

Analysis and synthesis of high-amplitude *Cis*-elements in the mammalian circadian clock

Yuichi Kumaki^{*†‡}, Maki Ukai-Tadenuma^{*}, Ken-ichiro D. Uno[§], Junko Nishio[§], Koh-hei Masumoto^{*¶}, Mamoru Nagano[¶], Takashi Komori[†], Yasufumi Shigeyoshi[¶], John B. Hogenesch^{||}, and Hiroki R. Ueda^{*§**}

^{*}Laboratory for Systems Biology and [§]Functional Genomics Unit, Center for Developmental Biology, RIKEN, 2-2-3 Minatojima-Minamimachi, Chuo-ku, Kobe 650-0047, Japan; [†]INTEC Systems Institute, Inc., 1-3-3 Shinsuna, Koto-ku, Tokyo 136-0075, Japan; [¶]Department of Anatomy and Neurobiology, Kinki University School of Medicine, 377-2 Ohno-Higashi, Osaka-Sayama, Osaka 589-8511, Japan; and ^{||}Institute for Translational Medicine and Therapeutics and the Department of Pharmacology, University of Pennsylvania School of Medicine, 810 Biomedical Research Building II/III, 421 Curie Boulevard, Philadelphia, PA 19104-6160

Edited by Joseph S. Takahashi, Northwestern University, Evanston, IL, and approved August 15, 2008 (received for review March 16, 2008)

Mammalian circadian clocks consist of regulatory loops mediated by Clock/Bmal1-binding elements, DBP/E4BP4 binding elements, and RevErbA/ROR binding elements. As a step toward system-level understanding of the dynamic transcriptional regulation of the oscillator, we constructed and used a mammalian promoter/enhancer database (<http://promoter.cdb.riken.jp/>) with computational models of the Clock/Bmal1-binding elements, DBP/E4BP4 binding elements, and RevErbA/ROR binding elements to predict new targets of the clock and subsequently validated these targets at the level of the cell and organism. We further demonstrated the predictive nature of these models by generating and testing synthetic regulatory elements that do not occur in nature and showed that these elements produced high-amplitude circadian gene regulation. Biochemical experiments to characterize these synthetic elements revealed the importance of the affinity balance between transactivators and transrepressors in generating high-amplitude circadian transcriptional output. These results highlight the power of comparative genomics approaches for system-level identification and knowledge-based design of dynamic regulatory circuits.

comparative genomics | promoter and enhancer database | synthetic biology | systems biology | transcription

The rapidly expanding number of sequenced mammalian genomes (1–3), annotated and cloned full-length cDNAs (4–6), transcriptional starts sites (TSSs) (7–9) and transcription factor binding sites (TFBSs) (10–12) has provided new opportunities to unravel the control of dynamic transcriptional programs. Comparative genomics approaches applying these resources have been used to identify target genes of specific biological pathways. These efforts used consensus sequence searches (13, 14), positional weight matrices (15), hidden Markov models (HMMs) (16), and specifically tailored algorithms (17, 18) to define candidate response elements and target genes in raw genomic sequence. Additionally, post hoc analysis employing evolutionary conservation (15, 16, 18) together with positional information of TSSs (15) and/or translational start sites (16) has helped to further define candidate elements and genes and greatly expanded our knowledge of transcriptional output regulation.

The mammalian circadian clock is an ideal system to apply these tools as it consists of integrated transcriptional regulatory loops that direct output through at least three types of transcriptional regulatory elements, the Clock/Bmal1-binding elements (E-box) (CACGTG) (19–21), DBP/E4BP4 binding elements (D-box) (TTATG[T/C]AA) (21–23), and RevErbA/ROR binding elements (RRE) ([A/T]A[A/T]NT[A/G]GGTCA) (15, 21, 24, 25). Several groups, including our own, have shown that approximately 5–10% of mammalian genes display circadian expression in central and peripheral clock tissues (26). However, for the most part, the transcriptional regulation of these thousands of clock-controlled genes has remained uncharacterized. We and others have used comparative genomics approaches to

analyze E-box (21, 27, 28), D-box (21), and RRE (15, 21), highlighting the importance of both their core consensus and flanking sequences (15, 21, 27, 28) in circadian gene control. In this study, we further extend comparative genomics approaches toward a system-level understanding of the dynamic transcriptional regulations of the mammalian circadian clock.

