

# Conservation of Regulatory Elements with Highly Diverged Sequences Across Large Evolutionary Distances

Mai H.Q. Phan<sup>1,2</sup>§, Tobias Zehnder<sup>2</sup>§, Fiona Puntieri<sup>2</sup>, Bai-Wei Lo<sup>2</sup>, Boris Lenhard<sup>3,4</sup>, Ferenc Mueller<sup>5</sup>  
Martin Vingron<sup>2</sup>, Daniel M. Ibrahim<sup>1,2\*</sup>

## Affiliations:

1 Berlin Institute of Health at Charité – Universitätsmedizin Berlin, Center for Regenerative Therapies, Charitéplatz 1, 10117 Berlin

2 Max Planck Institute for Molecular Genetics, Ihnestr. 63, 14195 Berlin

3 MRC London Institute of Medical Sciences, London, UK

4 Institute of Clinical Sciences, Faculty of Medicine, Imperial College London, Hammersmith Hospital Campus, London, UK

5 Institute of Cancer and Genomic Sciences, Birmingham Centre for Genome Biology, College of Medical and Dental Sciences, University of Birmingham, Birmingham, UK

§ co-first authors

\* Correspondence to: [daniel.ibrahim@bih-charite.de](mailto:daniel.ibrahim@bih-charite.de)

## Abstract

Embryonic gene expression is remarkably conserved across vertebrates as observed, for instance, in the developing hearts of chicken and mouse which diverged >300 million years ago. However, most *cis* regulatory elements (CREs) are highly divergent, which makes orthology tracing based on sequence similarity difficult, especially at larger evolutionary distances. Some evidence suggests functional conservation of CREs despite sequence divergence. However, it remains unclear how widespread such functional conservation might be. Here, we address this question by profiling the regulatory genome in the embryonic hearts of chicken and mouse at equivalent developmental stages. Gene expression and 3D chromatin structure show remarkable similarity, while the majority of CREs are non-alignable between the two species. To identify orthologous CREs independent of sequence alignability, we introduce a synteny-based strategy called Interspecies Point Projection (IPP). Compared to alignment-based approaches, IPP identifies up to 5-fold more putative orthologs in chicken, and up to 9-fold across distantly related vertebrates. We term these sequence-diverged orthologs indirectly conserved and characterize their functional conservation compared to sequence-alignable, directly conserved CREs. Indirectly and directly conserved elements show similar enrichment of functional chromatin signatures and cell-type specific enhancer sequence composition. Yet, shared transcription factor binding sites between orthologs are more heavily rearranged in indirectly conserved elements. Finally, we validate functional conservation of indirectly conserved chicken enhancers in mouse using *in vivo* reporter assays. Taken together, by overcoming the limitations of alignment-based methods our results reveal functional conservation of CREs across large evolutionary distances is more widespread than previously recognized.

## 1 **Introduction**

2 Embryonic organ development is driven by deeply conserved sets of transcription factors (TF) and  
3 signaling molecules that control tissue patterning, cell fates, and ultimately morphogenesis. Especially  
4 during the phylotypic stage, but also later during organogenesis many lineage and tissue-specific gene  
5 expression patterns are similar even between distantly related organisms (Irie and Kuratani 2011;  
6 Berthelot et al. 2017). One such example is the developing heart, where cellular patterning and  
7 morphological changes are deeply conserved across vertebrate lineages. The same key group of TFs  
8 expressed in the cardiac mesoderm is required for organogenesis, from the two-chambered heart in  
9 fish to the four-chambered hearts of birds and mammals (Olson 2006). Thus, the TFs of this gene  
10 regulatory network argue for a common genetic basis of embryonic development. Moreover, coding  
11 mutations in these genes have been shown account for ~45% of congenital heart disease (CHD) cases,  
12 the most common human birth defect (Pediatric Cardiac Genomics Consortium et al. 2013; Zaidi et al.  
13 2013; Jin et al. 2017). Many of the ~55% unsolved cases, not only for CHD, but also for other genetic  
14 diseases, might be caused by non-coding variants perturbing CREs of those developmental genes  
15 (Richter et al. 2020; Xiao et al. 2024).

16 However, many cis-regulatory elements (CREs) detected experimentally through DNA-accessibility or  
17 chromatin modifications are not sequence conserved (Visel et al. 2009; Villar et al. 2015), especially  
18 across larger evolutionary distances. For example, enhancers identified by chromatin modifications in  
19 embryonic heart tissues are poorly conserved (Blow et al. 2010). Similar observations were shown for  
20 TF binding sites (TFBS) in livers of five different vertebrate species (Schmidt et al. 2010). Yet, there are  
21 several examples for functionally conserved CREs in the absence of sequence conservation (Fisher et  
22 al. 2006; McGaughey et al. 2008; Madgwick et al. 2019). For example, the well-known *even-skipped*  
23 stripe 2 enhancer shows functional conservation amongst insects despite highly divergent sequences  
24 (Ludwig, Patel, and Kreitman 1998; Hare et al. 2008; Crocker and Stern 2017).

25 Determining orthologous CREs in distantly related species is complicated for several reasons. One, the  
26 rapid turnover in non-coding sequences constrains the effectiveness of pairwise alignments. Two,  
27 alignment-free methods struggle to accurately determine ortholog pairs. Alignment-free methods  
28 search for similar clusters of TFBS or “sequence words” as footprints of regulatory elements. (Sanges  
29 et al. 2006; Vinga 2014; Zielezinski et al. 2017, 2019). A more recent alignment-free approach,  
30 machine-learning algorithms, were successfully employed to identify cell-type specific enhancers  
31 across species. While this highlights the conservation of regulatory information independent of  
32 sequence alignability (Minnoye et al. 2020; Oh and Beer 2023; Kaplow et al. 2023; Kliesmete et al.  
33 2024), additional processing steps are needed to establish putative ortholog pairs. Three, the  
34 computational demands and availability of genome assemblies limit the use of multiple genome  
35 alignments, which is an alternative better suited to the task of orthology tracing across species. For  
36 example, the zebrafish ortholog of a human limb enhancer was identified indirectly through iterative

37 pairwise alignment between human and spotted gar, and between spotted gar and zebrafish (Braasch  
38 et al. 2016). More systematically, the use of one bridging species (*Xenopus*) helped to uncover  
39 hundreds of such “covert” ortholog pairs between human and zebrafish (Taher et al. 2011). Using  
40 Cactus multi-species alignments from hundreds of genomes, approaches like halliftover/HALPER  
41 (Paten et al. 2011; Hickey et al. 2013; Zhang et al. 2020; Armstrong et al. 2020) aim to trace orthology  
42 from genome sequences alone. However, in addition to the required computational infrastructure and  
43 the availability of genome assemblies, these approaches are currently not available for larger  
44 evolutionary distances (e.g. chicken-mouse).

45 Here we present an experimental-analytical framework to efficiently identify orthologous CREs  
46 combining two currently underutilized features – synteny and functional genomic data. In genomics,  
47 *synteny* describes the maintenance of colinear genomic sequences on chromosomes of different  
48 species (Engström et al. 2007; Kikuta et al. 2007). Not only genes are maintained in synteny;  
49 developmental genes are often flanked by conserved non-coding elements (CNEs), many of which act  
50 as enhancers (Siepel et al. 2005; Bejerano et al. 2005; Visel, Bristow, and Pennacchio 2007). Their  
51 syntenic arrangement reflects conserved regulatory environments that have been described as  
52 genomic regulatory blocks (GRBs)(Kikuta et al. 2007; Harmston et al. 2017). *Functional genomic data*,  
53 such as chromatin accessibility and histone modifications, are widely used to determine putative CREs  
54 in any tissue of interest. Given that the hearts of birds and mammals are evolutionary homologous  
55 structures, the active regulatory genome in both should be related. Therefore, experimentally  
56 identified CREs in both species might provide the genomic footprint of functionally conserved  
57 orthologs whose sequences have diverged to the point where alignment fails. We first use chromatin  
58 profiling from murine and chicken hearts at equivalent developmental stages to experimentally  
59 determine regulatory elements. We then apply Interspecies Point Projection (IPP), a synteny-based  
60 algorithm designed to map corresponding genomic locations in highly diverged genomes.

61 Using this strategy, we uncover thousands of previously hidden orthologous CREs based on their  
62 relative position in the genome and overcome current limitations. We term these sequence-diverged  
63 orthologs indirectly conserved, and validate their functional equivalence as compared to classical  
64 sequence-conserved elements. We find similar enrichment of chromatin marks at directly and  
65 indirectly conserved elements. Furthermore, using machine learning models and TFBS-driven analysis,  
66 we show that both classes display similar heart-enhancer specific sequence composition. Yet, shared  
67 TFBS are more heavily rearranged between indirectly conserved CRE pairs. Finally, we demonstrate  
68 their functional orthology using *in vivo* enhancer-reporter assays. Thereby we demonstrate a currently  
69 underrepresented widespread conservation of cis-regulatory elements with highly diverged sequences  
70 across large evolutionary distances.

## 71 Results

### 72 Identification of heart CREs from equivalent developmental stages in mouse and chicken

73 To identify the cis-regulatory elements driving gene expression at equivalent stages of heart  
 74 development, we generated comprehensive chromatin and gene expression profiles from embryonic  
 75 mouse and chicken hearts (ChIPmentation, ATAC-seq, RNA-seq, Hi-C) at E10.5/E11.5 and HH22/HH24  
 76 (Fig. 1a). To compare global gene expression profiles, we measured differentially expressed genes in  
 77 the heart vs. limb in both mouse and chicken (Fig. 1b). Consistent with previous reports (Olson 2006)  
 78 tissue-specific expression is conserved, including key TF genes specific for heart and limb development  
 79 (Fig. 1b, Fig. S1a). To characterize conservation of regulatory regions driving this expression, we first  
 80 estimated sequence conservation by alignment of open chromatin regions using LiftOver (Kuhn et al.  
 81 2009). Most mouse peaks in non-coding regions lacked sequence conservation in chicken, in stark

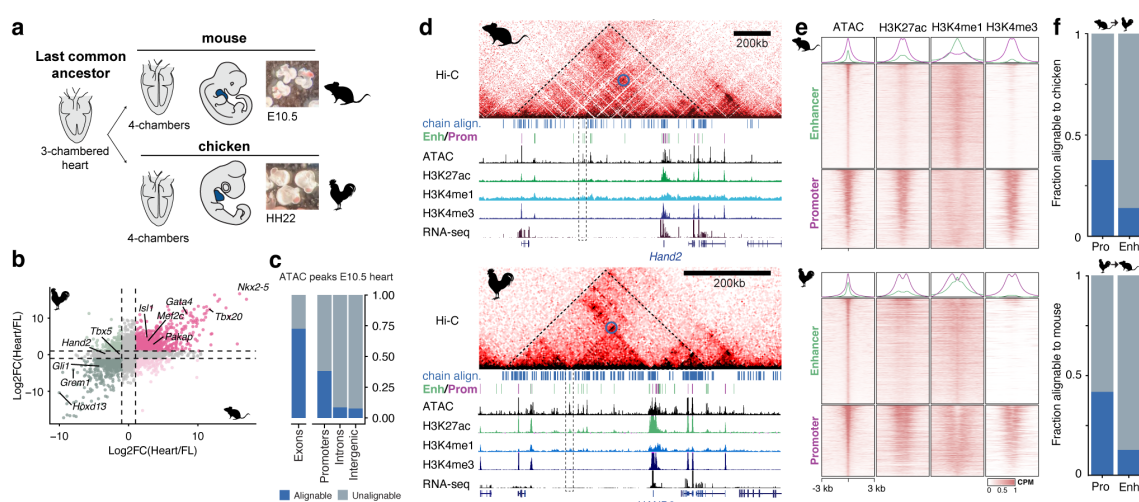


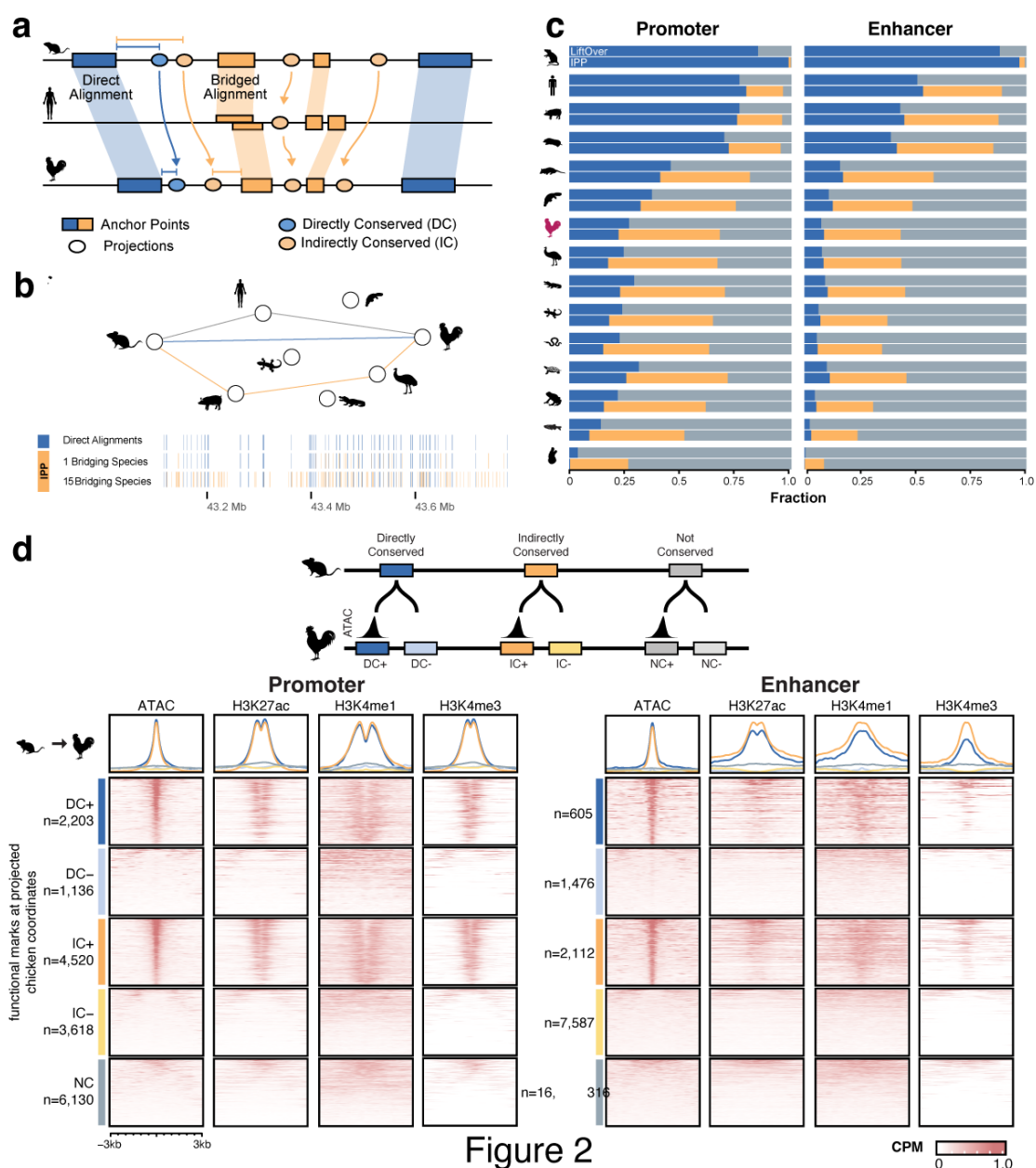
Figure 1

### Figure 1: Evolutionary conservation of gene expression and chromatin structure between mouse and chicken embryonic hearts despite divergent cis-regulatory elements

