Akim Demaille Étienne Renault Roland Levillain first.last@lrde.epita.fr

EPITA — École Pour l'Informatique et les Techniques Avancées

April 29, 2019

- Intermediate Representations
- 2 Memory Management
- 3 Translation to Intermediate Language
- 4 The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

- Intermediate Representations
 - Compilers Structure
 - Intermediate Representations
 - Tree
- 2 Memory Management
- Translation to Intermediate Language
- The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

Compilers Structure

- Intermediate Representations
 - Compilers Structure
 - Intermediate Representations
 - Tree
- 2 Memory Management
- Translation to Intermediate Language
- 4 The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

So many ends...

Ends:

```
front end analysis
middle end generic synthesis
back end specific synthesis
```

```
The gcc team suggests

front end name ("a front end").

front-end adjective ("the front-end interface").
```

Front Ends...

The front end is dedicated to analysis:

- lexical analysis (scanning)
- syntactic analysis (parsing)
- ast generation
- static semantic analysis (type checking, context sensitive checks)
- source language specific optimizations
- hir generation

... Back Ends

The back end is dedicated to specific synthesis:

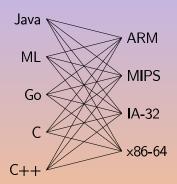
- instruction selection (mir to lir)
- register allocation
- assembly specific optimizations
- assembly code emission

... Middle Ends...

The middle end is dedicated to generic synthesis:

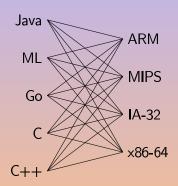
- stepwise refinement of hir to mir
- generic optimizations

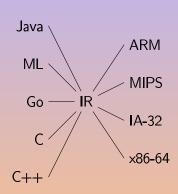
Retargetable Compilers





Retargetable Compilers





Other Compiling Strategies

- Intermediate language-based strategy: SmartEiffel, GHC
- Bytecode strategy: Java bytecode (JVM), CIL (.NET)
- Hybrid approaches: GCJ (Java bytecode or native code)
- Retargetable optimizing back ends: MLRISC, VPO (Very Portable Optimizer), and somehow C-- (Quick C--).
- Modular systems: LLVM (compiler as a library, centered on a typed IR). Contains the LLVM core libraries, Clang, LLDB, etc. Also:
 - VMKit: a substrate for virtual machines (JVM, etc.).
 - Emscripten: an LLVM-to-JavaScript compiler. Enables C/C++ to JS compilation.

Intermediate Representations (IR) are fundamental.

- Intermediate Representations
 - Compilers Structure
 - Intermediate Representations
 - Tree
- 2 Memory Management
- 3 Translation to Intermediate Language
- 4 The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

Format? Representation? Language?

Intermediate representation:

- a faithful model of the source program
- "written" in an abstract language, the intermediate language
- may have an external syntax
- may be interpreted/compiled (havm, byte code)
- may have different levels (gcc's Tree is very much like C).

What Language Flavor?

- Imperative?
 - Stack Based? (Java Byte-code)
 - Register Based? (gcc's rtl, tc's Tree)
- Functional? Most functional languages are compiled into a lower level language, eventually a simple λ -calculus.
- Other?

What Level?

A whole range of expressivities, typically aiming at making some optimizations easier:

• Keep array expressions?

Yes: adequate for dependency analysis and related optimizations,

No: Good for constant folding, strength reduction, loop invariant code motion, etc.

Keep loop constructs?

What level of machine independence?

Explicit register names?

Designing an Intermediate Representation



Intermediate-language design is largely an art, not a science.

— [Muchnick, 1997]

```
float a[20][10];
a[i][j+2];
```

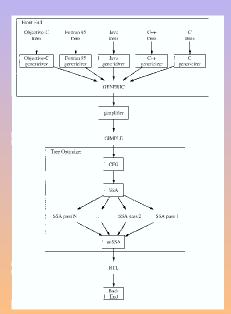
```
a[i][j+2];
t1 <- a[i,j+2]
```

float a[20][10];

```
float a[20][10];
. . .
a[i][j+2];
t1 <- a[i,j+2]
                        t1 < -j + 2
                        t2 <- i * 20
                        t3 < -t1 + t2
                        t4 <- 4 * t3
                        t5 <- addr a
                        t6 < - t5 + t4
                        t7 <- *t6
```

```
float a[20][10];
. . .
a[i][j+2];
t1 <- a[i,j+2]
                       t1 < -j + 2
                                             r1 < - [fp - 4]
                       t2 < -i * 20
                                             r2 < -r1 + 2
                       t3 < -t1 + t2
                                             r3 <- [fp - 8]
                       t4 <- 4 * t3
                                           r4 <- r3 * 20
                       t5 <- addr a
                                           r5 <- r4 + r2
                       t6 < - t5 + t4
                                           r6 <- 4 * r5
                       t7 <- *t6
                                            r7 <- fp - 216
                                             f1 < - [r7 + r6]
```

Different Levels: The GCC Structure



Stack Based: Java Byte-Code [Edwards, 2003]

```
class Gcd
  static public int gcd(int a, int b)
    while (a != b)
        if (a > b)
         a -= b;
        else
          b -= a:
    return a:
  static public int main(String[] arg)
    return gcd(12, 34);
```

Stack Based: Java Byte-Code

```
% gcj-3.3 -c gcd.java
                                    17: iload_1
% jcf-dump-3.3 -c gcd
                                    18: iload_0
                                    19: isub
Method name: "gcd" public static
                                    20: istore_1
Signature: 5=(int,int)int
                                    21: goto 0
Attribute "Code", length:66,
                                    24: iload_0
max_stack:2, max_locals:2,
                                    25: ireturn
                                   Attribute "LineNumberTable",
code_length:26
  0: iload_0
                                             length: 22, count: 5
  1: iload_1
                                     line: 5 at pc: 0
  2: if_icmpeq 24
                                     line: 7 at pc: 5
  5: iload 0
                                     line: 8 at pc: 10
  6: iload 1
                                     line: 10 at pc: 17
  7: if_icmple 17
                                     line: 12 at pc: 24
 10: iload_0
 11: iload_1
 12: isub
 13: istore 0
```

Stack Based [Edwards, 2003]

Advantages

- Trivial translation of expressions
- Trivial interpreters
- No pressure on registers
- Often compact

Disadvantages

- Does not fit with today's architectures
- Hard to analyze
- Hard to optimize

Stack Based [Edwards, 2003]

Advantages

- Trivial translation of expressions
- Trivial interpreters
- No pressure on registers
- Often compact

Disadvantages

- Does not fit with today's architectures
- Hard to analyze
- Hard to optimize

Stack Based: Examples

ucode, used in hp pa-risk, and mips, was designed for stack evaluation (HP 3000 is stack based).

Today it is less adequate.

mips translates it back and forth to triples for optimization.

hp converts it into sllic (Spectrum Low Level ir) [Muchnick, 1997].

Register Based: tc's Tree

Register Based: tc's Tree (1/4)

```
/* == High Level Intermediate representation. == */
# Routine: gcd
label 10
# Prologue
move temp t0 temp fp
move temp fp temp sp
move
 temp sp
  binop sub temp sp const 12
move
mem temp fp
    temp i0
move
  mem binop add temp fp const -4
  temp i1
move
  mem binop add temp fp const -8
  temp i2
```

Register Based: tc's Tree (2/4)

```
# Body
move temp rv
 eseq
 seq
  1abe1 12
  cjump ne mem binop add temp fp const -4
         mem binop add temp fp const -8
         name | 3 name | 1
  1abel 13
  seq
    cjump gt mem binop add temp fp const -4
          mem binop add temp fp const -8
          name 14 name 15
    label 4
    move mem binop add temp fp const -4
        binop sub mem binop add temp fp const -4
               mem binop add temp fp const -8
    jump name 16
```

Register Based: tc's Tree (3/4)

```
label 15
      move mem binop add temp fp const -8
           binop sub mem binop add temp fp const -8
                     mem binop add temp fp const -4
      label 16
    seq end
    jump name 12
    label 11
  seq end
  mem binop add temp fp const -4
# Epilogue
move temp sp temp fp
move temp fp temp t0
label end
```

Register Based: tc's Tree (4/4)

```
# Routine: main
label main
# Prologue
# Body
seq
  sxp
    call
      name print_int
      call name 10 temp fp const 42 const 51
      call end
    call end
  sxp
    const 0
seq end
# Epilogue
label end
```

How is the structure coded?

Addresses Expressions and instructions have names, or (absolute) addresses. (Stack based is a bit like a relative address).

- 2 address instructions? (*triples*)
- 3 address instructions? (quadruples)
- Tree Expressions and instructions are unnamed, related to each other
 - as nodes of trees
- dag Compact, good for local value numbering, but that's all.

How is the structure coded?

Addresses Expressions and instructions have names, or (absolute) addresses. (Stack based is a bit like a relative address).

- 2 address instructions? (triples)
- 3 address instructions? (quadruples)

Tree Expressions and instructions are unnamed, related to each other as nodes of trees

dag Compact, good for local value numbering, but that's all.

How is the structure coded?

- Addresses Expressions and instructions have names, or (absolute) addresses. (Stack based is a bit like a relative address).
 - 2 address instructions? (triples)
 - 3 address instructions? (quadruples)
 - Tree Expressions and instructions are unnamed, related to each other as nodes of trees
 - dag Compact, good for local value numbering, but that's all.

How is the structure coded?

- Addresses Expressions and instructions have names, or (absolute) addresses. (Stack based is a bit like a relative address).
 - 2 address instructions? (triples)
 - 3 address instructions? (quadruples)
 - Tree Expressions and instructions are unnamed, related to each other as nodes of trees
 - dag Compact, good for local value numbering, but that's all.

How is the structure coded?

- Addresses Expressions and instructions have names, or (absolute) addresses. (Stack based is a bit like a relative address).
 - 2 address instructions? (triples)
 - 3 address instructions? (quadruples)
 - Tree Expressions and instructions are unnamed, related to each other as nodes of trees
 - dag Compact, good for local value numbering, but that's all.

Quadruples vs. Triples [Muchnick, 1997]

```
L1: i <- i + 1
(2) i sto

t1 <- i + 1
(3) i + 1

t2 <- p + 4
(4) p + 4

t3 <- *t2
(5) *(4)

p <- t2
(6) p sto

t4 <- t1 < 10
(7) (3) <

*r <- t3
if t4 goto L1
(9) if (7)
```

Quadruples vs. Triples [Muchnick, 1997]

Register Based: gcc's rtl

```
int
gcd(int a, int b)
  while (a != b)
      if (a > b)
        a -= b;
      else
        b -= a;
  return a;
```

Register Based: gcc's rtl

```
(jump_insn 15 14 16 (set (pc)
        (if_then_else (ne (reg:CCZ 17 flags)
                (const_int 0 [0x0])
            (label_ref 18)
            (pc))) -1 (nil)
    (nil))
(jump_insn 16 15 17 (set (pc)
        (label ref 44)) -1 (nil)
    (nil))
(barrier 17 16 18)
(code label 18 17 19 4 "" "" [0 uses])
(note 19 18 20 NOTE_INSN_LOOP_END_TOP_COND)
(note 20 19 21 NOTE_INSN_DELETED)
(note 21 20 22 NOTE INSN DELETED)
```

```
(note 22 21 25 ("gcd.c") 6)
(insn 25 22 26 (set (reg:SI 60)
        (mem/f:SI (reg/f:SI 53 virtual-incoming-args) [0 a+0 S4 A32]))
    -1 (nil) (nil))
(insn 26 25 27 (set (reg:CCGC 17 flags)
        (compare:CCGC (reg:SI 60)
            (mem/f:SI (plus:SI (reg/f:SI 53 virtual-incoming-args)
                    (const_int 4 [0x4])) [0 b+0 S4 A32]))) -1 (nil)
    (nil))
(jump_insn 27 26 28 (set (pc)
        (if_then_else (le (reg:CCGC 17 flags)
                (const_int 0 [0x0]))
            (label_ref 34)
            (pc))) -1 (nil)
    (nil))
```

```
(note 28 27 30 ("gcd.c") 7)
(insn 30 28 31 (set (reg:SI 61)
        (mem/f:SI (plus:SI (reg/f:SI 53 virtual-incoming-args)
                (const_int 4 [0x4])) [0 b+0 S4 A32])) -1 (nil)
   (nil))
(insn 31 30 32 (parallel[
            (set (mem/f:SI (reg/f:SI 53 virtual-incoming-args)
                [0 a+0 S4 A32])
                (minus:SI (mem/f:SI (reg/f:SI 53 virtual-incoming-args)
                [0 a+0 S4 A32])
                    (reg:SI 61)))
           (clobber (reg:CC 17 flags))
       ] ) -1 (nil)
   (expr_list:REG_EQUAL (minus:SI (mem/f:SI (reg/f:SI 53
                                virtual-incoming-args) [0 a+0 S4 A32])
            (mem/f:SI (plus:SI (reg/f:SI 53 virtual-incoming-args)
                    (const_int 4 [0x4])) [0 b+0 S4 A32]))
        (nil)))
(jump_insn 32 31 33 (set (pc)
       (label_ref 39)) -1 (nil)
   (nil))
(barrier 33 32 34)
(code_label 34 33 35 5 "" "" [0 uses])
```

```
(note 35 34 37 ("gcd.c") 9)
(insn 37 35 38 (set (reg:SI 62)
        (mem/f:SI (reg/f:SI 53 virtual-incoming-args) [0 a+0 S4 A32]))
    -1 (nil) (nil))
(insn 38 37 39 (parallel[
            (set (mem/f:SI (plus:SI (reg/f:SI 53 virtual-incoming-args)
                        (const_int 4 [0x4])) [0 b+0 S4 A32])
                (minus:SI (mem/f:SI (plus:SI (reg/f:SI 53
                                              virtual-incoming-args)
                            (const_int 4 [0x4])) [0 b+0 S4 A32])
                    (reg:SI 62)))
           (clobber (reg:CC 17 flags))
       1) -1 (nil)
   (expr_list:REG_EQUAL (minus:SI (mem/f:SI (plus:SI (reg/f:SI
                                             53 virtual-incoming-args)
                    (const_int 4 [0x4])) [0 b+0 S4 A32])
            (mem/f:SI (reg/f:SI 53 virtual-incoming-args) [0 a+0 S4 A32]))
        (nil)))
(code_label 39 38 41 6 "" "" [0 uses])
(jump_insn 41 39 42 (set (pc)
       (label ref 10)) -1 (nil)
   (nil))
(barrier 42 41 43)
(note 43 42 44 NOTE_INSN_LOOP_END)
```

```
(note 45 44 46 ("gcd.c") 11)
(note 46 45 47 NOTE_INSN_DELETED)
(note 47 46 49 NOTE_INSN_DELETED)
(insn 49 47 51 (set (reg:SI 64)
        (mem/f:SI (reg/f:SI 53 virtual-incoming-args) [0 a+0 S4 A32])) -1 (nil)
    (nil))
(insn 51 49 52 (set (reg:SI 58)
        (reg:SI 64)) -1 (nil)
    (nil))
(jump_insn 52 51 53 (set (pc)
        (label_ref 56)) -1 (nil)
   (nil))
(barrier 53 52 54)
(note 54 53 55 NOTE_INSN_FUNCTION_END)
(note 55 54 59 ("gcd.c") 12)
(insn 59 55 60 (clobber (reg/i:SI 0 eax)) -1 (nil)
   (nil))
(insn 60 59 56 (clobber (reg:SI 58)) -1 (nil)
    (nil))
(code_label 56 60 58 1 "" "" [0 uses])
(insn 58 56 61 (set (reg/i:SI 0 eax)
        (reg:SI 58)) -1 (nil)
   (nil))
(insn 61 58 0 (use (reg/i:SI 0 eax)) -1 (nil)
```

Register Based [Edwards, 2003]

Advantages

- Suits today's architectures
- Clearer data flow

Disadvantages

- Harder to synthesize
- Less compact
- Harder to interpret

Register Based [Edwards, 2003]

Advantages

- Suits today's architectures
- Clearer data flow

Disadvantages

- Harder to synthesize
- Less compact
- Harder to interpret

Tree

- Intermediate Representations
 - Compilers Structure
 - Intermediate Representations
 - Tree
- 2 Memory Management
- Translation to Intermediate Language
- 4 The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

Tree [Appel, 1998]

A simple intermediate language:

- Tree structure (no kidding...)
- Unbounded number of registers (temporaries)
- Two way conditional jump

Tree: Grammar

```
\langle \text{Exp} \rangle ::= \text{"const"} \text{ int}
         "name" (Label)
         "temp" (Temp)
             "binop" (Oper) (Exp) (Exp)
             "mem" (Exp)
             "call" \langle \text{Exp} \rangle [{\langle \text{Exp} \rangle}] "call end"
             "eseq" (Stm) (Exp)
\langle Stm \rangle ::= "move" \langle Exp \rangle \langle Exp \rangle
         "sxp" (Exp)
         | "jump" \langle Exp \rangle [{\langle Label \rangle}]
            "cjump" (Relop) (Exp) (Exp) (Label) (Label)
          "seq" [{\langle Stm \rangle}] "seq end"
            "label" (Label)
(Oper) ::= "add" | "sub" | "mul" | "div" | "mod"
⟨Relop⟩ ::= "eq" | "ne" | "lt" | "gt" | "le" | "ge"
```

Tree Samples

```
\% echo '1 + 2 * 3' | tc -H -
/* == High Level Intermediate representation. == */
# Routine: Main Program
label Main
# Prologue
# Body
sxp
    binop add
        const 1
        binop mul
            const 2
            const 3
# Epilogue
label end
```

Tree Samples

```
% echo 'if 1 then print_int (1)' | tc -H -
# Routine: Main Program
label Main
# Prologue
# Body
seq
    cjump ne, const 1, const 0, name 11, name 12
    label 11
    sxp call name print_int, const 1
    jump name 13
    label 12
    sxp const 0
    label 13
seq end
# Epilogue
label end
```

Memory Management

- Intermediate Representations
- 2 Memory Management
 - Memory Management
 - Activation Blocks
 - Nonlocal Variables
- Translation to Intermediate Language
- 4 The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

Memory Management

- Intermediate Representations
- 2 Memory Management
 - Memory Management
 - Activation Blocks
 - Nonlocal Variables
- Translation to Intermediate Language
- 4 The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

Memory Hierarchy [Appel, 1998]

Different kinds of memory in a computer, with different performances:

Registers Small memory units built on the cpu (bytes, 1 cycle)

L1 Cache Last main memory access results (kB, 2-3 cycles)

L2 Cache (MB, 10 cycles)

Memory The usual ram (GB, 100 cycles)

Storage Disks (100GB, TB, > 1Mcycles)

Use the registers as much as possible.

Register Overflow

What if there are not enough registers? Use the main memory, but how? Recursion:

Without Each name is bound once. It can be statically allocated a single unit of main memory. (Cobol, Concurrent Pascal, Fortran (unless recursive)).

With A single name can be part of several concurrent bindings.

Memory allocation must be dynamic.

Register Overflow

What if there are not enough registers? Use the main memory, but how? Recursion:

- Without Each name is bound once. It can be statically allocated a single unit of main memory. (Cobol, Concurrent Pascal, Fortran (unless recursive)).
 - With A single name can be part of several concurrent bindings.

 Memory allocation must be dynamic.

Depending on the persistence, several models:

Global Global objects, whose liveness is equal to that of the program, are statically allocated (e.g., static variables in C)

Automatic Liveness is bound to that of the host function (e.g., auto variables in C)

Heap Liveness is independent of function liveness:

Garbage Collected

With or without new

(12 m. Completelle MT)

(lisp, Smalltalk, ML, Haskell, Tiger, Perl etc.

Depending on the persistence, several models:

Global Global objects, whose liveness is equal to that of the program, are statically allocated (e.g., static variables in C)

Automatic Liveness is bound to that of the host function (e.g., auto variables in C)

Heap Liveness is independent of function liveness:

- Garbage Collected
- With or without new
 - VIIII OF WITHOUT HEW
 - (lisp, Smalltalk, ML, Haskell, Tiger, Perl etc.).

Depending on the persistence, several models:

```
Global Global objects, whose liveness is equal to that of the program, are statically allocated (e.g., static variables in C)
```

Automatic Liveness is bound to that of the host function (e.g., auto variables in C)

Heap Liveness is independent of function liveness:

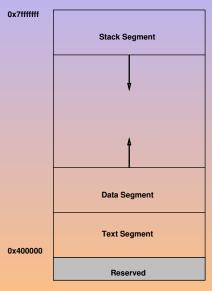
```
User Controlled
malloc/free (C), new/dispose (Pascal),
new/delete (C++) etc.

Garbage Collected
With or without new
(lisp, Smalltalk, ML, Haskell, Tiger, Perl etc.)
```

```
Depending on the persistence, several models:
    Global Global objects, whose liveness is equal to that of the program,
           are statically allocated
           (e.g., static variables in C)
Automatic Liveness is bound to that of the host function
           (e.g., auto variables in C)
     Heap Liveness is independent of function liveness:
           User Controlled
               malloc/free (C), new/dispose (Pascal),
               new/delete (C++) etc.
```

```
Depending on the persistence, several models:
    Global Global objects, whose liveness is equal to that of the program,
           are statically allocated
           (e.g., static variables in C)
Automatic Liveness is bound to that of the host function
           (e.g., auto variables in C)
     Heap Liveness is independent of function liveness:
           User Controlled
               malloc/free (C), new/dispose (Pascal),
               new/delete (C++) etc.
           Garbage Collected
               With or without new
                (lisp, Smalltalk, ML, Haskell, Tiger, Perl etc.).
```

spim Memory Model [Larus, 1990]



Stack Management

Function calls is a last-in first-out process, hence, it is properly represented by a stack.

Or...

"Call tree": the complete history of calls.

The execution of the program is its depth first traversal.

Depth-first walk requires a stack.

- Intermediate Representations
- 2 Memory Management
 - Memory Management
 - Activation Blocks
 - Nonlocal Variables
- Translation to Intermediate Language
- 4 The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

- In recursive languages, a single routine can be "opened" several times concurrently.
- An activation designates one single instance of execution.
- Automatic variables are bound to the liveness of the activation.
- Their location is naturally called activation block, or stack frame.

- In recursive languages, a single routine can be "opened" several times concurrently.
- An activation designates one single instance of execution.
- Automatic variables are bound to the liveness of the activation.
- Their location is naturally called *activation block*, or *stack frame*.

- In recursive languages, a single routine can be "opened" several times concurrently.
- An activation designates one single instance of execution.
- Automatic variables are bound to the liveness of the activation.
- Their location is naturally called activation block, or stack frame.

- In recursive languages, a single routine can be "opened" several times concurrently.
- An activation designates one single instance of execution.
- Automatic variables are bound to the liveness of the activation.
- Their location is naturally called activation block, or stack frame.

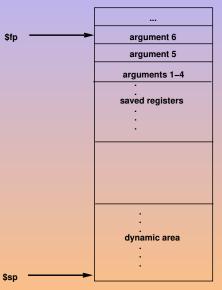
Activation Blocks Contents

Data to store on the stack:
 arguments incoming
local variables user automatic variables
return address where to return
saved registers the caller's environment to restore
temp compiler automatic variables, spills
static link when needed

Activation Blocks Layout

The layout is suggested by the constructor. Usually the layout is from earliest known, to latest.

Activation Blocks Layout on mips [Larus, 1990]





Frame & Stack Pointers

The stack of activation blocks is implemented as an array with frame pointer the inner frontier of the activation block stack pointer the outer frontier

Usually the stack is represented growing towards the bottom.

Flexible Automatic Memory

```
auto Static size, automatic memory.
malloc Dynamic size, persistent memory.
Automatic memory is extremely convenient...
int
open2(char* str1, char* str2, int flags, int mode)
{
   char name[strlen(str1) + strlen(str2) + 1];
   stpcpy(stpcpy(name, str1), str2);
   return open(name, flags, mode);
}
```

Flexible Automatic Memory

malloc is a poor replacement.

```
int
open2(char* str1, char* str2, int flags, int mode)
₹
  char* name
    = (char*) malloc(strlen(str1) + strlen(str2) + 1);
  if (name == 0)
    fatal("virtual memory exceeded");
  stpcpy(stpcpy(name, str1), str2);
  int fd = open(name, flags, mode);
 free(name):
 return fd;
```

Flexible Automatic Memory

alloca is a good replacement.

Advantages of alloca [Loosemore et al., 2003]

- Using alloca wastes very little space and is very fast.
 (It is open-coded by the GNU C compiler.)
- alloca does not cause memory fragmentation.
 Since alloca does not have separate pools for different sizes of block, space used for any size block can be reused for any other size.
- Automatically freed.
 Nonlocal exits done with longjmp automatically free the space allocated with alloca when they exit through the function that called alloca. This is the most important reason to use alloca.

Advantages of alloca [Loosemore et al., 2003]

- Using alloca wastes very little space and is very fast.
 (It is open-coded by the GNU C compiler.)
- alloca does not cause memory fragmentation.
 Since alloca does not have separate pools for different sizes of block, space used for any size block can be reused for any other size.
- Automatically freed.
 Nonlocal exits done with longjmp automatically free the space allocated with alloca when they exit through the function that called alloca. This is the most important reason to use alloca.

Advantages of alloca [Loosemore et al., 2003]

- Using alloca wastes very little space and is very fast.
 (It is open-coded by the GNU C compiler.)
- alloca does not cause memory fragmentation.
 Since alloca does not have separate pools for different sizes of block, space used for any size block can be reused for any other size.
- Automatically freed.
 Nonlocal exits done with longjmp automatically free the space allocated with alloca when they exit through the function that called alloca. This is the most important reason to use alloca.

Disadvantages of alloca [Loosemore et al., 2003]

- If you try to allocate more memory than the machine can provide, you
 don't get a clean error message.
 Instead you get a fatal signal like the one you would get from an
 infinite recursion; probably a segmentation violation.
- Some non-GNU systems fail to support alloca, so it is less portable However, a slower emulation of alloca written in C is available for use on systems with this deficiency.

Disadvantages of alloca [Loosemore et al., 2003]

- If you try to allocate more memory than the machine can provide, you
 don't get a clean error message.
 Instead you get a fatal signal like the one you would get from an
 infinite recursion; probably a segmentation violation.
- Some non-GNU systems fail to support alloca, so it is less portable. However, a slower emulation of alloca written in C is available for use on systems with this deficiency.

Arrays vs. Alloca [Loosemore et al., 2003]

- A variable size array's space is freed at the end of the scope of the name of the array.
 The space allocated with alloca remains until the end of the function.
- It is possible to use alloca within a loop, allocating an additional block on each iteration.

 This is impossible with variable sized arrays.

Arrays vs. Alloca [Loosemore et al., 2003]

- A variable size array's space is freed at the end of the scope of the name of the array.
 The space allocated with alloca remains until the end of the function.
- It is possible to use alloca within a loop, allocating an additional block on each iteration.
 - This is impossible with variable-sized arrays.

Implementing Dynamic Arrays & Alloca

- Playing with \$sp which makes \$fp mandatory.
- An additional stack (as with the C emulation of alloca).

