1  **A tripartite structure, the complex nuclear receptor element (cNRE), is a *cis-***

2  **regulatory module of viral origin required for atrial chamber preferential**

3  **gene expression**

4

5  Luana Nunes Santos[1,2,3], Ângela Maria da Souza Costa[1], Martin Nikolov[2],

6  Allysson Coelho Sampaio[3, 4], Frank E. Stockdale[5], Gang F Wang[ø], Hozana

7  Andrade Castillo[1,2], Mariana Bortoletto Grizante[1], Stefanie Dudczig[6], Michelle

8  Vasconcelos[3], Nadia Rosenthal[7,8], Patricia Regina Jusuf[6], Paulo de Oliveira[1],

9  Tatiana Guimarães de Freitas Matos[3], William Nikovits, JR.[5], Michael Schubert[9],

10  Mirana Ramialison[2,10, *] and José Xavier-Neto[3,11]

11

12  **Affiliations:**

13  [1] Brazilian Biosciences National Laboratory (LNBio), Brazilian Center of Research

14  in Energy and Materials (CNPEM), Campinas, SP, Brazil.

15  [2] Australian Regenerative Medicine Institute, Monash University, VIC Australia -

16  Systems Biology Institute Australia.

17  [3] Department of Cell and Developmental Biology, Institute of Biomedical

18  Sciences, University of São Paulo (USP), São Paulo, SP, Brazil.

19  [4] Faculdade Santa Marcelina - São Paulo, Brazil.

20  [5] Department of Medicine, Stanford University, Stanford, California, USA.

21  [6] School of BioSciences, University of Melbourne, Parkville, VIC, Australia.

22   [7] The Jackson Laboratory, Bar Harbor, 04609 Maine, USA.

23   [8] National Heart and Lung Institute, Imperial College London, London, United

24   Kingdom.

25   [9] Laboratoire de Biologie du Développement de Villefranche-sur-Mer, Institut de

26   la Mer de Villefranche, Sorbonne Université, CNRS, Villefranche-sur-Mer,

27   France.

28   [11] Murdoch Children's Research Institute, Parkville, VIC, Australia.

29   [10] Department of Morphology, Federal University of Ceará (UFC), Ceará, CE,

30   Brazil.

31

32   [*] Corresponding author.

33   [ø] in memory.

34

35   **E-mail addresses:**

36   Luana Nunes Santos: santosln89@gmail.com

37   Allysson Coelho Sampaio: allysson.sampaio@santamarcelina.edu.br

38   Ângela Maria da Souza Costa: angelamarsou@gmail.com

39   Frank E Stockdale: stockdale@stanford.edu

40   Hozana Andrade Castillo: hozana.castillo@monash.edu

41   Mariana Bortoletto Grizante: mari.grizante@gmail.com

42    Martin Nikolov: mnik0002@student.monash.edu

43    Michelle Vasconcelos: michellevasconcelos@gmail.com

44    Nadia Rosenthal: nadia.rosenthal@jax.org

45    Patricia Regina Jusuf: patricia.jusuf@unimelb.edu.au

46    Paulo de Oliveira: paulo.oliveira@lnbio.cnpem.br

47    Stefanie Dudczig: stefanie.dudczig@unimelb.edu.au

48    Tatiana Guimarães de Freitas Matos: tgfmatos@yahoo.com.br

49    William Nikovits JR: nikovits@comcast.net

50    Michael Schubert: michael.schubert@imev-mer.fr

51    Mirana Ramialison: mirana.ramialison@monash.edu

52    José Xavier-Neto: josexavierneto@gmail.com

53

## Abstract

55    Optimal cardiac function requires appropriate contractile proteins in each heart

56    chamber. Atria require slow myosins to act as variable reservoirs, while ventricles

57    demand fast myosin for swift pumping functions. Hence, myosin is under

58    chamber-biased *cis*-regulatory control to achieve this functional distribution.

59    Failure in proper regulation of myosin genes can lead to severe congenital heart

60    dysfunction. The precise regulatory input leading to cardiac chamber-biased

61    expression remains uncharted. To address this, we computationally and

molecularly dissected the quail Slow Myosin Heavy Chain III (SMyHC III) promoter that drives specific gene expression to the atria to uncover the regulatory information leading to chamber expression and understand their evolutionary origins. We show that SMyHC III gene states are autonomously orchestrated by a complex nuclear receptor *cis*-regulatory element (cNRE), a 32-bp sequence with hexanucleotide binding repeats. Using *in vivo* transgenic assays in zebrafish and mouse models, we demonstrate that preferential atrial expression is achieved by the combinatorial regulatory input composed of atrial activation motifs and ventricular repression motifs. Through comparative genomics, we provide evidence that the cNRE emerged from an endogenous viral element, most likely through infection of an ancestral host germline. Our study reveals an evolutionary pathway to cardiac chamber-specific expression.

**Keywords:** atrial expression, evolution, enhancer, viral insertion, nuclear receptors

## Introduction

Vertebrate chambered hearts are efficient pumps organized according to an ancient evolutionary paradigm that divides their circulatory functions into inflow and outflow working modules, the atria, and ventricles, respectively (Simões-Costa et al., 2005). Most current views on cardiac chamber development agree with the principles of epigenesis, in which higher-level structures arise through sequential morphogenetic steps. The gradual three-dimensional organization is accompanied by progressive restriction of cellular fates, from large, early

4

85    embryonic fields that initially display broad tissue potencies (e.g., epiblast or

86    mesoderm) to terminally differentiated cells, such as most heart cell types. The

87    highly varied cardiac cell fates originate from a combination of mosaic and

88    regulative development processes (José Xavier-Neto et al., 2012). Ultimately,

89    these developmental processes combine clues from the relative position of each

90    group of cells inside the embryo with the information imparted by patterning and

91    migration, which modify the initial relationships between cell progenitors and

92    further restrict fates.

93    Most efforts in the last twenty years have been dedicated to transpose the

94    blueprints of the above-mentioned cardiac ontogenetic events from the four-

95    dimension arena into the overlapping chemical (i.e., signaling) and genomic

96    spaces, with a strong focus on gene regulatory pathways ((e.g. (Erwin &

97    Davidson, 2002)). A clear description of ontological and epigenetic events

98    analogous to an engineering blueprint in all these axes of information is years

99    ahead. Moreover, it has become clear that, yet another dimension will have to be

100   accounted for. This additional dimension is evolutionary time, and its major arena

101   is, again, the genome. Rather than being a monotonous landscape with an

102   occasional mutation occurring here and there, the genome is a dynamic space,

103   frequently challenged by the insertion and/or proliferation of mobile genetic

104   elements such as viruses and transposons (Kazazian, 2004). All these wandering

105   elements have the potential to disrupt the genome, leading to insertion mutations,

106   inversions, and translocations, but, more importantly, from an evolutionary

107   standpoint, they also have the capacity to foster genome innovation, including the

108   assembly of entirely new gene regulatory networks.

109    To understand the establishment and the evolution of gene regulatory networks,

110    we have previously examined the Slow Myosin Heavy Chain III gene promoter

111    (SMyHC III) and determined that it is able to drive preferential atrial gene

112    expression (Bruneau et al., 2000; Nikovits et al., 1996; G F Wang et al., 1998;

113    Gang Feng Wang et al., 1996, 2001; J Xavier-Neto et al., 1999). The sequence

114    AGGACAAAGAGGGGA located from −801 to −787 bp upstream from the

115    transcription start site of SMyHC III contains two Hexad sequences (Hexads A

116    and B). Hexads A and B were previously identified as a dual putative Vitamin D

117    Receptor Element (VDRE) and a Retinoic Acid Receptor Responsive Element

118    (RARE) (G F Wang et al., 1998; Gang Feng Wang et al., 1996, 2001),

119    respectively, with a well-established ventricular inhibitory function associated with

120    the VDRE (Gang Feng Wang et al., 1996). Although correct, the current model

121    for selective atrial expression via ventricular repression requires an extension to

122    include atrium-specific activating elements present in the SMyHC III promoter (G

123    F Wang et al., 1998; Gang Feng Wang et al., 1996).  We sought to investigate

124    additional mechanisms for atrial specificity in the SMyHC III promoter using *in*

125    *silico* and in vivo approaches. We show that the atrial preference exhibited by the

126    SMyHC III promoter is manifested in avian, mammals, and teleost fishes, chiefly

127    on account of a low frequency repetitive 32 base pair genome element formed by

128    tandem reiterations of three purine-rich hexanucleotide repeats, here designated

129    as the complex Nuclear Receptor Element (cNRE). The cNRE is a versatile

130    regulator of selective cardiac chamber expression, switching from SMyHC III

131    activator to repressor functions according to atrium, or ventricular contexts,

132    respectively. We demonstrated that the combination of three Hexads A, B, and C

133    within the cNRE, provides an information processing platform that integrates

6

134    signals from different elements and that the cNRE is necessary and sufficient to

135    confer preferential atrial expression. Finally, using comparative genomics, we

136    provide evidence that the cNRE was associated with the SMyHC III gene through

137    infection of an ancestral host germline by an unknown virus resulting in

138    recombination into the genome of a Galliform bird ancestor at the root of the

139    Galliformes radiation in the Cretaceous, about 63 million years ago.

140    ## Results

141    **Identification of the cNRE as a tripartite structure with Hexads A, B, and C**

142    To investigate additional mechanisms of atrial specificity, we performed

143    computerized profiling of nuclear receptor binding sites in the SMyHC III

144    promoter. This study predicted a novel nuclear receptor binding Hexad (Hexad

145    C), adjacent to Hexads A and B known to act as ventricular repressor sequences

146    (aaggacaaagaggggacaaagAGGCGGaggt at -786 to -778 bp) (G F Wang et al.,

147    1998; Gang Feng Wang et al., 1996) (Figure 1A). The combination of these three

148    Hexads sequences (A +B +C) was designated as the complex Nuclear Receptor

149    Element (cNRE). We postulated that this novel tripartite nuclear receptor binding

150    contains the minimal information necessary and sufficient for specific atrial

151    expression in vertebrate embryos.

