

プログラミング言語 10
言語処理系 / Implementing Programming
Languages

田浦

目次 / Contents

- ① コンパイラの仕事 / What do compilers do basically?
- ② より一般の場合 / More general cases
- ③ 最小限のCからのコード生成器 / Implementing a minimum C compiler
- ④ 中間言語 / Intermediate Representation

言語処理系実装の形態

- **インタプリタ**: プログラムを解釈実行 (プログラムと入力から出力を直接計算)
- **トランスレータ**: プログラムを別の言語 (例: C) に翻訳
 - ▶ 例: OpenMP (Cの並列拡張) を C (+ Pthreads) に翻訳
- **コンパイラ**: プログラムを機械語に翻訳

Various forms of language implementation

- **interpreter**: interprets and executes programs (takes a program and an input; and computes the output)
- **translator**: translates programs into another language (e.g., C)
 - ▶ e.g. translate OpenMP (parallel extension to C) to C (+ Pthreads)
- **compiler**: translates programs into **a machine (assembly) code**

なぜ(今も)言語処理系を学ぶか

- 新ハードウェア用の処理系
 - ▶ GPU用のC/C++言語 (CUDA, OpenACC, OpenMP)
 - ▶ プロセッサの新しい命令セット (e.g., SIMD) への対応
 - ▶ 量子コンピュータ, 量子アニーラ
- 新汎用言語
 - ▶ Scala, Julia, Go, Rust, etc.
- 言語の拡張
 - ▶ 並列処理用拡張 (例: OpenMP, CUDA, OpenACC, Cilk)
 - ▶ ベクトル命令用拡張
 - ▶ 型システム拡張 (例: PyPy, TypeScript)
- 目的に特化した言語
 - ▶ 統計パッケージ (R, MatLab, etc.)
 - ▶ データ処理 (SQL, NoSQL, SPARQL, etc.)
 - ▶ 機械学習
 - ▶ 制約解消系, 定理証明系 (Coq, Isabelle, etc.)
 - ▶ アプリケーション用マクロ言語 (Visual Basic (MS Office 用), Emacs Lisp (Emacs), Javascript (ウェブブラウザ), etc.)

Why do you want to build a language, today?

- new hardware
 - ▶ C/C++ for GPUs (CUDA, OpenACC, OpenMP)
 - ▶ new instruction set (e.g., SIMD) of the processor
 - ▶ quantum computers, quantum annealers
- new general purpose languages
 - ▶ Scala, Julia, Go, Rust, etc.
- new extension
 - ▶ parallel processing (ex: OpenMP, CUDA, OpenACC, Cilk)
 - ▶ vector/SIMD processing
 - ▶ type system extension for safety (ex: PyPy, TypeScript)
- new special purpose (domain specific) languages
 - ▶ statistics (R, MatLab, etc.)
 - ▶ data processing (SQL, NoSQL, SPARQL, etc.)
 - ▶ deep learning
 - ▶ constraint solving, proof assistance (Coq, Isabelle, etc.)
 - ▶ macro (Visual Basic (MS Office), Emacs Lisp (Emacs), Javascript (web browser), etc.)

Contents

- ① コンパイラの仕事 / What do compilers do basically?
- ② より一般の場合 / More general cases
- ③ 最小限の C からのコード先生器 / Implementing a minimum C compiler
- ④ 中間言語 / Intermediate Representation

高水準言語 vs. 機械語

	高水準言語 (e.g., C)	機械語
制御構造	for, while, if, ...	≈ go to だけ
式	任意の入れ子	≈ $C = A \text{ op } B$ だけ
局所変数の数	いくらでも	≈ レジスタ数まで
局所変数の寿命	関数実行中	≈ 関数呼び出しまで

- これらのギャップを埋めるのがコンパイラ
- <https://www.felixcloutier.com/x86/index.html>
- https://wiki.cdot.senecacollege.ca/wiki/X86_64_Register_and_Instruction_Quick_Start

High level language vs. machine code

	high-level (e.g., C)	machine
control expression	for, while, if, ...	\approx jump (“go to”) only
the number of local variables	arbitrary nest	\approx C = A op B only
lifetime of a local variable	arbitrary	\approx only a fixed number of registers
	during the function call that defined it	\approx some registers live only up to the next function call

- compiler's main job is to fill those *gaps*
- <https://www.felixcloutier.com/x86/index.html>
- https://wiki.cdot.senecacollege.ca/wiki/X86_64_Register_and_Instruction_Quick_Start

コード生成 — 人間コンパイラ内観

- 例: 以下 (ちなみに \sqrt{c} を求めるニュートン法) をどう機械語にするか

```
1 double sq(double c, long n) {  
2     double x = c;  
3     for (long i = 0; i < n; i++) {  
4         x = x / 2 + c / (x + x);  
5     }  
6     return x;  
7 }
```

Code generation by hand — introspecting “human compiler”

- ex: how to convert the following (which finds \sqrt{c} by the Newton method) into machine language

```
1 double sq(double c, long n) {  
2     double x = c;  
3     for (long i = 0; i < n; i++) {  
4         x = x / 2 + c / (x + x);  
5     }  
6     return x;  
7 }
```

ステップ 1 — 制御構造を goto だけに

```
1 double sq(double c, long n) {  
2     double x = c;  
3     for (long i = 0; i < n; i++) {  
4         x = x / 2 + c / (x + x);  
5     }  
6     return x;  
7 }
```

```
⇒ 1 double sq(double c, long n) {  
2     double x = c;  
3     long i = 0;  
4     if (i >= n) goto Lend;  
5     Lstart:  
6     x = x / 2 + c / (2 * x);  
7     i++;  
8     if (i < n) goto Lstart;  
9     Lend:  
10    return x;  
11 }
```

Step 1 — make all controls “goto”s

```
1 double sq(double c, long n) {  
2     double x = c;  
3     for (long i = 0; i < n; i++) {  
4         x = x / 2 + c / (x + x);  
5     }  
6     return x;  
7 }
```

```
1 double sq(double c, long n) {  
2     double x = c;  
3     long i = 0;  
4     if (i >= n) goto Lend;  
5     Lstart:  
6     x = x / 2 + c / (2 * x);  
7     i++;  
8     if (i < n) goto Lstart;  
9     Lend:  
10    return x;  
11 }
```

ステップ 2 — 式を $C = A \text{ op } B$ に

```
1 double sq(double c, long n) {  
2     double x = c;  
3     long i = 0;  
4     if (i >= n) goto Lend;  
5     Lstart:  
6     x = x / 2 + c / (2 * x);  
7     i++;  
8     if (i < n) goto Lstart;  
9     Lend:  
10    return x;  
11 }
```

⇒

```
1 double sq3(double c, long n) {  
2     double x = c;  
3     long i = 0;  
4     if (!(i < n)) goto Lend;  
5     Lstart:  
6     double t0 = 2;  
7     double t1 = x / t0;  
8     double t2 = t0 * x;  
9     double t3 = c / t2;  
10    x = t1 + t3;  
11    i = i + 1;  
12    if (i < n) goto Lstart;  
13    Lend:  
14    return x;  
15 }
```

