

Semantics for EqMLton

Jordan Quinn

June 18, 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 [2]. 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.

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

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 [1]. 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, $R = \langle N, I, r \rangle$, as a closure. This using the semantics of the node r , $\llbracket R \rrbracket = \llbracket r \rrbracket$ where we use the semantics of E-node r . Then, the semantics of application, $R(X)$, where I and X have congruent types, is defined as the application of the closure, $\llbracket R \rrbracket$ on tuple X . I.e., $\llbracket R(X) \rrbracket = \llbracket R \rrbracket (\llbracket X \rrbracket)$ and $\llbracket X \rrbracket = \langle \llbracket x_1 \rrbracket, \dots, \llbracket x_n \rrbracket \rangle$.

We can also give a semantic meaning to a region well-typed region application of, $R = \langle N, I, r \rangle$, in E-Graph G , applied to a tuple via:

$$\llbracket R(X = \langle x_1, x_2, \dots, x_n \rangle) \rrbracket = \llbracket r \rrbracket_{G'}$$

where $G' = \langle N, L', C, E, \Gamma \rangle$ and $L' = L[i_1 \vdash x_1, i_2 \vdash x_2, \dots, i_n \vdash x_n]$ for $I = \langle i_1, i_2, \dots, i_n \rangle$.

2.3 CFG Skeleton

- Use the denotational semantics of a region to define operational semantics

$$\overline{\langle R, \sigma, G \rangle \rightarrow \langle \llbracket R \rrbracket_G, \sigma, G \rangle}$$

$$\overline{\langle R \langle x_1, \dots, x_n \rangle, \sigma, G \rangle \rightarrow \langle \llbracket R \langle x_1, \dots, x_n \rangle \rrbracket_G, \sigma, G \rangle}$$

- X_i is the parameter tuple for each region

$$\frac{\langle x_1, \sigma, G \rangle \rightarrow \langle x'_1, \sigma, G \rangle \dots \langle x_n, \sigma, G \rangle \rightarrow \langle x'_n, \sigma, G \rangle \langle \langle R_1 X_1, \sigma, G \rangle \rightarrow \langle R'_1, \sigma, G \rangle \dots}{\langle b \langle x_1, \dots, x_n \rangle = \text{call}_f(R_1 X_1, \dots, R_{n'} X_{n'}), \sigma \rangle \rightarrow f \langle R'_1, \dots, R'_{n'} \rangle}$$

$$\frac{\langle x_1, \sigma, G \rangle \rightarrow \langle x'_1, \sigma, G \rangle \dots \langle x_n, \sigma, G \rangle \rightarrow \langle x'_n, \sigma, G \rangle \langle \langle R_1 X_1, \sigma, G \rangle \rightarrow \langle R'_1, \sigma, G \rangle \dots}{\langle b \langle x_1, \dots, x_n \rangle = \text{goto}_f(R_1 X_1, \dots, R_{n'} X_{n'}), \sigma \rangle \rightarrow b \langle', \dots, R'_{n'} \rangle}$$

References

- [1] Matthew DellaNeve. *Equality Saturation for the MLton Compiler*. Capstone report, Rochester Institute of Technology, 2023.
- [2] Ross Tate, Michael Stepp, Zachary Tatlock, and Sorin Lerner. Equality saturation: a new approach to optimization. In *Logical Methods in Computer Science*, volume 7, 2011.