

# Register Allocation and Instruction Scheduling in Unison

**Roberto Castañeda Lozano – SICS, KTH**

joint work with:

G. Hjort Blindell – KTH, SICS

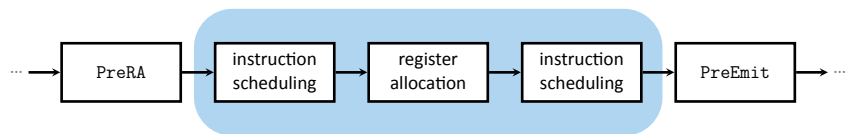
M. Carlsson – SICS

F. Drejhammar – SICS

C. Schulte – KTH, SICS

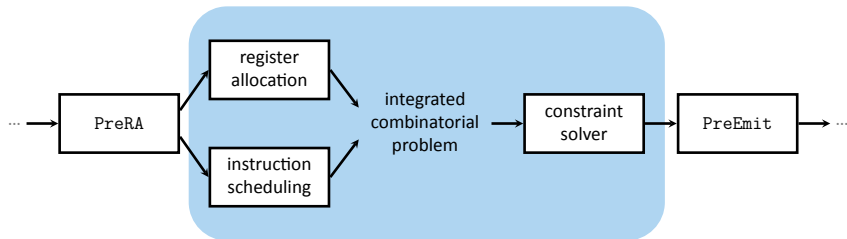


# Code Generation in LLVM



- Stages, heuristics
- Pros: compilation speed
- Cons: suboptimal, complex

# Introducing Unison



- Integration, combinatorial optimization
- Pros: simple, optimal
- Cons: compilation slowdown

# Unison Is Practical and Effective

- For LLVM Users
  - traditional LLVM for compile/debug cycle
  - LLVM + Unison for release builds
- For LLVM developers
  - evaluation of heuristics
  - identification of improvement opportunities

# 1 Optimal Approaches

2 Model

3 Results

4 Case Studies

5 Conclusion

# Earlier Optimal Approaches

- Global register allocation  
Local instruction scheduling
  - practical and effective
- Global instruction scheduling
  - scales up to medium-size problems
- Integrated optimal approaches
  - ignore essential register allocation subproblems
  - do not scale beyond small problems

# Integrated Optimal Approaches

approach	GL	register allocation								instr. sched.			max. size
		SP	RA	CO	SO	RP	LS	RM	MB	BD	MU	2D	
Wilson 1994	✓	✓	✓	✓	-	-	✓	-	-	✓	✓	-	30
Chang 1997	-	✓	-	-	✓	-	-	-	-	✓	✓	-	~ 10
Gebotys 1997	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	108
ICG 1999	-	✓	✓	✓	✓	-	✓	-	✓	✓	-	-	23
PROPAN 2000	✓	-	✓	-	-	-	-	-	✓	✓	✓	✓	42
Nagar. 2007	-	✓	✓	-	✓	-	✓	-	-	✓	-	-	?
OPTIMIST 2012	-	✓	-	-	✓	-	✓	-	✓	✓	✓	✓	100
<b>Unison</b>	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	605

- Few global approaches
- Few register allocation subproblems
- Low scalability
- Unison: All subproblems, better scalability
  - key: constraint programming

1 Optimal Approaches

**2 Model**

3 Results

4 Case Studies

5 Conclusion



# Model

- Register allocation
  - allocate temps to registers or memory
  - **register assignment**
  - spilling
  - coalescing
  - live range splitting
  - ...
- Instruction scheduling
- Connection: temp live ranges

# Register Assignment as Rectangle Packing

## Register Assignment

temp live ranges

temp size (16 bits, 32 bits, ...)

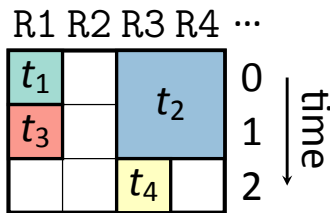
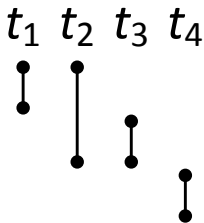
interfering temps cannot share registers

## Rectangle Packing

rectangles

rectangle width

rectangles cannot overlap



- Model based on [Pereira and Palsberg, 2008](#)

no-overlap( $\langle r_{t_1}, r_{t_1} + 1, ls_{t_1}, le_{t_1} \rangle$ ,  $\langle r_{t_2}, r_{t_2} + 2, ls_{t_2}, le_{t_2} \rangle$ ,  $\dots$ )

