

COS320: Compiling Techniques

Zak Kincaid

March 22, 2022

Logistics

- HW3 due **today**
- HW4 released today, due April 11th. You will implement a typechecker and translator for an extension of Oat.

Oat v2

- Specified by a (fairly large) type system
 - ~20 judgements, ~80 inference rules
 - Invest some time in making sure you understand how to read them
- Adds several features to the Oat language:
 - Memory safety
 - *nullable* and *non-null* references. Type system enforces no null pointer dereferences.
 - Run-time array bounds checking (like Java, OCaml)
 - Mutable record types
 - Subtyping
 - *ref <: ref?*: non-null references are a subtype of nullable references
 - Record subtyping: width but not depth (*why?*)

Compiling with Types

- Intrinsic view: an ill-typed program is not a program at all
- Compiler translates programs in the source language to programs in the target language
 - Well-typed source programs translate to well-typed target programs

- Intrinsic view: an ill-typed program is not a program at all
- Compiler translates programs in the source language to programs in the target language
 - Well-typed source programs translate to well-typed target programs
 - Compiler may reject ill-typed source programs

- Intrinsic view: an ill-typed program is not a program at all
- Compiler translates programs in the source language to programs in the target language
 - **Well-typed** source programs translate to **well-typed** target programs
 - Compiler may reject ill-typed source programs
 - Compiler must ensure that target program is well-typed

- Intrinsic view: an ill-typed program is not a program at all
- Compiler translates programs in the source language to programs in the target language
 - Well-typed source programs translate to well-typed target programs
 - Compiler may reject ill-typed source programs
 - Compiler must ensure that target program is well-typed
- IR may also have its own type system (LLVM)
 - Your backend does not check types, but does throw exceptions for (some) ill-typed programs
 - LLVM does check types: use `--clang` to check that your front-end produces type-correct code

We can think of compilation as translation of derivations of judgements from a source language to a target language

- Each kind of judgement has a different translation category. E.g.,
 - Well-formed types in source become well-formed types in target
 - Expressions in source become (operand, instruction list) pairs in target
 - ...
- Each inference rule corresponds to a case within that category

Oat v1 (HW3) – well-formed types

Judgements take the form:

- $\vdash t$: “ t is a well-formed type” (ty)
- $\vdash_r \mathit{ref}$: “ ref is a well-formed reference type” (rty)
- $\vdash_{rt} \mathit{rt}$: “ rt is a well-formed return type” (ret_ty)

$$\begin{array}{c} \text{TINT} \\ \hline \vdash \text{int} \end{array} \qquad \begin{array}{c} \text{TBOOL} \\ \hline \vdash \text{bool} \end{array} \qquad \begin{array}{c} \text{TREF} \\ \vdash_r \mathit{ref} \\ \hline \vdash \mathit{ref} \end{array}$$

$$\begin{array}{c} \text{RSTRING} \\ \hline \vdash_r \text{string} \end{array} \qquad \begin{array}{c} \text{RARRAY} \\ \vdash t \\ \hline \vdash_r t[] \end{array} \qquad \begin{array}{c} \text{RFUN} \\ \vdash t_1 \quad \dots \quad \vdash t_n \quad \vdash_{rt} \mathit{rt} \\ \hline \vdash_r (t_1, \dots, t_n) \rightarrow \mathit{rt} \end{array}$$

$$\begin{array}{c} \text{RTVOID} \\ \hline \vdash_{rt} \text{void} \end{array} \qquad \begin{array}{c} \text{RTTYP} \\ \vdash t \\ \hline \vdash_{rt} t \end{array}$$

LLVMlite well-formed types

Judgements take the form:

- $T \vdash t$: With named types T , t is a well-formed type
- $T \vdash_s t$: With named types T , t is a well-formed simple type
- $T \vdash_r t$: With named types T , t is a well-formed reference type
- $T \vdash_{rt} t$: With named types T , t is a well-formed return type

LLBOOL

$$\frac{}{T \vdash_s \text{i1}}$$

LLINT

$$\frac{}{T \vdash_s \text{i64}}$$

LLPTR

$$\frac{T \vdash_r \text{ref}}{T \vdash_s \text{ref}*}$$

LLTUPLE

$$\frac{T \vdash t_1 \quad \dots \quad T \vdash t_n}{T \vdash \{t_1, \dots, t_n\}}$$

