COS320: Compiling Techniques

Zak Kincaid

April 6, 2019



Intrinsic view: a term that cannot be typed is not a term at all

Compiler cannot translate terms that cannot be typed

· Well-typed programs in the source language translate to well-typed programs in the target

- If target language is typed, this imposes an additional burden:

language

- Intrinsic view: a term that cannot be typed is not a term at all
- Compiler cannot translate terms that cannot be typed
- If target language is typed, this imposes an additional burden: Well-typed programs in the source language translate to well-typed programs in the target language
- 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 (HW4) – well-formed types

Judgements take the form:

- $\vdash t$: "t is a well-formed type"
- \vdash_r ref: "ref is a well-formed reference type"
- $\vdash_{rt} rt$: "rt is a well-formed return type"

TINT		TBOOL			$\vdash_r \textit{ref}$	
$\overline{\vdash}$ int		⊢ bool			⊢ ref	
RSTRING	$\begin{array}{c} \textbf{RARRAY} \\ \vdash t \end{array}$		$\begin{matrix} RFUN \\ \vdash t_1 \end{matrix}$		$\vdash t_n$	$dash_{rt}$ rt
$\overline{\vdash_r}$ string	$\vdash_r t exttt{ iny extstyle ext$			$\vdash_r (t_1, .$	$, t_n) \rightarrow t$	rt
RTVoid				$\begin{array}{c} RTTYP \\ \vdash t \end{array}$		
	$\overline{\vdash_{rt}void}$			$\overline{\vdash_{rt} t}$		

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	LLINT	$LLPTR$ $T \vdash_{\tau} ref$	LLTUPLE $T \vdash t_1 \qquad$	$T \vdash t_n$	$\frac{T \vdash t}{$	LLSIMPLE $\vdash_s t$
$\overline{T dash_s ext{i1}}$	$\overline{T dash_s}$ i64	$\frac{T \vdash_r ref}{T \vdash_s ref*}$	$T \vdash \{t_1,,$	$\{t_1,, t_n\}$ $T \vdash$	$T \vdash [n \times t]$	$\overline{\vdash t}$

LLRTVoid	LLRTSIMPLE	LLRCHAR	LLRType	LLRFun			
	$T \vdash_s t$		$T \vdash t$	$T \vdash_{rt} \textit{rt}$	$T \vdash_s t_1$		$T \vdash_s t_n$
$\overline{T dash_{rt} void}$	$T \vdash_{rt} t$	$\overline{T \vdash_r \mathtt{i8}}$	$T \vdash_r t$		$T \vdash_r rt(t_1,$	$, t_n)$	

LLNamed

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

Translating well-formed types

- Each well-formed Oat type is translated to a well-formed LLVM type
 - types \rightarrow simple types
 - reference types → reference types
 - return types → return types
- Use [I⋅] to denote translation
 - I.e., $\llbracket \vdash \text{int} \rrbracket = \vdash_s \text{i64}$ denotes that the Oat type int is translated to the (simple) LLVMlite type i64

Translating well-formed types

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

TINT	TBOOL	TREF	
		$dash_r$ re	
<u>⊢ int</u>	⊢ bool	⊢ re	

Each has a corresponding case:

- Case TINT: $\llbracket \vdash \mathsf{int} \rrbracket = \vdash_s \mathsf{i64}$ (well-formed by LLInt)
- Case TBOOL: $\llbracket \vdash boo1 \rrbracket = \vdash_s i1$ (well-formed by LLBoo1)

Translating well-formed types

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

TINT	TBOOL	TREF $\vdash_r \mathit{re}$	
$\overline{\vdash}$ int	⊢ bool	$\frac{r}{\vdash}$ ref	

Each has a corresponding case:

- Case TINT: $\llbracket \vdash \text{ int} \rrbracket = \vdash_s \text{ i64 (well-formed by LLInt)}$
- Case TBOOL: $\llbracket \vdash boo1 \rrbracket = \vdash_s i1$ (well-formed by LLBoo1)
- Case TREF: By induction on the derivation, $\llbracket\vdash_r ref\rrbracket$ is a valid judgement of an LLVM reference type, say $\vdash_r t$

$$\frac{\mathsf{TREF}}{\vdash_r \mathit{ref}} \underset{\vdash}{\mathsf{ref}} \leadsto \frac{\mathsf{LLPTR}}{\vdash_r t \ (= \llbracket \vdash_r \mathit{ref} \rrbracket)} \\ \vdash t \ast$$

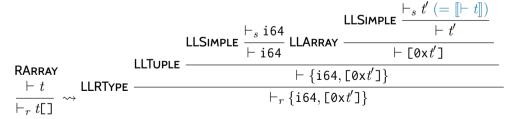
• I.e.,
$$\llbracket \vdash \mathit{ref} \rrbracket = \vdash_s t *$$
, where $\vdash_r t = \llbracket \vdash_r \mathit{ref} \rrbracket$

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.



