

Lecture 24

# **CIS 4521/5521: COMPILERS**

# Announcements

- HW6: Analysis & Optimizations
  - Alias analysis, constant propagation, dead code elimination, register allocation
  - Due: Wednesday, April 30<sup>th</sup> at 10:00pm
  - Posted test case due Tuesday, April 29<sup>th</sup> at 10:00pm
- Lecture Cancelled 4/24
  - Dr. Zdancewic will be out of town
- Final Exam:
  - According to registrar: Thursday, May 8<sup>th</sup> noon - 2:00pm
  - Coverage: emphasizes material since the midterm
  - Cheat sheet: one, hand-written, double-sided, letter-sized page of notes



# REGISTER ALLOCATION

# Moving Towards Register Allocation

- The OAT compiler currently generates as *many* temporary variables as it needs
  - These are the `%uids` you should be very familiar with by now.
- Current compilation strategy:
  - Each `%uid` maps to a stack location.
  - This yields programs with many loads/stores to memory.
  - Very inefficient.
- Ideally, we'd like to map as many `%uid`'s as possible into registers.
  - Eliminate the use of the `alloca` instruction?
  - Only 16 max registers available on 64-bit X86
  - `%rsp` and `%rbp` are reserved and some have special semantics, so really only 10 or 12 available
  - This means that a register must hold more than one slot
- When is this safe?

# Register Allocation Problem

- Given: an IR program that uses an unbounded number of temporaries
  - e.g. the uids of our LLVM programs
- Find: a mapping from temporaries to machine registers such that
  - program semantics is preserved (i.e. the behavior is the same)
  - register usage is maximized
  - moves between registers are minimized
  - calling conventions / architecture requirements are obeyed
- Stack Spilling
  - If there are  $k$  registers available and  $m > k$  temporaries are live at the same time, then not all of them will fit into registers.
  - So: "spill" the excess temporaries to the stack.

# Example of Register Allocation

source LLVM IR code

```
define i64 @baz(i64 %x1, i64 %x2, i64 %x3, i64 %x4,  
               i64 %x5, i64 %x6, i64 %x7, i64 %x8) {  
    %1 = add i64 %x1, %x2  
    %2 = add i64 %1, %x3  
    %3 = add i64 %2, %x4  
    %4 = add i64 %3, %x5  
    %5 = add i64 %4, %x6  
    %6 = add i64 %5, %x7  
    %7 = add i64 %6, %x8  
    ret i64 %7  
}
```

- Takes advantage of knowledge of calling conventions to avoid moving data.
  - e.g. %x1 is in %rdi, %x2 is in %rsi, etc.
  - leaves %x7 and %x8 on the stack
- Assigns %1 to %rsi  
    %2, %3, ... %7 all to %rdx

resulting x86 assembly

```
_baz:  
    pushq %rbp  
    movq %rsp, %rbp  
    subq $8, %rsp  
    movq %rcx, -8(%rbp)  
    addq %rdi, %rsi  
    addq %rsi, %rdx  
    addq -8(%rbp), %rdx  
    addq %r8, %rdx  
    addq %r9, %rdx  
    addq 16(%rbp), %rdx  
    addq 24(%rbp), %rdx  
    movq %rdx, %rax  
    movq %rbp, %rsp  
    popq %rbp  
    retq
```

# Exploiting Liveness

```
define i64 @baz(i64 %x1, i64 %x2, i64 %x3, i64 %x4, i64 %x5, i64 %x6, i64 %x7, i64 %x8) {
  _entry={%x1, %x2, %x3, %x4, %x5, %x6, %x7, %x8}
  IN : {%x1, %x2, %x3, %x4, %x5, %x6, %x7, %x8}
  %1 = add i64 %x1, %x2
  OUT: {%1, %x3, %x4, %x5, %x6, %x7, %x8}
  IN : {%1, %x3, %x4, %x5, %x6, %x7, %x8}
  %2 = add i64 %1, %x3
  OUT: {%2, %x4, %x5, %x6, %x7, %x8}
  IN : {%2, %x4, %x5, %x6, %x7, %x8}
  %3 = add i64 %2, %x4
  OUT: {%3, %x5, %x6, %x7, %x8}
  IN : {%3, %x5, %x6, %x7, %x8}
  %4 = add i64 %3, %x5
  OUT: {%4, %x6, %x7, %x8}
  IN : {%4, %x6, %x7, %x8}
  %5 = add i64 %4, %x6
  OUT: {%5, %x7, %x8}
  IN : {%5, %x7, %x8}
  %6 = add i64 %5, %x7
  OUT: {%6, %x8}
  IN : {%6, %x8}
  %7 = add i64 %6, %x8
  OUT: {%7}
  IN : {%7}
  ret i64 %7
  OUT: {}
}
```

Results of running the  
printanalysis tool from HW6  
on the example code from  
llprograms/call6.ll

Observe:

- %x2 and %1 are not live at the same time
- none of %2, ..., %7 are live at the same time

# Linear-Scan Register Allocation

Simple, greedy register-allocation strategy:

1. Compute liveness information:  $\text{live}(x)$ 
  - recall:  $\text{live\_out}(x)$  is the set of uids that are live on exit from  $x$ 's definition
2. Let  $\text{pal}$  be the set of usable registers
  - usually reserve a couple for spill code [our implementation uses `rax,rcx`]
3. Maintain "layout"  $\text{uid\_loc}$  that maps uids to locations
  - locations include registers and stack slots  $n$ , starting at  $n=0$
4. Scan through the program. For each instruction that defines a uid  $x$ 
  - $\text{used} = \{r \mid \text{reg } r = \text{uid\_loc}(y) \text{ s.t. } y \in \text{live\_out}(x)\}$
  - $\text{available} = \text{pal} - \text{used}$
  - If  $\text{available}$  is empty: *// no registers available, spill*  
 $\text{uid\_loc}(x) := \text{slot } n \ ; \ n = n + 1$
  - Otherwise, pick  $r$  in  $\text{available}$ : *// choose an available register*  
 $\text{uid\_loc}(x) := \text{reg } r$