- Reptilian and mammalian lineages convergently evolved fully separated 4-chambered hearts. E10.5/HH22 represent equivalent stages of heart formation
- Conservation of global gene expression (log<sub>2</sub> fold-change of heart vs. limb expressed genes) between mouse (E10.5) and chicken (HH22).
- ATAC-seq peaks (E10.5 heart) are mostly alignable (LiftOver (--minMatch = 0.1)) to chicken in coding, but not in non-coding regions.
- Syntenic regions at the *Hand2*/*HAND2* locus shows conserved 3D chromatin structure and histone modifications relative to the target gene despite different genomic size. Dashed triangle outline conserved TAD structure, blue circles/dashed rectangle show specific contacts to conserved enhancers. Blue ticks: conserved sequences, Green/Purple ticks: predicted promoters/enhancers.
- Signal enrichment (+/-3kb) of histone modifications at heart promoters and enhancers, centred on ATAC-seq peaks
- Fraction of alignable elements identified in e) with the chicken/mouse genome

82 contrast to those overlapping exons (Fig. 1c, Fig. S1b). We then used Hi-C and ChIPmentation data to  
 83 more comprehensively profile the regulatory genome. Hi-C showed global conservation of the 3D  
 84 genome in syntenic regions surrounding most developmental genes (Fig. 1d) and enrichment of  
 85 synteny breaks at TAD boundaries (Fig. S1). Syntenic regions surrounding developmental genes  
 86 showed comparable distribution of chromatin marks indicating that the position of regulatory

87 elements relative to their targets might be conserved (**Fig. 1d**). We used CRUP to predict active CREs  
 88 from typical histone modifications (Ramisch et al. 2019). To further refine our list, we integrated CRUP  
 89 predictions with chromatin accessibility and gene expression data, followed by stringent filtering (see  
 90 Methods) to have a high-confidence set of active enhancers and promoters for both species,



**Figure 2**

**Figure 2: A syntenic-based algorithm, Interspecies Point Projection (IPP), identifies thousands of putative sequence orthologs of mouse heart CREs with functional chromatin signatures in chicken**

- Syntenic-based proximity to direct/indirectly aligned regions determines orthology between features (e.g. ATAC-peak summits).
- Multi-species bridged alignments increase the number of anchor points in a representative region using 0, 1 and 15 bridging species
- IPP increases the number of putatively homologous regions from mouse to 15 other species used as bridging species (compare blue vs. orange portion). LiftOver alignments (top bar) are compared to IPP *Directly Conserved* and *Indirectly Conserved*. Increase is particularly high at greater evolutionary distances to non-mammalian species
- Classification of elements with or without conserved activity +/- . Signal enrichment at chicken genomic regions to which mouse E10.5 heart elements were projected.

91 minimizing the number of false-positive regions. In total, we called 20,252 promoters and 29,498  
92 enhancers in mouse and 14,806 and 21,641 in chicken hearts, respectively (**Fig. 1d, Fig. S1c**).

93 We then estimated the degree of sequence conservation for this high-confidence set of regulatory  
94 elements. Consistent with previous reports (Blow et al. 2010) less than 50% of promoters and only  
95 ~10% of enhancers were sequence-conserved between mouse and chicken (**Fig. 1f, Fig. S1d**). Thus, the  
96 lack of sequence alignability remains consistent, even when restricting the analysis to a stringently  
97 filtered set of enhancers and promoters, contrasting conserved gene expression patterns and 3D  
98 chromatin structure.

### 99 **A synteny-based strategy to identify orthologous genomic regions**

100 Because enhancer function can be maintained despite rapid turnover of underlying sequences, DNA  
101 sequence conservation alone likely underestimates conserved regulatory activity. To identify such  
102 conserved, non-alignable CREs we developed a synteny-based algorithm, Interspecies Point Projection  
103 (IPP) (Baranasic et al. 2022), designed to find orthologous regions independent of sequence divergence  
104 (see Supplemental Text and Fig. S2). The approach is based on conserved synteny. We assume any non-  
105 alignable element in one genome located between flanking blocks of alignable regions is located at  
106 the same relative position in another genome (**Fig. 2a**). Thus, for a given species pair we can interpolate  
107 the position of an element (e.g. an enhancer) relative to adjacent alignable regions, so-called *anchor*  
108 *points*. We refer to the interpolated coordinates in the target genome as *projections*. Because a larger  
109 distance to an *anchor point* reduces accuracy of projections, the second pillar of IPP involves optimizing  
110 the use of bridged/tunneled alignments (Taher et al. 2011; Braasch et al. 2016). IPP uses not one, but  
111 multiple bridging species, which increases the number of anchor points thereby minimizing this  
112 distance (**Fig. 2b**). With this, IPP classifies orthologous regions by their distance to a *bridged alignment*  
113 or *direct alignment*. Regions projected within 300bp to a *direct alignment* are defined as *directly*  
114 *conserved* (DC). Those further than 300bp to a *direct alignment* but which can be projected through  
115 bridged alignments we define as *indirectly conserved* (IC) regions if the summed distance to anchor  
116 points is less than 2.5kb. The remaining projections are defined as non-conserved (NC) (see **Fig. S2** and  
117 Supplemental Text for details and parameterization).

### 118 **IPP improves detection of orthologs between distantly related species**

119 To optimize our mouse-chicken projections, we selected a set of 16 species, consisting of mouse,  
120 chicken and 14 bridging species from the reptilian and mammalian lineages along with additional  
121 ancestral vertebrate/chordate genomes (see Methods). After building our collection of anchor points  
122 from pairwise alignments, we project our set of murine heart CREs to chicken and all bridging species  
123 to estimate their conservation at varying evolutionary distances. In parallel we used UCSC LiftOver to  
124 serve as a reference for sequence conservation for IPP projections. In practice, LiftOver performed  
125 similarly to IPP DC projections for all 15 species (**Fig. 2c**), with the exception that multiple mappings

126 can occur when ‘lifting’ the entire sequence. The proportion of mouse CREs classified as *directly*  
127 *conserved* (DC) reduces drastically with increasing evolutionary distances. While over 90% of CREs are  
128 conserved when comparing mouse to the closely related rat, this number drops to 50-70% within  
129 placental mammals and even more so to non-mammalian vertebrates. Specifically for chicken, only  
130 22% of all promoters and 10% of enhancers are sequence conserved (**Fig. 2c**).

131 By additionally identifying indirectly conserved (IC) regions, IPP increases the number of conserved  
132 elements in all species. Especially within distantly related vertebrates this increases by a factor of 3 to  
133 9, and substantially adds to the number of putatively conserved CREs (orange fraction, **Fig. 2c**). For the  
134 mouse-chicken comparison, the percentage of conserved promoters increases 3-fold (18,9%, DC) to  
135 65%, DC+IC), and for enhancers 5-fold (7,4% to 42%). With this, IPP pairs an additional 8,138 and 9,699  
136 promoters and enhancers with candidate ortholog sequences in chicken.

137 Unlike the synteny-based approach of IPP, other efforts to improve ortholog identification include  
138 hierarchical alignments, which are multiple-genome alignments guided by evolutionary relationships  
139 (Hickey et al. 2013; Zhang et al. 2020). We compared IPP with halliftover/HALPER (Zhang et al. 2020),  
140 which uses Cactus alignments from hundreds of mammalian genomes, for all placental mammals in  
141 our species collection (i.e. rat, human, pig, mole). Depending on parameterization, IPP performs  
142 similarly or better at identifying orthologous enhancers within this relatively short evolutionary  
143 distance. This indicates that ortholog can be traced across evolutionary distances by comparing  
144 hundreds of genome sequences. However, the synteny-based strategy of IPP achieves comparable  
145 detection rates using only 16 species and spans a greater evolutionary distance than currently available  
146 for hierarchical alignments.

147 Since IPP can project any set of genomic coordinates, we next used IPP on a set of limb enhancers we  
148 identified, as well as on two published datasets that reported low conservation between mouse and  
149 chicken: murine heart enhancers from (Blow et al. 2010), and a set of CEBP/A TFBSs in liver from  
150 (Schmidt et al. 2010). IPP increased the number of putative ortholog heart enhancers equivalent to  
151 that of our heart enhancers (**Fig. S2b**). Heart enhancers were slightly less well conserved (DC and IC)  
152 than limb enhancers (**Fig. S2c**), confirming general trends observed previously (Blow et al. 2010). For  
153 CEBP/A only 2% of murine peaks were reported to be conserved in chicken, and even less bound by  
154 CEBP/A in chicken livers (Schmidt et al. 2010). We re-analyzed the ChIP-seq data from mouse and  
155 chicken livers and confirmed that only a small fraction (5,7%) of mouse CEBP/A binding sites were  
156 directly conserved in chicken (DC) and just 173 of these sites overlapped with a CEBP/A peak in chicken  
157 (**Fig. S2f**). However, by including IC projections, we increased the number conserved CEBP/A sites to  
158 32% and found an additional set of 579 peaks that were also CEBP/A bound in chicken livers.

159 Taken together, IPP dramatically increases detection of orthologous genomic regions, particularly for  
160 larger evolutionary distances, uncovering a previously hidden set of conserved elements that can be  
161 investigated for their role in evolution and gene regulation.

162

### 163 **Indirectly and directly conserved CREs show a similar enrichment for functional chromatin marks**

164 The large additional number of IC regions suggests that up to 80% of conserved CREs might have gone  
165 undetected in most analyses to date. Since we collected functional genomic data from  
166 developmentally equivalent stages, we first profiled chromatin signal and compared how well the  
167 chromatin state at DC and IC predicted CREs is conserved in chicken. For DC CREs, we found that 66%  
168 mouse promoter and 29% enhancer projections overlapped an ATAC-seq peak in the chicken genome.  
169 Interestingly, these percentages were similar for IC CREs with 56% promoter and 26% enhancer  
170 projections, although the absolute numbers of IC promoters and enhancers is substantially higher than  
171 DC. We classified these regions with conserved activity as DC+/IC+ and those without an ATAC-seq  
172 peak at the projected site as DC-/IC- (**Fig. 2d**). Consistent with the ATAC-seq signal, DC+/IC+ CREs  
173 showed equivalent specific enrichment of H3K4me3 at promoters and H3K4me1 at enhancers,  
174 suggesting that the IC CREs identify the functional orthologs of murine heart CREs in the chicken  
175 genome (**Fig. 2d**). This similar enrichment of functional chromatin marks suggests that interpolated  
176 regions point to “functionally conserved” CREs in the target genome and that sequence homology is  
177 an incomplete indicator of conserved activity.

178

### 179 **SVM model robustly learns tissue-specific sequence features and independently validates IPP** 180 **performance**

181 Recently, machine learning (ML) methods have become a viable strategy to identify cell-type specific  
182 CREs in distantly related species, by virtue of their ability to capture complex sequence-function  
183 relationships without relying on strict sequence conservation (Minnoye et al. 2020; Oh and Beer 2023;  
184 Kliesmete et al. 2024; Kaplow et al. 2023). To test the regulatory content in the IPP projections in  
185 chicken, we first trained a gapped k-mer Support Vector Machine (gkm-SVM) model on mouse data to  
186 identify heart-specific enhancers. To learn predictive heart-specific enhancer vocabularies, we trained  
187 the SVM on aggregated tissue-specific ATAC-seq peaks from mouse embryonic heart outside promoter  
188 regions, against the background of non-overlapping peaks from non-heart cell/tissues (**Fig. 3a**, see  
189 Methods).

190 We then tested the model’s cross-species predictive power on the chicken enhancer regions we  
191 identified in the embryonic heart and forelimb (FL). The mouse-trained SVM correctly distinguished  
192 between heart-specific, shared, and FL-specific chicken enhancers (**Fig. 3a**). This validated that the  
193 features from mouse sequences are in fact predictive of heart-specific enhancers in chicken. A recent



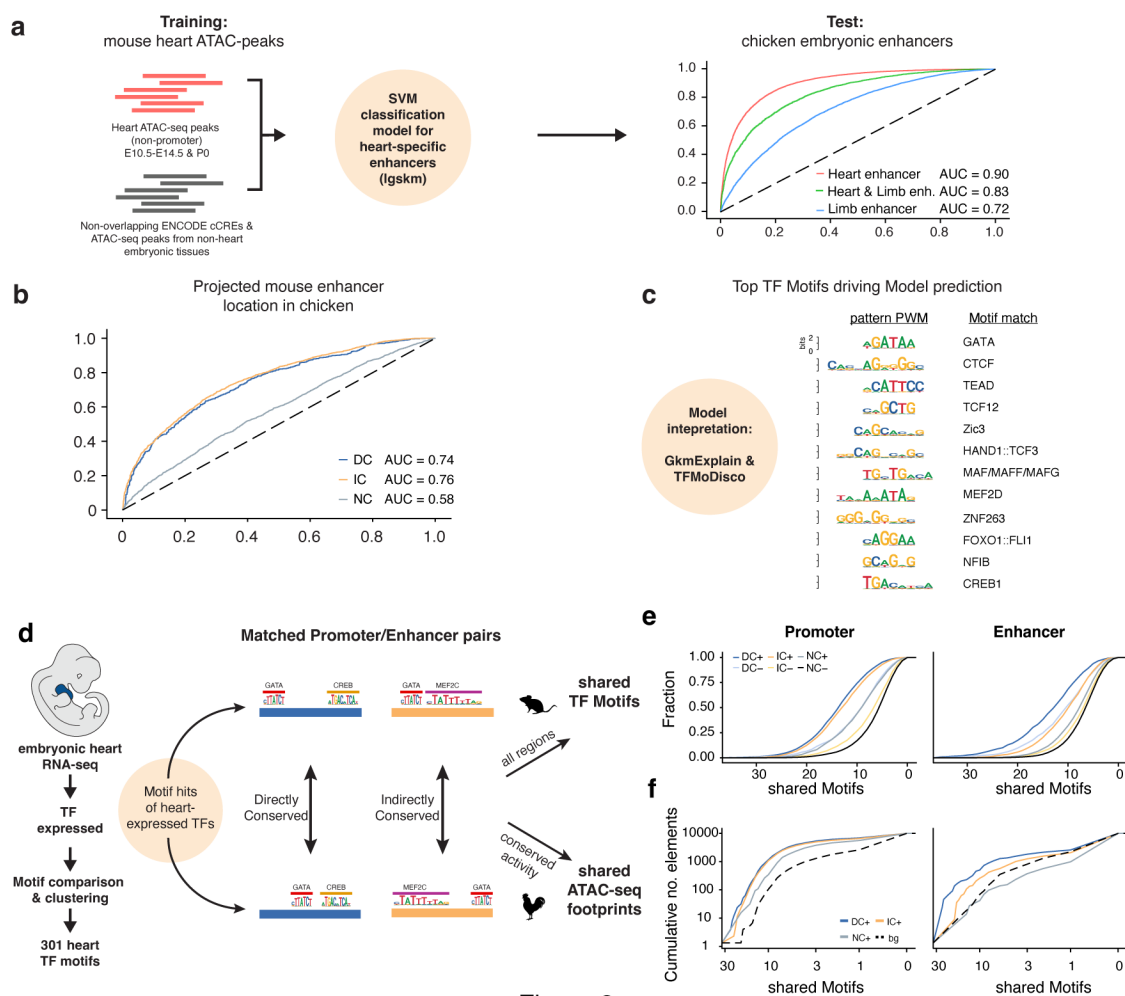


Figure 3

**Figure 3: In silico analysis of sequence composition and motif content of indirectly and directly conserved elements**

- Training of a Support Vector Machine (SVM) model to identify heart enhancers with independent data from public repositories. Positive set: embryonic heart/cardiomyocyte ATAC-seq peaks, Negative Set: non-overlapping ATAC-seq peaks from non-heart tissues. The model distinguishes heart- vs. limb-specific enhancers from chicken embryos. AUC: Area Under Curve
- Evaluation of DC+, IC+ and NC classified regions of the chicken genome by the SVM Model.
- TFMoDisco interpretation of the putative TFBS that contribute to model specificity. BS of several known heart-specific TFs contribute to model accuracy.
- Heart-expressed TFs identified from RNA-seq were consolidated to 301 motifs of heart specific TFs. Promoter/Enhancer pairs were screened for shared TFBS or ATAC-seq footprints
- DC+/IC+ promoters/enhancers share more heart TFBS than DC-/IC-, or non-conserved NC regions
- Functionally conserved DC and IC ATAC-seq peak pairs share more TF-footprints than NC ATAC-seq peak pairs or control pairs (bg = a non-paired ATAC-seq peak in the same TAD)

194 study found that tissue-specific CREs show a lower degree of sequence conservation than more  
 195 pleiotropic CREs (Kliesmete et al. 2024). We therefore evaluated SVM-predicted tissue-specificity of all  
 196 ATAC-Seq peaks from chicken embryonic hearts and noted a clear inverse relationship to sequence  
 197 conservation (**Fig. S3**). In other words, predicted heart-specific chicken regions (i.e. positive score) are  
 198 more sequence-divergent from mouse than more pleiotropic peaks, providing further evidence that  
 199 sequence alignability is a poor estimator of conserved regulatory activity.  
 200 Since IC elements exhibit similar degree of conserved activity to DC elements in terms of epigenomic  
 201 signatures (**Fig. 2d**), we next wanted to estimate conservation as defined by its shared tissue specificity

202 between species. We therefore compared the predicted tissue-specificity of mouse enhancers  
203 projected to orthologous chicken loci between DC, IC and NC elements. DC and IC projections were  
204 equally likely to be classified as heart-specific enhancers (AUC, DC=0.74, IC=0.76). NC projections,  
205 however, were less likely to be classified as heart enhancer (AUC=0.58) (**Fig. 3b**), further indicating  
206 conserved tissue-specific enhancer activity.

