

# RNAiFold2T: Constraint Programming design of thermo-IRES switches

Juan Antonio Garcia-Martin<sup>1</sup>, Ivan Dotu<sup>2</sup>, Javier Fernandez-Chamorro<sup>3</sup>,  
Gloria Lozano<sup>3</sup>, Jorge Ramajo<sup>3</sup>,  
Encarnacion Martinez-Salas<sup>3</sup>, Peter Clote<sup>1\*</sup>

1: Biology Department, Boston College, Chestnut Hill, MA 02467 (USA).

2: Research Programme on Biomedical Informatics (GRIB), Department of Experimental and Health Sciences, Universitat Pompeu Fabra. Dr. Aiguader 88. Barcelona, (Spain).

3: Centro de Biología Molecular Severo Ochoa, Consejo Superior de Investigaciones Científicas – Universidad Autónoma de Madrid, Nicolás Cabrera 1, 28049 Madrid (Spain).

## Abstract

**Motivation:** RNA *thermometers* (RNATs) are *cis*-regulatory elements that change secondary structure upon temperature shift. Often involved in the regulation of heat shock, cold shock and virulence genes, RNATs constitute an interesting potential resource in synthetic biology, where engineered RNATs could prove to be useful tools in biosensors and conditional gene regulation.

**Results:** Solving the 2-temperature inverse folding problem is critical for RNAT engineering. Here we introduce **RNAiFold2T**, the first Constraint Programming (CP) and Large Neighborhood Search (LNS) algorithms to solve this problem. Benchmarking tests of **RNAiFold2T** against existent programs (adaptive walk and genetic algorithm) inverse folding show that our software generates two orders of magnitude more solutions, thus allowing ample exploration of the space of solutions. Subsequently, solutions can be prioritized by computing various measures, including probability of target structure in the ensemble, melting temperature, etc. Using this strategy, we rationally designed two thermosensor internal ribosome entry site (*thermo*-IRES) elements, whose normalized cap-independent translation efficiency is approximately 50% greater at 42°C than 30°C, when tested in reticulocyte lysates. Translation efficiency is lower than that of the wild-type IRES element, which on the other hand is fully resistant to temperature shift-up. This appears to be the first purely computational design of functional RNA thermoswitches, and certainly the first purely computational design of functional thermo-IRES elements.

---

\*Corresponding author: [clote@bc.edu](mailto:clote@bc.edu)

**Availability:** RNAiFold2T is publicly available as part of the new release RNAiFold3.0 at <https://github.com/clotelab/RNAiFold> and <http://bioinformatics.bc.edu/clotelab/RNAiFold>, which latter has a web server as well. The software is written in C++ and uses OR-Tools CP search engine.

**Contact:** [clote@bc.edu](mailto:clote@bc.edu)

**Final version:** This article will appear in *Bioinformatics* 2016.

## 1 Introduction

RNA *thermometers* (RNATs), also known as *thermosensors*, are *cis*-regulatory elements that change secondary structure upon temperature shift. Examples include (1) repression of heat shock gene expression (ROSE) elements [28], that control the expression of small heat shock genes, such as *hspA* in *Bradyrhizobium japonicum* and *ibpA* in *Escherichia coli*, (2) FourU elements [40], such as the virulence factor LcrF in *Yersinia pestis*, (3) Hsp17 thermosensor [36, 18], which controls membrane integrity of the cyanobacterium *Synechocystis sp. PCC6803* under stress conditions, critical for photosynthetic activity. Additional examples are described in [17]. ROSE elements and FourU elements operate as temperature-sensitive, reversible zippers, while the *Listeria monocytogenes* *prfA* thermosensor [16], phage  $\lambda$  cIII thermoswitch [1] and *E. coli* CspA cold shock thermometer [3] operate in a switch-like fashion. Here, the helix of a zipper melts gradually with increasing temperature, returning to the original structure when temperature is reduced, while a switch consists of two mutually exclusive structures determined by temperature.

Several bioinformatics search methods exist to identify and predict candidate RNA thermometers. In [39] the database RNA-SURIBA (Structures of Untranslated Regions In Bacteria) was created; using regular expressions, particular structural motifs were detected in the minimum free energy (MFE) structure, as determined by *mfold* [45]. In contrast, the *RNAtips* web server [6] and the *RNAthermsw* [5] web server both rely on base pairing probabilities computed at different temperatures using *RNAfold* from the *Vienna RNA Package* [21].

For some time now, RNA thermosensors have been recognized as an attractive target for rational design [42, 41]. Indeed, within the broader context of synthetic biology, rationally designed thermometers could be used as a *thermogenetic* tool to control expression by temperature regulation (i.e. *on-demand protein translation*), or even as a multifunctional nanoscale devices to measure temperature in the context of hyperthermic treatment of cancer cells, imaging, or drug delivery [20].

In [27], synthetic (zipper) thermosensors were *manually* designed to sequester the Shine-Dalgarno (SD) sequence AAGGAG within a single stem-loop structure containing 4-9 base pairs, several of which contained 1-2 bulges of size 1. In [41], synthetic (switch) thermosensors were *computationally* designed to switch between a single stem-loop structure that sequesters the SD sequence GGAGG, and two shorter stem-loop structures where the SD sequence is found in the

apical loop of the second stem-loop. In particular, the 2-temperature inverse folding (adaptive walk) program **SwitchDesign** [12] was used to obtain 300 candidate sequences; only two candidate sequences survived after the application of several computational filters including the computation of melting curves with **RNAheat** [21]. Since neither of these sequence displayed any temperature-dependent control of a reporter gene (bgaB) fusion, the top candidate sequence was used as a template in two rounds of error-prone mutagenesis followed by selection, resulting in a successful thermosensor – see Figure 5 of [41]. In [29], a non temperature-dependent *riboswitch* was manually designed, which promotes cap-independent translation in wheat germ cell lysate only upon binding of the ligand theophylline. In [38], a synthetic theophylline riboswitch was designed by a computational pipeline including inverse folding and experimentally shown to perform transcriptional regulation in *Escherichia coli* Top10 cells. In [35], a *thermozyme* was created by fusing a *Salmonella* RNA thermometer (RNAT) to a hammerhead ribozyme, followed by *in vivo* screening – thus showing that naturally occurring hammerheads and RNATs appear to be modules that can be combined. In [15] small, heat-repressible RNA thermosensors (zippers) were manually designed in *E. coli*, which at low temperature sequester a cleavage site for RNaseE, and at high temperatures unfold to allow mRNA degradation.

Despite these impressive results, [17] state that: “RNATs have little, if any sequence conservation and are difficult to predict from genome sequences. ... Therefore, the bioinformatic prediction and rational design of functional RNATs has remained a major challenge”.

In this paper, we introduce the software **RNAiFold2T**, capable of solving the inverse folding problem for two or more temperatures, i.e. generating one or more RNA sequences whose minimum free energy (MFE) secondary structures at temperatures  $T_1$  and  $T_2$  [resp.  $T_1, \dots, T_m$ ] are user-specified target structures  $S_1$  and  $S_2$  [resp.  $S_1, \dots, S_m$ ], or which reports that no such solution exists. **RNAiFold2T** is unique in that it implements two different algorithms – *Constraint Programming* (CP) and *Large Neighborhood Search* (LNS). CP is an exact, non-heuristic method that uses an exhaustive yet efficient branch-and-prune process, and is the only currently available software capable of generating all solutions or determining that no solution exists (since there are possibly exponentially many solutions, a complete solution is feasible only for structures of modest size). LNS uses a local search heuristic, complemented with local calls of constraint programming to explore solutions of substructures of the target structures. We use **RNAiFold2T** to rationally design two *thermo*-IRES elements, whose normalized cap-independent translation efficiency is approximately 50% greater at 42°C than 30°C, when tested in reticulocyte lysates. We then benchmark **RNAiFold2T** with the only two other programs that solve the 2-temperature inverse folding problem: the adaptive walk **SwitchDesign** (SD) [12] and the genetic algorithm **Frnakenstein** (FRNA) [25]. **RNAiFold2T** CP generates two orders of magnitude more solutions than either, when all programs are run for 24 hours, while the number of distinct problem instances that can be solved by **RNAiFold2T** LNS within 30 [resp. 60] minutes for short [resp. long] target structures is roughly comparable for each program. Finally, by analyz-

ing existent RNATs found in the Rfam 12.0 database [26], we determine that naturally occurring RNATs appear not to be optimized for the *cost* function (equation (7) of [12]), and that both SD and FRNA appear to generate solutions whose *cost* function value is substantially lower than that of naturally occurring RNATs, in contrast to solutions returned by RNAiFold2T.

