

# Enum Translation with Type Constraints

Minwook Lee  
ryan-lee@kaist.ac.kr

## I. NOTATION AND ASSUMPTIONS

- We consider a type as a **set of values**. Therefore, `int` is equivalent to  $\mathbb{Z}$ .
  - Every set operation can be applied to types as well. For example,  $\text{int} \cup \epsilon = \mathbb{Z}$  and  $\epsilon \subseteq \mathbb{Z}$ .
  - Every unary or binary operation can be applied to types as well. For example:

$$\tau_1 \text{ binop } \tau_2 = \{v_1 \text{ binop } v_2 \mid v_1 \in \tau_1 \wedge v_2 \in \tau_2\}.$$

- Type constraint of the identifier  $x$  is of the form  $\llbracket x \rrbracket = \langle \lfloor x \rfloor, \lceil x \rceil \rangle$ .
  - $\lfloor x \rfloor$  is a set of possible values that can be assigned to the identifier.
  - $\lceil x \rceil$  is the required type for  $x$  based on the type annotation and its usage context.
  - Every set operation can be applied to type constraints as well. For example:
    - \*  $\llbracket x \rrbracket \subseteq \llbracket y \rrbracket$  iff  $\lfloor x \rfloor \subseteq \lfloor y \rfloor \wedge \lceil x \rceil \subseteq \lceil y \rceil$
    - \*  $\llbracket x \rrbracket \cup \llbracket y \rrbracket = \langle \lfloor x \rfloor \cup \lfloor y \rfloor, \lceil x \rceil \cup \lceil y \rceil \rangle$
    - \*  $\llbracket x \rrbracket \text{ binop } \llbracket y \rrbracket = \langle \lfloor x \rfloor \text{ binop } \lfloor y \rfloor, \lceil x \rceil \text{ binop } \lceil y \rceil \rangle$

## II. BASIC MIR

### A. Abstract Syntax

$$\begin{aligned} n &\in \mathbb{Z} \\ x, y &\in Id \\ f &:= \overline{B : b; \overline{x : \tau}} \\ b &:= \overline{s; t} \\ s &:= x = n \mid x = y \mid x = \text{unop } y \mid x = y_1 \text{ binop } y_2 \\ t &:= \text{goto } B \mid \text{switch } x \text{ } \overline{n : B} \mid \text{return} \\ \epsilon &\in Enum \\ \tau &:= \text{int} \mid \epsilon \end{aligned}$$

Fig. 1. Basic MIR Abstract Syntax

The abstract syntax of the basic MIR is shown in Fig 1. In MIR, a function  $f$  is represented with the type `Body`<sup>1</sup>, which denotes a control flow graph (CFG) consisting of multiple basic blocks. Each basic block  $b$  is represented with the type `BasicBlockData`<sup>2</sup> and consists of multiple statements followed by a single terminator. In a basic block, statements do not have any control flow effect, and only terminators affect the control flow. A basic block is associated with a unique identifier  $B$ , represented with the type `BasicBlock`<sup>3</sup>. A statement  $s$  has the type `Statement`<sup>4</sup>. While various kinds of statements exist in MIR, we only need to consider assignments. The left-hand side of an assignment is a variable  $x$ , and the right-hand side is either an integer  $n$ , a variable, a unary operation, or a binary operation. A terminator  $t$  has the type `Terminator`<sup>5</sup> and is either a jump to another basic block, a switch, or a return. When returning, the return value is always the value stored in the variable named `_0`. If a function has  $n$  parameters, their values are stored in the variables `_1`, `_2`, ..., `_n`. Note that a function declares the type of each variable as well. A type  $\tau$  is either `int` or a C `enum` type  $\epsilon$ .

### B. Generating Type Constraints

The rules for type constraints in basic MIR are shown in Fig 2.