152    **The cNRE is a transferable *cis*-regulatory agent necessary for driving atrial**

153    **expression**

154    To test the necessity of the cNRE for atrial specificity, we performed transient

155    expression assays in zebrafish embryos (**Figure 1B**). Two reporter constructs,

156    the quail *SMyHC III* promoter driving GFP (*SMyHC III*::GFP) and the quail

157   *SMyHC III* promoter in which the cNRE was deleted (*SMyHC IIIΔcNRE*::GFP),

158   were injected into the Tg(*vhmc*::mCherry) embryos (Jin et al., 2009). This line

159   exclusively expresses mCherry in the ventricle. We assessed the proportion of

160   embryos displaying GFP expression in only the atrium, ventricle, or both

161   chambers (**Figure 2A, B**). With the wild-type (WT) SMyHC III promoter construct,

162   37% of embryos (n= 33) displayed GFP expression in the atrium, and 53% (n=

163   45) expressed the reporter in both ventricular and atrial chambers (**Figure 2B**).

164   However, there was greater atrial-specific expression (37%) than ventricular-

165   specific expression (10%, n= 6) (**Figures 1B, 2B**). In contrast, atrial-specific

166   expression was statistically significantly reduced in the mutated quail promoter

167   *SMyHC IIIΔcNRE*::GFP reporter assays (21%, n= 9, p=0.0026), and

168   concomitantly, ventricular-specific expression was statistically significantly

169   increased (38%, n= 26, p=0.0038) (**Figures 1B, 2B**). There was no statistically

170   significant change in the proportion of embryos displaying non-chamber-specific

171   expression for both SMyHC *III*::GFP and the *SMyHC IIIΔcNRE*::GFP constructs

172   (53%, n= 82 and 41%, n= 55 respectively) (**Figure 2B, Figure S1**). Taking

173   together, these results suggest that the deletion of the cNRE from the *SMyHC III*

174   promoter reduces atrial-specificity and support the idea that the cNRE is

175   necessary for driving atrium-specific gene expression, through the release of

176   ventricular repression.

177   To support this interpretation, we performed transgenic assays in the *SMyHC*

178   *III*::HAP transgenic mouse line driving atrial-specific expression of the human

179   alkaline phosphatase reporter gene (HAP) (J Xavier-Neto et al., 1999) under the

180   control of the quail *SMyHC III* promoter (**Figure 1C, Figure S2**). We observed

181   that the deletion of a 72bp region encompassing the cNRE (*SMyHC*

182    *III*ΔcNRE::HAP) abrogated atrial expression (**Figure 1C, Figure S3**), further

183    supporting the notion that the cNRE is necessary for driving atrium-specific

184    expression. A statistically insignificant level of reporter expression was observed

185    in non-atrial regions in both the *SMyHC III*::HAP and the *SMyHC III*ΔcNRE::HAP

186    transgenic mouse lines (**Figure 1D-E**).

**The cNRE is a transferable *cis*-regulatory agent sufficient to drive atrial**

**expression**

189    To test the sufficiency of the cNRE for driving atrial specific expression, we

190    devised a conversion assay, which aimed at testing whether the cNRE is

191    sufficient to revert a pattern from ventricular to atrial activation. To do so, we

192    performed transient expression assays with the *vmhc* promoter driving GFP

193    expression (*vmhc*::GFP) (**Figure 1F**) and with a 5' fusion of five tandem repeats

194    of the cNRE to this *vmhc* reporter construct (*5xcNRE-vmhc*::GFP) (**Figure 1G**).

195    In 48 hpf zebrafish, we observed ventricule-specific GFP expression in most

196    transients injected with *vmhc*::GFP (73%, n= 39) (**Figure 1H**). There were no

197    embryos expressing GFP exclusively in the atrium. In contrast to the WT *vmhc*

198    promoter, with the *5xcNRE-vmhc* fusion construct (**Figure 1G**), we observed a

199    significant increase in the proportion of embryos showing both atrial and

200    ventricular GFP expression (n= 55), a decrease in the number of embryos

201    expressing GFP exclusively in the ventricle (46% of the embryos, n= 47), and

202    strikingly, a single embryo with reporter expression exclusively in the atrium

203    (0,97%, n= 1) (**Figure 1H**). These experiments demonstrated that, outside of its

204    native context in the SMyHC III promoter, the cNRE is sufficient to convert

205    ventricular to atrial expression.

206    Taking together, our results suggest that the cNRE is necessary and sufficient for

207    directing preferential atrial gene expression. Experiments in zebrafish suggest

208    that atrial specificity could be achieved through ventricular gene repression.

209    Hence, we sought to investigate the precise contribution of the binding elements

210    within the cNRE to the activation of atrial-specific gene expression or repression

211    of ventricular-specific gene expression.

212    **Combinatorial recruitment of Hexads A, B and C is essential for cNRE**

213    **activity *in vivo***

214    To understand the *cis*-regulatory composition of the cNRE we assessed the

215    contribution of Hexads A, B and C in driving atrial specific expression *in vivo*. We

216    used site-directed mutagenesis of individual Hexads in the *SMyHC* III promoter

217    used in zebrafish transient assays (*SMyHC* III*::GFP*) (**Figure 2**) and in mouse

218    mutant lines (*SMyHC* III*::HAP*) (**Figure 3, Figures S2, S3, S4, S5**). Reporter

219    expression driven by the mutated cNRE was compared to WT cNRE (**Figure 4A**),

220    or its complete deletion (**Figure 4B**). Mutation of Hexad A (Mut A) was obtained

221    by substituting the Hexad A sequence 5' AGGACA 3' for 5' GTCGAC 3' (**Figure**

222    **4C**). Dinucleotide substitutions were performed on Hexad B and Hexad C to

223    obtain Mut B and Mut C, respectively (**Figure 4D-E**). One dinucleotide

224    substitution in the spacer region between Hexads B and C was designed as a

225    non-Hexad control mutation (Mut S) (**Figure. 4F**). In zebrafish embryos, when

226    comparing Mut A to the WT *SMyHC* III*::GFP* promoter (**Figures 2, 4A, A'**), we

227    observed a statistically significant increase in the proportion of embryos with GFP

228    positive cells in both atrium and ventricle (from 53%, n=45, in WT to 81%, n=70,

229    in Mut A, p<0.001) and a statistically significant decrease in the proportion of

10

230  embryos with GFP-positive cells exclusively in the atrium (from 38%, n=33, in WT

231  to 8%, n=7, in Mut A, p<0.0001) (**Figure 2B**). Mut A thus leads to a significant

232  increase of GFP-positive cells in both atrium and ventricle, mainly by stimulating

233  GFP expression in the ventricle (**Figures 2B, 4C, C'**). We hence conclude that,

234  in zebrafish, Hexad A is required for ventricular repression. In contrast, in mouse

235  embryos, Mut A (**Figures 3A, 4C'', Figure S4**) had no effect on HAP expression

236  when compared to WT *SMyHC* III::HAP (**Figures 1C, 4A'', Figure S2**). Similar

237  to the effect of Mut A, for Mut B in zebrafish, we observed a statistically significant

238  increase in the proportion of GFP-positive cells in both chambers (from 53%,

239  n=45, in WT to 78%, n=62, in Mut B, p<0.0001) and a statistically significant

240  decrease in the proportion of embryos with GFP-positive cells exclusively in the

241  atrium (from 37%, n=33, in WT to 17%, n=9, in Mut B, p=0.0002) (**Figure 2B**).

242  For Mut B, we found a similar effect in mouse embryos, with conspicuous HAP

243  expression in the left ventricle and proximal outflow tract (**Figure 3C-H**).

244  Altogether, the experiments in both the zebrafish and the mouse support the

245  notion that Mut B released ventricular repression (**Figure 4D, D', D''**), thus

246  suggesting that Hexad B is required for repression of gene expression in the

247  ventricle. Mutation of Hexad C, Mut C, resulted in a statistically significant

248  increase of ventricular-specific expression (31%, n=16, p=0.0005) when

249  compared to the *SMyHC III* control (10%, n=6) in zebrafish embryos (**Figure 2B**).

250  This increase is accompanied by a decrease of the proportion of embryos with

251  GFP expression only in the atrium or in both the atrial and ventricular chambers.

252  Consistent with this result in zebrafish, we found a marked decrease of HAP

253  staining in mouse Mut C mutants (**Figure 3I, J, M, Figure S5B-D**), relative to WT

254  mice (**Figure 3K, Figure S5A**). Of note, compared to the WT, Mut C mutant mice

11

255    were characterized by a general reduction of HAP expression in 10.5 dpc hearts

256    (**Figure 3L, N, Figure S5A**). Taken together, the reduction of atrial reporter

257    expression of Mut C in both zebrafish and mouse (**Figures 4E, E', E''**) suggests

258    that Hexad C acts as an atrial activator. As control of the Hexad deletions, we

259    assessed the effect of a mutation in a spacer region outside the Hexads (Mut S).

260    This mutation is located between Hexads B and C, adjacent to Hexad B (**Figure**

261    **4F**). In zebrafish, Mut S did not trigger a significant change in atrial-specific

262    expression. However, we observed an increase of GFP-positive cells in the

263    ventricle (29%, n=8) compared to the *SMyHC III* control (10%, n=6), at the

264    expense of GFP expression in both chambers (29%, n=10, p=0.0388) compared

265    to the *SMyHC III* control (53%, n=45) (**Figure 2B**). This suggests that Mut S might

266    contribute to the repression of atrial expression (**Figures 4F, F'**). However, these

267    results could not be confirmed in mice, where Mut S had no effect on HAP

268    expression compared to the WT *SMyHC III*::HAP (**Figures 3O-P, 4F''**). In

269    conclusion, the phenotype displayed by Mut A and Mut B transgenic embryos is

270    consistent with the release of a highly stereotyped ventricular expression pattern.

271    These observations are consistent with those previously described after the

272    deletion of the whole VDRE/RARE motif (Hexad A and B) (G F Wang et al., 1998;

273    Gang Feng Wang et al., 1996). We further reveal that the Hexad C element in

274    the *SMyHC III* promoter is responsible for atrial-specific activation **(Figure 4)**.