## 2 Methods

### 2.1 Biochemical methods

***Thermo-IRES activity assay:*** Thermo-IRES constructs were created by replacement of wild-type nucleotides at positions 417-462 in domain 5 of the foot-and-mouth disease virus (FMDV) IRES element, whose secondary structure is depicted in Figure 1 of [10] and Figure 2 of [22]. Six computationally designed thermo-IRES elements were selected for validation, along with a negative control and the wild-type FMDV IRES element as positive control. Specifically, synthetic oligonucleotides containing the designed sequence (46 nts) in either positive or negative orientation were annealed in Tris 50 mM pH 7.5, NaCl 100 mM, MgCl<sub>2</sub> 10 mM, 15 min at 37°C and subsequently inserted into the HindIII and XhoI restriction sites of pBIC, which harbors the wild-type IRES, linearized with the same enzymes. Colonies that carried the correct insert were then selected, and prior to expression analysis, the nucleotide sequence of the entire length of each region under study was determined (Macrogen).

*In vitro* transcription was performed for 1 h at 37°C using T7 RNA polymerase, as described in [24]. RNA was extracted with phenol-chloroform, ethanol precipitated and then resuspended in TE. Using gel electrophoresis, the transcripts were checked for integrity. Equal amounts of the RNAs synthesized *in vitro* were translated in 70% rabbit reticulocyte lysate (RRL) (Promega) supplemented with <sup>35</sup>S-methionine (10  $\mu$ Ci), as described in [30]. Each experiment was independently repeated in triplicate, using the wild type RNA as a control in all assays. Luciferase (LUC) and chloramphenicol acetyl transferase (CAT) activities were measured for the bicistronic plasmid, as previously described [11]. In particular, intensity of the LUC band, as well as the CAT band, produced by each transcript was determined in a densitometer, and normalized against the intensity of LUC and CAT bands produced by the wild type RNA, set at 100%. Values represent the mean  $\pm$  SD.

Luciferase activity reflects the efficiency of IRES-dependent translation, while CAT activity reflects the efficiency of 5'-dependent translation; thus the ratio LUC/CAT was determined at 30°C and 42°C for wild-type FMDV IRES, the negative control and two thermo-IRES constructs. Selective 2'-hydroxyl acylation analyzed by primer extension (SHAPE) was performed for wild-type IRES element, as described in [9]

## 2.2 Computational methods

RNAiFold2T uses Constraint Programming (CP) to determine those sequences, whose minimum free energy (MFE) structure at temperature  $T_1$  [resp.  $T_2$ ] is identical to a user-specified target structure  $S_1$  [resp.  $S_2$ ]. The target structures  $S_1, S_2$  can also be hybridization complexes of two RNAs, rather than single secondary structures. CP performs a complete, exhaustive (branch and prune) exploration of the search space and therefore, it can return all possible solutions of the thermoswitch design problem or prove that no solution exists (given an unlimited amount of time). In addition to CP, RNAiFold2T also supports Large Neighborhood Search (LNS), a fast local (not complete) search metaheuristic that employs CP to exhaustively explore large neighborhoods of every candidate solution at each iteration step. Moreover, since it is written in C++ using the OR-Tools engine <https://github.com/google/or-tools>, together with plugins to Vienna RNA Package [21] and RNAstructure [32], the user can install and run RNAiFold2T locally, thus permitting much longer execution times than supported by our web server.

The overall methodology of RNAiFold2T is similar to its precursor, RNAiFold 2.0 [14, 13]; however, as explained below, there are a number of algorithmic details that are new and not present in RNAiFold 2.0 – decomposition tree for 2 or more target structures, novel constraints that are underlined below, variable (helix) and value heuristics that are proper only to RNAiFold2T, the introduction of two types of restart heuristic in LNS to ensure a good trade-off between exploration of promising regions of the search space versus the exploration of remote portions of the search space. RNAiFold2T cleanly separates all constraints from the CP or LNS solver, thus permitting our software to be extended to support any future desirable constraints. Current constraints supported by RNAiFold2T include the following (novel constraints underlined): full and/or partial target structures or hybridization complexes at two temperatures; a plug-in to determine MFE structure using either RNAfold [21] or RNAstructure [32] with the Turner99, Turner2004 [37] or Andronescu [2] energy model with dangle treatment (-d0, -d1, -d2, -d3 corresponding respectively to no dangle, max of 5' and 3' dangle, sum of 5' and 3' dangle, sum of 5' and 3' dangle with coaxial stacking); IUPAC nucleotide constraints, IUPAC amino acid constraints that require all returned RNA sequences to code specified peptides in one or more overlapping reading frames, structural compatibility and structural incompatibility constraints, etc. Additionally, RNAiFold2T can determine a user-specified number of solutions, all solutions (given sufficient time), or whether no solution exists. Indeed, memory requirements for RNAiFold2T are minimal, and since there are no memory leaks, the software can be run for weeks.

In developing a CP solution to a given problem the main tasks are to define the problem (specify variables, domains and constraints) and to define the search procedure. The extension of the model from RNAiFold2.0 is trivial and consists of adding new constraints for the helices corresponding to the second structure. The search procedure, however, must be adapted to the new difficulties imposed by the 2-temperature problem. New variable and value ordering

heuristics are needed in order to solve 2-temperature inverse folding efficiently (see Supplementary Tables 1,2). The algorithmic details related to the new search procedure are explained below:

### 2.2.1 Structure decomposition

As in other inverse folding methods, such as `RNAinverse` [21], `NUPACK` [43], etc., we rely on a decomposition tree of the structure into independent helices, called *extended helices* (EHs) and *extended helices with dangles* (EHwDs). See Supplementary Information (SI) for definitions of EH, EHwD and decomposition tree, and see Figure 1 for an illustration of the EHwD decomposition tree for a FourU thermometer. Decomposition trees play a special role in `RNAiFold2T` for the following reasons: (a) each node in the decomposition tree is a constraint, (b) the helix and variable heuristics (described later) cause the search tree to be searched in a specific order. To improve efficiency in solving multi-temperature inverse folding, we investigated various helix and value heuristics, which steer the search within a search space defined by a composite decomposition tree, comprising subtrees for each target structure at the corresponding temperature. The helix and value heuristics are new and not part of `RNAiFold 2.0`.

Consider the 65 nt FourU RNA thermometer whose MFE structures determined by `RNAfold` from `Vienna RNA Package 2.1.9` [21] at 37°C and 53°C are given in dot-bracket notation by

```
>CP000647.1/1773227-1773291)
12345678901234567890123456789012345678901234567890123456789012345
GGACAAAGCAAGUCCUUGCCUUUUGUUGAGCUUUUGAAUGAAUUAUUCAGGAGGUAAUUAUGGCAC
((.((((.....))))).(((C.(((((((.....))))))))).))))).))))).
.....(((C.....(((C.(((((((.....))))))))).))))).))))).
```

Let  $S_1$  [resp.  $S_2$ ] denote the MFE structure of this FourU thermometer at 37°C [resp. 53°C].  $S_1$  is identical to the consensus structure from Rfam 12.0 [26], as well as the structure displayed in Figure 1 of [40] for the FourU sequence taken from the 5'-UTR of the *Salmonella agsA* gene. Giving labels as described to nodes in the EHwD decomposition trees  $T_1, T_2$  respectively for  $S_1$  and  $S_2$ , we find the EHwD decomposition of  $S_1$  has EHwD 0 from positions 1-65, EHwD 1 from positions 1-19, and EHwD 2 from positions 23-58, while the EHwD decomposition of  $S_2$  has EHwD 4 from positions 1-65, EHwD 5 from positions 14-65, and EHwD 6 from positions 23-58. The lower portion of Figure 1 depicts the EHwD decomposition trees for temperatures  $T_1, T_2$ , joined together with a (dummy) root that corresponds to the solution returned by `RNAiFold2T`.