Summary of type translation

```
• [ | int | = | 64
       • [ | bool ] = | . i1
       • \llbracket \vdash ref \rrbracket = \vdash_s t*, where \vdash_r t = \llbracket \vdash_r ref \rrbracket
       • \llbracket \vdash_r \mathsf{string} \rrbracket = \vdash_r \mathsf{i8}
       • \llbracket \vdash_r t \llbracket \rrbracket \rrbracket = \vdash_r \{ \mathsf{i64}, \llbracket \mathsf{0x}t' \rrbracket \}, where \vdash_s t' = \llbracket \vdash t \rrbracket
      • \llbracket \vdash_r (t_1, ..., t_n) \rightarrow rt \rrbracket = \vdash_{rt} rt'(t'_1, ..., t'_n), where
                   • \vdash_{rt} rt = \llbracket \vdash_{rt} rt \rrbracket.
                   \cdot \vdash_{e} t'_{1} = \llbracket \vdash t_{1} \rrbracket, \dots \vdash_{e} t'_{n} = \llbracket \vdash t_{n} \rrbracket
       • \llbracket \vdash_{rt} \mathsf{void} \rrbracket = \vdash_{rt} \mathsf{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 HW4)
```

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' "
- $\Gamma \vdash opn : t$ "under type environment Γ , operand opn has type t"

ID NUM
$$\frac{\Gamma \vdash id : \Gamma(id)}{\Gamma \vdash n : \text{i64}} \ n \in \mathbb{Z}$$

ADD

$$\frac{\Gamma \vdash \textit{opn}_1 : \texttt{i64} \qquad \Gamma \vdash \textit{opn}_2 : \texttt{i64}}{T, \Gamma \vdash \%uid = \texttt{add i64 } \textit{opn}_1, \textit{opn}_2 \Rightarrow \Gamma\{\%uid \mapsto \texttt{i64}\}} \ \%uid \not\in \textit{dom}(\Gamma)$$

$$\frac{SEQ}{T, \Gamma \vdash s_1 \Rightarrow \Gamma' \qquad T, \Gamma' \vdash s_2 \Rightarrow \Gamma''}{T, \Gamma \vdash s_1, s_2 \Rightarrow \Gamma''}$$

$$\overline{T,\Gamma \vdash \epsilon \Rightarrow \Gamma}$$

...lots more

Well-typed expressions

$$\frac{\mathsf{ADD}}{\Gamma \vdash x : \Gamma(x)} \qquad \frac{\Gamma \vdash e_1 : \mathsf{int} \qquad \Gamma \vdash e_2 : \mathsf{int}}{\Gamma \vdash e_1 + e_2 : \mathsf{int}}$$

...

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

- A codestream judgement $\Gamma_{ll} \vdash s \Rightarrow \Gamma'_{ll}$, and
- An operand judgement $\Gamma'_{ll} \vdash opn : t_{ll}$

- 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 $[\Gamma \vdash x : t](\mathsf{ctxt})$, first let $(id, t*) = \mathsf{ctxt}(x)$, then:

 - Define [ctxt] to be the type environment associated with ctxt
 - Codestream: $[ctxt] \vdash \%uid = load \ t* \ opn \Rightarrow [ctxt] \{\%uid \mapsto t\}$ • Operand: $[ctxt]{\{\%uid \mapsto t\}} \vdash \%uid : t$

- Need a symbol table ctxt, which maps Oat identifiers to LLVMlite operand judgements
 - The operand associated with a variable x is a ${\it pointer}$ to the memory location associated with x
 - To compute $[\Gamma \vdash x : t](\mathsf{ctxt})$, first let $(id, t*) = \mathsf{ctxt}(x)$, then:
 - Define $[t_1, t_2, t_3]$ to be the time equivergence of each with a term
 - Define [ctxt] to be the type environment associated with ctxt
 - Codestream: $[\![\mathtt{ctxt}]\!] \vdash \%uid = \mathtt{load} \ t \!\!\!* \ opn \Rightarrow [\![\mathtt{ctxt}]\!] \{\%uid \mapsto t\}$ • Operand: $[\![\mathtt{ctxt}]\!] \{\%uid \mapsto t\} \vdash \%uid : t$
- How can translate $\Gamma \vdash e_1 + e_2 : \text{int (i.e., ADD)}$?

