

# CardioTF, a database of deconstructing transcriptional circuits in the heart system

Yisong Zhen

**Background.** Information on cardiovascular gene transcription is fragmented and far behind the present requirements of the systems biology field. To create a comprehensive source of data for cardiovascular gene regulation and to facilitate a deeper understanding of genomic data, the CardioTF database was constructed. The purpose of this database is to collate information on cardiovascular transcription factors (TFs), position weight matrices, and enhancer sequences discovered using the ChIP-seq method. **Methods.** The Naïve-Bayes algorithm was used to classify literature and identify all PubMed abstracts on cardiovascular development. The natural language learning tool GNAT was then used to identify corresponding gene names embedded within these abstracts. Local Perl scripts were used to integrate and dump data from public databases into the MariaDB management system (MySQL). In-house R scripts were written to analyze and visualize the results. **Results.** Known cardiovascular TFs from humans and human homologs from fly, *Ciona*, zebrafish, frog, chicken, mouse and human were identified and deposited in the database. Position weight matrices from Jaspar, hPDI, and UniPROBE databases were deposited in the database and can be retrieved using their corresponding TF names. Gene enhancer regions from various sources of ChIP-seq data were deposited into the database and were able to be visualized by graphical output. Besides biocuration, mouse homologs of the 81 core cardiac TFs were selected using a machine learning method and then by intersecting four independent data sources: RNA profiling, expert annotation, PubMed abstract and phenotype. **Discussion.** The CardioTF database can be used as a portal to construct transcriptional network of cardiac development. **Availability and Implementation.** Database URL: <http://www.cardiosignal.org/database/cardiotsf.html>

1 CardioTF, a database of deconstructing transcriptional circuits in  
2 the heart system

3

4

Yisong Zhen<sup>1</sup>

5

6 1. State Key Laboratory of Cardiovascular Disease, Fuwai Hospital, National Center for  
7 Cardiovascular Diseases Chinese Academy of Medical Sciences and Peking Union  
8 Medical College, Beijing, 100037, People's Republic of China

9

10 Corresponding Author:

11 Yisong Zhen

12 Beilishilu, 167, Beijing, 100037 P.R.China

13 Email address: [zhenyisong@fuwaihospital.org](mailto:zhenyisong@fuwaihospital.org)

14 **Abstract**

15 **Background.** Information on cardiovascular gene transcription is fragmented and far behind the  
16 present requirements of the systems biology field. To create a comprehensive source of data for  
17 cardiovascular gene regulation and to facilitate a deeper understanding of genomic data, the  
18 CardioTF database was constructed. The purpose of this database is to collate information on  
19 cardiovascular transcription factors (TFs), position weight matrices, and enhancer sequences  
20 discovered using the ChIP-seq method.

21 **Methods.** The Naïve-Bayes algorithm was used to classify literature and identify all PubMed  
22 abstracts on cardiovascular development. The natural language learning tool GNAT was then  
23 used to identify corresponding gene names embedded within these abstracts. Local Perl scripts  
24 were used to integrate and dump data from public databases into the MariaDB management  
25 system (MySQL). In-house R scripts were written to analyze and visualize the results.

26 **Results.** Known cardiovascular TFs from humans and human homologs from fly, *Ciona*,  
27 zebrafish, frog, chicken, mouse and human were identified and deposited in the database.  
28 Position weight matrices from Jaspar, hPDI, and UniPROBE databases were deposited in the  
29 database and can be retrieved using their corresponding TF names. Gene enhancer regions from  
30 various sources of ChIP-seq data were deposited into the database and were able to be  
31 visualized by graphical output. Besides biocuration, mouse homologs of the 81 core cardiac TFs  
32 were selected using a machine learning method and then by intersecting four independent data  
33 sources: RNA profiling, expert annotation, PubMed abstract and phenotype.

34 **Discussion.** The CardioTF database can be used as a portal to construct transcriptional network  
35 of cardiac development.

36 **Availability**                    **and**                    **Implementation.**                    Database                    URL:  
37 <http://www.cardiosignal.org/database/cardiotf.html>

## 38 **Introduction**

39 Heart disease is a leading cause of morbidity and mortality in both infants and adults [1,2].  
40 Insights into the cause of congenital heart diseases (CHDs) has led to the identification of  
41 mutations in essential cardiac transcription factors (TFs) [3]. At the opposite end of the temporal  
42 spectrum, adult cardiac disease can be traced to variants of gene regulatory sequences [4]. Thus,  
43 knowledge of TFs, their downstream targets, and regulatory genomic sequences involved in the  
44 heart transcriptional regulatory network will enhance our understanding of heart disease.

45 Although the vast amounts of data generated by high-throughput technologies is archived in  
46 databases such as ArrayExpress or GEO of NCBI [5,6], they do not contain cohesive knowledge  
47 and lack expert annotation. In addition, the field of cardiac development has experienced  
48 accelerated growth that can be attributed to the use of various animal models. However, to our  
49 present understanding, few efforts have been made to create a database which collects cardiac  
50 transcriptional information across species, thereby limiting the benefits from an evolutionary  
51 perspective to study heart development.

52 At present, two branching efforts have been made to archive and analyze the data. One is to  
53 construct small scale databases, like BloodChIP or CistromeMap, which are dedicated to  
54 collecting specific types of data [7, 8]. The other approach is to establish a number of consortia,  
55 like ENCODE, modENCODE, and Epigenomics Roadmap, which are created to generate huge  
56 amounts of raw data and archive them [9-11]. In addition to these projects, analysis and  
57 visualization software are valuable resources that lead to deeper understanding of the data, and  
58 facilitate the generation of novel hypotheses. Central databases, like Ensembl and UCSC also  
59 have search functions which allow browsing of the results generated by the consortia mentioned

60 above [12, 13]. However, there are currently few databases committed exclusively to  
61 cardiovascular development [14]. This prompted us to combine information about TFs, position  
62 weight matrices (PWMs), and ChIP-seq results and create a one-stop site for information on  
63 cardiovascular development, thus facilitating systems biology studies in transcriptional network  
64 regulation

65 CardioTF was therefore constructed to capture all transcriptional information relating to  
66 cardiovascular development. As a biocuration project, it documents TFs, PWM files and  
67 enhancers across species, including fly, *Ciona*, fish, frog, chicken, mouse and human. It also  
68 implements a search engine to query this information on the fly. In addition to the data-mining  
69 effort, core cardiac TFs are identified using a machine learning approach, which could be used  
70 as a roadmap to generate gene regulatory network of heart development.

