

Constructing Set Constraints for ReScript

서울대학교 전기·정보공학부 2018-12602 이준협

1 Definition of set expressions

| | | | |
|------|-------|--------------------|---------------------------------------------------------------------------------------------------|
| se | $::=$ | \emptyset | <i>empty set</i> |
| | | $-$ | <i>maximum set</i> |
| | | $()$ | <i>unit</i> |
| | | n | <i>integer</i> |
| | | b | <i>boolean</i> |
| | | $\lambda x.e$ | <i>function</i> |
| | | loc | <i>memory location</i> |
| | | V_e | <i>set variable corresponding to the possible values of e</i> |
| | | P_e | <i>set variable corresponding to the possible exn packets of e</i> |
| | | $body_V(se)$ | <i>values that se can spit out when se is applied to something</i> |
| | | $body_P(se)$ | <i>exn packets that se can spit out when se is applied to something</i> |
| | | $par(se)$ | <i>values that can be a parameter for se</i> |
| | | κ | <i>constructor</i> |
| | | l | <i>field of a record</i> |
| | | $con(\kappa, se)$ | <i>construct</i> |
| | | $exn(\kappa, se)$ | <i>exception</i> |
| | | $fld(se, l)$ | <i>contents of the field l of a record se</i> |
| | | $cnt(se)$ | <i>contents of a reference</i> |
| | | $bop(se, se)$ | <i>binary operators, where $bop \in \{+, -, \times, \div, =, <, >\}$</i> |
| | | $f_{(i)}^{-1}(se)$ | <i>projection onto the i-th argument of f</i> |
| | | $se \cup se$ | <i>union</i> |
| | | $se \cap se$ | <i>intersection</i> |
| | | \overline{se} | <i>complement</i> |

The definition of the conditional set expression needs clarification.

$$se_1 \Rightarrow se_2 := \begin{cases} \emptyset & (se_1 = \emptyset) \\ se_2 & (o.w) \end{cases}$$

The conditional set expression is a naive approximation for pattern matching. Consider the case when we want to match an expression against the record pattern with fields x and y . The constraints describing the record $r = \{x = 1, y = 2\}$ are $1 \subseteq fld(V_r, x) \wedge 2 \subseteq fld(V_r, y)$. To pattern-match r against $\{x, y\}$, we want the fields x and y of V_r to be nonempty. Thus, the value of “match r with $\{x, y\} \rightarrow e$ ” is $fld(V_r, x) \Rightarrow (fld(V_r, y) \Rightarrow V_e)$.

The conditional set expression can also be used to define conditional set constraints [1].

$$se \Rightarrow \bigwedge_{i=1}^n X_i \subseteq Y_i := \bigwedge_{i=1}^n (se \Rightarrow X_i) \subseteq Y_i$$

2 Constructing set constraints

Now we are in a position to define constraint construction rules for our ReScript-like language. Hopefully this would be reasonably fast when implemented and be accurate enough...

$$\begin{array}{c}
 \text{[UNIT, INT, BOOL]} \quad \frac{}{\triangleright c : V_e \supseteq c} \quad c = (), n, b \\
 \\
 \text{[APP]} \quad \frac{\triangleright e_1 : C_1 \quad \triangleright e_2 : C_2}{\triangleright e_1 e_2 : (V_e \supseteq body_V(V_{e_1})) \wedge (P_e \supseteq (body_P(V_{e_1}) \cup P_{e_1} \cup P_{e_2})) \wedge (par(V_{e_1}) \supseteq V_{e_2}) \wedge C_1 \wedge C_2} \\
 \\
 \text{[FN]} \quad \frac{\triangleright e' : C'}{\triangleright \lambda x.e' : (V_e \supseteq \lambda x.e') \wedge (body_V(V_e) \supseteq V_{e'}) \wedge (body_P(V_e) \supseteq P_{e'}) \wedge (par(V_e) \subseteq V_x) \wedge C'}
 \end{array}$$

