

Lecture Notes 3

Type Safety

Carlo Angiuli

B522: PL Foundations
January 26, 2026

In this lecture, we use the concepts introduced in the previous lecture to state and prove our first interesting theorem about a programming language: *type safety* (or *type soundness*) for a first-order language with booleans and numbers.

Type safety for first-order languages is covered in Chapters 4–6 of Harper [Har16], but we will consider a different language taken from Chapter 8 of Pierce [Pie02] because we’re covering topics in a different order from the textbook.

Remark 3.1. Today we will encounter several instances of inconsistent or unsettled terminology; I will flag these as they come up.

1 A slightly more interesting language

Our running example in the previous lecture was a simple programming language of boolean expressions. This lecture we’ll consider a language with both booleans and natural numbers.

Definition 3.2. The judgment $e \text{ tm}$ is defined by the following BNF grammar:

$$\begin{aligned} \text{Terms } e ::= & \text{ true } | \text{ false } | \text{ if}(e, e, e) \\ & | \text{ zero } | \text{ suc}(e) | \text{ pred}(e) | \text{ zero?}(e) \end{aligned}$$

We will define the meaning of this language shortly, but some helpful remarks: We’ve removed `not` to shorten our proofs; it’s definable in terms of `if`. `pred` is short for “predecessor,” as in the opposite of “successor”; it subtracts one. `zero?` is a boolean test of whether its (numerical) input is the number zero.

2 Operational semantics (dynamics)

At the end of the previous lecture, we defined the meaning of programs using a binary *evaluation* judgment $e \Downarrow e'$ stating that the expression e evaluates to the expression e' . We will do two things differently this time:

- Rather than defining a binary judgment for the “full” evaluation of a term, we will define a binary judgment for taking *a single step of computation*.
- We will define a unary judgment expressing that a term is *finished computing*.

The evaluation judgment $e \Downarrow e'$ from previous lecture is sometimes called *natural semantics* or *big-step operational semantics*; it corresponds to defining an interpreter by structural recursion.

Today’s judgment $e \rightarrow e'$ is often called *structural operational semantics* or *small-step operational semantics* (for obvious reasons). Small-step judgments are used more commonly than evaluation judgments; there’s nothing wrong with evaluation judgments, but small-step judgments are often easier to reason about, especially when proving type safety.

Remark 3.3. In both approaches, note that the intermediate and final stages of a computation are drawn from the same exact collection of terms as our input language. This is a convenient simplifying assumption when it works, but it doesn’t always. For example, your interpreters in C311/B521 evaluated lambdas to “closures,” which were not part of the input language.

Definition 3.4 (Values). For e tm, we define the judgment e val (“ e is a value”) by the following inference rules. We notate values using the metavariable v .

$$\frac{}{\text{true val}} \quad \frac{}{\text{false val}} \quad \frac{}{\text{zero val}} \quad \frac{v \text{ val}}{\text{suc}(v) \text{ val}}$$

Values are programs that are “finished computing.” Note that, like our set of programs, the set of values includes some “nonsense,” like $\text{suc}(\text{true})$.

Definition 3.5 (Small-step operational semantics). For e tm, we define the judgment $e \rightarrow e'$ (“ e steps to e' ”) by the inference rules in Figure 3.1.

The reader may notice that there are two distinct “kinds of rules” in small-step operational semantics. Rules like

$$\frac{}{\text{if}(\text{true}, e_2, e_3) \rightarrow e_2}$$

$$\begin{array}{c}
\overline{\text{if}(\text{true}, e_2, e_3) \mapsto e_2} \quad \overline{\text{if}(\text{false}, e_2, e_3) \mapsto e_3} \\
\overline{e_1 \mapsto e'_1} \quad \overline{e \mapsto e'} \\
\text{if}(e_1, e_2, e_3) \mapsto \text{if}(e'_1, e_2, e_3) \quad \text{suc}(e) \mapsto \text{suc}(e') \\
\overline{\text{pred(zero)} \mapsto \text{zero}}^* \quad \overline{v \text{ val}} \quad \overline{e \mapsto e'} \\
\text{pred}(\text{succ}(v)) \mapsto v \quad \text{pred}(e) \mapsto \text{pred}(e') \\
\overline{\text{zero?}(\text{zero}) \mapsto \text{true}} \quad \overline{v \text{ val}} \\
\text{zero?}(\text{succ}(v)) \mapsto \text{false} \\
\overline{e \mapsto e'} \\
\text{zero?}(e) \mapsto \text{zero?}(e')
\end{array}$$