## 71 **Materials & Methods**

### 72 **Comprehensive collection and annotation of cardiac TFs**

73 Cardiac TFs were previously defined as regulators of cardiac gene expression, which can  
74 impact the process of heart development, particularly the initiation and maintenance of the  
75 myocardium. By definition, cardiac TFs themselves should be involved in the steps of  
76 specification, determination, patterning, and differentiation that will result in a heart fate. We  
77 recruited cardiac TFs involved in the development of the epicardium and endocardium, both of  
78 which, along with the myocardium, should be from the *Mesp1*-expressing cell lineage [15].  
79 *Mesp1* is an indicator of pre-cardiac mesoderm, which, to our knowledge, is the earliest marker  
80 that labels all cardiac lineages. In *Ciona*, the *Mesp* enhancer specifies heart precursor cells [16].  
81 In the initial screen, we identified human TFs from previously published annotations [17]. We

82 used these human TFs, excluding human-specific TFs, as a reference to search for their  
83 homologs in other species, including fly, *Ciona*, zebrafish, chicken, frog and mouse [18].  
84 Human-specific TFs are defined as genes which have no homologs in the mouse genome. The  
85 NCBI HomoloGene database [19] was used as a reference to assess homologs between human  
86 and mouse/zebrafish. Other specie homologs of human were retrieved from their central  
87 databases, namely, FlyBase (Fly), BirdBase (Chicken), Aniseed (*Ciona*) and Xenbase (Frog)  
88 [20-23]. We also documented the expression status for mouse TFs from four independent  
89 sources which included annotation from the Cardiovascular Gene Ontology Annotation  
90 Initiative [24], Mouse Genome Database (MGI) genes with cardiovascular phenotype [25],  
91 PubMed abstract parsing results and RNA expression profiling results.

#### 92 **TFs from PubMed abstract parsing**

93 The Weinstein Cardiovascular Conference provides a platform for talks and posters on all  
94 aspects of heart development and congenital heart disease. We used Weinstein meeting abstracts  
95 from 2010 to 2013 (954 abstracts) as the positive group. We assumed that Weinstein-like  
96 abstracts deposited in PubMed are all from the cardiovascular community and focus on  
97 cardiovascular development. Abstracts from the negative group were from non-heart related  
98 journals, which were manually selected from the PMC Open Access Subset at NCBI. To choose  
99 the negative control journals, the following criteria were set. First, well-known cardiovascular  
100 journals were excluded, such as ‘Circulation’ and ‘Circulation Research’. Second, journals  
101 without key words ‘heart’ or ‘cardiac’ or ‘cardiovascular’ in their title were selected. Third,  
102 journals which are dedicated to the study of other organs or diseases, for example, ‘Neuron’ or  
103 ‘Cancer’ were selected. Fourth, other journal which are unlikely to publish articles about  
104 cardiovascular development and related topics, such as journals about plants or viruses, were

105 selected. Journals in the negative group contained research from across kingdoms and topics  
106 obviously in other fields, such as ‘Sleep\_Disord’ or ‘Toxicology’. All journal names in the  
107 negative group were saved in a file and uploaded onto the cardioTF server  
108 (negative\_set.journal.txt). The negative group includes 57080 abstracts. We split the data  
109 (positive and negative groups) into a training (80%) and test (20%) set. The Naïve-Bayes  
110 module from The Comprehensive Perl Archive Network (CPAN) was used with a local Perl  
111 script to classify Weinstein-like abstracts. We used the training set and adopted the 5×2 cross-  
112 validation proposed by Dietterich [26] to train and validate the data. The parameter ( the cutoff  
113 to decide whether an abstract is a true Weinstein-like) was selected based on average predictive  
114 performance which led to an accuracy (ACC) of 0.99. A wrapper function was implemented to  
115 parse the abstract and calculate the word frequencies. This function called two Perl modules  
116 (Lingua::EN::Splitter and Lingua::EN::StopWords) to extract words and perform text analysis.  
117 In doing so, the word frequency was forwarded to the algorithm as the only feature. A test set  
118 using optimized parameter was used to assess the algorithm’s final performance. All publication  
119 abstracts from 2008 to 2013 were downloaded to the local environment and analyzed by the  
120 algorithm. We targeted journals which had at least 6 publications classified as Weinstein-like  
121 abstracts in the six-year period (annual publication rate is  $\geq 1$ ). Then all abstracts from the  
122 targeted journal were downloaded. All abstracts were downloaded from journals that met the  
123 criteria. The selected abstracts were then processed by GNAT [27] using its default script  
124 (test100.sh) to recognize the gene name in mouse. PMID was recorded when the gene name  
125 matched with the name in the curated mouse TF set.

## 126 **RNA expression profiling data procession**

127 Affymetrix data (GSE1479) were processed by R using the MAS5 algorithm which provides a  
128 present call for each gene (see the script ExtractAffy.R at the Github) Gene expression status  
129 was defined as “on” if the gene was expressed in any microarray at selected developmental  
130 stages and had a present call. RNA-seq data (GSE47950 and GSE29184) were re-analyzed  
131 using the recommended protocol [28]. Any gene with an FPKM value greater than 1 was  
132 defined as expressed.

### 133 **Depositing PWM files**

134 The gene symbol was used as the unique identifier to link the original database ID to our local  
135 database primary key. A local Perl script was written to change their format to TRANSFAC  
136 style, which was used by our in-house CardioSignalScan program. PWM files were collected  
137 from Jaspar, UniPROBE and hPDI databases [29-31]. Users can retrieve their annotations by  
138 directing them to the respective database.

### 139 **Orthologs of TFs from model systems**

140 NCBI has its own collection of all ortholog gene collections using its unpublished algorithm.  
141 TFs from mouse, human and zebrafish are annotated by NCBI [32]. Frog, chicken and *Ciona*  
142 TF homolog annotations were downloaded from their central databases including Xenbase,  
143 BirdBase and ANISEED. Fly TFs, which have their counterparts in the human proteome, were  
144 annotated by the Inparanoid system [33]. Each TF collected in the database was assigned one  
145 treeID on the basis of its human counterpart. The treeID is equivalent to a TF family by the  
146 recommendation of TFClass.

### 147 **Enhancer curation: TF-ChIP and Histone-ChIP data processing**

148 Enhancer regions were defined by ChIP-seq signals. Peak calling was performed using the  
149 recommended pipeline [34]. In brief, sequencing reads were aligned to the mm10/hg19

150 reference genome using Bowtie/Bowtie2. mm9/hg9 represent the genome build assigned by  
151 UCSC. Index files for mm9/hg19 were downloaded from the iGenome project.