Results and Discussion

Prediction of Direct Clock Targets Through Utilization of the Mammalian Promoter/Enhancer Database. To generate a resource that facilitates identification of clock-controlled genes, we constructed a mammalian promoter/enhancer database (<http://promoter.cdb.riken.jp/>) by integrating information sources such as conserved non-coding regions, TSSs and TFBSs [supporting information (SI) Fig. S1 and SI Appendix]. Although excellent and similar databases exist such as DBTSS (8), CisView (29) and ECRbase (30), none were tailored to specifically identify clock gene targets and having local control of the database facilitated manipulation of the underlying data (see also SI Appendix). We then developed a comparative genomics strategy employing this database and profile HMMs using the HMMER software package (31). Profile HMMs are powerful tools to extract the statistical properties of input sequences by representing multiple sequences as a transition probability matrix marching from one position to the neighboring position. HMMs were built and calibrated on known functional clock-controlled elements experimentally verified in our previous (15, 21) or current studies (Fig. S2 and Table S1), consisting of 12 E-boxes, 10 D-boxes and 15 RREs (Table S2). Profile HMM searches to identify new clock-controlled elements from conserved non-coding regions between human and mouse identified 1,108 E-boxes, 2,314 D-boxes, and 3,288 RREs candidate elements (see the circadian section of the mammalian promoter/enhancer database: <http://promoter.cdb.riken.jp/circadian.html> for element lists). To set appropriate reporting thresholds, we used the match scores of known functional clock-controlled elements (*Material and Methods*). Predicted clock-controlled elements exhibited an un-biased distribution of chromosomal position spread over whole mouse genome (Fig. 1A and <http://promoter.cdb.riken.jp/circadian.html>).

Author contributions: J.B.H. and H.R.U. designed research; Y.K., M.U., K.D.U., J.N., K.M., M.N., and Y.S. performed research; Y.K. and T.K. contributed new analytic tools; Y.K. analyzed data; and Y.K., J.B.H., and H.R.U. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Freely available online through the PNAS open access option.

[†]Present address: Laboratory for Mammalian Epigenetic Studies, Center for Developmental Biology, RIKEN, 2-2-3 Minatojima-Minamimachi, Chuo-ku, Kobe 650-0047, Japan.

^{**}To whom correspondence should be addressed. E-mail: uedah-ky@umin.ac.jp.

This article contains supporting information online at www.pnas.org/cgi/content/full/0802636105/DCSupplemental.

© 2008 by The National Academy of Sciences of the USA

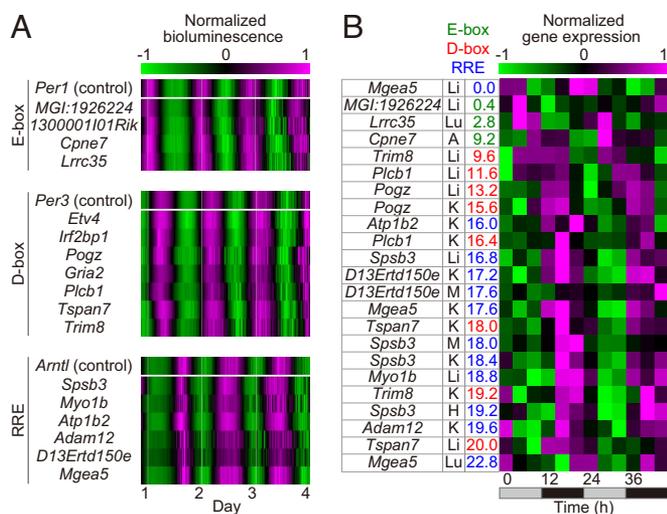


Fig. 2. Experimental validation of HMM-based predictions at cellular and organismal levels. (A) Circadian rhythms of bioluminescence from the predicted clock-controlled elements fused to the SV40 basic promoter driving a dLuc reporter in NIH 3T3 fibroblasts. Three known clock-controlled elements from clock genes (E-box of *Per1*, D-box of *Per3*, and RRE of *Arntl*) are used as positive controls. The bioluminescence data were detrended in baseline and amplitude and normalized so that their maximum, minimum, and average were set to 1, -1 , and 0, respectively. The colors in descending order from magenta to black to green represent the detrended bioluminescence. Columns represent time points and rows represent the predicted elements on the designated genes. (B) Circadian rhythms of temporal mRNA expression profiles of the predicted clock-controlled genes in mouse seven tissues ('A', 'B', 'H', 'K', 'Li', 'Lu' and 'M' for aorta, bone, heart, kidney, liver, lung, and muscle, respectively). An estimated peak time with color of type of predicted clock-controlled element (green, red, and blue for E-box, D-box, and RRE, respectively) is also indicated. The colors in descending order from magenta to black to green represent the normalized data (the average and standard deviation over 12-point time courses are 0.0 and 1.0, respectively). Columns represent time points, and rows represent the predicted clock-controlled genes in the designated tissues.

