## **DTTNU** | Norwegian University of Science and Technology

### **Compiler Construction**

Lecture 10: Context-sensitive analysis

**Michael Engel** 

### Overview

- Where are we standing now?
- There's more to languages than context-free grammars can describe...
  - From syntax to semantics
- Syntax-directed translation
  - Ad-hoc approach
  - Examples
  - A tiny (very imperfect) arithmetical expression to ARM assembly compiler



### Where are we standing now?

#### Source code



 Decides if input *token sequence* can be derived from the grammar

 $ext{expression} 
ightarrow ext{term} \{ (+|-) ext{term} \}$   $ext{term} 
ightarrow ext{factor} \{ (*|/) ext{factor} \}$   $ext{factor} 
ightarrow ext{'(' expression ')'}$  $| ext{ id } | ext{ number}$ 

Norwegian University of

Science and Technology

Compiler Construction 10: Context-sensitive analysis

id(v)

op(+)

id(x)

**Semantic** 

analysis

machine-level program

3

number(42



#### Semantic analysis

- *Name analysis* (check def. & scope of symbols)
- *Type analysis* (check correct type of expressions)
- Creation of *symbol tables* (map identifiers to their types and positions in the source code)

machine-level program

### **Beyond syntax: Example**

- Consider this C program
  - Which errors can you detect?
  - Which of these can be detected using a context-free grammar?



**D**NTNU

Norwegian University of Science and Technology

Compiler Construction 10: Context-sensitive analysis

**Semantic** 

### **Beyond syntax**

- All of these errors are "deeper than syntax"
  - There is a level of correctness that is *deeper than grammar*
  - To generate code, we need to understand its meaning!
- To generate code, the compiler needs to answer many questions, such as:
  - Is "x" a scalar, an array, or a function? Is "x" declared?
  - Are there names that are not declared? Declared but not used?
  - Which declaration of "x" does a given use reference?
  - Is the expression "x \* y + z" type-consistent?
  - In "a[i,j,k]", does a have three dimensions?
  - Where can "z" be stored? (register, local, global, heap, static)
  - In "f = 15", how should 15 be represented?
  - How many arguments does "bar()" take? What about "printf()"?
  - Does "\*p" reference the result of a "malloc()"?
  - Do "p" and "q" refer to the same memory location?
  - Is "x" defined before it is used?

All these are beyond the expressive power of a context-free grammar!



### **Context-sensitive analysis**

#### These questions are part of context-sensitive analysis

- Answers depend on values, not parts of the language
- Questions & answers involve non-local information
- Answers may involve computation

#### How can we answer these questions?

- Use formal methods
  - Context-sensitive grammars?
  - Attribute grammars? (attributed grammars?)
- Use ad-hoc techniques
  - Symbol tables
  - Ad-hoc code (action routines)

Norwegian University of

Science and Technology

In context-sensitive analysis, ad-hoc techniques are often used in practice

### Semantic analysis

For parsing and scanning, formal approaches won



### **Non-syntactical information**

#### Idea: Track the definitions of symbols in a global structure



Semantic

### Symbol tables

Semantic analysis

#### Which information is required to compile an instruction?

023 int x; ... 099 x = x + 1;

#### Line 99 might be translated to:

- 1. Read value from **memory location** of x
- 2. Add integer value 1 to this
- 3. Store value to **memory location** of x

It is convenient to store all this information in a table and link the nodes of the AST to this information



### Implementing symbol tables

# This linking requires finding the table entry of x every time that name is used

- We only get the name ( $\rightarrow$  scanner), so this is a text search problem
- We potentially have thousands of names when compiling a program

#### **Possible approaches:**

- Direct indexing: keep table where the index is a function of the text
   → limits number of identifiers to size of symbol table
- Linked list: keep a dynamic list, go through it and compare  $\rightarrow$  expensive searches for identifiers in the back of the list

• Hash table



### Symbol tables as hash tables

- An unpredictable, fixed-length code (hash value) can be computed from any length of identifier
- Elements stored in fixed-length array of linked lists

Science and Technology

• Search and compare only in the list where hash value matches



Semantic

### Advantage of hash tables

#### Hash tables are a good compromise

- Can dynamically grow with number of stored elements
- Constant time to find the right list to search
- If the hashing function distributes elements evenly, search time is divided by the number of lists
- Balance between static size limitation and list length can be adjusted depending on the data stored

#### However...

No implementation of hash tables directly available in C



# Ad-hoc syntax-directed translation Semantic analysis

#### Build on bottom-up, shift-reduce parser

- Associate a snippet of code with each production
- At each reduction, the corresponding snippet runs
- Allowing arbitrary code provides complete flexibility
  - Includes ability to do tasteless and bad things