Nonlocal Variables

- Intermediate Representations
- 2 Memory Management
 - Memory Management
 - Activation Blocks
 - Nonlocal Variables
- Translation to Intermediate Language
- 4 The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

escapes-n-recursion

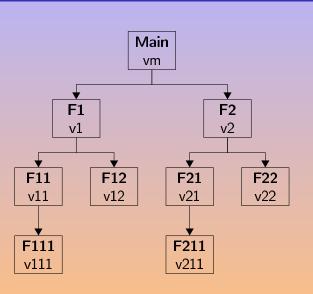
```
let function trace(fn: string, val: int) =
      (print(fn); print("("); print_int(val); print(") "))
    function one(input : int) =
    let function two() =
        (trace("two", input); one(input - 1))
    in
      if input > 0 then
        (two(); trace("one", input))
    end
in
  one(3); print("\n")
end
```

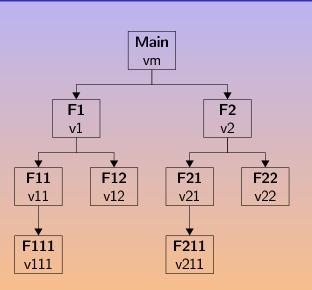
escapes-n-recursion

```
let function trace(fn: string, val: int) =
      (print(fn); print("("); print_int(val); print(") "))
    function one(input : int) =
    let function two() =
        (trace("two", input); one(input - 1))
    in
      if input > 0 then
        (two(); trace("one", input))
    end
in
  one(3); print("\n")
end
% tc -H escapes-n-recursion.tig > f.hir && havm f.hir
```

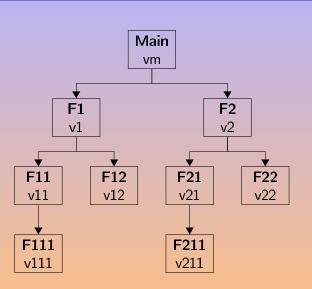
escapes-n-recursion

```
let function trace(fn: string, val: int) =
      (print(fn); print("("); print_int(val); print(") "))
    function one(input : int) =
    let function two() =
        (trace("two", input); one(input - 1))
    in
      if input > 0 then
        (two(); trace("one", input))
    end
in
  one(3); print("\n")
end
% tc -H escapes-n-recursion.tig > f.hir && havm f.hir
two(3) two(2) two(1) one(1) one(2) one(3)
```

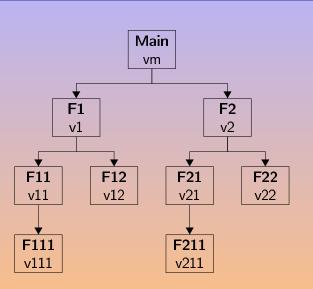




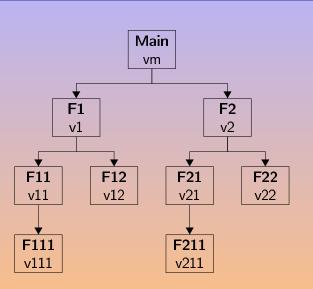
- Main uses vm
- Main calls F1
- F1 uses v1
- F1 uses vm, nor
- 5 F1 calls F11
- **6** F11 uses v11
- F11 uses v1
- 6 F11 uses vm
- F11 calls F12
 - U F12 Calls F1



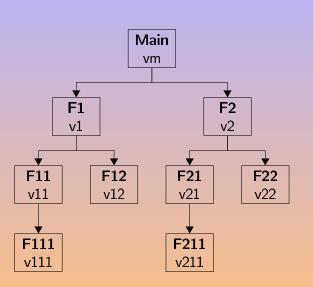
- Main uses vm
- Main calls F1
- 3 F1 uses v1
- local
- F1 calls F11F11 uses v11
- F11 uses v1
- **3** F11 uses vm
- F11 calls F12



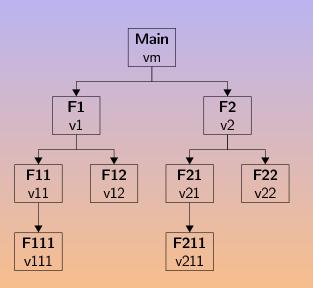
- Main uses vm
- 2 Main calls F1
- F1 uses v1
- 4 F1 uses vm, non local
- 5 F1 calls F11
- **6** F11 uses v11
- F11 uses v1
- **6** F11 uses vm
- F11 calls F12
- F12 calls F1



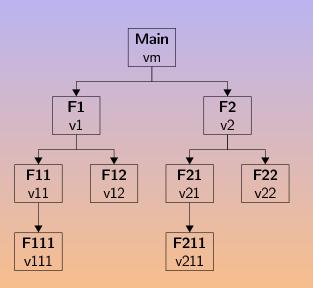
- Main uses vm
- Main calls F1
- F1 uses v1
- F1 uses vm, non
 local
- 5 F1 calls F11
- F11 uses v11
- **7** F11 uses v1
- F11 uses vm
- F11 calls F12
- F12 calls F1



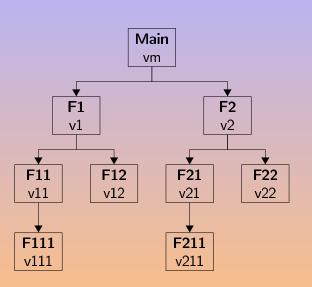
- Main uses vm
- Main calls F1
- 3 F1 uses v1
- 4 F1 uses vm, non local
- 5 F1 calls F11
- 9 F11 uses --1
- F11 uses tm
- F11 calls F19
- FII CallS F12
- F12 calls F1



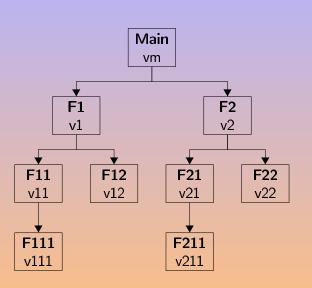
- Main uses vm
- Main calls F1
- 3 F1 uses v1
- 4 F1 uses vm, non local
- F1 calls F11
- **5** F11 uses v11
- F11 uses v1
- 6 F11 uses vm
- F11 calls F12
- F12 calls F1



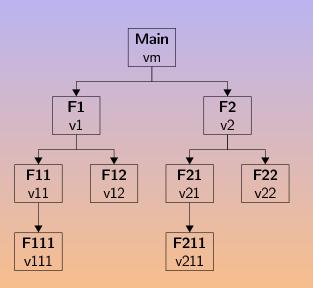
- Main uses vm
- Main calls F1
- 3 F1 uses v1
- F1 uses vm, non local
- F1 calls F11
- **o** F11 uses v11
- F11 uses v1
- F11 uses vm
- F11 calls F12
- F12 calls F1



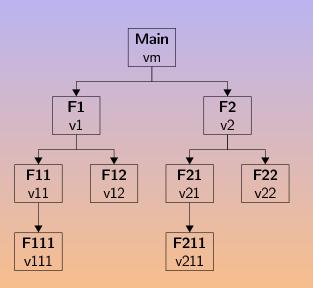
- Main uses vm
- Main calls F1
- F1 uses v1
- F1 uses vm, non local
- F1 calls F11
- **o** F11 uses v11
- F11 uses v1
- F11 uses vm
- F11 calls F12
- F12 calls F1



- Main uses vm
- Main calls F1
- 3 F1 uses v1
- F1 uses vm, non local
- F1 calls F11
- **o** F11 uses v11
- F11 uses v1
- F11 uses vm
- F11 calls F12
- F12 calls F1



- Main uses vm
- Main calls F1
- 3 F1 uses v1
- F1 uses vm, non local
- F1 calls F11
- **6** F11 uses v11
- F11 uses v1
- 3 F11 uses vm
- F11 calls F12
- F12 calls F1



- Main uses vm
- Main calls F1
- 3 F1 uses v1
- F1 uses vm, non local
- F1 calls F11
- **o** F11 uses v11
- F11 uses v1
- 3 F11 uses vm
- F11 calls F12
- F12 calls F1

The caller must provide the callee with its static link.

Caller	Callee	Static Link
Main	F1	$fp_{Main} = fp$
F1	F11	$fp_{F1} = fp$
F11	F12	$fp_{F1} = sl_{F11} = *fp_{F11} = *fp$
F12	F2	$fp_{Main} = sl_{F1} = *sl_{F12} = **fp_{F12} = **fp$
F2	F22	$fp_{F2} = fp$
F22	F11	fp _{F1} = ???

Assuming that the static link is stored at fp.

Higher Order Functions

```
let
  function addgen (a: int) : int -> int =
    let.
      function res (b: int) : int =
        a + b
    in
      res
    end
  var add50 := addgen (50)
in
  add50 (1)
end
```

Translation to Intermediate Language

- Intermediate Representations
- 2 Memory Management
- 3 Translation to Intermediate Language
 - Calling Conventions
 - Clever Translations
 - Complex Expressions
- 4 The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

Calling Conventions

- Intermediate Representations
- 2 Memory Management
- 3 Translation to Intermediate Language
 - Calling Conventions
 - Clever Translations
 - Complex Expressions
- The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

Calling Conventions at hir Level

You must:

- Preserve some registers (fp, sp)
- Allocate the frame
- Handle the static link (i0)
- Receive the (other) arguments (i1, i2...)

You don't:

- Save temporaries (havm has magic for recursion)
- Jump to the ra (this is not nice feature from havm)

havm Calling Conventions

```
let function gcd (a: int, b: int) : int = (...)
in print_int (gcd (42, 51)) end
# Routine: gcd
                                  # Body
label 10
                                  move temp rv
# Prologue
                                        eseq
move temp t2, temp fp
move temp fp, temp sp
                                         temp t0
                                  # Epilogue
move temp sp, temp sp - const 4
move mem temp fp, temp i0
                                  move temp sp, temp fp
move temp t0, temp i1
                                  move temp fp, temp t2
move temp t1, temp i2
                                  label end
                                  # Routine: Main Program
                                   label Main
                                   sxp call name print_int
                                            call name 10 temp fp
                                                 const 42 const 51
                                  label end
                                            ◆□ ▶ ◆□ ▶ ◆ ■ ● り へ ○
```

Clever Translations

- Intermediate Representations
- 2 Memory Management
- 3 Translation to Intermediate Language
 - Calling Conventions
 - Clever Translations
 - Complex Expressions
- The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

What is the right translation for $\alpha < \beta$, with α and β two arbitrary expressions?

```
cjump (\alpha < \beta, ltrue, lfalse)

eseq (seq (cjump (\alpha < \beta, ltrue, lfalse), label ltrue move temp t, const 1 jump lend label lfalse move temp t, const 0 label lend), temp t)
```

What is the right translation for $\alpha < \beta$, with α and β two arbitrary expressions?

```
① cjump (\alpha < \beta, ltrue, lfalse)
```

```
eseq (seq (cjump (\alpha < \beta, ltrue, lfalse), label ltrue move temp t, const 1 jump lend label lfalse move temp t, const 0 label lend), temp t)
```

```
3 seq (sxp (\alpha) sxp (\beta)
```

What is the right translation for $\alpha < \beta$, with α and β two arbitrary expressions?

```
① cjump (\alpha < \beta, ltrue, lfalse)
```

```
eseq (seq (cjump (\alpha < \beta, ltrue, lfalse), label ltrue move temp t, const 1 jump lend label lfalse move temp t, const 0 label lend), temp t)
```

```
3 seq (sxp (\alpha) sxp (\beta)
```

 \bigcirc cjump ($\alpha < \beta$, ltrue, lfalse)

What is the right translation for $\alpha < \beta$, with α and β two arbitrary expressions?

```
eseq (seq (cjump (\alpha < \beta, ltrue, lfalse), label ltrue move temp t, const 1 jump lend label lfalse move temp t, const 0 label lend),
```

```
3 seq (sxp (\alpha) sxp (\beta))
```

temp t)

What is the right translation for $\alpha < \beta$, with α and β two arbitrary expressions?

```
oldsymbol{0} cjump (lpha < eta, ltrue, lfalse)
```

```
eseq (seq (cjump (\alpha < \beta, ltrue, lfalse), label ltrue move temp t, const 1 jump lend label lfalse move temp t, const 0 label lend), temp t)
```

```
3 seq (sxp (\alpha) sxp (\beta))
```

```
1 if \alpha < \beta then ...
```

What is the right translation for $\alpha < \beta$, with α and β two arbitrary expressions?

```
1 cjump (\alpha < \beta, ltrue, lfalse)
```

```
eseq (seq (cjump (\alpha < \beta, ltrue, lfalse), label ltrue move temp t, const 1 jump lend label lfalse move temp t, const 0 label lend), temp t)
```

```
3 seq (sxp (\alpha) sxp (\beta))
```

```
① if \alpha < \beta then ...
```

- $\mathbf{a} := \alpha < \beta$
- \bigcirc ($\alpha < \beta$, ())

What is the right translation for $\alpha < \beta$, with α and β two arbitrary expressions?

```
1 cjump (\alpha < \beta, ltrue, lfalse)
```

```
2 eseq (seq (cjump (\alpha < \beta, ltrue, lfalse), label ltrue move temp t, const 1 jump lend label lfalse move temp t, const 0 label lend), temp t)
```

```
3 seq (sxp (\alpha) sxp (\beta))
```

```
① if \alpha < \beta then ...
```

$$\textbf{2} \text{ a := } \alpha < \beta$$



What is the right translation for $\alpha < \beta$, with α and β two arbitrary expressions?

```
1 cjump (\alpha < \beta, ltrue, lfalse)
```

```
2 eseq (seq (cjump (\alpha < \beta, ltrue, lfalse), label ltrue move temp t, const 1 jump lend label lfalse move temp t, const 0 label lend), temp t)
```

```
3 seq (sxp (\alpha) sxp (\beta))
```

- ① if $\alpha < \beta$ then ...
- $oldsymbol{2}$ a := $\alpha < \beta$
- \odot ($\alpha < \beta$, ()).

- The right translation depends upon the use.
 This is context sensitive!
- How to implement this?

- Don't forget to preserve the demands of higher levels...
- Eek.

- The right translation depends upon the *use*. This is context sensitive!
- How to implement this?
 - When entering an IfExp, warn "I want a condition";
 - then, depending whether it is an expression or a statement, warn "I want an expression" or "I want a statement".
- Don't forget to preserve the demands of higher levels...
- Eek.

- The right translation depends upon the *use*. This is context sensitive!
- How to implement this?
 - When entering an IfExp, warn "I want a condition",
 - then, depending whether it is an expression or a statement, warn "I want an expression" or "I want a statement".
- Don't forget to preserve the demands of higher levels...
- Eek.

- The right translation depends upon the *use*. This is context sensitive!
- How to implement this?
 - When entering an IfExp, warn "I want a condition",
 - then, depending whether it is an expression or a statement, warn "I want an expression" or "I want a statement".
- Don't forget to preserve the demands of higher levels...
- Eek.

- The right translation depends upon the *use*. This is context sensitive!
- How to implement this?
 - When entering an IfExp, warn "I want a condition",
 - then, depending whether it is an expression or a statement, warn "I want an expression" or "I want a statement".
- Don't forget to preserve the demands of higher levels...
- Eek.

- The right translation depends upon the *use*. This is context sensitive!
- How to implement this?
 - When entering an IfExp, warn "I want a condition",
 - then, depending whether it is an expression or a statement, warn "I want an expression" or "I want a statement".
- Don't forget to preserve the demands of higher levels...
- Eek.

Rather, delay the translation until the use is known (translate::Exp):

- Ex Expression shell, encapsulation of a proto value,
- Nx Statement shell, encapsulating a wannabe statement,
- Cx Condition shell, encapsulating a wannabe condition.

Rather, delay the translation until the use is known (translate::Exp):

- Ex Expression shell, encapsulation of a proto value,
- Nx Statement shell, encapsulating a wannabe statement,
- Cx Condition shell, encapsulating a wannabe condition.

Rather, delay the translation until the use is known (translate::Exp):

- Ex Expression shell, encapsulation of a proto value,
- Nx Statement shell, encapsulating a wannabe statement,
- Cx Condition shell, encapsulating a wannabe condition.

Exp	un_nx	un_ex	$un_cx(t, f)$
Ex(e)	sxp(e)	е	
Cx(a < b)			
Nx(s)			

Rather, delay the translation until the use is known (translate::Exp):

- Ex Expression shell, encapsulation of a proto value,
- Nx Statement shell, encapsulating a wannabe statement,
- Cx Condition shell, encapsulating a wannabe condition.

Exp	un_nx	un_ex	$un_cx(t, f)$
Ex(e)	sxp(e)	e	cjump(e \neq 0, t, f)
Cx(a < b)			
N×(s)			

Rather, delay the translation until the use is known (translate::Exp):

- Ex Expression shell, encapsulation of a proto value,
- Nx Statement shell, encapsulating a wannabe statement,
- Cx Condition shell, encapsulating a wannabe condition.

```
\begin{array}{c|ccccc} Exp & un\_nx & un\_ex & un\_cx \ (t, f) \\ \hline Ex(e) & sxp(e) & e & cjump(e \neq 0, t, f) \\ \hline Cx(a < b) & seq(sxp(a), sxp(b)) \\ \hline Nx(s) & & & & & & & & & \\ \hline \end{array}
```

Rather, delay the translation until the use is known (translate::Exp):

- Ex Expression shell, encapsulation of a proto value,
- Nx Statement shell, encapsulating a wannabe statement,
- Cx Condition shell, encapsulating a wannabe condition.

$$\begin{array}{c|ccccc} Exp & un_nx & un_ex & un_cx \ (t, f) \\ \hline Ex(e) & sxp(e) & e & cjump(e \neq 0, t, f) \\ Cx(a < b) & seq(sxp(a), sxp(b)) & eseq(t \leftarrow (a < b), t) \\ Nx(s) & & \end{array}$$

Rather, delay the translation until the use is known (translate::Exp):

- Ex Expression shell, encapsulation of a proto value,
- Nx Statement shell, encapsulating a wannabe statement,
- Cx Condition shell, encapsulating a wannabe condition.

Exp	un_nx	un_ex	$un_cx(t, f)$
Ex(e)	sxp(e)	е	cjump(e \neq 0, t, f)
Cx(a < b)	seq(sxp(a), sxp(b))	$eseq(t \leftarrow (a < b), t)$	cjump(a < b, t, f)
N×(s)			

Rather, delay the translation until the use is known (translate::Exp):

- Ex Expression shell, encapsulation of a proto value,
- Nx Statement shell, encapsulating a wannabe statement,
- Cx Condition shell, encapsulating a wannabe condition.

Exp	un_nx	un_ex	$un_cx(t, f)$
Ex(e)	sxp(e)	е	cjump(e \neq 0, t, f)
Cx(a < b)	seq(sxp(a), sxp(b))	$eseq(t \leftarrow (a < b), t)$	cjump(a < b, t, f)
N×(s)	S		

Rather, delay the translation until the use is known (translate::Exp):

- Ex Expression shell, encapsulation of a proto value,
- Nx Statement shell, encapsulating a wannabe statement,
- Cx Condition shell, encapsulating a wannabe condition.

Exp	un_nx	un_ex	$un_cx(t, f)$
Ex(e)	sxp(e)	e	cjump(e \neq 0, t, f)
Cx(a < b)	seq(sxp(a), sxp(b))	$eseq(t \leftarrow (a < b), t)$	cjump(a < b, t, f)
N×(s)	S	???	

Rather, delay the translation until the use is known (translate::Exp):

- Ex Expression shell, encapsulation of a proto value,
- Nx Statement shell, encapsulating a wannabe statement,
- Cx Condition shell, encapsulating a wannabe condition.

Exp	un_nx	un_ex	$un_cx(t, f)$
Ex(e)	sxp(e)	е	cjump(e \neq 0, t, f)
Cx(a < b)	seq(sxp(a), sxp(b))	$eseq(t \leftarrow (a < b), t)$	cjump(a < b, t, f)
N×(s)	S	???	???

if 11 < 22 | 22 < 33 then print_int(1) else print_int(0)

```
cjump ne
   eseq seq cjump 11 < 22 name 10 name 11
               label 10 move temp t0 const 1
                         jump name 12
               label 11 move temp t0
                         eseq seq move temp t1 const 1
                                  cjump 22 < 33 name 13 name 14
                                  label 14
                                  move temp t1 const 0
                                  label 13
                              seq end
                              temp t1
                        jump name 12
               label 12
         seq end
        temp t0
   const 0
   name 15
   name 16
label 15 sxp call name print_int const 1
          jump name 17
label 16
          sxp call name print_int const 0
          jump name 17
label 17
```

A Better Translation: Ix

```
seq
    cjump 11 < 22 name 13 name 14
    label 13
      cjump 1 <> 0 name 10 name 11
    label 14
      cjump 22 < 33 name 10 name 11
  seq end
label 10
  sxp call name print_int const 1
  jump name 12
label 11
  sxp call name print_int const 0
label 12
```

Complex Expressions

- Intermediate Representations
- 2 Memory Management
- 3 Translation to Intermediate Language
 - Calling Conventions
 - Clever Translations
 - Complex Expressions
- The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

Complex Expressions

- Array creation
- Record creation
- String comparison
- While loops
- For loops

While Loops

while condition do body

While Loops

while condition do body

```
test:
   if not (condition)
     goto done
   body
   goto test
done:
```

For Loops

```
for i := min to max
  do body
```

```
let i := min
    limit := max
in
  while i <= limit
  do
    (body; ++i)
end</pre>
```

For Loops

```
for i := min to max
do body
```

```
let i := min
    limit := max
in
  if (i > limit)
   goto end
loop:
    body
    if (i >= limit)
     goto end
    ++i
    goto loop
end:
```

Additional Features

- Bounds checking
- Nil checking
- ..

The Case of the Tiger Compiler

- Intermediate Representations
- 2 Memory Management
- 3 Translation to Intermediate Language
- 4 The Case of the Tiger Compiler
 - Translation in the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

Translation in the Tiger Compiler

- Intermediate Representations
- 2 Memory Management
- Translation to Intermediate Language
- 4 The Case of the Tiger Compiler
 - Translation in the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

Actors: The temp Module

temp::Temp temporaries are pseudo-registers.
Generation of fresh temporaries.

temp::Label Pseudo addresses, both for data and code.

Generation of fresh labels.

misc::endo_map<T> Mapping from T to T.

Used during register allocation.

Actors: The tree Module

Implementation of hir and lir.

```
/Tree/ /Exp/ Const (int)
               Name
                      (temp::Label)
               Temp (temp::Temp)
               Binop (Oper, Exp, Exp)
               Mem (Exp)
               Call (Exp, list<Exp*>)
               Eseq (Stm, Exp)
        /Stm/
               Move
                     (Exp, Exp)
               Sxp
                     (Exp)
                     (Exp, list<temp::Label>)
               Jump
               CJump (Relop, Exp, Exp, Label, Label)
               Seq
                      (list<Stm *>)
               Label (temp::Label)
```

Actors: The tree Module: Warnings

- temp::Temp is not tree::Temp.
 The latter aggregates one of the former.
 Similarly with Label.
- n-ary seq. (Unlike [Appel, 1998]).
- Sxp instead of Exp.

Actors: The frame Module

Access How to reach a "variable".

Abstract class with two concrete subclasses.

frame::In_Register
frame::In_Frame

Frame What "variables" a frame contains.

local_alloc(bool escapes_p) -> Access

Frames and (frame::) accesses are not aware of static links.

Actors: The translate Module

Access Static link aware version of frame::Access:
how to reach a variable, including non local: a frame::Access
and a translate::Level.

exp(Level use) -> Exp Tree expression
The location of this Access, from the use point of view.

Level Static link aware version of frame::Frame: what variables a frame contains, and where is its parent level.

fp(Level use) -> Exp Tree expression

The frame pointer of this Level, from the use point of view. Used for calls, and reaching frame resident temporaries.

Actors: The translate Module

translate::Exp

Prototranslation wrappers (Ex, Nx, Cx, Ix).

translate/translation.hh

Auxiliary functions used by the Translator.

translate::Translator

The translator.

lir: Low Level Intermediate Representation

- Intermediate Representations
- 2 Memory Management
- Translation to Intermediate Language
- The Case of the Tiger Compiler
- 5 lir: Low Level Intermediate Representation

- Structure
 No nested sequences.
- Expressions
 Assembly is imperative: there is no "expression"
- Calling Conventions
 A (high-level) call is a delicate operation, not a simple instruction.
- Machines provide "jump or continue" instructions.
- Limited Number of Registers
 From temps to actual registers

- Structure
 No nested sequences.
- Expressions
 Assembly is imperative: there is no "expression".
- Calling Conventions
 A (high-level) call is a delicate operation, not a simple instruction
- Two Way Conditional Jumps
 Machines provide "jump or continue" instructions.
- Limited Number of Registers
 From temps to actual registers.

- Structure
 No nested sequences.
- Expressions
 Assembly is imperative: there is no "expression".
- Calling Conventions
 A (high-level) call is a delicate operation, not a simple instruction.
- Two Way Conditional Jumps
 Machines provide "jump or continue" instructions
- Limited Number of Registers
 From temps to actual registers

- Structure
 No nested sequences.
- Expressions
 Assembly is imperative: there is no "expression".
- Calling Conventions
 A (high-level) call is a delicate operation, not a simple instruction.
- Two Way Conditional Jumps
 Machines provide "jump or continue" instructions.
- Limited Number of Registers
 From temps to actual registers

- Structure
 No nested sequences.
- Expressions
 Assembly is imperative: there is no "expression".
- Calling Conventions
 A (high-level) call is a delicate operation, not a simple instruction.
- Two Way Conditional Jumps
 Machines provide "jump or continue" instructions.
- Limited Number of Registers
 From temps to actual registers.

Linearization: Principle

- eseq and seq must be eliminated (except the outermost seq).
- Similar to cut-elimination: permute inner eseq and seq to lift them higher, until they vanish.
- A simple rewriting system.

```
eseq (s1, eseq (s2, e)) \rightarrow eseq (seq (s1, s2), e)

x (eseq (s, e)) \rightarrow seq (s, x (e))
```

seq (ss1, seq (ss2), ss3) \sim seq (ss1, ss2, ss3)

```
seq (ss1, seq (ss2), ss3) \sim seq (ss1, ss2, ss3) call (f, eseq (s, e), es)
```

```
seq (ss1, seq (ss2), ss3) \rightsquigarrow seq (ss1, ss2, ss3) call (f, eseq (s, e), es) \rightsquigarrow eseq (s, call (f, e, es))
```

```
seq (ss1, seq (ss2), ss3) \rightarrow seq (ss1, ss2, ss3) call (f, eseq (s, e), es) \rightarrow eseq (s, call (f, e, es)) binop (+, eseq (s, e1), e2)
```

```
seq (ss1, seq (ss2), ss3) \rightarrow seq (ss1, ss2, ss3)
call (f, eseq (s, e), es) \rightarrow eseq (s, call (f, e, es))
binop (+, eseq (s, e1), e2) \rightarrow eseq (s, binop (+, e1, e2))
```

```
seq (ss1, seq (ss2), ss3) \rightarrow seq (ss1, ss2, ss3)
call (f, eseq (s, e), es) \rightarrow eseq (s, call (f, e, es))
binop (+, eseq (s, e1), e2) \rightarrow eseq (s, binop (+, e1, e2))
binop (+, e1, eseq (s, e2))
```

```
seq (ss1, seq (ss2), ss3) \rightarrow seq (ss1, ss2, ss3)
call (f, eseq (s, e), es) \rightarrow eseq (s, call (f, e, es))
binop (+, eseq (s, e1), e2) \rightarrow eseq (s, binop (+, e1, e2))
binop (+, e1, eseq (s, e2)) \rightarrow eseq (s, binop (+, e1, e2))
```

- This transformation is invalid: it changes the semantics.
- How can it be solved?

```
binop (+, e1, eseq (s, e2)) \rightarrow eseq (s, binop (+, e1, e2))
```

• But what if s modifies the value of e1?

- This transformation is invalid: it changes the semantics.
- How can it be solved?

```
binop (+, e1, eseq (s, e2)) \rightarrow eseq (s, binop (+, e1, e2))
```

• But what if s modifies the value of e1?

- This transformation is invalid: it changes the semantics.
- How can it be solved?

Wrong

Right

Linearization: More Temporaries

When "de-expressioning" fresh temporaries are needed

More generally

This is extremely inefficient when not needed...

Linearization: Commutativity

- Save useless extra temporaries and moves.
- Problem: commutativity cannot be known statically.
 E.g., move (mem (t1), e) and mem (t2)
 commute iff t1 \(\neq \) t2.
- We need a conservative approximation,
 i.e., never say "commute" when they don't.
 E.g., "if e is a const then s and e definitely commute".

Linearization: Commutativity

- Save useless extra temporaries and moves.
- Problem: commutativity cannot be known statically.
 E.g., move (mem (t1), e) and mem (t2)
 commute iff t1 ≠ t2.
- We need a conservative approximation,
 i.e., never say "commute" when they don't.
 E.g., "if e is a const then s and e definitely commute".

Linearization: Commutativity

- Save useless extra temporaries and moves.
- Problem: commutativity cannot be known statically.
 E.g., move (mem (t1), e) and mem (t2)
 commute iff t1 ≠ t2.
- We need a conservative approximation,
 i.e., never say "commute" when they don't.
 E.g., "if e is a const then s and e definitely commute".

Call Normalization

Normalization of a call depends on the kind of the routine:

procedure then its parent must be an sxp
function then its parent must be a move (temp t, .)

This normalization is performed simultaneously with linearization.

Two Way Jumps

Obviously, to enable the translation of a cjump into actual assembly instructions, the "false" label must follow the cjump. How?

Two Way Jumps: Basic Blocks

Split the long outer seq into "basic blocks":

- a single entry: the first instruction
- a single (maybe multi-) exit: the last instruction

It may require

- a new label as first instruction, to which the prologue jumps
- new labels after jumps or cjumps
- a new jump from the last instruction to the epilogue.

Two Way Jumps: Traces

Start from the initial block, and "sew" each remaining basic block to this growing "trace".

- If the last instruction is a jump
 - if the "destination block" is available, add it
 - otherwise, fetch any other remaining block.
- If the last instruction is a cjump
 - If the false destination is available, push it
 - If the true destination is available, flip the cjump and push it,
 - otherwise, change the cjump to go to a fresh label, attach this label, and finally jump to the initial false destination.

Two Way Jumps: Optimizing Traces

Many jumps should be removable, but sometimes there are choices to make.

```
label prologue

Prologue.
jump name test
```

```
label test
cjump i <= N, body, done</pre>
```

```
label body
Body.
jump name test
```

```
label done
Epilogue
jump name end
```

Two Way Jumps: Optimizing Traces

label prologue
Prologue
jump name test

label prologue
Prologue
jump name test

label prologue
Prologue
jump name test

label test
cjump i > N,
 done, body

label test
cjump i <= N,
body, done

label body *Body* jump name test

label body
Body
jump name test

label done
Epilogue
jump name end

label test
cjump i <= N,
body, done</pre>

label done
Epilogue
jump name end

label body
Body
jump name test

Epilogue
jump name end

Two Way Jumps: Optimizing Traces

label prologue
 Prologue
jump name test

label test
cjump i > N,
 done, body

label body
Body
jump name test

label done

Epilogue
jump name end

label prologue
Prologue
jump name test

label test
cjump i <= N,
 body, done</pre>

abel done
Epilogue
ump name end

abel body

Body

ump name test

label prologue
Prologue
jump name test

label body *Body* jump name test

cjump i <= N,
body, done</pre>

label done

Epilogue
jump name er

<ロ > ← □ > ← □ > ← □ > ← □ = ← の へ ○

Two Way Jumps: Optimizing Traces

label prologue
 Prologue
jump name test

label prologue
 Prologue
jump name test

Prologue jump name

label test
cjump i > N,
 done, body

label test
cjump i <= N,
 body, done</pre>

label body
Body
jump name test

label done
Epilogue
jump name end

label done
Epilogue
jump name end

label body
Body
jump name test

label done

Epilogue
jump name e

<ロ > ← □ > ← □ > ← □ > ← □ = ← の へ ○

Two Way Jumps: Optimizing Traces

label prologue

Prologue
jump name test

label prologue
 Prologue
jump name test

label prologue
Prologue
jump name test

label test
cjump i > N,
 done, body

label test
cjump i <= N,
 body, done</pre>

label body
Body
jump name test

label body
Body
jump name test

label done

Epilogue

jump name end

label test
cjump i <= N,
 body, done</pre>

label done
Epilogue
jump name end

label body
Body
jump name test

label done

Epilogue
jump name end

<ロ>

Bibliography I

- Appel, A. W. (1998).

 Modern Compiler Implementation in C, Java, ML.

 Cambridge University Press.
- Edwards, S. (2003).

COMS W4115 Programming Languages and Translators.

http://www.cs.columbia.edu/~sedwards/classes/2003/w4115/.

larus, J. R. (1990).

SPIM S20: A MIPS R2000 simulator.

Technical Report TR966, Computer Sciences Department, University of Wisconsin–Madison.

Bibliography II

Loosemore, S., Stallman, R. M., McGrath, R., Oram, A., and Drepper, U. (2003).

The GNU C Library Reference Manual.

Free Software Foundation, 59 Temple Place – Suite 330, Boston, MA 02111-1307 USA, 0.10 edition.

Muchnick, S. (1997).

Advanced Compiler Design and Implementation.

Morgan Kaufmann Publishers.

Instruction Selection

Akim Demaille Étienne Renault Roland Levillain first.last@lrde.epita.fr

EPITA — École Pour l'Informatique et les Techniques Avancées

April 23, 2018

Instruction Selection

- Microprocessors
- 2 A Typical risc: mips
- 3 The EPITA Tiger Compiler
- 4 Instruction Selection
- 5 Instruction Selection

Microprocessors

- Microprocessors
- 2 A Typical risc: mips
- 3 The EPITA Tiger Compiler
- 4 Instruction Selection
- 5 Instruction Selection

Instruction set architecture is the structure of a computer that a machine language programmer (or a compiler) must understand to write a correct (timing independent) program for that machine IBM introducing 360 (1964)

The Instruction Set Architecture (ISA) is the part of the processor that is visible to the programmer or compiler writer.

