enhancer. *Development* 2003;130:4823–34.
- [22] Kurokawa D, Kiyonari H, Nakayama R, Kimura-Yoshida C, Matsuo I, Aizawa S. Regulation of *otx2* expression and its functions in mouse forebrain and midbrain. *Development* 2004;131:3319–31.
- [23] Kurokawa D, Takasaki N, Kiyonari H, Nakayama R, Kimura-Yoshida C, Matsuo I, et al. Regulation of *otx2* expression and its functions in mouse epiblast and anterior neuroectoderm. *Development* 2004;131:3307–17.
- [24] Yanagisawa H, Clouthier DE, Richardson JA, Charite J, Olson EN. Targeted deletion of a branchial arch-specific enhancer reveals a role of *dhand* in craniofacial development. *Development* 2003;130:1069–78.
- [25] Gu H, Marth JD, Orban PC, Mossmann H, Rajewsky K. Deletion of a DNA polymerase beta gene segment in T cells using cell type-specific gene targeting. *Science* 1994;265:103–6.
- [26] Vong LH, Ragusa MJ, Schwarz JJ. Generation of conditional *mef2c*loxP/loxP mice for temporal- and tissue-specific analyses. *Genesis* 2005;43:43–8.
- [27] Li Song D, Joyner AL. Two *pax2/5/8*-binding sites in *engrailed2* are required for proper initiation of endogenous mid-hindbrain expression. *Mech Dev* 2000;90:155–65.
- [28] Iwahori A, Fraidenraich D, Basilico C. A conserved enhancer element that drives *fgf4* gene expression in the embryonic myotomes is synergistically activated by *gata* and *bhlh* proteins. *Dev Biol* 2004;270:525–37.
- [29] Guyot B, Valverde-Garduno V, Porcher C, Vyas P. Deletion of the major *gata1* enhancer *hs 1* does not affect eosinophil *gata1* expression and eosinophil differentiation. *Blood* 2004;104:89–91.
- [30] Chen JC, Goldhamer DJ. The core enhancer is essential for proper timing of myod activation in limb buds and branchial arches. *Dev Biol* 2004;265:502–12.
- [31] Xiong N, Kang C, Raulet DH. Redundant and unique roles of two enhancer elements in the *trcg* locus in gene regulation and γ T cell development. *Immunity* 2002;16:453–63.
- [32] Lettice LA, Horikoshi T, Heaney SJ, van Baren MJ, van der Linde HC, Breedveld GJ, et al. Disruption of a long-range *cis*-acting regulator for *shh* causes preaxial polydactyly. *Proc Natl Acad Sci USA* 2002;99:7548–53.
- [33] Sagai T, Hosoya M, Mizushima Y, Tamura M, Shiroishi T. Elimination of a long-range *cis*-regulatory module causes complete loss of limb-specific *shh* expression and truncation of the mouse limb. *Development* 2005;132:797–803.
- [34] Lettice LA, Heaney SJ, Purdie LA, Li L, de Beer P, Oostra BA, et al. A long-range *shh* enhancer regulates expression in the developing limb and fin and is associated with preaxial polydactyly. *Hum Mol Genet* 2003;12:1725–35.
- [35] Stenson PD, Ball EV, Mort M, Phillips AD, Shiel JA, Thomas NS, et al. Human gene mutation database (hgmd): 2003 update. *Hum Mutat* 2003;21:577–81.
- [36] Gumucio DL, Heilstedt-Williamson H, Gray TA, Tarle SA, Shelton DA, Tagle DA, et al. Phylogenetic footprinting reveals a nuclear protein which binds to silencer sequences in the human γ and ϵ globin genes. *Mol Cell Biol* 1992;12:4919–29.
- [37] Tagle DA, Koop BF, Goodman M, Slightom JL, Hess DL, Jones RT. Embryonic ϵ and γ globin genes of a prosimian primate (*galago crassicaudatus*). Nucleotide and amino acid sequences, developmental regulation and phylogenetic footprints. *J Mol Biol* 1988;203:439–55.
- [38] Hardison RC, Oeltjen J, Miller W. Long human–mouse sequence alignments reveal novel regulatory elements: a reason to sequence the mouse genome. *Genome Res* 1997;7:959–66.
- [39] Elgar G, Sandford R, Aparicio S, Macrae A, Venkatesh B, Brenner S. Small is beautiful: comparative genomics with the pufferfish (*fugu rubripes*). *Trends Genet* 1996;12:145–50.
- [40] Götting B, Barton LM, Gilbert JG, Bench AJ, Sanchez MJ, Bahn S, et al. Analysis of vertebrate *scl* loci identifies conserved enhancers. *Nat Biotechnol* 2000;18:181–6.
- [41] Loots GG, Locksley RM, Blankespoor CM, Wang ZE, Miller W, Rubin EM, et al. Identification of a coordinate regulator of interleukins 4, 13, and 5 by cross-species sequence comparisons. *Science* 2000;288:136–40.