- 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 $[\Gamma \vdash x : t](\mathsf{ctxt})$, first let $(id, t*) = \mathsf{ctxt}(x)$, then:
 - Define [ctxt] to be the type environment associated with ctxt
 - Codestream: $[ctxt] \vdash \%uid = load \ t* \ opn \Rightarrow [ctxt] \{\%uid \mapsto t\}$ • Operand: $[ctxt]{\{\%uid \mapsto t\}} \vdash \%uid : t$
- How can translate $\Gamma \vdash e_1 + e_2 : \text{int (i.e., ADD)}$?
 - Let $(\llbracket \mathsf{ctxt} \rrbracket \vdash s_1 \Rightarrow \Gamma_1, \Gamma_1 \vdash opn_1 : \mathsf{i64}) = \llbracket e_1 \rrbracket (\mathsf{ctxt})$
 - Let ($\llbracket \mathsf{ctxt} \rrbracket \vdash s_2 \Rightarrow \Gamma_2, \Gamma_2 \vdash opn_2 : \mathsf{i64}$) = $\llbracket e_2 \rrbracket (\mathsf{ctxt})$
 - Codestream: $\Gamma_1 + \Gamma_2 \vdash s_1, s_2, \%uid = \text{add i64 } opn_1, opn_2) \Rightarrow (\Gamma_1 + \Gamma_2) \{\%uid \mapsto \text{i64}\}$
 - Operand: $(\Gamma_1 + \Gamma_2)$ { $\%uid \mapsto i64$ } $\vdash \%uid : i64$

Global initializers

One would expect the following coherence property:

If $\Gamma \vdash e: t$ translates to the codestream judgement $\Gamma_{ll} \vdash s \Rightarrow \Gamma'_{ll}$ and the operand judgement $\Gamma'_{ll} \vdash opn: t_{ll}$, then $[\![t]\!] = t_{ll}$

I.e., "compilation preserved types"

Global initializers

One would expect the following coherence property:

```
If \Gamma \vdash e:t translates to the codestream judgement \Gamma_{ll} \vdash s \Rightarrow \Gamma'_{ll} and the operand judgement \Gamma'_{ll} \vdash opn:t_{ll}, then [\![t]\!]=t_{ll}
```

- I.e., "compilation preserved types"
- There is some subtlety in making this work!
 - Global declaration of string / array constants must compile to types with known length
 - E.g., var x = int[]{0,1} translates to @x = global { 2, {0, 1} }}
 - Oat: x has type int[]
 - LLVM: $@x \text{ has type } \{ i64, [2xi64] \} (\neq [int[]] = \{ i64, [0xi64] \})$
 - cmp_exp_as is a variant of cmp_exp that ensures type preservation via bitcast.

Oat v2 (HW5)

- Specified by a (fairly large) type system
 - Invest some time in making sure you understand how to read the judgements and inference rules
- 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?)

Subtyping and type inference

$$\frac{\Gamma \vdash e : s \qquad \vdash s <: t}{\Gamma \vdash e : t}$$

- Challenge:
 - In the presence of the subsumption rule, a term may have more than one type (how can we
 infer types for a declaration like var x = exp?)
 - Subsumption destroys syntax-directed quality of the type system
- Solution:
 - Do not use subsumption! Integrate subtyping into other inference rules. E.g.,

$$\frac{\textit{Typ_CARR}}{\textit{H} \vdash t \qquad \textit{H}; \textit{G}; \textit{L} \vdash \textit{e}_1 : t \qquad ... \qquad \textit{H}; \textit{G}; \textit{L} \vdash \textit{e}_n : t}{\textit{H}; \textit{G}; \textit{L} \vdash \textit{new t} [] \{\textit{e}_1, ... \textit{e}_n\}}$$

Subtyping and type inference

$\frac{\Gamma \vdash e : s \qquad \vdash s <: t}{\Gamma \vdash e : t}$

- In the presence of the subsumption rule, a term may have more than one type (how can we infer types for a declaration like var x = exp?)
- Subsumption destroys syntax-directed quality of the type system
- Solution:
 - Do not use subsumption! Integrate subtyping into other inference rules. E.g.,

TYP_CARR

$$H; G; L \vdash \mathsf{new} \mathsf{t[]}\{e_1, ...e_n\}$$