275    **Evolutionary origin of the cNRE**

276    Given that the quail cNRE sequence drives preferential atrial expression in

277    different vertebrate taxa, including zebrafish and mice, we sought to define the

278    evolutionary origin of the cNRE. Based on the observation that the *SMyHC III*

279    gene is a representative of a clade of slow myosin's specific to birds (Chen et al.,

12

280    1997; Nikovits et al., 1996) (**Figure 5A**). In all galliform birds, at least one hit of

281    the 30-32 bp-long cNRE was found with a significant E-value (**Figure 5A**). The

282    cNRE described here, composed of the three Hexads A, B and C, is thus unique

283    to galliform birds and has likely originated in their last common ancestor,

284    approximately 63 million years ago (Kuhl et al., 2021). Interestingly, the cNRE

285    sequence strongly matches more than one locus in the genomes of various

286    galliform birds, including Japanese quail (*Coturnix japonica*), chicken (*Gallus*

287    *gallus*), pinnated grouse (*Tympanuchus cupido*), wild turkey (*Meleagris*

288    *gallopavo*), helmeted guineafowl (*Numida meleagris*), Indian peafowl (*Pavo*

289    *cristatus*), common pheasant (*Phasianus colchicus*) and Mikado pheasant

290    (*Syrmaticus mikado*) (**Table 1**). Since galliform birds are known to be highly

291    susceptible to viral integration (Holmes, 2011; Kaleta, 1990; McGeoch et al.,

292    2000; Morissette & Flamand, 2010; Nair, 2005), we hypothesized that the cNRE

293    might have a viral origin. To test this hypothesis, we systematically screened for

294    cNRE sequences in the avian viral database (see materials and methods). We

295    found statistically significant cNRE matches in the genomes of papillomaviruses,

296    paramyxoviruses, retroviruses and herpesviruses (**Figure 5B**). This suggests

297    that the cNRE might have originally been associated with the *SMyHC III* gene

298    through viral infection of an ancestral galliform host and that integration into the

299    genome occurred early during galliform radiation in the Cretaceous period, about

300    63 million years ago. Taken together, our study demonstrates that the cNRE is a

301    tripartite *cis*-regulatory element that confers cardiac chamber specificity and

302    arose during early diversification of galliform birds by genomic integration of a

303    virus-derived sequence.

304

13

305

306 **Table 1.** Species with multiple cNRE hits.

| Species | Number of hits | Sequence length | E-value | Genomic location | BLAST+ parameters used |
|---|---|---|---|---|---|
| Japanese quail (*Coturnix japonica*) | 2 | 31/31 | 0.004/0.049 | Unplaced genomic scaffold/Ch1 | 1/-3; 1/1 |
| Chicken (*Gallus gallus*) | 2 | 31/31 | 5e-05/0.019 | Ch19/Ch19 | 1/-1; 2/1 |
| Wild turkey (*Meleagris gallopavo*) | 2 | 31/30 | 0.002/0.005 | Whole genome shotgun sequence | 1/-3; 1/1 |
| Helmeted guineafowl (*Numida meleagris*) | 2 | 30/30 | 1e-08/3e-06 | Ch18/Ch18 | 1/-3; 5/2 |
| Indian blue peafowl (*Pavo cristatus*) | 2 | 31/32 | 9e-08/0.005 | Whole genome shotgun sequence | 1/-3; 1/1 |
| Common pheasant (*Phasianus colchicus*) | 2 | 30/31 | 3e-06/2e-04 | Whole genome shotgun sequence | 1/-3; 5/2 |
| Mikado pheasant (*Syrmaticus mikado*) | 2 | 31/30 | 1e-05/4e-05 | Whole genome shotgun sequence | 1/-3; 1/1 |
| Pinnated grouse (*Tympanuchus cupido*) | 2 | 31/31 | 8e-07/0.048 | Whole genome shotgun sequence | 1/-3; 5/2 |

307 (Ch) chromosome

308

309

14

## Discussion

310

311  Cardiac development in mammals, *i.e.*, the formation of a four-chambered heart,

312  is characterized by a series of complex morphogenetic movements and is

313  controlled by an intricate gene regulatory network (Waardenberg et al., 2014). In

314  this work, we have identified the complex Nuclear Receptors Element (cNRE), a

315  new, 32bp-long regulatory element located within the quail Slow Myosin Heavy

316  Chain III (*SMyHC* III) promoter. We further showed that the cNRE is necessary

317  and sufficient for driving reporter gene expression specifically in atrial cells of

318  mice and preferentially in atrial cells of zebrafish. We also demonstrated that this

319  specific atrial regulation mediated by the cNRE crosses species barriers (from

320  quail to mice and zebrafish), highlighting the importance of this new element for

321  understanding the evolution of gene regulation underlying the specification of

322  cardiac chambers during early embryonic development. Wang *et al.* (G F Wang

323  et al., 1998) previously demonstrated that the VDRE/RARE element, which

324  includes Hexads A and B of the cNRE, is responsible for ventricular repression

325  of *SMyHC* III promoter activity in avian and murine hearts. They also identified a

326  GATA-binding element in the *SMyHC III* promoter involved in activating

327  expression in both the atrium and the ventricle, but they failed to identify the DNA

328  sequence driving expression of the promoter in atrial cells. By searching for

329  potential novel nuclear receptor binding sites, we defined the cNRE as a 3'

330  expansion of the initial 17-bp long, adding a third Hexad. The 32-bp cNRE thus

331  contains three Hexads, A, B and C, and is responsible for the preferential activity

332  of the promoter in the atrium (Nikovits et al., 1996; G F Wang et al., 1998; Gang

333  Feng Wang et al., 1996, 2001). To functionally characterize the activity of the

334  cNRE in the heart, we took advantage of the power of transient expression

15

335    assays in zebrafish embryos as a complementary strategy to the generation of

336    stable transgenic lines in mice (Meyers, 2018). The evaluation of the effects of

337    point mutations in the cNRE in both zebrafish and mice ultimately allowed us to

338    assess the role of each of the three Hexads constituting the cNRE. We found that

339    cNRE can drive specific reporter gene expression in the atrium of both zebrafish

340    and mice, as deletion of the cNRE within the *SMyHC III* promoter abrogated the

341    atrial-specific expression driven by this promoter.

342    Of note, while the activity of the *SMyHC III* promoter is limited to the atrium in

343    mice, the promoter drives expression in both chambers of the zebrafish heart. It

344    might be that the teleost fish-specific whole genome duplication (REF) has

345    secondarily altered the regulatory landscape controlling heart development (Jin

346    et al., 2009). However, preferential atrial expression was nonetheless

347    significantly decreased in the zebrafish ΔcNRE mutant, and the multimerized

348    cNRE construct still clearly shifted the activity of the zebrafish *vmhc* promoter in

349    an atrial direction. We sought to identify which regions in the cNRE sequence are

350    critical to its activity by creating constructs containing mutations in the cNRE

351    sequence for transient and stable transgenic analyses in, respectively, zebrafish

352    and mice. When mutating the sequences of Hexad A and B in zebrafish, we

353    observed a decrease in atrial expression concomitant with an increase in the

354    number of individuals with expression in both chambers. In mice, in contrast,

355    mutation of Hexad A did not affect expression, while the mutation of Hexad B

356    increased expression in ventricular chambers. These data suggest that Hexad B,

357    and potentially Hexad A, act as ventricular repressors. This is consistent with

358    previous *in vitro* studies performed in quail atrial cells, where removal of Hexad

359    A and B led to increased reporter gene expression in ventricular cells (Gang Feng

360 Wang et al., 2001). Mutation of Hexad C resulted in a statistically significant

361 increase of ventricle-specific expression in zebrafish. We hypothesize that this

362 effect is due to the specific loss of atrial expression in embryos that would

363 normally be characterized by transgenic activity of the construct in both

364 chambers. This notion is consistent with our finding that mutation of Hexad C in

365 mice led to an abrogation of atrial expression. We thus conclude that Hexad C

366 plays an important role in atrial activation. In summary, we demonstrate that the

367 cNRE contains information needed for both atrial activation and ventricular

368 repression of the *SMyHC III* promoter. Further work will be required to define the

369 transcription factors binding to the cNRE sequence to regulate its activity.

370 Although the *SMyHC III* promoter is not conserved between species, it is capable

371 of driving atrium-specific expression in different animal models, including chicken,

372 mouse and zebrafish (Nikovits et al., 1996; G F Wang et al., 1998; Gang Feng

373 Wang et al., 1996, 2001; J Xavier-Neto et al., 1999). We postulate that the cNRE

374 acts as a dual *cis*-regulatory module integrating both activating and repressing

375 signals, some of which likely mediated by members of the nuclear receptor

376 superfamily via the VDRE/RARE binding sites. We propose that the cNRE *cis*-

377 regulatory module emerged from an endogenous viral element (Holmes, 2011),

378 a genomic remnant of a rare recombination event caused by a viral infection of

379 an ancestral host. Using comparative genomics, we provide evidence that this

380 virus-derived cNRE was associated with the *SMyHC III* gene in the last common

381 ancestor germline of Galliform birds in the Cretaceous period, about 63 million

382 years ago.

383

## Conclusion

Supported by *in vivo* experiments in zebrafish and mice, we demonstrate that the cNRE, a sequence of merely 32 bp, carries information to control both atrial activation and ventricular repression. We further provide evidence this *cis*-regulatory element is of viral origin. Our work thus highlights the evolution of specific regulatory motifs in a sequence that is present exclusively in avian genomes but can drive preferential atrial expression in different vertebrate taxa. Taken together, this study sheds light on the origin of enhancers and defines the minimum amount of information required for regulating gene expression in an atrial-specific fashion.

## Materials and Methods

### Bioinformatics profiling of nuclear receptor binding sites at the SMyHC III promoter

We devised a simulation approach to identify nuclear receptor binding sites (Hexads) in the SMyHC III promoter. The principle of the approach is based on the Poisson-Boltzmann theory and aims at calculating interaction energies between nuclear receptors and DNA as an approximation of their respective binding affinities. Protein/DNA complexes were assembled for molecular dynamics profiling by positioning three-dimensional (3D) structures of nuclear receptor DNA-binding domains on the 3D structure of the cNRE DNA. RXR, RAR, and VDR crystal structures are available at the Protein Data Bank (http://www.rcsb.org) with the codes 1DSZ, 1KB4 and 1BY4, respectively (Rastinejad et al., 2000; Shaffer & Gewirth, 2002; Zhao et al., 2000). The free

407 binding energy was calculated for protein/DNA complexes using the trajectories

408 obtained from the molecular dynamics profiling with the software MM-PBSA in

409 the AMBER package (Sekijima et al., 2003). Each complex (AB) was split into

410 two parts: the nuclear receptor structure (A) and the cNRE structure (B), and the

411 energy was calculated for the whole complex AB as well as for each part, A and

412 B. Binding energy differences were obtained according to: $\Delta\Delta G = \Delta GAB - (\Delta GA$

413 $+ \Delta GB)$. Binding free energies for all cNRE Hexads were plotted with reference

414 to its first nucleotide. Data were pooled and analyzed as a box plot to identify

415 values below the 10th percentile. These values were used to identify all potential

416 Hexads within the cNRE. To quantify the potential for nuclear receptor binding

417 within the cNRE we scored the number of times a given cNRE nucleotide was

418 part of a Hexad.