$$\begin{array}{c}
\text{[LET]} \frac{\triangleright e_1 : C_1 \quad \triangleright e_2 : C_2}{\triangleright \text{let } x = e_1 \text{ in } e_2 : (V_x \supseteq V_{e_1}) \wedge (V_e \supseteq V_{e_2}) \wedge (P_e \supseteq P_{e_1} \cup P_{e_2}) \wedge C_1 \wedge C_2} \\
\\
\text{[BOP]} \frac{\triangleright e_1 : C_1 \quad \triangleright e_2 : C_2}{\triangleright e_1 \text{ bop } e_2 : (V_e \supseteq \text{bop}(V_{e_1}, V_{e_2})) \wedge (P_e \supseteq P_{e_1} \cup P_{e_2}) \wedge \underbrace{(P_e \supseteq (V_{e_2} \cap 0 \Rightarrow \text{exn}(\text{Divide_by_zero}, ())))}_{\text{when bop is } \div} \wedge C_1 \wedge C_2} \\
\\
\text{[CON]} \frac{\triangleright e' : C'}{\triangleright \text{con } \kappa e' : (V_e \supseteq \text{con}(\kappa, V_{e'})) \wedge (P_e \supseteq P_{e'}) \wedge C'} \\
\\
\text{[EXN]} \frac{\triangleright e' : C'}{\triangleright \text{exn } \kappa e' : (V_e \supseteq \text{exn}(\kappa, V_{e'})) \wedge (P_e \supseteq P_{e'}) \wedge C'} \\
\\
\text{[DECON]} \frac{\triangleright e' : C'}{\triangleright \text{decon } \kappa e' : (V_e \supseteq \underbrace{\text{con}_{(2)}^{-1}(V_{e'} \cap \text{con}(\kappa, _))}_{\text{filter}} \cup \underbrace{\text{exn}_{(2)}^{-1}(V_{e'} \cap \text{exn}(\kappa, _))}_{\text{filter}}) \wedge (P_e \supseteq P_{e'}) \wedge C'} \quad \text{as in SML-NJ's Lambda} \\
\\
\text{[RECORD]} \frac{\triangleright e_i : C_i \ (1 \leq i \leq n)}{\triangleright \{l_1 = e_1, \dots, l_n = e_n\} : \bigwedge_{i=1}^n (\text{fld}(V_e, l_i) \supseteq V_{e_i}) \wedge (P_e \supseteq \bigcup_{i=1}^n P_{e_i}) \wedge \bigwedge_{i=1}^n C_i} \\
\\
\text{[FIELD]} \frac{\triangleright e' : C'}{\triangleright e'.l : (V_e \supseteq \text{fld}(V_{e'}, l)) \wedge (P_e \supseteq P_{e'}) \wedge C'} \\
\\
\text{[REF]} \frac{\triangleright e' : C'}{\triangleright \text{ref } e' : (V_e \supseteq \text{loc}) \wedge (\text{cnt}(\text{loc}) \supseteq V_{e'}) \wedge (P_e \supseteq P_{e'}) \wedge C'} \quad \text{new loc} \\
\\
\text{[UPDATE]} \frac{\triangleright e_1 : C_1 \quad \triangleright e_2 : C_2}{\triangleright e_1 := e_2 : (\text{cnt}(V_{e_1}) \supseteq V_{e_2}) \wedge (P_e \supseteq P_{e_1} \cup P_{e_2}) \wedge C_1 \wedge C_2} \\
\\
\text{[BANG]} \frac{\triangleright e' : C'}{\triangleright !e' : (V_e \supseteq \text{cnt}(V_{e'})) \wedge (P_e \supseteq P_{e'}) \wedge C'} \\