1 Optimal Approaches

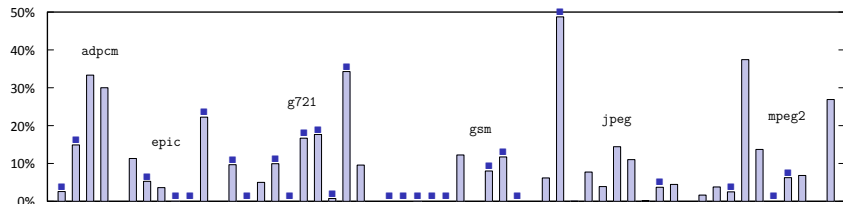
2 Model

**3 Results**

4 Case Studies

5 Conclusion

# Speedup over LLVM 3.8



- 50 MediaBench functions
- Hexagon V4 processor
- Provably optimal (■) for 54% of the functions
- Compilation time: from seconds to minutes

1 Optimal Approaches

2 Model

3 Results

**4 Case Studies**

5 Conclusion

# Disclaimer

The case studies are generated with LLVM 3.8 for  
Hexagon V4

# Disclaimer

opt is run with different optimization levels and  
some llc passes are disabled for simplicity  
(llc is always run with O3)

# Disclaimer

I am no expert on LLVM – just a humble user

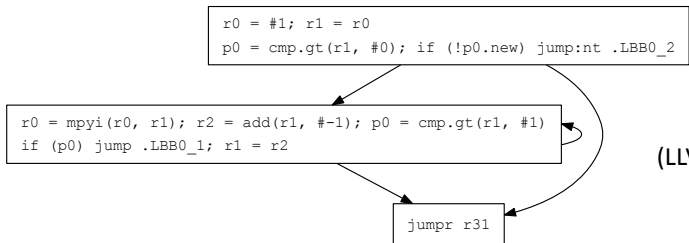


# Case Study: fac

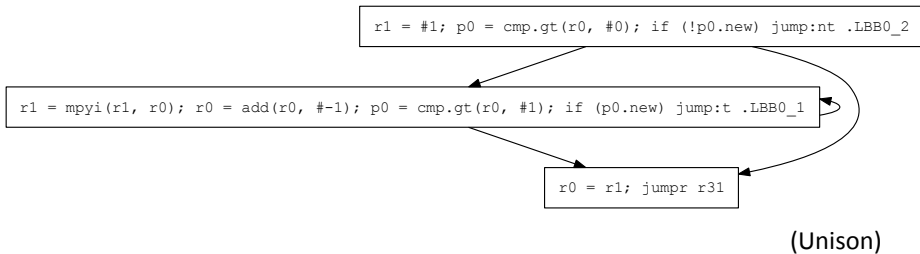
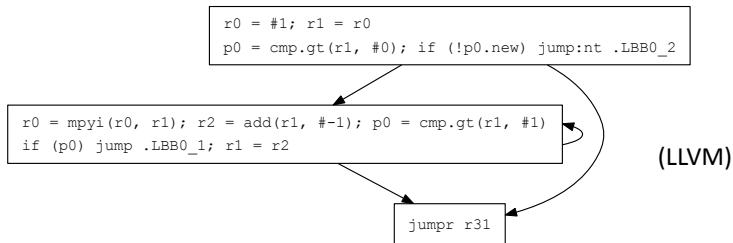
- Simple iterative factorial:

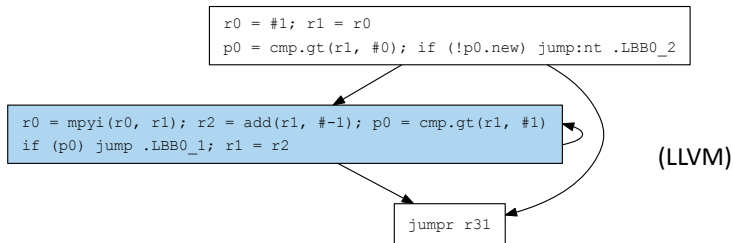
```
int fac(int n) {
    int f = 1;
    while (n > 0) {
        f = f * n;
        n--;
    }
    return f;
}
```

- Exposes opportunity for better coalescing
- Illustrates effect of integrated reasoning