What is an instruction set?

An instruction set specifies a processor functionality:

- what operations are supported
- what storage mechanisms are used
- how to access storage
- how to communicate program to processor

Technical aspect of instruction set

- o format: length, encoding
- operations: data type (floating or fixed point), number and kind of operands
- 3 storage:
 - internal: accumulator, stack, register
 - memory: address size, addressing modes
- control: branch condition, support for procedures, predication

What makes a good instruction set?

An instruction set specifies a processor functionnality:

- implementability: support for a (high performances) range of implementation
- programmability: easy to express program (by Humans before 80's, mostly by compiler nowadays)
- backward/forward compatibility: implementability & programmability across generation

cisc - Complex Instruction Set Chip

- large number of instructions (100-250)
- 6, 8, 16 registers, some for pointers, others for integer computation
- arithmetic in memory can be processed
- two address code
- many possible effects (e.g., self-incrementation)

cisc - Pros & Cons

Pros:

- Simplified compiler: translation from IR is straightforward
- Smaller assembly code than risc assembly code
- Fewer instructions will be fetched
- Special purpose register available: stack pointer, interrupt handling ...

Cons:

- Variable length instruction format
- Many instruction require many clock for execution
- Limiter number of general purpose register
- (often) new version of cisc include the subset of instructions of the previous version

Motivations for something else!

Though the CISC programs could be small in length, but number of bits of memory occupies may not be less

The complex instructions do not simplify the compilers: many clock cycles can be wasted to find the appropriate instruction.

risc architectures were designed with the goal of executing one instruction per clock cycle.

risc - Reduced Instruction Set Chip

- 32 generic purpose registers
- arithmetic only available on registers
- 3 address code
- load and store relative to a register (M[r + const])
- only one effect or result per instruction

risc – Pipeline 1/3

Pipelining is the overlapping the execution of several instructions in a pipeline fashion.

A pipeline is (typically) decomposed into five stages:

- Instruction Fetch (IF)
 - 2 Instruction Decode (ID)
 - Execute (EX)
 - Memory Access (MA)
 - Write Back (WB)

risc – Pipeline 2/3

```
inst1:
      IF
          ID
              EX
                   MA
                        WB
          IF
inst2:
               ID
                   EX
                        MA
                              WB
               IF
inst3:
                    ID EX
                              MA
                                   WB
inst4:
                    IF
                         ID EX
                                   MA
                                         WB
                         IF
                               ID
                                    ΕX
inst5:
                                         MA
                                              WB
```

The slowest stage determines the speed of the whole pipeline!

Ex introduces latency

- Register-Register Operation: 1 cycle
- Memory Reference: 2 cycles
- Multi-cycle Instructions (floating point): many cycles

risc – Pipeline 3/3

Data hazard: When an instruction depends on the results of a previous instruction still in the pipeline.

- inst1 write in \$s1 during WB
- inst1 read in \$s1 during ID

```
inst1: IF ID EX MA WB inst2: IF ID EX MA WB
```

inst2 must be split, causing delays...

other dependencies can appears

risc – Pros & Cons

Pros:

- Fixed length instructions: decoding is easier
- Simpler hardware: higher clock rate
- Efficient Instruction pipeline
- Large number of general purpose register
- Overlapped register windows to speed up procedure call and return
- One instruction per cycle

Cons:

- Minimal number of addressing modes: only Load and Store
- Relatively few instructions

Nowadays

- the classification pure-risc or pure-cisc is becoming more and more inappropriate and may be irrelevant
- modern processors use a calculated combination elements of both design styles
- decisive factor is based on a tradeoff between the required improvement in performance and the expected added cost
- Some processors that are classified as CISC while employing a number of RISC features, such as pipelining

ARM provides the advantage of using a CISC (in terms of functionality) and the advantage of an RISC (in terms of reduced code lengths).

Lessons to be learned

Implementability

Driven by technology: microcode, VLSI, FPGA, pipelining, superscalar, SIMD, SSE

Programmability

Driven by compiler technology

Sum-up

- Many non technical issues influence ISA's
- Best solutions don't always win (Intel X86)

Intel X86 (IA32)

- Introduced in 1978
- 8 × 32 bits "general" register
- variable length instructions (1–15 byte)
- long life to the king! 15 generations from Intel 8086 to Intel Kabylake

Intel's trick?

Decoder translates cisc into risc micro-operations

A Typical risc: mips

- Microprocessors
- 2 A Typical risc: mips
 - Integer Arithmetics
 - Logical Operations
 - Control Flow
 - Loads and Stores
 - Floating Point Operations
- 3 The EPITA Tiger Compiler
- Instruction Selection
- 5 Instruction Selection

mips Registers and Use Convention [Larus, 1990]

Name	Number	Usage
zero	0	Constant 0
at	1	Reserved for assembler
v0-v1	2–3	Expression evaluation and results of a function
a0-a3	4–7	Function argument 1–4
t0-t7	8–15	Temporary (not preserved across call)
s0-s7	16–23	Saved temporary (preserved across call)
t8-t9	24–25	Temporary (not preserved across call)
k0-k1	26–27	Reserved for OS kernel
gp	28	Pointer to global area
sp	29	Stack pointer
fp	30	Frame pointer
ra	31	Return address (used by function call)

Typical risc Instructions

The following slides are based on [Larus, 1990].

- The assembler translates pseudo-instructions (marked with † below).
- In all instructions below, Src2 can be
 - a register
 - an immediate value (a 16 bit integer).
- The immediate forms are included for reference.
- The assembler translates the general form (e.g., add) into the immediate form (e.g., addi) if the second argument is constant.

Integer Arithmetics

- Microprocessors
- 2 A Typical risc: mips
 - Integer Arithmetics
 - Logical Operations
 - Control Flow
 - Loads and Stores
 - Floating Point Operations
- The EPITA Tiger Compiler
- 4 Instruction Selection
- 5 Instruction Selection

Arithmetic: Addition/Subtraction

```
addi Rdest, Rsrc1, Imm
Addition Immediate (with overflow)

addu Rdest, Rsrc1, Src2
Addition (without overflow)

addiu Rdest, Rsrc1, Imm
Addition Immediate (without overflow)

Put the sum of the integers from Rsrc1 and Src2 (or Imm) into Rdest.

sub Rdest, Rsrc1, Src2
Subtract (with overflow)

subu Rdest, Rsrc1, Src2
Subtract (without overflow)
```

Put the difference of the integers from Rsrc1 and Src2 into Rdest.

Addition (with overflow)

add Rdest, Rsrc1, Src2

Arithmetic: Division

If an operand is negative, the remainder is unspecified by the mips architecture and depends on the conventions of the machine on which spim is run.

```
div Rsrc1, Rsrc2
divu Rsrc1, Rsrc2
```

Divide (signed)
Divide (unsigned)

Divide the contents of the two registers. Leave the quotient in register 10 and the

remainder in register hi.

Divide (signed, with overflow) †

div Rdest, Rsrc1, Src2
divu Rdest, Rsrc1, Src2

Divide (unsigned, without overflow) †

Put the quotient of the integers from Rsrc1 and Src2 into Rdest.

rem Rdest, Rsrc1, Src2

Remainder †

remu Rdest, Rsrc1, Src2

Unsigned Remainder †

Likewise for the the remainder of the division.

Arithmetic: Multiplication

```
mul Rdest, Rsrc1, Src2

mulo Rdest, Rsrc1, Src2

mulou Rdest, Rsrc1, Src2

Multiply (with overflow) †

mulou Rdest, Rsrc1, Src2

Unsigned Multiply (with overflow) †

Put the product of the integers from Rsrc1 and Src2 into Rdest.

mult Rsrc1, Rsrc2

Multiply

multu Rsrc1, Rsrc2

Unsigned Multiply

Multiply the contents of the two registers. Leave the low-order word of the product in register 10 and the high-word in register hi.
```

Arithmetic Instructions

abs Rdest, Rsrc

Absolute Value †
Put the absolute value of the integer from Rsrc in Rdest.

Negate Value (with overflow) †

 neg Rdest, Rsrc
 Negate Value (with overflow) †

 negu Rdest, Rsrc
 Negate Value (without overflow) †

Put the negative of the integer from Rsrc into Rdest.

Logical Operations

- Microprocessors
- 2 A Typical risc: mips
 - Integer Arithmetics
 - Logical Operations
 - Control Flow
 - Loads and Stores
 - Floating Point Operations
- The EPITA Tiger Compiler
- 4 Instruction Selection
- Instruction Selection

Logical Instructions

and Rdest, Rsrc1, Src2
andi Rdest, Rsrc1, Imm

AND

AND Immediate

Put the logical AND of the integers from Rsrc1 and Src2 (or Imm) into Rdest.

not Rdest, Rsrc

NOT †

Put the bitwise logical negation of the integer from Rsrc into Rdest.

Logical Instructions

```
nor Rdest, Rsrc1, Src2
```

NOR

Put the logical NOR of the integers from Rsrc1 and Src2 into Rdest.

or Rdest, Rsrc1, Src2

OR

ori Rdest, Rsrc1, Imm

OR Immediate

Put the logical OR of the integers from Rsrc1 and Src2 (or Imm) into Rdest.

xor Rdest, Rsrc1, Src2

XOR

xori Rdest, Rsrc1, Imm

XOR Immediate

Put the logical XOR of the integers from Rsrc1 and Src2 (or Imm) into Rdest.

Logical Instructions

```
rol Rdest, Rsrc1, Src2
                                                                  Rotate Left †
                                                                 Rotate Right †
ror Rdest, Rsrc1, Src2
 Rotate the contents of Rsrc1 left (right) by the distance indicated by Src2 and
 put the result in Rdest.
                                                              Shift Left Logical
sll Rdest, Rsrc1, Src2
                                                      Shift Left Logical Variable
sllv Rdest, Rsrc1, Rsrc2
                                                         Shift Right Arithmetic
sra Rdest, Rsrc1, Src2
srav Rdest, Rsrc1, Rsrc2
                                                 Shift Right Arithmetic Variable
                                                             Shift Right Logical
srl Rdest, Rsrc1, Src2
                                                    Shift Right Logical Variable
srlv Rdest, Rsrc1, Rsrc2
 Shift the contents of Rsrc1 left (right) by the distance indicated by Src2
 (Rsrc2) and put the result in Rdest.
```

Control Flow

- Microprocessors
- 2 A Typical risc: mips
 - Integer Arithmetics
 - Logical Operations
 - Control Flow
 - Loads and Stores
 - Floating Point Operations
- The EPITA Tiger Compiler
- 4 Instruction Selection
- Instruction Selection

Comparison Instructions

seq Rdest, Rsrc1, Src2
Set Rdest to 1 if Rsrc1 equals Src2, otherwise to 0.

Set Equal †

sne Rdest, Rsrc1, Src2

Set Not Equal †

Set Rdest to 1 if Rsrc1 is not equal to Src2, otherwise to 0.

Comparison Instructions

```
Set Greater Than Equal †
sge Rdest, Rsrc1, Src2
                                           Set Greater Than Equal Unsigned †
sgeu Rdest, Rsrc1, Src2
 Set Rdest to 1 if Rsrc1 \geq Src2, otherwise to 0.
                                                          Set Greater Than †
sgt Rdest, Rsrc1, Src2
                                                 Set Greater Than Unsigned †
sgtu Rdest, Rsrc1, Src2
 Set Rdest to 1 if Rsrc1 > Src2, otherwise to 0.
                                                       Set Less Than Equal †
sle Rdest, Rsrc1, Src2
sleu Rdest, Rsrc1, Src2
                                              Set Less Than Equal Unsigned †
 Set Rdest to 1 if Rsrc1 < Src2, otherwise to 0.
                                                               Set Less Than
slt Rdest, Rsrc1, Src2
                                                     Set Less Than Immediate
slti Rdest, Rsrc1, Imm
                                                      Set Less Than Unsigned
sltu Rdest, Rsrc1, Src2
                                           Set Less Than Unsigned Immediate
sltiu Rdest, Rsrc1, Imm
 Set Rdest to 1 if Rsrc1 < Src2 (or Imm), otherwise to 0.
```

Branch instructions use a signed 16-bit offset field: jump from -2^{15} to $+2^{15}-1$) instructions (not bytes). The jump instruction contains a 26 bit address field.

b label Branch instruction †

Unconditionally branch to label.

j label Jump

Unconditionally jump to label.

jal labelJump and Linkjalr RsrcJump and Link Register

Unconditionally jump to *label* or whose address is in Rsrc. Save the address of the next instruction in register 31.

jr Rsrc Jump Register

Unconditionally jump to the instruction whose address is in register Rsrc.

bczt label Branch Coprocessor z True
bczf label Branch Coprocessor z False
Conditionally branch to label if coprocessor z's condition flag is true (false).

Conditionally branch to *label* if the contents of Rsrc1 * Src2.

```
beq Rsrc1, Src2, label
bne Rsrc1, Src2, label
```

begz Rsrc, label

bnez Rsrc, label

Branch on Equal Branch on Not Equal

Branch on Equal Zero † Branch on Not Equal Zero †

Conditionally branch to *label* if the contents of Rsrc1 * Src2.

bge Rsrc1, Src2, label
Branch on Greater Than Equal †
bgeu Rsrc1, Src2, label
Branch on GTE Unsigned †
Branch on Greater Than Equal Zero
bgezal Rsrc, label
Branch on Greater Than Equal Zero And Link
Conditionally branch to label if the contents of Rsrc are greater than or equal to
0. Save the address of the next instruction in register 31.

bgt Rsrc1, Src2, label
bgtu Rsrc1, Src2, label
bgtz Rsrc, label

Branch on Greater Than † Branch on Greater Than Unsigned † Branch on Greater Than Zero

```
Conditionally branch to label if the contents of Rsrc1 are * to Src2.
                                                  Branch on Less Than Equal †
ble Rsrc1, Src2, label
                                                    Branch on LTE Unsigned †
bleu Rsrc1, Src2, label
                                               Branch on Less Than Equal Zero
blez Rsrc, label
                                  Branch on Greater Than Equal Zero And Link
bgezal Rsrc, label
bltzal Rsrc, label
                                                Branch on Less Than And Link
 Conditionally branch to label if the contents of Rsrc are greater or equal to 0 or
 less than 0, respectively. Save the address of the next instruction in register 31.
                                                        Branch on Less Than †
blt Rsrc1, Src2, label
                                               Branch on Less Than Unsigned †
bltu Rsrc1, Src2, label
bltz Rsrc, label
                                                     Branch on Less Than Zero
```

Exception and Trap Instructions

rfe

Return From Exception

Restore the Status register.

syscall

System Call

Register \$v0 contains the number of the system call provided by spim.

break n

Break

Cause exception n. Exception 1 is reserved for the debugger.

nop

No operation

Do nothing.

Loads and Stores

- Microprocessors
- 2 A Typical risc: mips
 - Integer Arithmetics
 - Logical Operations
 - Control Flow
 - Loads and Stores
 - Floating Point Operations
- The EPITA Tiger Compiler
- 4 Instruction Selection
- 5 Instruction Selection

Constant-Manipulating Instructions

li Rdest, imm

Load Immediate †

Move the immediate imm into Rdest.

lui Rdest, imm

Load Upper Immediate

Load the lower halfword of the immediate imm into the upper halfword of Rdest. The lower bits of the register are set to 0.

Load: Byte & Halfword

1b Rdest, address
Load Byte
1bu Rdest, address
Load Unsigned Byte
Load the byte at address into Rdest. The byte is sign-extended by the 1b, but
not the 1bu, instruction.

1h Rdest, address
Load Halfword
1hu Rdest, address
Load Unsigned Halfword
Load the 16-bit quantity (halfword) at address into register Rdest. The halfword

is sign-extended by the 1h, but not the **lhu**, instruction

Load: Word

```
Load Word
lw Rdest, address
 Load the 32-bit quantity (word) at address into Rdest.
lwcz Rdest. address
                                                        Load Word Coprocessor
 Load the word at address into Rdest of coprocessor z (0-3).
                                                                Load Word Left
lwl Rdest, address
lwr Rdest, address
                                                              Load Word Right
 Load the left (right) bytes from the word at the possibly-unaligned address into
 Rdest.
                                                     Unaligned Load Halfword †
ulh Rdest, address
                                           Unaligned Load Halfword Unsigned †
ulhu Rdest, address
 Load the 16-bit quantity (halfword) at the possibly-unaligned address into Rdest.
 The halfword is sign-extended by the ulh, but not the ulhu, instruction
                                                        Unaligned Load Word †
ulw Rdest, address
 Load the 32-bit quantity (word) at the possibly-unaligned address into Rdest.
```

Load Instructions

la Rdest, address Load Address †
Load computed address, not the contents of the location, into Rdest.

1d Rdest, address Load Double-Word † Load the 64-bit quantity at address into Rdest and Rdest + 1.

4ロト 4個ト 4 国 ト 4 国 ト 国 の 9 (で)

Store: Byte & Halfword

sb Rsrc, address
Store the low byte from Rsrc at address.

Store Byte

sh Rsrc, address
Store the low halfword from Rsrc at address.

Store Halfword

Store: Word

Store Word sw Rsrc, address Store the word from Rsrc at address.

Store Word Coprocessor swcz Rsrc, address

Store the word from Rsrc of coprocessor z at address.

Store Word Left swl Rsrc, address Store Word Right swr Rsrc, address

Store the left (right) bytes from Rsrc at the possibly-unaligned address.

Unaligned Store Halfword † ush Rsrc, address

Store the low halfword from Rsrc at the possibly-unaligned address.

Unaligned Store Word † usw Rsrc, address

Store the word from Rsrc at the possibly-unaligned address.

Store: Double Word

sd Rsrc, address Store Double-Word † Store the 64-bit quantity in Rsrc and Rsrc + 1 at address.

Data Movement Instructions

move Rdest, Rsrc

Move †

Move the contents of Rsrc to Rdest.

The multiply and divide unit produces its result in two additional registers, hi and lo (e.g., mul Rdest, Rsrc1, Src2).

mfhi Rdest

Move From hi

mflo Rdest

Move From lo

Move the contents of the hi (lo) register to Rdest.

mthi Rdest

Move To hi

mtlo Rdest

Move To lo

Move the contents Rdest to the hi (lo) register.

Data Movement Instructions

Coprocessors have their own register sets. These instructions move values between these registers and the CPU's registers.

mfc1.d Rdest, FRsrc1 Move Double From Coprocessor 1 † Move the contents of floating point registers FRsrc1 and FRsrc1 + 1 to CPU registers Rdest and Rdest + 1.

Floating Point Operations

- Microprocessors
- 2 A Typical risc: mips
 - Integer Arithmetics
 - Logical Operations
 - Control Flow
 - Loads and Stores
 - Floating Point Operations
- 3 The EPITA Tiger Compiler
- 4 Instruction Selection
- 5 Instruction Selection

mips Floating Point Instructions

- Floating point coprocessor 1 operates on single (32-bit) and double precision (64-bit) FP numbers.
- 32 32-bit registers \$f0-\$f31.
- Two FP registers to hold doubles.
- FP operations only use even-numbered registers including instructions that operate on single floats.
- Values are moved one word (32-bits) at a time by 1wc1, swc1, mtc1, and mfc1 or by the 1.s, 1.d, s.s, and s.d pseudo-instructions.
- The flag set by FP comparison operations is read by the CPU with its bc1t and bc1f instructions.

Floating Point: Arithmetics

Compute the \ast of the floating float doubles (singles) in FRsrc1 and FRsrc2 and put it in FRdest.

```
add.d FRdest, FRsrc1, FRsrc2
add.s FRdest, FRsrc1, FRsrc2
div.d FRdest, FRsrc1, FRsrc2
div.s FRdest, FRsrc1, FRsrc2
mul.d FRdest, FRsrc1, FRsrc2
mul.s FRdest, FRsrc1, FRsrc2
sub.d FRdest. FRsrc1. FRsrc2
sub.s FRdest, FRsrc1, FRsrc2
abs.d FRdest, FRsrc
abs.s FRdest, FRsrc
neg.d FRdest, FRsrc
neg.s FRdest, FRsrc
```

Floating Point Addition Double Floating Point Addition Single Floating Point Divide Double Floating Point Divide Single Floating Point Multiply Double Floating Point Multiply Single Floating Point Subtract Double Floating Point Subtract Single Floating Point Absolute Value Double Floating Point Absolute Value Single Negate Double Negate Single

Floating Point: Comparison

Compare the floating point double in FRsrc1 against the one in FRsrc2 and set the floating point condition flag true if they are *.

```
c.eq.d FRsrc1, FRsrc2
c.eq.s FRsrc1, FRsrc2
c.le.d FRsrc1, FRsrc2
c.le.s FRsrc1, FRsrc2
c.lt.d FRsrc1, FRsrc2
```

c.lt.s FRsrc1, FRsrc2

```
Compare Equal Double
Compare Less Than Equal Double
Compare Less Than Equal Single
```

Compare Less Than Double Compare Less Than Single

Floating Point: Conversions

Convert between (i) single, (ii) double precision floating point number or (iii) integer in FRsrc to FRdest.

```
cvt.d.s FRdest, FRsrc
cvt.d.w FRdest, FRsrc
cvt.s.d FRdest, FRsrc
cvt.s.w FRdest, FRsrc
cvt.w.d FRdest, FRsrc
```

cvt.w.s FRdest, FRsrc

```
Convert Single to Double
Convert Integer to Double
Convert Double to Single
Convert Integer to Single
```

Floating Point: Moves

mov.d FRdest, FRsrc

Move Floating Point Double

Move Floating Point Single

Move Floating Point Single

Move Floating Point Single

Move the floating float double (single) from FRsrc to FRdest.

s.d FRdest, address

Store Floating Point Double †

S.s FRdest, address

Store Floating Point Single †

Store the floating float double (single) in FRdest at address.

The EPITA Tiger Compiler

- Microprocessors
- 2 A Typical risc: mips
- 3 The EPITA Tiger Compiler
- 4 Instruction Selection
- 5 Instruction Selection

- mips is a nice assembly language
- mips is more modern
- mips is meaningful:

- spim is a portable mips emulator
- spim has a cool modern gui, xspim!

- mips is a nice assembly language
- mips is more modern
- mips is meaningful:

- spim is a portable mips emulator
- spim has a cool modern gui, xspim!

- mips is a nice assembly language
- mips is more modern
- mips is meaningful:
 - Million Instructions Per Second (10 mips, 1 mip)
 - Meaningless Indication of Processor Speed
 - Meaningless Information Provided by Salesmen
 - Meaningless Information per Second
 - Microprocessor without Interlocked Piped Stages
- spim is a portable mips emulator
- spim has a cool modern gui, xspim!

- mips is a nice assembly language
- mips is more modern
- mips is meaningful:
 - Million Instructions Per Second (10 mips, 1 mip)
 - Meaningless Indication of Processor Speed
 - Meaningless Information Provided by Salesmen
 - Meaningless Information per Second
 - Microprocessor without Interlocked Piped Stages
- spim is a portable mips emulator
- spim has a cool modern gui, xspim!

- mips is a nice assembly language
- mips is more modern
- mips is meaningful:
 - Million Instructions Per Second (10 mips, 1 mip)
 - Meaningless Indication of Processor Speed
 - Meaningless Information Provided by Salesmen
 - Meaningless Information per Second
 - Microprocessor without Interlocked Piped Stages
- spim is a portable mips emulator
- spim has a cool modern gui, xspim!

- mips is a nice assembly language
- mips is more modern
- mips is meaningful:
 - Million Instructions Per Second (10 mips, 1 mip)
 - Meaningless Indication of Processor Speed
 - Meaningless Information Provided by Salesmen
 - Meaningless Information per Second
 - Microprocessor without Interlocked Piped Stages
- spim is a portable mips emulator
- spim has a cool modern gui, xspim!

- mips is a nice assembly language
- mips is more modern
- mips is meaningful:
 - Million Instructions Per Second (10 mips, 1 mip)
 - Meaningless Indication of Processor Speed
 - Meaningless Information Provided by Salesmen
 - Meaningless Information per Second
 - Microprocessor without Interlocked Piped Stages
- spim is a portable mips emulator
- spim has a cool modern gui, xspim!

- mips is a nice assembly language
- mips is more modern
- mips is meaningful:
 - Million Instructions Per Second (10 mips, 1 mip)
 - Meaningless Indication of Processor Speed
 - Meaningless Information Provided by Salesmen
 - Meaningless Information per Second
 - Microprocessor without Interlocked Piped Stages
- spim is a portable mips emulator
- spim has a cool modern gui, xspim!

- mips is a nice assembly language
- mips is more modern
- mips is meaningful:
 - Million Instructions Per Second (10 mips, 1 mip)
 - Meaningless Indication of Processor Speed
 - Meaningless Information Provided by Salesmen
 - Meaningless Information per Second
 - Microprocessor without Interlocked Piped Stages
- spim is a portable mips emulator
- spim has a cool modern gui, xspim!

- mips is a nice assembly language
- mips is more modern
- mips is meaningful:
 - Million Instructions Per Second (10 mips, 1 mip)
 - Meaningless Indication of Processor Speed
 - Meaningless Information Provided by Salesmen
 - Meaningless Information per Second
 - Microprocessor without Interlocked Piped Stages
- spim is a portable mips emulator
- spim has a cool modern gui, xspim!

Cause = 0000000 BadVaddr = 00000000 Status= 00000000 HI General Registers = 000000000 R16 R16 (s0) = 0000000 R17 (s1) = 0000000 Register Display gp R13 (gp) Double Floating Point Registers 0.00000 Single Floating Point Registers quit load clear step set value Control **Buttons** help termina: mode **Text Segments** [0x00400000] 0x8fa40000 lw R4, 0(R29) [] [0x00400004] 0x27a50004 addiu R5, R29, 4 [] User and [0x00400008] 0x24a60004 addiu R6, R5, 4 [] Kernel [0x0040000c] 0x00041090 sll R2, R4, 2 [0x00400010] 0x00c23021 addu R6, R6, R2 Text 0x004000141 0x0c000000 ial 0x00000000 [1 Segments [0x00400018] 0x3402000a ori RO, RO, 10 [] 0x0040001cl 0x0000000c svscall **Data Segments** [0x100000001...[0x100100001 0x00000000 0x74706563 0x636f2000 0x6920646e Data and 0x7265746e Stack 0x100100401 Segments 0x100100601 SPIM Version 3.2 of January 14, 1990 SPIM Messages

A Sample: fact

```
/* Define a recursive function. */
let
  /* Calculate n! */
  function fact (n : int) : int =
    if n = 0
        then 1
        else n * fact (n - 1)
in
  print_int (fact (10));
  print ("\n")
end
```

```
# Routine: fact
                                        .data
10:
               $fp, -8 ($sp)
       SW
                                        14:
       move
               $fp, $sp
                                                .word 1
               $sp, $sp, 16
       sub
                                                .asciiz "\n"
               $ra, -12 ($fp)
       SW
                                        .text
               $a0, ($fp)
       SW
               $a1, -4 ($fp)
                                        # Routine: Main
       SW
               $t0, -4 ($fp)
15:
       lw
                                        t_main: sw
                                                         $fp, ($sp)
               $t0. 0. 11
       beq
                                                        $fp, $sp
                                                move
12:
       lw
               $a0, ($fp)
                                                sub
                                                         $sp, $sp, 8
       lw
               $t0, -4 ($fp)
                                                         $ra, -4 ($fp)
                                                SW
               $a1, $t0, 1
       sub
                                        17:
                                                        $a0, $fp
                                                move
       jal
               10
                                                li
                                                         $a1, 10
       lw
               $t0, -4 ($fp)
               $t0, $t0, $v0
       mul
                                                jal
                                                        10
13:
       move
              $v0, $t0
                                                        $a0, $v0
                                                move
               16
                                                jal
                                                        print_int
11:
       li
               $t0, 1
                                                la
                                                        $a0, 14
       j
               13
                                                jal
                                                        print
16:
       lw
               $ra, -12 ($fp)
                                       18:
                                                lw
                                                         $ra, -4 ($fp)
               $sp, $fp
       move
       lw
               $fp, -8 ($fp)
                                                move
                                                         $sp, $fp
       jr
               $ra
                                                         $fp, ($fp)
                                                lw
                                                         $ra
                                                jr
```

Nolimips (formerly Mipsy)

- Another mips emulator
- Interactive loop
- Unlimited number of \$x42 registers!

# Routine: fact 10: sw \$a0, (\$fp) # Routine: fact	(ten)
sw \$a1, -4 (\$fp) 10: sw \$fp, -8 move \$x11, \$s0	(Aph)
$\underset{\text{move}}{\text{move}}$ $\underset{\$x12}{\$x1}$, $\$s1$ $\underset{\text{move}}{\text{move}}$ $\$fp$, $\$s1$	p
move \$x13, \$s2 sub \$sp, \$s	
move \$x14. \$s3	
move \$x15, \$s4 sw \$ra, -1	2 (\$fp)
move \$x16, \$s5 sw \$a0, (\$:	fn)
move \$X17, \$Sb	_
move \$x18, \$s7 sw \$a1, -4	(\$fp)
15: lw \$x5, -4 (\$fp) 15: lw \$t0, -4	(\$fp)
beq \$x5, U, II	-
12: lw \$x6, (\$fp) beq \$t0, 0,	11
lw \$x8, -4 (\$fp) 12: lw \$a0, (\$:	fp)
sub \$x7, \$x8, 1 lw \$t0, -4	(¢fn)
move \$a1.\$x7	-
jal 10 sub \$a1, \$t	0, 1
move \$x3, \$v0 jal 10	
lw \$x10, -4 (\$ip)	(4.5.)
mul \$x9, \$x10, \$x3 lw \$t0, -4	(\$ip)
move \$x0, \$x9 mul \$t0, \$t	0, \$v0
13: move \$v0, \$x0	•
j 16 13: move \$v0, \$to	U
j 16	
16: move \$s0, \$x11	
move \$s1. \$x12	
move	
move \$s3, \$x14 l6: lw \$ra, -1	2 (\$fp)
move \$s4, \$x15	•
move \$s5, \$x16 move \$sp, \$f	р
move \$s6, \$x17 lw \$fp, -8	(\$fp)
move \$57, \$x10	
jr \$ra	

Instruction Selection

- Microprocessors
- 2 A Typical risc: mips
- 3 The EPITA Tiger Compiler
- 4 Instruction Selection
- 5 Instruction Selection

Nolimips (formerly Mipsy)

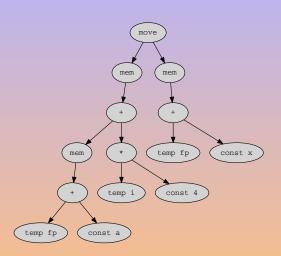
- Another mips emulator
- Interactive loop
- Unlimited number of \$x42 registers!