It can happen that the MFE structure of an extended helix does not agree with the target substructure, while that of the extended helix with dangles does. This is the reason that the EH decomposition trees of `RNAiFold 1.0` [14] were replaced by EHwD decomposition trees in `RNAiFold 2.0` [13], and why EHwD trees are used in `RNAiFold2T`.

In the CP and LNS search strategy, whenever a subsequence corresponding to a node of the decomposition tree has been instantiated, a check is made to determine whether the MFE structure of the subsequence is identical to the target substructure. In the case of structural disagreement, the instantiated subsequence is discarded and backtracking occurs. For any solution sequence

returned by `RNAiFold2T`, it follows that at temperature  $T_1$  [resp.  $T_2$ ], each subsequence of the solution that corresponds to a node in decomposition tree  $\mathcal{T}_1$  [resp.  $\mathcal{T}_2$ ] folds into the corresponding substructure of target  $S_1$  [resp.  $S_2$ ]. The top portion of Figure 1 illustrates the differences between extended helix (EH) and extended helix with dangle (EHwD) decomposition, and when a check is made for whether the target substructure agrees with the MFE structure of the instantiated subsequence. In the sequel, we refer to four types of elements in an EHwD:

- Dangling position: Unpaired position at any side of a helix.
- Unpaired position: Any other unpaired position.
- Closing base pair: Outermost base pair of a helix.
- Normal base pair: Any other base pair.

### 2.2.2 Heuristics for variable and value order

In a Constraint Programming (CP) algorithm, one typically specifies the order in which variables are instantiated (assigned), known as the *variable ordering heuristic*, as well as the order in which the values belonging to the domain of each variable are to be assigned, known as the *value ordering heuristic*. The variable ordering heuristic is divided into two levels: first, the order in which extended helices with dangles (EHwDs) are to be assigned, and second, the order in which nucleotide positions within helices are to be assigned. In this section, we also discuss *restarting heuristics* for the Large Neighbourhood Search (LNS) variant of `RNAiFold2T`.

**Variable ordering at the helix level:** In the search for thermosensors, there is often an overlap between EHwDs of structure  $S_1$  and those of structure  $S_2$  – this situation substantially complicates the task of finding an optimized order of exploration of the CP search space. In the leaf-to-root heuristic of `RNAiFold 2.0`, EHwD node  $H$  is explored before EHwD node  $H'$  if the height  $ht(H)$  of  $H$  is less than the height  $ht(H')$  of  $H'$ , or if  $ht(H) = ht(H')$  and  $H$  appears to the left of  $H'$  in the decomposition tree for the single target structure  $S$ . In contrast to this heuristic, `RNAiFold2T` implements two different approaches in order to find an adequate exploration ordering for the extended helices with dangles for two target structures  $S_1$  and  $S_2$ , whereby a high priority is given to solve those helices, whose sequence is determinant for other parts of the structure due to overlaps. Let  $N$  denote the number of nodes (EHwDs) in the decomposition tree for  $S_1$ , plus the number of nodes (EHwDs) in the decomposition tree for  $S_2$ . Suppose that  $H, H'$  are two distinct EHwDs belonging to  $S_1$  or  $S_2$ , where the outermost base pair of  $H$  [resp.  $H'$ ] is  $(i, j)$  [resp.  $(i', j')$ ]. Define the following relations:

- $includes(H, H')$  is 1 if  $i < i'$  and  $j > j'$ , i.e. interval  $[i, j]$  properly contains  $[i', j']$ ; otherwise  $includes(H, H')$  is 0.

- $overlap_1(H, H')$  is 1 if  $[i, j] \cap [i', j'] \neq \emptyset$ , or equivalently  $\max(i, i') \leq \min(j, j')$ ; otherwise  $overlap_1(H, H')$  is 0.
- $overlap_2(H, H')$  is the number of positions  $k$  in  $H$  and  $H'$ , for which  $k$  is base-paired in both  $H$  and  $H'$ .
- $degree_\alpha(H) = \sum_{i=0}^N overlap_\alpha(H, H_i)$ , for  $\alpha = 1, 2$ .
- For  $\alpha = 1, 2$ ,  $degreedist_\alpha(H, H')$  is equal to  $degree_\alpha(H) - degree_\alpha(H')$ , provided that  $degree_\alpha(H) > degree_\alpha(H')$ ; otherwise,  $degreedist_\alpha(H, H')$  is 0.
- For  $\alpha = 1, 2$ ,  $diff_\alpha(\sigma, H, H')$  is equal to  $degreedist_\alpha(H, H')$ , provided that  $\sigma(label(H)) < \sigma(label(H'))$ ; otherwise  $diff_\alpha(\sigma, H, H')$  is 0.

The value  $label(H)$  is defined in SI, and corresponds to visitation order in breadth-first traversal of tree  $\mathcal{T}$ , and  $\sigma : \{0, \dots, N-1\} \rightarrow \{0, \dots, N-1\}$  is a permutation that minimizes  $\sum_{i=0}^N \sum_{j=0}^N diff_\alpha(\sigma, H_i, H_j)$ , subject to the constraint that if  $H$  properly includes  $H'$ , then  $order(H) > order(H')$ . Note that the partition  $\sigma$  orders the EHwDs of  $S_1$  and  $S_2$  in order to minimize the total *incremental overlap*. Before the search for thermosensors begins, RNAiFold2T executes a very fast CP search to determine the optimal ordering permutation  $\sigma$ . Finally, by setting the index  $\alpha$  to either 1 or 2 in the definition of  $diff_\alpha(\sigma, H, H')$ , we obtain the first or second search heuristic. (Two additional helix ordering heuristics are explained in Supplementary Information.) An example of each helix order heuristic 1,2 is shown in Figure 1 – note that for even small structures, there are several differences in the helix exploration order when  $\alpha$  is 1 or 2.

**Variable ordering at the nucleotide level:** The second level of variable ordering heuristic deals with the exploration of nucleotide positions within a given EHwD structure. This second level of variable ordering can be stated as follows:

- First, non-outermost base-paired positions  $(x, y)$  of a given EHwD are instantiated from the innermost base pair to the outermost base pair.
- Second, unpaired positions in a given EHwD are grouped together in consecutive runs, and these runs are ordered from largest to smallest and then instantiated from left to right.
- Third, the outermost, closing base pair of a given EHwD is instantiated.
- Finally, dangling positions of a given EHwD (if any) are instantiated (note that not all EHwDs contain dangling positions).

The small example shown in Figure 2 graphically depicts the differences between constraint checks and variable ordering for both RNAiFold and RNAiFold2T. Supplementary Tables 1,2 describe and benchmark variable and value heuristics in RNAiFold2T.



**Value ordering:** Value ordering establishes the order in which values are assigned to variables – in our case, this means values GC,CG, AU, UA, GU, UG for base-paired positions and A,C,G,U for unpaired positions. The underlying idea for ordering domain values for variables for base-paired positions and unpaired positions is to allow the creation of thermodynamically stable helices and to take into account the nature of the overlap in overlapping positions. This requires specific value orderings for base-paired positions, depending on whether the position is a dangle, mismatch, normal or closing base pair of an EHwD for both targets  $S_1, S_2$ , or only one of the target structures. For each of these cases, we define different ordering heuristics, described in Supplementary Table 1.