Step 2 — flatten all nested expressions to “C = A op B”

```
1 double sq(double c, long n) {  
2     double x = c;  
3     long i = 0;  
4     if (i >= n) goto Lend;  
5 Lstart:  
6     x = x / 2 + c / (2 * x);  
7     i++;  
8     if (i < n) goto Lstart;  
9 Lend:  
10    return x;  
11 }
```

⇒

```
1 double sq3(double c, long n) {  
2     double x = c;  
3     long i = 0;  
4     if (!(i < n)) goto Lend;  
5 Lstart:  
6     double t0 = 2;  
7     double t1 = x / t0;  
8     double t2 = t0 * x;  
9     double t3 = c / t2;  
10    x = t1 + t3;  
11    i = i + 1;  
12    if (i < n) goto Lstart;  
13 Lend:  
14    return x;  
15 }
```

ステップ3 — 変数に機械語レベルでの変数(レジスタまたはメモリ)を割り当て

- 注: 浮動小数点数の定数は命令中には書けない/ cannot embed floating point constants in instructions

```
1  /* c : xmm0, n : rdi */
2  double sq3(double c, long n) {
3      double x = c;          /* x : xmm1 */
4      long i = 0;           /* i : rsi */
5      if (!(i < n)) goto Lend;
6  Lstart:
7      double t0 = 2;        /* t0 : xmm2 */
8      double t1 = x / t0;   /* t1 : xmm3 */
9      double t2 = t0 * x;   /* t2 : xmm4 */
10     double t3 = c / t2;   /* t3 : xmm5 */
11     x = t1 + t3;
12     i = i + 1;
13     if (i < n) goto Lstart;
14 Lend:
15     return x;
16 }
```


Step 3 —assign “machine variables” (registers or memory) to variables

- note: cannot write floating point constants in instructions

```
1  /* c : xmm0, n : rdi */
2  double sq3(double c, long n) {
3      double x = c;          /* x : xmm1 */
4      long i = 0;           /* i : rsi */
5      if (!(i < n)) goto Lend;
6      Lstart:
7      double t0 = 2;        /* t0 : xmm2 */
8      double t1 = x / t0;   /* t1 : xmm3 */
9      double t2 = t0 * x;   /* t2 : xmm4 */
10     double t3 = c / t2;   /* t3 : xmm5 */
11     x = t1 + t3;
12     i = i + 1;
13     if (i < n) goto Lstart;
14     Lend:
15     return x;
16 }
```

ステップ4 — 命令に変換

```
1 /* c : xmm0, n : rdi */
2 double sq3(double c, long n) {
3     # double x = c;          /*x:xmm1*/
4     movasq %xmm0,%xmm1
5     # long i = 0;          /*i:rsi*/
6     movq $0,%rsi
7     .Lstart:
8     # if (!(i < n)) goto Lend;
9     cmpq %rdi,%rsi # n - i
10    jle .Lend
11    # double t0 = 2;        /*t0:xmm2*/
12    movasq .L2(%rip),%xmm2
13    # double t1 = x / t0; /*t1:xmm3*/
14    movasq %xmm1,%xmm3
15    divq %xmm2,%xmm3
16    # double t2 = t0 * x; /*t2:xmm4*/
17    movasq %xmm0,%xmm4
18    mulsd xmm2,%xmm4
```

```
1 # double t3 = c/t2; /*t3:xmm5*/
2 movasq %xmm0,%xmm5
3 divsd %xmm4,%xmm5
4 # x = t1 + t3;
5 movasq %xmm3,%xmm1
6 addsd %xmm5,%xmm1
7 # i = i + 1;
8 addq $1,%rsi
9 # if (i < n) goto Lstart;
10 cmpq %rdi,%rsi # n - i
11 jl .Lstart
12 .Lend:
13 # return x;
14 movq %xmm1,%xmm0
15 ret
16 }
```

Step 4 — convert them to machine instructions

```
1 /* c : xmm0, n : rdi */
2 double sq3(double c, long n) {
3     # double x = c;          /*x:xmm1*/
4     movasd %xmm0,%xmm1
5     # long i = 0;          /*i:rsi*/
6     movq $0,%rsi
7     .Lstart:
8     # if (!(i < n)) goto Lend;
9     cmpq %rdi,%rsi # n - i
10    jle .Lend
11    # double t0 = 2;        /*t0:xmm2*/
12    movasd .L2(%rip),%xmm2
13    # double t1 = x / t0; /*t1:xmm3*/
14    movasd %xmm1,%xmm3
15    divq %xmm2,%xmm3
16    # double t2 = t0 * x; /*t2:xmm4*/
17    movasd %xmm0,%xmm4
18    mulsd xmm2,%xmm4
```

```
1     # double t3 = c/t2; /*t3:xmm5*/
2     movasd %xmm0,%xmm5
3     divsd %xmm4,%xmm5
4     # x = t1 + t3;
5     movasd %xmm3,%xmm1
6     addsd %xmm5,%xmm1
7     # i = i + 1;
8     addq $1,%rsi
9     # if (i < n) goto Lstart;
10    cmpq %rdi,%rsi # n - i
11    jl .Lstart
12    .Lend:
13    # return x;
14    movq %xmm1,%xmm0
15    ret
16 }
```

Contents

- ① コンパイラの仕事 / What do compilers do basically?
- ② より一般の場合 / More general cases
- ③ 最小限の C からのコード先生器 / Implementing a minimum C compiler
- ④ 中間言語 / Intermediate Representation

コード生成 — 一般的には難しいところ

- 気軽に各途中結果にレジスタを割り当てたが ...

```
1   double x = c;          /* x : xmm1 */
2   Lstart:
3   if (!(i < n)) goto Lend;
4   double t0 = 2;        /* t0 : xmm2 */
5   double t1 = x / t0;   /* t1 : xmm3 */
6   double t2 = t0 * x;   /* t2 : xmm4 */
7   double t3 = c / t2;   /* t3 : xmm5 */
```

- レジスタは足りなくなるかも知れない
- 多くのレジスタは関数呼び出しをまたがると破壊される
- オペランドレジスタが限定されている命令もある (e.g., 整数割り算の被除数は %rax, %rdx ≡ %rax, %rdx は割り算をまたがると破壊される)
- → 一般にはメモリ (スタック領域) も使う必要がある

Things are more complex in general ...

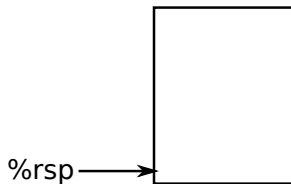
- we've liberally assign registers to intermediate results, but ...

```
1  double x = c;          /* x : xmm1 */
2  Lstart:
3  if (!(i < n)) goto Lend;
4  double t0 = 2;         /* t0 : xmm2 */
5  double t1 = x / t0;    /* t1 : xmm3 */
6  double t2 = t0 * x;    /* t2 : xmm4 */
7  double t3 = c / t2;    /* t3 : xmm5 */
```

- registers are finite (may run out)
- some registers are destroyed (i.e., values on them are lost) across a function call
- some instructions demand operands to be on specific registers (e.g., dividend of integer division must be on `rax` and `rdx` \equiv `rax` and `rdx` are destroyed across an integer division)
- \rightarrow you must use memory (“stack” region) as well