Figure 3.1: Definition of \mapsto .

are sometimes called *principal reductions*; they define how the primitive operations behave on values. On the other hand, *congruence reduction* rules like

$$\overline{e_1 \mapsto e'_1} \\
\text{if}(e_1, e_2, e_3) \mapsto \text{if}(e'_1, e_2, e_3)$$

define the language's *evaluation order*; in this language, `if` fully evaluates its first argument but not its second or third arguments. We say that the first argument of `if` is the *principal argument*.

Remark 3.6. There seems to be no standard terminology for principal reductions—see [this February 2021 Twitter thread](#) between Derek Dreyer, myself, and others—but this term is well-motivated by proof theory and was previously used by Benton et al. [Ben+93]. As for the congruence reductions, terminology varies; Harper [Har16] calls them *search transitions* and Pierce [Pie02] calls them *congruence rules*.

Remark 3.7. There are also two kinds of term operators: the ones that “make data” and the ones that “consume data.” `true`, `false`, `zero`, and `suc` make data, and are the values of our language. `if`, `pred`, and `zero?` consume data; each of them has a congruence reduction rule and several principal reductions.

Remark 3.8. Figure 3.1 contains one possibly dubious rule: $\text{pred}(\text{zero}) \mapsto \text{zero}$. We will revisit this rule later.

The following are “sanity check” lemmas that we can prove by rule induction.

Lemma 3.9. *If v val then v tm.*

Lemma 3.10. *If e tm and $e \mapsto e'$ then e' tm.*

Lemma 3.11 (Determinacy). *If $e \mapsto e'$ and $e \mapsto e''$ then $e' = e''$.*

Lemma 3.12 (Finality of values). *If v val then there is no e' such that $v \mapsto e'$.*

We write $e \nmapsto$ to mean that there exists no e' such that $e \mapsto e'$.

Remark 3.13. You might be tempted to also state the following lemma:

“If e tm then either $e \mapsto e'$ or e val.”

Unfortunately, that is **false**: terms such as $\text{zero?}(\text{true})$ and $\text{if}(\text{zero}, e_2, e_3)$ are neither values nor take a step, because their principal arguments are somehow “the wrong kind of thing.” Such terms are called *stuck*.

Definition 3.14. For e tm, we define the judgment $e \mapsto^* e'$ (“ e takes zero or more steps to e' ”) by the following inference rules:

$$\frac{}{e \mapsto^* e} \quad \frac{e \mapsto e' \quad e' \mapsto^* e''}{e \mapsto^* e''}$$

\mapsto^* is called the *reflexive transitive closure of \mapsto* .

Lemma 3.15 (Uniqueness of values). *If $e \mapsto^* v$ and $e \mapsto^* v'$ where v val and v' val, then $v = v'$.*

We typically don’t write derivation trees for \mapsto or \mapsto^* because they are quite awkwardly shaped. Instead, we typically write a sequence of single-step reductions like so:

```

suc(if(zero?(pred(suc(zero)))), suc(zero), zero))
  ↪ suc(if(zero?(zero), suc(zero), zero))
  ↪ suc(if(true, suc(zero), zero))
  ↪ suc(suc(zero))

```

The underlines are optional but indicate the subterm that is being simplified in the following step. A term that matches the left-hand side of a principal reduction is called a *reducible expression*, or *redex* for short.

Remark 3.16. The plural of redex is redexes, **not** redices.

3 Type system (statics): 2 Types 2 Furious

Because our language has two different kinds of data in it, it is possible for evaluation to get stuck if an operator that expects a number is given a boolean, or vice versa. We consider programs that eventually get stuck to be erroneous or nonsense, and would like to exclude them from consideration *before we even try to evaluate them*. We do this by syntactically characterizing which programs produce numerical values, and which programs produce boolean values.