**LNS restart heuristics:** Similarly as in the LNS variant of RNAiFold, the restart condition is a given amount of time, proportional to the length of the target structures, after which search is stopped and some variables are fixed in order to explore exhaustively a large neighbourhood of the current solution. After the first restart, full exploration of the remaining space with no solution found is also a restart condition.

In RNAiFold2T, when a restart is triggered, a set of positions is selected as candidates to be fixed. The MFE structure for each EHwD of the current candidate solutions is evaluated independently. If the MFE structure of an EHwD matches with the target structure at the desired temperature, and the MFE structure of all the overlapping helices in the second target structure (at the corresponding temperature) also matches, then all the EHwD positions are included in the set of candidates. When all the EHwDs have been evaluated, candidate positions are fixed with a probability of 0.9, and the set of candidate positions is stored.

Since the order of exploration is similar in each round, it could be possible that fixing similar parts of the sequence results in an exploration of almost the same region of the search space in subsequent searches, so two mechanisms are implemented to avoid this behavior: (1) In subsequent restarts, if the candidate positions to be fixed are the same as in the previous restart, then the probability of fixing positions decreases by 0.05, if not, then the initial probability of 0.9 is restored. (2) There is a hard restart (no nucleotide position is fixed) in the case that, after 10 restarts, the set of candidate positions remains unchanged, or if all possible solutions for the current subproblem have been explored.

In local search algorithms, there is always a trade-off between *exploitation* and *exploration*. Exploitation means focusing the search on promising regions, as reflected in our choice of probability 0.9 to remain close to currently instantiated portions of the sequence. Exploration means covering different, remote regions of the search space, as reflected in our choice to decrement the probability by 0.05 and our choice to perform a hard restart after 10 restarts.

## 2.3 Benchmarking

**Benchmarking data:** We retrieved the sequences of seven families of non-coding RNA thermometers from Rfam: RF00038, RF00433, RF00435, RF01766,

RF01795, RF01804, RF01832. These families include both cold and heat shock RNA thermometers, taken from diverse organisms including phages, prokaryotes and eukaryotes, with sequence length ranging from 60 nt to 450 nt. The benchmarking was divided into two groups: sequences shorter than and longer than 130 nt. For each sequence, we used **RNAfold** [21] with the Turner99 energy parameters [37] to determine the MFE structures at temperatures  $T_1$  and  $T_2$ , where temperatures were chosen (essentially) according to published experimental studies for each thermometer family [33, 4, 1] – in particular, we increased the temperature difference  $T_2 - T_1$  from the published values to ensure that **RNAfold** produced distinct structures at  $T_1$  and  $T_2$  if possible (see SI). Turner99 rather than Turner2004 energies were used, since it required less distortion from published temperatures  $T_1, T_2$ . All sequences, whose MFE structures at  $T_1$  and  $T_2$  were identical, were subsequently removed. See Supplementary Tables 3-5 and SI Excel files for benchmarking results, and Supplementary Table 6 [resp. SI Excel file] for a list of sequences, structures, and temperatures for sequences of length less than 130 nt [resp. greater than 130 nt]. The resulting benchmarking set includes all 5 *Lambda* phage CIII thermoregulator elements (Lambda.thermo), all 3 FourU thermometers, 11 of 13 repression of heat shock (ROSE) elements, 8 of 14 sequences from a second family of repression of heat shock (ROSE\_2) elements, 3 of 13 thermoregulators of PrfA virulence genes (PrfA), 4 of 6 HSP90 *cis* regulatory elements (HSP\_CRE), and 14 of 15 cold shock protein regulator sequences (CspA).

**Software used in benchmarking tests:** **SwitchDesign** (SD) [12], **Frnakenstein** (FRNA) [25], and **RNAiFold2T** were benchmarked on Rfam family RF01804 of *Lambda* phage CIII thermoregulator elements to determine the maximum number of distinct solutions over 24 hours, restarting when no new solution is found within 1 hour. Since neither **FRNA** nor **SD** output more than one solution, we made the following modifications of each program. The genetic algorithm **Frnakenstein** was run as many times as possible over 24 hours (each time with a time limit of 1 hour); we then output all sequences found in the most recent (internally stored) population which fold into the target structures  $S_1, S_2$  at temperatures  $T_1, T_2$ . **SD** returns a single sequence which minimizes a cost function described in [12]; thus we modified **SD** source code in order to test whether any sequence explored in the search was a solution. Sequences were checked at two different points in **SD**: when a new sequence is generated by a single mutation (**SD** update), and when a sequence is selected by minimization of the cost function (**SD** selected). In all cases, **SwitchDesign** was restarted if no new sequence was found in one hour.

For each solution set obtained, additional solutions were generated by testing all single point mutations of any solution returned. Additionally, a reference solution set was produced by running **RNAiFold2T** for several days.

### 3 Results

### 3.1 Design of thermo-IRES switches

Domain 5 of wild-type FMDV IRES element contains the 46 nucleotides AUAG-GUGACC GGAGGUCGGC ACCUUUCCUU UACAAUAAU GACCCU at positions 417-462, as shown in Figure 1 of [10] and Figure 2 of [22]. The domain 5 stem-loop at positions 419-440 and *unpaired* polypyrimidine tract (PPT) region at positions 441-447, are both known to be essential for IRES activity [19]. Using the following pipeline, two candidate thermo-IRES elements were tested, along with a negative control and a positive control (wild-type IRES).

1. As shown in Figures 3a and 3b, the *inactive* target structure  $S_1$  at  $T_1 = 30^\circ\text{C}$  was chosen to destroy domain 5 stem-loop and unpaired PPT region, while target *active* structure  $S_2$  at  $T_2 = 42^\circ\text{C}$  is the experimentally determined structure of wild-type domain 5 FMDV IRES element. Sequence constraints were chosen in accordance with conservation observed in a multiple alignment of 183 IRES elements [10], where we added AUG start codon at positions 47-49 (corresponding to IRES positions 463-465).

Inactive S1: .....(((((((.....))))))..))))).  
Active S2: ..(((((((.....)))))).....  
Constraints: NUAGGNGACCGNAGGCGGCNCNUUYYYYYYRNNNNNNNNNNNNNAUG

Using RNAiFold2T, 24,410 solutions were generated, although an additional 45,442 sequences were generated using variants of target  $S_1$ .

2. RNAiFold2T solutions were discarded if any of the following criteria were not met:
- Wild-type structures for domain 4 and domain 5 appear as stable substructures using **Vienna Package RNALfold -L 110 -T 42**.
  - Domain 4 appears as a stable substructure using **RNALfold -L 110 -T 30**
  - Probability  $Pr(S_2, T_2)$  of active conformation at  $T_2 = 42^\circ\text{C}$  exceeds 0.2
  - Probability  $Pr(S_1, T_1)$  of inactive conformation at  $T_1 = 30^\circ\text{C}$  exceeds 0.2
  - Probability of intended target structure at intended temperature is more than double that of unintended target, i.e.  $Pr(S_1, 42) / Pr(S_2, 42) < 0.5$  and  $Pr(S_2, 30) / Pr(S_1, 30) < 0.5$ .
3. Retained solutions were further filtered using various measures. For instance, candidate 1 (Seq1) had the highest value of  $A + B$ , where  $A$  is  $a \cdot (b \cdot Pr(S_1, 30) + c \cdot (Pr(S_1, 30) - Pr(S_2, 30)))$  and  $B$  is  $d \cdot (1 - Pr(S_1, 30)) + e \cdot (b \cdot Pr(S_2, 42) + c \cdot (Pr(S_2, 42) - Pr(S_1, 42)))$ , and  $a = 4, b = 0.5, c = 0.5, d = 2, e = 1$ . This measure was designed to select sequences where the probability of the intended target at the intended temperature is high, while probability of the unintended target is low. The measure is

weighted to increase the likelihood of not having the inactive conformation at 42°C. In contrast, Candidate 2 (Seq2) is one of two sequences satisfying  $Pr(S_2, 42) > 0.3$  and  $P(S_1, 30) > 0.3$ . See Supplementary Information (SI) for a spread sheet of measures and their values.