# Routine: fact 10: sw \$a0, (\$fp) # Routine: fact	(ten)
sw \$a1, -4 (\$fp) 10: sw \$fp, -8 move \$x11, \$s0	(Aph)
$\underset{\text{move}}{\text{move}}$ $\underset{\$x12}{\$x1}$, $\$s1$ $\underset{\text{move}}{\text{move}}$ $\$fp$, $\$s1$	р
move \$x13, \$s2 sub \$sp, \$s	
move \$x14. \$s3	
move \$x15, \$s4 sw \$ra, -1	2 (\$fp)
move \$x16, \$s5 sw \$a0, (\$:	fn)
move \$X17, \$Sb	_
move \$x18, \$s7 sw \$a1, -4	(\$fp)
15: lw \$x5, -4 (\$fp) 15: lw \$t0, -4	(\$fp)
beq \$x5, U, II	-
12: lw \$x6, (\$fp) beq \$t0, 0,	11
lw \$x8, -4 (\$fp) 12: lw \$a0, (\$:	fp)
sub \$x7, \$x8, 1 lw \$t0, -4	(¢fn)
move \$a1.\$x7	-
jal 10 sub \$a1, \$t	0, 1
move \$x3, \$v0 jal 10	
lw \$x10, -4 (\$ip)	(4.5.)
mul \$x9, \$x10, \$x3 lw \$t0, -4	(\$ip)
move \$x0, \$x9 mul \$t0, \$t	0, \$v0
13: move \$v0, \$x0	•
j 16 13: move \$v0, \$to	U
j 16	
16: move \$s0, \$x11	
move \$s1. \$x12	
move	
move \$s3, \$x14 l6: lw \$ra, -1	2 (\$fp)
move \$s4, \$x15	•
move \$s5, \$x16 move \$sp, \$f	р
move \$s6, \$x17 lw \$fp, -8	(\$fp)
move \$57, \$x10	
jr \$ra	

Instruction Selection

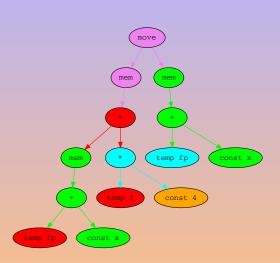
- Microprocessors
- 2 A Typical risc: mips
- 3 The EPITA Tiger Compiler
- 4 Instruction Selection
- 5 Instruction Selection

How would you translate
a[i] := x
where x is frame resident, and
i is not? [Appel, 1998]



Simple Instruction: Translation 1

```
load t17 <- M[fp + a]
addi t18 <- r0 + 4
mul t19 <- ti * t18
add t20 <- t17 + t19
load t21 <- M[fp + x]
store M[t20 + 0] <- t21
```



Tree Patterns

- Translation from Tree to Assembly corresponds to parsing a tree.
- Looking for a covering of the tree, using tiles
- The set of tiles corresponds to the instruction set









Tree Patterns

- Translation from Tree to Assembly corresponds to parsing a tree.
- Looking for a covering of the tree, using tiles.
- The set of tiles corresponds to the instruction set.









Tree Patterns

- Translation from Tree to Assembly corresponds to parsing a tree.
- Looking for a covering of the tree, using tiles.
- The set of tiles corresponds to the instruction set.



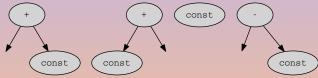




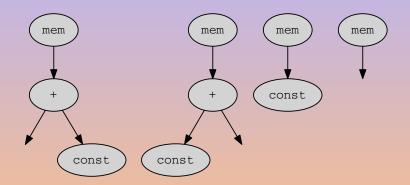


Tiles

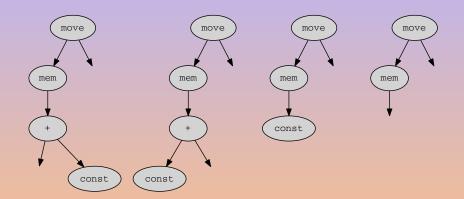
Missing nodes are plugs for *temporaries*: tiles read from temps, and create temps.



Some architectures rely on a special register to produce 0.

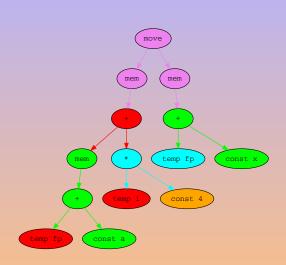


Tiles: Storing store $M[r_j + c] \leftarrow r_i$

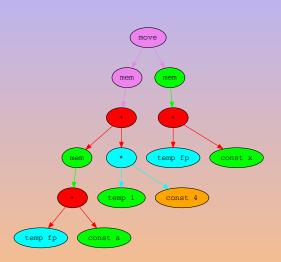


Simple Instruction: Translation 2

```
load t17 <- M[fp + a]
addi t18 <- r0 + 4
mul t19 <- ti * t18
add t20 <- t17 + t19
addi t21 <- fp + x
movem M[t20] <- M[t21]
```



Simple Instruction: Translation 3



- There is always a solution (provided the instruction set is reasonable)
- there can be several solutions
- given a cost function, some are better than others

- cpus are really layers of units that work in parallel
- Costs are therefore interrelated.

- There is always a solution (provided the instruction set is reasonable)
- there can be several solutions
- given a cost function, some are better than others

- cpus are really layers of units that work in parallel.
- Costs are therefore interrelated.

- There is always a solution (provided the instruction set is reasonable)
- there can be several solutions
- given a cost function, some are better than others:
 - some are locally better, optimal coverings (no fusion can reduce the cost).
 - some are globally better, optimum coverings.

- cpus are really layers of units that work in parallel
- Costs are therefore interrelated.

- There is always a solution (provided the instruction set is reasonable)
- there can be several solutions
- given a cost function, some are better than others:
 - some are locally better, *optimal coverings* (no fusion can reduce the cost),
 - some are globally better, optimum coverings.

- cpus are really layers of units that work in parallel.
- Costs are therefore interrelated.

- There is always a solution (provided the instruction set is reasonable)
- there can be several solutions
- given a cost function, some are better than others:
 - some are locally better, *optimal coverings* (no fusion can reduce the cost),
 - some are globally better, optimum coverings.

- cpus are really layers of units that work in parallel.
- Costs are therefore interrelated.

- There is always a solution (provided the instruction set is reasonable)
- there can be several solutions
- given a cost function, some are better than others:
 - some are locally better, *optimal coverings* (no fusion can reduce the cost),
 - some are globally better, optimum coverings.

- cpus are really layers of units that work in parallel.
- Costs are therefore interrelated.

Algorithms for Instruction Selection

Maximal Munch Find an optimal tiling.

- Top-down strategy.
- Cover the current node with the largest tile.
- Repeat on subtrees.
- Generate instructions in reverse-order after tile placement.

Dynamic Programming Find an optimum tiling.

- Bottom-up strategy.
- Assign cost to each node.
- Cost = cost of selected tile + cost of subtrees.
- Select a tile with minimal cost and recurse upward.
- Implemented by code generator generators (Twig, Burg, iBurg, MonoBURG, ...).

Algorithms for Instruction Selection

Maximal Munch Find an optimal tiling.

- Top-down strategy.
- Cover the current node with the largest tile.
- Repeat on subtrees.
- Generate instructions in reverse-order after tile placement.

Dynamic Programming Find an optimum tiling.

- Bottom-up strategy.
- Assign cost to each node.
- Cost = cost of selected tile + cost of subtrees.
- Select a tile with minimal cost and recurse upward.
- Implemented by code generator generators (Twig, Burg, iBurg, MonoBURG, ...).

Tree Matching

- The basic operation is the pattern matching.
- Not all the languages stand equal before pattern matching...

Tree Matching

- The basic operation is the pattern matching.
- Not all the languages stand equal before pattern matching. . .

... in Stratego

```
Select-swri :
  MOVE(MEM(BINOP(PLUS, e1, CONST(n))), e2) \rightarrow
  SEQ(MOVE(r2, e2), SEQ(MOVE(r1, e1), sw-ri(r2, r1, n)))
  where \langle \text{new-atemp} \rangle e1 \Rightarrow r1; \langle \text{new-atemp} \rangle e2 \Rightarrow r2
Select-swr:
  MOVE(MEM(e1), e2) \rightarrow SEQ(MOVE(r2, e2), SEQ(MOVE(r1, e1), sw-r(r2, r1)))
  where \langle \text{new-atemp} \rangle e1 \Rightarrow r1; \langle \text{new-atemp} \rangle e2 \Rightarrow r2
Select-nop :
  MOVE(TEMP(r), TEMP(r)) \rightarrow NUL
Select-nop :
  MOVE(REG(r), REG(r)) \rightarrow NUL
Select-mover :
  MOVE(TEMP(r), TEMP(t)) \rightarrow move(TEMP(r), TEMP(t)) where move(q) > (r, t)
Select-mover :
  MOVE(TEMP(r), REG(t)) \rightarrow move(TEMP(r), REG(t))
                                                                  where <not(eq)> (r, t)
Select-mover :
  MOVE(REG(r), TEMP(t)) \rightarrow move(REG(r), TEMP(t))
                                                                  where <not(eq)> (r, t)
Select-mover :
  MOVE(REG(r), REG(t)) \rightarrow move(REG(r), REG(t))
                                                                  where <not(eq)> (r, t)
```

... in Haskell: Ir.hs [Anisko, 2003]

```
module Ir (Stm (Move, Sxp, Jump, CJump, Seq, Label,
                LabelEnd, Literal),
           ...)
where
data Stm a =
    Move { ma :: a, lval :: Exp a, rval :: Exp a }
  | Sxp a (Exp a)
  | Jump a (Exp a)
  | CJump { cja :: a,
            rop :: Relop, cleft :: Exp a, cright :: Exp a,
            iftrue :: Exp a, iffalse :: Exp a }
  | Seq a [Stm a]
  | Label { la :: a,
            name :: String, size :: Int }
  | LabelEnd a
  | Literal { lita :: a,
              litname :: String, litcontent :: [Int] }
```

... in Haskell Eval.hs [Anisko, 2003]

```
module Eval (evalStm, ...)
where
import Ir
import Monad (Mnd, rfetch, rstore, rpush, rpop, ...)
import Result (Res (IntRes, UnitRes))
import Profile (profileExp, profileStm)
evalStm :: Stm Loc -> Mnd ()
evalStm stm@(Move loc (Temp _ t) e) =
    do (IntRes r) <- evalExp e</pre>
       verbose loc ["move", "(", "temp", t, ")", show r]
       profileStm stm
       rstore t r
evalStm stm@(Move loc (Mem _ e) f) =
    do (IntRes r) <- evalExp e</pre>
       (IntRes s) <- evalExp f
       verbose loc ["move", "(", "mem", show r, ")", show s]
       profileStm stm
       mstore r s
                                                   <ロ > < 回 > < 回 > < 亘 > く 亘 > ・ 亘 ・ り g (や
```

... in Haskell Low.hs [Anisko, 2003]

```
module Low (lowExp, lowStms)
where import ...
lowStms :: Int -> [Stm Ann] -> Mnd Bool
lowStms [] = return True
lowStms level
        ((CJump _ _ e f _ (Name _ s)) : (Label _ s' _) : stms)
        | s == s' =
    do a <- lowExp (level + 1) e</pre>
       b <- lowExp (level + 1) f
       c <- lowStms level stms
       return $ a && b && c
lowStms level (CJump l _ e f _ _ : stms) =
    do awarn l ["invalid cjump"]
       lowExp (level + 1) e
       lowExp (level + 1) f
       lowStms level stms
      return False
```

... in Haskell High.hs [Anisko, 2003]

```
module High (highExp, highStms)
where import ...
highStms :: Int -> [Stm Ann] -> Mnd Bool
highStms level ss =
    do a <- sequence $ map (highStm level) ss</pre>
       return (and a)
highStm :: Int -> Stm Ann -> Mnd Bool
highStm level (Move 1 dest src) =
    do a <- highExp (level + 1) dest</pre>
       a' <- case dest of
               Temp _ _ -> return True
               Mem _ _ -> return True
                        -> do awarn (annExp dest)
                                     ["invalid move destination:",
                                      show dest1
                              return False
       b <- highExp (level + 1) src
       return $ a && a' && b
                                                ◆□▶ ◆□▶ ◆三▶ ◆三 ◆90℃
```

... in C++

```
52 lines matching "switch\\|case\\|default\\|/" in buffer codegen.cc.
28:switch (stm.kind_get ())
30:
      case Tree::move_kind :
36:
          switch (dst->kind_get ())
38:
            case Tree::mem kind: // dst
41:
                // MOVE (MEM (...), ...)
42:
                switch (src.kind_get ())
44:
                    // MOVE (MEM (...), MEM (...))
45:
                  case Tree::mem_kind : // src
55:
                  default : // src
57:
                      // MOVE (MEM (...) , e1)
59:
                      switch (addr->kind_get ())
61:
                         case Tree::binop_kind : // addr
                             // MOVE (MEM (BINOP (..., ..., ...)) , e1)
63:
69:
                             switch (binop.oper_get ())
71:
                               case Binop::minus:
73:
                               case Binop::plus:
74:
                                 // MOVE (MEM (BINOP (+/-, e1, CONST (i))),
77:
                                 // MOVE (MEM (BINOP (+/-, CONST (i), e1)) ,
87:
                               default:
88:
                                 // MOVE (MEM (BINOP (♥.., ♥ ₹..., ₹...)) , ∞e1)
```

... in C++

```
case Node::move kind :
 DOWN_CAST (Move, move, stm);
 const Exp* dst = move.dst_get (); const Exp* src = move.src_get ();
 switch (dst->kind_get ()) {
    case Node::mem_kind : { // dst
      DOWN CAST (Mem. mem. *dst):
     // MOVE (MEM (...), ...)
     switch (src.kind_get ()) {
       // MOVE (MEM (...), MEM (...))
     case Node::mem kind: // src
     default : { // src
       // MOVE (MEM (...) , e1)
        const Exp* addr = dst.exp_get ();
        switch (addr->kind_get ()) {
        case Node::binop_kind : { // addr
          // MOVE (MEM (BINOP (..., ..., ...)) , e1)
          DOWN_CAST (Binop, binop, *addr);
          const Exp* binop_left = binop.left_get ();
          const Exp* binop_right = binop.right_get ();
          short sign = 1;
          switch (binop.oper get ()) {
          case Binop::minus: sign = -1:
          case Binop::plus:
           // MOVE (MEM (BINOP (+/-, e1, CONST (i))), e2)
           if (binop_right->kind_get () == Node::const_kind)
              std::swap (binop_left, binop_right);
           // MOVE (MEM (BINOP (+/-, CONST (i), e1)) , e2)
           if (binop left->kind get () == Node::const kind) {
              DOWN_CAST (Const, const_left, *binop_left);
              emit (_assembly->store_build (munchExp (src),
                                            munchExp (* binop_right), < -> + -> + -> + -> + -> + ->
```

... in C++

Break down long switches into smaller functions.

Twig, Burg, iBurg [Fraser et al., 1992]

Twig, Burg, iBurg [Fraser et al., 1992]

```
/* ... */
%%
stmt: ASGNI(disp,reg) = 4 (1);
stmt:
      reg = 5;
       ADDI(reg,rc) = 6 (1);
reg:
       CVCI(INDIRC(disp)) = 7 (1);
reg:
       IOI = 8:
reg:
       disp = 9 (1);
reg:
     ADDI(reg, con) = 10;
disp:
      ADDRLP = 11;
disp:
rc: con = 12;
rc: reg = 13;
con: CNSTI = 14;
con: IOI = 15;
%%
/* ... */
```

MonoBURG

```
binop: Binop(lhs : exp, rhs : Const)
  auto binop = tree.cast<Binop>();
  auto cst = rhs.cast<Const>();
  EMIT(IA32_ASSEMBLY
       .binop_build(binop->oper_get(), lhs->asm_get(),
                    cst->value_get(), tree->asm_get()));
binop: Binop(lhs : exp, rhs : exp)
  auto binop = tree.cast<Binop>();
  EMIT(IA32 ASSEMBLY
       .binop_build(binop->oper_get(), lhs->asm_get(),
                    rhs->asm_get(), tree->asm_get()));
```

Bibliography I

Anisko, R. (2003).

Havm.

http://tiger.lrde.epita.fr/Havm.

Appel, A. W. (1998).

Modern Compiler Implementation in C, Java, ML.

Cambridge University Press.

- Fraser, C. W., Hanson, D. R., and Proebsting, T. A. (1992).
 Engineering a simple, efficient code-generator generator.

 ACM Letters on Programming Languages and Systems, 1(3):213–226.
- Larus, J. R. (1990).

SPIM S20: A MIPS R2000 simulator.

Technical Report TR966, Computer Sciences Department, University of Wisconsin–Madison.

Liveness Analysis

Akim Demaille Étienne Renault Roland Levillain first.last@lrde.epita.fr

EPITA — École Pour l'Informatique et les Techniques Avancées

May 19, 2018

Liveness Analysis

- 1 Control Flow Graph
- 2 Liveness
- 3 Various Dataflow Analysis
- 4 Interference Graph

Control Flow Graph

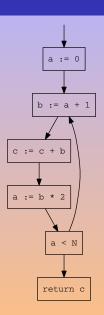
- 1 Control Flow Graph
- 2 Liveness
- 3 Various Dataflow Analysis
- 4 Interference Graph

Control Flow Graph [Appel, 1998]

```
a := 0
L1: b := a + 1
c := c + b
a := b * 2
if a < N goto L1
return c
```

Control Flow Graph [Appel, 1998]

```
a := 0
L1: b := a + 1
c := c + b
a := b * 2
if a < N goto L1
return c
```



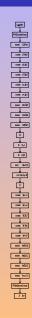
7.tig

1 + 2 * 3

7's Pre-Assembly

```
tc_main:
# Allocate frame
                                11:
                $x13, $ra
        move
        move
                $x5, $s0
                                                 $s0, $x5
                                         move
                $x6, $s1
                                                 $s1, $x6
        move
                                         move
                $x7, $s2
                                                 $s2, $x7
        move
                                        move
                $x8, $s3
                                                 $s3, $x8
        move
                                        move
                $x9, $s4
                                                 $s4, $x9
        move
                                         move
                $x10, $s5
                                                 $s5, $x10
        move
                                        move
                                                 $s6, $x11
                $x11, $s6
        move
                                        move
                $x12, $s7
                                                 $s7, $x12
        move
                                         move
10:
                                                 $ra, $x13
                                         move
        li.
                $x1, 1
                                # Deallocate frame
        li 
                $x2, 2
                                                 $ra
                                         jr
        mul
                $x3, $x2, 3
        add
                $x4, $x1, $x3
```

7's Flowgraph



7000.tig

1 | 2 & 3

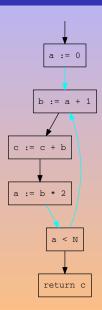
7000's Pre-Assembly

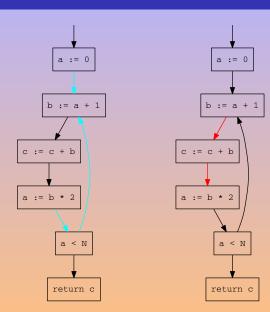
```
10:
tc_main:
                                         li
                                                 $x1, 1
# Allocate frame
                                         li 
                                                 $x5, 3
                $x6, $ra
        move
                                                 $x5, 0, 13
                                         bne
18:
                                 14:
               $x3, 1
        li
                                         li
                                                 $x1, 0
        bne
                $x3, 0, 15
                                 13:
16:
                                                $x0, $x1
                                         move
                $x4, 2
        li
                                                 12
        bne
                $x4, 0, 10
                                15:
11:
                                                 17
                $x0, 0
        li
                                19:
12:
                                         move $ra, $x6
17:
                                 # Deallocate frame
                19
                                                 $ra
                                         jr
```

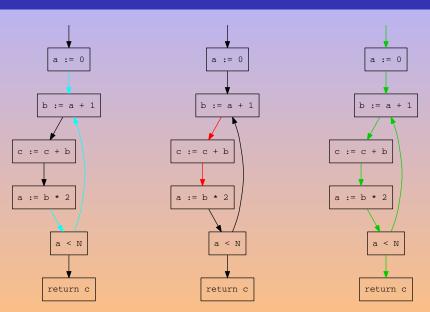
7000's Flowgraph



- 1 Control Flow Graph
- 2 Liveness
- 3 Various Dataflow Analysis
- 4 Interference Graph







Dataflow Equations for Liveness Analysis

$$in[n] = use[n] \cup (out[n] \setminus def[n])
out[n] = \bigcup_{s \in succ[n]} in[s]$$

	use	def	in	out	in	out	in	out	in	out
1		а								
2	а	Ь								
3	bc	С								
4	b	а								
5	а									
6	С									

	use	def	in	out	in	out	in	out
1		а						
2	а	b						
2	bc	С						
4	b	a						
5	a							
6	С							

$$\begin{array}{rcl} \operatorname{in}[n] & = & \operatorname{use}[n] \cup (\operatorname{out}[n] \setminus \operatorname{def}[n]) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

1st step

	use	def	in	out	in	out	in	out	in	out
1		a								
2	а	b	a							
3	bc	С	bc							
4	b	a	b							
5	а		а	а						
6	С		С							

	use	def	in	out	in	out	in	out
1		a						
2	а	b						
2	bc	С						
4	b	a						
5	a							
6	С							

$$\begin{array}{rcl} \operatorname{in}[n] & = & \operatorname{use}[n] \cup (\operatorname{out}[n] \setminus \operatorname{def}[n]) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

			1st	step	2nd	step				
	use	def	in	out	in	out	in	out	in	out
1		а				а				
2	а	b	a		a	bc				
3	bc	С	bc		bc	b				
4	b	a	b		b	а				
5	а		a	а	a	ac				
6	С		С		С					

	use	def	in	out	in	out	in	out
1		a						
2	а	b						
2	bc	С						
4	b	a						
5	a							
6	С							

$$\begin{array}{lcl} \operatorname{in}[n] & = & \operatorname{use}[n] \cup (\operatorname{out}[n] \setminus \operatorname{def}[n]) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

			1st	step	2nd	step	3rd	step		
	use	def	in	out	in	out	in	out	in	out
1		а				а		а		
2	а	b	a		a	bc	ac	bc		
3	bc	С	bc		bc	b	bc	b		
4	b	a	b		b	а	b	а		
5	а		a	а	a	ac	ac	ac		
6	С		С		С		С			

	use	def	in	out	in	out	in	out
1		a						
2	а	b						
2	bc	С						
4	b	a						
5	a							
6	С							

$$\begin{array}{rcl} \operatorname{in}[n] & = & \operatorname{use}[n] \cup \left(\operatorname{out}[n] \setminus \operatorname{def}[n]\right) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

			1st	step	2nd	step	3rd step		4th step	
	use	def	in	out	in	out	in	out	in	out
1		а				а		а		ac
2	а	b	a		a	bc	ac	bc	ac	bc
3	bc	С	bc		bc	b	bc	b	bc	С
4	b	a	b		b	а	b	а	b	ac
5	a		a	а	a	ac	ac	ac	ac	ac
6	С		С		С		С		С	

	use	def	in	out	in	out	in	out
1		а						
2	а	b						
2	bc	С						
4	b	a						
5	a							
6	С							

$$\begin{array}{rcl} \operatorname{in}[n] & = & \operatorname{use}[n] \cup (\operatorname{out}[n] \setminus \operatorname{def}[n]) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

			1st	step	2nd	step	3rd	step	4th	step
	use	def	in	out	in	out	in	out	in	out
1		а				а		а		ac
2	а	b	а		a	bc	ac	bc	ac	bc
3	bc	С	bc		bc	b	bc	b	bc	С
4	Ь	a	b		b	а	b	а	b	ac
5	a		а	а	a	ac	ac	ac	ac	ac
6	С		С		С		С		С	

5th step

			Juli	step				
	use	def	in	out	in	out	in	out
1		а	С	ac				
2	a	b	ac	bc				
3	bc	С	bc	b				
4	Ь	a	bc	ac				
5	а		ac	ac				
6	С		С					

$$\begin{array}{rcl} \operatorname{in}[n] & = & \operatorname{use}[n] \cup (\operatorname{out}[n] \setminus \operatorname{def}[n]) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

			1st step		2nd	2nd step		3rd step		4th step	
	use	def	in	out	in	out	in	out	in	out	
1		а				а		а		ac	
2	а	b	a		a	bc	ac	bc	ac	bc	
3	bc	С	bc		bc	b	bc	b	bc	С	
4	b	a	b		b	a	b	a	b	ac	
5	а		a	а	a	ac	ac	ac	ac	ac	
6	С		С		С		С		С		

			5th	5th step 6th s		step		
	use	def	in	out	in	out	in	out
1		а	С	ac	С	ac		
2	а	Ь	ac	bc	ac	bc		
3	bc	С	bc	b	bc	bc		
4	Ь	a	bc	ac	bc	ac		
5	a		ac	ac	ac	ac		
6	_		_		_			

$$\begin{array}{rcl} \operatorname{in}[n] & = & \operatorname{use}[n] \cup (\operatorname{out}[n] \setminus \operatorname{def}[n]) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

			1st step		2nd	2nd step		3rd step		4th step	
	use	def	in	out	in	out	in	out	in	out	
1		а				а		а		ac	
2	а	b	a		a	bc	ac	bc	ac	bc	
3	bc	С	bc		bc	b	bc	b	bc	С	
4	b	a	b		b	а	b	а	b	ac	
5	а		a	а	a	ac	ac	ac	ac	ac	
6	С		С		С		С		С		

			5th step		6th step		7th step	
	use	def	in	out	in	out	in	out
1		a	С	ac	С	ac	С	ac
2	a	b	ac	bc	ac	bc	ac	bc
3	bc	С	bc	b	bc	bc	bc	bc
4	b	a	bc	ac	bc	ac	bc	ac
5	a		ac	ac	ac	ac	ac	ac
6	С		С		С		С	

$$\begin{array}{rcl} \operatorname{in}[n] & = & \operatorname{use}[n] \cup (\operatorname{out}[n] \setminus \operatorname{def}[n]) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

Liveness Calculation (Forward)

			1st step		2nd step		3rd step		4th step	
	use	def	in	out	in	out	in	out	in	out
1		а				а		а		ac
2	а	Ь	a		a	bc	ac	bc	ac	bc
3	bc	С	bc		bc	b	bc	b	bc	С
4	b	a	b		b	а	b	а	b	ac
5	а		a	а	a	ac	ac	ac	ac	ac
6	С		С		С		С		С	

			5th step		6th step		7th step	
	use	def	in	out	in	out	in	out
1		a	С	ac	С	ac	С	ac
2	а	Ь	ac	bc	ac	bc	ac	bc
3	bc	С	bc	b	bc	bc	bc	bc
4	b	a	bc	ac	bc	ac	bc	ac
5	а		ac	ac	ac	ac	ac	ac
6	С		С		С		С	

$$\begin{array}{rcl} \operatorname{in}[n] & = & \operatorname{use}[n] \cup (\operatorname{out}[n] \setminus \operatorname{def}[n]) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

	use	def	out	in	out	in	out	in
6	С							
5	a							
4	b	a						
3	bc	С						
2	a	b						
1		a						

$$\begin{array}{lll} \operatorname{in}[n] & = & \operatorname{use}[n] \cup \left(\operatorname{out}[n] \setminus \operatorname{def}[n]\right) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

1st step

	use	def	out	in	out	in	out	in	
6	С			С					
5	а		С	ac					
4	b	a	ac	bc					
3	bc	С	bc	bc					
2	а	b	bc	ac					
1		a	ac	С					

$$\begin{array}{lll} \operatorname{in}[n] & = & \operatorname{use}[n] \cup (\operatorname{out}[n] \setminus \operatorname{def}[n]) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

			1st s	step	2nd	step		
	use	def	out	in	out	in	out	in
6	С			С		С		
5	а		С	ac	ac	ac		
4	b	a	ac	bc	ac	bc		
3	bc	с	bc	bc	bc	bc		
2	а	b	bc	ac	bc	ac		
1		a	ac	С	ac	С		

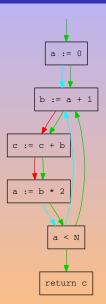
$$\begin{array}{lcl} \operatorname{in}[n] & = & \operatorname{use}[n] \cup (\operatorname{out}[n] \setminus \operatorname{def}[n]) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

			1st step		2nd step		3rd step	
	use	def	out	in	out	in	out	in
6	С			С		С		С
5	а		С	ac	ac	ac	ac	ac
4	Ь	a	ac	bc	ac	bc	ac	bc
3	bc	С	bc	bc	bc	bc	bc	bc
2	a	b	bc	ac	bc	ac	bc	ac
1		a	ac	С	ac	С	ac	С

$$\begin{array}{rcl} \operatorname{in}[n] & = & \operatorname{use}[n] \cup (\operatorname{out}[n] \setminus \operatorname{def}[n]) \\ \operatorname{out}[n] & = & \bigcup_{s \in \operatorname{succ}[n]} \operatorname{in}[s] \end{array}$$

Control Flow Graph [Appel, 1998]

```
a := 0
L1: b := a + 1
c := c + b
a := b * 2
if a < N goto L1
return c
```



Conservative Approximation

Suppose d a variable not used in the fragment of code

Another Solution

	use	def	out	in
1		а		
2	а	b		
3	bc	С		
4	b	а		
5	а			
6	С			

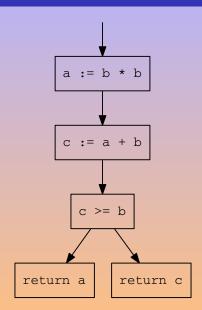
Conservative Approximation

Suppose d a variable not used in the fragment of code

Another Solution

	use	def	out	in
1		а	cd	acd
2	а	b	acd	bcd
3	bc	С	bcd	bcd
4	b	а	bcd	acd
5	а		acd	acd
6	С		С	

Conservative Approximation



ors.tig

1 | 2

ors' Flowgraph



ors' Liveness Graph



Various Dataflow Analysis

- 1 Control Flow Graph
- 2 Liveness
- 3 Various Dataflow Analysis
- 4 Interference Graph

Optimizing Compiler

- First step toward optimizing compilers
- How definitions and uses are related to each other
- What value a variable may have at a given point
- Constant propagation
- Common sub-expression elimination
- Copy propagation
- Dead Code Elimination

Constant propagation

An ambiguous definition is a statement that might or not assign a temporary t. For instance, a call may sometimes modifies t and sometimes not.