207 To better understand predictive sequence patterns learned by the model, we computed the  
208 contribution of individual nucleotides from input sequences to the SVM output classification with  
209 GkmExplain (Shrikumar, Prakash, and Kundaje 2019) and consolidated recurring high scoring patterns,  
210 or 'seqlets', into motifs (Shrikumar et al. 2018). Motifs discovered from mouse and chicken sequences  
211 largely overlap (**Fig. S3**), suggesting conserved enhancer vocabularies. In fact, known motifs of master  
212 regulators of heart development (e.g. GATA, TEAD and HAND) were most predictive of tissue specificity  
213 (**Fig. 3c**), further supporting the model's robustness in predicting heart-specific enhancers. Thus, this  
214 independent approach validates that the IPP projections of mouse enhancers faithfully identify heart-  
215 specific enhancer regions in the chicken genome.

216

#### 217 **Transcription factor binding site conservation as indicator of conserved CRE activity**

218 If IPP projections represent conserved pairs of CREs, these regions should share the same TFBS. Here,  
219 we can evaluate this both at the sequence- and chromatin level given our available data using TF motif  
220 scanning and ATAC-seq footprinting. We used our heart RNA-seq data to identify TFs expressed in the  
221 heart and curated a set of 301 heart TF motifs (**Fig. 3d**). We then calculated for every mouse-chicken  
222 ortholog pair how many TFBS were shared and plotted the results (**Fig. 3e**).

223 Overall, orthologous promoter regions shared more TFBS hits than enhancers. DC+/IC+ promoters  
224 were comparable in the number of shared TFBS and both were clearly distinguishable from DC-/IC-  
225 promoters (**Fig. 3e**). For enhancers, DC+ shared the most TFBSs, while IC+ enhancer pairs shared as  
226 many TFBS as DC- enhancers. Notably, CREs with conserved active chromatin marks (dark blue/orange  
227 lines) in all comparisons shared more TFBS than those without (light blue/orange lines), irrespective  
228 of direct or indirect conservation. This suggest that functionally conserved orthologs are more likely to  
229 retain regulatory information. Finally, we used our ATAC-seq data to compare shared TF footprints. We  
230 compared all DC/IC/NC pairs that had ATAC-seq signal in both genomes relative to background (non-  
231 orthologous ATAC-seq peaks within the same TAD, see Methods). Consistent with our TFBS motif  
232 results, all projection pairs outperformed control regions. DC and IC promoters were equal in the  
233 number of shared TF-footprints, while DC enhancers were overall slightly more likely to share TF-  
234 footprints than IC enhancers (**Fig. 4f**). These results confirm that IPP identifies orthologous pairs of

235 CREs with shared TFBS, representing a conserved sequence syntax that is independent of direct  
236 sequence conservation.

237

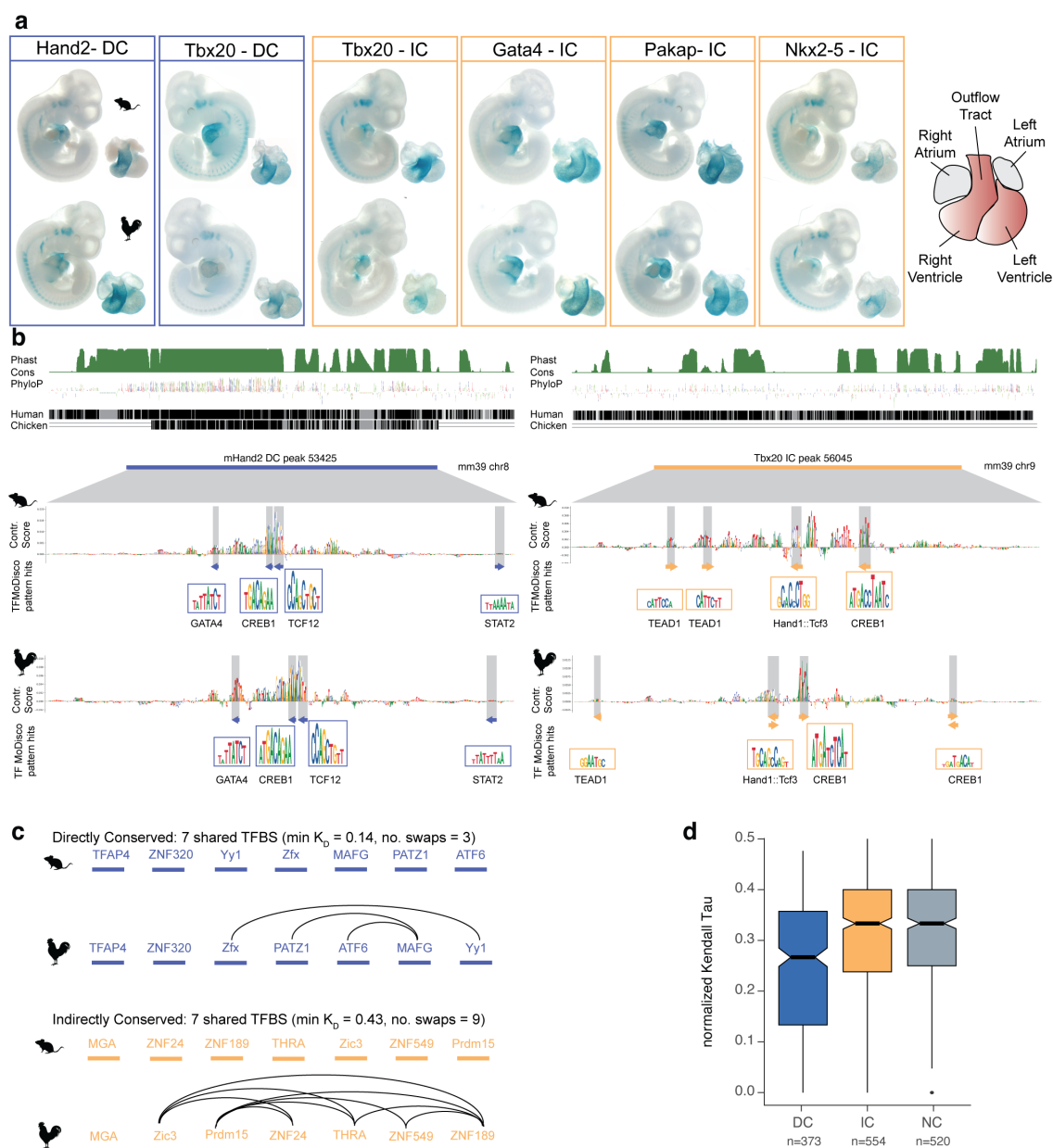
### 238 **Indirectly conserved heart enhancers from chicken drive conserved gene expression patterns in** 239 **mouse embryonic hearts**

240 Gene regulatory elements drive tissue and cell-type specific expression. Based on our analysis, directly  
241 and indirectly conserved elements are functionally conserved orthologs and should drive conserved  
242 expression patterns in the developing heart. To test this, we selected two pairs of DC and 4 pairs of IC  
243 enhancers and generated *in vivo* enhancer-reporter mice for each of these elements. We profiled  
244 enhancer activity using lacZ staining in E10.5 mouse embryos. All enhancer pairs drove conserved  
245 expression with remarkable specificity (**Fig. 4a**). Enhancers driving expression patterns in specific  
246 regions of the heart, such as the outflow tract and atrio-ventricular canal, were consistent with those  
247 from nearby genes (Hand2-DC, Tbx20-IC, Nkx2-5-IC) (Overbeek 1997; Srivastava and Olson 1997; Firulli  
248 et al. 1998; McFadden et al. 2000; Stennard et al. 2003; Prall et al. 2007). The same was true for the  
249 ventricle-specific expression of two other enhancers (Tbx20-DC, Gata4-IC) (Heikinheimo, Scandrett,  
250 and Wilson 1994). An indirectly conserved enhancer at the *Pakap* locus (Pakap-IC), which contains the  
251 *A-kinase anchoring protein 2* (*Akap2*) gene involved in general cardiomyocyte function (Maric et al.  
252 2021), drives broad expression in all cardiac tissues. We integrated scores obtained from the SVM  
253 model for all tested pairs. Many seqlets with high contribution scores to our enhancer prediction  
254 overlapped with predicted binding sites of key TFs (**Fig. 4b** shaded boxes and **Fig. S4**) and were shared  
255 between mouse and chicken CREs for each pair. These data show that the chicken IC enhancers we  
256 identify constitute *bona fide* orthologs to their mouse counterparts, regardless of sequence  
257 conservation.

### 258 **Indirectly conserved CREs show a higher degree of TFBS shuffling**

259 In all our analyses and validations, IC regions showed similar signatures of functional conservation to  
260 DC, despite lack of alignability. We therefore wanted to explore how the underlying DNA sequences  
261 may differ in ways they encode regulatory information. We hypothesized that for CRE pairs with a  
262 similar number of shared TFBSs, DC pairs would display a more conserved TFBS order within the  
263 element than IC pairs (**Fig. 4c**). For example, a DC and IC enhancer pair with both 7 shared TFBSs, show  
264 a more shuffled order between IC pairs, likely complicating alignment of the two sequences. To  
265 systematically evaluate this phenomenon, we calculated the Kendall-Tau rank distance for all enhancer  
266 pairs. The Kendall-Tau rank distance assesses the similarity between two ranked lists by measuring the  
267 number of swaps needed to change one list into the order of the other list (Qian and Yu 2019). We  
268 selected all functionally conserved enhancer pairs with at least 6 shared TFBSs and computed the  
269 normalized Kendall-Tau Distance for each pair (**Fig. 4 c,d**). DC enhancers exhibited a significantly lower

270 KD score (median = 0.27) than IC (median=0.33) and NC enhancers (median=0.33). Consequently,  
 271 conservation of an element's regulatory function is likely less dependent on exact sequence  
 272 conservation than on preserving the appropriate balance of TFBSs within the given element.



**Figure 4: Indirectly conserved heart enhancers from mouse and chicken drive conserved gene expression pattern *in vivo***

- Directly and Indirectly conserved enhancers from mouse (top) and chicken (bottom) drive highly similar expression patterns in the heart of E10.5 embryos. Individual enhancer show similar tissue-restricted or broad expression patterns.
- Sequence conservation scores (PhastCons/PhyloP) and direct alignments to human and chicken of the murine Hand2-DC and Tbx20-IC enhancer tested in a). SVM contribution scores and TF-Modisco Motif matches show conserved sequence features of the 500bp enhancer highlighting shared TF-Motif hits overlapping with seqlets.
- The different order of shared TFBSs in IC and DC enhancer pairs is reflected in the computed Kendall-Tau Distance,  $K_D$ .
- $K_D$  scores for all functionally conserved DC/IC/NC CRE enhancer pairs. Asterisks indicate the magnitude of effect size based on Cohen's  $d$ : small (\*,  $d < 0.2$ ), medium (\*\*,  $d \leq 0.5$ )

## 273 Discussion

274 Here, we show widespread conservation of functional gene regulatory elements in the absence of  
275 direct sequence conservation. By combining equivalent functional genomic data from two species, a  
276 synteny-based algorithm, and *in vivo* validation we reveal a substantial amount of previously hidden  
277 indirectly conserved elements functionally equivalent between mouse and chicken.

278 Identification of orthologous enhancers between distantly related species is an inherently difficult  
279 problem due to rapid enhancer evolution (Berthelot et al. 2017; Villar et al. 2015). While there have  
280 been several individual reports describing enhancers conserved in function rather than in sequence  
281 (Madgwick et al. 2019; Crocker and Stern 2017; Hare et al. 2008; Braasch et al. 2016; Fisher et al. 2006),  
282 a systematic evaluation of this phenomenon is challenging. Not only does it require algorithmic  
283 approaches that attempt to pair non-alignable sequences, but it also requires functional data that can  
284 be used to validate these predictions. By combining the synteny-based algorithm IPP with matching  
285 experimental data from two species, we were able to predict a large set of previously hidden indirectly  
286 conserved elements and demonstrate they are as likely to be functionally conserved as directly  
287 conserved elements. Our reanalysis of previous studies show that these likely 5-fold underestimate  
288 the number of chicken-conserved enhancers (Blow et al. 2010; Schmidt et al. 2010). While this does  
289 not change the general trend observed in these studies, the degree of underreported conserved  
290 regulatory elements changes the interpretation to which degree enhancers may evolve from neutral  
291 sequences (Galupa et al. 2023) and to which degree they are conserved. Our results indicate a degree  
292 of conservation invisible to current alignment-based measures. Thereby, our approach reconciles the  
293 apparent contrast between divergent non-coding genome sequences and other conserved features  
294 such as 3D chromatin structure and gene expression.

295 Rapidly diverging regulatory DNA allows adaptation of the regulatory genome during evolution but  
296 presents a major challenge for tracing the evolution of regulatory elements across species. Multiple  
297 sequence alignments and alignment-free algorithms are strategies to identify orthologous pairs of  
298 regulatory sequences, but are challenging, especially for large evolutionary distances. Efforts such as  
299 halliftover/HALPER (Zhang et al. 2020; Hickey et al. 2013) try to overcome this based on multiple  
300 alignment of hundreds of genomes. However, their performance is similar to IPP using only 16  
301 genomes, highlighting the potential of synteny as a proxy for conservation. Bridged/tunneled  
302 alignments (Taher et al. 2011; Baranasic et al. 2022) provide a viable strategy for orthologous CRE  
303 detection and have already indicated a higher degree of CRE conservation than commonly assumed.  
304 Our approach builds on the idea of bridged alignments and extends it in several ways. One, IPP  
305 implements multiple bridging species, which can be optimized for any pairwise comparison based on  
306 their specific phylogenetic relationships. Two, within the framework of conserved synteny, IPP

307 projections can assume orthology for any pair of regions between any two genomes, irrespective of  
308 their DNA sequence. Consequently, in non-syntenic regions, or between very distantly related  
309 genomes (Sanges et al. 2006) this strategy might miss orthologous elements. Nevertheless, IPP is a  
310 potent approach to identify putative orthologs for comparative studies at varying evolutionary  
311 distances provided the appropriate set of bridging species, in particular when combined with  
312 equivalent experimental data sets similar to our mouse and chicken heart data. Moreover,  
313 identification of indirectly conserved elements provides valuable information for interpretation of  
314 disease-associated non-coding variants in humans, for example in congenital heart disease (Richter et  
315 al. 2020; Xiao et al. 2024), and facilitates their functional characterization and testing in animal models.