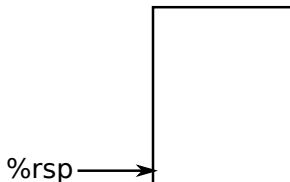
最も単純だが一般的な (コンパイラによる) コード生成の作戦

- 途中結果は一般にはメモリ (スタック領域) も使う必要がある
⇒ 「常に」スタック領域を使うのが単純



A simplest general strategy for code generation by a compiler

- in general, memory (stack) must be used to hold intermediate results \Rightarrow simply, “*always*” use stack
- a register is used only “temporarily” (to read an operand from memory, which is immediately used by an instruction)



レジスタ使用慣例 (ABI)

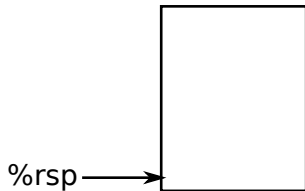
- 整数/ポインタの第 1-6 引数: rdi, rsi, rdx, rcx, r8, r9
- 浮動小数点数の引数は, xmm0, xmm1, ...
- 整数/ポインタの戻り値: rax
- rsp: 関数先頭でスタックの端を指し, そこには戻り番地が格納されている
- callee-save レジスタ: rbx, rbp, r12, r13, r14, r15 (関数呼び出しをまたがって保存 → 呼び出された関数がそれらを使う場合は保存してから使う)
- その他のレジスタは caller-save (関数呼び出しをまたがったら壊れると仮定してコードを生成)
- https://wiki.cdot.senecacollege.ca/wiki/X86_64_Register_and_Instruction_Quick_Start “general-purpose registers” を参照

Register usage conventions (ABI)

- the first six integer/pointers arguments : `rdi`, `rsi`, `rdx`, `rcx`, `r8`, `r9`
- floating point number arguments : `xmm0`, `xmm1`, ...
- an integer/pointer return value : `rax`
- `rsp` : points the end of the stack upon function entry, which holds the return address
- callee-save registers: `rbx`, `rbp`, `r12`, `r13`, `r14`, `r15`
(preserved across function calls → a function must save them before using (setting a value to) them)
- other registers are caller-save (a function must assume they are destroyed across function calls)
- see “general-purpose” registers in https://wiki.cdot.senecacollege.ca/wiki/X86_64_Register_and_Instruction_Quick_Start

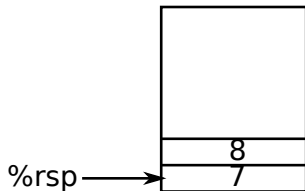
関数呼び出し時の動き

- `long f() { ... g(1,2,3,4,5,6,7,8); ... }`
- `f` 実行中



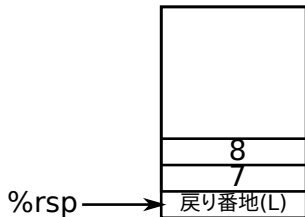
関数呼び出し時の動き

- `long f() { ... g(1,2,3,4,5,6,7,8); ... }`
- `call g` 実行直前 `rdi=1, rsi=2, rdx=3, rcx=4, r8=5, r9=6`



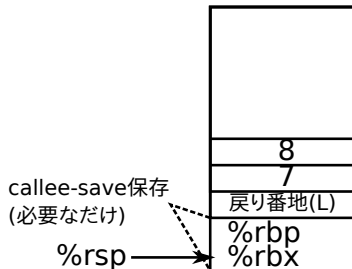
関数呼び出し時の動き

- `long f() { ... g(1,2,3,4,5,6,7,8); ... }`
- `call g` 実行直後 (g 開始)



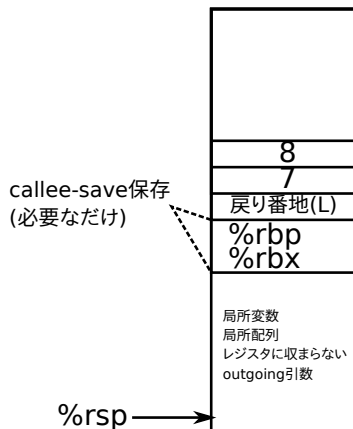
関数呼び出し時の動き

- `long f() { ... g(1,2,3,4,5,6,7,8); ... }`
- `g` が使う callee-save レジスタ保存



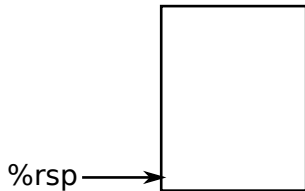
関数呼び出し時の動き

- `long f() { ... g(1,2,3,4,5,6,7,8); ... }`
- `g` 実行中



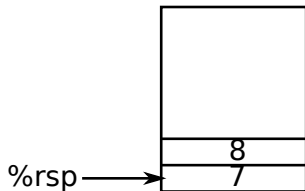
What happens upon function calls

- `long f() { ... g(1,2,3,4,5,6,7,8); ... }`
- during `f`



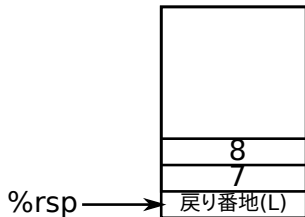
What happens upon function calls

- `long f() { ... g(1,2,3,4,5,6,7,8); ... }`
- right before “`call g`” `rdi=1, rsi=2, rdx=3, rcx=4, r8=5, r9=6`



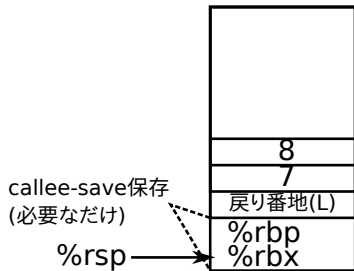
What happens upon function calls

- `long f() { ... g(1,2,3,4,5,6,7,8); ... }`
- right after “call g” (when g started)



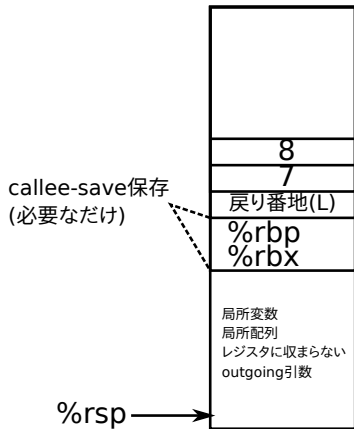
What happens upon function calls

- `long f() { ... g(1,2,3,4,5,6,7,8); ... }`
- save callee-save registers `g` uses



What happens upon function calls

- long f() { ... g(1,2,3,4,5,6,7,8); ... }
- during g



関数呼び出しを含むコード生成例

```
1 double integ(long n) {  
2     double x = 0;  
3     double dx = 1 / (double)n;  
4     double s = 0;  
5     for (long i = 0; i < n; i++) {  
6         s += f(x);  
7         x += dx;  
8     }  
9     return s * dx;  
10 }
```

Code generation including function calls

```
1 double integ(long n) {  
2     double x = 0;  
3     double dx = 1 / (double)n;  
4     double s = 0;  
5     for (long i = 0; i < n; i++) {  
6         s += f(x);  
7         x += dx;  
8     }  
9     return s * dx;  
10 }
```

