

Semantics for EqMLton

Jordan Quinn

June 25, 2025

1 Introduction

1.1 E-Nodes and E-Graphs

Our E-graph implementation will be defined similarly to E-PEGs in Tate's EqSat paper [?]. However, each node is annotated with the possible types it could represent. Therefore, our e-graphs will be 5-tuples, $\langle N, L, C, E, \Gamma \rangle$ where Γ is our typing context.

For nodes n_1 and n_2 in our e-graph, $\llbracket n_1 \rrbracket = \llbracket n_2 \rrbracket$ if $L(n_1) = L(n_2)$, $\Gamma \vdash n_1 : \tau$ and $\Gamma \vdash n_2 : \tau$, and, if for $C(n_1) = \langle c_{1,1}, c_{1,2}, \dots \rangle$, $C(n_2) = \langle c_{2,1}, c_{2,2}, \dots \rangle$, $n = |C(n_1)| = |C(n_2)|$, $\llbracket c_{1,i} \rrbracket = \llbracket c_{2,i} \rrbracket$ for each $i \in \{1, 2, \dots, n\}$.

1.2 Regions

A **region** defines a set of e-graph nodes representing a computation. These region values appear in the CFG-skeleton and essentially represent continuations.

Formally, a region, of an e-graph G , is a three-tuple, $R = \langle N, I, r \rangle$ where N is the set of nodes of G that compose the region, I is a list of parameter nodes in N that represent inputs to our region, and $r \in N$ is the root node for our region that represents the result of the region's computation.

Note that, for region $R = \langle N, I, r \rangle$,

$$\frac{\Gamma \vdash r : \tau \quad \Gamma \vdash I = \langle i_1, i_2, \dots, i_k \rangle : \langle \tau_1, \tau_2, \dots, \tau_k \rangle}{\Gamma \vdash R : \langle \tau_1, \tau_2, \dots, \tau_k \rangle \rightarrow \tau}$$

In other words, region R is treated as a continuation that takes a tuple containing arguments for each parameter and returns a value of type τ .

Then, for two regions, $R_1 = \langle N_1, I_1, r_1 \rangle$ and $R_2 = \langle N_2, I_2, r_2 \rangle$ in e-graph $G = \langle N, L, C, E, \Gamma \rangle$, where $N_1 = \langle n_{1,1}, n_{1,2}, \dots \rangle$ and $I_1 = \langle i_{1,1}, i_{1,2}, \dots \rangle$ and similar for N_2, I_2 . Then, if $\Gamma \vdash R_1 : \tau$ and $\Gamma \vdash R_2 : \tau$, we know the lists of parameters are the same length and typed in the same order. Then, define L' and e-graph $G' = \langle N, L', C, E, \Gamma \rangle$ such that $L' = L[i_{1,1} \vdash L(i_2, 1), i_{1,2} \vdash L(i_2, 2), \dots]$. This ensures that the labeling is the same for each parameter list and thus $\llbracket i_{1,i} \rrbracket_{G'} = \llbracket i_{2,i} \rrbracket_{G'}$ for each $i \in \{1, \dots, |I_1| = |I_2|\}$. Then $\llbracket R_1 \rrbracket_G = \llbracket R_2 \rrbracket_G$ if

$\llbracket r_1 \rrbracket_{G'} = \llbracket r_2 \rrbracket_{G'}$. In other words, two regions are semantically equivalent if their root nodes are semantically equivalent, ignoring the labels of parameter nodes in our regions.

1.3 CFG Skeleton

The CFG-skeleton is a directed graph that encodes the control-flow decisions taken by a program. It is directly related to MLton's SSA graph, except we are only interested in branch and merge points. Instead of representing computations as statements in the CFG-skeleton, as is done in MLton's SSA, we

The CFG-skeleton is a directed graph where nodes, called blocks, represent branch and merge points in a program, and edges represent control flow transfer. Blocks are related to those of MLton's SSA implementation. They allow for opaque parameters that are referred to within the block. However, where MLton's SSA places computation in statements within each block, our CFG-skeleton only contains control flow information, using regions to represent computations used in our control flow statements. Each block can be defined with one control flow statement. Therefore, it is perhaps more helpful to think of our CFG skeleton as a single root block, $S = b$, and consider it along with a function that returns child blocks, $C(b) = \{b_1, \dots, b_b\}$, calculated from the control flow statement that comprises each block. Then we can define all blocks in our CFG skeleton, S^* , as the closure of the child function.