152 Bowtie call

153 bowtie -m 2 -S -q -p 8

154 Peak calling was performed using the MACS peak calling algorithm [35].

155 MACS call

156 macs14 -t ERR231646.bam -c ERR231653.bam -g mm -n sham\_Anti\_H3k9ac

157 Detailed annotation of the peaks was performed by HOMER using annotatePeaks.pl package  
158 [36].

### 159 **Recognition of transcription factor binding sites (TFBSs) in enhancer**

160 CardioSignalScan was previously implemented to identify the transcription factor binding sites.  
161 However, this local program is brute-force solution which consumes computational time  
162 exponentially making unrealistic to scan sequences longer than 3000bp. This prompted us to  
163 choose MOODS [37] instead, which reduces the computational time proportionally to position  
164 weight matrix length. A wrapper module was written to calculate the threshold that gauges the  
165 match. The cutoff was empirically defined to be 0.75 (range from 0 to 1).

166 
$$\text{threshold} = \text{min\_log\_score} + (\text{max\_log\_score} - \text{min\_log\_score}) * \text{cutoff}$$

167 This step avoids using *P*-value to assess the significance of TFBS.

## 168 **Results**

### 169 **The database schema**

170 Our database uses MariaDB, a drop-in replacement for MySQL, as the database management  
171 system (DBMS). To address how information will be stored and how the elements will related to  
172 one another, we used the unified modeling language (UML) to describe the high-level database

173 model [38]. UML was originally developed as a graphical notation for describing software  
174 designs in an object-oriented style. It has been extended, and modified and is now a popular  
175 notation for describing database designs. Here we used UML instead of an entity/relationship  
176 diagram to design the relational database schema following modeling principles, such as  
177 faithfulness, avoiding redundancy, and simplicity counts [38] (Fig. 1). When possible, we used  
178 a composition distinguished by a line between two classes that ends in a solid diamond at one  
179 end. The diamond implies that the label at the end must be 1:1. For example, there is a  
180 composition from CardioTFmatrix to CardioTFCenter, which means that every matrix  
181 annotation row (PWM related information) belongs to exactly one row in CardioTFCenter (one  
182 type of TF may have more PWM records in a CardioTFmatrix table). A 1:1 label at the  
183 CardioTFCenter end is implied by a solid black diamond.

#### 184 **Web Interface and search engine**

185 CardioTF is a Perl website implemented using only Perl language to dynamically display the  
186 graphical output while querying the database in the backend (Fig. 2). To aid cardiovascular  
187 biologists, a search engine was created to allow users to tackle the following problems: (1)  
188 identify homology information for the queried TF across six species and linking to the  
189 corresponding central databases outside CardioTF; (2) identify PWM file union of three public  
190 databases regarding the queried TF; and (3) identify the enhancer regions revealed by ChIP-seq  
191 data of the queried gene. Thus, the database is able to perform the key functions required to  
192 construct a transcriptional network of heart development. **Cardiovascular TFs in the database**  
193 Wingender's annotation set [16] was used as a benchmark to recruit TFs across species. Human-  
194 specific TFs, defined as those with no orthologs in the mouse genome, were discarded because  
195 no model system could be used to verify their function *in vivo*. Therefore, 1200 mouse TFs were

196 collected. Other established animal models for cardiovascular development include fly, *Ciona*,  
197 zebrafish, frog, and chicken. TFs from these species were collected if they were homologs to the  
198 above mouse TFs. The distribution of TFs from different species is shown in Fig. 3. The  
199 expression status of mouse TFs was verified by four independent resources, namely RNA-seq  
200 data re-analysis [39], phenotype annotations from the MGI database [40], expert  
201 recommendation from the UK Cardiovascular Gene Annotation Initiative [24], and PubMed  
202 relevance from classification of Weinstein-like abstracts (see the subsequent section and Table  
203 S1).

#### 204 **Weinstein TFs from PubMed analysis**

205 We identified the journals that favored Weinstein-style papers, which were likely contain gene  
206 information one genes involved in cardiovascular development. As expected, after using a  
207 machine learning method, the journals we identified were among the 30 journals most relevant  
208 to developmental biology. Two of the journals (*Circ. Res.* and *J Mol Cell Cardiol.*) obviously  
209 publish research specifically in the area of heart system. (Table S1: CardioJournalDistribution).  
210 If normalized and ranked by publication rate, the above conclusion still holds true although two  
211 different heart journals (*Eur J Echocardiogr*, *Heart Rhythm*) are in the top 30 list (Table S1:  
212 CardioJournalDistribution\_norm) in this case. We then used GNAT, a tool that recognizes gene  
213 names in the literature, to recover all TFs mentioned in Weinstein-style abstracts because we  
214 assumed that these TFs are studied by researchers in the cardiovascular community (Fig. 4, Fig.  
215 S1 and Table S1).

#### 216 **PWM files collected in database**

217 Public databases for PWM files include UniPROBE, Jaspar, and hPDI, and they provide PWM  
218 files for TFs. Jaspar PWM files are curated from the published literatures whereas the other two

219 databases generate PWM files from experiments. Our database integrates these three sources,  
220 and the TF PWM can be queried on the basis of the TF name. Search results directly link to the  
221 original database through the PWM raw database key. The CardioTFmatrix class contains 904  
222 records, and these PWM files can be recognized by our local CardioSignalScan program to  
223 search for the motifs in genomics regions.

#### 224 **Core cardiovascular TFs**

225 The 1200 mouse TFs were included in the cardiac TF dataset as the entry point to initiate deep  
226 annotation. To define the core dataset of cardiac TFs, we intersected four independent sources  
227 of cardiovascular TF collections. Inclusion of these 81 TFs is supported by their expression  
228 status, phenotype annotation, expert recommendation and PubMed relevance (Table S1). We  
229 also performed DAVID functional analysis, and found that TFs are particularly enriched in  
230 cardiac muscle differentiation (Fig. S2). Is this set of gene the minimum requirement for cardiac  
231 development? Indeed, these four sources of supporting evidence indicate that these TFs genes  
232 play a key role in heart development. We wanted to determine if these TFs display specific  
233 expression patterns in heart development. A heatmap was generated using seven RNA-seq data  
234 sets of various development stages. This heatmap did not reveal any specific patterns (Fig. S3).  
235 In adult tissues, these TFs did not exhibit enriched expression in the adult heart. In case of TFs  
236 which are never expressed at any stage of heart development, no specific expression pattern was  
237 revealed by the boxplot assay (Fig. S4-S5).

#### 238 **Cardiovascular enhancers collected in this database version**

239 Few enhancers have already been verified by traditional biological experiments, for example, by  
240 using transgenic expression of isolated DNA fragment *in vivo* to analyze temporal-spatial  
241 patterns. Therefore, the ChIP-seq method provides a high-throughput approach to delineate

242 enhancer regions at the genome scale. A standard protocol was used to identify genome-wide  
243 locations of transcription/chromatin factor binding sites or histone modification sites from  
244 ChIP-seq data (Fig. S6). The present database houses 511,893 enhancer records, covering  
245 different developmental stages of the heart. The searching result of a single enhancer also  
246 provide user interface to scan the TFBSs. The binding matrices provided in the list are came  
247 from the core cardiac TFs. The output in the flat text format contains the matrix key of the TFs  
248 in the database and list of hits in the order.