$$\text{Types} \quad \tau ::= \text{bool} \mid \text{num}$$

Definition 3.17 (Type system). For e tm and τ ty, we define the judgment $e : \tau$ (“ e has type τ ”) by the following inference rules:

$$\begin{array}{cccc} \frac{}{\text{true} : \text{bool}} & \frac{}{\text{false} : \text{bool}} & \frac{}{\text{zero} : \text{num}} & \frac{e : \text{num}}{\text{suc}(e) : \text{num}} \\ \frac{e : \text{num}}{\text{pred}(e) : \text{num}} & \frac{e : \text{num}}{\text{zero?}(e) : \text{bool}} & \frac{e_1 : \text{bool} \quad e_2 : \tau \quad e_3 : \tau}{\text{if}(e_1, e_2, e_3) : \tau} & \end{array}$$

Exercise 3.18. Show there is no τ ty such that $\text{zero?}(\text{true}) : \tau$.

Lemma 3.19 (Uniqueness of types). *If $e : \tau$ and $e : \tau'$ then $\tau = \tau'$.*

Note that our type system does not refer to our operational semantics. Rather, it is a purely *static* analysis of the syntactic structure of terms: $e : \text{num}$ never gets stuck, and its only possible values are natural numbers; $e : \text{bool}$ never gets stuck, and its only possible values are `true` and `false`. In the next section, we will show that our analysis is *sound* (Theorem 3.29), but our analysis (like essentially all type systems) is not *complete*. That is, it is a conservative analysis that “misses” some programs that do actually compute boolean or numerical values.

Remark 3.20. Terms of type `bool` are not “booleans” and terms of type `num` are not “numbers”; they are programs that we expect to *compute* booleans or numbers.

Exercise 3.21. Show that `if(true, true, zero)` does not have type `bool` (in fact it has no type at all), but it steps to `true`.

4 Type safety = progress + preservation

The type safety theorem, also known as type soundness, expresses that our type system and operational semantics agree with one another. There are several ways

to state this result, but all of them have the upshot that programs of type `bool` compute boolean results, and programs of type `num` compute natural number results. Harper [Har16] formulates type safety as the following theorem:

Theorem 3.22 (Type safety à la Harper [Har16]).

1. *If $e : \tau$ and $e \mapsto e'$ then $e' : \tau$.*
2. *If $e : \tau$ then either e val or $e \mapsto e'$.*

The first clause, *preservation* (or *subject reduction*), says that \mapsto preserves typing. The second clause, *progress*, says that well-typed terms are never “stuck”: they either are values or can take a step.

There are two very famous slogans associated to this theorem. The first, popularized by Harper [Pie02, p. 95], is quite self-explanatory: “*type safety is progress + preservation*.” The second, due to Milner [Mil78], is “*well-typed expressions do not go wrong*.” Indeed, Theorem 3.22 implies that every well-typed program has either successfully finished evaluating, or steps to a well-typed program (at which point we are either done, or we take another successful step, *ad infinitum*).

Remark 3.23. Progress may sound like it alone implies that “well-typed programs don’t get stuck,” but it really says that well-typed programs aren’t *immediately* stuck. We need preservation to conclude that well-typed programs never become stuck during evaluation, immediately or otherwise.

However, the approach of using progress and preservation to establish type safety is decades newer than the concept of type safety itself; thus it is reasonable to argue that type safety “is” the idea that well-typed programs do not go wrong, demoting the status of progress + preservation to a pair of convenient lemmas.

Theorem 3.24 (Type soundness à la Wright and Felleisen [WF94]). *If $e : \tau$ and $e \mapsto^* e' \not\mapsto$, then $e' : \tau$ and e' val.*

The formulation in Theorem 3.24 is taken directly from Wright and Felleisen [WF94] who introduced the technique of progress and preservation in 1994, thereby allowing subsequent generations of PL researchers to avoid learning any denotational semantics. Progress and preservation were later famously used in the proof of type safety for Featherweight Java with generics [IPW01], and are now among the most famed and trusted tools in the PL researcher’s toolbox.

Returning to the main story, we prove type safety in three steps: canonical forms (Theorem 3.25), progress (Theorem 3.26), and preservation (Theorem 3.27).