316 Advances in machine learning have made it possible to predict the regulatory activity for any DNA  
317 sequence in a given cell type or tissue (Avsec et al. 2021; de Almeida et al. 2022; Reiter, de Almeida,  
318 and Stark 2023; de Almeida et al. 2023; Taskiran et al. 2023). Within mammals, models trained in one  
319 species can successfully predict activity in another (Minnoye et al. 2020; Kaplow et al. 2023; Kliesmete  
320 et al. 2024), but cannot match ortholog pairs. A recent study aimed to identify orthologous enhancers  
321 between human and mouse using a ML model, but requires syntenic regions as part of their algorithm  
322 to match orthologs (Oh and Beer 2023). Here we show that our SVM model trained in mouse can  
323 predict tissue-specific enhancers in chicken, highlighting the deep conservation of enhancer sequence  
324 syntax. Going beyond its predictive power, we use the model to independently validate IPP-projected  
325 regions in the chicken genome, demonstrating that indirectly conserved regions have sequence  
326 characteristics typical of heart enhancers. In the future, combination of both approaches can be a  
327 powerful strategy to study enhancer evolution. For example, IPP-identified pairs of orthologs can serve  
328 as training input for ML models to learn sequence changes compatible with functional conservation.

329 Sequence conservation of CREs, especially that of enhancers, displays a great level of heterogeneity  
330 ranging from ultra-conserved elements (Snetkova et al. 2021; Dickel et al. 2018; Snetkova et al. 2022)  
331 to the sequence-divergent IC elements we describe here. We show, however, that signals for functional  
332 conservation, in terms of chromatin signatures, encoded TFBS, and predicted tissue-specificity is  
333 relatively uncoupled from sequence conservation. As such, we imagine IPP to be an efficient approach  
334 to annotate orthologous CREs between species for example in single-cell ATAC-/ChIP-seq datasets from  
335 equivalent tissues, where cell types and expression programs are conserved, while the majority of  
336 CREs currently appear to be non-conserved.

337 Furthermore, the TFBS shuffling analysis suggests that CRE function may predominantly be maintained  
338 by TFBS composition. Consequently, conservation of an element's regulatory function is less  
339 dependent on exact sequence conservation than on preserving the appropriate balance of TFBSs  
340 within the given element. Given that we found thousands of IC elements between mouse and chicken,

341 the functional conservation of CREs across larger evolutionary distances is likely much more prevalent  
342 than currently appreciated.

343  
344

#### 345 **Code Availability**

346 The source code for Interspecies Point Projection from this study can be obtained from  
347 <https://github.com/tobiaszehnder/ipp>

348

#### 349 **Acknowledgements**

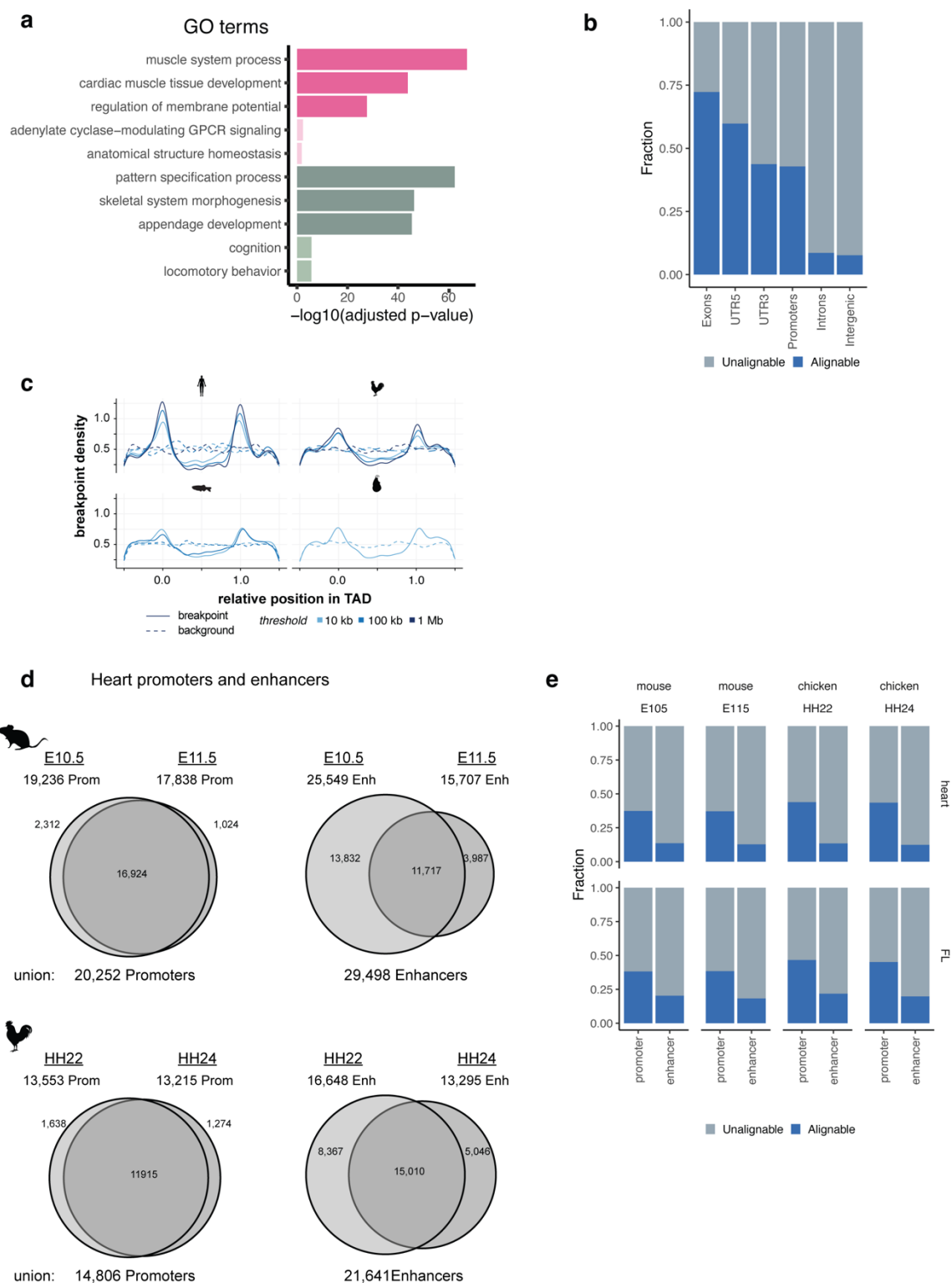
350 D.M.I. and M.P. were supported by funding from the DFG SPP 22.02 “3D Genome Architecture in  
351 Development and Disease” (IB139/1-1 and IB 139/6-1). Work in the Ibrahim Lab is supported by an  
352 ERC Starting Grant SYNREG (101076709). We thank the MPI-MG transgene facility and animal house  
353 for husbandry and members of Dominik Seelow and Martin Kircher’s labs for feedback on the Machine  
354 Learning analysis. We would like to thank Juliane Glaser, Alicia Madgwick and all members of the  
355 Ibrahim lab for feedback on the manuscript.

356

357

358

359

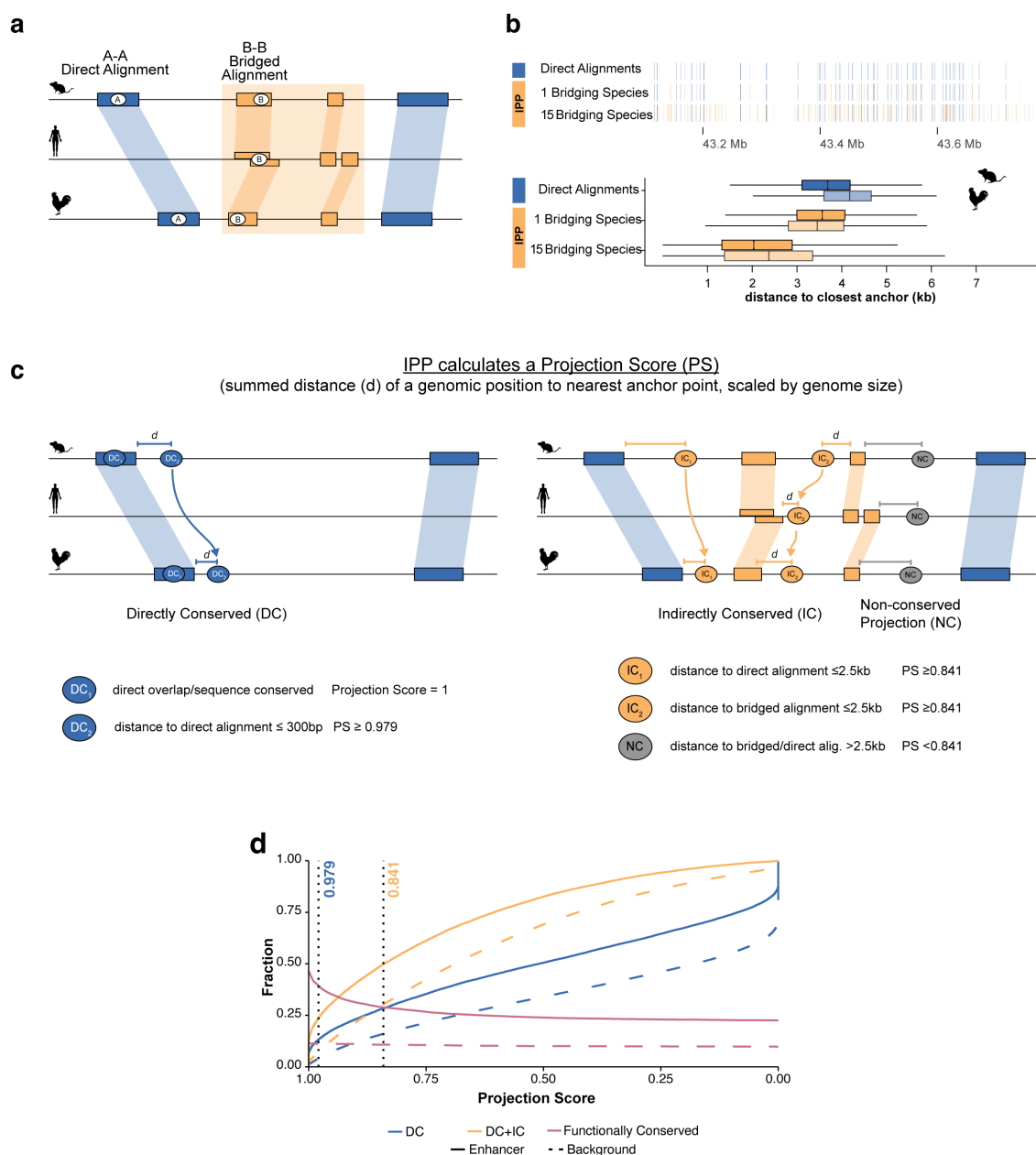


**Fig. S1**

**Figure S1: (a)** Gene Ontology (GO) annotations of differentially expressed genes (Heart vs. FL) in mouse and chicken. Dark pink = upregulated, both species. Dark green = downregulated, both species. Light pink = upregulated, mouse-only. Light green = upregulated, chicken-only. Grey = no differential expression. **(b)** Estimation of sequence alignability of ATAC-seq peaks from mouse embryonic heart at different annotated genomic locations. **(c)** Number of predicted promoters and enhancers from stage-specific and shared/union sets in both species. **(d)** Estimation of sequence alignability from stage-specific predicted promoters and enhancers from heart and FL in both species.

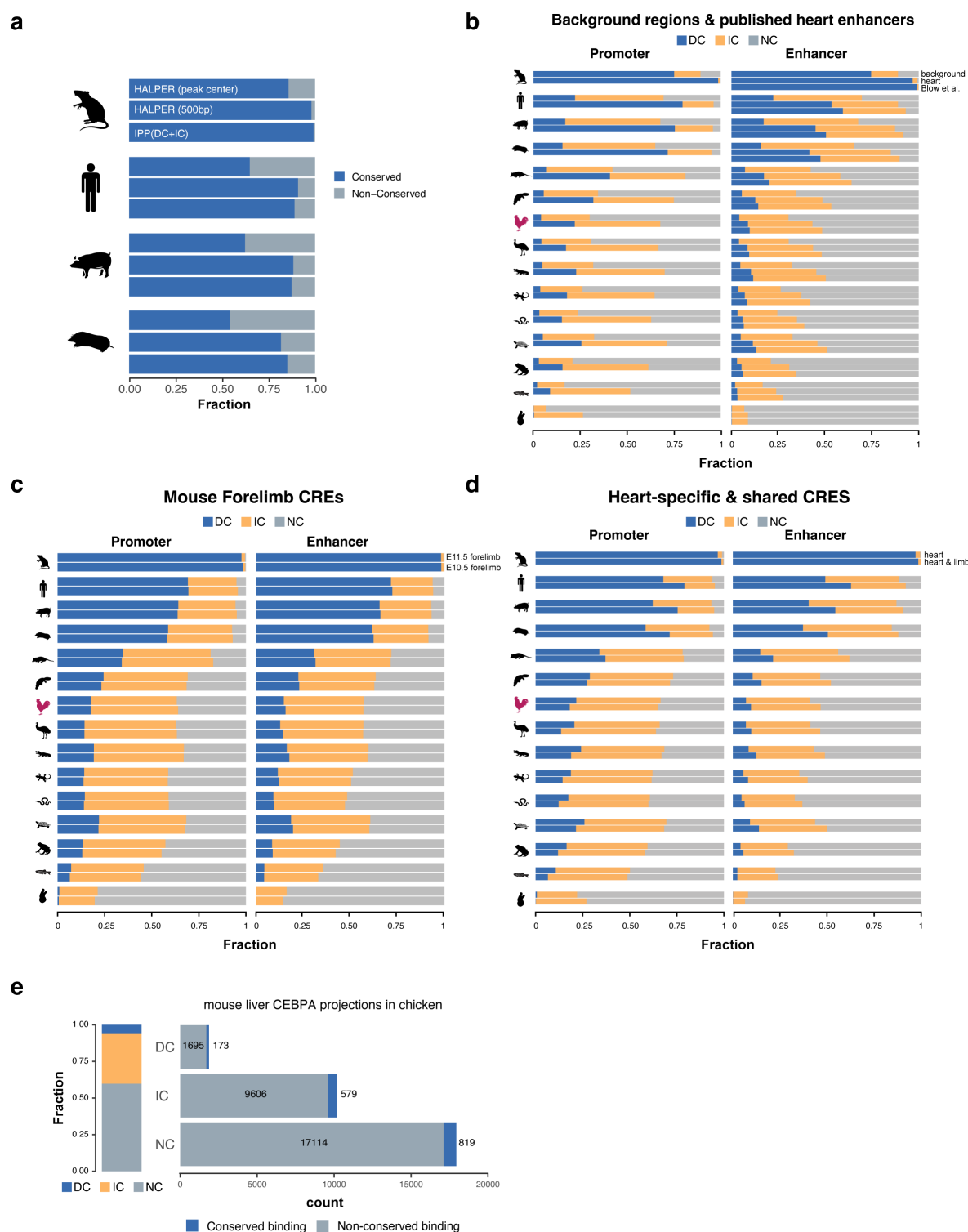
360  
361  
362  
363  
364  
365  
366  
367