- [42] Aparicio S, Chapman J, Stupka E, Putnam N, Chia JM, Dehal P, et al. Whole-genome shotgun assembly and analysis of the genome of *fugu rubripes*. *Science* 2002;297:1301–10.
- [43] Waterston RH, Lindblad-Toh K, Birney E, Rogers J, Abril JF, Agarwal P, et al. Initial sequencing and comparative analysis of the mouse genome. *Nature* 2002;420:520–62.
- [44] Boffelli D, Nobrega MA, Rubin EM. Comparative genomics at the vertebrate extremes. *Nat Rev Genet* 2004;5:456–65.
- [45] Aparicio S, Morrison A, Gould A, Gilthorpe J, Chaudhuri C, Rigby P, et al. Detecting conserved regulatory elements with the model genome of the Japanese puffer fish, *fugu rubripes*. *Proc Natl Acad Sci USA* 1995;92:1684–8.
- [46] Nobrega MA, Ovcharenko I, Afzal V, Rubin EM. Scanning human gene deserts for long-range enhancers. *Science* 2003;302:413.
- [47] Woolfe A, Goodson M, Goode DK, Snell P, McEwen GK, Vavouri T, et al. Highly conserved non-coding sequences are associated with vertebrate development. *PLoS Biol* 2005;3:e7.
- [48] Prabhakar S, Poulin F, Shoukry MI, Afzal V, Rubin EM, Couronne O, et al. Close sequence comparisons are sufficient to identify human *cis*-regulatory elements. *Genome Res* 2006;16:855–63.

- [49] International Human Genome Sequencing Consortium. Finishing the euchromatic sequence of the human genome. *Nature* 2004;431:931–45.
- [50] Buckingham M, Meilhac S, Zaffran S. Building the mammalian heart from two sources of myocardial cells. *Nat Rev Genet* 2005;6:826–35.
- [51] Kuo H, Chen J, Ruiz-Lozano P, Zou Y, Nemer M, Chien KR. Control of segmental expression of the cardiac-restricted ankyrin repeat protein gene by distinct regulatory pathways in murine cardiogenesis. *Development* 1999;126:4223–34.
- [52] Kuisk IR, Li H, Tran D, Capetanaki Y. A single *mef2* site governs desmin transcription in both heart and skeletal muscle during mouse embryogenesis. *Dev Biol* 1996;174:1–13.
- [53] Molkentin JD, Antos C, Mercer B, Taigen T, Miano JM, Olson EN. Direct activation of a *gata6* cardiac enhancer by *nkx2.5*: evidence for a reinforcing regulatory network of *nkx2.5* and *gata* transcription factors in the developing heart. *Dev Biol* 2000;217:301–9.
- [54] McFadden DG, Charite J, Richardson JA, Srivastava D, Firulli AB, Olson EN. A *gata*-dependent right ventricular enhancer controls *dhand* transcription in the developing heart. *Development* 2000;127:5331–41.
- [55] Dodou E, Verzi MP, Anderson JP, Xu SM, Black BL. *Mef2c* is a direct transcriptional target of *is11* and *gata* factors in the anterior heart field during mouse embryonic development. *Development* 2004;131:3931–42.
- [56] Phan D, Rasmussen TL, Nakagawa O, McAnally J, Gottlieb PD, Tucker PW, et al. *Bop*, a regulator of right ventricular heart development, is a direct transcriptional target of *mef2c* in the developing heart. *Development* 2005;132:2669–78.
- [57] Hu T, Yamagishi H, Maeda J, McAnally J, Yamagishi C, Srivastava D. *Tbx1* regulates fibroblast growth factors in the anterior heart field through a reinforcing autoregulatory loop involving forkhead transcription factors. *Development* 2004;131:5491–502.
- [58] Nobrega MA, Zhu Y, Plajzer-Frick I, Afzal V, Rubin EM. Megabase deletions of gene deserts result in viable mice. *Nature* 2004;431:988–93.
- [59] Pennacchio LA. Insights from human/mouse genome comparisons. *Mamm Genome* 2003;14:429–36.
- [60] Bejerano G, Pheasant M, Makunin I, Stephen S, Kent WJ, Mattick JS, et al. Ultraconserved elements in the human genome. *Science* 2004;304:1321–5.
- [61] Poulin F, Nobrega MA, Plajzer-Frick I, Holt A, Afzal V, Rubin EM, et al. In vivo characterization of a vertebrate ultraconserved enhancer. *Genomics* 2005;85:774–81.
- [62] Sandelin A, Bailey P, Bruce S, Engstrom PG, Klos JM, Wasserman WW, et al. Arrays of ultraconserved non-coding regions span the loci of key developmental genes in vertebrate genomes. *BMC Genomics* 2004;5:99.