$$\begin{aligned} b(x_1, \dots, x_k) ::= & \text{goto}_{b'}(R_1 \langle x_{1,1}, \dots, x_{1,k_1} \rangle, R_2 \langle x_{2,1}, \dots, x_{2,k_1} \rangle \dots) \\ & | \text{if } R \text{ then goto}_{b_t}(R_{1,1} \langle x_{1,1,1}, \dots, x_{1,1,k_{1,1}} \rangle, R_{1,2} \dots) \\ & \quad \text{else goto}_{b_f}(R_{2,1} \langle x_{2,1,1}, \dots, x_{1,k_{2,1}} \rangle, R_{2,2} \dots) \\ & | \text{return } R \langle x_1, \dots, x_{k_1} \rangle \end{aligned}$$

Here, block b has a single control flow statement that defines outgoing edges. In **goto**, we have one outgoing edge from b to b' . We are passing arguments computed by evaluating regions R_1, R_2 , etc. on their corresponding argument lists, $\langle x_{1,1}, \dots, x_{1,k'} \rangle$ for R_1 , etc. where each x is a parameter in $\{x_1, \dots, x_k\}$. Similar is done to show edges for the other control flow statements. Note that **if** is used as an example for notation; we will solely use **match** according to the next subsection.

1.3.1 Pattern Matching

The primary control flow consideration is handling pattern matching in our CFG skeleton. We consider three options: tagged tuples, first evaluation, and check-first evaluation.

- **Tagged Tuple Pattern Matching**

In this variation, we consider each constructor a tagged tuple. The first value in our tuple indicates which constructor we are using, and the following values are parameters for that constructor. We then would represent a block with pattern-matching control flow as follows:

$$\begin{aligned}
b(x_1, \dots, x_k) ::= & \text{match } R \langle x_{1,1}, \dots, x_{1,k'} \rangle \\
& '1 \Rightarrow \text{goto}_{b_1} (R_{1,1} \langle x_{1,1,1}, \dots, x_{1,1,k_{1,1}} \rangle, R_{1,2} \dots) \\
& '2 \Rightarrow \text{goto}_{b_2} (R_{2,1} \langle x_{2,1,1}, \dots, x_{2,1,k_{2,1}} \rangle, R_{2,2} \dots) \\
& \vdots \\
& 'n \Rightarrow \text{goto}_{b_n} (R_{n,1} \langle x_{n,1,1}, \dots, x_{n,1,k_{n,1}} \rangle, R_{n,2} \dots)
\end{aligned}$$

Where the left-hand side of each \Rightarrow , 'label', is a label describing a constructor for τ when $\Gamma \vdash R : \tau$. The right-hand side of the \Rightarrow potentially describes naming parts of each constructor or passing additional parameters needed in b_i .

- **First Evaluation Pattern Matching**

Here, we take advantage of the fact that this is a theoretical representation of a computation. We define our match to pass control onto the first \Rightarrow that succeeds in evaluating all of its arguments. These arguments are the same as in tagged tuple pattern matching. However, we specially define match to keep trying each case until all regions evaluate to a non- \perp value.

$$\begin{aligned}
b(x_1, \dots, x_k) ::= & \text{match } R \langle x_{1,1}, \dots, x_{1,k'} \rangle \\
& \Rightarrow \text{goto}_{b_1} (R_{1,1} \langle x_{1,1,1}, \dots, x_{1,1,k_{1,1}} \rangle, R_{1,2} \dots) \\
& \Rightarrow \text{goto}_{b_2} (R_{2,1} \langle x_{2,1,1}, \dots, x_{2,1,k_{2,1}} \rangle, R_{2,2} \dots) \\
& \vdots \\
& \Rightarrow \text{goto}_{b_n} (R_{n,1} \langle x_{n,1,1}, \dots, x_{n,1,k_{n,1}} \rangle, R_{n,2} \dots)
\end{aligned}$$

Note that here, we do not need the original R region in our match, and we can rely on our R_i s to construct our operand for us.