**Fig. S2**

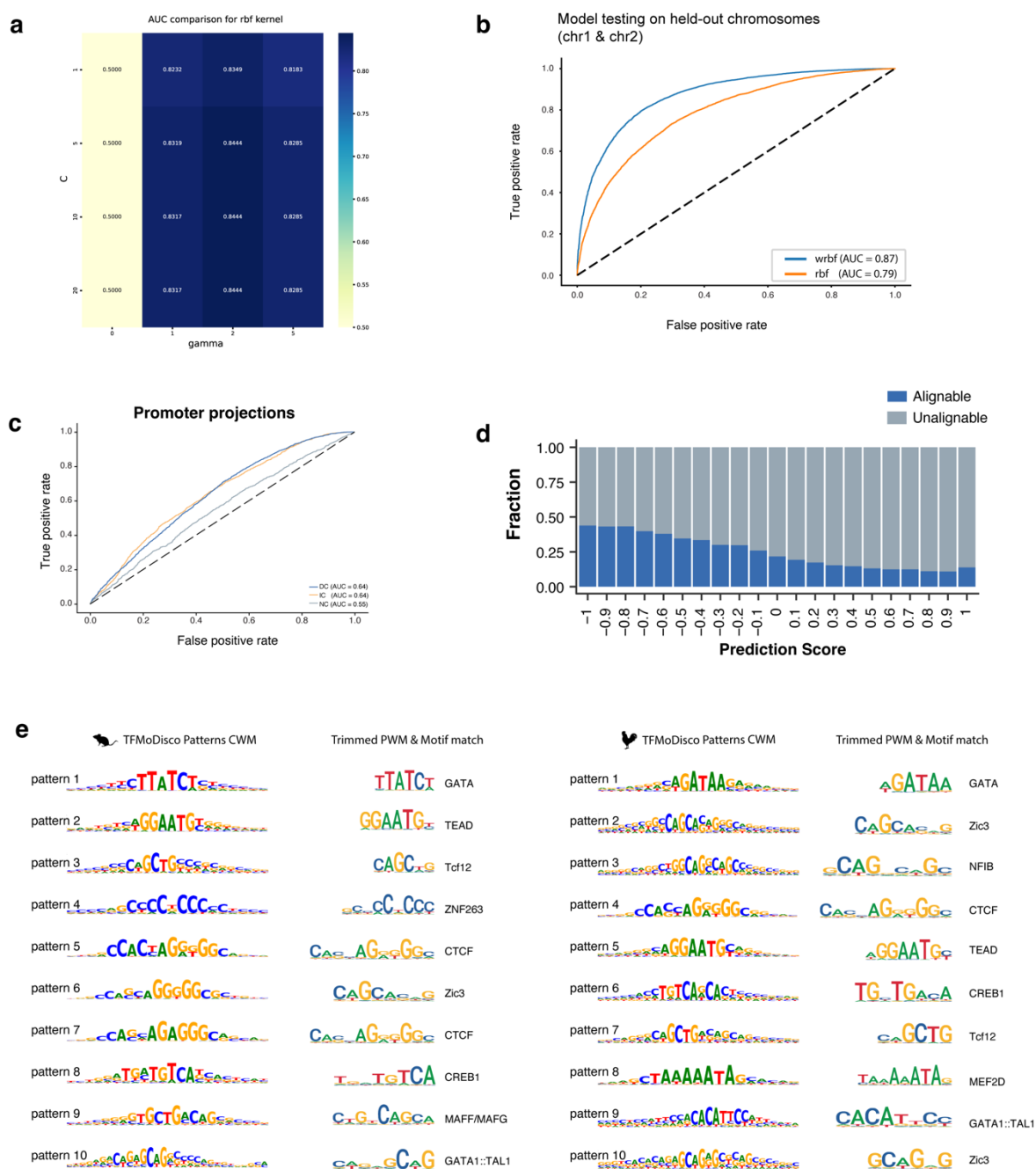
368  
 369 **Figure S2: Interspecies Point Projection combines bridged alignments and synteny to identify orthologous regions (a)**  
 370 Classification of direct and bridged alignments through the use of intermediate species (b) Increase in the number of anchor  
 371 points and distance to the nearest anchor points through multi-species bridged alignments. Comparison between 0, 1 and  
 372 15 bridging species (c) Classification of projections as directly and indirectly conserved. DC regions overlap a sequence  
 373 alignment or are  $\leq 300$ bp from a direct alignment. The distance of IC regions as  $>300$ bp but  $\leq 2.5$ kb from a direct or indirect  
 374 alignment. Regions with  $>2.5$ kb summed distance through the species graph from anchor points are classified as NC. (d)  
 375 Fractions of mouse enhancers identified as directly conserved (DC, blue) or either directly or indirectly conserved (DC + IC,  
 376 orange) as a function of the projection score threshold. Fraction of functionally conserved DC+IC elements as a function of  
 377 the projection score threshold (red). Solid lines = enhancers, dashed lines = randomly selected background regions. Dotted  
 378 vertical lines represent DC threshold score of 0.979 and IC of 0.841.  
 379



**Fig. S3**

**Figure S3** (a) IPP performance compared to halliftover/HALPER for mouse heart enhancer ortholog prediction in four placental mammals. (b-d) IPP projections for randomly selected genomic regions and published heart enhancers (Blow et al 2010) (b), forelimb CREs at E10.5 & E11.5 (c), and heart-specific or heart and limb CREs (d). (e) Fraction of directly/indirectly conserved mouse CEBP/A ChIP-seq peaks in the chicken genome. Blue fractions (right) show the number of conserved binding events (as determined by overlap with a CEBP/A ChIP-seq peak in chicken livers)

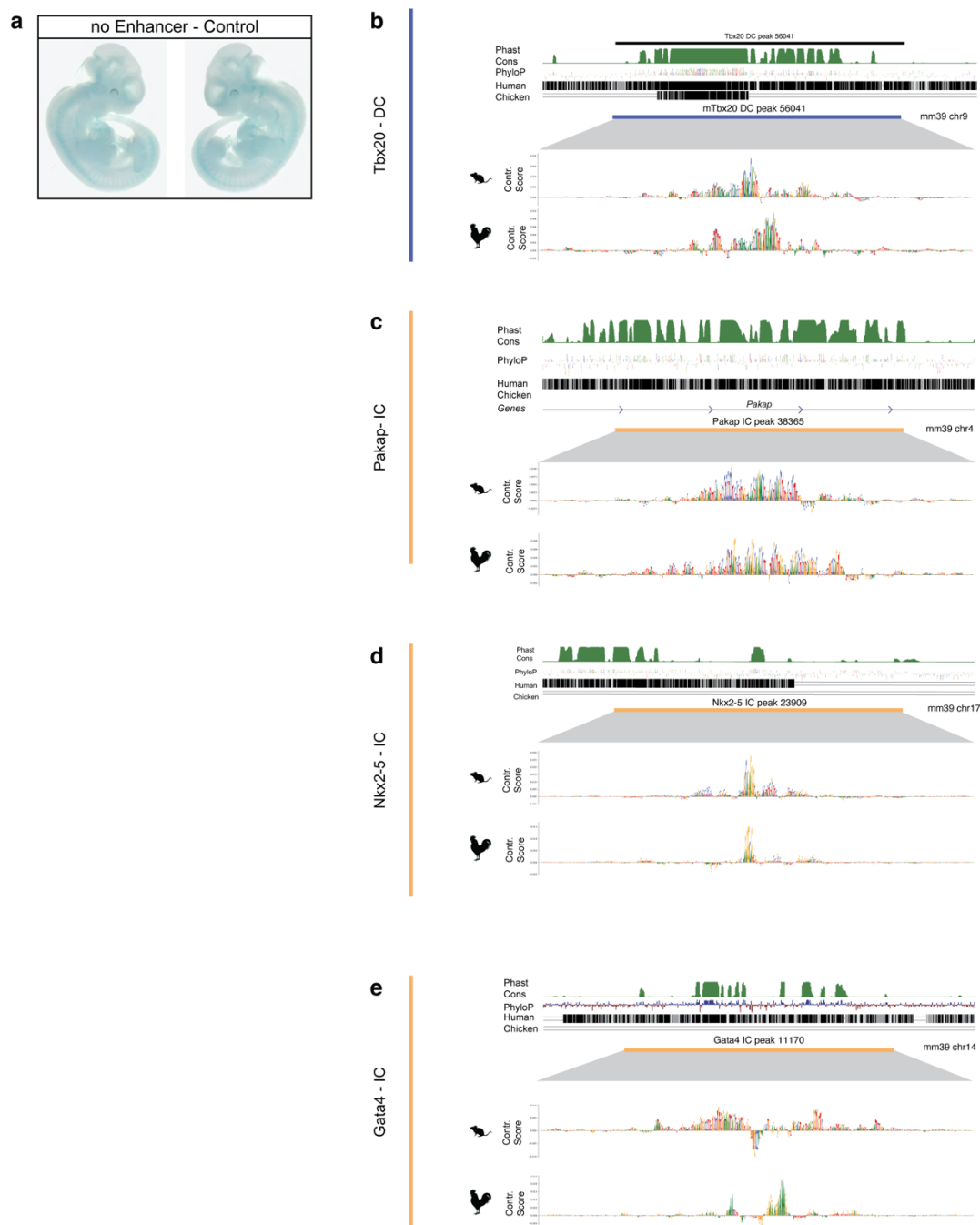
380  
 381  
 382  
 383  
 384  
 385  
 386  
 387  
 388



**Fig. S4**

**Figure S4 (a)** Parameter tuning to train the SVM with RBF kernel with a grid-search for parameters  $c$  and  $\gamma$  showing the calculated AUC after 5-fold cross validation. AUC = Area under the ROC curve. **(b)** ROC curves with computed AUC showing the performance of gkm-SVM with either RBF(rbf, orange) or weighted RBF(wrbf, blue) kernel on test data. The SVM was trained with the  $c$  &  $\gamma$  parameters chosen in (a). **(c)** ROC curves with computed AUC showing human-chicken interspecies prediction accuracy for different conservation classes of mouse promoters projected to chicken. **(d)** Estimation of sequence alignability as a function of SVM predicted tissue-specificity (as prediction score) for ATAC-Seq peaks from chicken embryonic heart. **(e)** Top 10 mouse (left) and chicken (right) patterns discovered by TF-MoDisco showing seqlet as CWM, trimmed and converted PWMs and their annotated JASPAR motif match.

389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401



**Fig. S5**

402  
 403 **Figure S5 (a)** Control for Enhancer reporter using a knock-in of the minimal promoter-lacZ without enhancer. Background  
 404 signal in somites along the antero-posterior axis. **(b-e)** Sequence conservation scores (PhastCons/PhyloP) and direct  
 405 alignments to human and chicken of all tested enhancers. SVM contribution scores show important sequence features of the  
 406 500bp enhancers.  
 407

## 408 Supplementary Text

409

410

### Interspecies Point Projection (IPP)

411 We project genomic point coordinates from a reference genome to a target genome by linear interpolation  
412 between blocks of pairwise sequence alignment, so called anchor points (Baranasic et al. 2022). Moreover, we  
413 use pairwise alignments between a set of bridging species to maximize anchor point density and thus optimize  
414 projection accuracy. This scenario is represented by a graph in which every node is a species, and the weighted  
415 edges represent the distance of a genomic coordinate to its anchor points between the nodes it connects (Fig.  
416 2). We established a distance scoring function that returns a score of 1 for a genomic location  $x$  overlapping an  
417 anchor point ( $|x - a| = 0$ ), and exponentially converges to zero with increasing distance  $|x - a|$ . For a single  
418 pairwise comparison, the function is defined as follows:

$$419 \quad f(x) = \exp\left(-\frac{d_{min}}{g_R s}\right), \quad (1)$$

420 with  $d_{min} = \min\{|x - a^{(1)}|, |x - a^{(2)}|\}$  denoting the distance of a genomic location  $x$  to its closest anchor  
421 point,  $g_R$  denoting the genome size of the reference species and  $s$  a scaling factor that can be tweaked to  
422 determine the decreasing rate of the function. For instance, we can set  $s$  by defining a distance half life  $d_h$  as the  
423 distance  $|x - a|$  at which the scoring function ought to return a value of 0.5:

$$424 \quad s = -\frac{d_h}{g_B \log(0.5)}. \quad (2)$$

425 All projections presented in this manuscript were computed using a distance half life of 10 kb.

426 For the score calculation in Equation 1, the distance is normalized by the genome size of the reference species  
427 ( $g_R$ ) of a pairwise comparison. In Equation 2, the scaling factor is normalized by the size of a basis genome ( $g_B$ )  
428 which we chose to be the mouse genome build mm39, allowing comparisons between projections from different  
429 reference species. In practice, this means that the distance scoring function decreases at equal rates for different  
430 reference genomes, however, these scores correspond to different distances based on the relative reference  
431 genome sizes. The function can thus be simplified to the following form:

$$432 \quad f(x) = 0.5^{\left(\frac{d_{min}}{d_h} \frac{g_B}{g_R}\right)}. \quad (3)$$

433

434 We can then compute the total distance score of a given path through the graph as the product of the score of  
435 all edges in that path. The length of a path is reciprocal to the distance scoring function, hence we can subtract  
436 the total score from 1 to obtain the path length  $l_p$ :

$$437 \quad l_p = 1 - \prod_{i \in p} f(x_i). \quad (4)$$

438 Finally, projection accuracy is optimized by finding the shortest path through the graph:

$$439 \quad \hat{p} = \arg \min_{p \in P} l_p, \quad (5)$$

440 with  $P$  denoting the set of all paths through the graph connecting the reference and the target species. Finding  
441 the shortest path through a graph can be solved using Dijkstra's Shortest Path Algorithm (Dijkstra 1959). We  
442 implemented the method in python and C++ and named it Interspecies Point Projection (IPP). IPP is publicly  
443 available at <https://github.com/tobiaszehnder/ipp>.

444

### 445 **Bridging species selection and pairwise alignment**

446 IPP relies on the additional anchor points provided by bridging species to map corresponding genomic locations  
447 between pair of divergent genomes (Fig. S2a,b). As such, the choice of bridging species depends on the specific  
448 comparison of interest. Here, for a mouse-chicken comparison, we selected mammalian species which have  
449 diverged from chicken after mouse (human, pig, mole, opossum, platypus), and those which have diverged from  
450 mouse after chicken (alligator, green anole, snake, turtle) (Supplementary Tab. 2). Additionally, we selected the  
451 rat and emu as two closely related species to mouse and chicken, respectively. Finally, we included frogs,  
452 zebrafish, and the sea squirt as outgroups.

453 Fasta files for all reference genome assemblies were obtained either from DNA Zoo or from NCBI were used as  
454 inputs for pairwise alignments with *lastal*. Chain files were then generated, preprocessed, and merged for each  
455 species pair before combined in one collection of large pairwise alignments and stored in a binary format. This  
456 collection of alignments consisting of the reference, target, and all bridging species is the necessary input for  
457 running IPP.

458

### 459 **Projection classification and distance score tuning**

460 IPP computes a score for every projection from one genome to another through the species graph. As described  
461 above, this score is a representation of the distance to the nearest anchor points, i.e. the higher the score, the  
462 shorter the distance and thus the more accurate the projection. We use this distance score as a threshold to  
463 classify projections into 3 classes: directly-conserved (DC), indirectly-conserved (IC), and non-conserved (NC)  
464 (Fig. S2c). An element is classified as DC if its projection score using only direct alignments was above this  
465 threshold, and as IC if their projection score from bridging alignment is also above such threshold. All remaining  
466 elements are then classified as NC. If ATAC-Seq data is available for the target genome, we further classify each  
467 projection by their functional conservation. Specifically, any projected point is classified as functionally  
468 conserved/'+' or non-functionally conserved/'-', if it is within or outside a 2.5kb distance from an ATAC-seq peak  
469 summit, respectively.