“goto”化と“C = A op B”化

```
1 double integ3(long n) {      /* n : 0(%rsp) */
2     double x = 0;           /* x : 8(%rsp) */
3     double t0 = 1;          /* t0 : 16(%rsp) */
4     double t1 = (double)n;   /* t1 : 24(%rsp) */
5     double dx = t0 / t1;     /* dx : 32(%rsp) */
6     double s = 0;           /* s : 40(%rsp) */
7     long i = 0;             /* i : 48(%rsp) */
8     if (!(i < n)) goto Lend;
9     Lstart:
10    double t2 = f(x);        /* t2 : 56(%rsp) */
11    s += t2;
12    x += dx;
13    i += 1;
14    if (i < n) goto Lstart;
15    Lend:
16    double t3 = s * dx;      /* t3 : 64(%rsp) */
17    return t3;
18 }
```

converting to “goto”s and “C = A op B”s

```
1 double integ3(long n) {      /* n : 0(%rsp) */
2   double x = 0;              /* x : 8(%rsp) */
3   double t0 = 1;             /* t0 : 16(%rsp) */
4   double t1 = (double)n;     /* t1 : 24(%rsp) */
5   double dx = t0 / t1;       /* dx : 32(%rsp) */
6   double s = 0;              /* s : 40(%rsp) */
7   long i = 0;                 /* i : 48(%rsp) */
8   if (!(i < n)) goto Lend;
9   Lstart:
10  double t2 = f(x);           /* t2 : 56(%rsp) */
11  s += t2;
12  x += dx;
13  i += 1;
14  if (i < n) goto Lstart;
15  Lend:
16  double t3 = s * dx;         /* t3 : 64(%rsp) */
17  return t3;
18 }
```


機械語 / Machine code

```
1 double integ3(long n) {
2     /* n : 0(%rsp) */
3     movq %rdi,0(%rsp)
4     # double x = 0;
5     /* x : 8(%rsp)*/
6     movsd .L0(%rip),%xmm0
7     movsd %xmm0,8(%rsp)
8     # double t0 = 1;
9     /* t0 : 16(%rsp)*/
10    movq $1,16(%rsp)
11    # double t1 = (double)n;
12    /* t1 : 24(%rsp)*/
13    cvtsi2sdq 0(%rsp),%xmm0
14    movsd %xmm0,24(%rsp)
15    # double dx = t0 / t1;
16    /* dx : 32(%rsp) */
17    movsd 16(%rsp),%xmm0
18    divsd 24(%rsp),%xmm0
19    movsd %xmm0,32(%rsp)
20    # double s = 0;
21    /* s : 40(%rsp) */
22    movsd .L0(%rip),%xmm0
23    movsd %xmm0,40(%rsp)
```

```
1     # long i = 0;
2     /* i : 48(%rsp) */
3     movq $0,48(%rsp)
4     # if (!(i < n)) goto Lend;
5     movq 0(%rsp),%rdi
6     cmpq 48(%rsp),%rdi # n - i
7     jle .Lend
8     .Lstart:
9     # double t2 = f(x);
10    /* t2 : 56(%rsp) */
11    movq 8(%rsp),%rdi
12    call f
13    movq %rax,56(%rsp)
14    # s += t2;
15    movq 40(%rsp),%xmm0
16    addsd 56(%rsp),%xmm0
17    movq %xmm0,40(%rsp)
18    # x += dx;
19    movsd 8(%rsp),%xmm0
20    addsd 32(%rsp),%xmm0
21    movsd %xmm0,8(%rsp)
```

機械語 / Machine code

```
1   # i += 1;
2   movq 48(%rsp),%rdi
3   addq $1,%rdi
4   movq %rdi,48(%rsp)
5   # if (i < n) goto Lstart;
6   movq 0(%rsp),%rdi
7   cmpq 48(%rsp),%rdi # n - i
8   jg .Lstart
9   .Lend:
10  movsd 40(%rsp),%xmm0
11  addsd 32(%rsp),%xmm0
12  addsd %xmm0,64(%rsp)
13  # return t3;
14  addsd 64(%rsp),%xmm0
15  ret
16 }
```

Contents

- ① コンパイラの仕事 / What do compilers do basically?
- ② より一般の場合 / More general cases
- ③ 最小限のCからのコード先生器 / Implementing a minimum C compiler
- ④ 中間言語 / Intermediate Representation

コード生成器 — 演習での前提

- 型は long (8 バイト整数) のみ
 - ▶ したがって typedef など無し
 - ▶ int もなし, 浮動小数点数もポインタもなし
 - ▶ 全部 long だから静的な型検査もいらぬ
- 大域変数もなし ⇒
 - ▶ プログラム = 関数定義のリスト
- ややこしい文は if, while, 複合文 ({ ... }) のみ
- 変数宣言は複合文の先頭で, 初期化の式もなし
- 以上は字句の定義 (cc_lex.mll), 文法の定義 (cc_parse.mly) に反映されている

Code generator — scope of the exercise

- all data types are `long` (8 byte integers)
 - ▶ no `typedefs`
 - ▶ no `ints`, floating point numbers, or pointers
 - ▶ everything is long, so `type checks` are unnecessary
- no global variables \Rightarrow
 - ▶ a program = list of function definitions
- supported complex statements are `if`, `while` and `compound statement` (`{ ... }`) only
- all variable definitions must come at the beginning of a block and initializes (`long x = expr`) are not supported
- they are expressed in token definitions (`cc_lex.mll`), and grammar definitions (`cc_parse.mly`)

プログラムの構成

- `cc_ast.ml` — 構文木定義
- `cc_parse.mly` — 文法定義
- `cc_lex.mll` — 字句定義
- `cc_cogen.ml` — 構文木からコード生成
- `cc.ml` — メイン

演習のほとんどの部分は, `cc_cogen.ml` で行われるだろう

Structure of the program

- `cc_ast.ml` — abstract syntax tree (AST) definition
- `cc_parse.mly` — grammar definition
- `cc_lex.mll` — lexer definition
- `cc_cogen.ml` — code generation from AST
- `cc.ml` — main driver

your work in this exercise will be mostly done in `cc_cogen.ml`

構文木定義 (cc_ast.ml) — 関数定義

- Cの関数定義の例

```
1 long f (long x, long y) {  
2     return x + y;  
3 }
```

- ⇒ 関数定義の構文木の定義

```
1 type definition =  
2     FUN_DEF of (type_expr * string * (type_expr * string) list * stmt)
```


AST (`cc_ast.ml`) — function definition

- an example C function definition

```
1 long f (long x, long y) {  
2     return x + y;  
3 }
```

- \Rightarrow AST definition for function definition

```
1 type definition =  
2   FUN_DEF of (type_expr * string * (type_expr * string) list * stmt)
```

構文木の定義 — 文

- if 文

```
1 if (x < y) { x++; return x; } else return y;
```

⇒

```
1 STMT_IF of (expr * stmt * stmt)
```

- 複合文

```
1 { long r; if (x < y) r = 10; else r = 20; }
```

⇒

```
1 STMT_COMPOUND of ((type_expr * string) list * stmt list)
```

AST — statements

- if statement

```
1 if (x < y) { x++; return x; } else return y;
```

⇒

```
1 STMT_IF of (expr * stmt * stmt)
```

- blocks

```
1 { long r; if (x < y) r = 10; else r = 20; }
```

⇒

```
1 STMT_COMPOUND of ((type_expr * string) list * stmt list)
```

構文木の定義 — 文

- while 文

```
1 while (i < n) { foo(i); i++; }
```

⇒

```
1 STMT_WHILE of (expr * stmt)
```

- ⇒ (諸々まとめた) 文の構文木の定義

```
1 type stmt =  
2   STMT_EMPTY  
3 | STMT_CONTINUE  
4 | STMT_BREAK  
5 | STMT_RETURN of expr (* e.g., return 123; *)  
6 | STMT_EXPR of expr (* e.g., f(x); *)  
7 | STMT_COMPOUND of ((type_expr * string) list * stmt list)  
8 | STMT_IF of (expr * stmt * stmt)  
9 | STMT_WHILE of (expr * stmt)
```