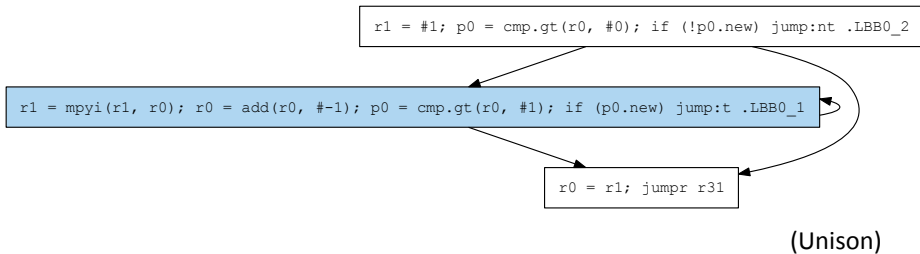


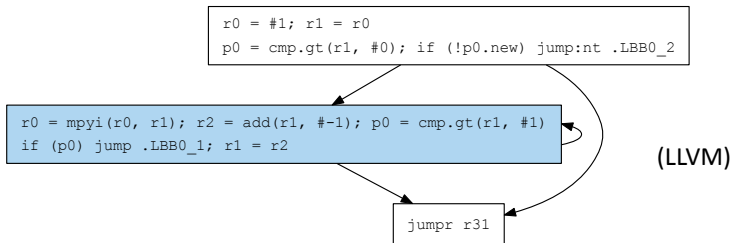
(LLVM)





LLVM's loop is twice as slow, why?





# After Simple Register Coalescing:

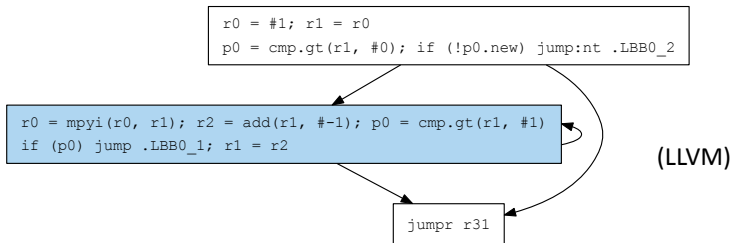
..

BB#1:

```

%vreg2 = A2_addi %vreg10, -1
%vreg9 = M2_mpyi %vreg9, %vreg10
%vreg8 = C2_cmpgti %vreg10, 1
%vreg10 = COPY %vreg2
J2_jumpr %vreg8, <BB#1>
  
```

..



# After Simple Register Coalescing:

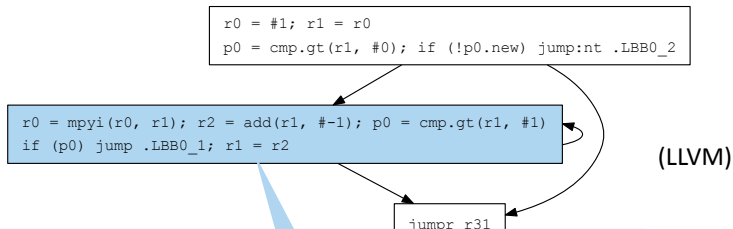
..

BB#1:

```
%vreg2 = A2_addi %vreg10, -1
%vreg9 = M2_mpyi %vreg9, %vreg10
%vreg8 = C2_cmpgti %vreg10, 1
%vreg10 = COPY %vreg2
J2_jumpr %vreg8, <BB#1>
```

..

**%vreg2 and %vreg10 not coalesced**



corresponding move requires a new bundle

# After Simple Register Coalescing:

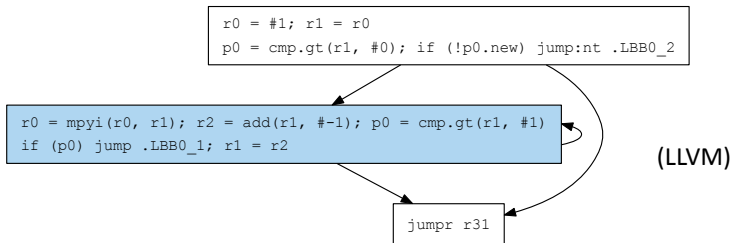
..

BB#1:

```

%vreg2 = A2_addi %vreg10, -1
%vreg9 = M2_mpyi %vreg9, %vreg10
%vreg8 = C2_cmpgti %vreg10, 1
%vreg10 = COPY %vreg2
J2_jumpr %vreg8, <BB#1>
  
```

..