<sup>1</sup>[https://doc.rust-lang.org/beta/nightly-rustc/rustc\\_middle/mir/struct.Body.html](https://doc.rust-lang.org/beta/nightly-rustc/rustc_middle/mir/struct.Body.html)

<sup>2</sup>[https://doc.rust-lang.org/beta/nightly-rustc/rustc\\_middle/mir/struct.BasicBlockData.html](https://doc.rust-lang.org/beta/nightly-rustc/rustc_middle/mir/struct.BasicBlockData.html)

<sup>3</sup>[https://doc.rust-lang.org/beta/nightly-rustc/rustc\\_middle/mir/struct.BasicBlock.html](https://doc.rust-lang.org/beta/nightly-rustc/rustc_middle/mir/struct.BasicBlock.html)

<sup>4</sup>[https://doc.rust-lang.org/beta/nightly-rustc/rustc\\_middle/mir/struct.Statement.html](https://doc.rust-lang.org/beta/nightly-rustc/rustc_middle/mir/struct.Statement.html)

<sup>5</sup>[https://doc.rust-lang.org/beta/nightly-rustc/rustc\\_middle/mir/struct.Terminator.html](https://doc.rust-lang.org/beta/nightly-rustc/rustc_middle/mir/struct.Terminator.html)

For the function  $f$ :

$$\overline{B} : \overline{b}; \overline{x} : \overline{\tau} \quad : \quad \lceil \overline{x} \rceil \subseteq \overline{\tau}$$

For the statement  $s$ :

$$\begin{aligned} x = n & \quad : \quad n \in \lfloor x \rfloor \\ x = y & \quad : \quad \llbracket y \rrbracket \subseteq \llbracket x \rrbracket \\ x = \text{unop } y & \quad : \quad \text{unop } \lfloor y \rfloor \subseteq \lfloor x \rfloor \\ x = y_1 \text{ binop } y_2 & \quad : \quad \lfloor y_1 \rfloor \text{ binop } \lfloor y_2 \rfloor \subseteq \lfloor x \rfloor \end{aligned}$$

Fig. 2. Basic MIR Type Constraints

### C. Resolving Type Constraints

We resolve the type constraints with the following principles:

- We always choose the smallest set for  $\lfloor x \rfloor$  that satisfies the type constraints.
  - This increases the coverage where the enum translation can be applied.
- We always choose the largest set for  $\lceil x \rceil$  that satisfies the type constraints.
  - This prevents the redundant enum translation.

### D. Examples

Assume  $\epsilon = \{0, 1\}$ .

1) *Simple case*: The following is an example of a basic MIR function:

$$x : \epsilon; y : \mathbb{Z}; \{y = 1; x = y; \text{return}\}$$

We can derive the following conditions:

$$\begin{aligned} x : \epsilon & \quad : \quad \lceil x \rceil \subseteq \epsilon \\ y : \mathbb{Z} & \quad : \quad \lceil y \rceil \subseteq \mathbb{Z} \\ y = 1 & \quad : \quad 1 \in \lfloor y \rfloor \\ x = y & \quad : \quad \llbracket y \rrbracket \subseteq \llbracket x \rrbracket \end{aligned}$$

By resolving the type constraints we have:

$$\llbracket x \rrbracket = \langle \{1\}, \epsilon \rangle \quad \llbracket y \rrbracket = \langle \{1\}, \epsilon \rangle.$$

Now we deduce that  $\lfloor x \rfloor \subseteq \lceil x \rceil$  and  $\lfloor y \rfloor \subseteq \lceil y \rceil$ . Therefore, we now safely modify the function to:

$$\dots; y : \epsilon; \dots$$

Also, the translation doesn't introduce the unsafe code since all required types are valid Rust enum types.

2) *Solving more cases with type constraints*: The following is an example of a basic MIR function showing how deriving value constraints would work in case where subtype inference could not be applied:

$$x : \epsilon; y : \mathbb{Z}; \{y = -1; x = -y; \text{return}\}$$

We can derive the following conditions:

$$\begin{aligned} x : \epsilon & \quad : \quad \lceil x \rceil \subseteq \epsilon \\ y : \mathbb{Z} & \quad : \quad \lceil y \rceil \subseteq \mathbb{Z} \\ y = -1 & \quad : \quad -1 \in \lfloor y \rfloor \\ x = -y & \quad : \quad -\lfloor y \rfloor \subseteq \lfloor x \rfloor \end{aligned}$$

By resolving the type constraints we have:

$$\llbracket x \rrbracket = \langle \{1\}, \epsilon \rangle \quad \llbracket y \rrbracket = \langle \{-1\}, \mathbb{Z} \rangle.$$

Now we deduce that all  $\lfloor x \rfloor \subseteq \lceil x \rceil$  and  $\lfloor y \rfloor \subseteq \lceil y \rceil$ . Therefore, we can translate the integer type to the Rust enum type. However, this procedure would introduce unsafe code on  $x = -y$  while translating the type of  $x$  to be  $\epsilon$ .

$$\begin{array}{lcl}
s & := & \dots \mid x = \&y \mid x = *y \mid *x = y \\
\tau & := & \dots \mid * \tau
\end{array}$$

Fig. 3. MIR Abstract Syntax with Pointers

### III. ADDING POINTERS

#### A. Abstract Syntax

The abstract syntax of MIR with pointers is shown in Fig 3.