## 249 **Discussion**

250 We identified 5442 TFs from six species, and integrate 904 PWM files from three PWM  
251 databases. We also collected 511,893 peak fragments for further analysis. The on-the-fly  
252 searching tool was implemented to match the core cardiac TFBSs in the specified enhancer  
253 sequence. Our database provides a framework where users can query homology information  
254 of various TFs across those species and PWM information corresponding to TFs and enhancers  
255 from high-throughput ChIP-seq data. Previously, we constructed the CardioSignal database  
256 which collates cardiac factors driving genes expressed in a tissue-specific or quantitative  
257 manner. Most enhancers archived in the database are expressed specifically in the myocardium.  
258 CardioTF is a complementary database that accumulates cardiac TFs expressed in the Mesp-1  
259 lineage, thus including TFs expressed in the epicardium, endocardium and myocardium. This  
260 information including PWMs can be used not only to find the features of enhancers in a  
261 machine learning approach (such the left-right patterning of the heart) but also to reconstruct  
262 regulatory network in systems biology.

263 We defined a core set of TFs using by four independent sources, namely, RNA profiling, expert  
264 curation, PubMed abstract parsing and phenotype annotation to support their roles in

265 cardiovascular development. In analyze PubMed abstracts, a machine learning approach was  
266 used to classify the Weinstein-style abstracts and then pick out the TFs embedded in those  
267 abstracts. Our text analysis, whether normalized by publication number or not, had a tendency  
268 of identifying journals favored by heart developmental biologists and journals that specialized  
269 in developmental biology (Table S1). The overlap of core cardiac TFs identified by these  
270 sources provide a roadmap for systems biology to construct transcriptional network of heart  
271 development. Current approaches by the Sperling group or the Pu group only report three to  
272 four TFs based on ChIP-seq data [41,42]. Similar approach by other genome biologists who  
273 tried to find cardiac enhancers on a genome scale have been reviewed elsewhere [43]. However,  
274 the information generated from these studies is well below our knowledge of these core  
275 cardiovascular TFs, which have multiple sources supporting their role in cardiovascular  
276 development.

277 We archived 1200 mouse TFs and wanted to determine at what stage of heart development they  
278 were expressed. Our preliminary analysis indicates that approximately 200 TFs have no  
279 evidence of their expression pattern, phenotype, expert recommendation, and PubMed abstracts.  
280 Whether these TF genes are expressed or play roles in heart disease requires further data  
281 analysis.

282 The database still lacks cell lineage-based expression profiling data, which will quantify the  
283 expression level of various TFs and thus construct a 4-D dynamic expression pattern *in vivo*.  
284 This information could be combined with cell lineage-based ChIP-seq data to create a super-  
285 resolution of enhancer tomography.

## 286 **Conclusions**

287 Modern translational medicine rests upon the progressive study of pathways and principles from  
288 model organisms such as yeast, fly, fish, and mice to clinical studies in humans. Therefore, we  
289 recruited TFs from six model organisms which are established models for research on  
290 cardiovascular development. The identification and collation of these well-annotated homologs  
291 in different animals enable investigators to interrogate more fundamental problems in heart  
292 development.

293 We hope that in the near future, single-cell sequencing data may provide comprehensive gene  
294 expression information with detailed temporal-spatial resolution, thereby providing insight into  
295 the enhancer network that contribute to the gene specific expression patterns. CardioTF is the  
296 initial step into the construction of a comprehensive transcriptional network reconstruction.

### 297 **Supplemental Information**

298 Additional file 1: Four independent sources of cardiac TF lists and preferred journals for  
299 Weinstein-like papers. Format: XLSX. Sheet 1 (CardioJournalDistribution) includes the top 100  
300 journals that favor Weinstein-like papers. Sheet 2 (NotHeartTFs) includes all TF names that  
301 have never appeared in the four sources of supporting evidence. Sheet 3 (MGI\_TFs) is the TF  
302 names from MGI database. Sheet 4 (PubMed\_TFs) lists TF names from Weinstein-like PubMed  
303 abstracts. Sheet 5 (Cardio\_GO\_UK\_TFs) includes TF names from UK Cardio-GO project.  
304 Sheet 5 (Cardio\_lineage\_TFs) includes TF names from RNA-seq or microarray analysis. Sheet  
305 6 (core\_TFs) defines the core 81 TFs known to be involved in heart development. Sheet 7  
306 (CardioJournalDistribution\_norm) includes journals ranked by publication rate (normalized  
307 version)

308 Additional file 2: Format: PDF. Figure S1. Confusion matrix and precision-recall curve. By  
309 convention, the class label of the minority class is positive (Weinstein abstracts), while the class

310 label of the majority class is negative (non-Weinstein-like abstracts). (A) The confusion matrix  
311 for a two-class problem. The first column shows the actual class label of the examples, and the  
312 first row presents their predicted class label. In the matrix, TP shows the true positive samples,  
313 FP shows the false positive samples, TN shows the true negative samples, and FN shows the  
314 false negative samples. (B) The precision-recall curve.

315 Additional file 3: Format: PDF. Figure S2. DAVID analysis of 81 core TFs. The results indicate  
316 that these TFs are truly associated with cardiac function by GO term enrichment analysis.

317 Additional file 4: Format: PDF. Figure S3. Heatmap of the 81 core cardiac transcriptional  
318 factors at the different stages of heart development.

319 Additional file 5: Format: PDF. Figure S4. RNA-seq expression patterns for the 81 core TFs  
320 across 13 adult tissues presented by the boxplot.

321 Additional file 6: Format: PDF. Figure S5. The expression profile across 13 adult tissues of TFs  
322 which are never expressed in the heart, as determined by our four sources supporting  
323 evidence.

324 Additional file 7: Format: PDF. Figure S6. The work flow of parsing ChIP-seq data and  
325 dumping it into the MySQL database.

### 326 **Availability and requirements**

327 CardioTF database is freely available on the web at  
328 <http://www.cardiosignal.org/database/cardiottf.html>.

### 329 **List of abbreviations**

330 TF: transcription factors; PWM: position weight matrix; UML: unified modeling language;

### 331 **Acknowledgements**

332 The author is grateful to Prof. Rutai Hui, Prof. Weinian Shou, Dr. Tingting Li and Dr. Jianxin  
333 Chen for their helpful comments. Special thanks to Prof. Paul Krieg (University of Arizona),  
334 Matija Brozovic (ANISEED), Prof. Carl J. Schmidt (BirdBase) and Zichao Sang (CardioSignal)  
335 for their comments or suggestions on data curation. I would like to thank the two anonymous  
336 reviewers and Dr. Rackham for their valuable comments and suggestions to improve the quality  
337 of the paper. I would also like to thank Dr. Rosannah C. Cameron at the Albert Einstein College  
338 of Medicine for her assistance editing the manuscript.

