

Assignment 4: Compiler for \mathcal{L}_{if}

(Deadline: 20.12.2021 09:00)

Exercise 1: Shrink Pass

Implement the pass `shrink` in `compiler.py` to remove `and` and `or` from the language by translating them to `if` expressions in \mathcal{L}_{if} (see end of this assignment sheet for concrete and abstract syntax):

$$\begin{aligned} e_1 \text{ and } e_2 &\Rightarrow e_2 \text{ if } e_1 \text{ else False} \\ e_1 \text{ or } e_2 &\Rightarrow \text{True if } e_1 \text{ else } e_2 \end{aligned}$$

Exercise 2: Remove Complex Operands Pass

Add cases for Boolean constants, comparisons, `if` expressions and statements to `rco_exp` and `rco_stmt` functions in `compiler.py`. Regarding `if`, it is particularly important to **not** replace its condition with a temporary variable because that would interfere with the generation of high-quality output in the `explicate_control` pass.

We add a new language form, the `Let` expression (its abstract syntax can be found in `utils.py`), to aid in the translation of `if` expressions. When we recursively process the two branches of the `if`, we generate temporary variables and their initializing expressions (see example below). However, these expressions may contain side effects and should only be executed when the condition of the `if` is true (for the “then” branch) or false (for the “else” branch). The `Let` provides a way to initialize the temporary variables within the two branches of the `if` expression. In general, the `Let(x, e_1, e_2)` form assigns the result of e_1 to the variable x , and then evaluates e_2 , which may reference x .

Example (of what we want to avoid):

$\begin{aligned} &(\text{input_int()} + 1) \setminus \\ &\text{if } x > 0 \text{ else } (42 + \text{input_int()}) \end{aligned} \Rightarrow$	$\begin{aligned} &\text{tmp0} = \text{input_int()} \\ &\text{tmp1} = \text{input_int()} \\ &(\text{tmp0} + 1) \setminus \\ &\text{if } x > 0 \text{ else } (42 + \text{tmp1}) \end{aligned}$
--	--

Exercise 3: Explicate Control Pass

Implement the `explicate_control` pass (translating $\mathcal{L}_{\text{if}}^{\text{mon}}$ to \mathcal{C}_{if} , grammar found at the end of this assignment) using the following four auxiliary functions (as found in your assignment stub):

`explicate_effect` generates code for expressions as statements, so their result is ignored and only their side effects matter. It has three parameters: 1) the expression to be compiled, 2) the already-compiled code for this expression’s *continuation*, that is, the list of statements that should execute after this expression, and 3) the dictionary of generated basic blocks. The `explicate_effect` function returns a list of \mathcal{C}_{if} statements and it may add to the dictionary of basic blocks. If the expression to be compiled is an `if` expression, we translate the two branches using `explicate_effect` and then translate the condition expression using `explicate_pred`, which generates code for the entire `if`.

explicate_assign generates code for expressions on the right-hand side of an assignment. It has four parameters: 1) the right-hand-side of the assignment, 2) the left-hand-side of the assignment (the variable), 3) the continuation, and 4) the dictionary of basic blocks. The **explicate_assign** function returns a list of \mathcal{C}_{if} statements and it may add to the dictionary of basic blocks.

explicate_pred generates code for an **if** expression or statement by analyzing the condition expression. It has four parameters: 1) the condition expression, 2) the generated statements for the “then” branch, 3) the generated statements for the “else” branch, and 4) the dictionary of basic blocks. The **explicate_pred** function returns a list of \mathcal{C}_{if} statements and it may add to the dictionary of basic blocks.

explicate_stmt generates code for statements. It has three parameters: 1) the statement to be compiled, 2) the code for its continuation, and 3) the dictionary of basic blocks. The **explicate_stmt** returns a list of statements and it may add to the dictionary of basic blocks.