Lemma 3.25 (Canonical forms). *Suppose $v : \tau$ and v val. Then:*

1. If $\tau = \text{bool}$, then $v = \text{true}$ or $v = \text{false}$.
2. If $\tau = \text{num}$, then $v \text{ nat}$ where the nat judgment is defined as:

$$\frac{}{\text{zero nat}} \qquad \frac{v \text{ nat}}{\text{suc}(v) \text{ nat}}$$

Proof. We use rule induction on $v : \tau$ to prove $P(v, \tau) = \text{If } v \text{ val then (1) if } \tau = \text{bool} \text{ then } v = \text{true or } v = \text{false, and (2) if } \tau = \text{num} \text{ then } v \text{ nat.}$ (The following proof is written out in extra detail, hopefully for clarity.)

- Case $\frac{}{\text{true} : \text{bool}}$:

The antecedent true val holds. For (1), the antecedent $\tau = \text{bool}$ holds, and $v = \text{true}$. For (2), the antecedent $\tau = \text{nat}$ is false so the statement holds vacuously.

- Case $\frac{}{\text{false} : \text{bool}}$:

Analogous to previous case.

- Case $\frac{}{\text{zero} : \text{num}}$:

zero val holds. (1) is vacuous; for (2), $\tau = \text{num}$ holds, and zero nat.

- Case $\frac{e : \text{num}}{\text{suc}(e) : \text{num}}$:

Suppose $\text{suc}(e)$ val; then by inversion, we have e val. (1) is vacuous because $\tau \neq \text{bool}$; for (2), $\tau = \text{num}$ holds, and we must show $\text{suc}(e)$ nat holds. By our inductive hypothesis $P(e, \text{num})$, if e val and $\text{num} = \text{num}$ then e nat. Thus e nat and so $\text{suc}(e)$ nat as required.

- Remaining cases (pred, zero?, if): The antecedent v val is false by inversion, so the statement holds vacuously. \square

Lemma 3.26 (Progress). *If $e : \tau$ then either e val or $e \mapsto e'$ for some e' .*

Proof. By rule induction on $e : \tau$.

- Case $\frac{}{\text{true} : \text{bool}}$:

True by true val.

- Case $\frac{\text{false}}{\text{false} : \text{bool}}$:

True by `false` val.

- Case $\frac{\text{zero}}{\text{zero} : \text{num}}$:

True by `zero` val.

- Case $\frac{e : \text{num}}{\text{suc}(e) : \text{num}}$:

We must show that $\text{suc}(e)$ val or $\text{suc}(e) \mapsto e'$. By our inductive hypothesis, either e val or $e \mapsto e''$. In the former case, $\text{suc}(e)$ val; in the latter case, $\text{suc}(e) \mapsto \text{suc}(e'')$.

- Case $\frac{e : \text{num}}{\text{pred}(e) : \text{num}}$:

This is never a value, so we will have to show $\text{pred}(e) \mapsto e'$. By our inductive hypothesis, either e val or $e \mapsto e''$. In the latter case, $\text{pred}(e) \mapsto \text{pred}(e'')$. In the former case, by Theorem 3.25 we have e nat, and complete the proof by inversion on e nat. If $e = \text{zero}$, then $\text{pred}(\text{zero}) \mapsto \text{zero}$. Otherwise, if $e = \text{suc}(v)$ where v nat, then $\text{pred}(\text{suc}(v)) \mapsto v$.

- Case $\frac{e : \text{num}}{\text{zero?}(e) : \text{bool}}$:

Similar to previous case.

- Case $\frac{e_1 : \text{bool} \quad e_2 : \tau \quad e_3 : \tau}{\text{if}(e_1, e_2, e_3) : \tau}$:

This is never a value, so we will have to show $\text{if}(e_1, e_2, e_3) \mapsto e'$. By our inductive hypothesis on e_1 , either e_1 val or $e_1 \mapsto e'_1$. In the latter case, $\text{if}(e_1, e_2, e_3) \mapsto \text{if}(e'_1, e_2, e_3)$. In the former case, by Theorem 3.25 we have $e_1 = \text{true}$ or $e_1 = \text{false}$. If $e_1 = \text{true}$ then $\text{if}(e_1, e_2, e_3) \mapsto e_2$; if $e_1 = \text{false}$ then $\text{if}(e_1, e_2, e_3) \mapsto e_3$. \square

Lemma 3.27 (Preservation). *If $e : \tau$ and $e \mapsto e'$ then $e' : \tau$.*

Proof. By rule induction on $e \mapsto e'$. We focus on a few representative cases.