LLARRAY

$$\frac{T \vdash t}{T \vdash [n \times t]} \quad n \in \mathbb{N}$$

LLSIMPLE

$$\frac{\vdash_s t}{\vdash t}$$

LLRTVOID

$$\frac{}{T \vdash_{rt} \text{void}}$$

LLRTSIMPLE

$$\frac{T \vdash_s t}{T \vdash_{rt} t}$$

LLRCHAR

$$\frac{}{T \vdash_r \text{i8}}$$

LLRTYPE

$$\frac{T \vdash t}{T \vdash_r t}$$

LLRFUN

$$\frac{T \vdash_{rt} \text{rt} \quad T \vdash_s t_1 \quad \dots \quad T \vdash_s t_n}{T \vdash_r \text{rt}(t_1, \dots, t_n)}$$

LLNAMED

$$\frac{}{T \vdash \%uid} \quad \%uid \in T$$

Translating well-formed types

- Each well-formed Oat type is translated to a well-formed LLVM type
 - types \rightarrow simple types (`cmp_ty`)
 - reference types \rightarrow reference types (`cmp_rty`)
 - return types \rightarrow return types (`cmp_ret_ty`)
- Use \rightsquigarrow to denote translation of derivations

Translating well-formed types

Suppose we have a well-formed type $\text{Oat type}, \vdash t$. There are three inference rules:

TINT

$$\frac{}{\vdash \text{int}}$$

TBOOL

$$\frac{}{\vdash \text{bool}}$$

TREF

$$\frac{\vdash_r \text{ref}}{\vdash \text{ref}}$$

Each has a corresponding case:

- $\left(\text{TINT} \frac{}{\vdash \text{int}} \right) \rightsquigarrow \left(\text{LLINT} \frac{}{\vdash_s \text{i64}} \right)$
- $\left(\text{TBOOL} \frac{}{\vdash \text{bool}} \right) \rightsquigarrow \left(\text{LLBOOL} \frac{}{\vdash_s \text{i1}} \right)$

Translating well-formed types

Suppose we have a well-formed type Oat type, $\vdash t$. There are three inference rules:

$$\frac{}{\vdash \text{int}} \text{TINT}$$

$$\frac{}{\vdash \text{bool}} \text{TBOOL}$$

$$\frac{\vdash_r \text{ref}}{\vdash \text{ref}} \text{TREF}$$

Each has a corresponding case:

- $\left(\frac{}{\vdash \text{int}} \text{TINT} \right) \rightsquigarrow \left(\frac{}{\vdash_s \text{i64}} \text{LLINT} \right)$
- $\left(\frac{}{\vdash \text{bool}} \text{TBOOL} \right) \rightsquigarrow \left(\frac{}{\vdash_s \text{i1}} \text{LLBOOL} \right)$
- $\left(\frac{\vdash_r \text{ref}}{\vdash \text{ref}} \text{TREF} \right) \rightsquigarrow \left(\frac{\vdash_r t}{\vdash_s t^*} \text{LLPTR} \right)$, where $(\vdash_r \text{ref}) \rightsquigarrow (\vdash_r t)$

Translating well-formed array types

- In Oat v2, arrays accesses are checked at runtime
- **Recall:** Can implement run-time array access checking by allocating additional memory at the beginning of the array to store its size
- In Oat v1, arrays accesses are unchecked, but for forwards-compatibility we represent arrays in the same way.

Translating well-formed array types

- In Oat v2, arrays accesses are checked at runtime
- Recall:** Can implement run-time array access checking by allocating additional memory at the beginning of the array to store its size
- In Oat v1, arrays accesses are unchecked, but for forwards-compatibility we represent arrays in the same way.

$$\begin{array}{c}
 \text{RARRAY} \\
 \frac{\vdash t}{\vdash_r t[]} \rightsquigarrow \\
 \text{LLRTYPE} \frac{\text{LLTUPLE} \frac{\text{LLSIMPLE} \frac{\text{LLINT} \frac{}{\vdash_s i64}}{\vdash_s i64} \quad \text{LLSIMPLE} \frac{\text{LLSIMPLE} \frac{\vdash_s t'}{\vdash t'}}{\vdash t'}}{\vdash [0xt']}}{\vdash \{i64, [0xt']\}}}{\vdash_r \{i64, [0xt']\}}
 \end{array}$$

where $\vdash t \rightsquigarrow \vdash_s t'$

Summary of type translation

Succinct notation: $\llbracket \vdash J \rrbracket = J'$ denotes that a derivation with root J translates to a derivation with root J'