### 339 **References:**

- 340 1. van der Linde D, Konings EE, Slager MA, Witsenburg M, Helbing WA, Takkenberg JJ,  
341 Roos-Hesselink JW...: **Birth prevalence of congenital heart disease worldwide: a systematic**  
342 **review and meta-analysis.** *J Am Coll Cardiol*, 2011, **58**(21):2241-7.
- 343 2. Celermajer DS, Chow CK, Marijon E, Anstey NM, Woo KS: **Cardiovascular disease in the**  
344 **developing world: prevalences, patterns, and the potential of early disease detection.** *J Am*  
345 *Coll Cardiol*, 2012, **60**(14):1207-16.
- 346 3. McCulley DJ, Black BL: **Transcription factor pathways and congenital heart disease.**  
347 *Curr Top Dev Biol*, 2012, **100**:253-77.
- 348 4. Smith JG, Newton-Cheh C: **Genome-wide association studies of late-onset cardiovascular**  
349 **disease.** *J Mol Cell Cardiol*, 2015, **83**:131-41.
- 350 5. Parkinson H, Sarkans U, Kolesnikov N, Abeygunawardena N, Burdett T, Dylag M, Emam I,  
351 Farne A, Hastings E, Holloway E, Kurbatova N, Lukk M, Malone J, Mani R, Pilicheva E,  
352 Rustici G, Sharma A, Williams E, Adamusiak T, Brandizi M, Sklyar N, Brazma A...:  
353 **ArrayExpress update--an archive of microarray and high-throughput sequencing-based**  
354 **functional genomics experiments.** *Nucleic Acids Res*, 2011,**39**(Database issue):D1002-4.

- 355 6. Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, Marshall KA,  
356 Phillippy KH, Sherman PM, Holko M, Yefanov A, Lee H, Zhang N, Robertson CL, Serova N,  
357 Davis S, Soboleva A.: **NCBI GEO: archive for functional genomics data sets--update.**  
358 *Nucleic Acids Res*, 2013, **41** (Database issue):D991-5.
- 359 7. Chacon D, Beck D, Perera D, Wong JW, Pimanda JE: **BloodChIP: a database of**  
360 **comparative genome-wide transcription factor binding profiles in human blood cells.**  
361 *Nucleic Acids Res*, 2014, **42**(Database issue):D172-7.
- 362 8. Qin B, Zhou M, Ge Y, Taing L, Liu T, Wang Q, Wang S, Chen J, Shen L, Duan X, Hu S, Li  
363 W, Long H, Zhang Y, Liu XS.: **CistromeMap: a knowledgebase and web server for ChIP-**  
364 **Seq and DNase-Seq studies in mouse and human.** *Bioinformatics*, 2012, **28** (10):1411-12.
- 365 9. Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, Aken BL,  
366 Barrell D, Zadissa A, Searle S, Barnes I, Bignell A, Boychenko V, Hunt T, Kay M, Mukherjee  
367 G, Rajan J, Despacio-Reyes G, Saunders G, Steward C, Harte R, Lin M, Howald C, Tanzer A,  
368 Derrien T, Chrast J, Walters N, Balasubramanian S, Pei B, Tress M, Rodriguez JM, Ezkurdia I,  
369 van Baren J, Brent M, Haussler D, Kellis M, Valencia A, Reymond A, Gerstein M, Guigó R,  
370 Hubbard TJ. **GENCODE: the reference human genome annotation for The ENCODE**  
371 **Project.** *Genome Res.*, 2012, **22**(9):1760-74.
- 372 10. Celniker SE, Dillon LA, Gerstein MB, Gunsalus KC, Henikoff S, Karpen GH, Kellis M, Lai  
373 EC, Lieb JD, MacAlpine DM, Micklem G, Piano F, Snyder M, Stein L, White KP, Waterston  
374 RH; modENCODE Consortium. Unlocking the secrets of the genome. *Nature.*, 2009,  
375 459(7249):927-30
- 376 11. Romanoski CE, Glass CK, Stunnenberg HG, Wilson L, Almouzni G. Epigenomics: Roadmap  
377 for regulation. *Nature.* 2015, 518(7539):314-6.

- 378 12. Mangan ME, Williams JM, Kuhn RM, Lathe WC 3rd. The UCSC Genome Browser: **What**  
379 **Every Molecular Biologist Should Know**. *Curr Protoc Mol Biol.*, 2014, 1;107:19.9.1-19.9.36.
- 380 13. Flicek P, Ahmed I, Amode MR, Barrell D, Beal K, Brent S, Carvalho-Silva D, Clapham P,  
381 Coates G, Fairley S, Fitzgerald S, Gil L, García-Girón C, Gordon L, Hourlier T, Hunt S,  
382 Juettemann T, Kähäri AK, Keenan S, Komorowska M, Kulesha E, Longden I, Maurel T,  
383 McLaren WM, Muffato M, Nag R, Overduin B, Pignatelli M, Pritchard B, Pritchard E, Riat HS,  
384 Ritchie GR, Ruffier M, Schuster M, Sheppard D, Sobral D, Taylor K, Thormann A, Trevanion S,  
385 White S, Wilder SP, Aken BL, Birney E, Cunningham F, Dunham I, Harrow J, Herrero J,  
386 Hubbard TJ, Johnson N, Kinsella R, Parker A, Spudich G, Yates A, Zadissa A, Searle SM.  
387 **Ensembl 2013**. *Nucleic Acids Res.*, 2013, 41(Database issue):D48-55.
- 388 14. Djordjevic D, Yang A, Zadoorian A, Rungrugeechooen K, Ho JW. How difficult is  
389 inference of mammalian causal gene regulatory networks? *PLoS One.*, 2014, 9(11):e111661.
- 390 15. Bondue A, Blanpain C: **Mesp1: a key regulator of cardiovascular lineage commitment**.  
391 *Circ Res*, 2010, **107**(12):1414-27.
- 392 16. Satou Y, Imai KS, Satoh N. The ascidian Mesp gene specifies heart precursor cells.  
393 *Development*. 2004, 131(11):2533-41.
- 394 17. Wingender E, Schoeps T, Dönitz J: **TFClass: an expandable hierarchical classification of**  
395 **human transcription factors**. *Nucleic Acids Res*, 2013, **41** (Database issue):D165-70.
- 396 18. Hutson MR, Kirby ML: **Model systems for the study of heart development and disease**.  
397 *Semin Cell Dev Biol*, 2007, **18**(1):1-2.
- 398 19. NCBI Resource Coordinators. **Database resources of the National Center for**  
399 **Biotechnology Information**. *Nucleic Acids Res*. 2015, 43(Database issue):D6-17.