- [63] Cooper GM, Stone EA, Asimenos G, Green ED, Batzoglu S, Sidow A. Distribution and intensity of constraint in mammalian genomic sequence. *Genome Res* 2005;15:901–13.
- [64] Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, et al. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res* 2005;15:1034–50.
- [65] Goodman M. The genomic record of humankind's evolutionary roots. *Am J Hum Genet* 1999;64:31–9.
- [66] Chimpanzee Sequencing and Analysis Consortium. Initial sequence of the chimpanzee genome and comparison with the human genome. *Nature* 2005;437:69–87.
- [67] Boffelli D, McAuliffe J, Ovcharenko D, Lewis KD, Ovcharenko I, Pachter L, et al. Phylogenetic shadowing of primate sequences to find functional regions of the human genome. *Science* 2003;299:1391–4.
- [68] Clark VJ, Cox NJ, Hammond M, Hanis CL, Di Rienzo A. Haplotype structure and phylogenetic shadowing of a hypervariable region in the *capn10* gene. *Hum Genet* 2005;117:258–66.
- [69] Margulies EH, Vinson JP, Miller W, Jaffe DB, Lindblad-Toh K, Chang JL, et al. An initial strategy for the systematic identification of functional elements in the human genome by low-redundancy comparative sequencing. *Proc Natl Acad Sci USA* 2005;102:4795–800.
- [70] Brudno M, Do CB, Cooper GM, Kim MF, Davydov E, Green ED, et al. Lagan and multi-lagan: efficient tools for large-scale multiple alignment of genomic DNA. *Genome Res* 2003;13:721–31.
- [71] Schwartz S, Zhang Z, Frazer KA, Smit A, Riemer C, Bouck J, et al. Pipmaker—a web server for aligning two genomic DNA sequences. *Genome Res* 2000;10:577–86.
- [72] Dubchak I, Brudno M, Loots GG, Pachter L, Mayor C, Rubin EM, et al. Active conservation of non-coding sequences revealed by three-way species comparisons. *Genome Res* 2000;10:1304–6.
- [73] Frazer KA, Pachter L, Poliakov A, Rubin EM, Dubchak I. Vista: computational tools for comparative genomics. *Nucleic Acids Res* 2004;32:W273–9.
- [74] Brudno M, Malde S, Poliakov A, Do CB, Couronne O, Dubchak I, et al. Global alignment: finding rearrangements during alignment. *Bioinformatics* 2003;19(Suppl 1):i54–62.
- [75] Loots GG, Ovcharenko I. Dcode.Org anthology of comparative genomic tools. *Nucleic Acids Res* 2005;33:W56–64.
- [76] Schwartz S, Kent WJ, Smit A, Zhang Z, Baertsch R, Hardison RC, et al. Human–mouse alignments with blastz. *Genome Res* 2003;13:103–7.
- [77] Bejerano G, Siepel AC, Kent WJ, Haussler D. Computational screening of conserved genomic DNA in search of functional non-coding elements. *Nat Methods* 2005;2:535–45.
- [78] Hinrichs AS, Karolchik D, Baertsch R, Barber GP, Bejerano G, Clawson H, et al. The ucsc genome browser database: update 2006. *Nucleic Acids Res* 2006;34:D590–8.
- [79] Blanchette M, Kent WJ, Riemer C, Elnitski L, Smit AF, Roskin KM, et al. Aligning multiple genomic sequences with the threaded blockset aligner. *Genome Res* 2004;14:708–15.
- [80] Ahituv N, Rubin EM, Nobrega MA. Exploiting human–fish genome comparisons for deciphering gene regulation. *Hum Mol Genet* 2004;13:R261–6. Spec No 2.
- [81] Kothary R, Clapoff S, Brown A, Campbell R, Peterson A, Rossant J. A transgene containing *lacZ* inserted into the dystonia locus is expressed in neural tube. *Nature* 1988;335:435–7.
- [82] Kothary R, Clapoff S, Darling S, Perry MD, Moran LA, Rossant J. Inducible expression of an *hsp68-lacZ* hybrid gene in transgenic mice. *Development* 1989;105:707–14.
- [83] Bard JL, Kaufman MH, Dubreuil C, Brune RM, Burger A, Baldock RA, et al. An internet-accessible database of mouse developmental anatomy based on a systematic nomenclature. *Mech Dev* 1998;74:111–20.
- [84] Pennacchio LA, Rubin EM. Genomic strategies to identify mammalian regulatory sequences. *Nat Rev Genet* 2001;2:100–9.
- [85] Shashikant CS, Kim CB, Borbely MA, Wang WC, Ruddell FH. Comparative studies on mammalian *hoxc8* early enhancer sequence reveal a baleen whale-specific deletion of a *cis*-acting element. *Proc Natl Acad Sci USA* 1998;95:15446–51.