Seq1 and Seq2 consist of the following 46 nt: Seq1 is AUAGGUGACC GGAGGGCGGC ACCUUUUUUC CAGAAAAGUA GUCGUC (15/46 positions differ from wild-type) and Seq2 is GUAGGUGACC GGAGGACGGC ACCUUUUUUC CAGAAAAGUA GUCGUC (16/46 positions differ from wild-type). Figure 3c shows that Seq1 and Seq2 displayed an increase of approximately 50% normalized IRES-dependent translation efficiency in RRL at 42°C versus 30°C. Seq1 and Seq2 IRES elements displayed about 20% lower normalized activity than the wild type IRES. Nonetheless, the wild type IRES was equally active at all temperatures tested (30, 37 and 42 °C).

### 3.2 Comparison of RNAiFold2T with other software

Here we compare the Constraint Programming (CP) and Large Neighborhood Search (LNS) programs of **RNAiFold2T** with the adaptive walk program **SwitchDesign** (SD) [12] and the genetic algorithm **Frnakenstein** (FRNA) [25]. Below, we describe benchmarking results for datasets of thermosensor target structures  $S_1$  resp.  $S_2$  at temperatures  $T_1$  resp.  $T_2$ , constructed as described in Methods. All benchmarking was carried out on a Core2Duo PC (2.8 GHz; 2 Gbyte memory; CentOS 5.5).

SD, the first algorithm capable of designing thermoswitches, achieves this by optimizing the cost function given in equation (7) of [12] – see also equation (1) in SI. In this context, we wanted to ascertain whether natural thermoswitches are optimized for this cost. Surprisingly, Figure 4 and Supplementary Figure 1 show that natural thermosensor sequences from Rfam appear *not* to be optimized for the cost function used in SD. In particular, SD and FRNA return solutions having substantially lower cost values (i.e. more optimal) than those of natural thermosensors, whose cost value appears to be the mean value returned by RNAiFold2T. Therefore, our benchmarking comparison of SD, Frnakenstein, and RNAiFold2T compares the number of problems solved by each algorithm within the same amount of time.

Supplementary Information (SI) Table 1 describes the value ordering heuristics used in RNAiFold2T; SI Table 2 benchmarks RNAiFold2T with respect to different helix ordering and value heuristics, using a cutoff time of 10 minutes. Since the data clearly demonstrates the superiority of *overlap<sub>2</sub>* helix ordering, this is taken as the default for all other benchmarks and for the web server. SI Table 3 [resp. SI Table 4] presents benchmarking data for Large Neighborhood Search (LNS) from RNAiFold2T and SwitchDesign and Frnakenstein, each with a cutoff time of 30 minutes, using Rfam thermosensor target structures of length less than 130 nt [resp. greater than 130 nt] – see Methods for construction of target structures and temperatures. RNAiFold2T has essentially the same performance as SD and FRNA for shorter sequences, while SD performs

better than other methods for longer sequences. Target structures  $S_1, S_2$  and temperatures  $T_1, T_2$  for this test are given in SI Table 5 displays the number of solutions for  $\lambda$  phage CIII thermoswitches from Rfam family RF01804. Constraint Programming (CP) from `RNAiFold2T`, adaptive walk `SwitchDesign` [12] and genetic algorithm `Frnakenstein` [25] were run on each thermosensor for 24 hours, forcing a restart if no new solution was found within 1 hour. Since both `SwitchDesign` and `Frnakenstein` return only a single solution, for this test we modified each as described in Methods, resulting in two versions of `SwitchDesign` and one of `Frnakenstein`, each of which returns two orders of magnitude less solutions than `RNAiFold2T`. SI Table 6 lists the target structures  $S_1, S_2$  and temperatures  $T_1, T_2$  of RNA thermometers of length less than 130 nt used in the benchmarking tests. The structures for thermometers of length 130-447 nt are too large for display, hence are available in an Excel file in Supplementary Information.

## 4 Discussion

In this paper, we introduce the software `RNAiFold2T`, capable of solving the inverse folding problem for two or more temperatures, i.e. generating one or more RNA sequences whose minimum free energy (MFE) secondary structures at temperatures  $T_1$  and  $T_2$  [resp.  $T_1, \dots, T_m$ ] are user-specified target structures  $S_1$  and  $S_2$  [resp.  $S_1, \dots, S_m$ ], or which reports that no such solution exists. `RNAiFold2T` is unique in that it implements two different algorithms – *Constraint Programming* (CP) and *Large Neighborhood Search* (LNS). CP is an exact, non-heuristic method that uses an exhaustive yet efficient branch-and-prune process, and is the only currently available software capable of generating all solutions or determining that no solution exists (since there are possibly exponentially many solutions, a complete solution is feasible only for structures of modest size). CP differs from what one might call a ‘brute-force’ approach only in that it relies on a highly efficient branch-and-prune search engine, that propagates the effects of currently instantiated variables held within a *constraint store*. LNS uses a local search heuristic, complemented with local calls of constraint programming to explore solutions of substructures of the target structures.

There exists an RNA sequence compatible with any two given secondary more than two structures [31] – nevertheless, `RNAiFold2T` solves inverse folding problem for more than two temperatures. `RNAiFold2T` is modular software, with a clear separation between search procedure and constraint descriptions, thus permitting the future addition of sequence and structural constraints. In its current form, `RNAiFold2T` includes constraints for *full* and/or *partial* target structures or hybridization complexes at two temperatures; a plug-in to use `RNAfold` or `RNAstructure` for MFE structure computation; IUPAC nucleotide constraints, IUPAC amino acid constraints that require all returned RNA sequences to code specified peptides in one or more overlapping reading frames, structural compatibility and structural incompatibility constraints, etc. These constraints support the design of temperature-sensitive selenocysteine insertion

(SECIS) elements, precursor microRNAs, and mRNA domains that are targeted by microRNAs, etc. Since Constraint Programming (CP) is not a heuristic, unlike other methods such as adaptive walk, genetic algorithm, etc., `RNAiFold2T` can in principle return all 2-temperature inverse folding solutions, or prove that none exist. The Large Neighborhood Search (LNS) algorithm of `RNAiFold2T` returns a single solution with approximately the same performance as state-of-the-art approaches `SD` and `FRNA`.

As with the synthetic hammerhead design in [8], our synthetic RNA design strategy consists of generating many solutions, which are prioritized for experimental validation by applying various computational filters. In our opinion, this strategy presents advantages over methods using `SD` or `NUPACK`, each of which returns a relatively small number of sequences that are optimized with respect to a single criterion – in the case of `SD`, this is the cost function [12], and in the case of `NUPACK`, this is ensemble defect [43].

In order to ascertain the viability of our approach of not committing to a particular cost function, we used the capabilities of `RNAiFold2T` to find hundreds of thousands of solutions to the 2-temperature inverse folding problem for target structures from Rfam family RF01804 ( $\lambda$  phage CIII thermoregulators). Figure 4 and SI Figure 1 show that the cost function value of (real) Rfam sequences is not close to the minimum, but rather close to the average of the distribution. Other figures can be found in Supplementary Information, where we investigated a variant of the cost function defined using ensemble defect. So, although `SD` benchmarking results in Supplementary Information indicate that cost function minimization is a good strategy to find sequences whose MFE structures at temperatures  $T_1$  resp.  $T_2$  are the target structures  $S_1$  resp.  $S_2$ , it appears that naturally occurring RNA thermoswitches are not optimized for the `SD` cost function. This observation may be important for the future design of functional synthetic thermoregulators.

