

Lecture Notes on Static Semantics

15-411: Compiler Design
Frank Pfenning, Jan Hoffmann*

Lecture 9a
September 26, 2019

1 Introduction

After lexing and parsing, a compiler will usually apply *elaboration* to translate the parse tree to a high-level intermediate form often called *abstract syntax*. Then we verify that the abstract syntax satisfies the requirements of the *static semantics*. Sometimes, there is some ambiguity whether a given condition should be enforced by the grammar, by elaboration, or while checking the static semantics. We will not be concerned with details of attribution, but how to describe and then implement various static semantic conditions. The principal properties to verify for C0 and the sublanguages discussed in this course are:

- *Initialization*: variables must be defined before they are used.
- *Proper returns*: functions that return a value must have an explicit return statement on every control flow path starting at the beginning of the function.
- *Types*: the program must be well-typed.

Type checking is frequently discussed in the literature, so we use initialization as our running example and discuss typing at the end, in [Section 9](#).

2 Abstract Syntax

We will use a slightly restricted form of the abstract syntax in [Lecture 9](#) on IR trees, with the addition of variable declaration with their scope. This fragment exhibits

*with small modifications in 2019

all the relevant features for the purposes of the present lecture.

Expressions $e ::= n \mid x \mid e_1 \oplus e_2 \mid e_1 \ \&\& \ e_2$
 Statements $s ::= \text{assign}(x, e) \mid \text{if}(e, s_1, s_2) \mid \text{while}(e, s)$
 $\mid \text{return}(e) \mid \text{nop} \mid \text{seq}(s_1, s_2) \mid \text{decl}(x, \tau, s)$

3 Definition and Use

Initialization guarantees that every variable is defined before it is used. It is natural to specify this in two parts: specify when a variable is *defined*, and specify when it is *used*. An error is signaled if we cannot show that every variable in the program is defined before it is used. As usual, this property is an *approximation* of what actual behaviors can be exhibited at runtime.

First, we define when a variable is used in an expression, written as $\text{use}(e, x)$. This is entirely straightforward, since we have a clear separation of expressions and statements in our language.

no rule for $\text{use}(n, x)$	$\frac{}{\text{use}(x, x)}$	no rule for $\text{use}(y, x), y \neq x$
$\frac{\text{use}(e_1, x)}{\text{use}(e_1 \oplus e_2, x)}$		$\frac{\text{use}(e_2, x)}{\text{use}(e_1 \oplus e_2, x)}$
$\frac{\text{use}(e_1, x)}{\text{use}(e_1 \ \&\& \ e_2, x)}$		$\frac{\text{use}(e_2, x)}{\text{use}(e_1 \ \&\& \ e_2, x)}$

We see already here that $\text{use}(e, x)$ is a so-called *may*-property: x *may* be used during the evaluation of x , but it is not guaranteed to be actually used. For example, the expression $y > 0 \ \&\& \ x > 0$ may or may not actually use x . The expression $\text{false} \ \&\& \ x > 0$ will actually never use x , and yet we flag it as possibly being used.

This is appropriate: we would like to raise an error if there is a possibility that an uninitialized variable may be used. Because determining this in general is undecidable, we need to approximate it. Our approximation essentially says that any variable occurring in an expression may be used. The rule above express this more formally.

For a language to be usable, it is important that the rules governing the static semantics are easy to understand for the programmer and have some internal coherence. While it might make sense to allow $\text{false} \ \&\& \ x > 0$ in particular, what is the general rule? Designing programming languages and their static semantics is difficult and requires a good balance of formal understanding of the properties of programming languages and programmer’s intuition.

Once we have defined use for expressions, we should consider statements. Does an assignment $x = e$ use x ? Our prior experience with liveness analysis for register allocation on abstract machine could would say: *only if it is used in e* . We stay consistent with this intuition and terminology and write $\text{live}(s, x)$ for the judgment that x is live in s . This means the value of x is relevant to the execution of s .