AST — statements

- while 文

```
1 while (i < n) { foo(i); i++; }
```

⇒

```
1 STMT_WHILE of (expr * stmt)
```

- ⇒ (putting them together) AST definition for statements

```
1 type stmt =  
2   STMT_EMPTY  
3 | STMT_CONTINUE  
4 | STMT_BREAK  
5 | STMT_RETURN of expr (* e.g., return 123; *)  
6 | STMT_EXPR of expr (* e.g., f(x); *)  
7 | STMT_COMPOUND of ((type_expr * string) list * stmt list)  
8 | STMT_IF of (expr * stmt * stmt)  
9 | STMT_WHILE of (expr * stmt)
```

構文木の定義 — 式

- 2項演算

1 `x + y + 1`

⇒

1 `EXPR_BIN_OP of bin_op * expr * expr`

注: 代入 ($a = b$) も 2項演算の一種 (Cの代入は, 文ではなく式)

- 関数呼び出し

1 `... f(x + 1, y + 2, z + 3) ...`

⇒

1 `EXPR_CALL of (string * expr list)`

AST — expressions

- binary operations

1 `x + y + 1`

⇒

1 `EXPR_BIN_OP of bin_op * expr * expr`

Note: assignment (`a = b`) is a kind of binary operation (C's assignment is not a statement but an expression)

- function call

1 `... f(x + 1, y + 2, z + 3) ...`

⇒

1 `EXPR_CALL of (string * expr list)`

構文木の定義 — 式

- ⇒ (諸々まとめた) 式の構文木の定義

```
1 type expr =  
2   EXPR_NUM of int      (* e.g., 3 *)  
3 | EXPR_VAR of string  (* e.g., x *)  
4 | EXPR_BIN_OP of bin_op * expr * expr  
5 | EXPR_UN_OP of un_op * expr (* e.g., -f(x) *)  
6 | EXPR_CALL of (string * expr list)
```


- \Rightarrow (putting them together) AST definitions for expressions

```
1 type expr =  
2   EXPR_NUM of int      (* e.g., 3 *)  
3 | EXPR_VAR of string  (* e.g., x *)  
4 | EXPR_BIN_OP of bin_op * expr * expr  
5 | EXPR_UN_OP of un_op * expr (* e.g., -f(x) *)  
6 | EXPR_CALL of (string * expr list)
```

コード生成 (cc_cogen.ml) — 基本スタイル

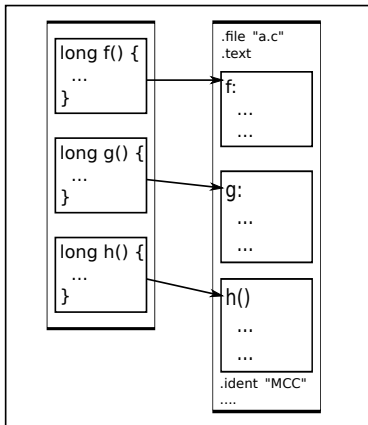
- 構文木 (AST) を受け取り, 対応する機械語 (「命令」のリスト) を返す
- ある構文木に対する機械語の生成 \approx その構成要素に対する機械語を適切に並べる
- プログラム全体 (program) \rightarrow 関数定義 (definition) \rightarrow 文 (stmt) \rightarrow 式 (expr)
- コード生成器のプログラムの見た目は, 多数の (1) 構文木に対するパターンマッチ (match) と (2) 子ノードに対する再帰呼出し

Code generation (`cc_cogen.ml`) — basic structure

- takes a parse tree (AST) and returns machine code (a list of instructions)
- generating machine code for an AST \approx arrange machine code for its components
- the program (`program`) \rightarrow function definition (`definition`) \rightarrow statement (`stmt`) \rightarrow expression (`expr`)
- code generator has lots of (1) pattern matching (`match`) against AST and (2) recursive calls to child trees

ファイル全体のコンパイル

- ≈ 関数毎にコンパイルしたものを連結



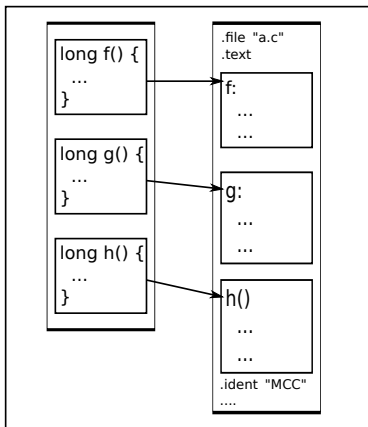
コード概形

```
1 let cogen_program defs ... =  
2   (gen_header ...)  
3   @ List.concat (List.map (fun def ->  
4     cogen_def def ...) defs)  
   (gen_trailer ...)
```

注: あくまで説明用の概形であって、
上記通りの形でなくてよい(こだわってはいけない)

Compiling an entire file

- \approx concatenate compilation of individual function definitions



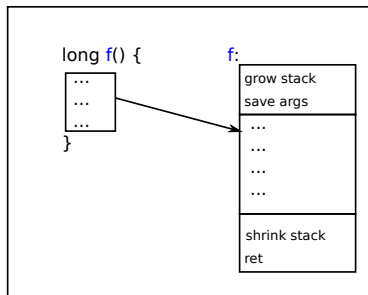
It will look like ...

```
1 let cogen_program defs ... =  
2   (gen_header ...)  
3   @ List.concat (List.map (fun def ->  
4     cogen_def def ...) defs)  
   (gen_trailer ...)
```

Note: the above is an outline for the illustration purpose you do not have to (should not) stick to

関数定義のコンパイル

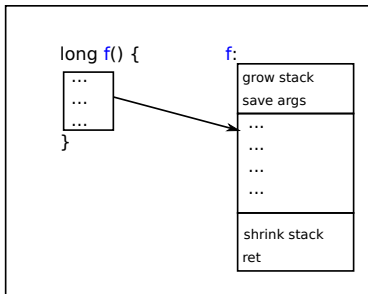
- \approx 文をコンパイルしたものの前後に, プロローグ (スタックを伸ばす, etc.), エピローグ (スタックを縮める, ret, etc.) をつける



```
1 let cogen_def def ... =  
2   match def with FUN_DEF(ret_type, f, params  
3     , body) =  
4     (gen_prologue def)  
5     @ (cogen_stmt body ...)  
6     @ (gen_epilogue def)
```