We don't have this problem for tiger due to excaping variables.

Don't loose optimisation! Consider every definiton as ambiguous

We need to define the set of definitions that reach the beginning and the end of each node.

- gen: when enter this statement, we know that we will reach the end of it
- kills: any statement that invalidates a gen
- begin[n]: which statements can reach the begining of statement n
- \bullet end[n]: which statements can reach the end of statement n

Reaching definition [Appel, 1998]

```
a := 5

c := 1

L1: if c > a goto L2

c := c + c

goto L1

L2: a := c - a

c := 0
```

	gen	kills	begin	end	begin	end	begin	end
1	1	6						
2	2	4,7						
3								
4	4	2,7						
5								
6	6	1						
7	7	2,4						

$$\operatorname{end}[n] = \operatorname{gen}[n] \cup \left(\operatorname{begin}[n] \setminus \operatorname{kills}[n]\right)$$

$$\operatorname{begin}[n] = \bigcup_{p \in \operatorname{pred}[n]} \operatorname{end}[p]$$

1st step

	gen	kills	begin	end	begin	end	begin	end
1	1	6		1				
2	2	4,7	1	1,2				
3			1,2	1,2				
4	4	2,7	1,2	1,4				
5			1,4	1,4				
6	6	1	1,2	2,6				
7	7	2,4	2,6	6,7				

$$\operatorname{end}[n] = \operatorname{gen}[n] \cup \left(\operatorname{begin}[n] \setminus \operatorname{kills}[n]\right)$$

$$\operatorname{begin}[n] = \bigcup_{p \in \operatorname{pred}[n]} \operatorname{end}[p]$$

			1st step		2nd step			
	gen	kills	begin	end	begin	end	begin	end
1	1	6		1		1		
2	2	4,7	1	1,2	1	1,2		
3			1,2	1,2	1,2,4	1,2,4		
4	4	2,7	1,2	1,4	1,2,4	1,4		
5			1,4	1,4	1,4	1,4		
6	6	1	1,2	2,6	1,2,4	2,4,6		
7	7	2,4	2,6	6,7	2,4,6	6.7		

$$\operatorname{end}[n] \ = \ \operatorname{gen}[n] \cup \left(\operatorname{begin}[n] \setminus \operatorname{kills}[n]\right)$$

$$\operatorname{begin}[n] = \bigcup_{p \in \operatorname{pred}[n]} \operatorname{end}[p]$$

			1st step		2nd step		3rd step	
	gen	kills	begin	end	begin	end	begin	end
1	1	6		1		1		1
2	2	4,7	1	1,2	1	1,2	1	1,2
3			1,2	1,2	1,2,4	1,2,4	1,2,4	1,2,4
4	4	2,7	1,2	1,4	1,2,4	1,4	1,2,4	1,4
5			1,4	1,4	1,4	1,4	1,4	1,4
6	6	1	1,2	2,6	1,2,4	2,4,6	1,2,4	2,4,6
7	7	2,4	2,6	6,7	2,4,6	6,7	2,4,6	6,7

$$\operatorname{end}[n] = \operatorname{gen}[n] \cup \left(\operatorname{begin}[n] \setminus \operatorname{kills}[n]\right)$$

$$\operatorname{begin}[n] = \bigcup_{p \in \operatorname{pred}[n]} \operatorname{end}[p]$$

Constant Propagation

- If we have a statement $d_1: t:=c$, with c constant, and another statement d_2 that uses t.
- t is constant
- if d_1 reaches d_2 and no other definition of t reaches d_2
- then we can rewrite d_2

In the previous example, only one definition of a reaches statement 3 so we can replace c>a by c>5.

Copy Propagation

- If we have a statement $d_1: t := z$, with z variable, and another statement d_2 that uses t.
- t is constant
- if d_1 reaches d_2 and no other definition of t reaches d_2 and the is no definition of t in all pathes between t and t
- then we can rewrite d_2

Good register allocator will automatically detects some such cases.

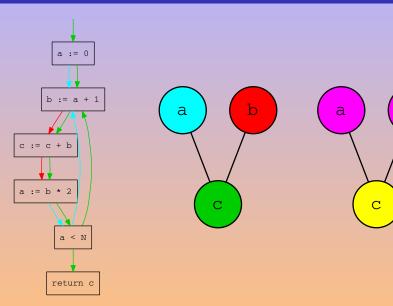
Optimizing compiler

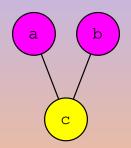
The removal of dead statements (or other optimizations) might introduce new dead statements.

To avoid the need for repeated global calculation, several strategies exist:

- Cutoff: perform no more than k round
- Cascading analysis: predict the cascade of effects of an optimization. Value numbering is a typical case of cascading analysis
- Incremental dataflow analysis: patch the dataflow after applying an optimization.

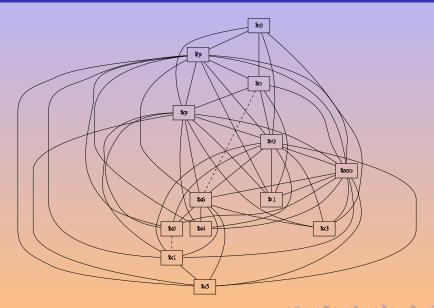
- 1 Control Flow Graph
- 2 Liveness
- 3 Various Dataflow Analysis
- 4 Interference Graph



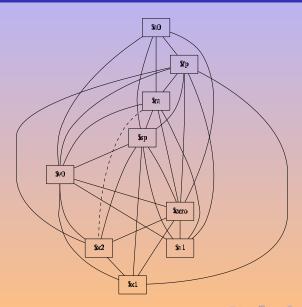


```
r1 := 0
L1: r1 := r1 + 1
    r2 := r2 + r1
    r1 := r1 * 2
    if r1 < N goto L1
    return r2
```





ors' Interference Graph

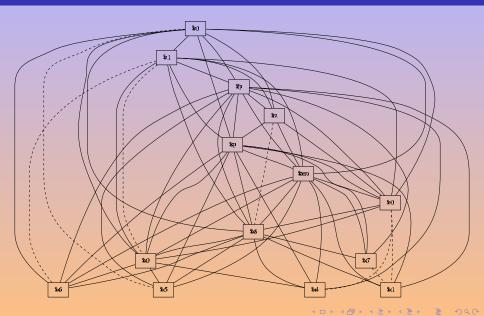


fact.tig

fact's Liveness Graph



fact's Interference Graph



Bibliography I

Appel, A. W. (1998).

Modern Compiler Implementation in C, Java, ML.

Cambridge University Press.

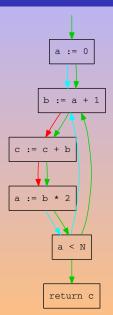
Akim Demaille Étienne Renault Roland Levillain first.last@lrde.epita.fr

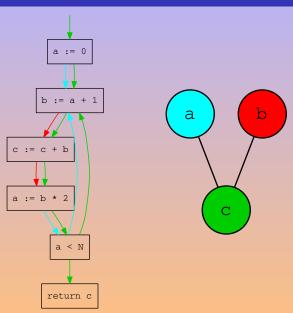
EPITA — École Pour l'Informatique et les Techniques Avancées

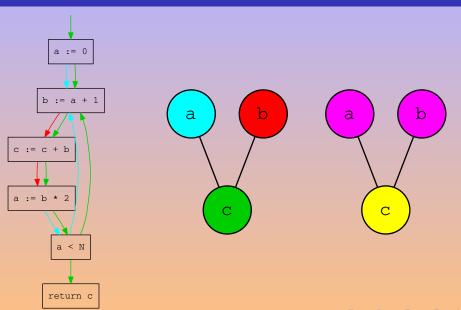
May 19, 2018

- 1 Interference Graph
- 2 Coloring by Simplification
- 3 Alternatives to Graph Coloring

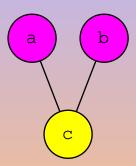
- Interference Graph
- 2 Coloring by Simplification
- 3 Alternatives to Graph Coloring

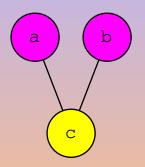






```
a := 0
L1: b := a + 1
c := c + b
a := b * 2
if a < N goto L1
return c
```





```
r1 := 0
L1: r1 := r1 + 1
r2 := r2 + r1
r1 := r1 * 2
if r1 < N goto L1
return r2
```

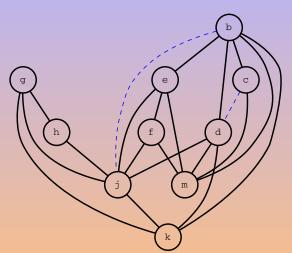
Coloring by Simplification

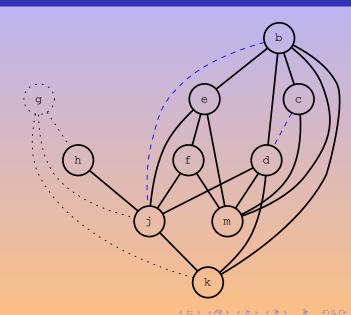
- Interference Graph
- 2 Coloring by Simplification
 - Spilling
 - Coalescing
 - Precolored Nodes
 - Implementation
- 3 Alternatives to Graph Coloring

Interference Graph [Appel, 1998]

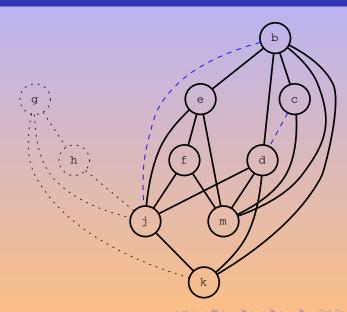
Four registers: r1, r2, r3, r4.

```
live in: k j
  g := [j + 12]
  h := k - 1
  f := g * h
  e := [j + 8]
  m := [j + 16]
  b := \lceil f \rceil
  c := e + 8
  d := c
  k := m + 4
  j := b
live out: d k j
```

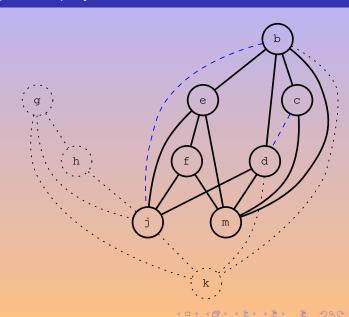




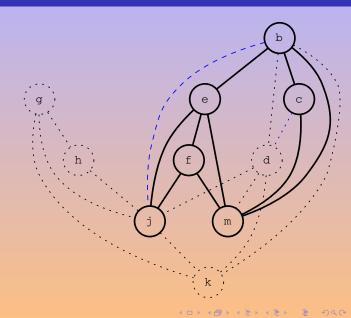
g



h g



k h



d k

h

g



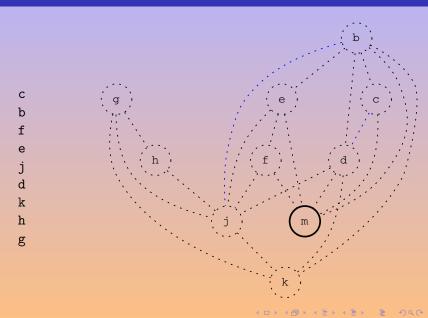
j d k h

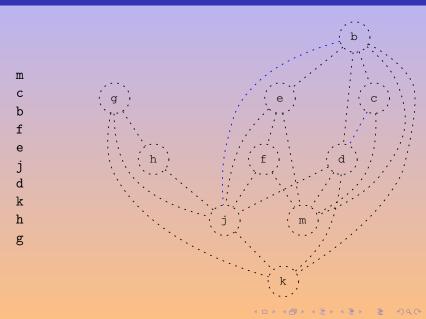


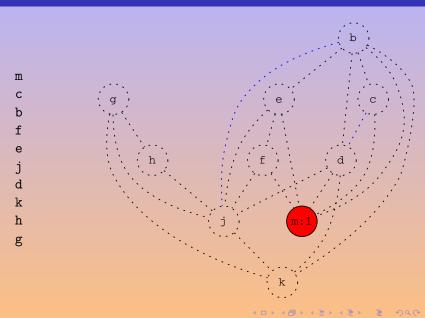
j d k h

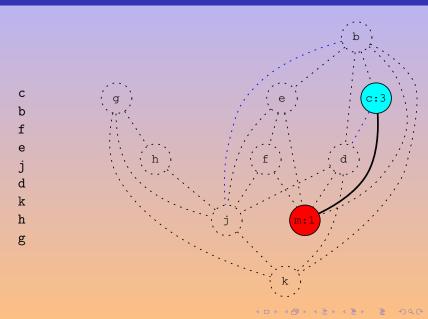


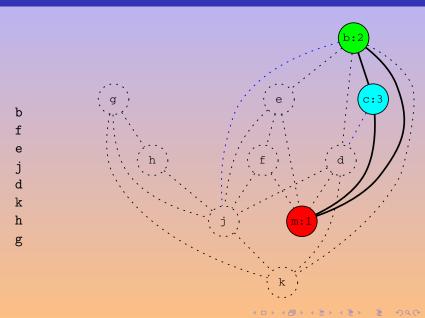


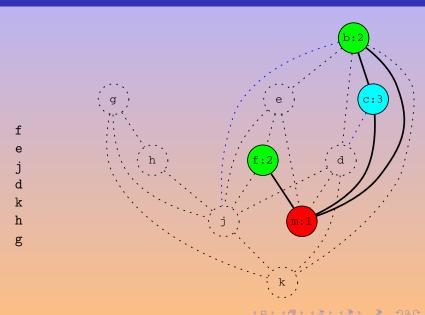


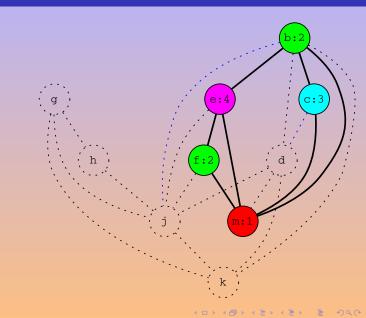






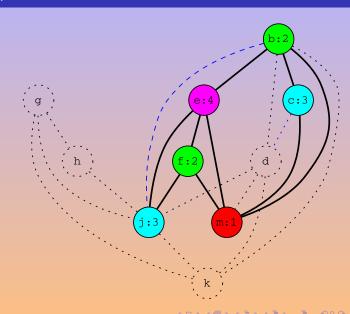




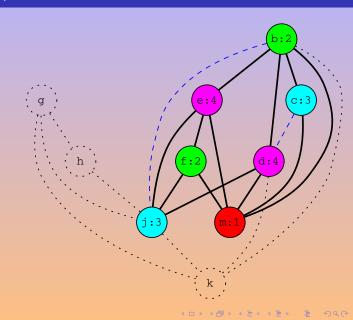


e j d k

g

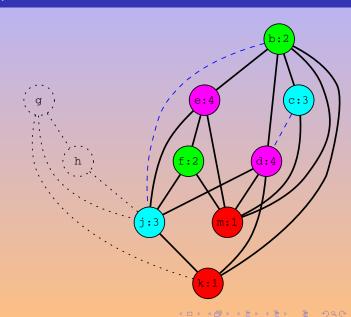


j d k h

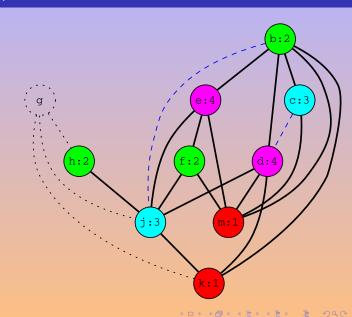


d k

g

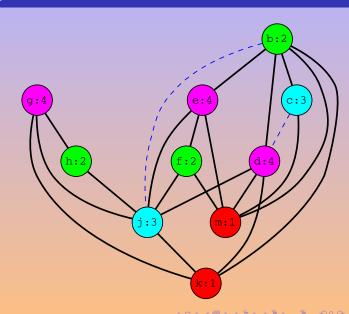


k h g



h

g



g

Result

```
live in: k j
  g := [j + 12]
 h := k - 1
  f := g * h
  e := [j + 8]
  m := [j + 16]
  b := [f]
  c := e + 8
  d := c
 k := m + 4
  i := b
live out: d k j
```

```
live in: r1 r3
  r4 := [r3 + 12]
  r2 := r1 - 1
  r2 := r4 * r2
 r4 := [r3 + 8]
  r1 := [r3 + 16]
 r2 := [r2]
  r3 := r4 + 8
 r4 := r3
 r1 := r1 + 4
 r3 := r2
live out: r4 r1 r3
```



build the conflict graph from the program

simplify the nodes with insignificant degree select (or color) while rebuilding the graph.



build the conflict graph from the program simplify the nodes with insignificant degree

select (or color) while rebuilding the graph.



build the conflict graph from the program simplify the nodes with insignificant degree select (or color) while rebuilding the graph.



build the conflict graph from the program simplify the nodes with insignificant degree select (or color) while rebuilding the graph.

Based on:

A.B. Kempe. On the Geographical problem of the four colors, Am. J. Math 2, 193–200, 1879.

[Appel, 1998, Matz, 2003]

Yes, but What Color? [Matz, 2003]

- Usually, first-fit (registers are ordered).
- Trying caller save first helps.
- Biased Coloring. [Briggs, 1992]
 Use a color already unavailable to our neighbors

Yes, but What Color? [Matz, 2003]

- Usually, first-fit (registers are ordered).
- Trying caller save first helps.
- Biased Coloring. [Briggs, 1992]
 Use a color already unavailable to our neighbors

Yes, but What Color? [Matz, 2003]

- Usually, first-fit (registers are ordered).
- Trying caller save first helps.
- Biased Coloring. [Briggs, 1992]
 Use a color already unavailable to our neighbors.

- Interference Graph
- 2 Coloring by Simplification
 - Spilling
 - Coalescing
 - Precolored Nodes
 - Implementation
- 3 Alternatives to Graph Coloring

A map can always be colored with 4 colors...

A map can always be colored with 4 colors. . .

But for graph coloring, there is no reason for:

- this simple heuristics to always find a solution,
- a solution to always exist...

A map can always be colored with 4 colors...

But for graph coloring, there is no reason for:

- this simple heuristics to always find a solution,
- a solution to always exist...

Not enough registers

$$t1 := t1 + t2$$

So use the stack

$$[sp + 4] := [sp + 4] + [sp + 8]$$

But use temporaries to do so!

Why should it solve the problem?

Not enough registers

$$t1 := t1 + t2$$

So use the stack

$$[sp + 4] := [sp + 4] + [sp + 8]$$

But use temporaries to do so!

• Why should it solve the problem?

Not enough registers

$$t1 := t1 + t2$$

So use the stack

$$[sp + 4] := [sp + 4] + [sp + 8]$$

• But use temporaries to do so!

Why should it solve the problem?

Not enough registers

$$t1 := t1 + t2$$

So use the stack

$$[sp + 4] := [sp + 4] + [sp + 8]$$

• But use temporaries to do so!

• Why should it solve the problem?

Register Allocation with Spills



spill when one cannot simplify, the (uses of the) temporary must be rewritten using the stack.

rebuild but then, the conflict graph is to be rewritten

[Appel, 1998, Matz, 2003]

- The simplification order does not matter
- The spilling order matters
- Spilling decreases the degree of the neighbors
- ... hence it enables additional simplifications
- ...so "first spilled, last served"
- ... therefore: spill cheap temporaries
 - few def/uses
- pay attention to loops

- The simplification order does not matter
- The spilling order matters
- Spilling decreases the degree of the neighbors
- ... hence it enables additional simplifications
- ... so "first spilled, last served
- ... therefore: spill cheap temporaries
 - few def/
- pay attention to loops

- The simplification order does not matter
- The spilling order matters
- Spilling decreases the degree of the neighbors
- ... hence it enables additional simplifications
- ... so "first spilled, last served"
- ... therefore: spill cheap temporaries
- pay attention to I

- The simplification order does not matter
- The spilling order matters
- Spilling decreases the degree of the neighbors
- ...hence it enables additional simplifications
- ...so "first spilled. last served"
- ... therefore: spill cheap temporaries

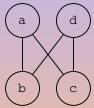
- The simplification order does not matter
- The spilling order matters
- Spilling decreases the degree of the neighbors
- ... hence it enables additional simplifications
- ...so "first spilled, last served"
- ... therefore: spill cheap temporaries

- The simplification order does not matter
- The spilling order matters
- Spilling decreases the degree of the neighbors
- ... hence it enables additional simplifications
- ...so "first spilled, last served"
- ... therefore: spill cheap temporaries
 - few def/uses
 - pay attention to loops

- The simplification order does not matter
- The spilling order matters
- Spilling decreases the degree of the neighbors
- ... hence it enables additional simplifications
- ...so "first spilled, last served"
- ... therefore: spill cheap temporaries
 - few def/uses
 - pay attention to loops

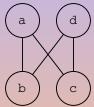
- The simplification order does not matter
- The spilling order matters
- Spilling decreases the degree of the neighbors
- ... hence it enables additional simplifications
- ...so "first spilled, last served"
- ... therefore: spill cheap temporaries
 - few def/uses
 - pay attention to loops

• We miss many opportunities to avoid the stack



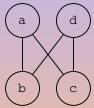
- Handle spills as if they were simplified (potential spills)
- then try to color them
- There might not be actual spills

We miss many opportunities to avoid the stack



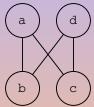
- Handle spills as if they were simplified (potential spills)
- then try to color them
- There might not be actual spills

We miss many opportunities to avoid the stack



- Handle spills as if they were simplified (potential spills)
- then try to color them
- There might not be actual spills

• We miss many opportunities to avoid the stack



- Handle spills as if they were simplified (potential spills)
- then try to color them
- There might not be actual spills

- Interference Graph
- 2 Coloring by Simplification
 - Spilling
 - Coalescing
 - Precolored Nodes
 - Implementation
- 3 Alternatives to Graph Coloring

- Some low-level form of copy propagation
- While building traces we tried to remove jumps
- While allocating registers, we try to remove moves

```
live-in: t2
t1 := ...
t2 := t1 + t2
t3 := t2
t4 := t1 + t3
t2 := t3 + t4
t1 := t2 - t4
```

- Some low-level form of copy propagation
- While building traces we tried to remove jumps
- While allocating registers, we try to remove moves

```
live-in: t2
t1 := ...
t2 := t1 + t2
t3 := t2
t4 := t1 + t3
t2 := t3 + t4
```

live-out: t1

- Some low-level form of copy propagation
- While building traces we tried to remove jumps
- While allocating registers, we try to remove moves

```
live-in: t2

t1 := ...

t2 := t1 + t2

t3 := t2

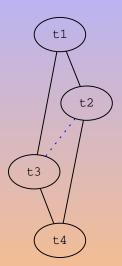
t4 := t1 + t3

t2 := t3 + t4

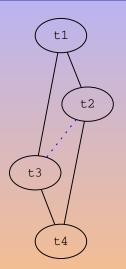
t1 := t2 - t4

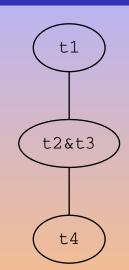
live-out: t1
```

Coalescing Improves the Coloralibility

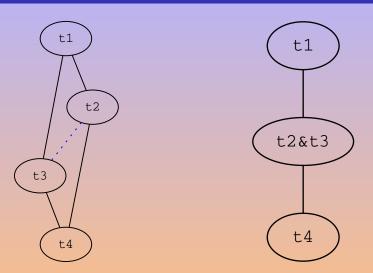


Coalescing Improves the Coloralibility





Coalescing Improves the Coloralibility



t1 and t4 have one neighbor less!

- Conservative Coalescing: don't make it harder.
- Coalesce a and b if

George's criterion is well suited for real registers

- Conservative Coalescing: don't make it harder.
- Coalesce a and b if

Briggs ab has fewer than k neighbors of significant degree. George every neighbor of a is

George's criterion is well suited for real registers

- Conservative Coalescing: don't make it harder.
- Coalesce a and b if
 - Briggs ab has fewer than k neighbors of significant degree.

• George's criterion is well suited for real registers

- Conservative Coalescing: don't make it harder.
- Coalesce a and b if

- of insignificant degreealready interfering with b
- George's criterion is well suited for real registers

- Conservative Coalescing: don't make it harder.
- Coalesce a and b if

- of insignificant degree
- already interfering with b
- George's criterion is well suited for real registers

- Conservative Coalescing: don't make it harder.
- Coalesce a and b if

- of insignificant degree
- already interfering with b
- George's criterion is well suited for real registers

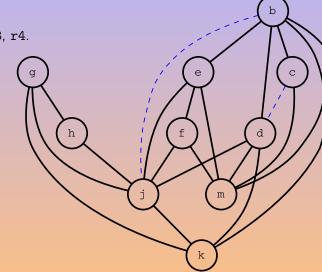
- Conservative Coalescing: don't make it harder.
- Coalesce a and b if

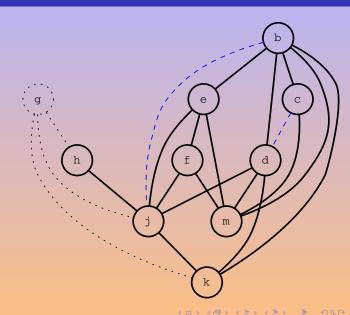
- of insignificant degree
- already interfering with b
- George's criterion is well suited for real registers

Interference Graph [Appel, 1998]

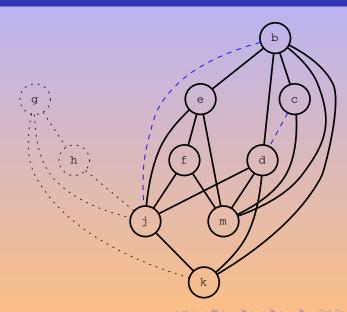
Four registers: r1, r2, r3, r4.

```
live in: k j
  g := [j + 12]
 h := k - 1
 f := g * h
  e := [i + 8]
 m := [j + 16]
  b := [f]
  c := e + 8
  d := c
  k := m + 4
  i := b
live out: d k j
```

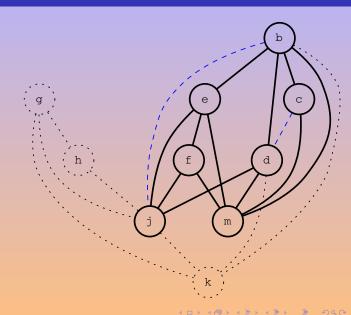




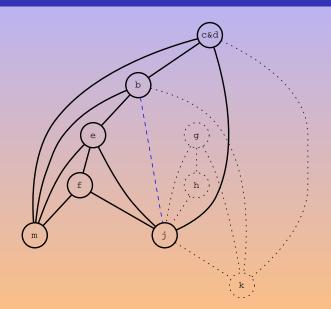
g



h g

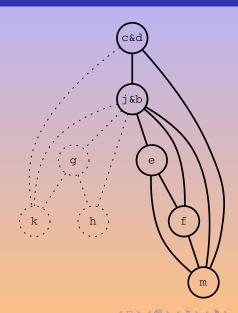


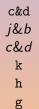
k h

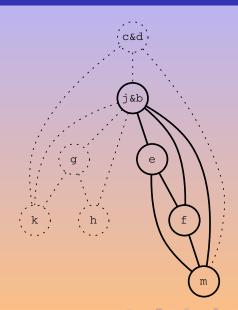


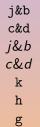
c&d k h

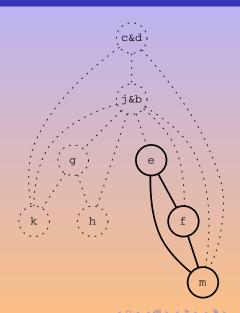


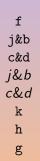


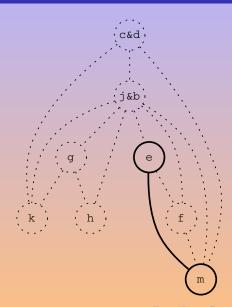


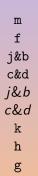


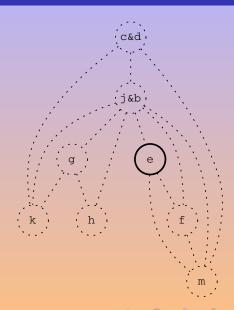


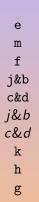


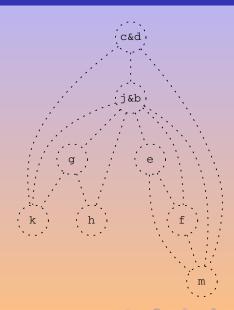


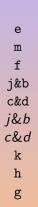


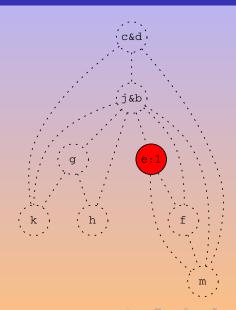


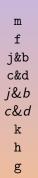


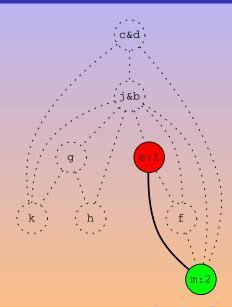


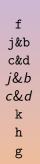


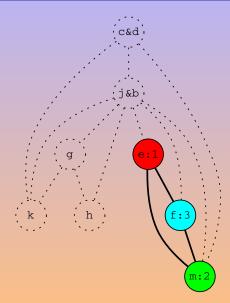


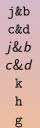


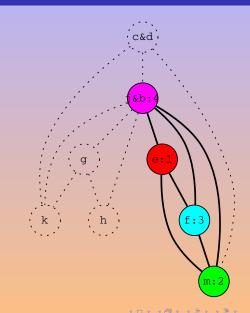


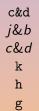


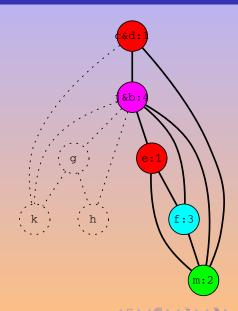


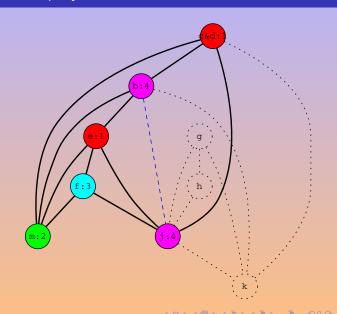




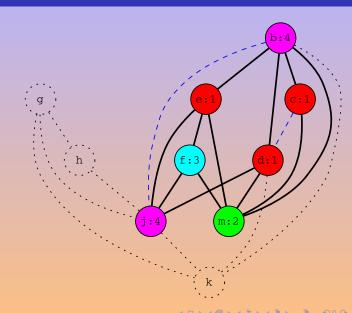




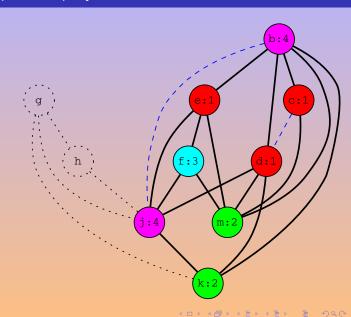




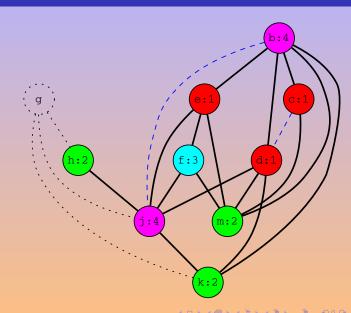
j&b c&d k h



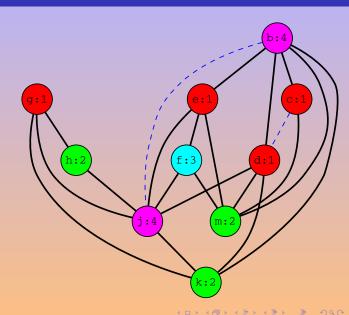
c&d k h



k h g



h g



g

Interference Graph: Result

```
live in: k j
  g := [j + 12]
 h := k - 1
  f := g * h
  e := [j + 8]
  m := [j + 16]
  b := [f]
  c := e + 8
  d := c
 k := m + 4
  j := b
live out: d k j
```

```
live in: r2 r4
  r1 := [r4 + 12]
 r2 := r2 - 1
  r3 := r1 * r2
  r1 := [r4 + 8]
  r2 := [r4 + 16]
 r4 := [r3]
  r1 := r1 + 8
# r1 := r1
r2 := r2 + 4
# r4 := r4
live out: r1 r2 r4
```

Precolored Nodes

- Interference Graph
- 2 Coloring by Simplification
 - Spilling
 - Coalescing
 - Precolored Nodes
 - Implementation
- 3 Alternatives to Graph Coloring

Some nodes are precolored: the real registers

```
the stack pointer ($sp)
the frame pointer ($fp)
the argument registers ($a0, $a1, etc.)
the return value ($v0, $v1)
the return address ($ra)
etc.
```

- They all interfere with each other
- They cannot be simplified (infinite degree)

- Some nodes are precolored: the real registers
 - the stack pointer (\$sp)
 - the frame pointer (\$fp)
 - the argument registers (\$a0, \$a1, etc.)
 - the return value (\$v0, \$v1)
 - the return address (\$ra)
 - etc.
- They all interfere with each other
- They cannot be simplified (infinite degree)

- Some nodes are precolored: the real registers
 - the stack pointer (\$sp)
 - the frame pointer (\$fp)
 - the argument registers (\$a0, \$a1, etc.)
 - the return value (\$v0, \$v1)
 - the return address (\$ra)
 - etc.
- They all interfere with each other
- They cannot be simplified (infinite degree)

- Some nodes are precolored: the real registers
 - the stack pointer (\$sp)
 - the frame pointer (\$fp)
 - the argument registers (\$a0, \$a1, etc.)
 - the return value (\$v0, \$v1)
 - the return address (\$ra)
 - etc.
- They all interfere with each other
- They cannot be simplified (infinite degree)

- Some nodes are precolored: the real registers
 - the stack pointer (\$sp)
 - the frame pointer (\$fp)
 - the argument registers (\$a0, \$a1, etc.)
 - the return value (\$v0, \$v1)
 - the return address (\$ra)
 - etc.
- They all interfere with each other
- They cannot be simplified (infinite degree)

- It just rocks!
 - Caller Save Def'd by calls.
 - Callee Save Def'd at entry, used at exit of functions.
- Register pressure will push temporaries live across calls into callee save

- It just rocks!
 - Caller Save Def'd by calls.
 - Callee Save Def'd at entry, used at exit of functions.
- Register pressure will push temporaries live across calls into callee save

- It just rocks!
 - Caller Save Def'd by calls.
 - Callee Save Def'd at entry, used at exit of functions.
- Register pressure will push temporaries live across calls into callee save

- It just rocks!
 - Caller Save Def'd by calls.
 - Callee Save Def'd at entry, used at exit of functions.
- Register pressure will push temporaries live across calls into callee save.

