

SAT Coding Expl.: State Transitions

[0, 0, 2, 2, 2, 2, 1, 1, 1, 1]
[1, 1, 2, 2, 2, 2, 1, 0, 0, 1]
[1, 1, 2, 2, 0, 0, 1, 2, 2, 1]
[1, 0, 0, 2, 1, 2, 1, 2, 2, 1]
[1, 2, 1, 2, 1, 2, 1, 2, 0, 0]

- unknowns: $x_{t,p}$ where
 $t = \text{time}$, $p = \text{position}$
- obvious initial/final condition, transitions:

$$\bigwedge_t \bigvee_{a,b} \bigwedge_p \left(\begin{array}{l} p \notin \{a, a+1, b, b+1\} \Rightarrow x_{t,p} = x_{t+1,p} \\ \wedge x_{t,a} = x_{t+1,b} \wedge x_{t,a+1} = x_{t+1,b+1} \\ \wedge x_{t+1,a} = 0 \wedge x_{t+1,a+1} = 0 \\ \wedge x_{t,b} = 0 \wedge x_{t,b+1} = 0 \end{array} \right)$$

improve to $\bigwedge_t \exists r, s \in \{1, 2\} (\bigvee_a \bigwedge_p \dots) \wedge (\bigvee_b \bigwedge_p \dots)$

complete source code (100 lines) <http://dfa.imn.htwk-leipzig.de/cgi-bin/gitweb.cgi?p=biosat.git;a=blob;f=rewriting/C.hs;hb=HEAD>

Johannes Waldmann () Constraint Programming for Secondary Struct October 5, 2012 9 / 18

SAT encoding for Sec. Struc. Pred.

model: disjoint circular matchings in graphs

input: $G = (V, E)$ where $V = \text{positions in RNA string}$,
 $E = \text{set of all possible base pairs}$; number $k \in \mathbb{N}$
output: sequence M_1, \dots, M_k with $M_i \subseteq E$
such that

- $M := \bigcup_i M_i$ is a matching (each $v \in V$ is incident to at most one edge in M)
- each M_i is circular (no crossing edges w.r.t. the ordering on V)

each M_i is an edge set, thus a relation, thus a boolean matrix $M_i : V \times V \rightarrow \mathbb{B}$

the unknowns of the constraint system are the entries of these matrices.

Johannes Waldmann () Constraint Programming for Secondary Struct October 5, 2012 10 / 18

Related Work

- Unyanee Poolsap, Yuki Kato, and Tatsuya Akutsu: *Prediction of RNA secondary structure with pseudoknots using integer programming*, BMC Bioinformatics. 2009; 10(Suppl 1): S38.
- Ganesh et al.: *Lynx: A Programmatic SAT Solver for the RNA-Folding Problem*, SAT'12 using the direct encoding (l^4 clauses) for the non-crossing condition

Johannes Waldmann () Constraint Programming for Secondary Struct October 5, 2012 11 / 18

Encoding details

- union: $M = \bigcup_i M_i$
 $M(p, q) := \bigvee_i M_i(p, q)$
- possible base pairs $M \subseteq E$:
 $\bigwedge \{ \neg M(p, q) \mid (w[p], w[q]) \notin \{AU, UA, CG, GC, GU, UG\} \}$
- M is matching:
 $\bigwedge \{ \neg (M(p, q) \wedge M(q, r)) \mid p \neq r \}$
- M_i is circular (non-crossing):
 $\bigwedge \{ \neg (M_i(p, q) \wedge M_i(r, s)) \mid p < r < q < s \}$
- number of variables: $l^2 k$, formula size: $\Theta(l^4 k)$.

Johannes Waldmann () Constraint Programming for Secondary Struct October 5, 2012 12 / 18

Encoding for CYK parsing

... to reduce the l^4 formula size

use table (relation) T with specification

$T(p, q) \iff w[p..q]$ is correctly parenthesized:

- $T(p, p) \iff p \notin \text{domain } M \cup \text{range } M$
- $T(p, q) \iff (M(p, q) \wedge T(p+1, q-1)) \vee \bigvee_h T(p, h) \wedge T(h+1, q)$
- $M(p, q) \Rightarrow T(p, q), M(1, l)$

l^2 variables, l^3 formula size

source code: <http://dfa.imn.htwk-leipzig.de/cgi-bin/gitweb.cgi?p=biosat.git;a=blob;f=ssp/code/SSP/Graph/Encode.hs;hb=HEAD>

Note: cannot apply CYK to the original problem, since we need to guess the type of parentheses. (this parsing problem is NP-hard)

Johannes Waldmann () Constraint Programming for Secondary Struct October 5, 2012 13 / 18

Encoding Numeric Valuation

For valuation of (M_1, \dots, M_k) , consider *stacks* (groups of parallel edges in $M = \bigcup_i M_i$)

- Define $S : V \rightarrow \mathbb{B}$ by
 $S(p) := \bigvee_q M(p, q) \wedge M(p+1, q-1)$,
- count number of 1 in $(S(1), \dots, S(l))$ by repeated binary addition (using half adder/full adder circuits represented as constraint systems)
- compare with a given bound
 $v \geq B \iff \exists d : v = B + d$

Johannes Waldmann () Constraint Programming for Secondary Struct October 5, 2012 14 / 18

Solving the Optimization Problem

- write the constraint system $C(P, S, V) = \text{“}S \text{ is an admissible solution for problem } P \text{ with value } \geq V\text{”}$
- to find $\max\{V \mid \exists S : C(P, S, V)\}$, determine a finite feasible range for V (e.g., $0 \dots \text{length of input}$)
- use iteration $V = 0, 1, 2, \dots$ or bisection $V = l/2, 3l/4, \dots$

Johannes Waldmann () Constraint Programming for Secondary Struct October 5, 2012 15 / 18

Prototype Implementation

of Secondary Structure Prediction with fixed number of parenthesis types
is proof-of-concept, as a basis for experimentation:

- source: `git`:
[//dfa.imn.htwk-leipzig.de/srv/git/biosat](http://dfa.imn.htwk-leipzig.de/srv/git/biosat)
- using Haskell library `Satchmo`
<https://github.com/jwaldmann/satchmo> to generate SAT constraint system and decode result
- solver: <https://github.com/niklasso/minisat>

Johannes Waldmann () Constraint Programming for Secondary Struct October 5, 2012 16 / 18

Program Inversion

Constraint system $C(P, S, V) =$

“ S is an admissible solution for problem P
with value $\geq V$ ”

can be used for:

- ▶ given P, V , determine S
e.g., RNA parsing (sec. struct. pred.)
- ▶ given S, V , determine P
e.g., RNA design

Conclusion/Claims

- ▶ constraint programming is *easy*: especially for non-programmers, since it is *declarative*
- ▶ constraint programming is *powerful*:
use generic *domain-specific* solver for *application-specific* program/problem
- ▶ constraint programming is *flexible*:
easily add/remove/change/invert constraints
(much easier than change an application-specific algorithm)
- ▶ write the constraint program in an EDSL
(*embedded domain specific language*)
that takes care of encoding and decoding