Before we specify liveness, we should specify when a variable is *defined*. This is because, for example, the variable x is not live before the statement $x = 3$, because its current value does not matter for this statement, or any subsequent statement. We write $\text{def}(s, x)$ is the execution of statement s will define x . This is an example of a *must*-property: we want to be sure that whenever s executes (and completes normally, without returning from the current function or raising an exception of some form), the x has been defined.

$$\frac{}{\text{def}(\text{assign}(x, e), x)} \quad \text{no rule for } \text{def}(\text{assign}(y, e), x), y \neq x$$

$$\frac{\text{def}(s_1, x) \quad \text{def}(s_2, x)}{\text{def}(\text{if}(e, s_1, s_2), x)} \quad \text{no rule for } \text{def}(\text{while}(e, s), x)$$

The last two rules clearly illustrate that $\text{def}(s, x)$ is a *must*-property: A conditional only defines a variable if is it defined along both branches, and a while loop does not define any variable (since the body may never be executed).

$$\text{no rule for } \text{def}(\text{nop}, x) \quad \frac{\text{def}(s_1, x)}{\text{def}(\text{seq}(s_1, s_2), x)} \quad \frac{\text{def}(s_2, x)}{\text{def}(\text{seq}(s_1, s_2), x)}$$

$$\frac{\text{def}(s, x) \quad y \neq x}{\text{def}(\text{decl}(y, \tau, s), x)}$$

The side condition on the last rule applies because s is the scope of y . It makes sense to maintain the invariant that the defined variables in a statement are a subset of the declared ones. Otherwise, the program would be invalid anyway. However, in C and C0 we require that a variable x is declared at most once on a control-flow path. So we know that y cannot be in scope in the tree above of the declaration. Hence it is also not defined above its declaration.

A strange case arises for return statement. Since a return statement will never pass the control flow to the next instruction, any subsequent statements are unreachable. It is therefore permissible to claim that *all* variables currently in scope have been defined. We capture this by simply stating that $\text{return}(e)$ defines any variable.

$$\frac{}{\text{def}(\text{return}(e), x)}$$

4 Liveness

We now lift the $\text{use}(e, x)$ property to statements, written as $\text{live}(s, x)$ (x is live in s). Liveness is again a *may*-property.

$$\frac{\text{use}(e, x)}{\text{live}(\text{assign}(y, e), x)}$$

$$\frac{\text{use}(e, x)}{\text{live}(\text{if}(e, s_1, s_2), x)} \quad \frac{\text{live}(s_1, x)}{\text{live}(\text{if}(e, s_1, s_2), x)} \quad \frac{\text{live}(s_2, x)}{\text{live}(\text{if}(e, s_1, s_2), x)}$$

We observe that liveness is indeed a *may*-property, since a variable is live in a conditional if is used in the condition or live in one or more of the branches. Similarly, if a variable is live in the body of a loop, it is live before because the loop body *may* be executed.

$$\frac{\text{use}(e, x)}{\text{live}(\text{while}(e, s), x)} \quad \frac{\text{live}(s, x)}{\text{live}(\text{while}(e, s), x)}$$

$$\frac{\text{use}(e, x)}{\text{live}(\text{return}(e), x)} \quad \text{no rule for } \text{live}(\text{nop}, x) \quad \frac{\text{live}(x, s) \quad y \neq x}{\text{live}(\text{decl}(y, \tau, s), x)}$$

In some way the most interesting case is a sequence of statements, $\text{seq}(s_1, s_2)$. If a variable is live in s_2 it is only live in the composition if it is not defined in s_1 !

$$\frac{\text{live}(s_1, x)}{\text{live}(\text{seq}(s_1, s_2), x)} \quad \frac{\neg \text{def}(s_1, x) \quad \text{live}(s_2, x)}{\text{live}(\text{seq}(s_1, s_2), x)}$$

5 Initialization

Given liveness, we can now say when proper initialization is violated: If a variable is live at the site of its declaration. That means that its value would be used somewhere before it is defined. Assume we have a program p and we write “ s in p ” if s is a statement appearing in p . The following rule captures the general condition.

$$\frac{\text{decl}(x, \tau, s) \text{ in } p \quad \text{live}(s, x)}{\text{error}}$$

Unlike the previous rules in the lecture, this one should be read from the premises to the conclusion. In this way it is similar to our rules for liveness from [Lecture 4](#).

This brings out an important distinction when we try to convert the specification rules into an implementation. We have to decide if the rules should be read

from the premises to the conclusion, or from the conclusion to the premises. Sometimes, the same property can be specified in different directions. For example, we can define a predicate `init` which verifies that all variables are properly initialized and which works from the conclusion to the premises with the following schema.

$$\frac{}{\text{init}(\text{nop})} \qquad \frac{\text{init}(s_1) \quad \text{init}(s_2)}{\text{init}(\text{seq}(s_1, s_2))}$$

$$\frac{\text{init}(s) \quad \neg \text{live}(s, x)}{\text{init}(\text{decl}(x, \tau, s))} \qquad (\text{other rules omitted})$$

The omitted rules just verify each substatement so that all declarations in the program are checked in the end.