470 Initially, we set the score threshold at 0.99 which, given Equation 2, represents a maximum distance to the next  
471 anchor point of ~150 bps for DC elements. For IC elements, this additionally means that the sum of distances  
472 from the query element to an anchor point at all intermediate projections is  $\leq$  150bp. While ensuring the  
473 confidence of projections, this rather stringent cutoff implies a certain level of false negatives within NC, i.e.  
474 projections with a projection score below the cutoff that are nevertheless pointing to the correct ortholog.  
475 Furthermore, taking into practical consideration that IPP only maps a single base-pair of an element between  
476 genomes, such stringency likely results in an underestimation of conservation of the element of interest.

477 We then seek to tune this threshold parameter, which is ultimately a trade-off between specificity and sensitivity.  
478 In other words, relaxing the distance threshold will result in more elements being classified as conserved with a  
479 higher likelihood of such classification being a false positive (i.e. a projection pointing to a non-orthologous  
480 region. We took advantage of available ATAC-Seq data in the chicken forelimb as an independent tissue model,  
481 and determine if and how the proportion of functionally conserved elements changes as we relax the cut-off  
482 score. We observe a clear drop in the fraction of functionally conserved elements at high projection scores (i.e.  
483  $\geq$ 0.9) from 38% to below 27% of all conserved enhancers (Fig. S2d). This sharp change in proportion appears to  
484 plateau at lower score thresholds. Indeed, even with dramatically forgiving cutoffs, just over 1/5th of all  
485 projections is putatively functionally conserved at every score threshold below 0.75. Importantly, this trend is  
486 not reflective of the spatial distribution of open chromatin, as only ~10% of randomly selected background  
487 regions reside within open chromatin after projection (Fig. S2d).

488 One can take advantage of such a relaxed approach to identify putative *functionally* conserved orthologous  
489 enhancers (e.g. projected elements that are residing in open chromatin), providing an additional layer of  
490 functional validation. Given the availability of equivalent functional datasets, we decided to relax this cut-off and  
491 used 2 different distance cut-offs for DC and IC classifications. Specifically, an element with a projection score of  
492 0.979 (~300bp distance) using only direct alignments is classified as DC. For IC classification, we used a score cut-  
493 off of 0.841, which is equivalent to a summed distance of 2.5kb through all intermediate projections. These  
494 projections are filtered - as before - for those overlapping open chromatin regions to select for putatively  
495 functionally conserved elements. This permits the detection of functional orthologs in highly dynamic genomic  
496 neighborhoods where sequence alignments are sparse, with the potential cost of a higher false discovery rate.

497

## 498 **Materials and Methods**

### 499 **Biological samples**

500 C57/BL6 inbred mice were used for timed mating and fertilized SPF eggs (Valo Biomedica) were incubated at 38°C  
501 50-55% humidity. Embryonic hearts and forelimbs from mouse and chicken embryos (E10.5, E11.5 and HH22,  
502 HH24) were dissected and further processed for sequencing libraries preparation. Each experiment was  
503 performed in biological replicates.

504

### 505 **Sample and Sequencing libraries preparation**

#### 506 a. RNA-seq

507 For RNA-seq, dissociated chicken embryonic heart cells were snap-frozen in liquid N<sub>2</sub>. RNA was extracted using  
508 the Qiagen RNeasy-Mini Kit according to manufacturer's instructions. Ribosomal RNA was depleted before library  
509 preparation with the Kapa HyperPrep Kit and sequenced on a Novaseq2 100 bp paired-end reads. RNA-seq  
510 experiments were performed in duplicates.

511

#### 512 b. ATAC-seq

513 ATAC-seq libraries were prepared using the Omni-ATAC protocol from 50k cells per replicate. Embryonic tissues  
514 were dissociated into single-cell suspension, washed with cold PBS, and lysed in fresh lysis buffer (10mM TrisCl  
515 pH7.4, 10mM NaCl, 3mM MgCl<sub>2</sub>, 0.1% (v/v) Igepal CA-630) on ice. Tn5 transposition for lysed nuclei was done  
516 for 30 min at 37° C, and DNA was then purified using the MinElute Reaction Cleanup kit (Qiagen) kit.  
517 Nextera indexing primers were added during library amplification from purified DNA, where the number of cycles  
518 were determined by qPCR as described. After double-sided size selection, we verified the expected nucleosomal  
519 fragment distribution with a BioAnalyzer or TapeStation. DNA concentration of libraries were measured with  
520 Qubit HS before sequencing on a Novaseq2 (Illumina) using 100bp paired-end reads.

## 521

### 522 c. ChIPmentation

523 ChIPmentation libraries were prepared as previously described (Schmidl et al. 2015). Briefly, dissociated cells  
524 were first filtered through a 100µm (embryonic heart) or 70µm (limb) MACS® SmartStrainer before fixation with  
525 1% MeOH-free formaldehyde (Thermo Scientific: 28906) in PBS on ice for 10 minutes. Fixed cells were first  
526 quenched using glycine and then lysed on ice in lysis buffer (10mM Tris pH 8.0, 100mM NaCl, 1mM EDTA pH 8.0,  
527 0.5mM EGTA, 0.1% Sodium deoxycholate, 0.5% N-lauroylsarcosine) before shearing with a Covaris E220 for a  
528 fragment distribution of 200-700bp. Antibodies were incubated overnight at 4C, followed by  
529 immunoprecipitation with protein G beads (id). After beads washing, transposition/'tagmentation' reaction with  
530 the Tn5 transposase was done at 37C for 5min. Beads were then again washed before overnight reverse  
531 crosslinking with Proteinase K. DNA was then purified using the MinElute Reaction Cleanup kit (Qiagen).  
532 Libraries were indexed and amplified similarly as previously described for ATAC-Seq libraries. The number of PCR  
533 cycles for each library was estimated using Ct values as determined by qPCR (where number of cycles = rounded  
534 up Ct value +1). After amplification, DNA was cleaned up with AmPureXP beads, and then checked on a  
535 TapeStation D5000 HS for size distribution. Size selection was then carried out accordingly. The concentration of  
536 final eluted DNA was measured using Qubit HS and checked again on a TapeStation D5000HS. All libraries were  
537 sequenced on a Novaseq2 (Illumina) using 100bp paired-end reads.

### 538

### 539 d. Hi-C

540 In situ Hi-C libraries were prepared as previously described (Schöpflin et al. 2022). Briefly, 3C libraries were  
541 digested with DpnII, and digested ends were marked with biotin-14-dATP. DNA was sheared with an S-Series 220  
542 Covaris to 300-600bp fragments before biotin pull-down using Dynabeads MyOne Streptavidin T1 beads.  
543 Sheared DNA ends were then repaired with T4 DNA polymerase and the Klenow fragment of DNA polymerase I,  
544 and subsequently and phosphorylated with T4 Polynucleotide Kinase NK. Sequencing adaptors were then added,  
545 and libraries were indexed via PCR amplification (4–8 cycles) using the NEBNext Ultra II Q5 Master Mix. PCR  
546 clean-up and size selection were done with AmPureXP beads before 100bp paired-end sequencing on a  
547 Novaseq2.

## 548

## 549

## 550 Data processing

### 551 a. RNA

552 We processed all RNA-seq libraries with STARv2.7.9a using reference genome sequences and annotations from  
553 GENCODE (vM32, primary) for mouse and Ensembl (GRCg7b) for chicken. We obtained gene-level counts for  
554 each sample with `--quantMode geneCounts`. In addition to in-house chicken heart RNA-seq libraries, we similarly  
555 processed the following publicly available datasets: mouse heart E10.5 & E11.5 (ENCODE3), chicken FL HH22 &  
556 HH24 ([GSE164737](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE164737), (Jhanwar et al. 2021). TPM values were computed from gene-level counts, where gene length  
557 is estimated as the sum of all exon lengths.

### 558 b. ATAC-seq & ChIPmentation

559 For ATAC-seq & ChIPmentation samples, Nextera Tn5 adaptor sequences were trimmed from fastq reads using  
560 `cutadapt` before further processing. Reads were aligned to appropriate reference genomes (mm10, mm39, or  
561 galGal6) using `bowtie2 v2.3.5.1` where the maximum fragment length set was either 1000bp (ATAC) or 700bp  
562 (ChIPmentation). Duplicated reads were then removed using `MarkDuplicates` (Picard v2.23.4). Finally, reads were  
563 further sorted and filtered using `samtools v1.10` to remove unmapped reads, low quality reads (MAPQ < 10),  
564 and mitochondrial reads. Filtered bam files from replicates were merged to generate bigwig files. We used  
565 `bamCoverage` (deepTools) with CPM normalization and bin size of either 1 for ATAC or 10 for ChIPmentation.

566 Peak calling for ATAC-seq data from replicates was done with Genrich v0.6.1 in ATAC mode ‘-j’ with default  
567 parameters. (<https://github.com/jsh58/Genrich>).

#### 568 c. HiC

569 Reads handling were done using Juicer v1.6.0 CPU version (Durand NC, et al. 2016). Specifically, alignment was  
570 done with BWA-MEM v0.7.17 to reference genome galGal6. Only read pairs with MAPQ > 30 were included in  
571 the final contact maps. Processing was done separately for each replicate, and output filtered de-duplicated read  
572 pairs were merged. Contact matrices were balanced with Knight-Ruiz normalization (Knight PA & Ruiz D. 2012)  
573 before visualization.

574

### 575 Data analysis

#### 576 a. Comparative differential expression analysis

577 Raw gene-level counts from heart and limb samples at both stages were used as input for differential analysis  
578 with DESeq2 [v1.36](Love, Huber, and Anders 2014). We obtain a set of differentially expressed genes in the  
579 heart relative to limb in both stages, accounting for the effects for biological replicates. To aid visualization and  
580 gene ranking for Gene Ontology (GO) analysis, effect size shrinkage was done for the coefficient modeling tissue-  
581 specificity (i.e. *tissue\_heart\_vs\_limb*).

582 Gene orthology annotations were obtained from Ensembl databases GRCm39 for mouse and GRCg7b for chicken.  
583 Duplicated annotations were filtered to retain only those with the highest GOC score. Only one-to-one  
584 orthologous genes (OGs) were used for all comparative analysis.

585 Gene Ontology (GO) analysis was done using R package clusterProfiler [v4.4.4] (Wu et al. 2021). Over-  
586 representation GO analysis of OGs was done given a background gene set of all detectably expressed mouse  
587 genes (i.e. raw counts  $\geq 10$ ). For statistical testing, testing gene-set sizes were set from a minimum of 5 to a  
588 maximum of 100 genes to allow focusing of specific biological processes (i.e. BP) over more general terms.

589

#### 590 b. Estimation of sequence alignability

591 To estimate conservation by means of sequence alignability, we used UCSC LiftOver as implemented within R  
592 package *rtracklayer* for reciprocal mapping between mouse and chicken genomes. Chain files for mm39 and  
593 galGal6 were obtained from UCSC before importing into R using *rtracklayer*. For mapping, we used the default  
594 parameter settings (minMatch=0.1) and allowed for multiple mapping (i.e. one-to-many) between query and  
595 target.

596

#### 597 c. Enhancer & Promoter prediction

598 Histone profiles (i.e. H3K27ac, H3K4me1, H3K4me3) from merged replicates were used to predict candidate  
599 regulatory regions using the enhancer prediction tool CRUP (Ramisch et al. 2019). In brief, CRUP computes the  
600 probability score for each 100bp bin in the entire genome to be an active enhancer element. Combining these  
601 probabilities and normalized histone signal values (i.e. mono:tri ratios), bins are filtered and merged into either  
602 promoter-like or enhancer-like regions.

603 To define active promoter regions, we intersected defined promoter-like regions with all TSS of actively  
604 transcribed genes (i.e. counts  $\geq 1$ TPM). Counts values were obtained as described previously for expression  
605 analysis. Once promoters are defined, we finalized the set of active enhancers by filtering enhancer-like regions  
606 by their accessibility from called ATAC peaks. Finally, those falling within 2kb of a predicted promoter are  
607 removed from the final set of active enhancers. The numbers of enhancers and promoters can be found in Sup.  
608 Tab. 1 and the bed files under GSE263587, GSE263753, GSE263755, GSE263783.

609

610

### 611 TFBS Motif and Foot-printing Analysis

#### 612 a. Reference motif collection

613 We obtained TF motif models from the JASPAR 2022 database (core vertebrate, non-redundant) and  
614 systematically curate this database to be used as a reference for all TFBS-based analysis. From over 700 JASPAR  
615 TF motifs, we filtered for those with detectable expression in the mouse embryonic heart by integrating RNA-  
616 seq counts (described above). Detectable expression is defined as having counts of  $\geq 1$  TPM in both replicates,  
617 in either stage E10.5 or E11.5 (n=520). From these, we further consolidate the reference collection by filtering



618 out redundant motifs based on sequence similarity within the same annotated TF family. Specifically, within each  
619 TF family, motifs are ranked by their informational content score before pair-wise comparison with others in the  
620 same family using the *compare\_motifs* function from R package *universalmotif*. Finally, motifs with lower  
621 informational content score and a similarity score > 0.9 (score of 1 = identical sequence) are discarded from the  
622 final reference set (n=301).

#### 623 b. Motif scanning

624 To characterize the TFBS composition of CREs, we searched for motif matches from the curated collection using  
625 FIMO implemented through R package *memes* (Grant, Bailey, and Noble 2011) with default parameters. DNA  
626 sequences were obtained from annotation packages *BSgenome.Mmusculus.UCSC.mm39*  
627 *BSgenome.Ggallus.UCSC.galGal6* for mouse and chicken, respectively. Motifs scanning was done within a 500bp  
628 window centered by ATAC peak summit or projected point. Peak centering by summit was done for projected  
629 regions in chicken only for functionally conserved elements (i.e. DC+ & IC+). Finally, any overlapping hits from  
630 the same motifs are discarded, keeping the match with higher score.

#### 631 c. ATAC-seq foot-printing

632 Aligned ATAC-seq reads from biological replicates were merged to be used as input for ATAC-seq footprinting  
633 analysis using TOBIAS (Bentsen et al. 2020) [v0.3.3]. The genomic regions of interests to be foot-printed were:  
634 (1) the union set of predicted enhancers and promoters in mouse and chicken hearts, and (2) all called chicken  
635 ATAC-seq peaks. Briefly, we used TOBIAS to correct for Tn5 bias before footprint scores were calculated at  
636 genomic regions of interest. Finally, we used our curated set of TFBS motifs as reference to predict TF binding.  
637 TOBIAS output from different stages were merged, and overlapping regions of predicted binding from the same  
638 TF were merged similarly to motif hits as described. Finally, quantification of shared footprints was done similarly  
639 to the motifs analysis previously described.

#### 640 d. Quantification of motifs and TFBS sharing between pairs of orthologous CREs

641 To quantify the similarity between mouse CREs and their corresponding chicken orthologs as determined by IPP,  
642 we determine the total number of shared motifs and TF-binding (i.e. TFBS) between every mouse-chicken pair  
643 of sequences. As a negative control, we also compare the number of shared motifs and TFBS between a mouse  
644 sequence and non-orthologous, i.e. background genomic region. Specifically, for every mouse sequence with a  
645 chicken projection overlapping an ATAC-seq peak (i.e DC+/IC+/NC+), another ATAC-seq peak (if possible, within  
646 the same TAD) is randomly selected as its non-ortholog.

647

### 648 Classification model for heart-specific enhancers

#### 649 a. Training strategy and data preparation

650 Our classification model is a Support Vector Machine (SVM) with a center-weighted radial basis gapped k-mer  
651 kernel function (wrbfgkm) (implemented at <https://github.com/kundajelab/lsgkm-svr>) (Ghandi et al. 2014; Lee  
652 2016). All datasets used for model training are processed bulk ATAC-seq data either obtained from ENCODE or  
653 in-house (as described above). To learn predictive features of heart-specific enhancers, we construct the positive  
654 set to include called ATAC-seq peaks from mouse hearts at 6 developmental stages (**in-house**: E10.5 & E11.5,  
655 **ENCODE**: E12.5-E14.5 & P0). All regions are centered at peak summit and extending 250bp on either side.  
656 Additionally, to ensure the model learns enhancer-specific regulatory features, regions within 2kb of an  
657 annotated mouse promoters (from EPD3 database) were removed from the final training set (n~65k).