In designing thermo-IRES elements, we solved the 2-temperature inverse folding problem depicted in Figure 3, where the AUG start codon was appended at the 3' end of the 46 IUPAC constraint mask NUAGGNGACC GNAG-GNCGGC NCNUUYYYYY YRNNNNNNNN NNNNNN. This was done since we assumed that the AUG start codon was located at position 463, as it occurs in the viral RNA. However, experimental validation was performed by using a construct containing the 35 nt spacer GAGCUCGAGC UUGGCAU-UCC GGUACUGUUG GUAAA at the 3' end of the designed 46 nt sequences, followed by the luciferase AUG start codon. It is possible that the secondary structure of the spacer might have rendered the polypyridine tract less accessible, thus explaining the low efficiency of thermo-IRES constructs Seq1 and Seq2. Another issue is that IRES elements are complex molecules, both in sequence and structure, which respond to more host factors than previously reported RNA thermometers [23]. Secondary structure predictions could differ from the physical structure due to unmodeled protein-RNA interactions, possibly unreliable free energy parameters at temperatures different than 37°C, and due to variations in structure prediction software as shown in Figure 5. Panels (a) resp. (b) of that figure show that the polypyrimidine tract (PPT)

is unpaired in the secondary structure of domain 5 of wild-type IRES at 37°C (hence consistent with experimental data), when determined by **RNAfold** (without SHAPE) resp. **RNAsc** (with SHAPE) [44], while panels (c) and (d) show the PPT is paired in the structure determined by **RNAstructure**, both with and without SHAPE [7] (hence inconsistent with experimental data). Due to this variation in MFE structure prediction, it may be advisable in future synthetic RNA design projects to run **RNAiFold2T** using plug-ins for both **RNAfold** and **RNAstructure**, and to select sequences which have the same desired target structures as predicted by both algorithms. Another issue is that for many RNA thermometers, the minimum free energy (MFE) structures at low temperature  $T_1$  and high temperature  $T_2$  are almost identical, where  $T_1, T_2$  are the temperatures for which a conformational change is reported in the literature. For instance, the MFE structures for the ROSE element CP000009.1/1450710-1450627 at 0°C and 60°C are nearly identical. This could be an important issue for certain software [6, 5] that predict RNA thermometers, but is of less importance in the synthetic design of thermometers, where it is possible to exaggerate the temperature difference  $T_2 - T_1$ . Perhaps future experimental work will provide more reliable energy parameters at temperatures different than 37°C, an issue that affects both RNA and DNA melting temperature predictions [34].

## 5 Conclusions

We present the software **RNAiFold2T** for the multi-temperature inverse folding problem, used to design functional thermoswitches. The CP variant of **RNAiFold2T** returns two orders of magnitude more solutions than other software, while the LNS variant, which returns a single solution, exhibits comparable performance with that of existent methods. The software design of **RNAiFold2T** currently supports a much greater variety of user-defined structural and sequence constraints than other methods, and moreover can be extended to support future constraints.

Naturally occurring RNA thermometers in heat-shock and virulence genes, as well as all previously known instances of rationally designed thermometers control translation initiation by sequestering the Shine-Dalgarno (SD) ribosomal binding sequence at low temperatures, whereby the SD becomes accessible at high temperatures due to the melting of a hairpin. In contrast, we have used **RNAiFold2T** to design a temperature-regulated internal ribosomal entry site (IRES) element by ensuring the presence at high temperatures of the domain 5 stem-loop and downstream single-stranded pyrimidine (Py) tract, both located upstream of the functional initiation codon and both known to be important for IRES functionality [23]. At low temperatures, our thermo-IRES element is designed to adopt a conformation that down-regulates protein product by disrupting both the domain 5 stem-loop and sequestering the Py tract. We showed that our rationally designed thermo-IRES elements are functional, where the cap-independent translational efficiency is approximately 50% higher at 42°C than at 30°C; however, since the focus of this paper is primarily to

describe a new method, we have not taken steps to improve efficiency using error-prone mutagenesis and selection.

## 6 Funding

Research of the Clote Lab was supported by National Science Foundation grant DBI-1262439. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation. Research of the Martinez-Salas Lab was supported by the Spanish Ministry of Economy and Competitiveness (MINECO) [CSD2009-00080, BFU2011-25437, BFU2014-54564] and by an Institutional Grant from Fundación Ramón Areces.

## References

- [1] S. Altuvia, D. Kornitzer, D. Teff, and A. B. Oppenheim. Alternative mRNA structures of the cIII gene of bacteriophage lambda determine the rate of its translation initiation. *J. Mol. Biol.*, 210(2):265–280, November 1989.
- [2] M. Andronescu, A. Condon, H. H. Hoos, D. H. Mathews, and K. P. Murphy. Efficient parameter estimation for RNA secondary structure prediction. *Bioinformatics*, 23(13):i19–i28, July 2007.
- [3] W. Bae, B. Xia, M. Inouye, and K. Severinov. Escherichia coli CspA-family RNA chaperones are transcription antiterminators. *Proc. Natl. Acad. Sci. U.S.A.*, 97(14):7784–7789, July 2000.
- [4] S. Chowdhury, C. Ragaz, E. Kreuger, and F. Narberhaus. Temperature-controlled structural alterations of an RNA thermometer. *J. Biol. Chem.*, 278(48):47915–47921, November 2003.
- [5] A. Churkin, A. Avihoo, M. Shapira, and D. Barash. RNAThermsw: direct temperature simulations for predicting the location of RNA thermometers. *PLoS. One.*, 9(4):e94340, 2014.
- [6] A. Chursov, S. J. Kopetzky, G. Bocharov, D. Frishman, and A. Shneider. RNAtips: Analysis of temperature-induced changes of RNA secondary structure. *Nucleic. Acids. Res.*, 41(Web):W486–W491, July 2013.
- [7] K. E. Deigan, T. W. Li, D. H. Mathews, and K. M. Weeks. Accurate SHAPE-directed RNA structure determination. *Proc. Natl. Acad. Sci. U.S.A.*, 106(1):97–102, January 2009.
- [8] I. Dotu, J. A. Garcia-Martin, B. L. Slinger, V. Mechery, M. M. Meyer, and P. Clote. Complete RNA inverse folding: computational design of functional hammerhead ribozymes. *Nucleic. Acids. Res.*, 42(18):11752–11762, February 2015.



- [9] I. Dotu, G. Lozano, P. Clote, and E. Martinez-Salas. Using RNA inverse folding to identify IRES-like structural subdomains. *RNA. Biol.*, 10(12):1842–1852, December 2013.
- [10] N. Fernandez, O. Fernandez-Miragall, J. Ramajo, A. Garca-Sacristan, N. Bellora, E. Eyra, C. Briones, and E. Martinez-Salas. Structural basis for the biological relevance of the invariant apical stem in IRES-mediated translation. *Nucleic Acids Res.*, 39:8572–8585, 2011.
- [11] J. Fernandez-Chamorro, D. Pineiro, J. M. Gordon, J. Ramajo, R. Francisco-Velilla, M. J. Macias, and E. Martinez-Salas. Identification of novel non-canonical RNA-binding sites in Gemin5 involved in internal initiation of translation. *Nucleic. Acids. Res.*, 42(9):5742–5754, May 2014.
- [12] C. Flamm, I.L. Hofacker, S. Mauer-Stroh, P.F. Stadler, and M. Zehl. Design of multi-stable RNA molecules. *RNA*, 7:254–265, 2001.
- [13] J. A. Garcia-Martin, I. Dotu, and P. Clote. RNAiFold 2.0: a web server and software to design custom and Rfam-based RNA molecules. *Nucleic. Acids. Res.*, 43(W1):W513–W521, July 2015.
- [14] J.A. Garcia-Martin, P. Clote, and I. Dotu. RNAiFold: A constraint programming algorithm for RNA inverse folding and molecular design. *Journal of Bioinformatics and Computational Biology*, 11(2):1350001, 2012.
- [15] A. Hoynes-O’Connor, K. Hinman, L. Kirchner, and T. S. Moon. De novo design of heat-repressible RNA thermosensors in *E. coli*. *Nucleic. Acids. Res.*, 43(12):6166–6179, July 2015.
- [16] J. Johansson, P. Mandin, A. Renzoni, C. Chiaruttini, M. Springer, and P. Cossart. An RNA thermosensor controls expression of virulence genes in *Listeria monocytogenes*. *Cell*, 110(5):551–561, September 2002.
- [17] J. Kortmann and F. Narberhaus. Bacterial RNA thermometers: molecular zippers and switches. *Nat. Rev. Microbiol.*, 10(4):255–265, April 2012.
- [18] J. Kortmann, S. Szodrok, J. Rinnenthal, H. Schwalbe, and F. Narberhaus. Translation on demand by a simple RNA-based thermosensor. *Nucleic. Acids. Res.*, 39(7):2855–2868, April 2011.
- [19] R. Kuhn, N. Luz, and E. Beck. Functional analysis of the internal translation initiation site of foot-and-mouth disease virus. *J. Virol.*, 64(10):4625–4631, October 1990.
- [20] J. Lee and J.A. Kotov. Thermometer design at the nanoscale. *Nano Today*, 2(1):48–51, 2007.
- [21] R. Lorenz, S. H. Bernhart, C. Höner zu Siederdissen, H. Tafer, C. Flamm, P. F. Stadler, and I. L. Hofacker. Viennarna Package 2.0. *Algorithms. Mol. Biol.*, 6:26, 2011.