419 **Generation of mouse reporter lines**

420 We generated mice containing mutations in critical nucleotides of Hexads A, B

421 and C (Mut A, B and C, respectively) as well as a non-Hexad control mutation

422 (Mut S). The constructs were synthesized with Agilent QuikChange II XL Site-

423 Directed Mutagenesis Kit® (Cat #200521) following the manufacturer's

424 instructions and using the primers described in Table 2. Constructs were

425 subsequently sequenced to confirm the substitutions. The Hexad A mutation (Mut

426 A) was generated by mutagenesis of this Hexad's nucleotides to a Sal I restriction

427 enzyme site (5'-AGGACA-3' to 5'-GTCGAC-3'). The mutation in Hexad B (Mut B)

428 were point mutations of the two first nucleotides of this Hexad (5'-AGGGGA-3' to

429 5'-TTGGGA-3'). Similarly, the mutation of Hexad C (Mut C) consisted of point

430 mutations of the two first nucleotides of this Hexad (5'-GGCGGA-3' to 5'-

431 TTCGGA-3'). The non-Hexad control mutation (Mut S) was obtained by mutating

432 two nucleotides in the putative nucleotide spacer between Hexads B and C (5'-

433 AGGGGAcaaagaGGCGGA-3' to 5'-AGGGGAttaagaGGCGGA-3'). *SMyHC*

434 *III*::HAP transgenic lines 1 and 5 have previously been described (J Xavier-Neto

435 et al., 1999), and four new *SMyHC III*::HAP transgenic lines were generated (2,

436 6, 27 and 29). Three mutB::HAP (5, 17 and 19) and seven mutC::HAP (3, 5, 7,

437 14, 15, 23 and 24) transgenic lines were established.


438 **HAP staining and histology**


439 For HAP staining and paraffin sections, embryos and hearts were handled as

440 described in (J Xavier-Neto et al., 1999). For HAP assays tissues were

441 homogenized in a lysis buffer containing 0.2% Triton X-100 with a mini bouncer.

442 HAP activity in cardiac tissues was measured with Phospha-Light™, a

443 chemiluminescent assay from PerkinElmer. HAP activity was normalized relative

444 to protein concentration.


445 **Zebrafish wild-type and stable transgenic Tg(*vmhc*::mCherry) lines used in**

446 **the transient experiments.**


447 To generate the Tg(*vmhc*::mCherry) line, we used a promoter fragment of 1.9 kb

448 upstream of the *vmhc* gene, as previously described (Jin et al., 2009; Zhang &

449 Xu, 2009). The PCR-amplified fragment was cloned into pT2AL200R150G

450 (courtesy of Dr. Koichi Kawakami) using the Xho I and Hind III restriction sites.

451 The eGFP, between the Cla I and BamH I restriction sites, was substituted by

452 mCherry. Injected embryos were raised to adulthood and an F2 generation was

453 established. To generate the *Tol2-SMyHC III*::GFP construct, the 840-bp

454 upstream regulatory sequence of quail *SMyHC III* was excised at the Sma I and

455 Hind III restriction sites from the *SMyHC III* pGl3 plasmid and cloned into the

456 pT2AL200R150G using the Xho I and Hind III restriction sites. Deletion of the 32-

457 bp cNRE from the *SMyHC III* promoter (*SMyHC III*::ΔcNRE) was obtained by

458 digesting the pGl3 plasmid with Xho I and Hind III followed by subsequent cloning

459 into pT2AL200R150G at the Sma I and Hind III restriction sites. For the mutational

460 analyses of the *SMyHC III* promoter, specific mutations of the *Tol2-SMyHC*

461 *III*::GFP plasmid were generated using the Agilent QuikChange II XL Site-

462 Directed Mutagenesis Kit® (Cat #200521) following the manufacturer's

463 instructions. Primers are listed in Table 2.

464 **Table 2.** Oligonucleotides (in sense and antisense) used for Hexad mutations.

| Hexad mutations | Oligonucleotide (sense) | Oligonucleotide (antisense) |
|---|---|---|
| **Mut A** | 5'-gaGTCGACaagaggggacaaagaggcggaggt-3' | 5'-acctccgcctctttgtcccctcttGTCGACtc-3' |
| **Mut B** | 5'-cttgcgaaggacaaagTTgggacaaagaggcggag-3' | 5'-ctccgcctctttgtcccAActttgtccttcgcaag-3' |
| **Mut C** | 5'-aggggacaaagaTTcggaggtggggctgg-3' | 5'-ccagcccccacctccgAAtctttgtcccctc-3' |
| **Mut S** | 5'-gaaggacaaagagggggaTTaagaggcggaggt-3' | 5'-acctccgcctcttAAtcccctctttgtccttc-3' |

465

466 To construct the chimeric promoter *5xcNRE-vmhc*, we cloned five tandem

467 repeats of the cNRE into the Xho I site of the *Tol2-vmhc*::GFP vector. The

468 5xcNRE sequence was obtained by annealing the following two oligonucleotides:

469 5'-

470 CTAGGAAGGACAAAGAGGGGACAAAGAGGCGGAGGTGAAGGACAAAGAG

471 GGGACAAAGAGGCGGAGGTGAAGGACAAAGAGGGGACAAAGAGGCGGA

472 GCTGAAGGACAAAGAGGGGACAAAGAGGCGGAGGTGAAGGACAAAGAGG

473 GGACAAAGAGGCGGAGGTCTCGAGA-3'                (sense)        and        5'-

474 GATCTCTCGAGACCTCCGCCTCTTTGTCCCCTCTTTGTCCTTCACCTCCGC

475 CTCTTTGTCCCCTCTTTGTCCTTCACCTCCGCCTCTTTGTCCCCTCTTTGTC

476 CTTCACCTCCGCCTCTTTGTCCCCTCTTTGTCCTTCACCTCCGCCTCTTTGT

477 CCCCTCTTTGTCCTTC-3' (antisense). For the transient transgenic assays, all

478 *Tol2*-based constructs were co-injected (~1nL) into the cytoplasm of one-cell

479 stage embryos with transposase mRNA, which was transcribed from the pCS-TP

480 vector using the mMESSAGE mMACHINE SP6 Kit (Ambion). The master mix for

481 injections was freshly prepared with 125 ng of the plasmid of interest, 175 ng of

482 transposase mRNA, 1 µL of 0,5% phenol red and water to complete the final

483 volume to 5 µL (Suster et al., 2009). All constructs were microinjected in at least

484 two independent experiments. Zebrafish embryos were staged and maintained

485 at 28.5°C, as previously described (Westerfield, 2000), and analyzed at 48 hpf.

486 **Image analyses and processing**

487 All imaging analyses were performed under a NIKON SMZ 25 fluorescent

488 stereomicroscope, and confocal imaging was carried out using a Leica SP8

489 microscope.

**Statistical tests**

We used either a chi-square test or an unpaired t-test (non-parametric when appropriate), and, for both analyses, a 95% confidence value was used to assess significance ($p < 0.05$). The data were presented in column graphs with standard deviation using GraphPad Prism 6. Each experimental group was compared to its respective control group.

**Searching animal genomes for cNRE-like sequences**

To search different animal genomes for cNRE-like sequences, the command-line version of BLAST (BLAST+, version 2.5.0) was used. Animal taxa included in the analyses were those with a genome assembly available on NCBI and with phylogenetic information available on TimeTree.org. The search was conducted using blastn-short, as this allowed for a search tailored to a query sequence as short as the cNRE. All parameters were set to default, except for word size (which was always set to 4 to maximize search accuracy) and gapopen, gapextend, reward and penalty, all of which were changed according to the search parameter stringency (Table 3).

**Table 3.** Search parameters for cNRE sequence queries.

| Parameter stringency | Gapopen | Gapextend | Reward | Penalty |
|---|---|---|---|---|
| **Lenient** | 2 | 1 | 1 | -1 |
| **Moderate** | 1 | 1 | 1 | -3 |
| **Stringent** | 5 | 2 | 1 | -3 |

508     While a gapopen/gapextend ratio of 1/-3 (in the moderate and stringent

509     parameter sets) yielded a sequence conservation of 99%, a ratio of 1/-1 (in the

510     lenient parameter set) resulted in a sequence conservation of 75%. Hits matching

511     at least 30 bp of the 32-bp cNRE were retained.

512     Hit sequences against animal genomes were individually extracted from the

513     databases, since BLAST+ is a local alignment algorithm that only yields parts of

514     hit regions within a given genome. Each animal genome was thus loaded into R

515     to obtain the exact position of the cNRE hit from the BLAST+ alignment. If the

516     best hit for the cNRE was on the minus strand, the genome sequence was

517     reverse complemented. Each sequence was subsequently extracted from the

518     genome based on the genomic location of the initial BLAST+ hit and manually

519     refined to allow a full alignment of the cNRE with a given genome sequence. Of

520     note, the *Coturnix japonica* cNRE sequence was the only one, where gaps

521     needed to be introduced to properly align the cNRE. The gaps were manually

522     introduced into the extracted sequence using R to match the sequence in the

523     BLAST+ results. All extracted sequences were compiled into a single FASTA file

524     and aligned using CLUSTAL OMEGA. The output alignment file was processed

525     in Jalview to create the final alignment figure.

526     **Searching through the viral genomes**

527     All available avian viruses in the NCBI virus database (as of February 19, 2020)

528     were scanned for the cNRE sequence. For this, the R function pairwise alignment

529     (in the Biostrings package, version 2.52.0) was used to create, for each virus, a

530     local-global alignment score of the single best hit for the cNRE sequence and the

531     cNRE reverse complement sequence. The "pattern" was set to the viral genome,

532  the "subject" to the cNRE (or the cNRE reverse complement), the "type" to "local-

533  global" and all other settings were set to default. To calculate the p-value

534  corresponding to a viral hit, 10,000 32 bp-long random DNA sequences matching

535  the cNRE length were generated and the pairwise alignment of each sequence

536  was determined against each viral genome. The same set of random sequences

537  was used to determine the p-value of every virus analyzed. If a cNRE (or a cNRE

538  reverse complement) hit in a viral genome had a local-global pairwise alignment

539  score of more than 95% of the alignment score of the random sequences against

540  this viral genome, this hit was considered as statistically significant (**Figure S6**).