658 For model training, we construct the negative set such that the model can accurately learn the sequence features  
659 determining whether an enhancer/CRE is heart-specific. First, to limit confounding factors, we generated a 10-  
660 fold null set of from random genomic loci. From these regions, we filtered for those overlapping any annotated  
661 ENCODE candidate CREs or ATAC-seq peaks from 5 non-heart embryonic organs (limbs, mid-/fore/hind-brain,  
662 liver, E12.5) and mESCs. Finally, those within a 2-kb overlap of any regions from the positive set were removed  
663 (n=70k).

664 All negative sets of GC- and repeats-matched sequences were generated using the *genNullSeqs* function from R  
665 package *gkmSVM* (Ghandi et al. 2014, 2016). Repeats-masked genomic sequences were obtained from custom  
666 masked *BSgenome* data packages for mm10, mm39 or galGal6.

#### 667 b. Hyperparameter tuning & performance evaluation

668 As a measure for classification performance, the area under the ROC curve (AUC) was computed and visualized.  
669 For parameter tuning, a grid search for *C* and *g* parameters for wrbfgkm-kernel was done using a 5-fold cross  
670 validation for each combination of *C* = 1, 5, 10, 20 and *g* = 0, 1, 2, 5 (=16 conditions). The best performing

671 parameter set ( $c=10$ ,  $g=2$ ) as determined by its calculated AUC was chosen for model training. The final model  
672 was tested on positive vs. negative regions on held-out chromosome 1 & 2.

### 673 c. Model prediction on chicken CREs and projections

674 Our heart-enhancer SVM model trained on mouse sequences was used to classify: (1) identified chicken  
675 enhancer and promoter sequences (described previously) from heart and FL, and (2) sequences mouse CREs at  
676 the projected chicken regions from by IPP. For each prediction, the negative set generated as described previously  
677 consists of GC- and repeats-matched regions. Additionally, only projected regions overlapping an ATAC-seq peaks  
678 (i.e. DC+ or IC+) were included in the analysis. AUROCs were computed to evaluate the model's performance on  
679 these regions.

### 680 d. Model interpretation and de novo motifs discovery

681 We used GkmExplain (Shrikumar, Prakash, and Kundaje 2019) (implemented at  
682 <https://github.com/kundajelab/lsgkm-svr>) to interpret the model's classification. GkmExplain computes the  
683 contribution score at each nucleotide to the SVM classification in all input sequences, i.e. its *importance score*.  
684 For each sequence, this importance score was computed by element-wise multiplication of the one-hot encoded  
685 sequence matrix by its hypothetical importance score. Scores were visualized using the *visualization* module  
686 from Python package *modisco*.

687 Computed hypothetical score was then normalized by the ratio of original importance scores and sum of all  
688 hypothetical scores having the same sign. Normalization allows the score to better reflect the importance of a  
689 specific base at each position thereby reducing noise for subsequent motif discovery with TF-Modisco (Shrikumar  
690 et al. 2018) (implemented at <https://github.com/jmschrei/tfmodisco-lite>). Computed and normalized scores  
691 from GkmExplain from: (1) mouse positive test set ( $n=9k$ ), and (2) heart-specific chicken enhancers ( $n=15k$ ) were  
692 used as input for two separate TF-Modisco runs. Similar positive sequence patterns from these were then merged  
693 for the final set of predictive sequence patterns and stored as PWM motifs. Flanking positions with information  
694 content  $< 0.5$  were trimmed from the PWMs before being annotated with known motifs using TOMTOM (Gupta  
695 et al. 2007) with our TF motifs collection as reference.

696

### 697 Quantification of motifs shuffling

698 To quantify the degree of motifs shuffling, we measure the Kendall tau distance ( $K_d$ ) between pairs of reference  
699 mouse sequence and its corresponding chicken orthologs. The Kendall tau distance metric measures the  
700 similarity between two ranking lists by counting the number of transpositions, or swaps, needed between pair  
701 of ranks in one list to achieve the same order from another list. The more similar the two lists are, the smaller  
702 the distance. A pair of mouse-chicken sequences is considered two ranking lists of motifs, where the order of  
703 shared motifs is the ranks. Here, we encode only the 5'-3' order of motifs for mouse sequences as reference and  
704 compare them to both orientations of the chicken sequences.

705 To ensure we faithfully encode the specific order of motifs as ranks, shared motifs obtained previously are further  
706 processed to filter out largely overlapping occurrences from different motifs (minimum overlap of 8bp), again  
707 keeping the hits with the highest mapping score. Additionally, to ensure unique rankings, runs of hits from the  
708 same motif are considered a singular match. Any sequence containing  $>1$  noncontiguous hits from the same  
709 motif (e.g. A,B,C,A,D) is stored as a matrix of ranking lists, where each row represents a unique ranking order  
710 (e.g. 1-A,B,C,D and 2- B,C,A,D). Using R package *rankdist* (Qian and Yu 2019), we compute the normalized  $K_D$   
711 between all unique ranking lists for a mouse-chicken pair, which accounts for varying number of shared motifs  
712 (i.e. list length). Finally, assuming the fewest possible changes have occurred during evolution, we take the  
713 smallest computed  $K_d$  value for every pairwise comparison and compared between conservation classes DC, IC,  
714 and NC. We also described the effect size of sequence conservation on sequence shuffling by computing Cohen's  
715  $d$  using R package *effsize* (Torchiano 2016).

716

717

### 718 In vivo enhancer-reporter assays

719 Transgenic mice carrying the individual mouse or chicken enhancers tested in this study were generated using a  
720 site-specific integration protocol, modified for mouse embryonic stem cells (mESCs). The PhiC31 system used  
721 (Chi et al. 2019) allows precise recombination between two att sites: the attP site inserted in a safe harbour  
722 genomic locus and the attB site in a donor vector. Genomic regions and primers used for generation of Enhancer  
723 Reporters can be found in Supplementary Table 2.

724 First, a master mESC line was established in which an Hsp68::LacZ expression cassette containing the attP site  
725 was inserted into a safe harbour locus (H11) via CRISPR/Cas9 using FuGENE technology (Promega). To create the

726 donor vectors, we cloned each individual enhancer in a vector containing the attB site and a puromycin (Sigma-  
727 Aldrich, P8833) selection marker using Gibson cloning. Subsequently, each resulting donor vector was co-  
728 transfected with the PhiC31 plasmid into the master line using Lipofectamine LTX (Invitrogen), following the  
729 manufacturer's guidelines. The enhancer-reporter mESC lines were cultured and embryos were generated via  
730 tetraploid complementation (Artus and Hadjantonakis 2011). At embryonic day E10.5, the embryos were  
731 harvested and processed for LacZ staining. Briefly, the embryos were kept in the dark at 37°C in LacZ staining  
732 buffer supplemented with 0.5 mg/ml X-gal, 5 mM potassium ferrocyanide and 5 mM potassium ferricyanide.  
733 When the desired staining was achieved, the embryos were washed several times in PBS and then fixed with 4%  
734 PFA/PBS supplemented with 0.2% glutaraldehyde and 5mM EDTA for long-term storage at 4°C. Embryos were  
735 imaged using a SteREO Discovery.V12 microscope with CL9000 cold light source and a Leica DFC420 digital  
736 camera. The embryo genotyping was performed by PCR using primers spanning the expected 5' and 3'  
737 integration junctions to confirm correct integration of the enhancers.

738

739

740

## 741 References

- 742 Almeida, Bernardo P. de, Franziska Reiter, Michaela Pagani, and Alexander Stark. 2022. “DeepSTARR  
743 Predicts Enhancer Activity from DNA Sequence and Enables the de Novo Design of Synthetic  
744 Enhancers.” *Nature Genetics* 54 (5): 613–24.
- 745 Almeida, Bernardo P. de, Christoph Schaub, Michaela Pagani, Stefano Secchia, Eileen E. M. Furlong,  
746 and Alexander Stark. 2023. “Targeted Design of Synthetic Enhancers for Selected Tissues in the  
747 *Drosophila* Embryo.” *Nature*, December, 1–5.
- 748 Armstrong, Joel, Glenn Hickey, Mark Diekhans, Ian T. Fiddes, Adam M. Novak, Alden Deran, Qi Fang, et  
749 al. 2020. “Progressive Cactus Is a Multiple-Genome Aligner for the Thousand-Genome Era.”  
750 *Nature* 587 (7833): 246–51.
- 751 Artus, Jérôme, and Anna-Katerina Hadjantonakis. 2011. “Generation of Chimeras by Aggregation of  
752 Embryonic Stem Cells with Diploid or Tetraploid Mouse Embryos.” In *Transgenic Mouse  
753 Methods and Protocols*, edited by Marten H. Hofker and Jan van Deursen, 693:37–56. Methods  
754 in Molecular Biology. Humana Press.
- 755 Avsec, Žiga, Vikram Agarwal, Daniel Visentin, Joseph R. Ledsam, Agnieszka Grabska-Barwinska, Kyle R.  
756 Taylor, Yannis Assael, John Jumper, Pushmeet Kohli, and David R. Kelley. 2021. “Effective Gene  
757 Expression Prediction from Sequence by Integrating Long-Range Interactions.” *Nature  
758 Methods* 18 (10): 1196–1203.
- 759 Baranasic, Damir, Matthias Hörtenhuber, Piotr J. Balwierz, Tobias Zehnder, Abdul Kadir Mukarram,  
760 Chirag Nepal, Csilla Várnai, et al. 2022. “Multiomic Atlas with Functional Stratification and  
761 Developmental Dynamics of Zebrafish Cis-Regulatory Elements.” *Nature Genetics* 54 (7): 1037–  
762 50.
- 763 Bejerano, G., A. C. Siepel, W. J. Kent, and D. Haussler. 2005. “Computational Screening of Conserved  
764 Genomic DNA in Search of Functional Noncoding Elements.” *Nature Methods* 2 (7): 535–45.
- 765 Bentsen, Mette, Philipp Goymann, Hendrik Schultheis, Kathrin Klee, Anastasiia Petrova, René  
766 Wiegandt, Annika Fust, et al. 2020. “ATAC-Seq Footprinting Unravels Kinetics of Transcription  
767 Factor Binding during Zygotic Genome Activation.” *Nature Communications* 11 (1): 4267.
- 768 Berthelot, Camille, Diego Villar, Julie E. Horvath, Duncan T. Odom, and Paul Flicek. 2017. “Complexity  
769 and Conservation of Regulatory Landscapes Underlie Evolutionary Resilience of Mammalian  
770 Gene Expression.” *Nature Ecology & Evolution* 2 (1): 152–63.
- 771 Blow, Matthew J., David J. McCulley, Zirong Li, Tao Zhang, Jennifer A. Akiyama, Amy Holt, Ingrid Plajzer-  
772 Frick, et al. 2010. “ChIP-Seq Identification of Weakly Conserved Heart Enhancers.” *Nature  
773 Genetics* 42 (9): 818–22.
- 774 Braasch, Ingo, Andrew R. Gehrke, Jeremiah J. Smith, Kazuhiko Kawasaki, Tereza Manousaki, Jeremy  
775 Pasquier, Angel Amores, et al. 2016. “The Spotted Gar Genome Illuminates Vertebrate  
776 Evolution and Facilitates Human-Teleost Comparisons.” *Nature Genetics* 48 (4): 427–37.
- 777 Chi, Xiuling, Qi Zheng, Ruhong Jiang, Ruby Yanru Chen-Tsai, and Ling-Jie Kong. 2019. “A System for Site-  
778 Specific Integration of Transgenes in Mammalian Cells.” *PLoS One* 14 (7): e0219842.
- 779 Crocker, J., and D. L. Stern. 2017. “Functional Regulatory Evolution Outside of the Minimal Even-  
780 Skipped Stripe 2 Enhancer.” *Development* 144 (17): 3095–3101.
- 781 Dickel, Diane E., Athena R. Ypsilanti, Ramón Pla, Yiwen Zhu, Iros Barozzi, Brandon J. Mannion, Yupar S.  
782 Khin, et al. 2018. “Ultraconserved Enhancers Are Required for Normal Development.” *Cell* 172  
783 (3): 491-499.e15.
- 784 Dijkstra, E. W. 1959. “A Note on Two Problems in Connexion with Graphs.” *Numerische Mathematik* 1  
785 (1): 269–71.
- 786 Engström, Pär G., Shannan J. Ho Sui, Oyvind Drivenes, Thomas S. Becker, and Boris Lenhard. 2007.  
787 “Genomic Regulatory Blocks Underlie Extensive Microsynteny Conservation in Insects.”  
788 *Genome Research* 17 (12): 1898–1908.
- 789 Firulli, A. B., D. G. McFadden, Q. Lin, D. Srivastava, and E. N. Olson. 1998. “Heart and Extra-Embryonic  
790 Mesodermal Defects in Mouse Embryos Lacking the BHLH Transcription Factor Hand1.” *Nature  
791 Genetics* 18 (3): 266–70.

- 792 Fisher, Shannon, Elizabeth A. Grice, Ryan M. Vinton, Seneca L. Bessling, and Andrew S. McCallion. 2006.  
793 "Conservation of RET Regulatory Function from Human to Zebrafish without Sequence  
794 Similarity." *Science (New York, N.Y.)* 312 (5771): 276–79.
- 795 Galupa, Rafael, Gilberto Alvarez-Canales, Noa Otilie Borst, Timothy Fuqua, Lautaro Gandara, Natalia  
796 Misunou, Kerstin Richter, et al. 2023. "Enhancer Architecture and Chromatin Accessibility  
797 Constrain Phenotypic Space during Drosophila Development." *Developmental Cell* 58 (1): 51-  
798 62.e4.
- 799 Ghandi, Mahmoud, Dongwon Lee, Morteza Mohammad-Noori, and Michael A. Beer. 2014. "Enhanced  
800 Regulatory Sequence Prediction Using Gapped K-Mer Features." *PLoS Computational Biology*  
801 10 (7): e1003711.
- 802 Ghandi, Mahmoud, Morteza Mohammad-Noori, Narges Ghareghani, Dongwon Lee, Levi Garraway,  
803 and Michael A. Beer. 2016. "GkmSVM: An R Package for Gapped-Kmer SVM." *Bioinformatics*  
804 (*Oxford, England*) 32 (14): 2205–7.
- 805 Grant, Charles E., Timothy L. Bailey, and William Stafford Noble. 2011. "FIMO: Scanning for  
806 Occurrences of a given Motif." *Bioinformatics* 27 (7): 1017–18.
- 807 Gupta, Shobhit, John A. Stamatoyannopoulos, Timothy L. Bailey, and William Stafford Noble. 2007.  
808 "Quantifying Similarity between Motifs." *Genome Biology* 8 (2): R24.
- 809 Hare, Emily E., Brant K. Peterson, Venky N. Iyer, Rudolf Meier, and Michael B. Eisen. 2008. "Sepsid Even-  
810 Skipped Enhancers Are Functionally Conserved in Drosophila despite Lack of Sequence  
811 Conservation." *PLoS Genetics* 4 (6): e1000106.
- 812 Harmston, Nathan, Elizabeth Ing-Simmons, Ge Tan, Malcolm Perry, Matthias Merckenschlager, and  
813 Boris Lenhard. 2017. "Topologically Associating Domains Are Ancient Features That Coincide  
814 with Metazoan Clusters of Extreme Noncoding Conservation." *Nature Communications* 8 (1):  
815 441.
- 816 Heikinheimo, M., J. M. Scandrett, and D. B. Wilson. 1994. "Localization of Transcription Factor GATA-4  
817 to Regions of the Mouse Embryo Involved in Cardiac Development." *Developmental Biology*  
818 164 (2): 361–73.
- 819 Hickey, Glenn, Benedict Paten, Dent Earl, Daniel Zerbino, and David Haussler. 2013. "HAL: A  
820 Hierarchical Format for Storing and Analyzing Multiple Genome Alignments." *Bioinformatics*  
821 (*Oxford, England*) 29 (10): 1341–42.
- 822 Irie, Naoki, and Shigeru Kuratani. 2011. "Comparative Transcriptome Analysis Reveals Vertebrate  
823 Phylotypic Period during Organogenesis." *Nature Communications* 2 (1): 248.
- 824 Jhanwar, Shalu, Jonas Malkmus, Jens Stolte, Olga Romashkina, Aimée Zuniga, and Rolf Zeller. 2021.  
825 "Conserved and Species-Specific Chromatin Remodeling and Regulatory Dynamics during  
826 Mouse and Chicken Limb Bud Development." *Nature Communications* 12 (1): 5685.
- 827 Jin, Sheng Chih, Jason Homsy, Samir Zaidi, Qiongshi Lu, Sarah Morton, Steven R. DePalma, Xue Zeng,  
828 et al. 2017. "Contribution of Rare Inherited and de Novo Variants in 2,871 Congenital Heart  
829 Disease Probands." *Nature Genetics* 49 (11): 1593–1601.
- 830 Kaplow, Irene M., Alyssa J. Lawler, Daniel E. Schäffer, Chaitanya Srinivasan, Heather H. Sestili, Morgan  
831 E. Wirthlin, Badoi N. Phan, et al. 2023. "Relating Enhancer Genetic Variation across Mammals  
832 to Complex Phenotypes Using Machine Learning." *Science (New York, N.Y.)* 380 (6643):  
833 eabm7993.
- 834 Kikuta, H., M. Laplante, P. Navratilova, A. Z. Komisarczuk, P. G. Engström, D. Fredman, A. Akalin, et al.  
835 2007. "Genomic Regulatory Blocks Encompass Multiple Neighboring Genes and Maintain  
836 Conserved Synteny in Vertebrates." *Genome Research* 17 (5): 545–55.
- 837 Kliesmete, Zane, Peter Orchard, Victor Yan Kin Lee, Johanna Geuder, Simon M. Krauß, Mari Ohnuki,  
838 Jessica Jocher, Beate Vieth, Wolfgang Enard, and Ines Hellmann. 2024. "Evidence for  
839 Compensatory Evolution within Pleiotropic Regulatory Elements." *BioRxiv*.  
840 <https://doi.org/10.1101/2024.01.10.575014>.