# After Simple Register Coalescing:

..

BB#1:

```

%vreg9 = M2_mpyi %vreg9, %vreg10
%vreg8 = C2_cmpgti %vreg10, 1
%vreg2 = A2_addi %vreg10, -1
%vreg10 = COPY %vreg2
J2_jumpt %vreg8, <BB#1>

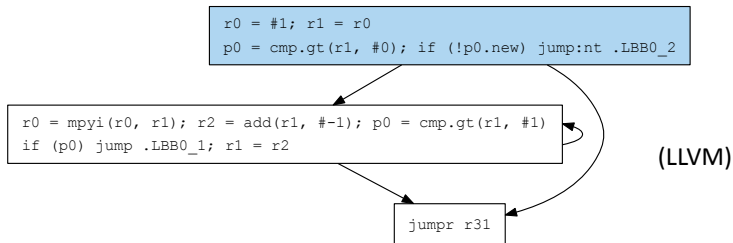
```

..

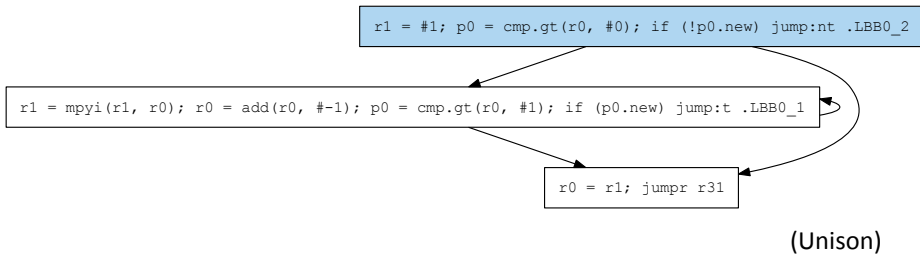


postponing A2\_addi enables coalescing %vreg2, %vreg10





## Unison's initialization is one cycle faster, why?



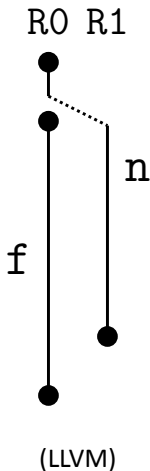
```
int fac(int n) {  
    int f = 1;  
    while(n > 0) {  
        f = f * n;  
        n--;  
    }  
    return f;  
}
```

Calling convention: argument, return value in R0

```
int fac(int n) {  
    int f = 1;  
    while(n > 0) {  
        f = f * n;  
        n--;  
    }  
    return f;  
}
```

However `n` and `f` interfere: move required

```
int fac(int n) {  
    int f = 1;  
    while(n > 0) {  
        f = f * n;  
        n--;  
    }  
    return f;  
}
```

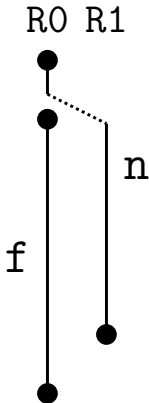


LLVM moves `n` in the initialization block

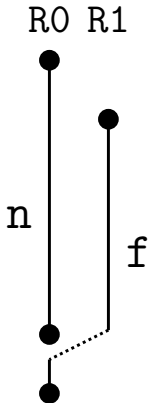
```

int fac(int n) {
  int f = 1;
  while(n > 0) {
    f = f * n;
    n--;
  }
  return f;
}

```

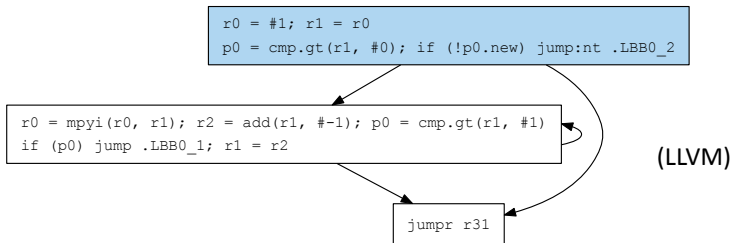


(LLVM)

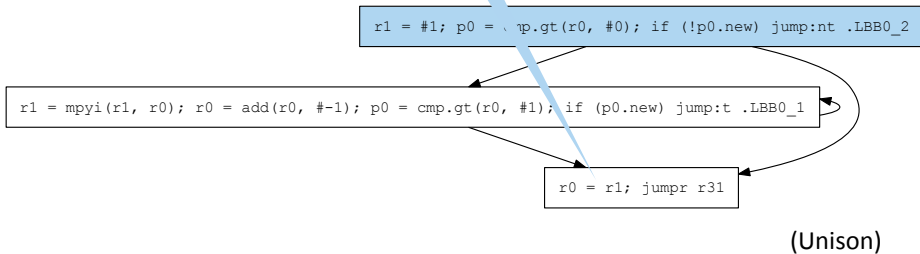


(Unison)

Unison moves `f` in the return block: does it matter?



