

Xi Type System Specification

Computer Science 4120
Cornell University

Version of October 2, 2011

Changes

- October 2: Made ARRAYDECL more general to match up with the language spec.
- September 30: Fixed bugs reported in recitation and some more. Tuple expressions and types may have 0 or 1 elements.
- September 29: Fixed syntax in ARRASSIGN.

Types

The Xi type system uses a somewhat bigger set of types than can be expressed explicitly in the source language:

$$\begin{aligned} \tau &::= \text{int} \\ &| \text{bool} \\ &| \tau[] \\ R &::= \text{unit} \mid \text{void} \\ T &::= \tau \\ &| (\tau_1, \tau_2, \dots, \tau_n) \quad (n \geq 0) \\ &| R \\ \sigma &::= \text{var } \tau \\ &| T_1 \rightarrow T_2 \end{aligned}$$

Ordinary types expressible in the language are denoted by the metavariable τ , which can be `int`, `bool`, or an array type.

The metavariable R represent the outcome of evaluating a statement. It can be either `unit` or `void`.

The type `unit` is used in various ways. It is the type of ordinary statements. It gives a type to the left-hand side of pattern-matching assignments that use the `_` placeholder, which lets their handling be integrated directly into the type system. The `unit` type is also used to represent the result type of procedures and the type of statements that may complete normally and permit the following statement to execute.

The type `void` is the type of statements such as `return` that pass control to a following statement. It should not be confused with the C type `void`, which is actually closer to `unit`.

The metavariable T denotes an expanded notion of type that includes the parameter types and return types of procedures and functions. It may be an ordinary type, a tuple type, `unit`, or `void`.

Tuple types represent the parameters or return types of functions that take multiple argument or return multiple results.

The set σ is used to represent typing environment entries, which can either be normal variables (bound to `var τ` for some type τ) or functions (bound to `$\tau \rightarrow \tau'$` where $\tau' \neq \text{unit}$), or procedures (bound to

$\tau \rightarrow \text{unit}$), where the “result type” (unit) indicates that the procedure result contains no information other than that the procedure call terminated.

Subtyping

The subtyping relation on T is the least partial order consistent with this rule:

$$\overline{\tau \leq \text{unit}}$$

However, there is not (for now) a subsumption rule, so subtyping matters only where it appears explicitly.

Type-checking expressions

To type-check expressions, we need to know what bound variables and functions are in scope; this is represented by the typing context Γ , which maps names x to types σ .

The judgment $\Gamma \vdash e : \tau$ is the rule for the type of an expression; it states that with bindings Γ we can conclude that e has the type τ .

We use the metavariable symbols x or f to represent arbitrary identifiers, n to represent a numeric constant, *string* to represent a literal string constant, and *char* to represent a literal character constant. Using these conventions, the expression typing rules are:

$$\overline{\Gamma \vdash n : \text{int}} \quad \overline{\Gamma \vdash \text{true} : \text{bool}} \quad \overline{\Gamma \vdash \text{false} : \text{bool}} \quad \overline{\Gamma \vdash \text{string} : \text{int}[]} \quad \overline{\Gamma \vdash \text{char} : \text{int}}$$

$$\frac{\Gamma(x) = \text{var } \tau}{\Gamma \vdash x : \tau} \quad \frac{\Gamma \vdash e_1 : \text{int} \quad \Gamma \vdash e_2 : \text{int} \quad \oplus \in \{+, -, /, *, \%\}}{\Gamma \vdash e_1 \oplus e_2 : \text{int}}$$

$$\frac{\Gamma \vdash e : \text{int}}{\Gamma \vdash -e : \text{int}} \quad \frac{\Gamma \vdash e_1 : \text{int} \quad \Gamma \vdash e_2 : \text{int} \quad \ominus \in \{==, !=, <, <=, >, >=\}}{\Gamma \vdash e_1 \ominus e_2 : \text{bool}}$$