- Case $\frac{\text{if}(\text{true}, e_2, e_3) \mapsto e_2}{\text{if}(\text{true}, e_2, e_3) : \tau}$:

By inversion on the typing judgment, if $\text{if}(\text{true}, e_2, e_3) : \tau$ then $e_2 : \tau$ and $e_3 : \tau$. Thus in particular $e_2 : \tau$ as required.

- Case $\frac{e_1 \mapsto e'_1}{\text{if}(e_1, e_2, e_3) \mapsto \text{if}(e'_1, e_2, e_3)}$:

By inversion on the typing judgment, if $\text{if}(e_1, e_2, e_3) : \tau$ then $e_1 : \text{bool}$, $e_2 : \tau$, and $e_3 : \tau$. By our inductive hypothesis, $e'_1 : \text{bool}$; the result follows by the typing rule for if.

- Case $\frac{v \text{ val}}{\text{pred}(\text{suc}(v)) \mapsto v}$:

By two inversions on typing, if $\text{pred}(\text{suc}(v)) : \text{num}$ then $v : \text{num}$. \square

Exercise 3.28. Prove the remaining cases of Theorem 3.27.

It follows that if a well-typed term terminates, it terminates in a value of the same type. In other words, programs of type `bool` compute booleans and programs of type `num` compute natural numbers.

Corollary 3.29 (Type soundness). *If $e : \tau$ and $e \mapsto^* e' \mapsto$, then:*

- If $\tau = \text{bool}$, then $e' = \text{true}$ or $e' = \text{false}$.
- If $\tau = \text{num}$, then $e' \text{ nat}$.

Proof. By rule induction on $e \mapsto^* e'$.

- Case $\frac{}{e \mapsto^* e}$:

By Theorem 3.26, either $e \text{ val}$ or $e \mapsto e'$, but the latter is impossible by our hypothesis $e \mapsto$. The result follows by Theorem 3.25.

- Case $\frac{e \mapsto e' \quad e' \mapsto^* e''}{e \mapsto^* e''}$:

By Theorem 3.27, $e' : \tau$. The result follows by our inductive hypothesis. \square

Remark 3.30. Type safety tells us that well-typed terms do not get stuck, but it allows for the possibility that a well-typed term will get into an infinite loop without ever reaching a value. For this particular language it is easy to establish that all programs (even ill-typed ones!) terminate—every \mapsto rule shrinks the size of the term—but in more interesting languages we will need to resort to more sophisticated techniques to prove termination.

5 Evaluation contexts

Add more words in this section.

We can rephrase the operational semantics in Figure 3.1 to group all the congruence reductions into a single rule. Also called *contextual dynamics*, *context-sensitive reduction semantics*, and *reduction semantics*.

Evaluation contexts $\mathcal{E} ::= \circ \mid \text{if}(\mathcal{E}, e, e) \mid \text{suc}(\mathcal{E}) \mid \text{pred}(\mathcal{E}) \mid \text{zero?}(\mathcal{E})$

Definition 3.31. For any evaluation context \mathcal{E} and e tm, we define the instantiation of \mathcal{E} with e , written $\mathcal{E}\{e\}$, by structural recursion:

$$\begin{aligned}\circ\{e\} &= e \\ \text{if}(\mathcal{E}, e_2, e_3)\{e\} &= \text{if}(\mathcal{E}\{e\}, e_2, e_3) \\ \text{suc}(\mathcal{E})\{e\} &= \text{suc}(\mathcal{E}\{e\}) \\ \text{pred}(\mathcal{E})\{e\} &= \text{pred}(\mathcal{E}\{e\}) \\ \text{zero?}(\mathcal{E})\{e\} &= \text{zero?}(\mathcal{E}\{e\})\end{aligned}$$

Isolate the principal reductions $e \xrightarrow{p} e'$ and define \xrightarrow{p} in one rule.

$$\begin{array}{c} \frac{e \xrightarrow{p} e'}{\mathcal{E}\{e\} \xrightarrow{p} \mathcal{E}\{e'\}} \quad \frac{}{\text{if}(\text{true}, e_2, e_3) \xrightarrow{p} e_2} \quad \frac{}{\text{if}(\text{false}, e_2, e_3) \xrightarrow{p} e_3} \\ \\ \frac{}{\text{pred}(\text{zero}) \xrightarrow{p} \text{zero}}^* \quad \frac{v \text{ val}}{\text{pred}(\text{suc}(v)) \xrightarrow{p} v} \\ \\ \frac{}{\text{zero?}(\text{zero}) \xrightarrow{p} \text{true}} \quad \frac{v \text{ val}}{\text{zero?}(\text{suc}(v)) \xrightarrow{p} \text{false}} \end{array}$$