- [22] G. Lozano, N. Fernandez, and E. Martinez-Salas. Magnesium-dependent folding of a picornavirus IRES element modulates RNA conformation and eIF4G interaction. *FEBS J.*, 281(16):3685–3700, August 2014.
- [23] G. Lozano and E. Martinez-Salas. Structural insights into viral IRES-dependent translation mechanisms. *Curr. Opin. Virol.*, 12:113–120, June 2015.
- [24] G. Lozano, A. Trapote, J. Ramajo, X. Elduque, A. Grandas, J. Robles, E. Pedroso, and E. Martinez-Salas. Local RNA flexibility perturbation of the IRES element induced by a novel ligand inhibits viral RNA translation. *RNA. Biol.*, 12(5):555–568, 2015.
- [25] R. B. Lyngso, J. W. Anderson, E. Sizikova, A. Badugu, T. Hyland, and J. Hein. Frnakenstein: multiple target inverse RNA folding. *BMC. Bioinformatics*, 13:260, 2012.
- [26] E. P. Nawrocki, S. W. Burge, A. Bateman, J. Daub, R. Y. Eberhardt, S. R. Eddy, E. W. Floden, P. P. Gardner, T. A. Jones, J. Tate, and R. D. Finn. Rfam 12.0: updates to the RNA families database. *Nucleic Acids Res.*, 0(O):O, November 2014.
- [27] J. Neupert, D. Karcher, and R. Bock. Design of simple synthetic RNA thermometers for temperature-controlled gene expression in *Escherichia coli*. *Nucleic. Acids. Res.*, 36(19):e124, November 2008.
- [28] A. Nocker, T. Hausherr, S. Balsiger, N. P. Krstulovic, H. Hennecke, and F. Narberhaus. A mRNA-based thermosensor controls expression of rhizobial heat shock genes. *Nucleic. Acids. Res.*, 29(23):4800–4807, December 2001.
- [29] A. Ogawa. Rational design of artificial riboswitches based on ligand-dependent modulation of internal ribosome entry in wheat germ extract and their applications as label-free biosensors. *RNA.*, 17(3):478–488, March 2011.
- [30] D. Pineiro, N. Fernandez, J. Ramajo, and E. Martinez-Salas. Gemin5 promotes IRES interaction and translation control through its C-terminal region. *Nucleic. Acids. Res.*, 41(2):1017–1028, January 2013.
- [31] C. Reidys, P.F. Stadler, and P. Schuster. Generic properties of combinatorial maps: neutral networks of RNA secondary structures. *Bull Math Biol.*, 59(2):339–397, 1997.
- [32] J. S. Reuter and D. H. Mathews. RNAstructure: software for RNA secondary structure prediction and analysis. *BMC. Bioinformatics*, 11:129, 2010.

- [33] J. Rinnenthal, B. Klinkert, F. Narberhaus, and H. Schwalbe. Modulation of the stability of the Salmonella fourU-type RNA thermometer. *Nucleic. Acids. Res.*, 39(18):8258–8270, October 2011.
- [34] I. Rouzina and V. A. Bloomfield. Heat capacity effects on the melting of DNA. 1. General aspects. *Biophys. J.*, 77(6):3242–3251, December 1999.
- [35] A. Saragliadis, S. S. Krajewski, C. Rehm, F. Narberhaus, and J. S. Hartig. Thermozyms: Synthetic RNA thermometers based on ribozyme activity. *RNA. Biol.*, 10(6):1010–1016, June 2013.
- [36] Z. Torok, P. Goloubinoff, I. Horvath, N. M. Tsvetkova, A. Glatz, G. Balogh, V. Varvasovszki, D. A. Los, E. Vierling, J. H. Crowe, and L. Vigh. Synechocystis HSP17 is an amphitropic protein that stabilizes heat-stressed membranes and binds denatured proteins for subsequent chaperone-mediated refolding. *Proc. Natl. Acad. Sci. U.S.A.*, 98(6):3098–3103, March 2001.
- [37] D. H. Turner and D. H. Mathews. NNDB: the nearest neighbor parameter database for predicting stability of nucleic acid secondary structure. *Nucleic. Acids. Res.*, 38(Database):D280–D282, January 2010.
- [38] M. Wachsmuth, S. Findeiss, N. Weissheimer, P. F. Stadler, and M. Morl. De novo design of a synthetic riboswitch that regulates transcription termination. *Nucleic. Acids. Res.*, 41(4):2541–2551, February 2013.
- [39] T. Waldminghaus, L. C. Gaubig, and F. Narberhaus. Genome-wide bioinformatic prediction and experimental evaluation of potential RNA thermometers. *Mol. Genet. Genomics.*, 278(5):555–564, November 2007.
- [40] T. Waldminghaus, N. Heidrich, S. Brantl, and F. Narberhaus. Fouru: a novel type of RNA thermometer in Salmonella. *Mol. Microbiol.*, 65(2):413–424, July 2007.
- [41] T. Waldminghaus, J. Kortmann, S. Gesing, and F. Narberhaus. Generation of synthetic RNA-based thermosensors. *Biol. Chem.*, 389(10):1319–1326, October 2008.
- [42] M. Wieland and J. S. Hartig. RNA quadruplex-based modulation of gene expression. *Chem. Biol.*, 14(7):757–763, July 2007.
- [43] J. N. Zadeh, B. R. Wolfe, and N. A. Pierce. Nucleic acid sequence design via efficient ensemble defect optimization. *J. Comput. Chem.*, 32(3):439–452, February 2011.
- [44] K. Zarringhalam, M. M. Meyer, I. Dotu, J. H. Chuang, and P. Clote. Integrating chemical footprinting data into RNA secondary structure prediction. *PLoS. One.*, 7(10):e45160, 2012.
- [45] M. Zuker. Mfold web server for nucleic acid folding and hybridization prediction. *Nucleic Acids Res.*, 31(13):3406–3415, 2003.



```

      .....((((((.....)))))).. T: 30°
      ((((((.....))))))..... T: 10°
1  GNNNNNNNNNNNGNNNNNCNNNNNN
2  GNNNNNNNNNNNGNNNNNCNNNNNN
3  GNNNNNNNNNNCGGNNNNNCCGNNNN
4  GNNNNNNNNNGCGGNNNNNCCGCNNN
5  GNNNNNNNNCGCGGNNNNNCCGCNN
6  GNNNNNGNCGCGGCNNNNCCGCGNN
7  GNNNGNCGCGGCNNNNCCGCGNN
8  GNNUGNCGCGGCCANNCCGCGNN
9  GNGUGNCGCGGCCACNCCGCGNN
10 GNGUGNCGCGGCCACUCCGCGNN
11 GNGUGNCGCGGCCACUCCGCGGN
12 GNGUGNCGCGGCCACUCCGCGGA
      .....((((((.....)))))).. ←MFE EHwD check (30°C)
13 GGGUGNCGCGGCCACUCCGCGGA
14 GGGUGAGCCGCGGCCACUCCGCGGA
      ((((((.....))))))..... ←MFE EHwD check (10°C)
      .....((((((.....)))))).. ←MFE EHwD check (30°C)
      ((((((.....))))))..... ←MFE EHwD check (10°C)