RRE, respectively (Fig. 2A and Fig. S3 A–D). The remaining sequences generated weak, low amplitude circadian transcriptional activity or were arrhythmic (Fig. S3 E–G). To supplement our observed 40% prediction success, we constructed 14 reporters containing conserved low-scoring E-boxes and found only one exhibited high-amplitude oscillations (Fig. S3H and Table S3). This result indicates that our observed 40% prediction success is suggestively higher than expected ($P = 0.075$, Fisher's exact test). Taken in sum, these results demonstrate utility of this approach in finding elements within structural genes that dictate rhythmic transcription.

If these *in vitro* validated 17 elements (4 E-boxes, 7 D-boxes, and 6 RREs) play a prominent role in gene regulation *in vivo*, we would predict that the endogenous transcripts for these genes would likely oscillate in a circadian fashion. To test this, we harvested mRNA from seven tissues (aorta, bone, heart, kidney, liver, lung, and muscle) isolated from mice entrained to a 12:12 light:dark cycle and then released to free run in constant darkness. Using quantitative PCR assays, we measured expression profiles from our predicted clock-controlled elements, and evaluated their rhythmicity using a statistical method based on analysis of variance (ANOVA) followed by curve fitting to a cosine wave. These experiments revealed that circadian expression profiles ($P < 0.03$) for 13 genes (76%): 3 E-box controlled genes, 4 D-box controlled genes and 6 RRE controlled genes, respectively, with a consistent order of peak time (4.1, 15.5, and 18.8 for mean value of the peak time of putative E-box, D-box, and RRE-controlled genes, respectively) (Fig. 2B; See also

<http://promoter.cdb.riken.jp/circadian.html> for detailed data). For those genes that did not confirm circadian rhythmicity, the average level of expression was lower, implying mRNA detection was limiting for these genes. Collectively, these *in vitro* and *in vivo* experiments suggest that many predicted E-box, RRE, and D-box containing genes are *bona fide* first-order clock-controlled genes.

Design and Validation of the Synthetic Regulatory Elements. One of the goals of systems biology is the synthesis of knowledge and the generation of testable (and tested) hypotheses. We reasoned that if our HMMs truly represented the functional response elements of these three transcription factor complexes, then synthetic regulatory elements derived from these models should mediate rhythmic transcription as well. To test this idea, we emitted sequences from the E-box, D-box, and RRE models, respectively, and filtered out those that naturally exist in either the human or mouse genomes. Furthermore, to not unduly focus our attention on outliers, we required that all candidates adhere to the consensus rules for each element, CACGTG for the E-box (19), TTATGTAA for the D-box (22), and [A/T]A[A/T]NT[A/G]GGTCA for the RRE (24). For the remaining sequences, we chose each one of the highest and lowest scoring synthetic representatives for three types of elements and named them "high-scoring" and "low-scoring" elements, respectively (Fig. 3A). We tested these elements in a synthetic reporter system as above (Fig. 3B). All three "high-scoring" elements showed high-amplitude circadian transcriptional activity equivalent to known elements from canonical clock genes (E-box of *Per1*, D-box of *Per3*, and RRE of *Arntl* are used as 1.0, respectively) (21) (Fig. 3C). On the other hand, the "low-scoring" elements emitted from the HMMs showed very low-amplitude transcriptional activity, despite the presence of "consensus" E-box, RRE, or D-box core sequences (Fig. 3C). These results show the utility of this comparative genomics approach in synthetic design of dynamic *cis*-acting elements, as well as highlight the contribution of flanking sequences in generating high-amplitude rhythmicity.

Investigating the Contribution of Flanking Sequences. Using these synthetic elements, we next attempted to explore the contribution of E-box flanking regions to identify critical residues that modulate amplitude and rhythmicity. We clustered their nucleotide sequences, and, interestingly, found two patterns of high-amplitude E-box flanking sequences adjacent to the core CACGTG element (Fig. S4). However, these positions do not absolutely dictate high-amplitude rhythmicity, as some elements that meet these rules exhibit lower-amplitude oscillations, possibly because they exhibit much higher GC content. In either case, these experimental results also imply that amplitude information is encoded in specific residues adjacent to the core consensus element and further strengthen the previous reports by other groups on the importance of flanking sequence of E-box (27, 28, 35, 36). Interestingly, the identified patterns in this study partly overlap with the computational models based on the evolutionarily conserved E-box structure from insects to mammals (27).