Conflicts

```
Minimize the conflicts ("pressure") with hard regs. Source and sink.
# Routine: fact
                       # def $s0, $s1...
10:
                       # def: $x11 use: $s0
  move $x11, $s0
                    # def: $x12 use: $s1
  move $x12, $s1
  . . .
16:
  move $s0, $x11
                     # def: $s0 use: $x11
  move $s1, $x12
                       # def: $s1 use: $x12
  . . .
                       # use: $fp, $ra, $sp,
                       # ... $v0. $zero
```

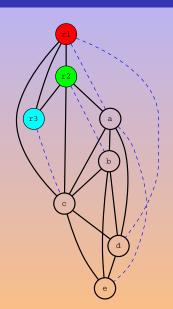
Example [Appel, 1998]

```
int
f (int a, int b)
 int d = 0;
 int e = a;
  do
   d += b;
   --e;
 } while (e > 0);
 return d;
```

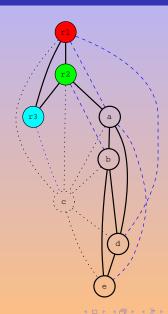
```
enter:
  c := r3
  a := r1
  b := r2
  d := 0
  e := a
loop:
  d := d + b
  e := e - 1
  if e > 0 goto loop
  r1 := d
  r3 := c
  return
# liveout: r1, r3
```

Example

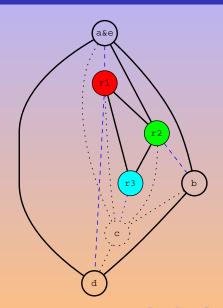
```
enter:
  c := r3
  a := r1
  b := r2
  d := 0
  e := a
loop:
  d := d + b
  e := e - 1
  if e > 0 goto loop
  r1 := d
  r3 := c
  return
# liveout: r1, r3
```



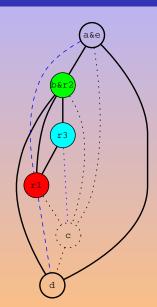




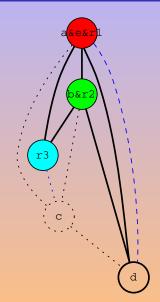




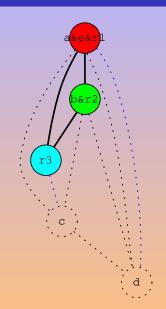
b&*r*2 *a*&*e* c



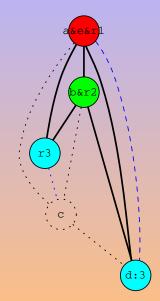
a&e&r1 b&r2 a&e c



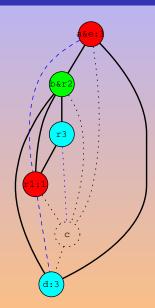
d a&e&r1 b&r2 a&e c



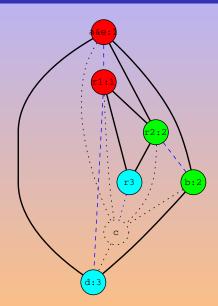
d a&e&r1 b&r2 a&e c



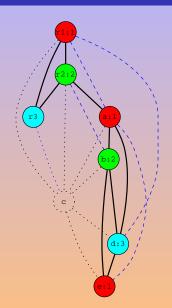
a&e&r1 b&r2 a&e c



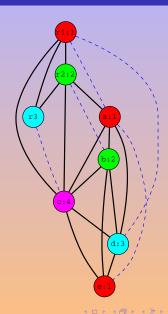




а&е с







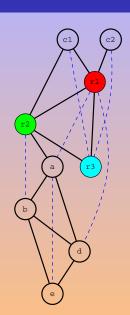
Spilling

```
enter:
  c := r3
  a := r1
  b := r2
  d := 0
  e := a
loop:
  d := d + b
  e := e - 1
  if e > 0 goto loop
  r1 := d
 r3 := c
  return
# liveout: r1, r3
```

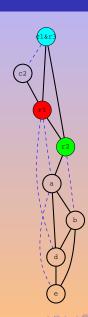
```
enter:
 c1 := r3
  [sp+8] := c1
 a := r1
 b := r2
 d := 0
 e := a
loop:
 d := d + b
 e := e - 1
 if e > 0
    goto loop
 r1 := d
 c2 := [sp+8]
 r3 := c2
 return
# liveout: r1, r3
```

Example

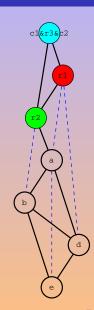
```
enter:
 c1 := r3
  [sp+8] := c1
     := r1
 b := r2
 d := 0
   := a
loop:
 d := d + b
 e := e - 1
 if e > 0
    goto loop
 r1 := d
 c2 := [sp+8]
 r3 := c2
 return
 liveout: r1, r3
```



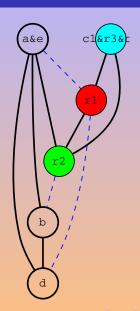
c1&r3



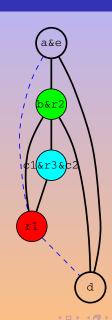
*c*1&*r*3&*c*2 *c*1&*r*3



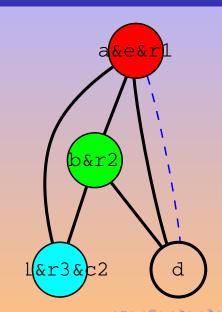
a&e c1&r3&c2 c1&r3



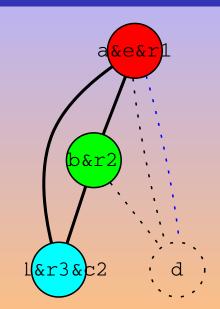
b&r2 a&e c1&r3&c2 c1&r3



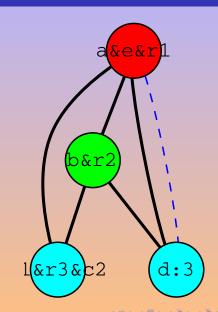
a&e&r1 b&r2 a&e c1&r3&c2 c1&r3



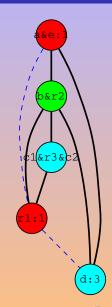
d
a&e&r1
b&r2
a&e
c1&r3&c2
c1&r3



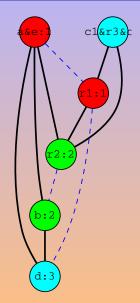
a&e&r1 b&r2 a&e c1&r3&c2 c1&r3



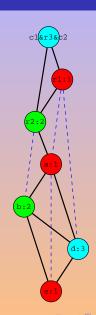
a&e&r1 b&r2 a&e c1&r3&c2 c1&r3



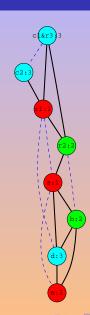
b&r2 a&e c1&r3&c2 c1&r3



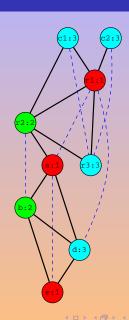
a&e c1&r3&c2 c1&r3



*c*1&*r*3&*c*2 *c*1&*r*3







Result

```
enter:
                     enter:
                                        enter:
 c1 := r3
                      r3 := r3
 [sp+8] := c1
                      [sp+8] := r3
                                          [sp+8] := r3
 a := r1
                      r1 := r1
 b := r2
                      r2 := r2
 d := 0
                      r3 := 0
                                         r3 := 0
                      r1 := r1
 e := a
loop:
                    loop:
                                        loop:
 d := d + b
                      r3 := r3 + r2
                                         r3 := r3 + r2
 e := e - 1
                      r1 := r1 - 1
                                         r1 := r1 - 1
 if e > 0
                      if r1 > 0
                                         if r1 > 0
  goto loop
                      goto loop
                                         goto loop
 r1 := d
                      r1 := r3
                                         r1 := r3
 c2 := [sp+8]
                      r3 := [sp+8]
                                         r3 := [sp+8]
 r3 := c2
                      r3 := r3
 return
                      return
                                         return
# liveout: r1, r3 # liveout: r1, r3 # liveout: r1, r3
```

- Interference Graph
- Coloring by Simplification
 - Spilling
 - Coalescing
 - Precolored Nodes
 - Implementation
- 3 Alternatives to Graph Coloring

- Naive implementation is quadratic

- Naive implementation is quadratic
- Lower with heavy use of worklists
- Queries on the conflict graph

• For more information, see [Appel, 1998]

- Naive implementation is quadratic
- Lower with heavy use of worklists
- Queries on the conflict graph
 - Iterate over neighbors, hence adjacency list
 - Existence of an edge between two nodes, hence bit matrix
- For more information, see [Appel, 1998].

- Naive implementation is quadratic
- Lower with heavy use of worklists
- Queries on the conflict graph
 - Iterate over neighbors, hence adjacency list
 - Existence of an edge between two nodes, hence bit matrix
- For more information, see [Appel, 1998].

- Naive implementation is quadratic
- Lower with heavy use of worklists
- Queries on the conflict graph
 - Iterate over neighbors, hence adjacency list
 - Existence of an edge between two nodes, hence bit matrix.
- For more information, see [Appel, 1998]

- Naive implementation is quadratic
- Lower with heavy use of worklists
- Queries on the conflict graph
 - Iterate over neighbors, hence adjacency list
 - Existence of an edge between two nodes, hence bit matrix.

Use both!

• For more information, see [Appel, 1998]

- Naive implementation is quadratic
- Lower with heavy use of worklists
- Queries on the conflict graph
 - Iterate over neighbors, hence adjacency list
 - Existence of an edge between two nodes, hence bit matrix.

Use both!

• For more information, see [Appel, 1998].

Alternatives to Graph Coloring

- 1 Interference Graph
- 2 Coloring by Simplification
- 3 Alternatives to Graph Coloring

Register Allocation for Trees

```
Can be done during instruction selection with maximal munch
function SimpleAlloc (t)
  for each nontrivial tile u child of t
    SimpleAlloc (u)
  for each nontrivial tile u child of t
    n := n - 1
  n := n + 1
  assign rn to (the root of) t
[Appel, 1998]
```

Bibliography I

- Appel, A. W. (1998).

 Modern Compiler Implementation in C, Java, ML.

 Cambridge University Press.
- Briggs, P. (1992).

 Register Allocation via Graph Coloring.

 PhD thesis, Rice University, Houston, Texas.
- Matz, M. (2003).

 Design and Implementation of a Graph Coloring Register Allocator for gcc.

pages 151-169.

Instruction scheduling

Akim Demaille, Etienne Renault, Roland Levillain

May 19, 2018

Table of contents

- Dependencies
- Dependency graph
- 3 Instruction Pipeline
- 4 Minimizing stalls
- 5 Loops unrolling
- 6 Managing caches

Dependencies analysis 1/2

Two instructions are independent they can be permuted without altering the consistency

Dependencies analysis 1/2

Two instructions are independent they can be permuted without altering the consistency

The 3 following instructions are independent

 $\begin{array}{lll} \text{inst}_1: & a \leftarrow 42 \\ \text{inst}_2: & b \leftarrow 51 \\ \text{inst}_3: & c \leftarrow 0 \end{array}$

Dependencies analysis 1/2

Two instructions are independent they can be permuted without altering the consistency

The 3 following instructions are independent

$$\begin{array}{ll} \mathsf{inst}_1: & \mathsf{a} \leftarrow \mathsf{42} \\ \mathsf{inst}_2: & \mathsf{b} \leftarrow \mathsf{51} \\ \mathsf{inst}_3: & \mathsf{c} \leftarrow \mathsf{0} \end{array}$$

• inst₁, inst₂ and inst₃ can then be reordered

Dependencies analysis 2/2

Two instructions are dependent if the first one needs to be executed before the second one.

Dependencies analysis 2/2

Two instructions are dependent if the first one needs to be executed before the second one.

 The 3 following instructions are dependent, i.e. no reordering is possible!

 $\begin{array}{ll} \mathsf{inst}_1: & \mathsf{a} \leftarrow \mathsf{42} \\ \mathsf{inst}_2: & \mathsf{b} \leftarrow \mathsf{a} + \mathsf{51} \\ \mathsf{inst}_3: & \mathsf{c} \leftarrow \mathsf{b} \times \mathsf{12} \end{array}$

Dependencies analysis 2/2

Two instructions are dependent if the first one needs to be executed before the second one.

 The 3 following instructions are dependent, i.e. no reordering is possible!

$$\begin{array}{ll} \mathsf{inst}_1: & \mathsf{a} \leftarrow \mathsf{42} \\ \mathsf{inst}_2: & \mathsf{b} \leftarrow \mathsf{a} + \mathsf{51} \\ \mathsf{inst}_3: & \mathsf{c} \leftarrow \mathsf{b} \times \mathsf{12} \end{array}$$

- Two kind of dependencies:
 - ▶ Data dependencies: the instruction manipulates a "variable" computed by another instruction.
 - ▶ **Instruction dependencies**: the instruction is a "cjump", the next instruction depends of the "cjump".

Read after Write (RAW)

An instruction reads from a location after an earlier instruction has written to it.

Read after Write (RAW)

An instruction reads from a location after an earlier instruction has written to it.

 $inst_1$: lw \$2, 0(\$4) $inst_2$: addi \$6, \$2, 42

Read after Write (RAW)

An instruction reads from a location after an earlier instruction has written to it.

 $inst_1$: lw \$2, 0(\$4) $inst_2$: addi \$6, \$2, 42

 $inst_1$ and $inst_2$ cannot be permuted, otherwise $inst_2$ would read an old value for \$2

Write after Read (WAR)

An instruction writes to a location after an earlier instruction has read from it.

Write after Read (WAR)

An instruction writes to a location after an earlier instruction has read from it.

 $inst_1: lw $2, 0($4)$

 $inst_2$: addi \$4, \$12, 42

Write after Read (WAR)

An instruction writes to a location after an earlier instruction has read from it.

inst₁: 1w \$2, 0(\$4) inst₂: addi \$4, \$12, 42

 $inst_1$ and $inst_2$ cannot be permuted, otherwise $inst_1$ would read a new value for \$4

Write after Write (WAW)

An instruction writes to a location after an earlier instruction has written to it.

Write after Write (WAW)

An instruction writes to a location after an earlier instruction has written to it.

inst₁: add \$1, \$2, \$3 inst₂: add \$1, \$5, \$6

Write after Write (WAW)

An instruction writes to a location after an earlier instruction has written to it.

 $inst_1$: add \$1, \$2, \$3 $inst_2$: add \$1, \$5, \$6

 inst_1 and inst_2 cannot be permuted, otherwise inst_1 would write an old value in \$1

Why and When reordering?

We would like to reorder the instructions within each basic block in a way which:

- preserves the dependencies between those instructions (and hence the correctness of the program)
- achieves the minimum possible number of pipeline stalls, i.e. two instructions simultaneously in the pipeline manipulates same data, registers, etc.

The two problems can be addressed separately (whew!).

Preserving and computing dependencies?

We construct a directed acyclic graph (DAG) to represent the dependencies between instructions:

- For each instruction in the basic block, create a corresponding vertex in the graph
- For each dependency between two instructions, create a corresponding (annotated) edge in the graph. Note that this edge is annotated.











 i_1 i_3

 i_2 i_4

 (i_1)

 (i_3)

 $\left(i_{5}\right)$

 $\left(i_{2}\right)$

 $\left(i_{4}\right)$

 (i_1)

 (i_3)

 $\left(i_{5}\right)$

 (i_2)

 (i_4)

 (i_6)

 $\left(i_{1}\right)$

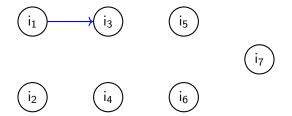
 (i_3)

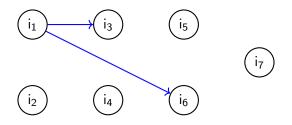
 $\left(i_{5}\right)$

 (i_2)

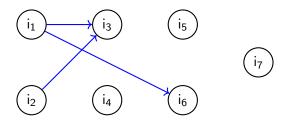
 (i_4)

 (i_6)

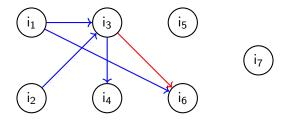


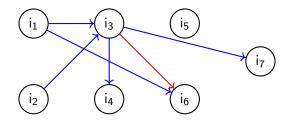


```
i_1: lw $1,0($10) | i_4: sw $3,12($10) | i_7: sw $3,16($10) | i_2: lw $2,4($10) | i_5: lw $4,8($10) | i_5: add $3,$1,$2 | i_6: add $3,$1,$4
```



```
i_1: lw $1,0($10) | i_4: sw $3,12($10) | i_7: sw $3,16($10) | i_2: lw $2,4($10) | i_5: lw $4,8($10) | i_5: add $3,$1,$2 | i_6: add $3,$1,$4
```

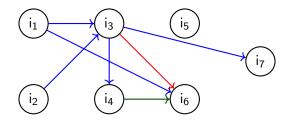


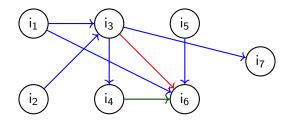


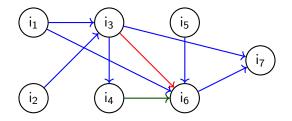
Type of dependency: RAW, WAW, WAR

10 / 57

```
i_1: lw $1,0($10) | i_4: sw $3,12($10) | i_7: sw $3,16($10) | i_2: lw $2,4($10) | i_5: lw $4,8($10) | i_5: add $3,$1,$2 | i_6: add $3,$1,$4
```







Type of dependency: RAW, WAW, WAR

10 / 57

Preserving dependencies: Critical Path 1/2

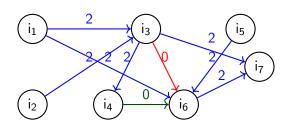
The critical path represents the longest path between two nodes. We add delays (weights) to edges:

- 0 for WAW and WAR dependencies
- 2 for RAW dependencies with memory access
- 1 for other RAW dependencies

Preserving dependencies: Critical Path 1/2

The critical path represents the longest path between two nodes. We add delays (weights) to edges:

- 0 for WAW and WAR dependencies
- 2 for RAW dependencies with memory access
- 1 for other RAW dependencies



Preserving dependencies: Critical Path 2/2

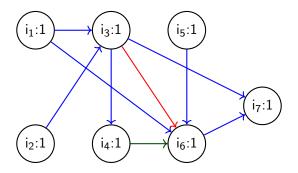
Any (reverse) topological sort of this DAG (i.e. any linear ordering of the vertices which keeps all the edges "pointing forwards") will maintain the dependencies and hence preserve the correctness of the program.

Algorithm:

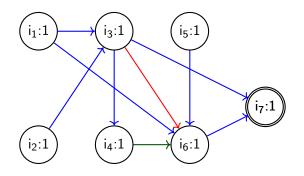
- Associate a weight 1 to all "instruction node"
- For all nodes n_i in topological postorder
 - ▶ If n; is not a leaf
 - * For all nodes n_j in $succ(n_i)$ do $n_i.weight \leftarrow max(n_i.weight, n_j.weight+ delay(n_i, n_j))$

Remember "important" edges during computations, they will form the critical path.

Delays: blue arrows 2, red and green 0



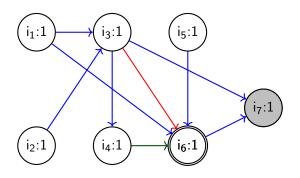
Delays: blue arrows 2, red and green 0



i₇ doesn't have successors, skip it!

13 / 57

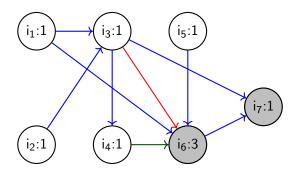
Delays: blue arrows 2, red and green 0



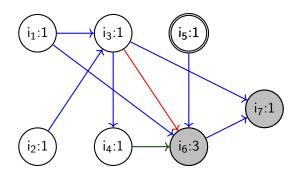
 $delay(i_6, i_7)=2 > 1$, change i_6 weight!

13 / 57

Delays: blue arrows 2, red and green 0



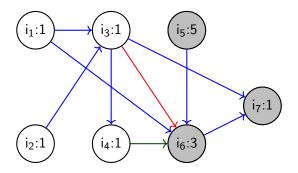
Delays: blue arrows 2, red and green 0



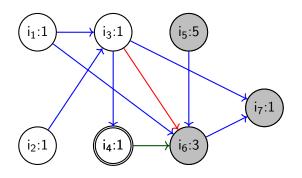
 $delay(i_5, i_6)=2 > 1$, change i_5 weight!

13 / 57

Delays: blue arrows 2, red and green 0



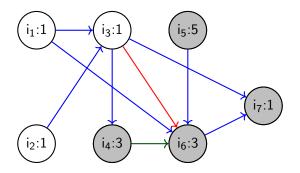
Delays: blue arrows 2, red and green 0



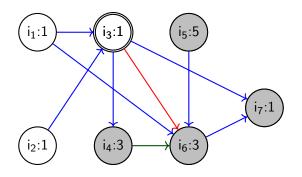
 $i_6.weight=3 > 1$, change i_4 weight!

13 / 57

Delays: blue arrows 2, red and green 0



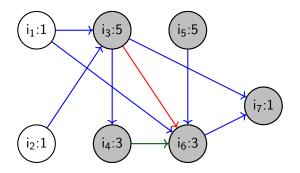
Delays: blue arrows 2, red and green 0



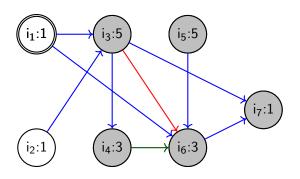
 $delay(i_3, i_4) + i_4.weight=3 > 1$, change i_3 weight!

May 19, 2018

Delays: blue arrows 2, red and green 0



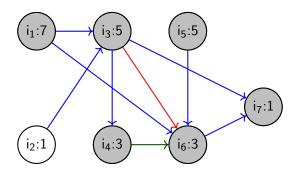
Delays: blue arrows 2, red and green 0



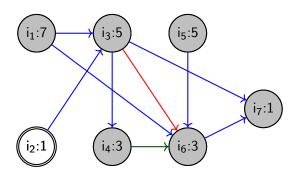
 $delay(i_1, i_3) + i_3.weight=7 > 1$, change i_1 weight!

13 / 57

Delays: blue arrows 2, red and green 0



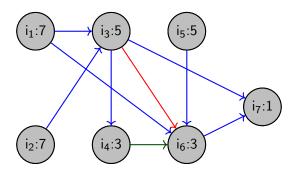
Delays: blue arrows 2, red and green 0



 $delay(i_2, i_3) + i_3.weight=7 > 1$, change i_2 weight!

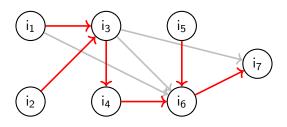
◆ロト ◆御 ト ◆ 恵 ト ◆ 恵 ・ からぐ

Delays: blue arrows 2, red and green 0



13 / 57

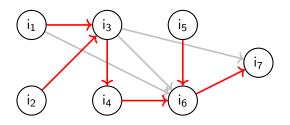
So many orders . . . with one critial path



 $i_1, i_2, i_3, i_4, i_5, i_6, i_7$ $i_1, i_2, i_5, i_3, i_4, i_6, i_7$ i₁,i₂,i₃,i₅,i₄,i₆,i₇ i₂,i₁,i₅,i₃,i₄,i₆,i₇ i₅,i₁,i₂,i₃,i₄,i₆,i₇ $i_2, i_1, i_3, i_5, i_4, i_6, i_7$ $i_1, i_5, i_2, i_3, i_4, i_6, i_7$ $i_5, i_2, i_1, i_3, i_4, i_6, i_7$

i₂,i₁,i₃,i₄,i₅,i₆,i₇ i₂,i₅,i₁,i₃,i₄,i₆,i₇

So many orders . . . with one critial path



 $i_1, i_2, i_3, i_4, i_5, i_6, i_7$ $i_1, i_2, i_5, i_3, i_4, i_6, i_7$ i₁,i₂,i₃,i₅,i₄,i₆,i₇ i₂,i₁,i₅,i₃,i₄,i₆,i₇ i₅,i₁,i₂,i₃,i₄,i₆,i₇ i₂,i₁,i₃,i₅,i₄,i₆,i₇ i₁,i₅,i₂,i₃,i₄,i₆,i₇ i₅,i₂,i₁,i₃,i₄,i₆,i₇ i₂,i₁,i₃,i₄,i₅,i₆,i₇ i₂,i₅,i₁,i₃,i₄,i₆,i₇

All these permutations respect dependencies but is there a best instruction scheduling?

Performances and Pipeline

Not all orders are equivalents!

- Some dependencies can bring hazards that slow down performances inside of the pipeline
- Hazard occurs when:
 - 1 instruction requires the previous instruction has finished
 - 2 instructions need the same data at the same time: one of the two is blocked

Instructions Pipeline

The microprocessor (MIPS) contains 5 stages:

- IF: Instruction Fetch
- ID: Instruction Decode. Read operands from registers, compute the address of the next instruction
- EX Execute instructions requiring the ALU
- ME Read/write into Memory
- WB Write Back. Results are written into registers.

	cycle ₁	cycle ₂	cycle ₃	cycle ₄	cycle ₅	cycle ₆	cycle ₇	cycle ₈	cycle ₉
instr ₁	IF	ID	EX	ME	WB				i
instr ₂		IF	ID	EX	ME	WB			
instr ₃		I	IF	ID	EX	ME	WB		;
instr ₄	 	! 		IF	ID	EX	ME	WB	1
instr ₅	1	l ı	l		IF	ID	EX	ME	WB

Hazard: RAW dependencies 1/2

Some instruction requires a result computed by a previous one!

Consider the following example:

	cycle ₁	cycle ₂	cycle ₃	cycle ₄	cycle ₅	cycle ₆	cycle ₇
lw \$2, 0(\$4)	IF	ID	EX	ME	WB	'	
addi \$5, \$2, 10		IF	ID		EX	ME	WB

- 1w produces its result into \$2 during the ME stage
- ADDI requires \$2 for the EX stage
- In this example, 1 stall (cycle 4)

The goal of risc architectures is to produce one per cycle!