It does! Unison's move can be scheduled in parallel




```
r0 = #1; r1 = r0
p0 = cmp.gt(r1, #0); if (!p0.new) jump:nt .LBB0_2
```

```
r0 = mpyi(r0, r1); r2 = add(r1, #-1); p0 = cmp.gt(r1, #1)
if (p0) jump .LBB0_1; r1 = r2
```

```
jumpr r31
```

(LLVM)



What if `cmp.gt` could choose `r0` during scheduling?

```
r1 = #1; p0 = cmp.gt(r0, #0); if (!p0.new) jump:nt .LBB0_2
```

```
r1 = mpyi(r1, r0); r0 = add(r0, #-1); p0 = cmp.gt(r0, #1); if (p0.new) jump:t .LBB0_1
```

```
r0 = r1; jumpr r31
```

(Unison)

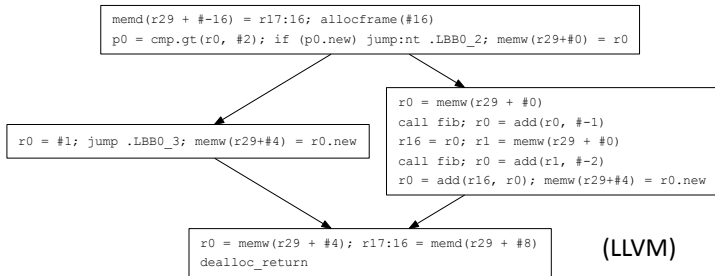
# Case Study: fib

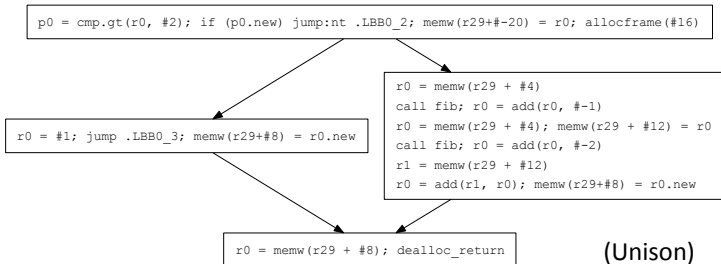
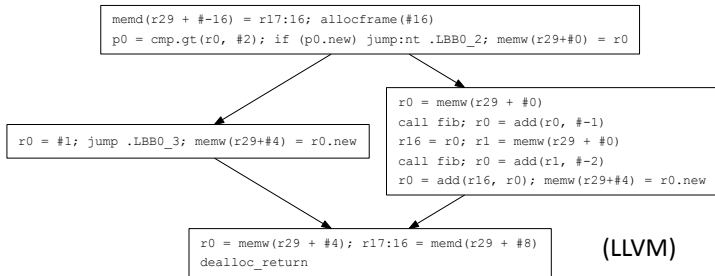
- Simple recursive *Fibonacci*:

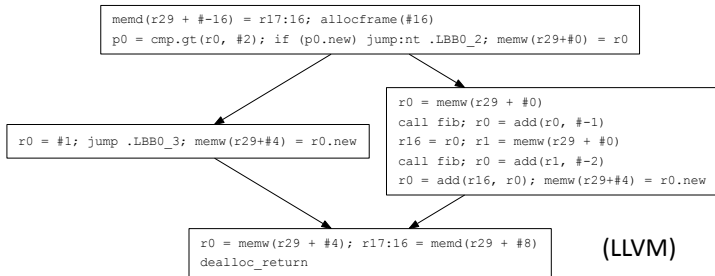
```
int fib(int n) {  
    if(n <= 2) {  
        return 1;  
    }  
    return fib(n-1) + fib(n-2);  
}
```