We denote the abstract cell  $c \in \text{Cell}$  as the set of all locations while  $\llbracket c \rrbracket = \langle \lfloor c \rfloor, \lceil c \rceil \rangle$  denotes the type constraint of the value which is located on  $c$ . The pointer type  $*\tau$  denotes the set of cells that can be dereferenced to a value of type  $\tau$ . If  $x$  can be dereferenced to a value of type  $\tau_1$  and  $\tau_1 \subseteq \tau_2$ , then  $x$  can also be dereferenced to a value of type  $\tau_2$ . This means that  $*\tau_1 \subseteq *\tau_2$ .

#### B. Generating Type Constraints

We choose an Andersen-style pointer analysis, since we do not expect real-world C code to use extensive pointer operations that would require flow-sensitive analysis for enum types.

The rules for type constraints in MIR with pointers are shown in Fig 4.

For the statement  $s$ :

$$\begin{array}{lcl}
& & \dots \\
x = \&y & : & \forall \tau, \quad y \subseteq \lfloor x \rfloor \quad \wedge \quad \tau \subseteq \lceil y \rceil \Rightarrow *\tau \subseteq \lceil x \rceil \\
x = *y & : & \forall \tau \forall c, \quad c \in \lfloor y \rfloor \Rightarrow \llbracket c \rrbracket \subseteq \llbracket x \rrbracket \quad \wedge \quad *\tau \subseteq \lceil y \rceil \Rightarrow \tau \subseteq \lceil x \rceil \\
*x = y & : & \forall \tau \forall c, \quad c \in \lfloor x \rfloor \Rightarrow \llbracket y \rrbracket \subseteq \llbracket c \rrbracket \quad \wedge \quad \tau \subseteq \lceil y \rceil \Rightarrow *\tau \subseteq \lceil x \rceil
\end{array}$$

Fig. 4. MIR Type Constraints with Pointers

#### C. Examples

Assume  $\epsilon = \{0, 1\}$ . The following is an example of a basic MIR function with pointers:

$$x : *\mathbb{Z}; y : \mathbb{Z}; z : \epsilon; \{y = 1; x = \&y; z = *x; \text{return}\}$$

Each statement and variable declaration generates the following type constraints:

$$\begin{array}{lcl}
x : *\mathbb{Z} & : & \lceil x \rceil \subseteq *\mathbb{Z} \\
y : \mathbb{Z} & : & \lceil y \rceil \subseteq \mathbb{Z} \\
z : \epsilon & : & \lceil z \rceil \subseteq \epsilon \\
y = 1 & : & 1 \in \lfloor y \rfloor \\
x = \&y & : & \forall \tau, \quad y \subseteq \lfloor x \rfloor \quad \wedge \quad \tau \subseteq \lceil y \rceil \Rightarrow *\tau \subseteq \lceil x \rceil \\
z = *x & : & \forall \tau \forall c, \quad c \in \lfloor x \rfloor \Rightarrow \llbracket c \rrbracket \subseteq \llbracket z \rrbracket \quad \wedge \quad *\tau \subseteq \lceil x \rceil \Rightarrow \tau \subseteq \lceil z \rceil
\end{array}$$

Now we resolve the type constraints as follows:

$$\begin{aligned}
\lfloor y \rfloor &= \{1\} \\
\lfloor x \rfloor &= \{y\} \\
\llbracket y \rrbracket &\subseteq \llbracket z \rrbracket \\
\tau \subseteq \lceil y \rceil &\Rightarrow *\tau \subseteq \lceil x \rceil \Rightarrow \tau \subseteq \lceil z \rceil \subseteq \epsilon \\
\therefore \lceil x \rceil &\subseteq *\epsilon \quad \wedge \quad \lceil y \rceil \subseteq \epsilon \quad \wedge \quad \lceil z \rceil \subseteq \epsilon
\end{aligned}$$

By resolving the type constraints we have:

$$\llbracket x \rrbracket = \langle \{y\}, *\epsilon \rangle \quad \llbracket y \rrbracket = \langle \{1\}, \epsilon \rangle \quad \llbracket z \rrbracket = \langle \{1\}, \epsilon \rangle.$$