Hazard: RAW dependencies 2/2

Consider now the following example:

	cycle ₁	cycle ₂	cycle ₃	cycle ₄	cycle ₅	cycle ₆	cycle ₇	cycle ₈	
lw \$2, 0(\$4)	IF	ID	EX	ME	WB				ı
addi \$5, \$2, 10		IF	ID		¥ΕΧ	ME	WB		1
add \$12, \$9, \$11			IF		ID	EX	ME	WB	

Hazard: RAW dependencies 2/2

Consider now the following example:

	cycle ₁	cycle ₂	cycle ₃	cycle ₄	cycle ₅	cycle ₆	cycle ₇	cycle ₈
lw \$2, 0(\$4)	IF	ID	EX	ME	WB			
addi \$5, \$2, 10		IF	ID		¥ΕΧ	ME	WB	
add \$12, \$9, \$11			IF		ID	EX	ME	WB

Let's look ... instruction 3 is independent from the others

18 / 57

Hazard: RAW dependencies 2/2

Consider now the following example:

	cycle ₁	cycle ₂	cycle ₃	cycle ₄	cycle ₅	cycle ₆	cycle ₇	cycle ₈
lw \$2, 0(\$4)	IF	ID	EX	ME	WB			
addi \$5, \$2, 10		IF	ID		¥ΕΧ	ME	WB	1
add \$12, \$9, \$11			IF		ID	EX	ME	WB

Let's look ... instruction 3 is independent from the others so we can change the order!

	cycle ₁	cycle ₂	cycle ₃	cycle ₄	cycle ₅	cycle ₆	cycle ₇	cycle ₈	
lw \$2, 0(\$4)	IF	ID	EX	MĘ	WB			' 	İ
add \$12, \$9, \$11		IF	ID	EX	ME	WB		l	1
addi \$5, \$2, 10		l	IF	ID	$ ho_{ m EX}$	ME	WB		1

Hazard: WAW dependencies

Two instructions write in the same register!

Consider the following example:

	cycle ₁	cycle ₂	cycle ₃	cycle ₄	cycle ₅	cycle ₆
addi \$5, \$11, 42	IF	ID	EX	ME	WB	'
addi \$5, \$2, 10		IF	ID	EX	ME	₩B

WAW do not produce stalls! (even when writing in the same memory address)

Hazard: WAR dependencies

One instruction writes where a previous one reads!

Consider the following example:



WAR do not produce stalls!

Back to the example – without scheduling

	c ₁	c ₂	С3	C4	C5	С6	C7	C8	Cg .	c ₁₀	c ₁₁	c ₁₂	c ₁₃
i_1	IF	ID	EX	ME	WB			! !		' 			
i ₂	ı	IF	ID	EX	$^{ m ME}$	WB		! !		l j	ı i		1
i ₃	l	ı	IF	ID		EX	$^{ m ME}$	WB		l I	l I		·
i ₄	l I	l	ı	IF		ID	EX	ME	WB				<u> </u>
i ₅	i	' 	I	1		IF	ID	EX	ME	WB			<u> </u>
i ₆	I	I	1	l I	1		IF	ID		EX	ME	WB	
i ₇			l		 -			IF		ID	EX	ME	WB

Back to the example – without scheduling

	c ₁	c ₂	С3	C4	С5	С6	C7	C8	Cg .	c ₁₀	c ₁₁	c ₁₂	C ₁₃
i_1	IF	ID	EX	ME	WB		! !	! !		' 			' '
i_2	I	IF	ID	EX	$^{ m ME}$	WB		! !		l j	ı		1
i ₃	l	i	IF	ID		EX	ME	WB		l I	l i		
i ₄	 	l		IF		ID	EX	ME	WB				' '
i ₅	I	! 	I			IF	ID	EX	ME	WB			<u>'</u> 1
i ₆	I	I	I	l i		Ī	IF	ID		EX	ME	WB	ļ
i ₇	l	l	l			l	1	IF		ID	EX	ME	WB
	 	 	 	l I		 	IF						V

Without scheduling: 2 dependencies, 2 stalls, 13 cycles!

Minimizing Stalls - First approach

Each time we emit the next instruction, we should try to choose one which

- P₁ does not conflict with the previous emitted instruction
- P₂: is most likely to conflict if first of a pair (e.g. prefer lw to add)
- P₃: is as far away as possible (along paths in the DAG) from an instruction which can validly be scheduled last

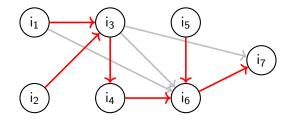
Minimizing Stalls - First approach

Each time we emit the next instruction, we should try to choose one which

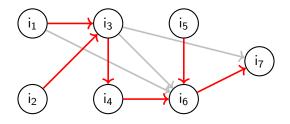
- P₁ does not conflict with the previous emitted instruction
- P2: is most likely to conflict if first of a pair (e.g. prefer lw to add)
- P₃: is as far away as possible (along paths in the DAG) from an instruction which can validly be scheduled last

Algorithm:

- Compute the dependency graph
- While the list of candidate instructions is not empty
 - If one instruction satisfies P₁, P₂, and P₃: remove it from the list and emit it.
 - * Remove the instruction from the DAG and insert the newly minimal elements into the candidate list.
 - Otherwise emit a nop instruction



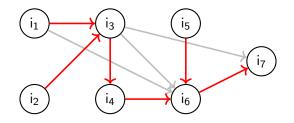
Candidates = $\{i_1, i_2, i_5\}$ Final Order =



Candidates = $\{i_1, i_2, i_5\}$ Final Order =

Choose i_1 since it satisfies P_1 , P_2 and P_3

```
i_1: lw $1,0($10) | i_4: sw $3,12($10) | i_7: sw $3,16($10) | i_2: lw $2,4($10) | i_5: lw $4,8($10) | i_5: add $3,$1,$2 | i_6: add $3,$1,$4
```

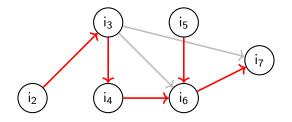


```
Candidates = \{i_1, i_2, i_5\}
Final Order = i_1
```

Choose i_1 since it satisfies P_1 , P_2 and P_3

10.40.45.45. 5 20.0

```
i_1: lw $1,0($10) | i_4: sw $3,12($10) | i_7: sw $3,16($10) | i_2: lw $2,4($10) | i_5: lw $4,8($10) | i_5: add $3,$1,$2 | i_6: add $3,$1,$4
```

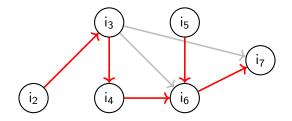


```
Candidates = \{i_1, i_2, i_5\}
Final Order = i_1
```

Choose i_1 since it satisfies P_1 , P_2 and P_3

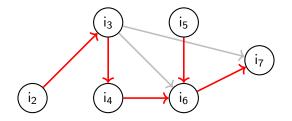
4 D > 4 A > 4 B > 4 B > 9 Q O

```
i_1: lw $1,0($10) | i_4: sw $3,12($10) | i_7: sw $3,16($10) | i_2: lw $2,4($10) | i_5: lw $4,8($10) | i_5: add $3,$1,$2 | i_6: add $3,$1,$4
```



```
Candidates = \{i_2, i_5\}
Final Order = i_1
```

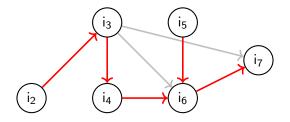
```
i_1: lw $1,0($10) | i_4: sw $3,12($10) | i_7: sw $3,16($10) | i_2: lw $2,4($10) | i_5: lw $4,8($10) | i_5: add $3,$1,$2 | i_6: add $3,$1,$4
```



```
Candidates = \{i_2, i_5\}
Final Order = i_1
```

Choose i_2 since it satisfies P_1 , P_2 and P_3

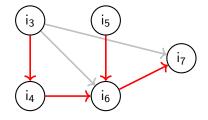
40 40 40 40 40 40 40 40 40



```
Candidates = \{i_2, i_5\}
Final Order = i_1, i_2
```

Choose i_2 since it satisfies P_1 , P_2 and P_3

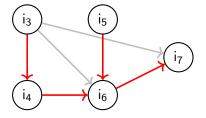
4 D > 4 A > 4 B > 4 B > B 900



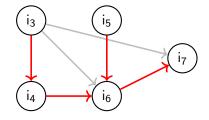
```
Candidates = \{i_2, i_5\}
Final Order = i_1, i_2
```

Choose i_2 since it satisfies P_1 , P_2 and P_3

4 D > 4 A > 4 B > 4 B > B 900



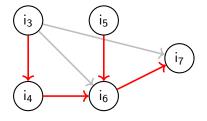
```
Candidates = \{i_5, i_3\}
Final Order = i_1, i_2
```



```
Candidates = \{i_5, i_3\}
Final Order = i_1, i_2
```

Choose i_5 since it satisfies P_1 , P_2 and P_3

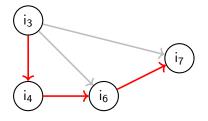
40 40 40 40 40 40 40 40 40



```
Candidates = \{i_5, i_3\}
Final Order = i_1, i_2, i_5
```

Choose i_5 since it satisfies P_1 , P_2 and P_3

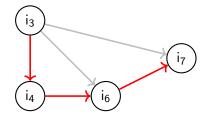
4 D > 4 A > 4 B > 4 B > B 900



```
Candidates = \{i_5, i_3\}
Final Order = i_1, i_2, i_5
```

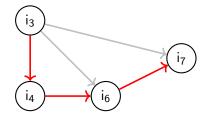
Choose i_5 since it satisfies P_1 , P_2 and P_3

4 D > 4 A > 4 B > 4 B > 9 Q O



```
Candidates = \{i_3\}
Final Order = i_1, i_2, i_5
```

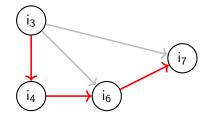
26 / 57



```
Candidates = \{i_3\}
Final Order = i_1, i_2, i_5
```

Choose i_3 since it satisfies P_1 , P_2 and P_3

40.40.43.43. 3 900

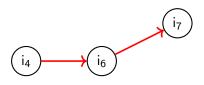


```
Candidates = \{i_3\}
Final Order = i_1, i_2, i_5, i_3
```

Choose i_3 since it satisfies P_1 , P_2 and P_3

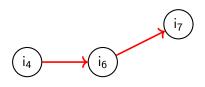
4 D > 4 A > 4 B > 4 B > B 900

```
i_1: lw $1,0($10) | i_4: sw $3,12($10) | i_7: sw $3,16($10) | i_2: lw $2,4($10) | i_5: lw $4,8($10) | i_5: add $3,$1,$2 | i_6: add $3,$1,$4
```

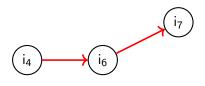


```
Candidates = \{i_3\}
Final Order = i_1, i_2, i_5, i_3
```

Choose i_3 since it satisfies P_1 , P_2 and P_3



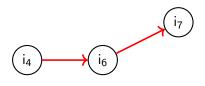
```
Candidates = \{i_4\}
Final Order = i_1, i_2, i_5, i_3
```



```
Candidates = \{i_4\}
Final Order = i_1, i_2, i_5, i_3
```

Choose i_4 since it satisfies P_1 , P_2 and P_3

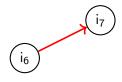
```
i_1: lw $1,0($10) | i_4: sw $3,12($10) | i_7: sw $3,16($10) | i_2: lw $2,4($10) | i_5: lw $4,8($10) | i_5: add $3,$1,$2 | i_6: add $3,$1,$4
```



```
Candidates = \{i_4\}
Final Order = i_1, i_2, i_5, i_3, i_4
```

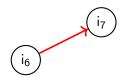
Choose i_4 since it satisfies P_1 , P_2 and P_3

4 D > 4 P > 4 B > 4 B > 9 Q C

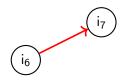


```
Candidates = \{i_4\}
Final Order = i_1, i_2, i_5, i_3, i_4
```

Choose i_4 since it satisfies P_1 , P_2 and P_3



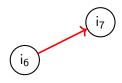
```
Candidates = \{i_6\}
Final Order = i_1, i_2, i_5, i_3, i_4
```



```
Candidates = \{i_6\}
Final Order = i_1, i_2, i_5, i_3, i_4
```

Choose i_6 since it satisfies P_1 , P_2 and P_3

4 D > 4 P > 4 B > 4 B > B = 900



```
Candidates = \{i_6\}
Final Order = i_1, i_2, i_5, i_3, i_4, i_6
```

Choose i_6 since it satisfies P_1 , P_2 and P_3

4 D > 4 A > 4 B > 4 B > 9 Q O



```
Candidates = \{i_6\}
Final Order = i_1, i_2, i_5, i_3, i_4, i_6
```

Choose i_6 since it satisfies P_1 , P_2 and P_3

◆□▶ ◆□▶ ◆■▶ ◆■▶ ● ◆9.0°

 (i_7)

```
Candidates = \{i_7\}
Final Order = i_1, i_2, i_5, i_3, i_4, i_6
```

 (i_7)

```
Candidates = \{i_7\}
Final Order = i_1, i_2, i_5, i_3, i_4, i_6
```

Choose i_7 since it satisfies P_1 , P_2 and P_3

 (i_7)

```
Candidates = \{i_7\}
Final Order = i_1, i_2, i_5, i_3, i_4, i_6, i_7
```

Choose i_7 since it satisfies P_1 , P_2 and P_3

```
Candidates = \{i_7\}
Final Order = i_1, i_2, i_5, i_3, i_4, i_6, i_7
```

Choose i_7 since it satisfies P_1 , P_2 and P_3

Final Order = i_1 , i_2 , i_5 , i_3 , i_4 , i_6 , i_7

	c_1	c_2	c ₃	c ₄	c ₅	c ₆	c ₇	c ₈	c ₉	c ₁₀	c ₁₁	
i_1	IF	ID	EX	ME	WB				l I	! 	' '	ı
i_2		IF	ID	EX	ME	WB			I	ı		I
i ₅	ļ i		IF	ID	EX	ME	WB		1	l		1
i_3			l	IF	ID	EX	ME	WB		l		ı
i ₄			l	1	IF	ID	EX	ME	WB		<u> </u>	I
i ₆	۱ ۱		l	ı	ı	IF	ID	EX	$^{ m ME}$	WB		1
i ₇			 -				IF	ID	EX	ME	WB	

```
i_1: lw $1,0($10) | i_4: sw $3,12($10) | i_7: sw $3,16($10)
i_2: lw $2,4($10) | i_5: lw $4,8($10)
i_3: add $3,$1,$2 | i_6: add $3,$1,$4
```

Final Order = i_1 , i_2 , i_5 , i_3 , i_4 , i_6 , i_7

	c ₁	c ₂	c ₃	c ₄	c ₅	c ₆	c ₇	c ₈	c ₉	c ₁₀	c ₁₁
i_1	IF	ID	EX	ME	WB				' 		' '
i ₂		IF	ID	EX	ME	WB			I	i	1
i ₅			IF	ID	EX	ME	WB		 	1	
i ₃	l I		l	IF	ID	EX	ME	WB			<u>'</u>
i ₄			l	ī	IF	ID	EX	ME	WB		1
i ₆	l		l	I	ı	IF	ID	EX	ME	WB	1
i ₇			 -	l			IF	ID	EX	ME	WB

With scheduling: still 2 dependencies but 0 stalls and 11 cycles!

A word on scheduling strategies

- Sometimes we cannot avoid some stalls
- Computing the critical path can be smarter:
 - ▶ Rather than attributing 1 as weight to every instruction, we can adjust according to the real time of executing the instruction
 - We can take advantages of the number of successors
 - ... many yet-to-be-define heuristics!
- Computing the DAG of dependencies can be done in $O(n^2)$ by scanning backwards through the basic block and adding edges as dependencies arise

31 / 57

A word on performances

We can statically compute instructions per cycle IPC= $\frac{\text{nb instructions}}{\text{nb cycles}}$, to evaluate 2 possible scheduling.

In the previous example:

- without scheduling IPC= $\frac{7}{13}$ = 0.53
- with scheduling IPC= $\frac{7}{11}$ = 0.63 (better!)

We can also statically compute cycle per instructions: $CPI = \frac{1}{IPC}$.

The CPI lower bound is $\frac{\sum \alpha \times \beta}{\text{nb instructions}}$, avec α is the number of instructions for a given instruction type and β the associated cost.

Can we do better?

Consider the following code (representing a basic block):

Can we do better?

Consider the following code (representing a basic block):

	c ₁	C2	С3	C4	C5	c6	C7	. C8	C9	C10	C11	C12	C13	C14	C ₁₅	C16	
i ₁	IF	ID	EX	ME	WB					ļ				'			
i ₂		IF	ID		EX	ME	WB]						l			
i ₃	l		IF					ID	EX	ME	WB			ı		1	
i4	ı	I	i					IF	ID	EX	ME	WB] ,				
i ₅							1		IF				ID	EX	ME	WB	

Can we do better?

Consider the following code (representing a basic block):

	c ₁	c ₂	С3	C4	C5	С6	C7	С8	C9	C ₁₀	c ₁₁	c ₁₂	c ₁₃	C ₁₄	C ₁₅	c ₁₆
i ₁	IF	ID	EX	ME	WB				'		1		1	'	ı	
i ₂		IF	ID		EX	ME	WB			ı			I	l	l	1
iз	I		IF					ID	EX	ME	WB			ı	l	1 1
i4	l .	I	Ī					IF	ID	EX	ME	WB	1		ı	
i ₅									IF				ID	EX	ME	WB
-				'	'	'	'									

16 cycles for 5 instructions that are all dependent!

$$IPC = 0.31$$



Loop Unrolling

- Replicate loop body to expose more parallelism
- Reduces loop-control overhead

Loop Unrolling

- Replicate loop body to expose more parallelism
- Reduces loop-control overhead

At high level, it can be seen as following:

Without Loop Unrolling	With Loop Unrolling
int i;	int i;
for $(i = 0; i < 100; ++i)$	for $(i = 0; i < 100; i+=5)$
for $(i = 0; i < 100; ++i)$ tab[i] = tab[i] +42;	for (i = 0; i < 100; i+=5) tab[i] = tab[i] +42;
	tab[i+1] = tab[i+1] + 42;
	tab[i+2] = tab[i+2] + 42;
	tab[i+3] = tab[i+3] + 42;
	tab[i+4] = tab[i+4] + 42;

34 / 57

Loop Unrolling

- Replicate loop body to expose more parallelism
- Reduces loop-control overhead

At high level, it can be seen as following:

Without Loop Unrolling	With Loop Unrolling
int i;	int i;
for $(i = 0; i < 100; ++i)$	for $(i = 0; i < 100; i+=5)$
tab[i] = tab[i] +42;	tab[i] = tab[i] +42;
	tab[i+1] = tab[i+1] + 42;
	tab[i+2] = tab[i+2] + 42;
	tab[i+3] = tab[i+3] + 42;
	tab[i+4] = tab[i+4] + 42;

Special care must be taken for pre and post loops operations (as well as intra-loop dependencies)

```
٦w
                   $t0, 0($s1)
                                   # t0=array element
    Loop:
                   $t0, $t0, $s2 # add scalar in s2
            addu
                   $t0, 0($s1)
                                   # store result
            SW
i<sub>4</sub>:
            addi
                   $s1, $s1,-4 # decrement pointer
                   $s1, $0, Loop # branch s1!=0
           bne
           lw
                   $t0, 0($s1)
                                   # t0=array element
    Loop:
                   $t0, $t0, $s2 # add scalar in s2
            addu
                   $t0, 0($s1)
                                   # store result
            SW
            addi
                   $s1, $s1,-4
                                   # decrement pointer
                   $s1, $0, Loop # branch s1!=0
            bne
                   $t0, 0($s1) # t0=array element
    Loop:
           lw
                   $t0, $t0, $s2 # add scalar in s2
            addu
                   $t0, 0($s1) # store result
            SW
i<sub>14</sub>:
            addi
                   $s1, $s1,-4 # decrement pointer
                   $s1, $0, Loop
                                   # branch s1!=0
            bne
```

First duplicate N times the body of the loop!

35 / 57

```
Loop:
           lw
                   $t0, 0($s1) # t0=array element
                   $t0, $t0, $s2 # add scalar in s2
            addu
            SW
                   $t0, 0($s1) # store result
            addi
                   $s1, $s1,-4 # decrement pointer
                   $t0, 0($s1)
                                   # t0=array element
           ٦w
                   $t0, $t0, $s2 # add scalar in s2
            addıı
                   $t0, 0($s1)
                                   # store result
            SW
                   $s1, $s1,-4
                                   # decrement pointer
           addi
           lw
                   $t0, 0($s1) # t0=array element
                   $t0, $t0, $s2 # add scalar in s2
           addu
                   $t0, 0($s1) # store result
            SW
i<sub>14</sub>:
                   $s1, $s1,-4 # decrement pointer
            addi
                   $s1, $0, Loop # branch s1!=0
           bne
```

Remove redundant labels and jump (by supposing that we are able to do it!)

```
Loop:
           lw
                   $t0, 0($s1) # t0=array element
                   $t0, $t0, $s2 # add scalar in s2
           addu
                   $t0, 0($s1) # store result
           SW
i<sub>4</sub>:
           addi
                   $s1, $s1,-4 # decrement pointer
                   $t1, 0($s1) # t0=array element
           lw
           addu
                   $t1, $t1, $s2 # add scalar in s2
                   $t1, 0($s1)
                                   # store result
           SW
           addi
                   $s1, $s1,-4 # decrement pointer
           ٦w
                   $t2, 0($s1) # t0=array element
                   $t2, $t2, $s2 # add scalar in s2
           addu
           SW
                   $t2, 0($s1) # store result
i<sub>14</sub>:
           addi
                   $s1, $s1,-4 # decrement pointer
                   $s1, $0, Loop # branch s1!=0
           bne
```

Use other temporaries name when possible!

```
$s1, $s1,-12 # decrement pointer
i4:
    Loop:
           addi
                  $t0, 0($s1) # t0=array element
           lw
                  $t0, $t0, $s2 # add scalar in s2
           addu
                  $t0, 0($s1) # store result
           SW
           ٦w
                  $t1, 4($s1) # t0=array element
                  $t1, $t1, $s2 # add scalar in s2
           addıı
                  $t1, 4($s1) # store result
           SW
                  $t2, 8($s1) # t0=array element
           lw
           addıı
                  $t2, $t2, $s2 # add scalar in s2
                  t2, t2, t3 t4 store result
           SW
                  $s1, $0, Loop # branch s1!=0
           bne
```

Grab redundant operation and merge them carefully!

38 / 57

```
\$s1, \$s1, -12 # decrement pointer for N=3
Loop:
       addi
              $t0, 0($s1)
                              # t0=array element
       lw
       lw
              $t1, 4($s1) # t1=array element
       lw
              $t2, 8($s1) # t2=array element
              $t0, $t0, $s2 # add scalar in s2
       addu
              $t1, $t1, $s2 # add scalar in s2
       addu
              $t2, $t2, $s2 # add scalar in s2
       addu
              $t0, 0($s1)
                              # store result
       SW
              $t1, 4($s1) # store result
       SW
              $t2, 8($s1) # store result
       SW
              $s1, $0, Loop # branch s1!=0
       bne
```

Schedule the instructions and renumber them (and update comments)!

Pros & Cons

- We avoid a lot of conditional jumps (and many stall hence)
- We require 19 cycles for 11 instructions: IPC=0.57 (a lot better than the previous 0.31)
- This trick allows to have more independent instructions to insert, and thus, less stalls!
- But we have now a prologue and an epilogue: i.e., two more basic blocks
- Require more temporaries: register allocation will be harder!
- Try it by yourself in gcc -funroll-loops

A very last word on Branch Hazards 1/2

- Conditional jumps often introduce delays since we cannot pre-fetch instrcutions
 - ► Branch Outcome and Branch Target Address are ready at the end of the EX stage (3th stage)
 - Conditional branches are solved when PC is updated at the end of the ME stage (4th stage)
- Can we avoid them?

We only know i_{next} at cycle 5!

	c_1	c_2	с3	C4	c ₅	С6	c ₇	C8	Cg
bne \$1,\$2, loop	IF	ID	EX	ME	WB			' ' 	1
nop	l I	IF	ID	EX	ME	WB			1
nop	ı	ı	IF	ID	EX	ME	WB		i
nop	l I	' 	ı	IF	ID	EX	ME	WB	1
i _{next}	1	l			IF	ID	EX	ME	WB

41 / 57

A very last word on Branch Hazards 2/2

- X delayed slot: the X instructions after a branch are systematically executed
- The original SPARC and MIPS processors each used a single branch delay slot to eliminate single-cycle stalls after branches
- We need branch prediction... but nowadays, most of processors do it for us (and use slt...)!
- Some architectures have bypass between stages to avoid stalls

Avoid as possible floating points and jumps!

A very last word on Branch Hazards 2/2

- X delayed slot: the X instructions after a branch are systematically executed
- The original SPARC and MIPS processors each used a single branch delay slot to eliminate single-cycle stalls after branches
- We need branch prediction... but nowadays, most of processors do it for us (and use slt...)!
- Some architectures have bypass between stages to avoid stalls

Avoid as possible floating points and jumps!

"Do you program in mips?" she asked. "nop", he said.

Stalls due to caches

When the processor processor needs to access a data:

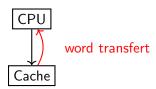
- If data is in cache: with a cost of 3 cycles
- Otherwise: with a cost of 100 cycles

Stalls due to caches

When the processor processor needs to access a data:

- If data is in cache: with a cost of 3 cycles
- Otherwise: with a cost of 100 cycles

CACHE HIT

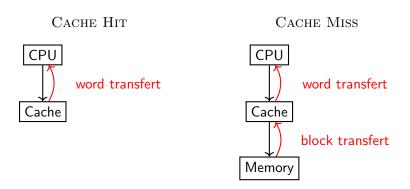


43 / 57

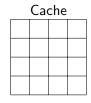
Stalls due to caches

When the processor processor needs to access a data:

- If data is in cache: with a cost of 3 cycles
- Otherwise: with a cost of 100 cycles



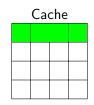
Cache Fundamentals 1/2



0×1	0×5	0×9	0×13	0×17	0x21

Memory

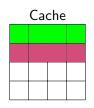
Cache Fundamentals 1/2

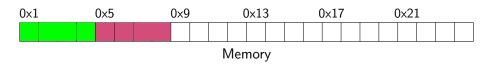


0×1	()x5		0×9				0×13			0×17			0×21				
							N	Лer	nor	у								

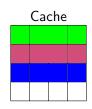
Access to adress 0x1, 4 words are fetched

Cache Fundamentals 1/2



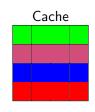


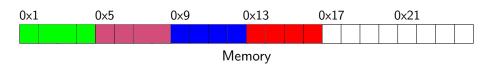
Access to adress 0x5, 4 words are fetched



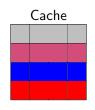


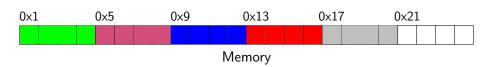
Access to adress 0x9, 4 words are fetched



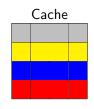


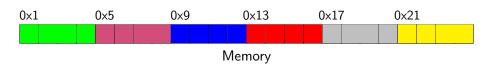
Access to adress 0x13, 4 words are fetched





Access to adress 0x17, 4 words are fetched First line of cache is replaced!





Access to adress 0x21, 4 words are fetched Second line of cache is replaced!

Many strategies to put data into the cache:

- Direct Mapping:
 - ► The address is decomposed in 3 parts: tag (8b), line (22b), and word(2b)
 - Each block of main memory maps to only one cache line, i.e. block-size = cache-line-size
 - ► Simple, Inexpensive, and fixed location for given block
- Associative Mapping:
 - ▶ A main memory block can load into any line of cache
 - Memory address is interpreted as tag and word
 - Tag uniquely identifies block of memory
 - Each block of main memory maps to only one cache line, i.e. block-size = cache-line-size
 - Complex, Expensive, and no-fixed location for given block

Prefetching

Fetch the data before it is needed (i.e. pre-fetch) by the program

- Eliminate cache misses
- Involves predicting which address will be needed in the future (as for branch prediction)
- In contrast to branch prediction:
 - incorrect prefetched data will simply not be used
 - there is no need for state recovery

Locality

- Locality is the principle that future memory accesses are near past accesses
- Memories take advantage of two types of locality
 - ► Temporal locality, i.e. near in time: we will often access the same data again very soon
 - Spatial locality, i.e. near in space/distance: our next access is often very close to our last access (or recent accesses)

Locality

- Locality is the principle that future memory accesses are near past accesses
- Memories take advantage of two types of locality
 - ► Temporal locality, i.e. near in time: we will often access the same data again very soon
 - Spatial locality, i.e. near in space/distance: our next access is often very close to our last access (or recent accesses)

Some Instruction Set Architecture (ISA) allows to pre-fetch some data: i.e., Humans or compilers has to insert (take advantage) of these instructions

Loops optimisations

We have already seen loops-unrolling to avoid stalls inside of the processor. Other techniques exist to avoid stalls due to cache:

- Loop Fission
- Loop interchanging
- Tabular Grouping
- Loop blocking
- Loop reversal
- Loop tiling
- . . .

48 / 57

Consider the following code, and direct mapping strategy:

```
int A[1024]; int B[1024]; int C[1024]; for (int i = 1; i<1024; ++i) {  A[i] = B[i]; \\ C[i] = C[i-1] + 1; \\ \}
```

49 / 57

Consider the following code, and direct mapping strategy:

```
int A[1024]; int B[1024]; int C[1024]; for (int i = 1; i<1024; ++i) {  A[i] = B[i]; \\ C[i] = C[i-1] + 1; \\ \}
```

Fetch A[i], A[i+1], A[i+2] and A[i+3]



Consider the following code, and direct mapping strategy:

```
int A[1024]; int B[1024]; int C[1024]; for (int i=1; i{<}1024; +{+}i) { A[i] = B[i]; \\ C[i] = C[i{-}1] + 1; }
```

Fetch B[i], B[i + 1], B[i + 2] and B[i + 3]



Consider the following code, and direct mapping strategy:

```
int A[1024]; int B[1024]; int C[1024]; for (int i = 1; i<1024; ++i) {  A[i] = B[i]; \\ C[i] = C[i-1] + 1;  }
```

Fetch C[i], C[i+1], C[i+2] and C[i+3]



Consider the following code, and direct mapping strategy:

```
int A[1024]; int B[1024]; int C[1024]; for (int i = 1; i<1024; ++i) {  A[i] = B[i]; \\ C[i] = C[i-1] + 1; \\ \}
```

Fetch C[i-1] will probably conflict



- Hopefully A[i], B[i] and C[i] will not conflict in the cache
- but ... C[i-1] will probably!