$$\frac{\Gamma \vdash e : \text{bool}}{\Gamma \vdash !e : \text{bool}} \quad \frac{\Gamma \vdash e_1 : \text{bool} \quad \Gamma \vdash e_2 : \text{bool} \quad \ominus \in \{==, !=, \&, \}}{\Gamma \vdash e_1 \ominus e_2 : \text{bool}} \quad \frac{\Gamma \vdash e : \tau[]}{\Gamma \vdash \text{length } e : \text{int}}$$

$$\frac{\Gamma \vdash e_1 : \tau[] \quad \Gamma \vdash e_2 : \tau[] \quad \ominus \in \{==, !=\}}{\Gamma \vdash e_1 \ominus e_2 : \text{bool}} \quad \frac{\Gamma \vdash e_1 : \tau \quad \dots \quad \Gamma \vdash e_n : \tau \quad n \geq 0}{\Gamma \vdash (e_1, \dots, e_n) : \tau[]}$$

$$\frac{\Gamma \vdash e_1 : \tau[] \quad \Gamma \vdash e_2 : \tau[]}{\Gamma \vdash e_1 + e_2 : \tau[]} \quad \frac{\Gamma \vdash e_1 : \tau[] \quad \Gamma \vdash e_2 : \text{int}}{\Gamma \vdash e_1[e_2] : \tau} \quad \frac{\Gamma(f) = (\tau_1, \dots, \tau_n) \rightarrow \tau' \quad \Gamma \vdash e_i : \tau_i \quad (\forall i \in 1..n)}{\Gamma \vdash f(e_1, \dots, e_n) : \tau'}$$

Type-checking statements

To type-check statements, we need all the information used to type-check expressions, plus the types of procedures, which are included in Γ . In addition, we extend the domain of Γ a little to include two special symbols, ρ and β . To check the return statement we need to know what the return type of the current function is or if it is a procedure. Let this be denoted by $\Gamma(\rho)$, which is some type τ if the statement is part of a function, or unit if the statement is in a procedure. For **break** statements, we also need to check whether we are inside a loop, which we will denote as $\Gamma(\beta)$, which is unit if we are inside a loop and void if we are not. Since statements include declarations, they can also produce new variable bindings, resulting in an updated typing context which we will denote as Γ' . To update typing contexts, we write $\Gamma[x \mapsto \tau]$, which is an environment exactly like Γ except that it maps x to τ . We use the metavariable s to denote

a statement, so the main typing judgment for statements has either the form $\Gamma \vdash s : \text{unit}, \Gamma'$ or the form $\Gamma \vdash s : \text{void}, \Gamma'$.

Most of the statements are fairly straightforward, and do not change Γ . However, statements like `break` and `return` are a bit tricky because their type is `void`.

$$\begin{array}{c}
\frac{\Gamma \vdash e : \text{bool} \quad \Gamma \vdash s : R, \Gamma'}{\Gamma \vdash \text{if } (e) s : \text{unit}, \Gamma} \text{ (IF)} \qquad \frac{\Gamma \vdash e : \text{bool} \quad \Gamma \vdash s_1 : R_1, \Gamma' \quad \Gamma \vdash s_2 : R_2, \Gamma''}{\Gamma \vdash \text{if } (e) s_1 \text{ else } s_2 : \text{lub}(R_1, R_2), \Gamma} \text{ (IFELSE)} \\
\frac{\Gamma \vdash e : \text{bool} \quad \Gamma[\beta \mapsto \text{unit}] \vdash s : R, \Gamma'}{\Gamma \vdash \text{while } (e) s : \text{unit}, \Gamma} \text{ (WHILE)} \\
\frac{\Gamma \vdash s_1 : \text{unit}, \Gamma_1 \quad \Gamma_1 \vdash s_2 : \text{unit}, \Gamma_2 \quad \dots \quad \Gamma_{n-2} \vdash s_{n-1} : \text{unit}, \Gamma_{n-1} \quad \Gamma_{n-1} \vdash s_n : R, \Gamma_n}{\Gamma \vdash \{s_1; s_2; \dots; s_n\} : R, \Gamma_n} \text{ (BLOCK)} \\
\frac{\Gamma(f) = \tau \rightarrow \text{unit} \quad \Gamma \vdash e : \tau}{\Gamma \vdash f(e_1, \dots, e_n) : \text{unit}, \Gamma} \text{ (PRCALL)} \qquad \frac{\Gamma(\beta) = \text{unit}}{\Gamma \vdash \text{break} : \text{void}, \Gamma} \text{ (BREAK)} \\
\frac{\Gamma(\rho) = \text{unit}}{\Gamma \vdash \text{return} : \text{void}, \Gamma} \text{ (RETURN)} \qquad \frac{\Gamma(\rho) = \tau \neq \text{unit} \quad \Gamma \vdash e : \tau}{\Gamma \vdash \text{return } e : \text{void}, \Gamma} \text{ (RETVAl)}
\end{array}$$