6 From Judgments to Functions

We now focus on the special case that the inference rules are to be read bottom-up, starting with the judgments we ultimately want to verify. Consider `init(s)`. When we start, `s` is known and we are trying to determine if there is a deduction of `init(s)` given the rules we have put down. If there is such a deduction, we succeed. If not, we issue an error message. We can model this as a function returning a boolean, or a function returning no interesting value but raising an exception in case the property is violated.

$$\text{init} : \text{stm} \rightarrow \text{bool}$$

Now each of the rules becomes a case in the function definition.

$$\begin{aligned} \text{init}(\text{nop}) &= \top \\ \text{init}(\text{seq}(s_1, s_2)) &= \text{init}(s_1) \wedge \text{init}(s_2) \\ \text{init}(\text{decl}(x, \tau, s)) &= \text{init}(s) \wedge \neg \text{live}(s, x) \\ &\dots \end{aligned}$$

Here we assume a boolean constant \top (for *true*) and boolean operators conjunction $A \wedge B$ and negation $\neg A$ in the functional language; later we might use disjunction $A \vee B$ and falsehood \perp . When we call `live(s, x)` we assume that it is a similar function.

$$\text{live} : \text{stm} \times \text{var} \rightarrow \text{bool}$$

This function is now a transcription of the rules for the `live` judgment. In this process we sometimes have to combine multiple rules into a single case of the function definition (as, for example, for `seq(s1, s2)`).

$$\begin{aligned}
\text{live}(\text{nop}, x) &= \perp \\
\text{live}(\text{seq}(s_1, s_2), x) &= \text{live}(s_1, x) \vee (\neg \text{def}(s_1, x) \wedge \text{live}(s_2, x)) \\
\text{live}(\text{decl}(y, \tau, s), x) &= y \neq x \wedge \text{live}(x, s) \\
&\dots
\end{aligned}$$

We still have to write functions for predicates $\text{def}(s, x)$ and $\text{use}(e, x)$, but these are a straightforward exercise now.

$$\begin{aligned}
\text{def} &: \text{stm} \times \text{var} \rightarrow \text{bool} \\
\text{use} &: \text{exp} \times \text{var} \rightarrow \text{bool}
\end{aligned}$$

The whole translation was relatively straightforward, primarily because the rules were well-designed, and because we always had enough information to just write a boolean function.

7 Maintaining Set of Variables

What we have done above is a perfectly adequate implementation of initialization checking. But we might also try to rewrite it in order limit the number of traversals of the statements. For example, in

$$\text{live}(\text{seq}(s_1, s_2), x) = \text{live}(s_1, x) \vee (\neg \text{def}(s_1, x) \wedge \text{live}(s_2, x))$$

we may traverse s_1 twice: once to check if x is live in s_1 , and once to see if x is defined in s_1 . In general, we might traverse statements multiple times, namely for each variable declaration in whose scope it lies. This in itself is not a performance bug, but let's see how one might change it.

One way this can often be done is to notice that for any statement s , there could be multiple variables x such that $\text{live}(s, x)$ or $\text{def}(s, x)$ holds. We can try to combine these into a *set*. We denote a set of variables with δ and define the following two judgments:

- $\text{init}(\delta, s, \delta')$: assuming all the variables in δ are defined when s is reached, no uninitialized variable will be referenced and after its execution all the variables in δ' will be defined.
- $\text{use}(\delta, e)$: e will only reference variables defined in δ .

As a notational convention, we isolate assumptions on the left-hand side of a turnstile symbol, writing:

- $\delta \vdash s \Rightarrow \delta'$ for $\text{init}(\delta, s, \delta')$.
- $\delta \vdash e$ for $\text{use}(\delta, e)$.