541  We retained top 2,000 viral genomes, based on the pairwise alignment scores of

542  the cNRE (or cNRE reverse complement). Of these, only the best scoring virus

543  of each viral family was chosen for further analysis. This measure ensured that

544  no viral species or genus was overrepresented in the list of hits. Of the 2,000 viral

545  genomes, we thus obtained 27 unique virus families and 4 unclassified viral hits.

546  The p-value of these 31 viruses was subsequently determined, with viruses

547  returning a p-value of 0.05 or less having been included in **Figure 5**. These hit

548  sequences of viral genomes were saved in a single FASTA file and the FASTA

549  file was processed using CLUSTAL OMEGA to produce a multiple sequence

550  alignment. If a hit in a virus genome was to the cNRE reverse complement, the

551  viral hit sequence was reverse complemented before it was included in the

552  FASTA file. The output file was then opened in Jalview and processed to highlight

553  different levels of sequence conservation.

554

555

**Acknowledgments**

We thank Dr. Koichi Kawakami for providing the plasmid used to produce transgenic zebrafish. We thank the members of the Ramialison group and Akriti Varshney, Gulrez Chahal, and Julian Stolper for their feedback and support and Jeannette Hallab, Jeanette Rientjes, Ekaterina Salimova for their assistance with experimental troubleshooting and design. We thank the members of the Monash Bioinformatics Platform for invaluable advice, especially Stuart Archer, Adele Barugahare, Paul Harrison, David Powell, Michael See and Nick Wong. We thank the ARMI FishCore staff and Melbourne University Fish Facility. We also thank Lucas Buscaratti for his drawing created with BioRender.com that helped us to illustrate the different mutations on the *SMyHC III* promoter comparing mice and zebrafish.

**Declarations**

- **Ethics approval and consent to participate**

We confirm that all relevant ethical guidelines have been followed, and any necessary IRB and/or ethics committee approvals have been obtained. The present study was approved by Ethics Committee on Animal Use, Institute of Biomedical Sciences, University of São Paulo (CEUA-ICB/USP), by Ethics Committee on Animal Use, Brazilian Center for Research in Energy and Materials (CEUA-CNPEM) protocol number 24, by the University of Melbourne guidelines and local ethics committee and by the Monash University Animal Ethics Committee.

579   ● **Availability of data and materials**

580   The datasets used and/or analyzed during the current study are available from

581   the corresponding authors on reasonable request, and available in the GitHub

582   repository: https://github.com/mart-nik/cNRE_project.git.

583   ● **Competing interests**

584   The authors declare that they have no competing interests.

585   ● **Funding**

586   This work was supported in part by FAPESP grants 00/04082-1, 03/06555-2,

587   15/12549-2, and 18/09839-7 and by the Coordenação de Aperfeiçoamento de

588   Pessoal de Nível Superior – Brasil (CAPES) – Finance Code 001, by CNPq grant

589   481983/2013-9 and by the Ceara State Scientist-in-chief program of Fundação

590   Cearense de Apoio ao Desenvolvimento Científico e Tecnológico (FUNCAP

591   08908197/2019) to José Xavier Neto. MR is supported by Grants from the

592   Australian Research Council and the NHMRC. The Australian Regenerative

593   Medicine Institute is supported by grants from the State Government of Victoria

594   and the Australian Government. MS is funded by the CNRS.

595   ● **Authors' contributions**

596   LNS generated transgenic zebrafish, contributed to discussions, and wrote this

597   manuscript. MN and MR performed bioinformatics analyses, contributed to

598   discussions and co-wrote the manuscript. AMSC, PJ and SD contributed to

599   transgenic analyses. HAC, MV, TGFM and ACS carried out experiments on

600    transgenic mice. MBG and PO performed bioinformatics analyses. FES, NR,

601    GFW, WK and MS contributed to the conception of this study. JXN was

602    responsible for the project design, mouse transgenics, scientific discussions and

603    manuscript writing. All authors have read and approved the final manuscript.

604

605    **References**

606    Bruneau, B. G., Bao, Z.-Z., Tanaka, M., Schott, J.-J., Izumo, S., Cepko, C. L.,

607        Seidman, J. G., & Seidman, C. E. (2000). Cardiac Expression of the

608        Ventricle-Specific Homeobox Gene Irx4 Is Modulated by Nkx2-5 and

609        dHand. *Developmental Biology*, *217*(2), 266–277.

610        https://doi.org/10.1006/DBIO.1999.9548

611    Chen, Q., Moore, L. A., Wick, M., & Bandman, E. (1997). Identification of a

612        genomic locus containing three slow myosin heavy chain genes in the

613        chicken. *Biochimica et Biophysica Acta - Gene Structure and Expression*,

614        *1353*(2), 148–156. https://doi.org/10.1016/S0167-4781(97)00067-5

615    Erwin, D. H., & Davidson, E. H. (2002). *The last common bilaterian ancestor*.

616    Holmes, E. C. (2011). The evolution of endogenous viral elements. In *Cell Host*

617        *and Microbe* (Vol. 10, Issue 4, pp. 368–377). Cell Press.

618        https://doi.org/10.1016/j.chom.2011.09.002

619    Jin, D., Ni, T. T., Hou, J., Rellinger, E., & Zhong, T. P. (2009). Promoter analysis

620        of ventricular myosin heavy chain (vmhc) in zebrafish embryos.

621        *Developmental Dynamics : An Official Publication of the American*

622    *Association of Anatomists*, *238*(7), 1760–1767.

623    https://doi.org/10.1002/dvdy.22000

624    Kaleta, E. F. (1990). Herpesviruses of birds - a review. In *Avian Pathology* (Vol.

625    19, Issue 2, pp. 193–211). Avian Pathol.

626    https://doi.org/10.1080/03079459008418673

627    Kazazian, H. H. (2004). Mobile Elements: Drivers of Genome Evolution. In

628    *Science* (Vol. 303, Issue 5664, pp. 1626–1632). American Association for

629    the Advancement of Science. https://doi.org/10.1126/science.1089670

630    Kuhl, H., Frankl-Vilches, C., Bakker, A., Mayr, G., Nikolaus, G., Boerno, S. T.,

631    Klages, S., Timmermann, B., & Gahr, M. (2021). An Unbiased Molecular

632    Approach Using 3'-UTRs Resolves the Avian Family-Level Tree of Life.

633    *Molecular Biology and Evolution*, *38*(1).

634    https://doi.org/10.1093/MOLBEV/MSAA191

635    McGeoch, D. J., Dolan, A., & Ralph, A. C. (2000). Toward a Comprehensive

636    Phylogeny for Mammalian and Avian Herpesviruses. *Journal of Virology*,

637    *74*(22), 10401–10406. https://doi.org/10.1128/jvi.74.22.10401-10406.2000

638    Meyers, J. R. (2018). Zebrafish: Development of a Vertebrate Model Organism.

639    *Current Protocols in Essential Laboratory Techniques*, *16*(1).

640    https://doi.org/10.1002/cpet.19

641    Morissette, G., & Flamand, L. (2010). Herpesviruses and Chromosomal

642    Integration. *Journal of Virology*, *84*(23), 12100–12109.

643    https://doi.org/10.1128/jvi.01169-10

29

644  Nair, V. (2005). Evolution of Marek's disease - A paradigm for incessant race

645      between the pathogen and the host. *Veterinary Journal*, *170*(2), 175–183.

646      https://doi.org/10.1016/j.tvjl.2004.05.009

647  Nikovits, W., Wang, G. F., Feldman, J. L., Miller, J. B., Wade, R., Nelson, L., &

648      Stockdale, F. E. (1996). Isolation and characterization of an avian slow

649      myosin heavy chain gene expressed during embryonic skeletal muscle

650      fiber formation. *The Journal of Biological Chemistry*, *271*(29), 17047–

651      17056. http://www.ncbi.nlm.nih.gov/pubmed/8663323

652  Rastinejad, F., Wagner, T., Zhao, Q., & Khorasanizadeh, S. (2000). Structure of

653      the RXR-RAR DNA-binding complex on the retinoic acid response element

654      DR1. *The EMBO Journal*, *19*(5), 1045–1054.

655      https://doi.org/10.1093/emboj/19.5.1045

656  Sekijima, M., Motono, C., Yamasaki, S., Kaneko, K., & Akiyama, Y. (2003).

657      Molecular dynamics simulation of prion protein by large scale cluster

658      computing. *Lecture Notes in Computer Science (Including Subseries*

659      *Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*,

660      *2858*, 476–485. https://doi.org/10.1007/978-3-540-39707-6_43

661  Shaffer, P. L., & Gewirth, D. T. (2002). Structural basis of VDR–DNA

662      interactions on direct repeat response elements. *The EMBO Journal*, *21*(9),

663      2242. https://doi.org/10.1093/EMBOJ/21.9.2242

664  Simões-Costa, M. S., Vasconcelos, M., Sampaio, A. C., Cravo, R. M., Linhares,

665      V. L., Hochgreb, T., Yan, C. Y. I., Davidson, B., & Xavier-Neto, J. (2005).

666      The evolutionary origin of cardiac chambers. *Developmental Biology*,

667      *277*(1), 1–15. https://doi.org/10.1016/J.YDBIO.2004.09.026

668    Suster, M. L., Kikuta, H., Urasaki, A., Asakawa, K., & Kawakami, K. (2009).

669      Transgenesis in Zebrafish with the Tol2 Transposon System. In *Methods in*

670      *molecular biology (Clifton, N.J.)* (Vol. 561, pp. 41–63).

671      https://doi.org/10.1007/978-1-60327-019-9_3

672    Waardenberg, A. J., Ramialison, M., Bouvere, R., & Harvey, R. P. (2014).

673      Genetic networks governing heart development. *Cold Spring Harbor*

674      *Perspectives in Medicine*, *4*(11), 1–23.

675      https://doi.org/10.1101/cshperspect.a013839

676    Wang, G F, Nikovits, W., Schleinitz, M., & Stockdale, F. E. (1998). A positive

677      GATA element and a negative vitamin D receptor-like element control atrial

678      chamber-specific expression of a slow myosin heavy-chain gene during

679      cardiac morphogenesis. *Molecular and Cellular Biology*, *18*(10), 6023–

680      6034. https://doi.org/10.1128/MCB.18.10.6023

681    Wang, Gang Feng, Nikovits, W., Bao, Z. Z., & Stockdale, F. E. (2001). Irx4

682      Forms an Inhibitory Complex with the Vitamin D and Retinoic X Receptors

683      to Regulate Cardiac Chamber-specific slow MyHC3 Expression. *Journal of*

684      *Biological Chemistry*, *276*(31), 28835–28841.