Exercise 4: Select Instructions

Adapt the **select_instructions** pass (now translating \mathcal{C}_{if} to $\text{x86}_{\text{if}}^{\text{Var}}$, see end of this assignment for grammar; note the new arg terminal **ByteReg**). Some useful examples:

True	\Rightarrow	1	False	\Rightarrow	0
$var = \text{not } var$	\Rightarrow	xorq \$1, var			
$var = \text{not } atm$	\Rightarrow	movq arg, var			
		xorq \$1, var			
$var = (atm_1 == atm_2)$	\Rightarrow	cmpq arg₂, arg₁			
		sete %al			
		movzbq %al, var			
goto ℓ	\Rightarrow	jmp ℓ			
if $atm_1 == atm_2$:		cmpq arg₂, arg₁			
goto ℓ_1	\Rightarrow	je ℓ_1			
else:		jmp ℓ_2			
goto ℓ_2					

Regarding the **return** statement, treat it as an assignment to the **rax** register followed by a jump a new block labelled **conclusion**, which you will fill later in the Prelude and Conclusion pass.

Exercise 5: Register Allocation

Register allocation for $\text{x86}_{\text{if}}^{\text{Var}}$ requires you to perform liveness analysis on blocks separately. To perform liveness analysis on a block, we must know the live-after set for the last instruction in the block. If there are no successors (no jumps to other blocks), the live-after set is empty. If there are successors, their live-before set must be calculated first.

5.1: Control Flow Graph

Implement the function `cfg` in `register_allocation.py` to generate a control flow graph for a given dictionary of blocks. The control flow graph should be a directed graph where there is an edge from v to v' if there is a jump in block v to block v' .

In your stub, the result of this function is then transposed and topologically sorted to generate an order for performing liveness analysis for blocks.

5.2: Arg/Read/Write Locations

Extend your auxiliary functions for the new type of instructions found in `x86Varif`. You can omit `Jump` and `JumpIf`, they will be handled in 5.3.

5.3: Liveness Analysis

Update your `uncover_live` function: Iterate through the blocks in the order calculated in 5.1. Save the live-before set of a block in the `live_before_block` dictionary.

Regarding `Jump` and `JumpIf` instructions: The locations that are live before a `Jump` should be the locations in L_{before} at the target of the jump. Liveness analysis for `JumpIf` is particularly interesting because, during compilation, we do not know which way a conditional jump will go. So we do not know whether to use the live-before set for the following instruction or the live-before set for the block associated with the *label*. However, there is no harm to the correctness of the generated code if we classify more locations as live than the ones that are truly live during one particular execution of the instruction. Thus, we can take the union of the live-before sets from the following instruction and from the mapping for *label* in `live_before_block`.

5.4: Interference Graph

Adapt `build_interference` to the new language.

Exercise 6: Patch Instructions & Prelude and Conclusion

6.1: Patch Instructions

The new instructions `cmpq` and `movzbq` have some special restrictions that need to be handled in the `patch_instructions` pass. The second argument of the `cmpq` instruction must not be an immediate value (such as an integer). So if you are comparing two immediates, we recommend inserting a `movq` instruction to put the second argument in `rax`. As usual, `cmpq` may have at most one memory reference. The second argument of the `movzbq` must be a register.

6.2: Prelude and Conclusion

The generation of the `main` function with its prelude and conclusion must change to accommodate how the program now consists of one or more basic blocks. After the prelude in `main`, jump to the `start` block. Place the conclusion in a basic block labelled with `conclusion`.

Concrete and Abstract Syntax of \mathcal{L}_{lf} , \mathcal{C}_{lf} , x86_{lf}

$exp ::= int \mid \text{input_int}() \mid -exp \mid exp + exp \mid exp - exp \mid (exp)$
$stmt ::= \text{print}(exp) \mid exp$
$exp ::= var$
$stmt ::= var = exp$
$cmp ::= == \mid != \mid < \mid <= \mid > \mid >=$
$exp ::= \text{True} \mid \text{False} \mid exp \text{ and } exp \mid exp \text{ or } exp \mid \text{not } exp$
$\quad \mid exp \text{ cmp } exp \mid exp \text{ if } exp \text{ else } exp$
$stmt ::= \text{if } exp: stmt^+ \text{ else: } stmt^+$
$\mathcal{L}_{\text{lf}} ::= stmt^*$

Figure 1: The concrete syntax of \mathcal{L}_{lf} .