- Exposes opportunity for better spilling

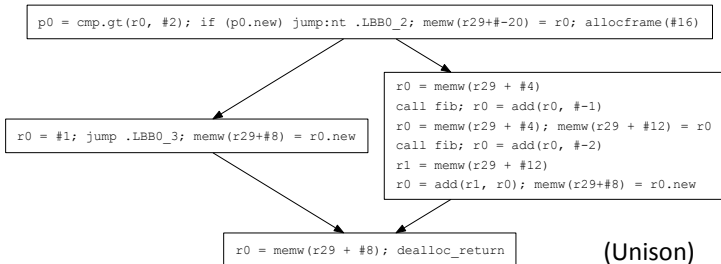


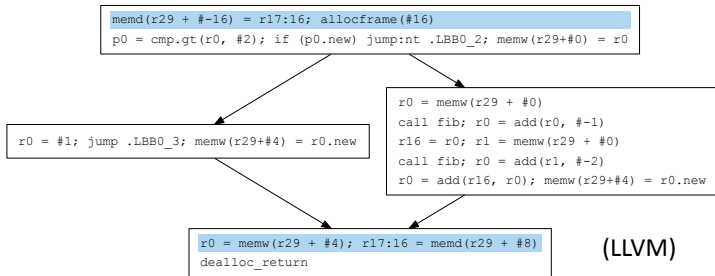




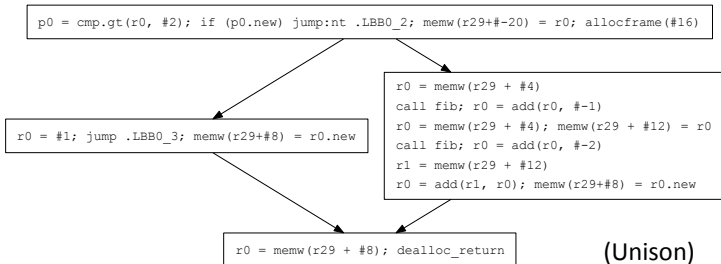


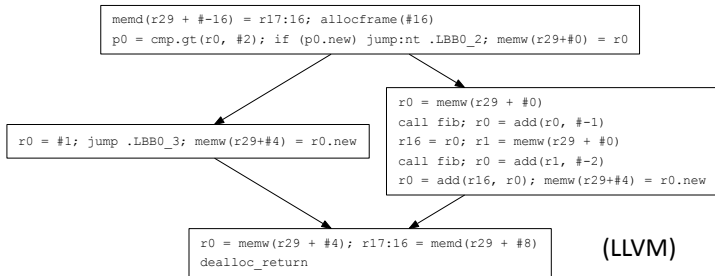
## recursive case requires spilling, where to spill?



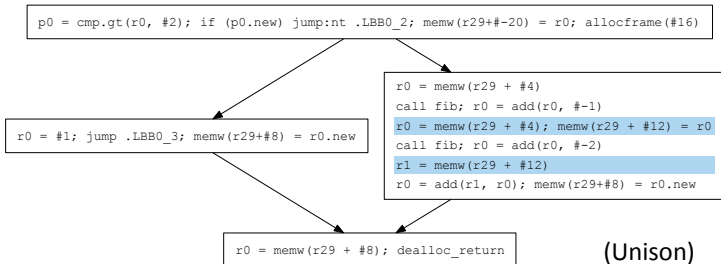


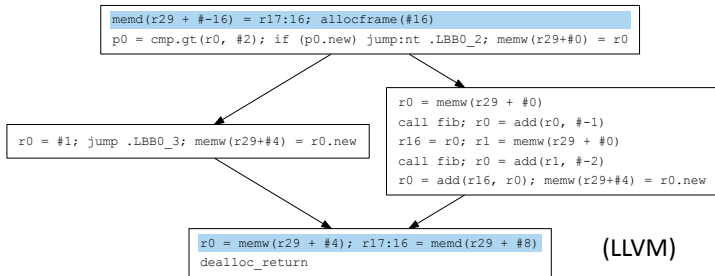
## LLVM frees a callee-saved register (always)