- **Check First Pattern Matching**

This differs from the first two in that it first checks a predicate region, R_{p_i} , for each case to ensure that it can be deconstructed according to the pattern in the original program. The RHS is the same as the previous two versions.

$$\begin{aligned}
b(x_1, \dots, x_k) ::= & \\
& \text{match } R \langle x_{1,1}, \dots, x_{1,k'} \rangle \\
& R_{p_1} \langle x_{1,p,1}, \dots, x_{1,p,k_{1,p}} \rangle \Rightarrow \text{goto}_{b_1} (R_{1,1} \langle x_{1,1,1}, \dots, x_{1,1,k_{1,1}} \rangle, R_{1,2} \dots) \\
& R_{p_2} \langle x_{2,p,1}, \dots, x_{2,p,k_{2,p}} \rangle \Rightarrow \text{goto}_{b_2} (R_{2,1} \langle x_{2,1,1}, \dots, x_{2,1,k_{2,1}} \rangle, R_{2,2} \dots) \\
& \vdots \\
& R_{p_n} \langle x_{n,p,1}, \dots, x_{n,p,k_{n,p}} \rangle \Rightarrow \text{goto}_{b_n} (R_{n,1} \langle x_{n,1,1}, \dots, x_{n,1,k_{n,1}} \rangle, R_{n,2} \dots)
\end{aligned}$$

Note that here, we do not need the original R region in our match, and we can rely on our R_{p_i} s to construct our operand for us.

1.4 Equality Saturation for MLton

Our IR used for EqSat in MLton will then be a two-tuple, notably $R = \langle S, G \rangle$ where S is a CFG-skeleton and G is a PEG-style E-graph, likely derivative of FEG.

2 Formal Semantics

2.1 E-Graphs and E-Nodes

The semantics of an E-node is defined similarly to in FPEG [?]. That being, for E-Node n ,

$$\llbracket n \rrbracket = L(n) \llbracket C(n) \rrbracket$$

An E-graph, however, does not have a specific semantic meaning.

2.2 Regions

We define the semantics of a region, $R = \langle N, I, r \rangle$, with a list of arguments to R , $X = \langle x_1, \dots, x_n \rangle$ similarly to the semantics of E-graphs given by Tate, et al. [?]. Notably, we place the arguments passed to the region, which share a type with the list of parameters, $\Gamma(X) \sim \Gamma(I)$, into the e-graph as constant nodes that replace the parameter nodes. This is done via $r' = r[i_1 \mapsto \llbracket x_1 \rrbracket, \dots, i_n \mapsto \llbracket x_n \rrbracket]$. Then, we say that the semantics of $R(X) = R \langle x_1, \dots, x_n \rangle$ is $\llbracket R(X) \rrbracket = \llbracket r' \rrbracket$. Note that $R(X) \equiv R X \equiv R \langle x_1, \dots, x_n \rangle$ for convenience.

2.3 CFG Skeleton

Our CFG can be defined using the grammar in figure ?? . However, it is useful to think of the set of all block labels, bs . In this grammar, the list of arguments, X , to block b , is described by the list of types in T .

$$\begin{aligned}
B &::= \langle b, X, T, t \rangle \\
b &::= \text{block_label} \\
T &::= \langle \tau_1, \dots, \tau_n \rangle \\
t &::= \text{goto}_b(R_1 X_1, \dots, R_n X_n) \\
& \quad | \text{call}_f(R_1 X_1, \dots, R_n X_n) \\
& \quad | \text{raise}(R_1 X_1, \dots, R_n X_n) \\
& \quad | \text{return}(R_1 X_1, \dots, R_n X_n) \\
& \quad | \text{match}(R_p X_p; \\
& \quad \quad c_1 \Rightarrow \text{goto}_{b_1}(R_{1,1} X_{1,1}, \dots, R_{1,n_1} X_{1,n_1}, \\
& \quad \quad \dots, \\
& \quad \quad c_n \Rightarrow \text{goto}_{b_n}(R_{n,1} X_{n,1}, \dots, R_{n,n_n} X_{n,n_n} \\
& \quad \quad [; \Rightarrow \text{goto}_{b_d}(R_{d,1} X_{d,1}, \dots, R_{d,n_d} X_{d,n_d}])) \\
R &::= \text{region_label} \\
X &::= \langle x_1, \dots, x_n \rangle \\
x &::= \text{argument_label} \\
c &::= \text{constructor_label}
\end{aligned}$$

Figure 1: Grammar for a CFG

Type	Meta-Var	Purpose
<i>block</i>	<i>b</i>	CFG block
<i>arg</i>	<i>x</i>	Argument to a block or function
<i>region</i>	<i>R</i>	Region in the E-graph
<i>value</i>	<i>v</i>	An SML value, primitive or constructor
<i>return</i>	<i>r</i>	Ok/Err Ok or exceptional return
<i>store</i>	<i>S</i>	Store containing variable, value pairs