- 400 20. Attrill H, Falls K, Goodman JL, Millburn GH, Antonazzo G, Rey AJ, Marygold SJ; FlyBase  
401 consortium. **FlyBase: establishing a Gene Group resource for *Drosophila melanogaster*.**  
402 *Nucleic Acids Res.* 2016, 44(D1):D786-92.
- 403 21. Karpinka JB, Fortriede JD, Burns KA, James-Zorn C, Ponferrada VG, Lee J, Karimi K,  
404 Zorn AM, Vize PD.: **Xenbase, the *Xenopus* model organism database; new virtualized**  
405 **system, data types and genomes.** *Nucleic Acids Res*, 2015, **43** (Database issue):D756-63.
- 406 22. Schmidt CJ, Romanov M, Ryder O, Magrini V, Hickenbotham M, Glasscock J, McGrath  
407 S, Mardis E, Stein LD.: **Gallus GBrowse: a unified genomic database for the chicken.**  
408 *Nucleic Acids Res*, 2008, **36** (Database issue):D719-23.
- 409 23. Tassy O, Dauga D, Daian F, Sobral D, Robin F, Khoueiry P, Salgado D, Fox V, Caillol D,  
410 Schiappa R, Laporte B, Rios A, Luxardi G, Kusakabe T, Joly JS, Darras S, Christiaen L,  
411 Contensin M, Auger H, Lamy C, Hudson C, Rothbacher U, Gilchrist MJ, Makabe KW, Hotta K,  
412 Fujiwara S, Satoh N, Satou Y, Lemaire P.: **The ANISEED database: digital representation,**  
413 **formalization, and elucidation of a chordate developmental program.** *Genome Res*, 2010,  
414 **20**(10):1459-68.
- 415 24. Khodiyar VK, Hill DP, Howe D, Berardini TZ, Tweedie S, Talmud PJ, Breckenridge R,  
416 Bhattarcharya S, Riley P, Scambler P, Lovering RC.: **The representation of heart**  
417 **development in the gene ontology.** *Dev Biol*, 2011, **354**(1):9-17.
- 418 25. Blake JA, Bult CJ, Eppig JT, Kadin JA, Richardson JE: Mouse Genome Database Group.  
419 **The Mouse Genome Database: integration of and access to knowledge about the**  
420 **laboratory mouse.** *Nucleic Acids Res*, 2014, **42** (Database issue):D810-17.
- 421 26. Dietterich TG: **Approximate statistical tests for comparing supervised classification**  
422 **learning algorithms.** *Neural Computation*, 1998, **10**:1895-923.

- 423 27. Hakenberg J, Plake C, Leaman R, Schroeder M, Gonzalez G: **Inter-species normalization**  
424 **of gene mentions with GNAT.** *Bioinformatics*, 2008, **24**(16):126-32.
- 425 28. Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, Pimentel H, Salzberg SL, Rinn  
426 JL, Pachter L.: **Differential gene and transcript expression analysis of RNA-seq**  
427 **experiments with TopHat and Cufflinks.** *Nat Protoc*, 2012, **7**(3):562-78.
- 428 29. Mathelier A, Zhao X, Zhang AW, Parcy F, Worsley-Hunt R, Arenillas DJ, Buchman S,  
429 Chen CY, Chou A, Ienasescu H, Lim J, Shyr C, Tan G, Zhou M, Lenhard B, Sandelin A,  
430 Wasserman WW.: **JASPAR 2014: an extensively expanded and updated open-access**  
431 **database of transcription factor binding profiles.** *Nucleic Acids Res*, 2014, **42** (Database  
432 issue):D142-7.
- 433 30. Hume MA, Barrera LA, Gisselbrecht SS, Bulyk ML: **UniPROBE, update 2015: new tools**  
434 **and content for the online database of protein-binding microarray data on protein-DNA**  
435 **interactions.** *Nucleic Acids Res*, 2015, **43**(Database issue):D117-22.
- 436 31. Xie Z, Hu S, Blackshaw S, Zhu H, Qian J: **hPDI: a database of experimental human**  
437 **protein-DNA interactions.** *Bioinformatics*, 2010, **26**(2):287-9.
- 438 32. Altenhoff AM, Dessimoz C: **Phylogenetic and functional assessment of orthologs**  
439 **inference projects and methods.** *PLoS Comput Biol*, 2009, **5**(1): e1000262.
- 440 34. Sonnhammer EL, Östlund G: **InParanoid 8: orthology analysis between 273 proteomes,**  
441 **mostly eukaryotic.** *Nucleic Acids Res*, 2015, **43** (Database issue):D234-39.
- 442 34. Bailey T, Krajewski P, Ladunga I, Lefebvre C, Li Q, Liu T, Madrigal P, Taslim C, Zhang  
443 J.: **Practical guidelines for the comprehensive analysis of ChIP-seq data.** *PLoS Comput Biol*,  
444 2013, **9**(11):e1003326.