6 Well-typed programs *can* go wrong

“Type errors” like $\text{zero?}(\text{true})$ are a major source of going wrong in a programming language, but sometimes things can go wrong even with the guardrails of a type system. Consider division by zero: there’s no number for division by zero to step to, but statically ruling it out would require a type system that can detect which numerical expressions cannot be zero at runtime.

This may suggest that type safety does not hold for languages with division, but in fact we can refine the statement of type safety to account for this. The idea

is to add a second kind of “outcome” to our language: in addition to evaluating to a value, a program may evaluate to a well-defined *error* state. This is distinct from a program failing to have a next state. We can think of the well-defined error states as *checked errors*, or errors that a user is forewarned might occur, and stuck terms as encountering *unchecked errors*, ones that indicate a gap in the language itself.

To see an example of this in action, let’s introduce a new judgment $e \text{ err}$ indicating that e evaluates no further because it has encountered the “subtraction from zero” error. We can adjust the contextual dynamics of the previous section by deleting the principal reduction labeled \star and replacing it with two new rules:

$$\frac{}{\text{pred(zero)} \xrightarrow{p} \text{zero}} \star \quad \rightsquigarrow \quad \frac{}{\text{pred(zero)} \text{ err}} \quad \frac{e \text{ err}}{\mathcal{E}\{e\} \text{ err}}$$

The first rule says that pred(zero) raises the “subtraction from zero” error, and the second rule propagates this error through evaluation contexts.

The statements and proofs of the canonical forms and preservation lemmas remain unchanged, but we must adjust progress to account for this error:

Lemma 3.32 (Progress). *If $e : \tau$ then either $e \text{ val}$ or $e \text{ err}$ or $e \xrightarrow{} e'$ for some e' .*

Type safety then still tells us that programs never get stuck: the possible outcomes are now that a program terminates successfully with a value, or terminates unsuccessfully with the “subtraction from zero” error, or gets into an infinite loop.⁵

We will study this further in the next problem set.

References

- [Ben+93] Nick Benton, Gavin Bierman, Valeria de Paiva, and Martin Hyland. “Linear λ -calculus and categorical models revisited”. In: *Computer Science Logic*. Ed. by E. Börger, G. Jäger, H. Kleine Bünning, S. Martini, and M. M. Richter. Berlin, Heidelberg: Springer Berlin Heidelberg, 1993, pp. 61–84. ISBN: 978-3-540-47890-4.
- [Har16] Robert Harper. *Practical Foundations for Programming Languages*. Second Edition. Cambridge University Press, 2016. ISBN: 9781107150300. doi: [10.1017/CBO9781316576892](https://doi.org/10.1017/CBO9781316576892).
- [IPW01] Atsushi Igarashi, Benjamin C. Pierce, and Philip Wadler. “Featherweight Java: a minimal core calculus for Java and GJ”. In: *ACM Transactions on Programming Languages and Systems* 23.3 (May 2001), pp. 396–450. ISSN: 0164-0925. doi: [10.1145/503502.503505](https://doi.org/10.1145/503502.503505).

⁵Again, there are no infinite loops in this language, although type safety does not tell us this.

[Mil78] Robin Milner. “A theory of type polymorphism in programming”. In: *Journal of Computer and System Sciences* 17.3 (1978), pp. 348–375. ISSN: 0022-0000. doi: [10.1016/0022-0000\(78\)90014-4](https://doi.org/10.1016/0022-0000(78)90014-4).

[Pie02] Benjamin C. Pierce. *Types and Programming Languages*. MIT Press, 2002. ISBN: 0-262-16209-1.

[WF94] A.K. Wright and M. Felleisen. “A Syntactic Approach to Type Soundness”. In: *Information and Computation* 115.1 (Nov. 1994), pp. 38–94. ISSN: 0890-5401. doi: [10.1006/inco.1994.1093](https://doi.org/10.1006/inco.1994.1093).