\\
\text{We define an auxiliary function for generating constraints out of pattern matching. If we want to figure out the constraint for the value } X \text{ of "match } Y \text{ with } p \rightarrow e": \\
\\
\text{case}(X, Y, p, e) := \begin{cases} (Y \subseteq V_x) \wedge (V_e \subseteq X) & (p = x) \\ (\text{fld}(Y, l_1) \Rightarrow \dots \Rightarrow \text{fld}(Y, l_n) \Rightarrow V_e) \subseteq X & (p = \{l_i\}_{i=1}^n) \\ (\kappa \cap (\text{con}_{(1)}^{-1}(Y) \cup \text{exn}_{(1)}^{-1}(Y))) \Rightarrow V_e \subseteq X & (p = \kappa) \\ (Y \cap c) \Rightarrow V_e \subseteq X & (p = \text{constant}) \\ V_e \subseteq X & (p = _) \end{cases} \\
\\
\text{[CASE]} \frac{\triangleright e' : C' \quad \triangleright e_i : C_i \ (1 \leq i \leq n)}{\triangleright \text{case } e' \ (p_i \rightarrow e_i)_{i=1}^n : \bigwedge_{i=1}^n \text{case}(V_{e'}, V_{e'}, p_i, e_i) \wedge (P_e \supseteq \bigcup_{i=1}^n P_{e_i}) \wedge C' \wedge \bigwedge_{i=1}^n C_i} \\
\\
\text{[HANDLE]} \frac{\triangleright e' : C' \quad \triangleright e_i : C_i \ (1 \leq i \leq n)}{\triangleright \text{handle } e' \ (p_i \rightarrow e_i)_{i=1}^n : (P_e \supseteq (P_{e'} \cap \bigcap_{i=1}^n \overline{\text{exn}(p_i, _)})) \cup \bigcup_{i=1}^n P_{e_i} \wedge (V_e \supseteq V_{e'} \cup \bigcup_{i=1}^n V_{e_i}) \wedge C' \wedge \bigwedge_{i=1}^n C_i} \\
\\
\text{[RAISE]} \frac{\triangleright e' : C'}{\triangleright \text{raise } e' : (P_e \supseteq V_{e_1} \cup P_{e_1}) \wedge C'} \\
\\
\text{[FOR]} \frac{\triangleright e_1 : C_1 \quad \triangleright e_2 : C_2 \quad \triangleright e_3 : C_3}{\triangleright \text{for } x \ e_1 \ e_2 \ e_3 : \begin{aligned} & (> (V_x, V_{e_1}) \cap \text{true} \Rightarrow () \subseteq V_e) \\ & \wedge (> (V_x, V_{e_1}) \cap \text{false} \Rightarrow (C_3 \wedge (V_x \supseteq +(V_x, 1)) \wedge (P_e \supseteq P_{e_3}))) \\ & \wedge (V_x \supseteq V_{e_1}) \wedge (P_e \supseteq P_{e_1} \cup P_{e_2}) \wedge C_1 \wedge C_2 \end{aligned}} \\
\\
\text{[WHILE]} \frac{\triangleright e_1 : C_1 \quad \triangleright e_2 : C_2}{\triangleright \text{while } e_1 \ e_2 : \begin{aligned} & (V_{e_1} \cap \text{true}) \Rightarrow (C_2 \wedge (P_e \supseteq P_{e_2})) \\ & \wedge ((V_{e_1} \cap \text{false}) \Rightarrow () \subseteq V_e) \wedge (P_e \supseteq P_{e_1}) \wedge C_1 \end{aligned}}
\end{array}$$