Compiling a function definition

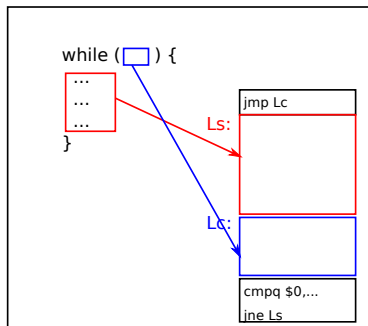
- \approx compile the body (statement); put prologue (grow the stack, etc.) and epilogue (shrink the stack, ret, etc.)



```
1 let cogen_def def ... =  
2   match def with FUN_DEF(ret_type, f, params  
3     , body) =  
4     (gen_prologue def)  
5     @ (cogen_stmt body ...)  
6     @ (gen_epilogue def)
```

文のコンパイル (例: while 文)

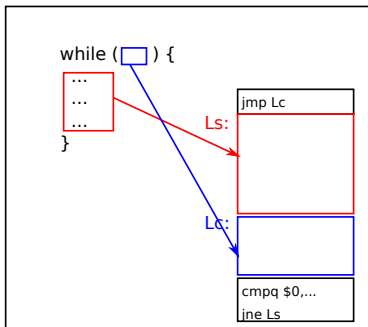
- \approx 条件式, 本体をコンパイルしたものを以下のように配置.
ループの継続判定コードをつける



```
1 let rec cogen_stmt stmt ... =  
2   match stmt with  
3     ...  
4   | Cc_ast.STMT_WHILE(cond, body) ->  
5     let cond_op, cond_insns = cogen_expr cond  
6       ... in  
7     let body_insns = cogen_stmt stmt ... in  
8     let ... in  
9     [ jmp Lc;  
10      Ls ]  
11    @ body_insns  
12    [ Lc ]  
13    @ cond_insns @  
14    [ cmpq $0, cond_op;  
      jne Ls ]
```


Compiling a statement (e.g., **while** statement)

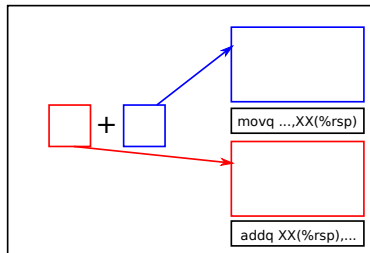
- \approx place compilation of the condition expression and the body as follows. add a conditional to determine if the loop continues



```
1 let rec cogen_stmt stmt ... =  
2   match stmt with  
3     ...  
4   | Cc_ast.STMT_WHILE(cond, body) ->  
5     let cond_op, cond_insns = cogen_expr cond  
6       ... in  
7     let body_insns = cogen_stmt stmt ... in  
8     let ... in  
9     [ jmp Lc;  
10      Ls ]  
11    @ body_insns  
12    [ Lc ]  
13    @ cond_insns @  
14    [ cmpq $0, cond_op;  
      jne Ls ]
```

式のコンパイル (算術演算)

- \approx 引数をそれぞれコンパイル; 演算命令

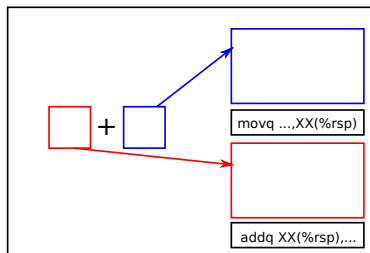


```
1 let rec cogen_expr expr ... =
2   match expr with
3   ...
4   | Cc_ast.EXPR_BIN_OP(op, e0, e1) ->
5     let insns1,op1 = cogen_expr e1 ... in
6     let insns0,op0 = cogen_expr e0 ... in
7     let m = スタック上のスロット in
8     ((insns1
9      @ [ movq op1,m ]
10     @ insns0
11     @ [ op m,op0 ]), (* op0 = op0 op m *))
12     op0)
13   | ...
```

- 注: `movq XX(%rsp),...` は第一オペランドの結果をスタックに格納し, 第二オペランドの評価中に壊されることがないようにしている
- 最も単純な方式 = 「すべての中間結果をスタックに格納する」に沿った方式

Compiling an expression (arithmetic)

- \approx compile the arguments; an arithmetic instruction



```
1 let rec cogen_expr expr ... =  
2   match expr with  
3     ...  
4   | Cc_ast.EXPR_BIN_OP(op, e0, e1) ->  
5     let insns1,op1 = cogen_expr e1 ... in  
6     let insns0,op0 = cogen_expr e0 ... in  
7     let m = a slot on the stack in  
8       ((insns1  
9         @ [ movq op1,m ]  
10        @ insns0  
11        @ [ op m,op0 ]), (* op0 = op0 op m *)  
12        op0)  
13    | ...
```

- Remark: `movq XX(%rsp), ...` saves the first operand, ensuring it won't be destroyed during the evaluation of the second
- remember we are following the simplest strategy = “save all intermediate results on the stack”

式のコンパイル (比較演算)

- $A < B$ は,
 - ▶ $A < B$ ならば 1
 - ▶ $A \geq B$ ならば 0という値を持つ式
- 式が許される任意の場所に現れうることに注意
 - ▶ $z = x < y$, $(x < y) + z$, $f(x < 1)$ のような式も許される (if や while の条件部分に来るとは限らない)
- アセンブリ言語でこれを生成する命令は?
 - ① 条件分岐
 - ② 条件つき set 命令. 例:

```
1 movq $0,%rax
2 cmpq %rdi,%rsi
3 setle %al
```

で, $\%rsi - \%rdi \leq 0$ (less-than-or-equal) ならば, $\%al$ ($\%rax$ の下 8 bit) に 1 がセットされる

Compiling an expression (comparison)

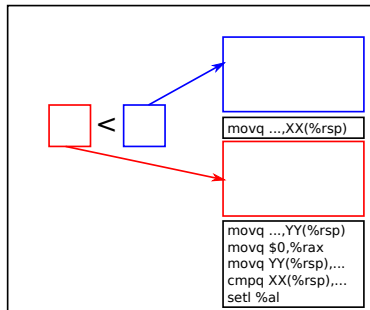
- $A < B$ is an expression that evaluates to
 - ▶ 1 if $A < B$
 - ▶ 0 if $A \geq B$
- no single instruction exactly does this
- note that they can appear anywhere expression can
 - ▶ $z = x < y$, $(x < y) + z$, and $f(x < 1)$ are allowed (they do not necessarily appear in condition expression of `if` or `while`)
- how to do it in assembly code?
 - 1 conditional branch
 - 2 *conditional set instruction*. e.g.,

```
1 movq $0,%rax
2 cmpq %rdi,%rsi
3 setle %al
```

will set `%al` (the lowest 8 bits of `%rax`) to 1 when `%rsi - %rdi` ≤ 0 (less-than-or-equal)

式のコンパイル (比較演算)

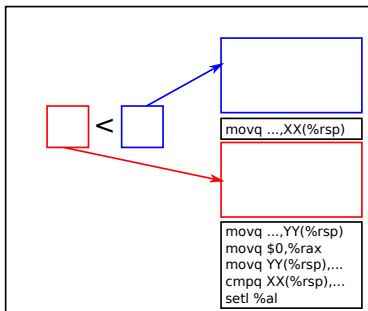
- \approx < の引数をそれぞれコンパイル; 比較; 条件付き set



```
1 let rec cogen_expr expr ... =
2   match expr with
3     ...
4   | Cc_ast.EXPR_CMP_OP(op, e0, e1) ->
5     let insns1,op1 = cogen_expr e1 ... in
6     let insns0,op0 = cogen_expr e0 ... in
7     let m0 = スタック上のスロット in
8     let m1 = スタック上のスロット in
9     ...
10    ((insns1
11      @ [ movq op1,m1 ]
12      @ insns0
13      @ [ movq op0,m0;
14          movq $0,%rax;
15          movq m0,op0;
16          cmpq m1,op0;
17          setop rax ]
18      op0)
19    | ...
```