685      https://doi.org/10.1074/jbc.M103716200

686    Wang, Gang Feng, Nikovits, W., Schleinitz, M., & Stockdale, F. E. (1996). Atrial

687      chamber-specific expression of the slow myosin heavy chain 3 gene in the

688      embryonic heart. *Journal of Biological Chemistry*, *271*(33), 19836–19845.

689      https://doi.org/10.1074/jbc.271.33.19836

690     Westerfield, M. (2000). *The Zebrafish Book: A guide for the Laboratory Use of*

691        *Zebrafish (Danio rerio)* (I. of N. U. of Oregon (ed.); 4 ed.). Institute of

692        Neuroscience. University of Oregon.

693     Xavier-Neto, J, Neville, C. M., Shapiro, M. D., Houghton, L., Wang, G. F.,

694        Nikovits, W., Stockdale, F. E., & Rosenthal, N. (1999). A retinoic acid-

695        inducible transgenic marker of sino-atrial development in the mouse heart.

696        *Development (Cambridge, England)*, *126*(12), 2677–2687.

697        http://www.ncbi.nlm.nih.gov/pubmed/10331979

698     Xavier-Neto, José, Trueba, S. S., Stolfi, A., Souza, H. M., Sobreira, T. J. P.,

699        Schubert, M., & Castillo, H. A. (2012). An Unauthorized Biography of the

700        Second Heart Field and a Pioneer/Scaffold Model for Cardiac

701        Development. In *Current Topics in Developmental Biology* (Vol. 100, pp.

702        67–105). Academic Press Inc. https://doi.org/10.1016/B978-0-12-387786-

703        4.00003-8

704     Zhang, R., & Xu, X. (2009). Transient and transgenic analysis of the zebrafish

705        ventricular myosin heavy chain (vmhc) promoter: an inhibitory mechanism

706        of ventricle-specific gene expression. *Developmental Dynamics : An Official*

707        *Publication of the American Association of Anatomists*, *238*(6), 1564–1573.

708        https://doi.org/10.1002/dvdy.21929

709     Zhao, Q., Chasse, S. A., Devarakonda, S., Sierk, M. L., Ahvazi, B., &

710        Rastinejad, F. (2000). Structural basis of RXR-DNA interactions. *Journal of*

711        *Molecular Biology*, *296*(2), 509–520.

712        https://doi.org/10.1006/JMBI.1999.3457

713

714  Bruneau, B. G., Bao, Z.-Z., Tanaka, M., Schott, J.-J., Izumo, S., Cepko, C. L.,

715      Seidman, J. G., & Seidman, C. E. (2000). Cardiac Expression of the

716      Ventricle-Specific Homeobox Gene Irx4 Is Modulated by Nkx2-5 and

717      dHand. *Developmental Biology*, *217*(2), 266–277.

718      https://doi.org/10.1006/DBIO.1999.9548

719  Chen, Q., Moore, L. A., Wick, M., & Bandman, E. (1997). Identification of a

720      genomic locus containing three slow myosin heavy chain genes in the

721      chicken. *Biochimica et Biophysica Acta - Gene Structure and Expression*,

722      *1353*(2), 148–156. https://doi.org/10.1016/S0167-4781(97)00067-5

723  Erwin, D. H., & Davidson, E. H. (2002). *The last common bilaterian ancestor*.

724  Holmes, E. C. (2011). The evolution of endogenous viral elements. In *Cell Host*

725      *and Microbe* (Vol. 10, Issue 4, pp. 368–377). Cell Press.

726      https://doi.org/10.1016/j.chom.2011.09.002

727  Jin, D., Ni, T. T., Hou, J., Rellinger, E., & Zhong, T. P. (2009). Promoter analysis

728      of ventricular myosin heavy chain (vmhc) in zebrafish embryos.

729      *Developmental Dynamics : An Official Publication of the American*

730      *Association of Anatomists*, *238*(7), 1760–1767.

731      https://doi.org/10.1002/dvdy.22000

732  Kaleta, E. F. (1990). Herpesviruses of birds - a review. In *Avian Pathology* (Vol.

733      19, Issue 2, pp. 193–211). Avian Pathol.

734      https://doi.org/10.1080/03079459008418673

735 Kazazian, H. H. (2004). Mobile Elements: Drivers of Genome Evolution. In

736  *Science* (Vol. 303, Issue 5664, pp. 1626–1632). American Association for

737  the Advancement of Science. https://doi.org/10.1126/science.1089670

738 Kuhl, H., Frankl-Vilches, C., Bakker, A., Mayr, G., Nikolaus, G., Boerno, S. T.,

739  Klages, S., Timmermann, B., & Gahr, M. (2021). An Unbiased Molecular

740  Approach Using 3'-UTRs Resolves the Avian Family-Level Tree of Life.

741  *Molecular Biology and Evolution*, *38*(1).

742  https://doi.org/10.1093/MOLBEV/MSAA191

743 McGeoch, D. J., Dolan, A., & Ralph, A. C. (2000). Toward a Comprehensive

744  Phylogeny for Mammalian and Avian Herpesviruses. *Journal of Virology*,

745  *74*(22), 10401–10406. https://doi.org/10.1128/jvi.74.22.10401-10406.2000

746 Meyers, J. R. (2018). Zebrafish: Development of a Vertebrate Model Organism.

747  *Current Protocols in Essential Laboratory Techniques*, *16*(1).

748  https://doi.org/10.1002/cpet.19

749 Morissette, G., & Flamand, L. (2010). Herpesviruses and Chromosomal

750  Integration. *Journal of Virology*, *84*(23), 12100–12109.

751  https://doi.org/10.1128/jvi.01169-10

752 Nair, V. (2005). Evolution of Marek's disease - A paradigm for incessant race

753  between the pathogen and the host. *Veterinary Journal*, *170*(2), 175–183.

754  https://doi.org/10.1016/j.tvjl.2004.05.009

755 Nikovits, W., Wang, G. F., Feldman, J. L., Miller, J. B., Wade, R., Nelson, L., &

756  Stockdale, F. E. (1996). Isolation and characterization of an avian slow

757  myosin heavy chain gene expressed during embryonic skeletal muscle

758    fiber formation. *The Journal of Biological Chemistry*, *271*(29), 17047–

759    17056. http://www.ncbi.nlm.nih.gov/pubmed/8663323

760    Rastinejad, F., Wagner, T., Zhao, Q., & Khorasanizadeh, S. (2000). Structure of

761    the RXR-RAR DNA-binding complex on the retinoic acid response element

762    DR1. *The EMBO Journal*, *19*(5), 1045–1054.

763    https://doi.org/10.1093/emboj/19.5.1045

764    Sekijima, M., Motono, C., Yamasaki, S., Kaneko, K., & Akiyama, Y. (2003).

765    Molecular dynamics simulation of prion protein by large scale cluster

766    computing. *Lecture Notes in Computer Science (Including Subseries*

767    *Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*,

768    *2858*, 476–485. https://doi.org/10.1007/978-3-540-39707-6_43

769    Shaffer, P. L., & Gewirth, D. T. (2002). Structural basis of VDR–DNA

770    interactions on direct repeat response elements. *The EMBO Journal*, *21*(9),

771    2242. https://doi.org/10.1093/EMBOJ/21.9.2242

772    Simões-Costa, M. S., Vasconcelos, M., Sampaio, A. C., Cravo, R. M., Linhares,

773    V. L., Hochgreb, T., Yan, C. Y. I., Davidson, B., & Xavier-Neto, J. (2005).

774    The evolutionary origin of cardiac chambers. *Developmental Biology*,

775    *277*(1), 1–15. https://doi.org/10.1016/J.YDBIO.2004.09.026

776    Suster, M. L., Kikuta, H., Urasaki, A., Asakawa, K., & Kawakami, K. (2009).

777    Transgenesis in Zebrafish with the Tol2 Transposon System. In *Methods in*

778    *molecular biology (Clifton, N.J.)* (Vol. 561, pp. 41–63).

779    https://doi.org/10.1007/978-1-60327-019-9_3

780    Waardenberg, A. J., Ramialison, M., Bouvere, R., & Harvey, R. P. (2014).

781    Genetic networks governing heart development. *Cold Spring Harbor*

782    *Perspectives in Medicine*, *4*(11), 1–23.

783    https://doi.org/10.1101/cshperspect.a013839

784    Wang, G F, Nikovits, W., Schleinitz, M., & Stockdale, F. E. (1998). A positive

785    GATA element and a negative vitamin D receptor-like element control atrial

786    chamber-specific expression of a slow myosin heavy-chain gene during

787    cardiac morphogenesis. *Molecular and Cellular Biology*, *18*(10), 6023–

788    6034. https://doi.org/10.1128/MCB.18.10.6023

789    Wang, Gang Feng, Nikovits, W., Bao, Z. Z., & Stockdale, F. E. (2001). Irx4

790    Forms an Inhibitory Complex with the Vitamin D and Retinoic X Receptors

791    to Regulate Cardiac Chamber-specific slow MyHC3 Expression. *Journal of*

792    *Biological Chemistry*, *276*(31), 28835–28841.

793    https://doi.org/10.1074/jbc.M103716200

794    Wang, Gang Feng, Nikovits, W., Schleinitz, M., & Stockdale, F. E. (1996). Atrial

795    chamber-specific expression of the slow myosin heavy chain 3 gene in the

796    embryonic heart. *Journal of Biological Chemistry*, *271*(33), 19836–19845.

797    https://doi.org/10.1074/jbc.271.33.19836

798    Westerfield, M. (2000). *The Zebrafish Book: A guide for the Laboratory Use of*

799    *Zebrafish (Danio rerio)* (I. of N. U. of Oregon (ed.); 4 ed.). Institute of

800    Neuroscience. University of Oregon.

801    Xavier-Neto, J, Neville, C. M., Shapiro, M. D., Houghton, L., Wang, G. F.,

802    Nikovits, W., Stockdale, F. E., & Rosenthal, N. (1999). A retinoic acid-

803    inducible transgenic marker of sino-atrial development in the mouse heart.

804    *Development (Cambridge, England)*, *126*(12), 2677–2687.