#### IV. ADDING STRUCTS AND ARRAYS

##### A. Abstract Syntax

The abstract syntax of MIR with structs and arrays is shown in Fig 5. The array type  $\tau[]$  denotes the set of arrays whose elements are of type  $\tau$ . i.e., if a variable  $x$  has the type  $\tau[]$  and  $y$  has the type  $\mathbb{Z}$ , then  $x[y] \subseteq \tau$ .

$$\begin{aligned} F &\in Field \\ S &\in Struct \\ s &:= \dots \mid x.F = y \mid x = y.F \mid x[x'] = y \mid x = y[y'] \\ \tau &:= \dots \mid S \mid \tau[] \\ p &:= \overline{S\{F : \tau\}}; \bar{f} \end{aligned}$$

Fig. 5. MIR Abstract Syntax with Structs and Arrays

##### B. Generating Type Constraints

The rules for type constraints in MIR with structs and arrays are shown in Fig 6.

For the struct definition  $S\{\overline{F : \tau}\}$ :

$$S\{\overline{F : \tau}\} : \forall F \in S, \Rightarrow [S.F] \subseteq \tau$$

For the statement  $s$ :

$$\begin{aligned} &\dots \\ x : S &: \forall F \in S, \llbracket x.F \rrbracket \subseteq \llbracket S.F \rrbracket \\ x.F = y &: \llbracket y \rrbracket \subseteq \llbracket x.F \rrbracket \\ x = y.F &: \llbracket y.F \rrbracket \subseteq \llbracket x \rrbracket \\ x[x'] = y &: \tau \subseteq \llbracket y \rrbracket \Rightarrow \tau[] \subseteq \llbracket x \rrbracket \quad \wedge \quad \forall \tau, \tau \subseteq \llbracket y \rrbracket \Rightarrow \tau[] \subseteq \llbracket x \rrbracket \\ x = y[y'] &: \tau[] \subseteq \llbracket y \rrbracket \Rightarrow \tau \subseteq \llbracket x \rrbracket \quad \wedge \quad \forall \tau, \tau[] \subseteq \llbracket y \rrbracket \Rightarrow \tau \subseteq \llbracket x \rrbracket \end{aligned}$$

Fig. 6. MIR Type Constraints with Structs and Arrays

#### V. ADDING FUNCTION CALLS

##### A. Abstract Syntax

The abstract syntax of MIR with function calls is shown in Fig 7.

$$\begin{aligned} g &\in Function \\ p &:= \overline{S\{F : \tau\}}; \overline{g : f} \\ t &:= \dots \mid x = g(\bar{y}); B \end{aligned}$$

Fig. 7. MIR Abstract Syntax with Function Calls

We denote the local variable  $x$  in the function  $g$  as  $g.x$ . For the function call  $x = g(\bar{y})$ , we assume that  $y_r$  is the local variable in  $g$  function which is assigned as the return value. Also, we assume that  $\bar{x}_p$  is the list of local variables in the current context, which is assigned to each parameter  $\bar{y}$ . The procedure of function application can be described as follows:

- We assign each value in  $\bar{x}_p$  to the corresponding parameter in  $g$ . i.e.,  $\bar{g.y} = \bar{x}_p$ .
- Evaluate the function  $g$  with the current context.
- Then assign the return value  $g.y_r$  to  $x$ .

### B. Generating Type Constraints

The rules for type constraints in MIR with function calls are shown in Fig 8.

$$\begin{array}{c} \dots \\ x = g(\bar{y}); B \quad : \quad \overline{\llbracket x_p \rrbracket} \subseteq \overline{\llbracket g.y \rrbracket} \quad \wedge \quad \llbracket g.y_r \rrbracket \subseteq \llbracket x \rrbracket \end{array}$$

Fig. 8. MIR Type Constraints with Function Calls

## VI. ADDING FUNCTION POINTERS

### A. Abstract Syntax

The abstract syntax of MIR with function pointers is shown in Fig 9.

$$\begin{array}{lcl} s & := & \dots \mid x = g \\ t & := & \dots \mid x = y(\overline{y'}); B \\ \tau & := & \dots \mid \bar{\tau} \rightarrow \tau' \end{array}$$

Fig. 9. MIR Abstract Syntax with Function Pointers