The function *lub* is defined as $\text{lub}(R, R) = R$ and $\text{lub}(\text{unit}, R) = \text{lub}(R, \text{unit}) = \text{unit}$. Therefore the type of an `if` is `void` only if both branches have that type.

Assignments require checking the left-hand side to make sure it is assignable:

$$\frac{\Gamma(x) = \text{var } \tau \quad \Gamma \vdash e : \tau}{\Gamma \vdash x = e : \text{unit}, \Gamma} \text{ (ASSIGN)} \qquad \frac{\Gamma \vdash e_1 : \tau[] \quad \Gamma \vdash e_2 : \text{int} \quad \Gamma \vdash e_3 : \tau}{\Gamma \vdash e_1[e_2] = e_3 : \text{unit}, \Gamma} \text{ (ARRASSIGN)}$$

Declarations are the source of new bindings. Three kinds of declarations can appear in the source language: regular variable declarations, tuple declarations, and function/procedure declarations. We are only concerned with the first two kinds within a function body. To handle tuples, we define a declaration d that can appear within a tuple:

$$d ::= x : \tau \mid _$$

and define functions *typeof*(d) and *varsof*(d) as follows: $\text{typeof}(x : \tau) = \tau$ and $\text{typeof}(_) = \text{unit}$, and $\text{varsof}(x : \tau) = \{x\}$ and $\text{varsof}(_) = \emptyset$. Using these notations, we have the following rules:

$$\frac{x \notin \text{dom}(\Gamma)}{\Gamma \vdash x : \tau : \Gamma[x \mapsto \tau]} \text{ (VARDECL)} \qquad \frac{x \notin \text{dom}(\Gamma) \quad \Gamma \vdash e : \tau}{\Gamma \vdash x : \tau = e : \Gamma[x \mapsto \tau]} \text{ (VARINIT)}$$

The following rule for array declarations with specified sizes is intended to capture the essence of type checking them, though we have sacrificed a bit of formal precision for the sake of readability. The declared variable is added to the typing context with a type with the same number of array dimensions.

$$\frac{\frac{x \notin \text{dom}(\Gamma) \quad \Gamma \vdash e_i : \text{int} \quad (\forall i \in 1..n)}{\Gamma \vdash x : \tau[e] [] : \Gamma[x \mapsto \tau [] []]} \text{ (ARRAYDECL)} \quad \frac{\Gamma(f) = (\tau'_1, \dots, \tau'_m) \rightarrow (\tau_1, \dots, \tau_n) \quad \Gamma \vdash e_i : \tau'_i \quad (\forall i \in 1..m) \quad \tau_i \leq \text{typeof}(d_i) \quad (\forall i \in 1..n)}{\text{dom}(\Gamma) \cap \text{varsof}(d_i) = \emptyset \quad (\forall i \in 1..n) \quad \text{varsof}(d_i) \cap \text{varsof}(d_j) = \emptyset \quad (\forall i, j \in 1..n | j \neq i)} \text{ (TUPLEDECL)}}{\Gamma \vdash d_1, \dots, d_n = f(e_1, \dots, e_m) : \Gamma[x_i \mapsto \text{typeof}(d_i) \quad (\forall i \in 1..n, x_i \mid \text{varsof}(d_i) = \{x_i\})} \text{ (TUPLEDECL)}$$

The final premise in rule `TUPLEDECL` prevents shadowing by ensuring that $\text{dom}(\Gamma)$ and all of the $\text{varsof}(d_i)$ are disjoint from each other.