805    http://www.ncbi.nlm.nih.gov/pubmed/10331979

806    Xavier-Neto, José, Trueba, S. S., Stolfi, A., Souza, H. M., Sobreira, T. J. P.,

807    Schubert, M., & Castillo, H. A. (2012). An Unauthorized Biography of the

808    Second Heart Field and a Pioneer/Scaffold Model for Cardiac

809    Development. In *Current Topics in Developmental Biology* (Vol. 100, pp.

810    67–105). Academic Press Inc. https://doi.org/10.1016/B978-0-12-387786-

811    4.00003-8

812    Zhang, R., & Xu, X. (2009). Transient and transgenic analysis of the zebrafish

813    ventricular myosin heavy chain (vmhc) promoter: an inhibitory mechanism

814    of ventricle-specific gene expression. *Developmental Dynamics : An Official*

815    *Publication of the American Association of Anatomists*, *238*(6), 1564–1573.

816    https://doi.org/10.1002/dvdy.21929

817    Zhao, Q., Chasse, S. A., Devarakonda, S., Sierk, M. L., Ahvazi, B., &

818    Rastinejad, F. (2000). Structural basis of RXR-DNA interactions. *Journal of*

819    *Molecular Biology*, *296*(2), 509–520.

820    https://doi.org/10.1006/JMBI.1999.3457

821

822    **Figure legends**

823    **Figure 1. The cNRE is necessary and sufficient to drive expression in atrial**

824    **cells. A)** Schematic representation of the *SMyHC III* promoter sequence

825    highlighting the position of the cNRE sequence. **B)** Confocal image in frontal

826    views, anterior is to the top of a representative zebrafish embryos. Exclusive

827    ventricular expression is demonstrated by overlapping GFP expression driven by

828    ΔcNRE and stable mCherry fluorescence driven by the ventricular stable line. **C)**

829    Frontal views, anterior is to the top, mouse embryos. *SMyHC III*::HAP isolated

830    heart at 10.5 dpc showing intense, dark blue, atrial coloring indicative of high HAP

831    expression and *SMyHC IIIΔcNRE*::HAP isolated heart at 10.5 dpc showing

832    absence of HAP expression. **D)** HAP assays in homogenates of atrial and non-

833    atrial cardiac tissues in *SMyHC III*::HAP (n= 18), p<0.0001. **E)** HAP assays in

834    homogenates of atrial and non-atrial cardiac tissues in the *SMyHC*

835    *IIIΔcNRE*::HAP mutant (n= 16), p=0.0013. **F)** Confocal image in lateral views,

836    anterior is to the left of a representative zebrafish embryo. Exclusive ventricular

837    GFP expression is observed at 48 hpf when injected with the *vmhc* promoter and

838    **G)** a representative embryo expressing GFP in both heart chambers at 48 hpf

839    when injected with the *5xcNRE-vmhc* construct. **H)** Graphical analysis of

840    chamber expression patterns of the cohort of embryos injected with *vmhc* or

841    *5XcNRE-vmhc* promoter constructs. (at) atrium. (vt) ventricle. chi-square test,

842    p<0.05, comparing *vmhc::*GFP and *5xcNRE-vmhc::*GFP embryos for each

843    condition. Scale bars are 30 μm.

844    **Figure 2. Mutational analysis of the *SMyHC III* promoter in zebrafish reveals**

845    **a dual role in atrial activation and ventricular repression. A)** Representative

846    panel of GFP expression patterns in cardiac chambers of zebrafish embryos in

847    lateral views, anterior to the left, injected with *SMyHC III*::GFP. **B)** Graphic

848    representation of GFP chamber expression patterns of cohorts of embryos

849    injected with *SMyHC III*::GFP, *SMyHC IIIΔcNRE*::GFP and constructs containing

850    point mutations in the cNRE Hexads A, B and C (Mut A, B and C, respectively)

851    as well as a non-Hexad control mutation (Mut S). Embryos were analyzed at 48

852  hpf and classified into three categories of cardiac expression patterns: exclusive

853  atrium (at), exclusive ventricular (vt) and atrium plus ventricular (at+vt). chi-

854  square test, p<0.05, comparing *SMyHC III* to each mutation and condition.

855  **Figure 3. Point mutations in Hexads B and C affect HAP expression in the**

856  **heart. A)** Strategy for the mutation of Hexad A (Mut A). **B)** Mouse line 14 (Mut A)

857  at 10.5 dpc, showing atrial-specific expression of HAP. **C)** Strategy for the

858  mutation of Hexad B (Mut B). **D)** Mouse line 17 (Mut B) at 10.5 dpc, showing

859  expression of HAP. **E-G)** Time course (10.5 dpc to 12.5 dpc) of cardiac

860  expression in both chambers (atrium and ventricle) in mouse line 17 (Mut B). **H)**

861  Comparison of HAP assays in homogenates of atrial and non-atrial cardiac

862  tissues from line 17 (Mut B) in 10.5 dpc embryos (n= 38). **I)** Strategy for the

863  mutation of Hexad C (Mut C). **J)** Mouse line 5 (Mut C) at 9.5 dpc. **K)** Isolated

864  heart from the *SMyHC III*::HAP line at 10.5 dpc. **L)** Isolated heart from a wild-type

865  littermate at 10.5 dpc. **M)** Isolated heart from mouse line 5 (Mut C) at 10.5 dpc.

866  **N)** HAP assays in homogenates of atrial and non-atrial cardiac tissues from line

867  5 (Mut C) in 10.5 dpc embryos (n= 15). **O)** Strategy for the mutation of the non-

868  Hexad control (Mut S) in the spacer sequence between Hexads B and C. **P)**

869  Representative mouse line 14 (Mut S) at 10.5 dpc, showing atrial-specific HAP

870  expression.

871  **Figure 4. Comparison of cNRE mutations between zebrafish and mice. A)**

872  *SMyHC III* promoter in **A')** zebrafish and **A'')** mice. **B)** cNRE deletion in **B')**

873  zebrafish and **B'')** mice. **C)** Mutation of Hexad A (Mut A) in **C')** zebrafish and **C'')**

874  mice. **D)** Mutation of Hexad B in **D')** zebrafish and **D'')** mice. **E)** Mutation of Hexad

875  C in **E')** zebrafish and **E'')** mice. **F)** Mutation of the spacer sequence between

876 Hexads B and C (Mut S) in **F')** zebrafish and **F'')** mice. Zebrafish hearts in lateral

877 views and mouse hearts in frontal view. (at) atrium. (vt) ventricle.

878 **Figure 5. Evolutionary conservation of the cNRE. A)** Phylogenetic tree of

879 cNRE sequences in American alligator (*Alligator mississippiensis*), Chinese

880 alligator (*Alligator sinensis*), ostrich (*Struthio camelus*), pigeon (*Columba livia*),

881 hoatzin (*Opisthocomus hoazin*), zebra finch (*Taeniopygia guttata*), helmeted

882 guineafowl (*Numida meleagris*), pinnated grouse (*Tympanuchus cupido*), black

883 grouse (*Lyrurus tetrix*), wild turkey (*Meleagris gallopavo*), chicken (*Gallus gallus*),

884 Indian peafowl (*Pavo cristatus*), golden pheasant (*Chrysolophus pictus*), common

885 pheasant (*Phasianus colchicus*), Mikado pheasant (*Syrmaticus mikado*),

886 Japanese quail (*Coturnix japonica*), Chinese bamboo partridge (*Bambusicola*

887 *thoracicus*), scaled quail (*Callipepla squamata*), bobwhite (*Colinus virginianus*),

888 mallard (*Anas platyrhynchos*), eastern spot-billed duck (*Anas zonorhyncha*),

889 tufted duck (*Aythya fuligula*), Muscovy duck (*Cairina moschata*), pink-footed

890 goose (*Anser brachyrhynchus*), swan goose (*Anser cygnoides*) and Canada

891 goose (*Branta canadensis*). For each species, the sequence displaying the

892 statistically most significant hit with the original cNRE is shown, along with the

893 corresponding E-value. **B)** Alignment of the original cNRE with sequences

894 identified in different viruses, along with their respective p-values. (Hex) Hexads.

895 **Supplementary Figure 1. GFP expression in cardiac and non-cardiac**

896 **tissues of zebrafish embryos at 48 hpf following injection with different**

897 **constructs of the *SMyHC III* promoter.** Graphic representation of GFP

898 expression in cardiac and other tissues of cohorts of embryos injected with

899 different *SMyHC III* promoter constructs and analyzed at 48 hpf. chi-square test,

900 $p < 0.05$, comparing *SMyHC III* to each mutation and condition.

901 **Supplementary Figure 2. *SMyHC III*::HAP embryos expressing HAP**

902 **specifically in the atrium. A-D)** HAP expression at 10.5 dpc in mouse lines 1,

903 5, 6 and 27.

904 **Supplementary Figure 3. Mutation of the distal 72 bp fragment containing**

905 **the cNRE sequence alters atrial HAP expression driven by the *SMyHC III***

906 **promoter. A)** Strategy for the deletion of the distal 72 bp of the *SMyHC III*

907 promoter. **B)** *SMyHC IIIΔcNRE*::HAP at 10.5 dpc in mouse line 30. **C)** *SMyHC*

908 *IIIΔcNRE*::HAP at 10.5 dpc in mouse line 11. **D)** Isolated embryonic heart of

909 *SMyHC IIIΔcNRE*::HAP mouse line 11 at 9.5 dpc. **E)** Isolated embryonic heart of

910 *SMyHC IIIΔcNRE*::HAP mouse line 11 at 12.5 dpc. (*) distal outflow.

911 **Supplementary Figure 4. Mutation of Hexad A does not affect atrial**

912 **specificity of the *SMyHC III* promoter. A)** Mouse line 1 with mutation of Hexad

913 A (*mutA*::HAP) at 10.5 dpc. **B)** Mouse line 25 with mutation of Hexad A

914 (*mutA*::HAP) at 9.5 dpc.

915 **Supplementary Figure 5. Mutation of Hexad C affects both atrial and**

916 **ventricular expression of the *SMyHC III* promoter. A)** A wild-type non-

917 transgenic embryo displaying background HAP staining at 10.5 dpc. **B)** Mouse

918 line 3 with mutation of Hexad C (*mutC*::HAP) at 10.5 dpc. **C)** Mouse line 7 with

919 mutation of Hexad C (*mutC*::HAP) at 10.5 dpc. **D)** Mouse line 14 with mutation of

920 Hexad C (*mutC*::HAP) at 10.5 dpc.

921 **Supplementary Figure 6. The R-based pairwise alignment method for**

922 **calculating p-values for viral hits.** 1000 random sequences were produced and

923 their pairwise alignment scores were compared to those obtained with the cNRE.