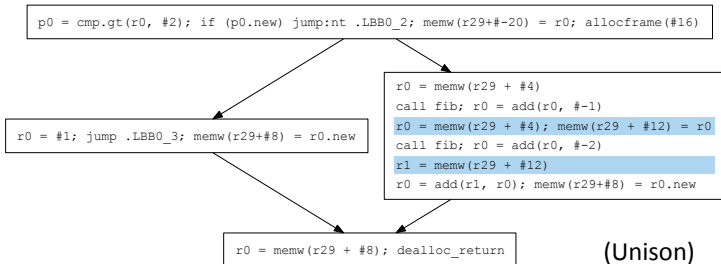


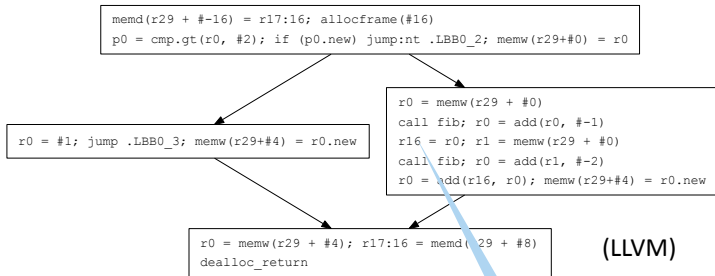
## Unison spills the value directly (recursive case)





## LLVM spills twice as much as Unison, why?





LLVM's register allocator thinks using r16 has no cost

# After Virtual Register Rewriter:

```

..
BB#0:

S2_storeri_io <fi#1>, 0,
%P0 = C2_cmpgti %R0, 2
J2_jumprt %P0, <BB#2>

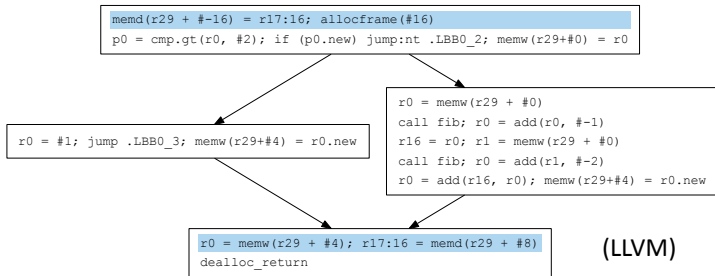
..
BB#3:

%R0 = L2_loadri_io <fi#0>, 0

JMPret %R31

..

```



# After Prologue/Epilogue Insertion & Frame Finalization:

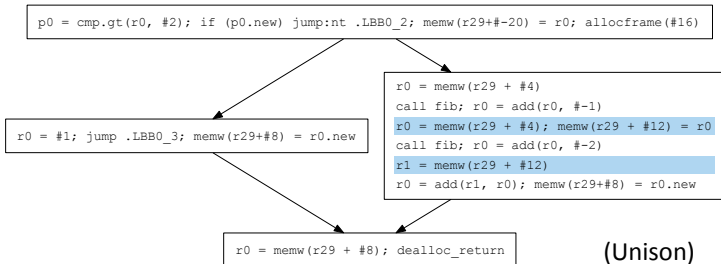
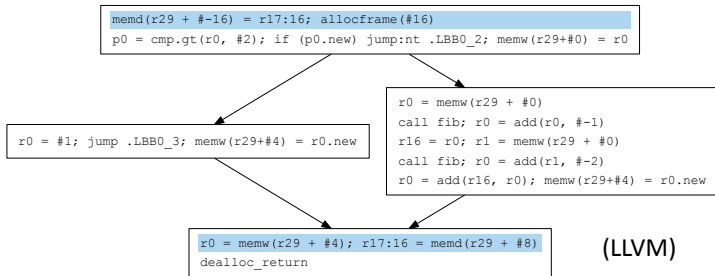
```

..
BB#0:
    S2_allocframe 16
    S2_storerd_io %R29, 8, %D8
    S2_storeri_io %R29, 0, %R0
    %P0 = C2_cmpgti %R0, 2
    J2_jumtp %P0, <BB#2>

..
BB#3:
    %R0 = L2_loadri_io %R29, 4
    %D8 = L2_loadrd_io %R29, 8
    L4_return
..
  