- $\llbracket \vdash \text{int} \rrbracket = \vdash_s \text{i64}$
- $\llbracket \vdash \text{bool} \rrbracket = \vdash_s \text{i1}$
- $\llbracket \vdash \text{ref} \rrbracket = \vdash_s t^*$, where $\vdash_r t = \llbracket \vdash_r \text{ref} \rrbracket$
- $\llbracket \vdash_r \text{string} \rrbracket = \vdash_r \text{i8}$
- $\llbracket \vdash_r t[\] \rrbracket = \vdash_r \{ \text{i64}, [\text{0x}t'] \}$, where $\vdash_s t' = \llbracket \vdash t \rrbracket$
- $\llbracket \vdash_r (t_1, \dots, t_n) \rightarrow \text{rt} \rrbracket = \vdash_{rt} \text{rt}'(t'_1, \dots, t'_n)$, where
 - $\vdash_{rt} \text{rt} = \llbracket \vdash_{rt} \text{rt} \rrbracket$,
 - $\vdash_s t'_1 = \llbracket \vdash t_1 \rrbracket, \dots, \vdash_s t'_n = \llbracket \vdash t_n \rrbracket$
- $\llbracket \vdash_{rt} \text{void} \rrbracket = \vdash_{rt} \text{void}$
- $\llbracket \vdash_{rt} t \rrbracket = \vdash_{rt} t$, where $\vdash_s t = \llbracket \vdash t \rrbracket$

(see: `cmp_ty`, `cmp_rty`, `cmp_ret_ty` in HW3)

Well-formed codestreams

Judgements take the form

- $\Gamma \vdash s \Rightarrow \Gamma'$: “under type environment Γ , code stream s is well-formed and results in type environment Γ' ”
- $\Gamma \vdash \text{opn} : t$: “under type environment Γ , operand opn has type t ”

ID

$$\frac{}{\Gamma \vdash \text{id} : t} \quad \Gamma(\text{id}) = t$$

NUM

$$\frac{}{\Gamma \vdash n : \text{i64}} \quad n \in \mathbb{Z}$$

ADD

$$\frac{\Gamma \vdash \text{opn}_1 : \text{i64} \quad \Gamma \vdash \text{opn}_2 : \text{i64}}{\Gamma \vdash \%_{\text{uid}} = \text{add i64 opn}_1, \text{opn}_2 \Rightarrow \Gamma\{\%_{\text{uid}} \mapsto \text{i64}\}} \quad \%_{\text{uid}} \notin \text{dom}(\Gamma)$$

SEQ

$$\frac{\Gamma \vdash s_1 \Rightarrow \Gamma' \quad \Gamma' \vdash s_2 \Rightarrow \Gamma''}{\Gamma \vdash s_1, s_2 \Rightarrow \Gamma''}$$

BASE

$$\frac{}{\Gamma \vdash \epsilon \Rightarrow \Gamma}$$

... lots more

Well-typed expressions

$$\text{VAR} \quad \frac{}{\Gamma \vdash x : \Gamma(x)}$$

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

...

Expression compilation (`cmp_exp`) translates a type judgement $\Gamma \vdash e : t$ to

- A **codestream judgement** $\Gamma_u \vdash s \Rightarrow \Gamma'_u$, and
- An **operand judgement** $\Gamma'_u \vdash \text{opn} : t_u$

How can translate $\Gamma \vdash x : t$ (i.e., VAR)?

How can translate $\Gamma \vdash x : t$ (i.e., VAR)?

- Need a symbol table `ctxt`, which maps Oat identifiers to LLVMlite operand judgements
 - The operand associated with a variable x is a *pointer* to the memory location associated with x

How can translate $\Gamma \vdash x : t$ (i.e., VAR)?

- Need a symbol table ctxt , which maps Oat identifiers to LLVMlite operand judgements
 - The operand associated with a variable x is a *pointer* to the memory location associated with x
- To compute $\llbracket \Gamma \vdash x : t \rrbracket(\text{ctxt})$, first let $(id, t^*) = \text{ctxt}(x)$, then:
 - Define $\llbracket \text{ctxt} \rrbracket$ to be the (LLVM) type environment associated with ctxt
 - $\llbracket \epsilon \rrbracket = \epsilon$ (empty context translates to empty context)
 - $\llbracket \text{ctxt}, x \mapsto (id, t) \rrbracket = \Gamma_u, id \mapsto t$, where $\llbracket \text{ctxt} \rrbracket = \Gamma_u$
 - Codestream: $\llbracket \text{ctxt} \rrbracket \vdash \%uid = \text{load } t, t^* id \Rightarrow \llbracket \text{ctxt} \rrbracket \{ \%uid \mapsto t \}$
 - Operand: $\llbracket \text{ctxt} \rrbracket \{ \%uid \mapsto t \} \vdash \%uid : t$

How can translate $\Gamma \vdash x : t$ (i.e., VAR)?