Compiling an expression (comparison)

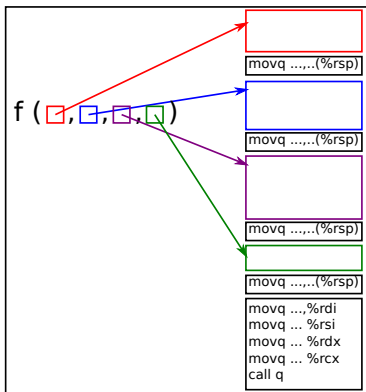
- \approx compile the arguments; compare; conditional set



```
1 let rec cogen_expr expr ... =
2   match expr with
3     ...
4   | Cc_ast.EXPR_CMP_OP(op, e0, e1) ->
5     let insns1,op1 = cogen_expr e1 ... in
6     let insns0,op0 = cogen_expr e0 ... in
7     let m0 = a slot on the stack in
8     let m1 = a slot on the stack in
9       ...
10    ((insns1
11      @ [ movq op1,m1 ]
12      @ insns0
13      @ [ movq op0,m0;
14          movq $0,%rax;
15          movq m0,op0;
16          cmpq m1,op0;
17          setop rax ]
18      op0)
19    | ...
```

式のコンパイル (関数呼び出し)

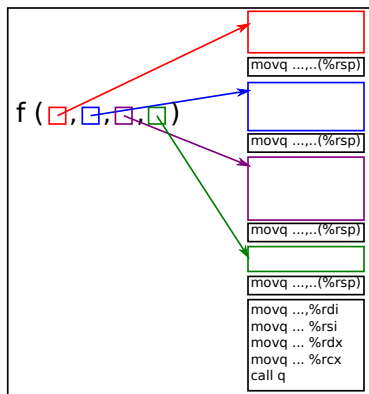
- ≈ 引数をそれぞれコンパイル; 引数を所定の位置に並べる;
call 命令



```
1 let rec cogen_expr expr ... =  
2   match expr with  
3     ...  
4   | Cc_ast.EXPR_CALL(f, args) ->  
5     let insns, arg_vars = cogen_exprs args  
6       env var_idx in  
       ((insns @ (make_call f arg_vars)), rax)
```


Compiling an expression (function call)

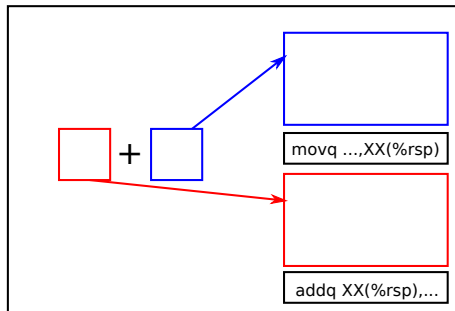
- \approx compile all arguments; put them to positions specified by ABI; a call instruction



```
1 let rec cogen_expr expr ... =
2   match expr with
3   ...
4   | Cc_ast.EXPR_CALL(f, args) ->
5     let insns, arg_vars = cogen_exprs
6       args env var_idx in
7     ((insns @ (make_call f arg_vars)),
8      rax)
```

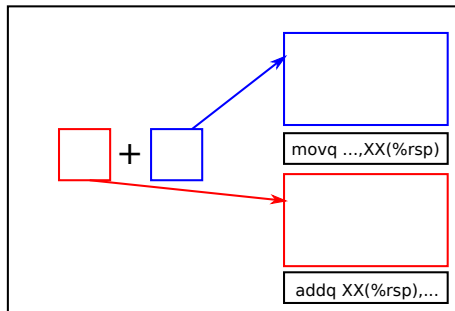
見過ごしてきた詳細

- 「部分式の値」や「変数の値」を保存しておく場所の決め方
- 以下の XX の決め方



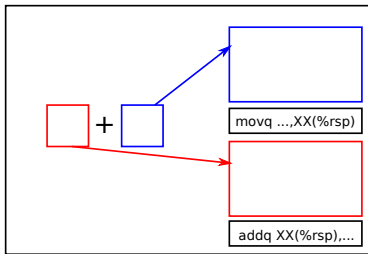
Details we have been leaving out

- how to determine locations to save values of *subexpressions* and *variables*
- that is, how to determine XX below



「部分式の値」の格納場所の決め方

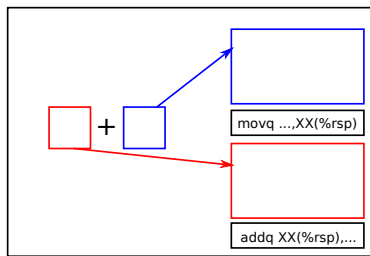
- `cogen_expr` に、空き領域の先頭を示す引数 (v) を渡す
意味: (`cogen_expr E v ...`) は, E を評価する命令列を生成;
それは $(\%rsp + v)$ 以上のアドレスのみを使う (破壊する)
- \rightarrow A + B の命令列中で, B の結果は $v(\%rsp)$ に保存
- A は $v + 8$ 以降を使う



```
1 let rec cogen_expr expr v =
2   match expr with
3     ...
4   | Cc_ast.EXPR_BIN_OP(op, e0, e1) ->
5     let insns1,op1 = cogen_expr e1 v ... in
6     let insns0,op0 = cogen_expr e0 (v + 8)
7       ... in
8     let m = v(%rsp) in
9     ((insns1
10      @ [ movq op1,m ]
11      @ insns0
12      @ [ op m,op0 ]), (* op0 = op0 op m *))
13   | ...
```

Determining where to save subexpressions

- `cogen_expr` receives a value (v) pointing to the free space spec: `cogen_expr E v ...` generates instructions that evaluate E using (destroying) only addresses above $(\%rsp + v)$
- \rightarrow when evaluating $A + B$, save B at $v(\%rsp)$
- let A use $v + 8$ and higher addresses



```
1 let rec cogen_expr expr v =
2   match expr with
3     ...
4   | Cc_ast.EXPR_BIN_OP(op, e0, e1) ->
5     let insns1,op1 = cogen_expr e1 v ... in
6     let insns0,op0 = cogen_expr e0 (v + 8)
7       ... in
8     let m = v(%rsp) in
9     ((insns1
10      @ [ movq op1,m ]
11      @ insns0
12      @ [ op m,op0 ]), (* op0 = op0 op m *))
13   | ...
```

「変数の値」の格納場所

- 例:

```
1  if (...) {  
2    long a, b, c;  
3    ...  
4  }
```

- 変数 a , b , c をスタック上に格納する必要がある
- 部分式の保存場所とほぼ同じ問題
- \rightarrow `cogen_stmt` にも, 空き領域の先頭を示す引数 v を渡す
意味: `cogen_stmt S v ...` は, S を実行する命令列を生成; それは $(\%rsp + v)$ 以上のアドレスのみを使う (破壊する)
- \rightarrow 例えば $a \mapsto v(\%rsp)$, $b \mapsto v + 8(\%rsp)$, $c \mapsto v + 16(\%rsp)$ に格納