```

■ Current assignment  
■ Constraint propagation

Figure 2: Example of CP search for target structures  $S_1$  (top) at temperature 30°C, and  $S_2$  (bottom) at temperature 10°C. Undetermined positions are assigned (red) in the following order: base pairs are instantiated in lines 1,2,3,4,5,13; unpaired positions in lines 6,7,8,9,10,14; a closing base pair in line 11; a dangling position in line 12. When the left nucleotide of a base pair is instantiated to C [resp. A], then propagation of the base pairing constraint reduces the domain of possible values of the right nucleotide of the base pair from  $\{A, C, G, U\}$  to the set  $\{G\}$  [resp.  $\{U\}$ ] – this happens in lines 1,6-9. For simplification, this example assumes that correct value assignments always occur at each search step.

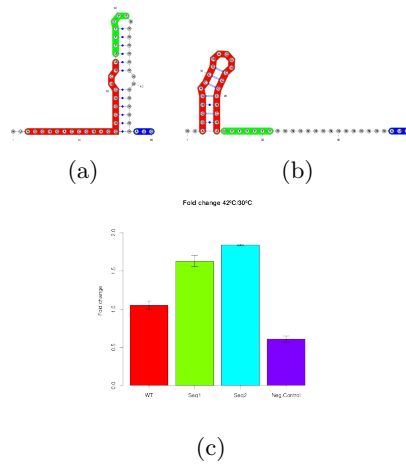


Figure 3: *(a,b)* Target structures  $S_1$  at temperature 30°C (a) and  $S_2$  at temperature 42°C (b) for domain 5 thermo-IRES element with added AUG codon (blue). IUPAC sequence constraints are determined from an alignment of 183 IRES sequences as shown in Figure 1 of [10]. Domain 5 stem-loop (positions 3-24 in red) and *unpaired* polypyrimidine tract (PPT positions 25-32 in green) are known to be *essential* for IRES activity [19]. Target structure  $S_1$  was designed to sequester the PPT at low temperature, thus creating a thermo-IRES which should be functional only at high temperature. *(c)* Ratio of normalized IRES activity at 42°C over that at 30°C for wild-type FMDV IRES, a negative control, and two thermosensors designed using RNAiFold2T.

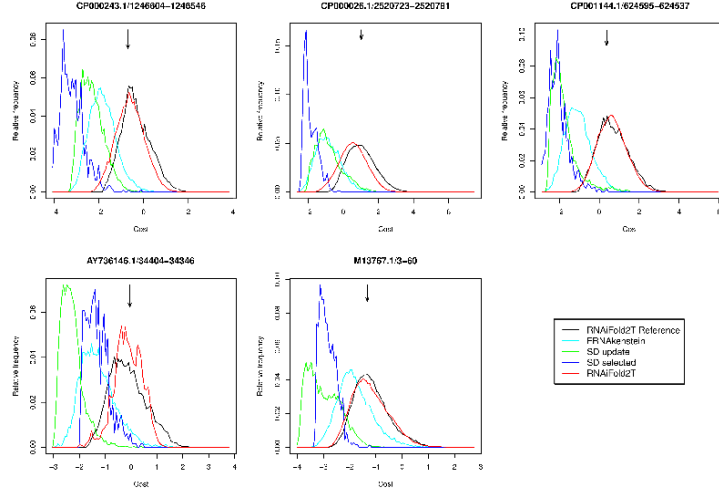


Figure 4: Relative frequency of the cost function optimized by **SwitchDesign**, for solutions returned by **RNAiFold2T**, **SwitchDesign** and **Frnakenstein**, given target structure  $S_1$  [resp.  $S_2$ ] at temperature  $T_1$  [resp.  $T_2$ ] for  $\lambda$  phage CIII thermoregulators from Rfam family RF01804. This figure is a more generous representation of the data from **SwitchDesign** and **Frnakenstein**, since all single point mutant solutions have been added to the raw output. (SI Figure 1 presents histograms for the raw output of these programs. The *reference* distribution for *RNAiFold2T Reference* (black curve), was produced by running *RNAiFold2T* for several days. Remaining curves are for **Frnakenstein** (light green), **SwitchDesign** (dark green and purple) and *RNAiFold2T* (red). Arrows indicate cost values for the real  $\lambda$  phage CIII thermoregulators from Rfam RF01804. Distribution for SD and FRNA without additional single point mutants shown in SI. SI Figures 1,3 show clearly that cost function values for Rfam sequences approximately equal the reference distribution mean.

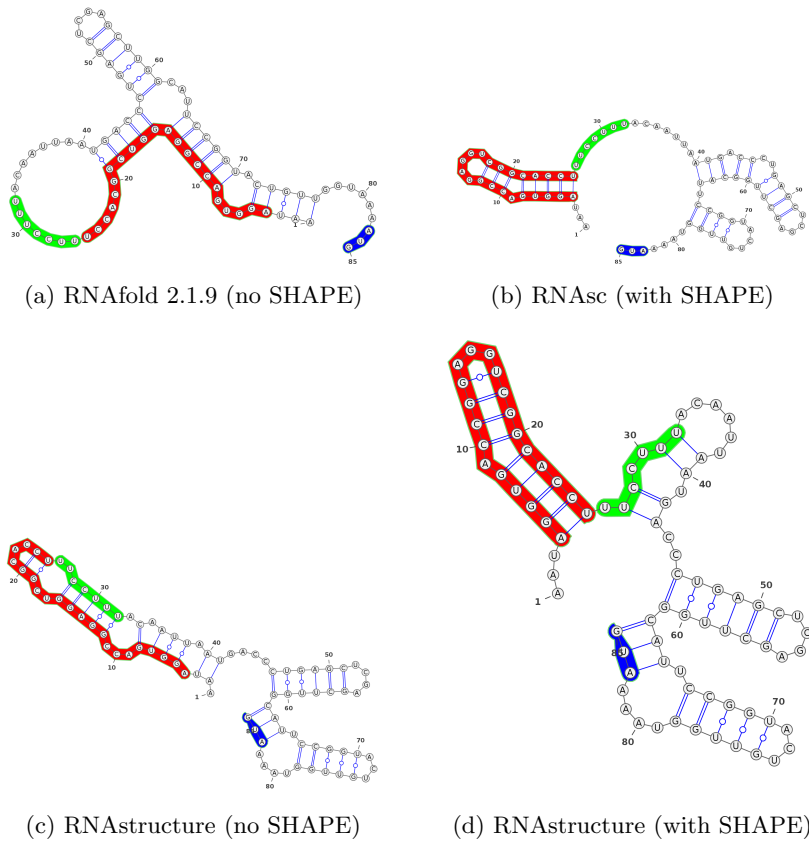


Figure 5: Secondary structure predictions of domain 5 of wild-type FMDV IRES, including the first functional AUG start codon, both with and without integration of probing data using selective 2'-hydroxyl acylation analyzed by primer extension (SHAPE). (a) **RNAfold 2.1.9**, which does not incorporate SHAPE data. (b) **RNAsc** [44], which penalizes deviations from SHAPE data for both base paired and loop regions. (c) **RNAstructure** [32] without SHAPE data. (d) **RNAstructure** [7], which penalizes deviations from SHAPE data only for base paired regions. The polypyrimidine tract (PPT) is unbound in only structures (a) from **RNAfold 2.1.9** and (b) from **RNAsc**, compatible with the full IRES structure appearing in Figure 1 of [10].