- Need a symbol table ctxt , which maps Oat identifiers to LLVMlite operand judgements
 - The operand associated with a variable x is a *pointer* to the memory location associated with x
- To compute $\llbracket \Gamma \vdash x : t \rrbracket(\text{ctxt})$, first let $(id, t^*) = \text{ctxt}(x)$, then:
 - Define $\llbracket \text{ctxt} \rrbracket$ to be the (LLVM) type environment associated with ctxt
 - $\llbracket \epsilon \rrbracket = \epsilon$ (empty context translates to empty context)
 - $\llbracket \text{ctxt}, x \mapsto (id, t) \rrbracket = \Gamma_u, id \mapsto t$, where $\llbracket \text{ctxt} \rrbracket = \Gamma_u$
 - Codestream: $\llbracket \text{ctxt} \rrbracket \vdash \%uid = \text{load } t, t^* id \Rightarrow \llbracket \text{ctxt} \rrbracket \{ \%uid \mapsto t \}$
 - Operand: $\llbracket \text{ctxt} \rrbracket \{ \%uid \mapsto t \} \vdash \%uid : t$

How can we translate $\Gamma \vdash e_1 + e_2 : \text{int}$ (i.e., ADD)?

How can translate $\Gamma \vdash x : t$ (i.e., VAR)?

- Need a symbol table ctxt , which maps Oat identifiers to LLVMlite operand judgements
 - The operand associated with a variable x is a *pointer* to the memory location associated with x
- To compute $\llbracket \Gamma \vdash x : t \rrbracket(\text{ctxt})$, first let $(id, t^*) = \text{ctxt}(x)$, then:
 - Define $\llbracket \text{ctxt} \rrbracket$ to be the (LLVM) type environment associated with ctxt
 - $\llbracket \epsilon \rrbracket = \epsilon$ (empty context translates to empty context)
 - $\llbracket \text{ctxt}, x \mapsto (id, t) \rrbracket = \Gamma_u, id \mapsto t$, where $\llbracket \text{ctxt} \rrbracket = \Gamma_u$
 - Codestream: $\llbracket \text{ctxt} \rrbracket \vdash \%uid = \text{load } t, t^* id \Rightarrow \llbracket \text{ctxt} \rrbracket \{ \%uid \mapsto t \}$
 - Operand: $\llbracket \text{ctxt} \rrbracket \{ \%uid \mapsto t \} \vdash \%uid : t$

How can we translate $\Gamma \vdash e_1 + e_2 : \text{int}$ (i.e., ADD)?

- Let $(\llbracket \text{ctxt} \rrbracket \vdash s_1 \Rightarrow \Gamma_1, \Gamma_1 \vdash \text{opn}_1 : \text{i64}) = \llbracket e_1 \rrbracket(\text{ctxt})$
- Let $(\Gamma_1 \vdash s_2 \Rightarrow \Gamma_2, \Gamma_2 \vdash \text{opn}_2 : \text{i64}) = \llbracket e_2 \rrbracket(\text{ctxt})$
- Codestream: $\llbracket \text{ctxt} \rrbracket \vdash s_1, s_2, \%uid = \text{add i64 } \text{opn}_1, \text{opn}_2 \Rightarrow \Gamma_2 \{ \%uid \mapsto \text{i64} \}$
- Operand: $\Gamma_2 \{ \%uid \mapsto \text{i64} \} \vdash \%uid : \text{i64}$

Summary

- *Semantic analysis* phase takes AST as input, constructs symbol table and performs well-formedness checks

Summary

- *Semantic analysis* phase takes AST as input, constructs symbol table and performs well-formedness checks
- Well-formedness derivations can impact compilation. E.g.,
 - `x.field` gets compiled differently depending on the type of `x`
 - We may have to emit bitcasts for uses of subsumption

Summary

- *Semantic analysis* phase takes AST as input, constructs symbol table and performs well-formedness checks
- Well-formedness derivations can impact compilation. E.g.,
 - `x.field` gets compiled differently depending on the type of `x`
 - We may have to emit bitcasts for uses of subsumption
- Compiler translates derivations of well-formedness judgements in the source language to derivations of well-formedness judgements in the target language
 - In an implementation, this viewpoint *implicit*
 - Don't need to do all the bookkeeping involved in manipulating derivations
 - *But* it is helpful for understanding how to organize the translation
 - E.g., `cmp_exp` returns a triple `L1.ty * L1.operand * stream`
In a sense: infers derivations in the source language “on the way down”
builds derivations in the target language “on the way up”
Only remembers the type of the operand (used in some compilation rules).