Figure 2: Table of metavariables used in semantics

true
false

Figure 3: Semantic-level atomics

block(b) = choose $B = \langle b', X, T, t \rangle \in bs$ such that $b \equiv b'$
sig(b) = **block**(b) _{T}
params(b) = **block**(b) _{X} : **sig**(b)
children(b) = $\{b' \in \mathbf{block}(b)_t\}$
is_prim(v) = **true** if v is a primitive type in MLton else **false**
c_label(v) = v if **is_prim**(v) else v_1

Figure 4: Semantic-level functions

$$\begin{array}{c}
\text{RETURN} \frac{\llbracket R(X) \rrbracket_s = v}{\langle \mathbf{return} \ R(X), s \rangle \Downarrow \langle \mathbf{Ok}(v), s \rangle} \\
\\
\text{GOTO} \frac{\begin{array}{c} \llbracket R_1 X_1 \rrbracket_s = v_1 \quad \dots \quad \llbracket R_n X_n \rrbracket_s = v_n \\ \langle b, s[\mathbf{params}(b) \mapsto \langle v_1, \dots, v_n \rangle] \rangle \Downarrow \langle r, s' \rangle \end{array}}{\langle \mathbf{goto}_b(R_1 X_1, \dots, R_n X_n), s \rangle \Downarrow \langle r, s' \rangle} \\
\\
\text{RAISE} \frac{\llbracket R_1 X_1 \rrbracket_s = v_1 \quad \dots \quad \llbracket R_n X_n \rrbracket_s = v_n}{\langle \mathbf{raise} \ (R_1 X_1, \dots, R_n X_n), s \rangle \Downarrow \langle \mathbf{Err}(v_1, \dots, v_n), s \rangle} \\
\\
\text{CALL} \frac{\llbracket R_1 X_1 \rrbracket_s = v_1 \quad \dots \quad \llbracket R_n X_n \rrbracket_s = v_n \quad \llbracket f(v_1, \dots, v_n) \rrbracket_s = r}{\langle \mathbf{call}_f(R_1 X_1, \dots, R_n X_n), s \rangle \Downarrow \langle r, s \rangle} \\
\\
\text{MATCH-CASE} \frac{\begin{array}{c} \llbracket R_p X_p \rrbracket_s = v_p \quad \llbracket c_i \rrbracket \sim \llbracket \mathbf{c_label}(v_p) \rrbracket \quad \langle \mathbf{goto}_{b_i}(R_{i,1} X_{i,1}, \dots), s \rangle \Downarrow \langle r, s' \rangle \\ \langle \mathbf{match}(R_p X_p; \dots, c_i \Rightarrow \mathbf{goto}_{b_i}(R_{i,1} X_{i,1}, \dots R_{i,n_i} X_{i,n_i}), \dots; \dots), s \rangle \Downarrow \langle r, s' \rangle \end{array}}{\langle \mathbf{match}(R_p X_p; \dots, c_i \Rightarrow \mathbf{goto}_{b_i}(R_{i,1} X_{i,1}, \dots R_{i,n_i} X_{i,n_i}), \dots; \dots), s \rangle \Downarrow \langle r, s' \rangle} \\
\\
\text{MATCH-DEFAULT} \frac{\begin{array}{c} \llbracket R_p X_p \rrbracket_s = v_p \\ \forall_{i \in \{1, \dots, n\}} \llbracket c_i \rrbracket \not\sim \llbracket \mathbf{c_label}(v_p) \rrbracket \quad \langle \mathbf{goto}_{b_d}(R_{d,1} X_{d,1}, \dots), s \rangle \Downarrow \langle r, s' \rangle \end{array}}{\langle \mathbf{match}(R_p X_p; c_1 \Rightarrow \dots, \dots, c_n \Rightarrow \dots; \mathbf{goto}_{b_d}(R_{d,1} X_{d,1}, \dots)), s \rangle \Downarrow \langle r, s' \rangle}
\end{array}$$

Figure 5: Semantics for CFG skeleton

In ??, there are a few additional considerations. First, $\llbracket RX \rrbracket_s$ is the same as $\llbracket RV \rrbracket$ where $V = \langle s(x_1), \dots, s(x_n) \rangle$. Then, MATCH-CASE defines the semantics for when we successfully match based on the constructor while MATCH-DEFAULT covers when we do not match on the constructor and execute the default case. Note that we can assume there will be a default case when necessary as we are translating from a well-formed MLton SSA program.

An important consideration is function arguments and globals. These are available to every block in the CFG regardless of where they appear. As such, we start with each function argument and global in our initial store, s .

2.4 CFG Properties

In our CFG skeleton, if we consider each `goto` assignment, as it saves this to the store, we can consider our CFG skeleton to have the dominance property, where each use of a variable is dominated by its definition.