High-Amplitude Oscillations Require Appropriate Affinity Balance Between Activators and Repressors. To explore the properties of these elements that result in high amplitude oscillations, we took a simplified molecular modeling and experimental approach. First, we assumed concentrations of activators and repressors were within similar ranges (see also SI Appendix Discussion in more general cases). We further hypothesized that flanking region DNA sequence impacted DNA-binding affinity of clock gene regulators and therefore altered amplitude. We further hypothesized that tightly binding sequences would have higher amplitudes of circadian oscillation. To test this notion, we analyzed the DNA-binding affinity of activators and repressors

other circadian response elements? We hypothesized that these differences might be encoded at the protein sequence level of the DNA-binding domains of activators and repressors. Interestingly, the DNA binding domains of transactivators and transrepressors of the RRE and D-box are more similar to each other (65% identity and 81% homology for RRE regulators, and 44% identity and 69% homology for D-box regulators) than those of E-box transactivators and transrepressor (22~30% identity and 55~57% homology) (Table S4). Based on these findings, we speculate that the DNA-binding domains for transactivators and transrepressors of the RRE and D-box have evolved similar affinities. In contrast, the evolutionarily and structurally divergent regulators of E-boxes, 108 bHLH proteins including several families of activators and repressors, as well as the unrelated *period* and *cryptochrome* gene families, may have required the co-evolution of specific DNA-binding domains and E-box sequences with specific flanking regions to generate higher amplitude rhythmicity.

Conclusion

In summary, we have applied a comparative genomics strategy to the understanding of a dynamic transcriptional regulatory system, the mammalian circadian clock. Our informatics strategy employs a model-based search with excellent statistical properties, the evolutionary conservation of putative transcriptional regulatory elements across mouse and human non-coding regions, and statistical evaluation of false discovery rates in each prediction. Experimental validation of this strategy *in vitro* and *in vivo* using real-time monitoring of transcriptional activity and quantitative PCR assay has led to the identification of dozens of novel clock-controlled genes and the elements that likely dictate their rhythmicity. High-scoring conserved E-boxes (mean HMM-score = 16.15) had a 40% rate of validation, while low-scoring conserved E-boxes (mean HMM-score = 2.5143) had a 7.1% probability of generating high-amplitude rhythmicity in reporter assays. Linear interpolation from these two numbers generates an estimate of approximately 347 novel conserved E-boxes that likely confer circadian rhythmicity (see also *SI Appendix*). Moreover, to demonstrate their predictive nature, we have taken these *in silico* models and designed synthetic elements that exhibit high-amplitude transcriptional rhythmicity as well as the best canonical regulatory elements. Furthermore, experimental measurement and *in silico* analysis of affinity of regulators to synthetic elements revealed the importance of the appropriate affinity balance between activators and repressors for high-amplitude rhythmicity. Surprisingly, for E-box sequences, lower affinity DNA element generates higher amplitude rhythms. The experimental, analytical, and synthetic approaches discussed here are especially timely as genomics tools are increasingly uncovering the complexity and flexibility of transcriptional regulatory circuits. We predict the general themes and resources reported here will enhance understanding of the biology mediated by complex and dynamic transcriptional regulation including the mammalian circadian clock.

Materials and Methods

Detailed information on the construction of the mammalian promoter/enhancer database, determination of distance from TSSs for natural and randomly positioned elements, calculation of FDR for putative elements, animals, genome sequences, oligonucleotide sequences, plasmid constructions, quantitative PCR, rhythmicity analysis of real-time bioluminescence data, amplitude analysis of real-time bioluminescence data, rhythmicity analysis of quantitative PCR data, over representation analysis of clock-controlled genes, estimation of the number of high-amplitude E-boxes, microarray expression data analysis of genes with predicted clock-controlled elements, affinity analysis of competitive DNA binding data, and *in silico* analysis of affinity to amplitude mechanism are available in *SI Appendix*.