#### To make this work

- Need names for attributes of each symbol on LHS & RHS
- Typically, one attribute passed through parser + arbitrary code (structures, globals, statics, ...)
- Yacc introduced \$\$, \$1, \$2, ... \$n, left to right
- Need an evaluation scheme
- Fits nicely into LR(1) parsing algorithm



Similar ideas work for

top-down parsers

### Example: expression grammar

```
1 Block → Block Assign
        Assign
2
3 Assign→ ident = Expr
₄ Expr → Expr + Term
     Expr – Term
5
       Term
6
7 Term → Term × Factor
         Term ÷ Factor
8
     | Factor
9
10 Factor → "(" Expr ")"
         number
11
         ident
12
```

```
{ cost = cost + COST(store); }
{ cost = cost + COST(add); }
{ cost = cost + COST(sub); }
{ cost = cost + COST(mult); }
{ cost = cost + COST(div); }
{ cost = cost + COST(loadImm); }
{ i = hash(ident);
 if (table[i].loaded == false) {
    cost = cost + COST(load);
    table[i].loaded = true; }}
```

Semantic

analysis

Introduce the cost of

expressions to grammar

### One thing was missing...





#### Before parser can reach *Block*, it must reduce *Init*

- Reduction by Init sets cost to zero
- We split the production to create a reduction in the middle
  - for the sole purpose of hanging an action there
  - This trick has many uses

### That wasn't chicken yacc..

Start : Block Block : Block Assign Assign Expr '-' Term Term : Term '\*' Factor Term Term '/' Factor Factor Factor: '(' Expr ')' number ident

```
{ printf("Cost: %d\n", $$); }
                              \{ \$\$ = \$1 + \$2; \}
                              \{ \$\$ = \$1; \}
Assign: ident '=' Expr { $$ = cost(STORE) + $3; }
Expr : Expr '+' Term { $$ = $1 + cost(ADD) + $3; }
                              { $$ = $1 + cost(SUB) + $3; }
                              \{ \$\$ = \$1; \}
                              \{ \$\$ = \$1 + cost(MULT) + \$3; \}
                              \{ \$\$ = \$1 + cost(DIV) + \$3; \}
                              \{ \$\$ = \$1; \}
                              \{ \$\$ = \$2; \}
                              { $$ = cost(LOADIMM); }
                              { int i = hash(ident);
                                if (table[i].loaded == 0) {
                                  $$ = $$ + cost(LOAD);
                                  table[i].loaded = 1;
                                else $$ = 0;
```

Complete yacc+lex code is online

**Semantic** 

### Use case example: timing, energy

Semantic analysis

- How long does a piece of code take to execute?
- How much energy will the code consume?

#### 3.5 Divide and Multiply Instructions

 $\Box$  NTNU

Much more complex to assess for modern high-end CPUs (due to superscalarity, pipelines, caches, ...)



Operand 2

### Example: building an AST



#### So far, our syntax tree was only implicit – we need to operate on it

- Assume constructors for each node
- Assume stack holds pointers to nodes
- Assume yacc-like syntax

1 Start :	Expr	{ \$\$ = \$1; }
2 Expr :	Expr <b>'+'</b> Term	{    \$\$ = MakeAddNode(\$1, \$3);    }
3	Expr <b>'-'</b> Term	{    \$\$ = MakeSubNode(\$1, \$3);    }
4	Term	$\{ $$ = $1; \}$
5 Term :	Term <b>'*'</b> Factor	{    \$\$ = MakeMultNode(\$1, \$3);    }
6	Term '/' Factor	{    \$\$ = MakeDivNode(\$1, \$3);    }
7	Factor	$\{ $$ = $1; \}$
<pre>8 Factor:</pre>	' <b>('</b> Expr <b>')'</b>	{ \$\$ = \$2; }
9	number	<pre>{ \$\$ = MakeNumberNode(token); }</pre>
10	ident	<pre>{ \$\$ = MakeIdentNode(token); }</pre>



### Example: emitting ARM assembly

Early simple compilers derived machine code directly from AST

- We won't do it this way later need more optimization opportunities
- Still a nice example (if the CPU instructions fit this scheme)
- Assume that NxReg() returns a CPU register number

Start	: Expr	{ \$\$ = \$1; }
Ехрт	: Expr <b>'+'</b> Te	<pre>rm { \$\$ = NxReg(); Emit("add", \$\$, \$1, \$3); }</pre>
	Expr <b>'-'</b> Te	<pre>rm { \$\$ = NxReg(); Emit("sub", \$\$, \$1, \$3); }</pre>
	Term	$\{ \$\$ = \$1; \}$
Term	: Term '*' Fa	<pre>ctor { \$\$ = NxReg(); Emit("mul", \$\$, \$1, \$3); }</pre>
	Term <b>'/'</b> Fa	<pre>ctor { \$\$ = NxReg(); Emit("div", \$\$, \$1, \$3); }</pre>
	Factor	$\{ $$ = $1; \}$
Factor	: '(' Expr ')	{ \$\$ = \$2; }
	number	{    \$\$ = NxReg();    EmitLI("mov", \$\$, yylval);    }
	ident	<pre>{ \$\$ = NxReg(); EmitLD("ldr", \$\$, yytext); }</pre>