Solution

Divide the loop into two:

- Less pressure on cache
- We can now insert padding to avoid conflicts

```
int A[1024]; padding[xx]; int B[1024]; int C[1024]; for (int i=1; i<1024; ++i) A[i] = B[i]; for (int i=1; i<1024; ++i) C[i] = C[i-1] + 1;
```

Try it by yourself in gcc -ftree-loop-distribution

Loop interchanging 1/2

Consider the following code, and direct mapping cache:

```
int A[1024][1024];
for (int j = 1; j < 1024; ++j)
for (int i = 1; i < 1024; ++i)
A[j][i] = A[j][i] * 42;
```

In Fortran, the elements of an array are stored in memory contiguously by column, and the original loop iterates over rows, potentially creating at each access a cache miss

A B C
is stored A D B E C F

Loop interchanging 1/2

Consider the following code, and direct mapping cache:

```
int A[1024][1024];
for (int i = 1; i < 1024; ++i)
  for (int i = 1; i < 1024; ++i)
    A[i][i] = A[i][i] * 42;
```



In Fortran, the elements of an array are stored in memory contiguously by column, and the original loop iterates over rows, potentially creating at is stored

each access a cache miss

A D B E C F

Loop interchanging 1/2

Consider the following code, and direct mapping cache:

```
int A[1024][1024];
for (int j = 1; j < 1024; ++j)
for (int i = 1; i < 1024; ++i)
A[j][i] = A[j][i] * 42;
```

Fetch A[
$$j + 1$$
][i], A[$j + 2$][i], A[$j + 3$][i], and A[$j + 4$][i]



In Fortran, the elements of an array are stored in memory contiguously by column, and the original loop iterates over rows, potentially creating at each access a cache miss

A B C
is stored A D B E C F

Loop interchanging 2/2

Solution

This transformation switches the positions of one loop that is tightly nested within another loop.

```
int A[1024][1024];
for (int i = 1; i < 1024; ++i)
for (int j = 1; j < 1024; ++j)
A[j][i] = A[j][i] * 42;
```

Legal if the outermost loop does not carry any data dependence Try it by yourself in gcc -floop-interchange

Consider the following code, and direct mapping cache:

```
int A[1024]; int B[1024];
for (int j = 1; j < 1024; ++j)
A[j] = B[j] * 42;
```

Consider the following code, and direct mapping cache:

int A[1024]; int B[1024]; for (int
$$j = 1$$
; $j < 1024$; $++j$) $A[j] = B[j] * 42$;



Consider the following code, and direct mapping cache:

int A[1024]; int B[1024]; for (int
$$j = 1$$
; $j < 1024$; $++j$) $A[j] = B[j] * 42$;



Consider the following code, and direct mapping cache:

$$\begin{array}{l} \text{int A[1024]; int B[1024];} \\ \text{for (int } j=1; \ j{<}1024; \ +{+}j) \\ \text{A[}j] = \text{B[}j] * 42; \end{array}$$



Dynamic allocation does not allow padding. In the worst case, two miss per iterations

Solution

Group the two tabular into one

```
struct twoval {int A; int B};
struct twoval R[1024];
for (int j = 1; j < 1024; ++j)
R[j].A = R[j].B * 42;
```

Avoid a lot of caches miss!

Very hard for compiler to detect such cases

Loop Blocking

Consider the code below.

```
int A[1024][1024]; int B[1024][1024]; for (int i = 1; i < 1024; ++i) for (int j = 1; j < 1024; ++j) A[i][j] = B[i][j];
```

- If A and B are aligned we may encounter problems.
- Similar problems occur when processing images: A[i][j] = B[i-1][j-1] + B[i-1][j] + B[i-1][j+1] + B[i][j-1] + B[i][j] + B[i-1][j+1] + B[i+1][j] + B[i+1][j] + B[i+1][j+1];
- In this latter case, padding is complicated...

Loop Blocking

Solution

Try to work with data that fit in memory!

```
\begin{split} &\text{int A}[1024][1024]; \text{ int B}[1024][1024]; \\ &\text{for (int i = 1; i < 1024; i += B)} \\ &\text{for (int j = 1; j < 1024; j += B)} \\ &\text{for (int ii = 1; ii < min(1024, ii + B - 1); ii += B)} \\ &\text{for (int jj = 1; jj < min(1024, ii + B - 1); jj += B)} \\ &\text{A}[i][j] = B[i][j]; \end{split}
```

56 / 57

Summary

- stalls in the processor can come from many reasons
 - from data dependencies between instructions
 - from instruction dependencies
 - from cache and memory
- modern compiler hardly try to reduce them
 - by using Instruction Level Parallelism (i.e, to have a lot of independent instructions)
 - all these optimization must occur before register allocation (which is the final step)
 - When writing a compiler, you must know the target processor by heart!
- Caches can be shared between many processors!

Garbage Collection

Akim Demaille, Etienne Renault, Roland Levillain

June 4, 2019

TYLA

1 / 35

Table of contents

- Motivations and Definitions
- 2 Reference Counting Garbage Collection
- Mark and Sweep Garbage Collection
- 4 Stop and Copy Garbage Collection
- 6 Hybrid Approaches

TYLA

Garbage Collection 1/2

- Fisrt apparition in LISP, 1959, McCarthy
- Garbage collection is the automatic reclamation of computer storage (heap) at runtime
- Automatic memory management
 - New/malloc doesn't need delete/free anymore
 - Necessary for fully modular programming.
 Otherwise some modules are responsible for allocation while others are responsible for deallocation.
 - ► No more memory leaks
 - Avoid dangling-pointers/references.
 Reclaiming memory too soon is no more possible

Garbage Collection 2/2

- Quite expensive relative to explicit heap management
 - ► Slow running programs down by (very roughly) 10 percent...
 - ▶ ... But sometime cheaper or competitive
 - ► Fair comparison is difficult since explicit deallocation affects the structure of programs in ways that may themselves be expensive
- Possible reduction of heap fragmentation
- Functional and logic programming languages generally incorporate garbage collection because their unpredictable execution patterns
- D, Python, Caml, Effeil, Swift, C#, Go, Java, Haskell, LISP, Dylan, Prolog, etc.

TYLA

What is Garbage?

- An object is called garbage at some point during execution if it will never be used again.
- What is garbage at the indicated points?

```
int main() {
 Object x, y;
 x = new Object();
 y = new Object();
 /* Point A */
 x.doSomething();
 y.doSomething();
 /* Point B */
 y = new Object();
 /* Point C */
```

Approximating Garbage

- In general, it is undecidable whether an object is garbage
- An object is reachable if it can still be referenced by the program.

Goals

Detect and reclaim unreachable objects

Basics of a Garbage Collector

- Distinguishing the live objects from the garbage ones
- Reclaiming the garbage object' storage

TYLA Garbage Collection June 4, 2019 7 / 35

Basics of a Garbage Collector

- Objects from the garbage ones
- Reclaiming the garbage object' storage

We focus on built-in garbage collectors so that:

- allocation routines performs special actions
 - reclaim memory
 - emit specific code to recognize object format
 - etc.
- explicit calls to the deallocator are unnecessary
 - the allocator will call it on-time
 - the objects will be automatically destroyed

TYLA

Different kind of GC

- Incremental techniques:
 - allow garbage collection to proceed piecemeal while application is running
 - my provide real-time garantees
 - can be generalized into concurrent collections
- Generationnal Schemes
 - improve efficiency/locality by garbage collecting a smaller area more often
 - avoid overhead due to long time objects
 - rely on pause to collect data

Table of contents

- Motivations and Definitions
- 2 Reference Counting Garbage Collection
- Mark and Sweep Garbage Collection
- 4 Stop and Copy Garbage Collection
- 6 Hybrid Approaches

9 / 35

TYLA Garbage Collection June 4, 2019



Intuition

• Maintain for each object a counter to the references to this object

Intuition

- Maintain for each object a counter to the references to this object
- Each time a reference to the object is created, increase the pointed-to object's counter

Intuition

- Maintain for each object a counter to the references to this object
- Each time a reference to the object is created, increase the pointed-to object's counter
- Each time an existing reference to an object is eliminated, the counter is decremented

Intuition

- Maintain for each object a counter to the references to this object
- Each time a reference to the object is created, increase the pointed-to object's counter
- Each time an existing reference to an object is eliminated, the counter is decremented
- When the object counter equals zero, the memory can be reclaimed

Caution

When an object is destructed:

Transitive reclamation can be deferred by maintaining a list of freed objects

TYLA

Caution

When an object is destructed:

examines pointer fields

Transitive reclamation can be deferred by maintaining a list of freed objects

Caution

When an object is destructed:

- examines pointer fields
- for any references R contained by this object, decrement reference counter of R

Transitive reclamation can be deferred by maintaining a list of freed objects

TYLA

Caution

When an object is destructed:

- examines pointer fields
- for any references R contained by this object, decrement reference counter of R
- If the reference counter of R becomes 0, reclaim memory

Transitive reclamation can be deferred by maintaining a list of freed objects

TYLA

```
class LinkedList {
  LinkedList next = null;
}
int main() {
```

```
class LinkedList {
  LinkedList next = null;
}
int main() {
  LinkedList head = new LinkedList;
```

TYLA

```
head 1
```

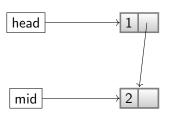
```
class LinkedList {
   LinkedList next = null;
}
int main() {
   LinkedList head = new LinkedList;
   LinkedList mid = new LinkedList;
```

```
head 1
```



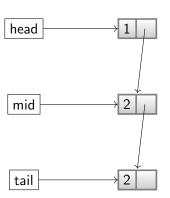
```
class LinkedList {
                                       head
  LinkedList next = null;
int main() {
LinkedList head = new LinkedList;
LinkedList mid = new LinkedList;
                                       mid
LinkedList tail = new LinkedList;
                                        tail
```

```
class LinkedList {
   LinkedList next = null;
}
int main() {
   LinkedList head = new LinkedList;
   LinkedList mid = new LinkedList;
   LinkedList tail = new LinkedList;
   head.next = mid;
```

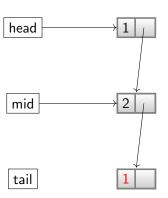




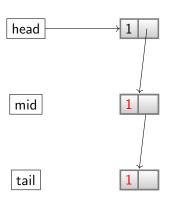
```
class LinkedList {
  LinkedList next = null;
int main() {
LinkedList head = new LinkedList;
LinkedList mid = new LinkedList;
LinkedList tail = new LinkedList;
head.next = mid;
mid.next = tail;
```



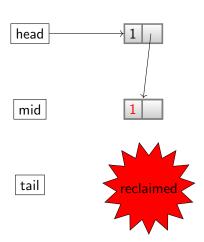
```
class LinkedList {
  LinkedList next = null;
int main() {
LinkedList head = new LinkedList;
LinkedList mid = new LinkedList;
LinkedList tail = new LinkedList;
head.next = mid;
mid.next = tail;
mid = tail = null;
```



```
class LinkedList {
  LinkedList next = null;
int main() {
LinkedList head = new LinkedList;
LinkedList mid = new LinkedList;
LinkedList tail = new LinkedList;
head.next = mid;
mid.next = tail;
mid = tail = null;
```



```
class LinkedList {
  LinkedList next = null;
int main() {
LinkedList head = new LinkedList;
LinkedList mid = new LinkedList;
LinkedList tail = new LinkedList;
head.next = mid;
mid.next = tail;
mid = tail = null;
head.next.next = null;
```



```
class LinkedList {
  LinkedList next = null;
}
int main() {
  LinkedList head = new LinkedList;
  LinkedList mid = new LinkedList;
  LinkedList tail = new LinkedList;
  head.next = mid;
  mid.next = tail;
  mid = tail = null;
  head.next.next = null;
  head = null;
  tail
```

```
class LinkedList {
                                        head
  LinkedList next = null;
int main() {
LinkedList head = new LinkedList;
LinkedList mid = new LinkedList;
                                        mid
LinkedList tail = new LinkedList;
head.next = mid;
mid.next = tail;
mid = tail = null;
head.next.next = null;
head = null;
```



tail

```
class LinkedList {
                                        head
  LinkedList next = null;
int main() {
LinkedList head = new LinkedList;
LinkedList mid = new LinkedList;
                                         mid
LinkedList tail = new LinkedList;
head.next = mid;
mid.next = tail;
mid = tail = null;
head.next.next = null;
                                         tail
head = null;
```

```
class LinkedList {
  LinkedList next = null;
int main() {
LinkedList head = new LinkedList;
LinkedList mid = new LinkedList;
LinkedList tail = new LinkedList;
head.next = mid;
mid.next = tail;
mid = tail = null;
head.next.next = null;
head = null;
```

head

mid



tail

```
class LinkedList {
   LinkedList next = null;
}
int main() {
```

If the objects create a directed cycle, the objects references counters will never reduced to zero.

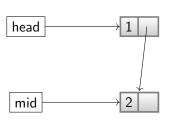
```
class LinkedList {
  LinkedList next = null;
}
int main() {
  LinkedList head = new LinkedList;
```

}

```
class LinkedList {
  LinkedList next = null;
                                       head
int main() {
LinkedList head = new LinkedList;
LinkedList mid = new LinkedList;
LinkedList tail = new LinkedList;
                                        mid
                                        tail
```

If the objects create a directed cycle, the objects references counters will never reduced to zero.

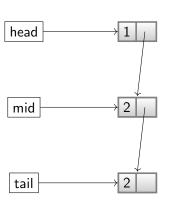
```
class LinkedList {
   LinkedList next = null;
}
int main() {
   LinkedList head = new LinkedList;
   LinkedList mid = new LinkedList;
   LinkedList tail = new LinkedList;
   head.next = mid;
```



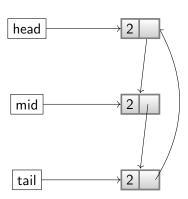


}

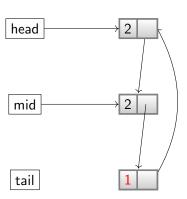
```
class LinkedList {
   LinkedList next = null;
}
int main() {
   LinkedList head = new LinkedList;
   LinkedList mid = new LinkedList;
   LinkedList tail = new LinkedList;
   head.next = mid;
   mid.next = tail;
```



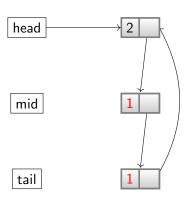
```
class LinkedList {
  LinkedList next = null;
int main() {
LinkedList head = new LinkedList;
LinkedList mid = new LinkedList;
LinkedList tail = new LinkedList;
head.next = mid;
mid.next = tail:
tail.next = head;
```



```
class LinkedList {
  LinkedList next = null;
int main() {
LinkedList head = new LinkedList;
LinkedList mid = new LinkedList;
LinkedList tail = new LinkedList:
head.next = mid;
mid.next = tail;
tail.next = head;
tail = null;
```



```
class LinkedList {
  LinkedList next = null;
int main() {
LinkedList head = new LinkedList:
LinkedList mid = new LinkedList;
LinkedList tail = new LinkedList:
head.next = mid;
mid.next = tail;
tail.next = head;
tail = null:
mid = null;
```



```
class LinkedList {
  LinkedList next = null;
                                        head
int main() {
LinkedList head = new LinkedList:
LinkedList mid = new LinkedList;
LinkedList tail = new LinkedList:
                                         mid
                                                           1
head.next = mid;
mid.next = tail;
tail.next = head;
tail = null:
mid = null;
                                         tail
head = null;
```

Pros and Cons

Pros:

- Easy to implement: perl, Firefox
- Can be implemented on top of explicit memory management librairies (shared_ptr)
- Interleaved with running time
- Small overage per unit of program execution
- Transitive reclamation can be deferred by maintaining a list of freed objects
- Real-time requierements: no halt of the system. Necessary for application where response-time is critical

Cons:

- A whole machine word per object
- When the number of references to an object overflows, the counter is set to the maximum and the memory will never be reclaimed
- Problem with cycles
- Efficiency: cost relative to the running program

Table of contents

- Motivations and Definitions
- 2 Reference Counting Garbage Collection
- Mark and Sweep Garbage Collection
- 4 Stop and Copy Garbage Collection
- 6 Hybrid Approaches

Analysis

• Reference counting tries to find unreachable objects by finding objects without incoming references

These references have been forgotten!



TYLA

16 / 35

Analysis

- Reference counting tries to find unreachable objects by finding objects without incoming references
- These references have been forgotten!

We have to trace the lifetime of objects

Intuition

Given knowledge of what's immediately accessible, find everything reachable in the program

The root set is the set of memory locations in the program that are known to be reachable

Graph Problem

Simply do a graph search starting at the root set:

- Any objects reachable from the root set are reachable
- Any objects not reachable from the root set are not reachable

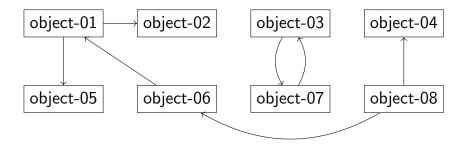
4□ > 4□ > 4□ > 4 = > = 90

How to obtain the root set?

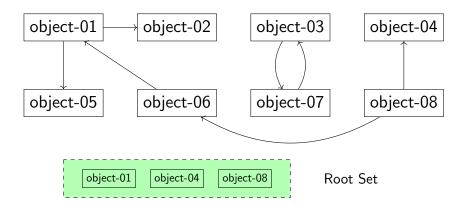
- static reference variables
- references registered through librairies (JNI, for instance)
- For each threads:
 - local variables
 - current method(s) arguments
 - stack
 - etc.

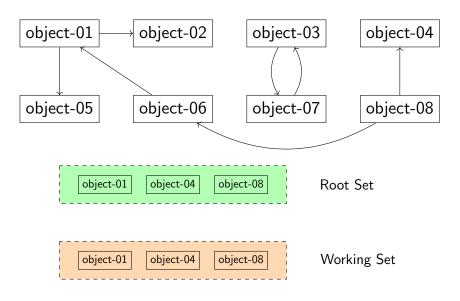
Mark-and-Sweep: the Algorithm

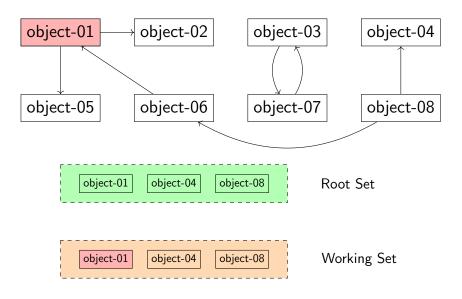
- Marking phase: Find reachable objects
 - Add the root set to a worklist
 - While the worklist isn't empty
 - ★ Remove an object from the worklist
 - If it is not marked, mark it and add to the worklist all objects reachable from that object
- Sweeping phase: Reclaim free memory
 - If that object isn't marked, reclaim its memory
 - If the object is marked, unmark it

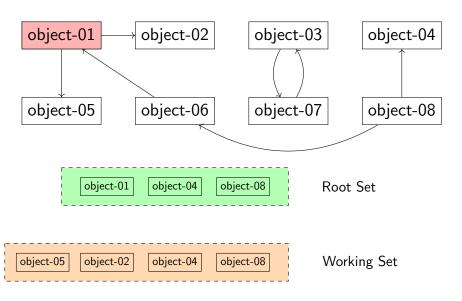


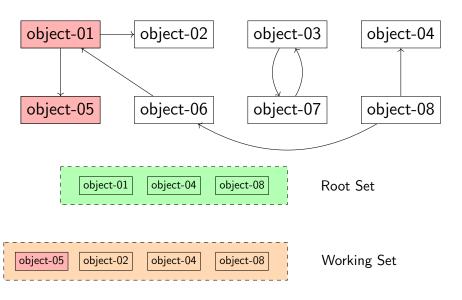
20 / 35

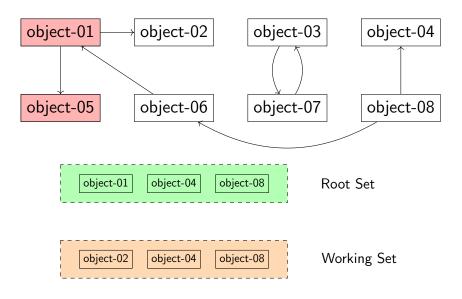


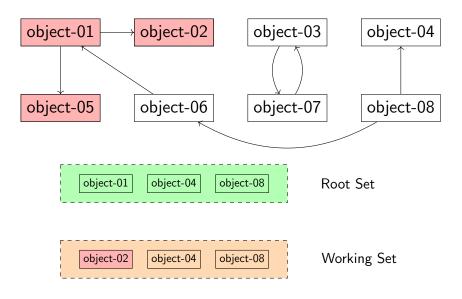


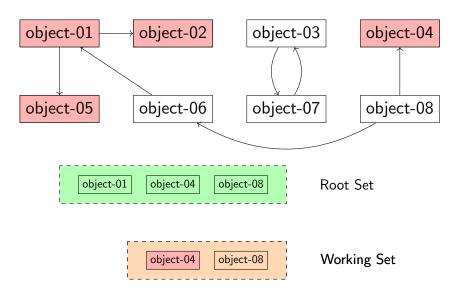


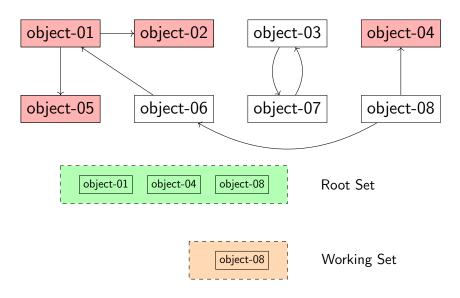


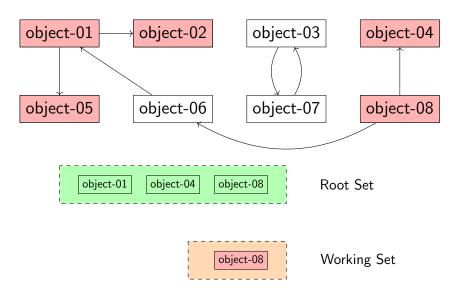


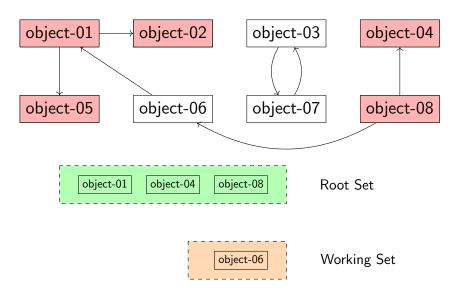


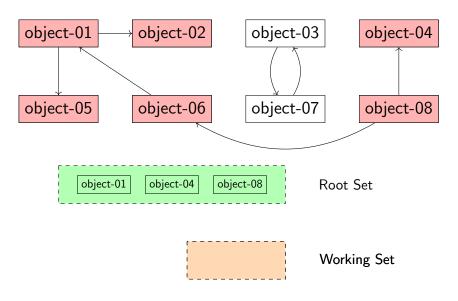


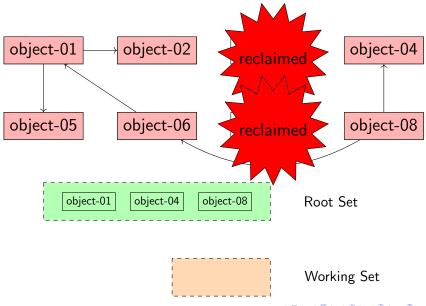












How to sweep?

Sweeping requires to know where are unreacheable objets !

Heap:

object-01
object-02
object-03
object-04
object-05
object-06
object-07
object-08

object-07 object-08

How to sweep?

Sweeping requires to know where are unreacheable objets !

Heap:

object-01 object-02 object-03 object-04 object-05 object-06 object-07

object-08

Just remove from the heap all non-marked objects

Problems

- Runtime proportional to number of allocated objects
 - Sweep phase visits all objects to free them or clear marks
- Work list requires lots of memory
 - ▶ Amount of space required could potentially be as large as all of memory
 - Can't preallocate this space

Pros and Cons

Pros:

- Can free cyclic references
- 1 bits per state
- Runtime can be proportional to the number of reachable objects (Baker's algorihtm)

Cons:

- Stop the world algorithm with possibly huge pauses times
- Memory Fragmentation
- Need to walk the whole heap

Table of contents

- Motivations and Definitions
- 2 Reference Counting Garbage Collection
- Mark and Sweep Garbage Collection
- 4 Stop and Copy Garbage Collection
- 6 Hybrid Approaches

Analysis

- Locality can be improved
 - After garbage collection, objects are no longer closed in memory
- Allocation speed can be improved
 - After garbage collection, the free list of the allocator must be walked.

The Sweep Phase can be improved

Zone 1										
Zone 2										

• Split memory in two pieces

Zone 1											
Zone 2											

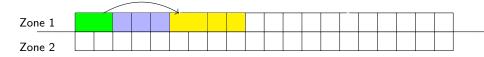
- Split memory in two pieces
- Allocate memory in the first zone

TYLA

Zone 1											
Zone 2											

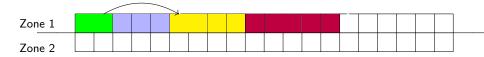
- Split memory in two pieces
- Allocate memory in the first zone





- Split memory in two pieces
- Allocate memory in the first zone

TYLA



- Split memory in two pieces
- Allocate memory in the first zone

TYLA

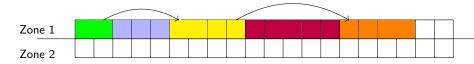


- Split memory in two pieces
- Allocate memory in the first zone



- Split memory in two pieces
- Allocate memory in the first zone
- When running out-of-space in the first zone: Garbage Collect!

TYLA



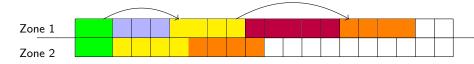
- Split memory in two pieces
- Allocate memory in the first zone
- When running out-of-space in the first zone: Garbage Collect!
- Explore only reachable references from the root set (here only green object)



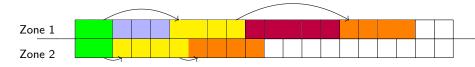
- Split memory in two pieces
- Allocate memory in the first zone
- When running out-of-space in the first zone: Garbage Collect!
- Explore only reachable references from the root set (here only green object)
- Copy objects



- Split memory in two pieces
- Allocate memory in the first zone
- When running out-of-space in the first zone: Garbage Collect!
- Explore only reachable references from the root set (here only green object)
- Copy objects



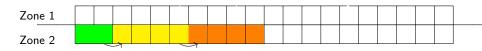
- Split memory in two pieces
- Allocate memory in the first zone
- When running out-of-space in the first zone: Garbage Collect!
- Explore only reachable references from the root set (here only green object)
- Copy objects



- Split memory in two pieces
- Allocate memory in the first zone
- When running out-of-space in the first zone: Garbage Collect!
- Explore only reachable references from the root set (here only green object)
- Copy objects
- Update References & Root set



- Split memory in two pieces
- Allocate memory in the first zone
- When running out-of-space in the first zone: Garbage Collect!
- Explore only reachable references from the root set (here only green object)
- Copy objects
- Update References & Root set
- Clean zone 1 (Constant time)

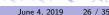


- Split memory in two pieces
- Allocate memory in the first zone
- When running out-of-space in the first zone: Garbage Collect!
- Explore only reachable references from the root set (here only green object)
- Copy objects
- Update References & Root set
- Clean zone 1 (Constant time)
- Swap zone 1 and 2 (Now allocation will happen in zone 2)

→ロト →同ト → 三ト → 三 → りゅべ



- Split memory in two pieces
- Allocate memory in the first zone
- When running out-of-space in the first zone: Garbage Collect!
- Explore only reachable references from the root set (here only green object)
- Copy objects
- Update References & Root set
- Clean zone 1 (Constant time)
- Swap zone 1 and 2 (Now allocation will happen in zone 2)
- Allocate the object that have provoqued the GC



Implementation

- Partition memory into two regions: the old space and the new space.
- Keep track of the next free address in the new space.
- To allocate n bytes of memory:
- If n bytes space exist at the free space pointer, use those bytes and advance the pointer.
- Otherwise, do a copy step. To execute a copy step:
- For each object in the root set:
 - Copy that object over to the start of the old space.
 - Recursively copy over all objects reachable from that object.
- Adjust the pointers in the old space and root set to point to new locations.
- Exchange the roles of the old and new spaces.

Problems

How to adjust pointers in the copied objects correctly?



TYLA

Problems

How to adjust pointers in the copied objects correctly?

- Have each object contain a extra space for a forwarding pointer
- First, do a complete bitwise copy of the object
- Next, set the forwarding pointer of the original object to point to the new object
 - Follow the pointer to the object it references
 - Replace the pointer with the pointee's forwarding pointer

Pros and Cons

Pros:

- Compact the Heap
- Allocation only increments a pointer
- No sweep

Cons:

- Smaller Heap
- Copy
- Reference adjusting

Table of contents

- Motivations and Definitions
- 2 Reference Counting Garbage Collection
- Mark and Sweep Garbage Collection
- 4 Stop and Copy Garbage Collection
- 6 Hybrid Approaches

Analysis

The best garbage collectors in use today are based on a combination of smaller garbage collectors

Objects Die Young

Most objects have extremely short lifetimes

Optimize garbage collection to reclaim young objects rapidly while spending less time on older objects

Generational Garbage Collector

- Partition memory into several generations
- Objects are always allocated in the first generation.
- When the first generation fills up, garbage collect it.
 - Runs quickly; collects only a small region of memory.
- Move objects that survive in the first generation long enough into the next generation.
- When no space can be found, run a full (slower) garbage collection on all of memory.

Garbage Collection in Java

- Split the Heap in 3 zones: eden, survivors and tenured
- New objects are allocated using a modified stop-and-copy collector in the Eden space.
- When Eden runs out of space, the stop-and-copy collector moves its elements to the survivor space.
- Objects that survive long enough in the survivor space become tenured and are moved to the tenured space.
- When memory fills up, a full garbage collection (perhaps mark-and-sweep) is used to garbage-collect the tenured objects

Garbage Collection in C

- Boehm GC
- Mark and Sweep
- Conservative
- Consider all program variables as root set
- ullet Easy to combine with C

Bibliography

• Uniprocessor Garbage Collection , Paul R. Wilson

35 / 35

YLA Garbage Collection June 4, 2019