From the previous rules we develop the following:

$$\begin{array}{c}
 \frac{}{\delta \vdash \text{nop} \Rightarrow \delta} \qquad \frac{\delta \vdash s_1 \Rightarrow \delta_1 \quad \delta_1 \vdash s_2 \Rightarrow \delta_2}{\delta \vdash \text{seq}(s_1, s_2) \Rightarrow \delta_2} \\
 \\
 \frac{\delta \vdash e}{\delta \vdash \text{assign}(x, e) \Rightarrow \delta \cup \{x\}} \qquad \frac{\delta \vdash e \quad \delta \vdash s_1 \Rightarrow \delta_1 \quad \delta \vdash s_2 \Rightarrow \delta_2}{\delta \vdash \text{if}(e, s_1, s_2) \Rightarrow \delta_1 \cap \delta_2} \\
 \\
 \frac{\delta \vdash e \quad \delta \vdash s \Rightarrow \delta'}{\delta \vdash \text{while}(e, s) \Rightarrow \delta} \qquad \frac{\delta \vdash s \Rightarrow \delta'}{\delta \vdash \text{decl}(y, \tau, s) \Rightarrow \delta' - \{y\}} \\
 \\
 \frac{\delta \vdash e}{\delta \vdash \text{return}(e) \Rightarrow \{x \mid x \text{ in scope}\}}
 \end{array}$$

It is worth reading these rules carefully to make sure you understand them. The last one is somewhat problematic, since we don't have enough information to know which declarations we are in the scope of. We should generalize our judgment to $\gamma ; \delta \vdash s \Rightarrow \delta'$, where γ is the set of all variables currently in scope. Then the last rule might become

$$\frac{\delta \vdash e}{\gamma ; \delta \vdash \text{return}(e) \Rightarrow \gamma}$$

and we would systematically add γ to all other judgments. We again leave this as an exercise.

In these judgments we have traded the complexity of traversing statements multiple times with the complexity of maintaining variables sets.

8 Modes of Judgments

If we consider the judgment $\delta \vdash e$ there is nothing new to consider: we would translate this to a function

$$\text{use} : \text{set var} \times \text{exp} \rightarrow \text{bool}$$

On the other hand, it does not work to translate $\delta \vdash s \Rightarrow \delta'$ as

$$\text{init} : \text{set var} \times \text{stm} \times \text{set var} \rightarrow \text{bool}$$

This is because, in general, we do not know δ' before we start out. We need to *compute* it as part of building the deduction! So we need to implement

$$\text{init} : \text{set var} \times \text{stm} \rightarrow \text{bool} \times \text{set var}$$

In order to handle $\text{return}(e)$, we probably should also pass in a second set of declared variables or a context. We could also avoid returning a boolean by just returning an optional set of defined variables, or raise an exception in case we discover a variable that is used but not defined.

Examining the rules shows that we will need to be able to add variables and to remove variables from sets, as well as compute intersections. Otherwise, the code should be relatively straightforward.

Before we actually start this coding, we should go over the inference rules to make sure we always have enough information to compute the output δ' given the inputs δ and s . This is the purpose of *mode checking*. Let's go over one example:

$$\frac{\delta \vdash s_1 \Rightarrow \delta_1 \quad \delta_1 \vdash s_2 \Rightarrow \delta_2}{\delta \vdash \text{seq}(s_1, s_2) \Rightarrow \delta_2}$$

Initially, we know the input δ and $s = \text{seq}(s_1, s_2)$. This means we also know s_1 and s_2 . We cannot yet compute δ_2 , since the required input δ_1 in the second premise is unknown. But we can compute δ_1 from the first premise since we know δ and s_1 and this point. This gives us δ_1 and we can now compute δ_2 from the second premise and return it in the conclusion.

$$\text{init}(\delta, \text{seq}(s_1, s_2)) = \text{let } \delta_1 = \text{init}(\delta, s_1) \text{ in } \text{init}(\delta_1, s_2)$$

9 Typing Judgments¹

Arguably the most important judgment on programs is whether they are well-typed. The *context* (or *type environment*) Γ assigns types to variables. This is often defined explicitly with

$$\Gamma ::= \cdot \mid \Gamma, x:\tau$$

where we write $\Gamma(x) = \tau$ when $x:\tau$ in Γ . The typing judgment for expressions

$$\Gamma \vdash e : \tau$$

verifies that the expression e is well-typed with type τ , assuming the variables are typed as prescribed by Γ . Most of the rules are straightforward; we show a couple.