Semantic

analysis

We omit

### Example: emitting ARM assembly

Emit, EmitLI and EmitLD print assembler instructions

NxReg should return *free* (unused) register number We will run out of registers for complex

```
expressions!
int NxReq(void) {
  static int reg = 0;
  if (reg > 11) { reg = 0; return reg; } // wraparound if > 12 registers used!
  return req++;
void EmitLD(char *op, int rd, char *adr) { // emit memory load from address "adr"
  printf("\tldr r%d, =%s\n", rd, adr);
  printf("\t%s r%d, [r%d]\n", op, rd, rd);
void EmitLI(char *op, int rd, int val) { // emit load of constant value "val"
 printf("\t%s r%d, #%d\n", op, rd, val);
void Emit(char *op, int rd, int rs1, int rs2) { // emit given arithmetic instrn.
  printf("\t%s r%d, r%d, r%d\n", op, rd, rs1, rs2);
```

**Semantic** 

### **Example: compiler output**

Input: (z-3)\*x +5 Input: (z-3)\*x+5 \$ echo "(z-3)\*x+5" | ./compile \$ echo "(z-3)\*x)+5" | ./compile 1dr r0, =z1dr r0, =zldr r0, [r0] // r0 = z ldr r0, [r0] mov r1, #3 // r1 = 3 mov r1, #3 sub r2, r0, r1 // r2 = z-3 sub r2, r0, r1 Directly generating code ldr r3, =x ldr r3, =x during parsing  $\rightarrow$ ldr r3, [r3] // r3 = x ldr r3, [r3] partial assembler code mul r4, r2, r3 // r4 = (z-3)\*xmul r4, r2, r3 is being emitted! mov r5, #5 // r5 = 5 syntax error: ) add r6, r4, r5 // r6 = (z-3)\*x+5

#### **ARM** instruction overview:

ldr	rd,	= Z				load	address of memory location z into reg. rd
ldr	rd,	[rs	]			load	contents of memory at addr. rs into rd
mov	rd,	#va	1			сору	numerical value val into register rd
(add	d sub	o mu	l div)	rd, rs1,	<b>rs2</b> –	exec	cute rd = rs1 (+ - * /) rs2



### Example: register wraparound

Input: (a+(b+(c+(d+e))))\*x

\$ echo "(a+(b+(c+(d+e)	)))*x"   ./compile
ldr r0, =a	
ldr r0, [r0]	// r0 = a
ldr r1, =b	
ldr r1, [r1]	// r1 = b
ldr r2, =c	
ldr r2, [r2]	// r2 = c
ldr r3, =d	
ldr r3, [r3]	// r3 = d
ldr r4, =e	
ldr r4, [r4]	// r4 = e
add r5, r3, r4	// r5 = d+e
add r0, r2, r5	// r0 = (d+e)+c
add r0, r1, r0	A CONTRACTOR OF THE OWNER OWNER OWNER OF THE OWNER OWNE
add r1, r0, r0	
ldr r2, =x	
ldr r2, [r2]	
mul r3, r1, r2	

Number of registers in NxReg() reduced to 5 here to make example shorter!

A real compiler needs a method for *register allocation* 

- assign values to *free* registers
- when running out of registers, *spill* (save to memory) register contents and *restore* them when needed later
- efficient register allocation is complex – as we will see later

No more unused registers: wraparound! r0 is overwritten here Value of "a" is lost → incorrect result!



### What's next?

Semantic analysis

- A quick look at attribute grammars
- Some insight into type systems and type analysis

#### References

[1] ARM Cortex-A57 Software Optimization Guide http://infocenter.arm.com/help/topic/com.arm.doc.uan0015b/ Cortex\_A57\_Software\_Optimization\_Guide\_external.pdf

[2] Kerstin Eder and John P. Gallagher, Energy-Aware Software Engineering, DOI: 10.5772/65985

https://www.intechopen.com/books/ict-energy-concepts-for-energy-efficiency-and-sustainability/energyaware-software-engineering

[3] Peter Marwedel, slide set on Embedded System Evaluation and Validation: WCET analysis (sl. 14 ff.) <u>https://ls12-www.cs.tu-dortmund.de/daes/media/documents/staff/marwedel/es-book/slides11/es-marw-5.1-evaluation.pdf</u>

[4] ARM Instruction Set reference guide https://static.docs.arm.com/100076/0100/ arm\_instruction\_set\_reference\_guide\_100076\_0100\_00\_en.pdf