924 A viral hit was considered as statistically significant, when the pairwise alignment

925    score of the cNRE sequence was more than 95% higher than the score obtained
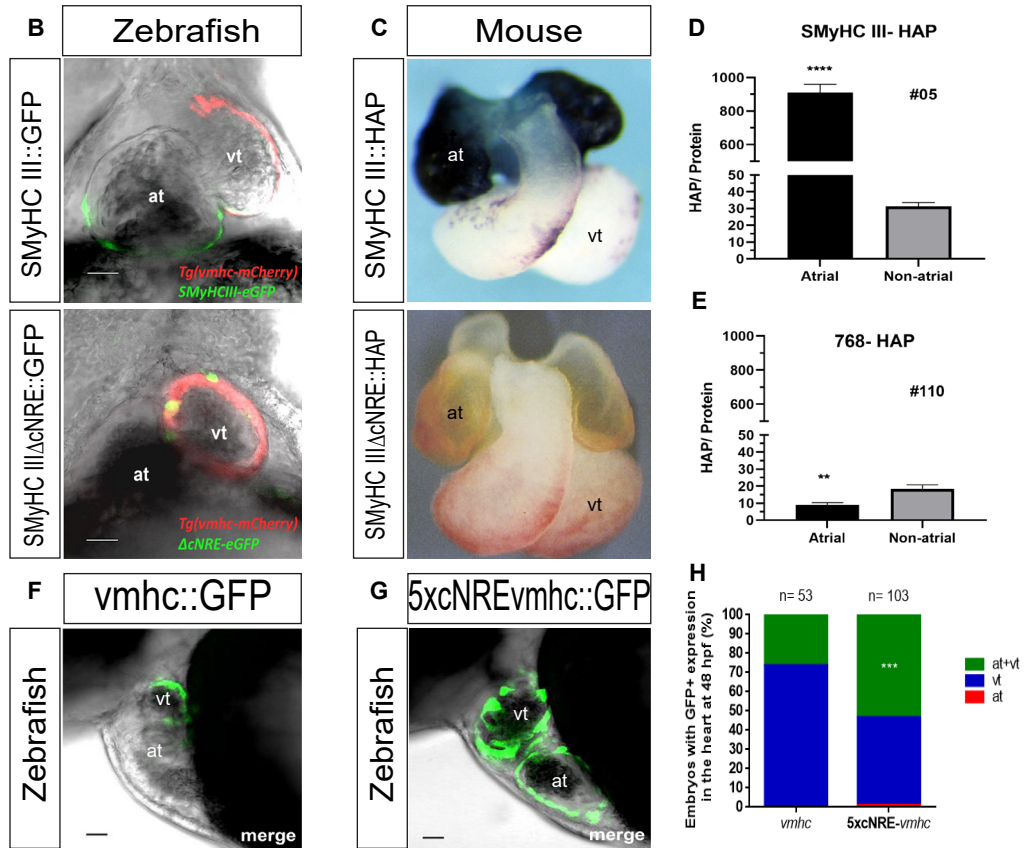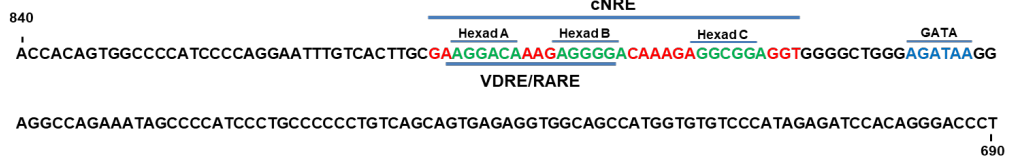
926    with the random sequences.

# A

```
840
|
ACCACAGTGGCCCCATCCCCAGGAATTTGTCACTTGCGAAGGACAAAGAGGGGGACAAAGAGGCGGAGGTGGGGCTGGGAGATAAGG
```

cNRE

Hexad A    Hexad B    Hexad C    GATA

VDRE/RARE

```
AGGCCAGAAATAGCCCCATCCCTGCCCCCCTGTCAGCAGTGAGAGGTGGCAGCCATGGTGTGTCCCATAGAGATCCACAGGGACCCT
                                                                                    |
                                                                                   690
```



**B** Zebrafish

SMyHC III::GFP

*Tg(vmhc-mCherry)*
*SMyHCIII-eGFP*

SMyHC IIIΔcNRE::GFP

*Tg(vmhc-mCherry)*
*ΔcNRE-eGFP*

**C** Mouse

SMyHC III::HAP

SMyHC IIIΔcNRE::HAP

**D** SMyHC III- HAP

#05

**E** 768- HAP

#110

**F** vmhc::GFP

Zebrafish

merge

**G** 5xcNREvmhc::GFP

Zebrafish

merge

**H**

n= 53    n= 103

at+vt
vt
at

*vmhc*    5xcNRE-*vmhc*

**Fig. 1 Santos *et al.,* 2021**

**A**

| Atrium | Ventricle | Atrium + Ventricle |

**B**

Fig. 2 Santos *et al.,* 2021

| MUT A | MUT B | MUT C | MUT S |

**A** Hexad A   Hexad B   Hexad C
GAAGGACACAAGAGGGGACAAAGAGGCGGAGGT
GTCGAT

**C** Hexad A   Hexad B   Hexad C
GAAGGACACAAGAGAGGGGACAAAGAGGCGGAGGT
TT

**I** Hexad A   Hexad B   Hexad C
GAAGGACACAAGAGGGGACAAAGAGGCGGAGGT
TT

**O** Hexad A   Hexad B   Hexad C
GAAGGACAAAGAGGGGACAAAGAGGCGGAGGT
TT

**B** 10.5 dpc   # 14

**D** 10.5 dpc   # 17

**J** 9.5 dpc   # 05

**P** 10.5 dpc   # 14

**E** 10.5 dpc   **F** 11.5 dpc   **G** 12.5 dpc

**K** 10.5 dpc   SMyHC3::HAP   **L** 10.5 dpc   littermate wild-type   **M** 10.5 dpc   mutC-HAP

**H** mutB- HAP
#17
HAP/ Protein
****
Atrial   Non-atrial

**N** mutC- HAP
#05
HAP/ Protein
****
Atrial   Non-atrial

**Fig. 3 Santos *et al.*, 2021**

| SMyHC III promoter Constructs | Zebrafish 48 hpf | Mouse 10.5 dpc |
|---|---|---|

**A** / **A'** / **A"**

cNRE — SMyHC III promoter — reporter

**B** / **B'** / **B"**

GAAGGACAAAGAGGGGACAAAGAGGCCGGAGGT
*Zebrafish deletion of 32 bp (cNRE)
**Mouse deletion of 72 bp (40 bp +cNRE)

**C** / **C'** / **C"**

GA AGGACA AAGAGGGGACAAAGAGGCGGAGGT
Hexad A
↓ Sal-I
GTCGAC
Mut A

**D** / **D'** / **D"**

GAAGGACAAAG AGGGGA CAAAGAGGCGGAGGT
Hexad B
↓
TT
Mut B

**E** / **E'** / **E"**

GAAGGACAAAGAGGGGACAAAGA GG CG GA GGT
Hexad C
↓
TT
Mut C

**F** / **F'** / **F"**

GAAGGACAAAG AGGGGA CA AAGA GGCGGA GGT
Hexad B    Hexad C
↓
TT
Mut S

Mut A, Mut B, Mut C, Mut S

**Fig. 4 Santos *et al.,* 2021**

**A**



**B**



Fig. 5  Santos *et al.*, 2021

**Fig. Sup. 1 Santos *et al.,* 2021**

**SMyHC III::HAP**

A  10.5 dpc  #01

B  10.5 dpc  #05

C  10.5 dpc  #06

D  10.5 dpc  #27

Fig. Sup. 2 Santos *et al.,* 2021

**A**

72 bp 768 bp HAP

**B** 10.5 dpc #30

**C** 10.5 dpc #11

**D** 9.5 dpc * #11

**E** 12.5 dpc * #11

Fig. Sup. 3  Santos *et al.*, 2021

## mutA::HAP



**Fig. Sup. 4 Santos *et al.,* 2021**

**A** 10.5 dpc

non-transgenic WT

**B** 10.5 dpc

mutC::HAP #03

**C** 10.5 dpc

mutC::HAP #07

**D** 10.5 dpc

mutC::HAP #14

**1. Make 1000 random sequences.**

**2. Align them.**

Random sequence    Random sequence    cNRE    Random sequence    Random sequence

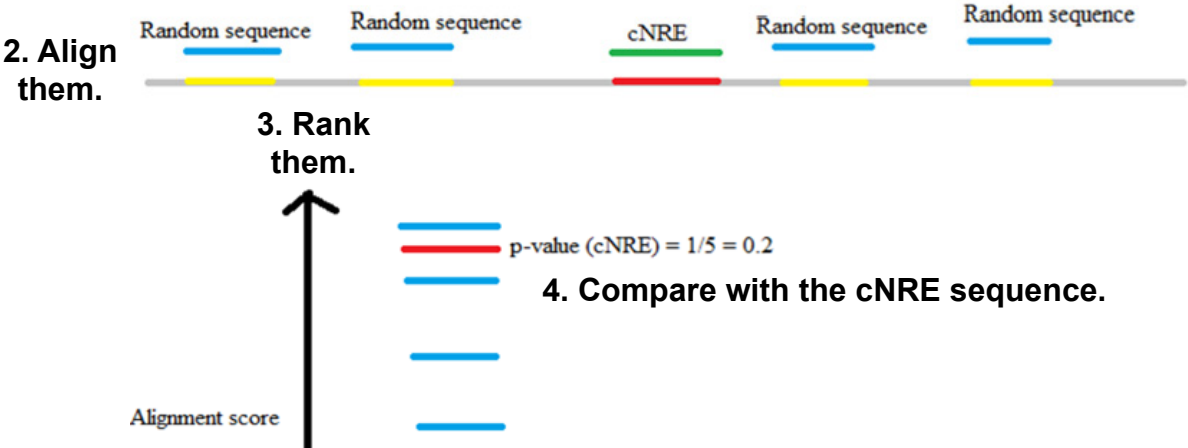**3. Rank them.**

Alignment score

p-value (cNRE) = 1/5 = 0.2

**4. Compare with the cNRE sequence.**

**Fig. Sup. 6 Santos *et al.,* 2021**