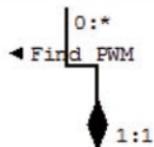
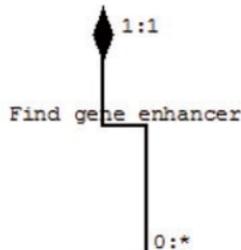
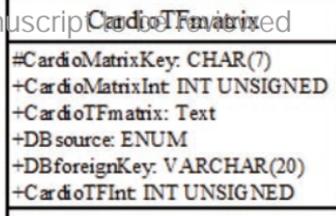
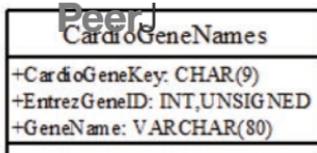
- 445 35. Feng J, Liu T, Qin B, Zhang Y, Liu XS: **Identifying ChIP-seq enrichment using MACS.**  
446 *Nat Protoc*,2012, **7**(9):1728-40.
- 447 36. Heinz S, Benner C, Spann N, Bertolino E, Lin YC, Laslo P, Cheng JX, Murre C, Singh H,  
448 Glass CK: **Simple combinations of lineage-determining transcription factors prime cis-**  
449 **regulatory elements required for macrophage and B cell identities.** *Mol Cell*,2010,  
450 **38**(4):576-89.
- 451 37. Korhonen J, Martinmäki P, Pizzi C, Rastas P, Ukkonen E. **MOODS: fast search for**  
452 **position weight matrix matches in DNA sequences.** *Bioinformatics.*, 2009, 25(23):3181-2.
- 453 38. Ullman JD, Widom J. A first course in database systems, 3rd. Beijing. Pearson Education  
454 Asia Limited & China Machine Press; 2008. p.140-147.
- 455 39. Shen Y, Yue F, McCleary DF, Ye Z, Edsall L, Kuan S, Wagner U, Dixon J, Lee L,  
456 Lobanenkov VV, Ren B. **A map of the cis-regulatory sequences in the mouse genome.**  
457 *Nature.*, 2012, 488(7409):116-20.
- 458 40. Wamstad JA, Alexander JM, Truty RM, Shrikumar A, Li F, Eilertson KE, Ding H, Wylie JN,  
459 Pico AR, Capra JA, Erwin G, Kattman SJ, Keller GM, Srivastava D, Levine SS, Pollard KS,  
460 Holloway AK, Boyer LA, Bruneau BG. **Dynamic and coordinated epigenetic regulation of**  
461 **developmental transitions in the cardiac lineage.** *Cell.* 2012, 151(1):206-20.
- 462 41. Schlesinger J, Schueler M, Grunert M, Fischer JJ, Zhang Q, Krueger T, Lange M, Tönjes M,  
463 Dunkel I, Sperling SR: **The cardiac transcription network modulated by Gata4, Mef2a,**  
464 **Nkx2.5, Srf, histone modifications, and microRNAs.** *PLoS Genet*, 2011, **7**(2):e1001313.
- 465 42. He A, Kong SW, Ma Q, Pu WT: **Co-occupancy by multiple cardiac transcription factors**  
466 **identifies transcriptional enhancers active in heart.** *Proc Natl Acad Sci U S A*,  
467 2011, **108**(14):5632-37.

- 468 43. Wamstad JA, Wang X, Demuren OO, Boyer LA. **Distal enhancers: new insights into heart**  
469 **development and disease.** *Trends Cell Biol.* ,2014, 24(5):294-302.

**Figure 1**(on next page)

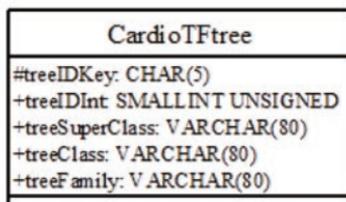
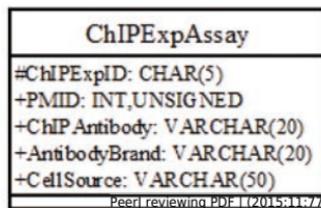
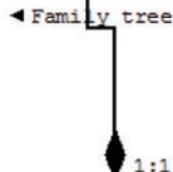
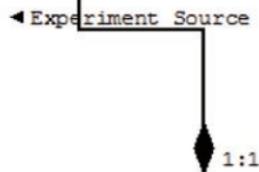
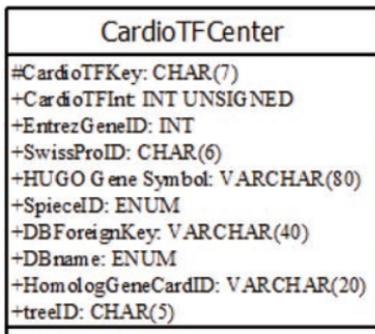
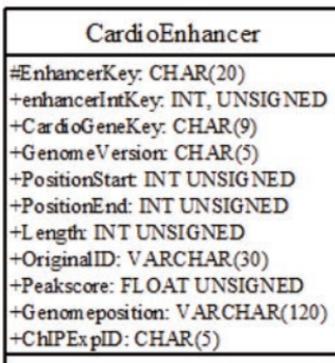
Unified modeling language diagram for the Cardio-TF database design.

The six graphical notation boxes represent of six major classes, namely CardioTFmatrix, CardioTFCenter, CardioTree, CardioGeneNames, CardioEnhancer, and ChIPExpAssay. These classes are analogous to entity/relationship sets. Each class has only two sections, one for the class name and one for the attributes. The attribute of each class is associated with the type used in MariaDB. The “#” in front of an attribute indicates that it’s visibility is “protected”, thereby making it a primary key. These classes faithfully represent the real world.



1:1

0:1



**Figure 2**(on next page)

The search engine and the web interface of the database.

(A) The search engine was implemented to perform three functions: querying TFs, their PWMs and gene enhancers (B) Web graphical output of Gata4 enhancers in mouse. Black lines indicate the enhancer region found by the ChIP-seq scanning program. (C) Query results for the GATA4 TF across species. TFs are listed and indexed according to their database identifiers.