$binaryop ::= \text{Add}() \mid \text{Sub}()$
$unaryop ::= \text{USub}()$
$exp ::= \text{Constant}(int) \mid \text{Call}(\text{Name}('input_int'), [])$
$\quad \mid \text{UnaryOp}(unaryop, exp) \mid \text{BinOp}(exp, binaryop, exp)$
$stmt ::= \text{Expr}(\text{Call}(\text{Name}('print'), [exp])) \mid \text{Expr}(exp)$
$exp ::= \text{Name}(var)$
$stmt ::= \text{Assign}([\text{Name}(var)], exp)$
$boolop ::= \text{And}() \mid \text{Or}()$
$unaryop ::= \text{Not}()$
$cmp ::= \text{Eq}() \mid \text{NotEq}() \mid \text{Lt}() \mid \text{LtE}() \mid \text{Gt}() \mid \text{GtE}()$
$bool ::= \text{True} \mid \text{False}$
$exp ::= \text{Constant}(bool) \mid \text{BoolOp}(boolop, [exp, exp])$
$\quad \mid \text{Compare}(exp, [cmp], [exp]) \mid \text{IfExp}(exp, exp, exp)$
$stmt ::= \text{If}(exp, stmt^+, stmt^+)$
$\mathcal{L}_{\text{lf}} ::= \text{Module}(stmt^*)$

Figure 2: The abstract syntax of \mathcal{L}_{lf} .

```

atm ::= int | var | bool
exp ::= atm | input_int() | atm binaryop atm | unaryop atm
      | atm cmp atm
stmt ::= print(exp) | exp
      | var = exp | return exp | goto label
      | if atm cmp atm: goto label else: goto label
Clf ::= (label: stmt*)...

```

Figure 3: The concrete syntax of \mathcal{C}_{lf} .

```

atm ::= Constant(int) | Name(var) | Constant(bool)
exp ::= atm | Call(Name('input_int'), [])
      | BinOp(atm, binaryop, atm) | UnaryOp(unaryop, atm)
      | Compare(atm, [cmp], [atm])
stmt ::= Expr(Call(Name('print'), [exp])) | Expr(exp)
      | Assign([Name(var)], exp) | Return(exp) | Goto(label)
      | If(Compare(atm, [cmp], [atm]), [Goto(label)], [Goto(label)])
Clf ::= CProgram({label: stmt*, ...})

```

Figure 4: The abstract syntax of \mathcal{C}_{lf} .

```

bytereg ::= ah | al | bh | bl | ch | cl | dh | dl
arg ::= $int | %reg | int(%reg) | %bytereg
cc ::= e | ne | l | le | g | ge
instr ::= addq arg, arg | subq arg, arg | negq arg | movq arg, arg |
      callq label | pushq arg | popq arg | retq | jmp label |
      label: instr | xorq arg, arg | cmpq arg, arg |
      setcc arg | movzbq arg, arg | jcc label
x86lf ::= .globl main
      main: instr...

```

Figure 5: The concrete syntax of x86_{lf} .

```

bytereg ::= 'ah' | 'al' | 'bh' | 'bl' | 'ch' | 'cl' | 'dh' | 'dl'
arg      ::= Immediate(int) | Reg(reg) | Deref(reg,int) | ByteReg(bytereg)
cc       ::= 'e' | 'ne' | 'l' | 'le' | 'g' | 'ge'
instr    ::= Instr('addq',[arg,arg]) | Instr('subq',[arg,arg])
              | Instr('movq',[arg,arg]) | Instr('negq',[arg])
              | Callq(label,int) | Retq() | Instr('pushq',[arg])
              | Instr('popq',[arg]) | Jump(label)
              | Instr('xorq',[arg,arg]) | Instr('cmpq',[arg,arg])
              | Instr('set'+cc,[arg]) | Instr('movzbq',[arg,arg])
              | JumpIf(cc,label)
x86lf   ::= X86Program({label : instr*, ... })

```

Figure 6: The abstract syntax of x86_{lf}.