- 841 Kuhn, R. M., D. Karolchik, A. S. Zweig, T. Wang, K. E. Smith, K. R. Rosenbloom, B. Rhead, et al. 2009.  
842 "The UCSC Genome Browser Database: Update 2009." *Nucleic Acids Research* 37 (Database  
843 issue): D755-61.
- 844 Lee, Dongwon. 2016. "LS-GKM: A New Gkm-SVM for Large-Scale Datasets." *Bioinformatics (Oxford,  
845 England)* 32 (14): 2196-98.
- 846 Love, Michael, Wolfgang Huber, and Simon Anders. 2014. "Moderated Estimation of Fold Change and  
847 Dispersion for RNA-Seq Data with DESeq2." *Genome Biology* 15 (12): 550.
- 848 Ludwig, M. Z., N. H. Patel, and M. Kreitman. 1998. "Functional Analysis of Eve Stripe 2 Enhancer  
849 Evolution in Drosophila: Rules Governing Conservation and Change." *Development  
850 (Cambridge, England)* 125 (5): 949-58.
- 851 Madgwick, Alicia, Marta Silvia Magri, Christelle Dantec, Damien Gailly, Ulla-Maj Fiuza, Léo Guignard,  
852 Sabrina Hettinger, Jose Luis Gomez-Skarmeta, and Patrick Lemaire. 2019. "Evolution of  
853 Embryonic Cis-Regulatory Landscapes between Divergent Phallusia and Ciona Ascidiars."  
854 *Developmental Biology* 448 (2): 71-87.
- 855 Maric, Darko, Aleksandra Paterek, Marion Delaunay, Irene Pérez López, Miroslav Arambasic, and Dario  
856 Diviani. 2021. "A-Kinase Anchoring Protein 2 Promotes Protection against Myocardial  
857 Infarction." *Cells (Basel, Switzerland)* 10 (11): 2861.
- 858 McFadden, D. G., J. Charité, J. A. Richardson, D. Srivastava, A. B. Firulli, and E. N. Olson. 2000. "A GATA-  
859 Dependent Right Ventricular Enhancer Controls DHAND Transcription in the Developing  
860 Heart." *Development (Cambridge, England)* 127 (24): 5331-41.
- 861 McGaughey, David M., Ryan M. Vinton, Jimmy Huynh, Amr Al-Saif, Michael A. Beer, and Andrew S.  
862 McCallion. 2008. "Metrics of Sequence Constraint Overlook Regulatory Sequences in an  
863 Exhaustive Analysis at Phox2b." *Genome Research* 18 (2): 252-60.
- 864 Minnoye, Liesbeth, Ibrahim Ihsan Taskiran, David Mauduit, Maurizio Fazio, Linde Van Aerschot, Gert  
865 Hulselmans, Valerie Christiaens, et al. 2020. "Cross-Species Analysis of Enhancer Logic Using  
866 Deep Learning." *Genome Research* 30 (12): 1815-34.
- 867 Oh, Jin Woo, and Michael A. Beer. 2023. "Gapped-Kmer Sequence Modeling Robustly Identifies  
868 Regulatory Vocabularies and Distal Enhancers Conserved between Evolutionarily Distant  
869 Mammals." *BioRxiv*. <https://doi.org/10.1101/2023.10.06.561128>.
- 870 Olson, E. N. 2006. "Gene Regulatory Networks in the Evolution and Development of the Heart."  
871 Overbeek, P. A. 1997. "Right and Left Go DHAND and EHAND." *Nature Genetics*. Springer Science and  
872 Business Media LLC.
- 873 Paten, Benedict, Dent Earl, Ngan Nguyen, Mark Diekhans, Daniel Zerbino, and David Haussler. 2011.  
874 "Cactus: Algorithms for Genome Multiple Sequence Alignment." *Genome Research* 21 (9):  
875 1512-28.
- 876 Pediatric Cardiac Genomics Consortium, Bruce Gelb, Martina Brueckner, Wendy Chung, Elizabeth  
877 Goldmuntz, Jonathan Kaltman, Juan Pablo Kaski, et al. 2013. "The Congenital Heart Disease  
878 Genetic Network Study: Rationale, Design, and Early Results." *Circulation Research* 112 (4):  
879 698-706.
- 880 Prall, Owen W. J., Mary K. Menon, Mark J. Solloway, Yusuke Watanabe, Stéphane Zaffran, Fanny Bajolle,  
881 Christine Biben, et al. 2007. "An Nkx2-5/Bmp2/Smad1 Negative Feedback Loop Controls Heart  
882 Progenitor Specification and Proliferation." *Cell* 128 (5): 947-59.
- 883 Qian, Zhaozhi, and Philip L. H. Yu. 2019. "Weighted Distance-Based Models for Ranking Data Using the  
884 R Package Rankdist." *Journal of Statistical Software* 90 (5).  
885 <https://doi.org/10.18637/jss.v090.i05>.
- 886 Ramisch, Anna, Verena Heinrich, Laura V. Glaser, Alisa Fuchs, Xinyi Yang, Philipp Benner, Robert  
887 Schöpflin, et al. 2019. "CRUP: A Comprehensive Framework to Predict Condition-Specific  
888 Regulatory Units." *Genome Biology* 20 (1): 227.
- 889 Reiter, Franziska, Bernardo P. de Almeida, and Alexander Stark. 2023. "Enhancers Display Constrained  
890 Sequence Flexibility and Context-Specific Modulation of Motif Function." *Genome Research* 33  
891 (3): 346-58.

- 892 Richter, Felix, Sarah U. Morton, Seong Won Kim, Alexander Kitaygorodsky, Lauren K. Wasson, Kathleen  
893 M. Chen, Jian Zhou, et al. 2020. "Genomic Analyses Implicate Noncoding de Novo Variants in  
894 Congenital Heart Disease." *Nature Genetics* 52 (8): 769–77.
- 895 Sanges, Remo, Eva Kalmar, Pamela Claudiani, Maria D'Amato, Ferenc Muller, and Elia Stupka. 2006.  
896 "Shuffling of Cis-Regulatory Elements Is a Pervasive Feature of the Vertebrate Lineage."  
897 *Genome Biology* 7 (7): R56.
- 898 Schmidl, Christian, André F. Rendeiro, Nathan C. Sheffield, and Christoph Bock. 2015. "ChIPmentation:  
899 Fast, Robust, Low-Input ChIP-Seq for Histones and Transcription Factors." *Nature Methods* 12  
900 (10): 963–65.
- 901 Schmidt, D., M. D. Wilson, B. Ballester, P. C. Schwalie, G. D. Brown, A. Marshall, C. Kutter, et al. 2010.  
902 "Five-Vertebrate ChIP-Seq Reveals the Evolutionary Dynamics of Transcription Factor Binding."  
903 *Science* 328 (5981): 1036–40.
- 904 Schöpflin, Robert, Uirá Souto Melo, Hossein Moeinzadeh, David Heller, Verena Laupert, Jakob  
905 Hertzberg, Manuel Holtgrewe, et al. 2022. "Integration of Hi-C with Short and Long-Read  
906 Genome Sequencing Reveals the Structure of Germline Rearranged Genomes." *Nature*  
907 *Communications* 13 (1): 6470.
- 908 Shrikumar, Avanti, Eva Prakash, and Anshul Kundaje. 2019. "GkmExplain: Fast and Accurate  
909 Interpretation of Nonlinear Gapped k-Mer SVMs." *Bioinformatics (Oxford, England)* 35 (14):  
910 i173–82.
- 911 Shrikumar, Avanti, Katherine Tian, Žiga Avsec, Anna Shcherbina, Abhimanyu Banerjee, Mahfuza  
912 Sharmin, Surag Nair, and Anshul Kundaje. 2018. "Technical Note on Transcription Factor Motif  
913 Discovery from Importance Scores (TF-MoDISco) Version 0.5.6.5." *ArXiv [Cs.LG]*. arXiv.  
914 <http://arxiv.org/abs/1811.00416>.
- 915 Siepel, A., G. Bejerano, J. S. Pedersen, A. S. Hinrichs, M. Hou, K. Rosenbloom, H. Clawson, et al. 2005.  
916 "Evolutionarily Conserved Elements in Vertebrate, Insect, Worm, and Yeast Genomes."  
917 *Genome Research* 15 (8): 1034–50.
- 918 Snetkova, Valentina, Len A. Pennacchio, Axel Visel, and Diane E. Dickel. 2022. "Perfect and Imperfect  
919 Views of Ultraconserved Sequences." *Nature Reviews. Genetics* 23 (3): 182–94.
- 920 Snetkova, Valentina, Athena R. Ypsilanti, Jennifer A. Akiyama, Brandon J. Mannion, Ingrid Plajzer-Frick,  
921 Catherine S. Novak, Anne N. Harrington, et al. 2021. "Ultraconserved Enhancer Function Does  
922 Not Require Perfect Sequence Conservation." *Nature Genetics* 53 (4): 521–28.
- 923 Srivastava, D., and E. N. Olson. 1997. "Knowing in Your Heart What's Right." *Trends in Cell Biology* 7  
924 (11): 447–53.
- 925 Stennard, Fiona A., Mauro W. Costa, David A. Elliott, Scott Rankin, Saskia J. P. Haast, Donna Lai, Lachlan  
926 P. A. McDonald, et al. 2003. "Cardiac T-Box Factor Tbx20 Directly Interacts with Nkx2-5, GATA4,  
927 and GATA5 in Regulation of Gene Expression in the Developing Heart." *Developmental Biology*  
928 262 (2): 206–24.
- 929 Taher, Leila, David M. McGaughey, Samantha Maragh, Ivy Aneas, Seneca L. Bessling, Webb Miller,  
930 Marcelo A. Nobrega, Andrew S. McCallion, and Ivan Ovcharenko. 2011. "Genome-Wide  
931 Identification of Conserved Regulatory Function in Diverged Sequences." *Genome Research* 21  
932 (7): 1139–49.
- 933 Taskiran, Ibrahim I., Katina I. Spanier, Hannah Dickmanken, Niklas Kempynck, Alexandra Pančíková,  
934 Eren Can Ekşi, Gert Hulselmans, et al. 2023. "Cell-Type-Directed Design of Synthetic  
935 Enhancers." *Nature*, December, 1–9.
- 936 Torchiano, Marco. 2016. *Effsize - a Package for Efficient Effect Size Computation*. Zenodo.  
937 <https://doi.org/10.5281/ZENODO.1480624>.
- 938 Villar, Diego, Camille Berthelot, Sarah Aldridge, Tim F. Rayner, Margus Lukk, Miguel Pignatelli, Thomas  
939 J. Park, et al. 2015. "Enhancer Evolution across 20 Mammalian Species." *Cell* 160 (3): 554–66.
- 940 Vinga, Susana. 2014. "Information Theory Applications for Biological Sequence Analysis." *Briefings in*  
941 *Bioinformatics* 15 (3): 376–89.

- 942 Visel, Axel, Matthew J. Blow, Zirong Li, Tao Zhang, Jennifer A. Akiyama, Amy Holt, Ingrid Plajzer-Frick,  
943 et al. 2009. "ChIP-Seq Accurately Predicts Tissue-Specific Activity of Enhancers." *Nature* 457  
944 (7231): 854–58.
- 945 Visel, Axel, James Bristow, and Len A. Pennacchio. 2007. "Enhancer Identification through Comparative  
946 Genomics." *Seminars in Cell & Developmental Biology* 18 (1): 140–52.
- 947 Wu, Tianzhi, Erqiang Hu, Shuangbin Xu, Meijun Chen, Pingfan Guo, Zehan Dai, Tingze Feng, et al. 2021.  
948 "ClusterProfiler 4.0: A Universal Enrichment Tool for Interpreting Omics Data." *Innovation*  
949 (*Cambridge (Mass.)*) 2 (3): 100141.
- 950 Xiao, Feng, Xiaoran Zhang, Sarah U. Morton, Seong Won Kim, Youfei Fan, Joshua M. Gorham, Huan  
951 Zhang, et al. 2024. "Functional Dissection of Human Cardiac Enhancers and Noncoding de  
952 Novo Variants in Congenital Heart Disease." *Nature Genetics* 56 (3): 420–30.
- 953 Zaidi, Samir, Murim Choi, Hiroko Wakimoto, Lijiang Ma, Jianming Jiang, John D. Overton, Angela  
954 Romano-Adesman, et al. 2013. "De Novo Mutations in Histone-Modifying Genes in Congenital  
955 Heart Disease." *Nature* 498 (7453): 220–23.
- 956 Zhang, Xiaoyu, Irene M. Kaplow, Morgan Wirthlin, Tae Yoon Park, and Andreas R. Pfenning. 2020.  
957 "HALPER Facilitates the Identification of Regulatory Element Orthologs across Species."  
958 *Bioinformatics (Oxford, England)* 36 (15): 4339–40.
- 959 Zielezinski, Andrzej, Hani Z. Girgis, Guillaume Bernard, Chris-Andre Leimeister, Kujin Tang, Thomas  
960 Dencker, Anna Katharina Lau, et al. 2019. "Benchmarking of Alignment-Free Sequence  
961 Comparison Methods." *Genome Biology* 20 (1): 144.
- 962 Zielezinski, Andrzej, Susana Vinga, Jonas Almeida, and Wojciech M. Karlowski. 2017. "Alignment-Free  
963 Sequence Comparison: Benefits, Applications, and Tools." *Genome Biology* 18 (1): 186.