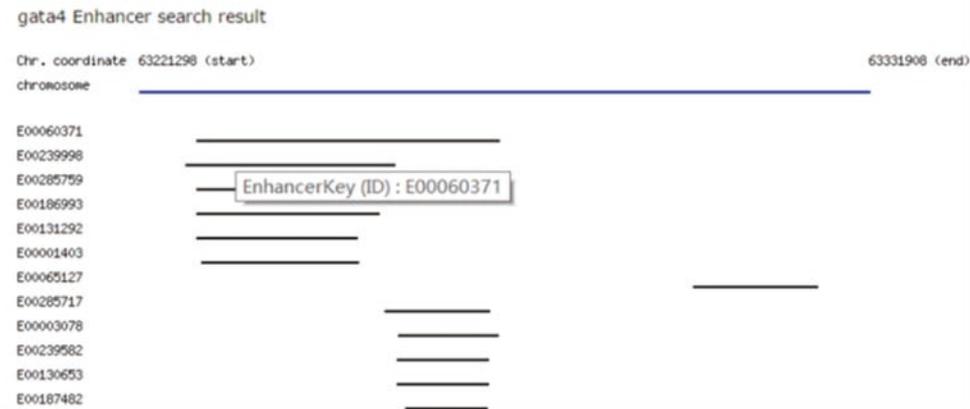
A

CardioTF Module: Cardiac Transcriptional Factor (TF) ▾

OPTIONS: Cardiac Gene Name ▾

QUERY: Gata4 submit

B



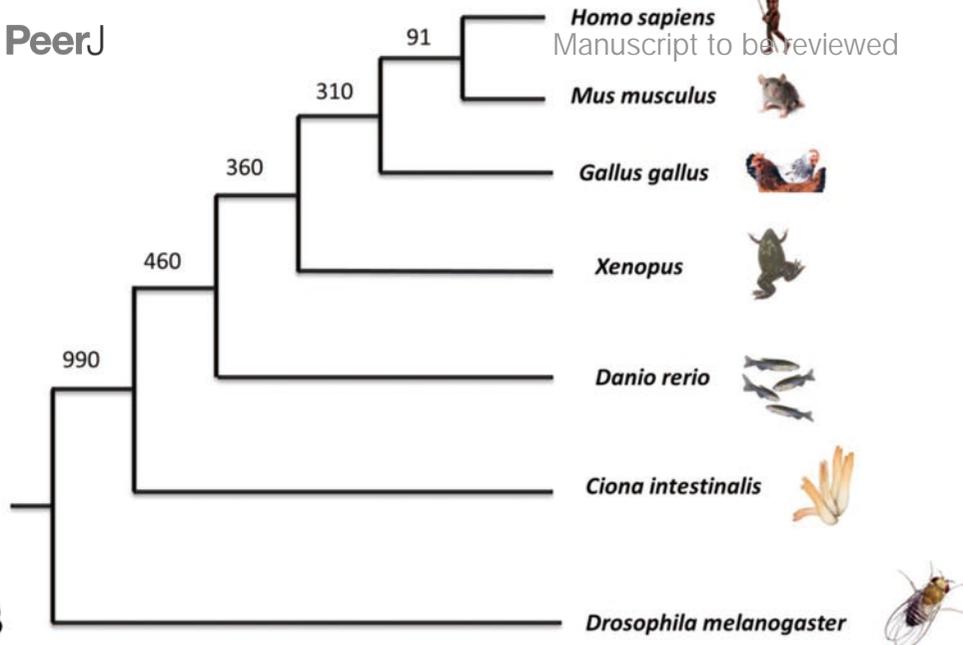
C

CardioTF Identifier	Gene Symbol	Species Name	EntrezGeneID	SwissProID	DB_link
TF02887	<i>gata4</i>	<i>Danio rerio</i>	30483	B8JKU1	ZFIN
TF01887	GATA4	<i>Homo sapiens</i>	2626	P43694	NCBI
TF04086	<i>gata4</i>	<i>Xenopus (Silurana) tropicalis</i>	549703	Unknown	XenBase
TF04714	GATA4	<i>Gallus gallus</i>	396392	Unknown	BirdBase
TF05400	PNR	<i>Drosophila melanogaster</i>	44849	P52168	FlyBase
TF01888	Gata4	<i>Mus musculus</i>	14463	Q08369	MGJ

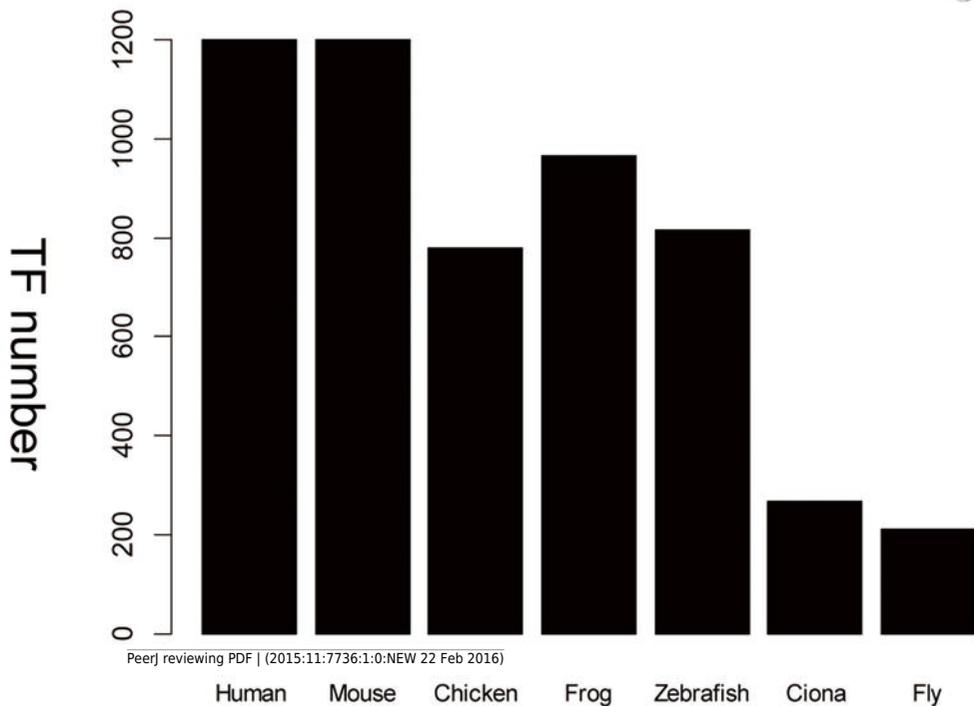
**Figure 3**(on next page)

TF distribution across species in the database.

(A) Phylogenetic tree showing the main animal models commonly used in heart development research and their evolutionary relationship. The divergence times in millions of years ago (Mya) are shown on the basis of multigene and multiprotein studies. Branch lengths are not proportional to time (B) Distribution TFs across six species. All TFs have homologs in their human counterparts. The unit of Y-axis is TF number.



B



**Figure 4**(on next page)

Machine learning protocol to select TFs described in Weinstein-like papers.

(A) The pipeline to select TF gene symbols from Weinstein PubMed abstracts. First, Naïve-Bayes module was used to select Weinstein-like papers from PubMed abstracts. Second, GNAT, a software that recognizes gene symbols, was used to identify all TF names from these Weinstein-like papers. (B) ROC curve and prediction performance judged by sensitivity, precision and F1 score.

A

Natural language  
processing Module:  
GNAT



Machine learning tool -  
PerlModule: NaiveBayes

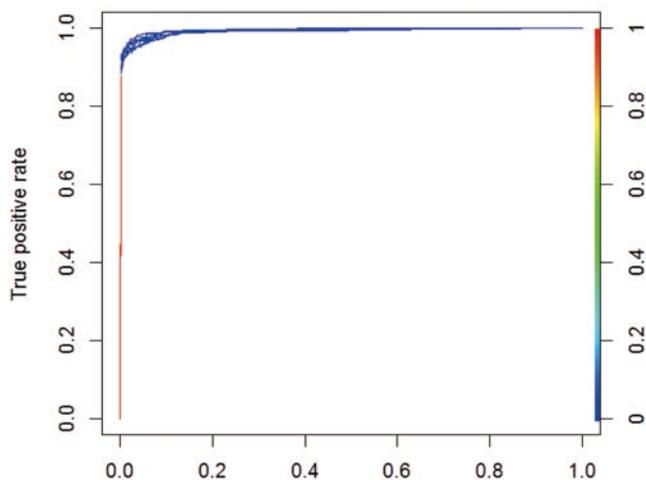


Pick up Weinstein style  
Cardiovascular PubMed  
abstracts



Gene symbol/Entrez GeneID

B



Sensitivity	0.9316
-------------	--------

Precision	0.9516
-----------	--------

F1 Score	0.9415
----------	--------