Top-level declarations

At the top level of the program, we need to figure out the types of procedures and functions, and make sure their bodies are well-typed. Since mutual recursion is supported, this needs to be done in two passes. First, we use the judgment $\Gamma \vdash fd : \Gamma'$ to state that the function or procedure declaration fd extends top-level bindings Γ to Γ' :

$$\frac{f \notin \text{dom}(\Gamma)}{\Gamma \vdash f(x:\tau) : \tau' = s : \Gamma[f \mapsto \tau \rightarrow \tau']}$$

$$\frac{f \notin \text{dom}(\Gamma) \quad n \geq 2 \quad \Gamma' = \Gamma[f \mapsto (\tau_1, \tau_2, \dots, \tau_n) \rightarrow \tau_r]}{\Gamma \vdash f(x_1:\tau_1, x_2:\tau_2, \dots, x_n:\tau_n) : \tau_r = s : \Gamma'}$$

$$\frac{f \notin \text{dom}(\Gamma)}{\Gamma \vdash f(x:\tau) = s : \Gamma[f \mapsto \tau \rightarrow \text{unit}]}$$

$$\frac{f \notin \text{dom}(\Gamma) \quad n \geq 2 \quad \Gamma' = \Gamma[f \mapsto (\tau_1, \tau_2, \dots, \tau_n) \rightarrow \text{unit}]}{\Gamma \vdash f(x_1:\tau_1, x_2:\tau_2, \dots, x_n:\tau_n) = s : \Gamma'}$$

The second pass over the program is captured by the judgment $\Gamma \vdash fd \text{ def}$, which defines how to check well-formedness of each function definition against a top-level environment Γ , ensuring that parameters do not shadow anything and that the body is well-typed. The body of a function definition must have type `void`, which ensures that the function body does not fall off the end without returning. We treat procedures just like functions that return the unit type. Therefore their bodies are allowed to have type `unit`.

$$\frac{x \notin \text{dom}(\Gamma) \quad \Gamma[x \mapsto \tau, \rho \mapsto \tau', \beta \mapsto \text{void}] \vdash s : \text{void}, \Gamma'}{\Gamma \vdash f(x:\tau) : \tau' = s \text{ def}}$$

$$\frac{|\text{dom}(\Gamma) \cup \{x_1, \dots, x_n\}| = |\text{dom}(\Gamma)| + n \quad \Gamma[x_1 \mapsto \tau_1, \dots, x_n \mapsto \tau_n, \rho \mapsto \tau', \beta \mapsto \text{void}] \vdash s : \text{void}, \Gamma'}{\Gamma \vdash f(x_1:\tau_1, \dots, x_n:\tau_n) : \tau' = s \text{ def}}$$

$$\frac{x \notin \text{dom}(\Gamma) \quad \Gamma[x \mapsto \tau, \rho \mapsto \text{unit}, \beta \mapsto \text{void}] \vdash s : \text{unit}, \Gamma'}{\Gamma \vdash f(x:\tau) = s \text{ def}}$$

$$\frac{|\text{dom}(\Gamma) \cup \{x_1, \dots, x_n\}| = |\text{dom}(\Gamma)| + n \quad \Gamma[x_1 \mapsto \tau_1, \dots, x_n \mapsto \tau_n, \rho \mapsto \text{unit}, \beta \mapsto \text{void}] \vdash s : \text{unit}, \Gamma'}{\Gamma \vdash f(x_1:\tau_1, \dots, x_n:\tau_n) = s \text{ def}}$$

Checking a program

Using the previous judgments, we can define when an entire program $fd_1 \dots fd_n$ is well-formed, written $\vdash fd_1 \dots fd_n \text{ prog}$:

$$\frac{\emptyset \vdash d_1 : \Gamma_1 \quad \Gamma_1 \vdash d_2 : \Gamma_2 \quad \dots \quad \Gamma_{n-1} \vdash d_n : \Gamma_n \quad \Gamma_n \vdash d_1 \text{ def} \quad \Gamma_n \vdash d_2 \text{ def} \quad \dots \quad \Gamma_n \vdash d_n \text{ def}}{\vdash fd_1 fd_2 \dots fd_n \text{ prog}}$$