Real-Time Circadian Reporter Assays. Real-time circadian assays were performed as previously described (40) with the following modifications. NIH

3T3 cells (American Type Culture Collection) were grown in DMEM (Invitrogen) supplemented with 10% FBS (JRH Biosciences) and antibiotics (25 units ml⁻¹ penicillin, 25 µg ml⁻¹ streptomycin; Invitrogen). Cells were plated at 5 × 10⁴ cells per well in 24-well plates 24 h before transfection. Cells were transfected with 0.32 µg of plasmids in total (0.13 µg reporter plasmid and 0.19 µg empty plasmid) per well using FuGENE6 (Roche Applied Science) according to the manufacturer's instructions. After 72 h, medium in each well was replaced with 500 µl of culture medium (DMEM/10% FBS) supplemented with 10 mM Hepes (pH 7.2), 0.1 mM luciferin (Promega), antibiotics and 0.01 µM forskolin (nacalai tesque). Bioluminescence was measured with photomultiplier tube (PMT) detector assemblies (Hamamatsu Photonics). The modules and cultures were maintained in a darkroom at 30 °C and interfaced with computers for continuous data acquisition until 96 h after forskolin stimulation. Photons were counted 2 min at 24-min intervals.

Construction, Search, and Design of Putative *cis*-Acting Elements. A HMM is a statistical model in which the target system is assumed to be a Markov process with unknown parameters. A HMM describes a probability distribution over input training sequences, i.e., probabilities of the state transition and emission. The extracted model can be used to find the probability of query sequence that is a product of all transition and emission probabilities at training sequences. Nucleotide sequences for known functional clock-controlled elements, 12 E-boxes (18 bp), 10 D-boxes (24 bp), and 15 RREs (23 bp), experimentally verified in previous (21) and current studies (Table S1 and Fig. S2), were used as a training dataset to construct HMMs. We also attempted to construct an HMM for the E'-box, but were unable (i.e., positive controls exhibited poor scores) because of the small number of experimentally validated E'-box (only three: *Per2*, *Bhlhb3*, and *Cry1*) and the relatively short core consensus sequence of the E-box. Thus, we did not use an E'-box HMM in this study. The lengths of these known functional elements were based on our previous experiments (21) and these were sufficient to produce circadian transcriptional activity in circadian reporter assays. These sites were aligned without gaps according to the direction of consensus sequences (TTATGT/C)AA for the D-box; ref. 22), [A/T]A[A/T]NT[A/G]GGTCA for the RRE; ref. 24). Because the consensus sequence for E-box is palindromic (CACGTG; ref. 19), we generated all possible alignments by changing sequence directions (forward and reverse) and selected one alignment as described below. These alignments were used to build HMMs using hmmt program in the HMMER 1.8.4 software package (31) with default parameters (using sim annealing, starting kT for sim annealing run as 5.0, and multiplier for sim annealing as 0.95). We used the older version 1.8.4 package (the current version is 2.3.2) in this study because the version 2 series was optimized for analysis of protein sequences. Following construction, models were used to search genomic regions for putative clock-controlled elements using the hmmls program with default parameters (by using threshold matches score to report as 0) except use '-c' option only for bidirectionally search. The average score was used in the search for the conserved elements between human and mouse. To select only one alignment for each E-box, we constructed 2048 HMMs of all possible alignments, and calculated match scores of 12 known E-box sequences in directional HMMER search. We selected the alignment that generated the highest average match score for further work.

To design the "high-scoring" and "low-scoring" sequence of clock-controlled elements, bidirectional HMMER searches were performed against all possible sequences of the same lengths as training dataset (18 bp for E-box, 24 bp for D-box, and 23 bp for RRE) that contain ordinary consensus sequence at the center (CACGTG for E-box; ref. 19; TTATGTAA for D-box; ref. 22, [A/T]A[A/T]NT[A/G]GGTCA for RRE; ref. 24), then filtered out those that naturally exist in either the human or mouse genome. The sequence of the highest and lowest score was selected as the "high-scoring" and "low-scoring" sequences, respectively. All HMMER searches, except the directional search in the selection of E-box alignments, were performed bidirectionally. The higher score was adopted if match scores were obtained for both directions at the same position. The training data are available in Table S2. The HMMs are publicly available on the circadian section of the mammalian promoter/enhancer database: <http://promoter.cdb.riken.jp/circadian.html>.

Competitive DNA Binding Assays. *In vitro* transcription/translation of Flag-tagged mouse protein from pMU2-*Arntl*, pMU2-*Clock*, pMU2-*Bhlhb2*, pMU2-*Dbp*, pMU2-*Nfil3*, pMU2-*Nr1d1*, and pMU2-*Rora* were performed with TNT T7 Quick Coupled Transcription/Translation System (Promega) according to the manufacturer's specifications. *In vitro* transcribed/translated *Arntl* and *Clock* proteins were mixed in equal volume. The complementary oligonucleotides of three tandem repeats sequence of designed and control *cis*-acting elements, which were labeled with biotin on 5'-end or non-labeled (for competitor)