3 Conditional expressions : a good approximation?

Case expressions are the cause for inaccuracy in approximating the program states. Take a look at the for loop.

```

let for x = match x > e2 with
    | true -> ()
    | false -> e3; for (x + 1)
in
for e1

```

The above program is equivalent to `for x = e1 to e2 do e3 done`. Translating this to set constraints is difficult as case statements partition the program states. That is, “`e3; for (x + 1)`” is evaluated under the constraint that $\triangleright (V_x, V_{e_2}) \subseteq \text{false}$ and “`()`” is evaluated under the constraint that $\triangleright (V_x, V_{e_2}) \subseteq \text{true}$.

These constraints obviously cannot be and-ed together, as the partitions are mutually disjoint. In other words, $x > e_2$ cannot evaluate to both true and false at the same time. Thus, it is straightforward that each set expression must have different “versions” of itself in each partition. Each case statement creates a “parallel world” where some set constraint becomes true.

The conditional expression approximation is very coarse, as can be observed in the [FOR] rule in section 2. Assume that $e_1 = 1$, $e_2 = 5$, $e_3 = ()$. Since $1 \subseteq V_x$, $\text{false} \subseteq \triangleright (V_x, 5) \subseteq \triangleright (V_x, V_{e_2})$, so the conditional expression is activated. Then $V_x \supseteq \triangleright (V_x, 1)$ becomes valid. Then the least model must map V_x to $\{x \in \mathbb{Z} \mid x \geq 1\}$. This overshoots the possible values that x may have, which is $\{1, 2, 3, 4, 5, 6\}$.

Why does this happen? It is because in the condition $V_x \supseteq \triangleright (V_x, 1)$, the V_x -s in different sides are in different partitions. The V_x on the left-hand side can be in either partition, as it is not matched against any pattern yet. However, the V_x on the right-hand side must be in the partition where $x > e_2$ is matched against `false`.

Then when are conditional expressions successful? In the case when $e_1 = 5$ and $e_2 = 1$, for example, the conditional expression eliminates the treacherous condition $V_x \supseteq \triangleright (V_x, 1)$. In the context of exception analysis, $P_e \supseteq P_{e_3}$ is not considered at all, so the analysis is a bit more accurate. That is, only when the program doesn’t lay a foot on the wrong path is when conditional expressions succeed.

The obvious problem is that programs are often written so that all cases are reachable, maybe except for the case when some value is divided by a nonzero constant. Then the question is whether conditional expressions in the [CASE, FOR, WHILE] rules actually decrease false alarms. When implementing the analyzer it might not be a bad idea to eliminate the conditional expressions.

If the conditional expressions are eliminated, then the case function changes to

$$\text{case}(X, Y, p, e) := \begin{cases} (Y \subseteq V_x) \wedge (V_e \subseteq X) & (p = x) \\ V_e \subseteq X & (o.w) \end{cases}$$

and the [FOR, WHILE] rules become

$$\begin{aligned} \text{[FOR]} & \frac{\triangleright e_1 : C_1 \quad \triangleright e_2 : C_2 \quad \triangleright e_3 : C_3}{\triangleright \text{for } x \ e_1 \ e_2 \ e_3 : \begin{aligned} & (V_e \supseteq ()) \wedge (V_x \supseteq \triangleright (V_x, 1) \cup V_{e_1}) \\ & \wedge (P_e \supseteq P_{e_1} \cup P_{e_2} \cup P_{e_3}) \wedge C_1 \wedge C_2 \wedge C_3 \end{aligned}} \\ \text{[WHILE]} & \frac{\triangleright e_1 : C_1 \quad \triangleright e_2 : C_2}{\triangleright \text{while } e_1 \ e_2 : (V_e \supseteq ()) \wedge (P_e \supseteq P_{e_1} \cup P_{e_2}) \wedge C_1 \wedge C_2} \end{aligned}$$

TODO:

1. How to solve the set constraints (first in the case when there are no conditional expressions) : based on the work of Nevin Heintze
2. Maybe, how to formulate the “parallel worlds” approach?

References

- [1] Alexander Aiken. “Introduction to set constraint-based program analysis”. In: *Science of Computer Programming* 35.2 (1999), pp. 79–111. ISSN: 0167-6423. DOI: [https://doi.org/10.1016/S0167-6423\(99\)00007-6](https://doi.org/10.1016/S0167-6423(99)00007-6).