```

but r16 is callee-saved  
and must be preserved





```
memd(r29 + #-16) = r17:16; allocframe(#16)
p0 = cmp.gt(r0, #2); if (p0.new) jump:nt .LBB0_2; memw(r29+#0) = r0
```

```
r0 = #1; jump .LBB0_3; memw(r29+#4) = r0.new
```

```
r0 = memw(r29 + #0)
call fib; r0 = add(r0, #-1)
r16 = r0; r1 = memw(r29 + #0)
call fib; r0 = add(r1, #-2)
r0 = add(r16, r0); memw(r29+#4) = r0.new
```

```
r0 = memw(r29 + #4); r17:16 = memd(r29 + #8)
dealloc_return
```

(LLVM)

could LLVM handle callee-saved spilling earlier?

```
p0 = cmp.gt(r0, #2); if (p0.new) jump:nt .LBB0_2; memw(r29+#-20) = r0; allocframe(#16)
```

```
r0 = #1; jump .LBB0_3; memw(r29+#8) = r0.new
```

```
r0 = memw(r29 + #4)
call fib; r0 = add(r0, #-1)
r0 = memw(r29 + #4); memw(r29 + #12) = r0
call fib; r0 = add(r0, #-2)
r1 = memw(r29 + #12)
r0 = add(r1, r0); memw(r29+#8) = r0.new
```

```
r0 = memw(r29 + #8); dealloc_return
```

(Unison)



# Case Study: chol

- Complex Cholesky decomposition (simplified)
- Illustrates the need for accurate information
- Exposes opportunity for better freq. estimation

# Case Study: chol

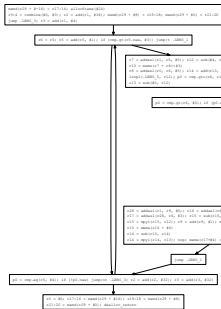
```
typedef struct complex {
    int re;
    int im;
} complex;

int chol(const complex A[4][4], complex L[4][4]) {
    for (int i = 0; i < 4; i++) {
        int f = 1/L[i][i].re;
        for (int j = i+1; j < 4; j++) {
            complex q = {0, 0};
            for (int k = 0; k < i; k++) {
                q.re += L[i][k].re * L[j][k].re;
                q.im -= L[i][k].re * L[j][k].im;
                q.re += L[i][k].im * L[j][k].im;
                q.im += L[i][k].im * L[j][k].re;
            }
            L[j][i].re = f * (A[i][j].re - q.re);
            L[j][i].im = -f * (A[i][j].im - q.im);
        }
    }
    return 0;
}
```









(LLVM)

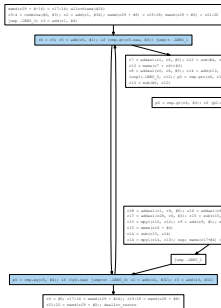


(Unison)

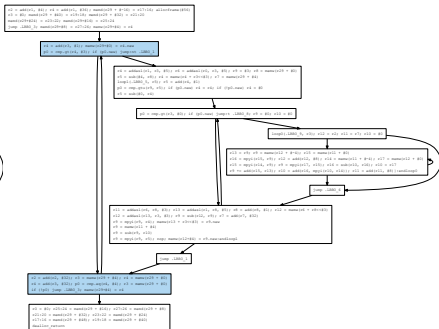
Unison optimizes that basic block (6 vs. 4 cycles)







(LLVM)



(Unison)

so much that Unison's code performs worse!

```

typedef struct complex {
    int re;
    int im;
} complex;

int chol(const complex A[4][4], complex L[4][4]) {
    for (int i = 0; i < 4; i++) {
        int f = 1/L[i][i].re;
        for (int j = i+1; j < 4; j++) {
            complex q = {0, 0};
            for (int k = 0; k < i; k++) {
                q.re += L[i][k].re * L[j][k].re;
                q.im -= L[i][k].re * L[j][k].im;
                q.re += L[i][k].im * L[j][k].im;
                q.im += L[i][k].im * L[j][k].re;
            }
            L[j][i].re = f * (A[i][j].re - q.re);
            L[j][i].im = -f * (A[i][j].im - q.im);
        }
    }
    return 0;
}

```



could LLVM's freq. estimation consider loop counts?

1 Optimal Approaches

2 Model

3 Results

4 Case Studies

**5 Conclusion**

# Unison Is Practical and Effective

- Integrated
  - register allocation
  - instruction scheduling
- Simple, optimal, slower
- For LLVM Users
  - traditional LLVM for compile/debug cycle
  - LLVM + Unison for release builds
- For LLVM developers
  - evaluation of heuristics
  - identification of improvement opportunities
    - coalescing, scheduling, spilling, frequency estimation
  - but: aggressive optimization requires accuracy

`unison-code.github.io`