Locations to hold variables

- ex:

```
1  if (...) {  
2    long a, b, c;  
3    ...  
4  }
```

- we need to hold **a**, **b**, **c** on the stack
- the problem is almost identical to saving values of subexpressions
- → `cogen_stmt` also takes v pointing to the beginning of the free space
spec: `cogen_stmt S v ...` generates instructions to execute S ;
they use (destroy) only addresses above $(\%rsp + v)$
- → e.g., hold $a \mapsto v(\%rsp)$, $b \mapsto v + 8(\%rsp)$,
 $c \mapsto v + 16(\%rsp)$

環境: 変数の格納場所の情報

- 「変数の値の格納場所」は、式に変数が出現した際にそれを取り出せる必要がある
 - ▶ 例: $x + 1$ をコンパイルするには, x の格納場所を知る必要
- 「変数 \mapsto 格納場所」の写像を管理するデータ構造 (環境) を作り, `cogen_stmt`, `cogen_expr` はそれを受け取るようにする
- 複合文 (`{ ... }`) の先頭で変数宣言が行われた時に, 環境に新たな写像を追加
- 環境は, 連想リスト (`List.assoc`) を用いて簡単に作れる

Environment: records where variables are held

- when a variable occurs in an expression, we need to get the location that holds the variable
 - ▶ ex: to compile `x + 1`, we need to know where `x` is held
- make a data structure that holds a mapping “variable \mapsto location” (*environment*) and pass it to `cogen_stmt` and `cogen_expr`
- when new variables are declared at the beginning of a compound statement (`{ ... }`), add new mappings to it
- an environment can be easily built with [association list](#) (`List.assoc`)

cogen_expr は環境を受け取る

```
1 let rec cogen_expr expr env v =  
2   match expr with  
3     ...  
4   | Cc_ast.EXPR_VAR(x) ->  
5     let loc = env_lookup x env in  
6     ([ movq loc,... ], ...)  
7   | ...
```

- `env_lookup x env` は、環境 `env` 中から `x` の格納位置を探す

cogen_expr receives an environment

```
1 let rec cogen_expr expr env v =  
2   match expr with  
3     ...  
4   | Cc_ast.EXPR_VAR(x) ->  
5     let loc = env_lookup x env in  
6     ([ movq loc,... ], ...)  
7   | ...
```

- `env_lookup x env` searches environment `env` for `x` and returns its location

cogen_stmt も環境を受け取る

```
1 let rec cogen_stmt expr env v =
2   match expr with
3     ...
4   | Cc_ast.STMT_COMPOUND(decls, stmts) ->
5     let env',v' = env_extend decls env v in
6     cogen_stmts stmts env' v' ...
7   | ...
```

- `env_extend decls env v` は,
 - ▶ 変数宣言 `decls` で宣言された変数に格納場所を割り当て (v , $v + 8$, $v + 16$, ...),
 - ▶ それらを環境 `env` に登録
 - ▶ 新しい環境 `env'` と空き領域 v' を返す

cogen_stmt receives an environment too

```
1 let rec cogen_stmt expr env v =  
2   match expr with  
3     ...  
4   | Cc_ast.STMT_COMPOUND(decls, stmts) ->  
5     let env',v' = env_extend decls env v in  
6     cogen_stmts stmts env' v' ...  
7   | ...
```

- `env_extend decls env v`

- ▶ assign locations (v , $v + 8$, $v + 16$, ...) to variables declared in *decls*
- ▶ register them in *env*
- ▶ return the new environment *env'* and the new free space *v'*

環境の実装

- 環境は (変数名, 格納場所) のリスト
- $loc = env_lookup\ x\ env$
環境 env において, 変数 x の格納場所 loc を返す (cf. `List.assoc`)
- $env' = env_add\ x\ loc\ env$
 env に, $x \mapsto loc$ が追加された環境 env' を返す ($(x, loc) :: env$)
- これを元に `env_extend decls env v` を作るのは演習問題

Implementing environment

- an environment is a list of (variable name, location)
- $loc = \text{env_lookup } x \text{ env}$
returns x 's location in environment env (cf. `List.assoc`)
- $env' = \text{env_add } x \text{ loc env}$
returns a new environment env' which has a new mapping $x \mapsto loc$ in addition to env ($(x, loc) :: env$)
- implementing `env_extend decls env v` based on this is your exercise

Contents

- ① コンパイラの仕事 / What do compilers do basically?
- ② より一般の場合 / More general cases
- ③ 最小限の C からのコード先生器 / Implementing a minimum C compiler
- ④ 中間言語 / Intermediate Representation

中間言語 (IR)

- 原理的には「構文木, 環境」から直接アセンブリ言語を出すことも可能だが, 色々な理由で, アセンブリ言語と似ているが少し違う「中間言語 (Intermediate Representation; IR)」を通すことが普通
- IR vs. アセンブリ

	IR	機械語
制御構造	≈ go to だけ	≈ go to だけ
式	≈ $C = A \text{ op } B$ だけ	≈ $C = A \text{ op } B$ だけ
局所変数の数	いくらでも	≈ レジスタ数まで
局所変数の寿命	関数実行中	≈ 関数呼び出しまで

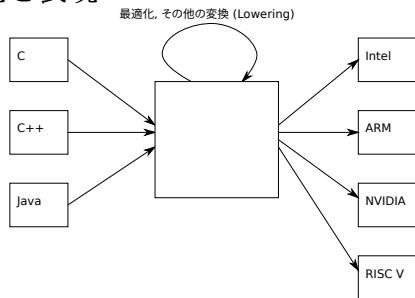
Intermediate representation (IR)

- while it is in theory possible to directly emit assembly code from a pair (AST, environment), most compilers first emit “Intermediate Representation; IR” which are similar to but different from the assembly code
- IR vs. assembly

	IR	assembly
control expression	\approx go to	\approx go to
the number of local variables	$\approx C = A \text{ op } B$ only arbitrary	$\approx C = A \text{ op } B$ only \approx only a fixed number of registers
lifetime of a local variable	during the function call that defined it	\approx some registers live only up to the next function call

中間言語 (IR) の存在理由

- 入力 → IR の変換を易しくする (任意個の, 寿命が関数呼び出しをまたがる変数を利用可能にする)
- 複数の入力言語の実装を容易にする — C, C++, Java, etc. で
入力言語 → IR 以外は共通
- 複数プロセッサへの実装を容易にする — Intel, ARM, etc. で,
IR → 機械語 以外は共通
- 最適化 — 「IR → IR の変換」または「IR → 機械語の変換」
で種々の最適化を表現



Why IR?

- simplify program \rightarrow IR conversion (allow arbitrary number of local variables, live across function calls)
- make it easy to implement multiple input languages — C, C++, Java, etc. can share everything but input \rightarrow IR
- make it easy to target multiple processors — Intel, ARM, etc. can share everything but IR \rightarrow machine code
- optimization — represent various optimizations as “IR \rightarrow IR” or “IR \rightarrow machine code” transformation

最適化, 其他の変換 (Lowering)