$$\frac{\Gamma(x) = \tau}{\Gamma \vdash x : \tau} \quad \frac{}{\Gamma \vdash n : \text{int}} \quad \frac{}{\Gamma \vdash \text{true} : \text{bool}} \quad \frac{}{\Gamma \vdash \text{false} : \text{bool}}$$

$$\frac{\Gamma \vdash e_1 : \text{int} \quad \Gamma \vdash e_2 : \text{int}}{\Gamma \vdash e_1 + e_2 : \text{int}} \quad \frac{\Gamma \vdash e_1 : \text{bool} \quad \Gamma \vdash e_2 : \text{bool}}{\Gamma \vdash e_1 \ \&\& \ e_2 : \text{bool}}$$

¹Not covered in lecture

Typing for statements is slightly more complex. Statements are executed for their effects, but statements in the body of a functions also ultimately return a value. We write

$$\Gamma \vdash s : [\tau]$$

to express that s is well-typed in context Γ . If s returns (using a $\text{return}(e)$ statement), then e must be of type τ . We use this to check that no matter how a function returns, the returned value is always of the correct type.

$$\frac{\Gamma(x) = \tau' \quad \Gamma \vdash e : \tau'}{\Gamma \vdash \text{assign}(x, e) : [\tau]} \quad \frac{\Gamma \vdash e : \text{bool} \quad \Gamma \vdash s_1 : [\tau] \quad \Gamma \vdash s_2 : [\tau]}{\Gamma \vdash \text{if}(e, s_1, s_2) : [\tau]}$$

$$\frac{\Gamma \vdash e : \text{bool} \quad \Gamma \vdash s : [\tau]}{\Gamma \vdash \text{while}(e, s) : [\tau]} \quad \frac{\Gamma \vdash e : \tau}{\Gamma \vdash \text{return}(e) : [\tau]}$$

$$\frac{}{\Gamma \vdash \text{nop} : [\tau]} \quad \frac{\Gamma \vdash s_1 : [\tau] \quad \Gamma \vdash s_2 : [\tau]}{\Gamma \vdash \text{seq}(s_1, s_2) : [\tau]}$$

$$\frac{\Gamma, x:\tau' \vdash s : [\tau]}{\Gamma \vdash \text{decl}(x, \tau', s) : [\tau]}$$

In the last rule for declarations, we might prohibit shadowing of variables by requiring that $x \notin \text{dom}(\Gamma)$. Alternatively, we could stipulate that the rightmost occurrence of x in Γ is the one considered when calculating $\Gamma(x)$. It is also possible that we already know that no conflict can occur, since shadowing may have been ruled out during elaboration already.

10 Modes for Typing

When implementing type-checking, we need to decide on a *mode* for the judgment. Clearly, we want the context Γ and the expression e or statement s to be known, but what about the type?

We first look at expression typing, $\Gamma \vdash e : \tau$. Can we always know τ ? Perhaps in our small language fragment from this lecture, but not in L3. For example, if we check an expression $e_1 == e_2 : \text{bool}$, we may know it is of type bool but we do not know the types of e_1 or e_2 (they could be bool or int). Similarly, if we have an expression used as a statement, we do not know the type of expression. Therefore, we should implement a function that takes the context Γ and e as input and *synthesizes* a type τ such that $\Gamma \vdash e : \tau$ (if such a type exists, and fails otherwise). The resulting type τ can be then be compared to a given type if that is known. Of course, you should go through the rules and verify that one can indeed always synthesize a type.

For the typing of statements $\Gamma \vdash s : [\tau]$ the situation is slightly different. Because τ is the return type of the function in which s occurs, we will know τ instead of having to synthesize it. We say we *check* the statement s against the return type τ .

Therefore, if we assume that functions raise an exception if an expression or statement is not well-typed, we might have functions such as

$$\begin{aligned}\text{syn_exp} & : \text{ctx} \times \text{exp} \rightarrow \text{tp} \\ \text{chk_stm} & : \text{ctx} \times \text{stm} \times \text{tp} \rightarrow \text{unit}\end{aligned}$$

For convenience, we might also write a function

$$\text{chk_exp} : \text{ctx} \times \text{exp} \times \text{tp} \rightarrow \text{unit}$$

where $\text{chk_exp}(e, \tau)$ would simply synthesize a type τ' for e and compare it to τ .

Questions

1. Write out the rules for proper returns: along each control flow path starting at the beginning of a function, there must be a return statement. Clearly, this is a *must*-property.
2. Write some cases in the functions for type-checking.