# For HW6

- HW 6 implements two naive register allocation strategies:
  - **none**: spill all registers
  - **greedy**: uses (a slightly "nerfed" version of) linear scan
- Also offers choice of liveness
  - **trivial**: assume all variables are live everywhere
  - **dataflow**: use the dataflow algorithms
- Your job: do "better" than these.
  - To beat "greedy" on small programs – it is necessary to take into account the calling conventions
- Quality Metric - *lower* score is better:
  - total number of memory accesses  
(**Ind2** and **Ind3** operands, **Push/Pop**)
  - ties broken by total number of instructions
- Linear scan is OK (and can be optimal if it gets lucky)
  - but... how can we do better?

# Linear Scan vs. "nerfed" Linear Scan

source LLVM IR code

```
define i64 @baz(i64 %x1, i64 %x2, i64 %x3, i64 %x4,  
               i64 %x5, i64 %x6, i64 %x7, i64 %x8) {  
    %1 = add i64 %x1, %x2  
    %2 = add i64 %1, %x3  
    %3 = add i64 %2, %x4  
    %4 = add i64 %3, %x5  
    %5 = add i64 %4, %x6  
    %6 = add i64 %5, %x7  
    %7 = add i64 %6, %x8  
    ret i64 %7  
}
```

```
_baz:  
    pushq %rbp  
    movq %rsp, %rbp  
    subq $8, %rsp  
    movq %rcx, -8(%rbp)  
    addq %rdi, %rsi  
    addq %rsi, %rdx  
    addq -8(%rbp), %rdx  
    addq %r8, %rdx  
    addq %r9, %rdx  
    addq 16(%rbp), %rdx  
    addq 24(%rbp), %rdx  
    movq %rdx, %rax  
    movq %rbp, %rsp  
    popq %rbp  
    retq
```

Linear Scan

resulting x86 assembly

```
_baz:  
    pushq    %rbp  
    movq     %rsp, %rbp  
    subq     $8, %rsp  
    movq     %rcx, -8(%rbp)  
    movq     %rdi, %r10  
    addq     %rsi, %r10  
    movq     %r10, %rsi  
    addq     %rdx, %rsi  
    movq     %rsi, %rdx  
    addq     -8(%rbp), %rdx  
    movq     %rdx, %rsi  
    addq     %r8, %rsi  
    movq     %rsi, %rdx  
    addq     %r9, %rdx  
    movq     %rdx, %rsi  
    addq     16(%rbp), %rsi  
    movq     %rsi, %rdx  
    addq     24(%rbp), %rdx  
    movq     %rdx, %rax  
    movq     %rbp, %rsp  
    popq     %rbp  
    retq
```

"nerfed"  
Linear Scan



# GRAPH COLORING

# Register Allocation

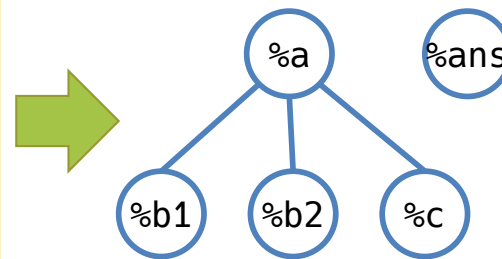
## Basic process:

1. Compute liveness information for each temporary.
2. Create an *interference graph*:
  - Nodes are temporary variables.
  - There is an edge between node  $n$  and  $m$  if  $n$  is live at the same time as  $m$
3. Try to color the graph
  - Each color corresponds to a register
4. In case step 3. fails, “spill” a register to the stack and repeat the whole process.
5. Rewrite the program to use registers

# Interference Graphs

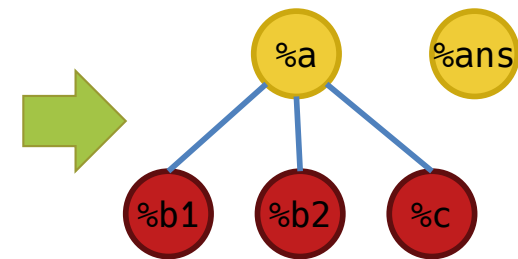
- Nodes of the graph are `%uids`
- Edges connect variables that *interfere* with each other
  - Two variables interfere if their live ranges intersect (i.e. there is an edge in the control-flow graph across which they are both live).
- Register assignment is a *graph coloring*.
  - A graph coloring assigns each node in the graph a color (register)
  - Any two nodes connected by an edge must have different colors.
- Example:

```
// live = {%a}
%b1 = add i32 %a, 2
// live = {%a,%b1}
%c = mult i32 %b1, %b1
// live = {%a,%c}
%b2 = add i32 %c, 1
// live = {%a,%b2}
%ans = mult i32 %b2, %a
// live = {%ans}
return %ans;
```



Interference Graph

Even better: take calling conventions into account and put %ans in %rax.



2-Coloring of the graph  
red = r8  
yellow = r9

# Register Allocation Questions

- Can we efficiently find a  $k$ -coloring of the graph whenever possible?
  - Answer: in general the problem is NP-complete (it requires search)
  - But, we can do an efficient approximation using heuristics.
- How do we assign registers to colors?
  - If we do this in a smart way, we can eliminate redundant MOV instructions.
- What do we do when there aren't enough colors/registers?
  - We have to use stack space, but how do we do this effectively?

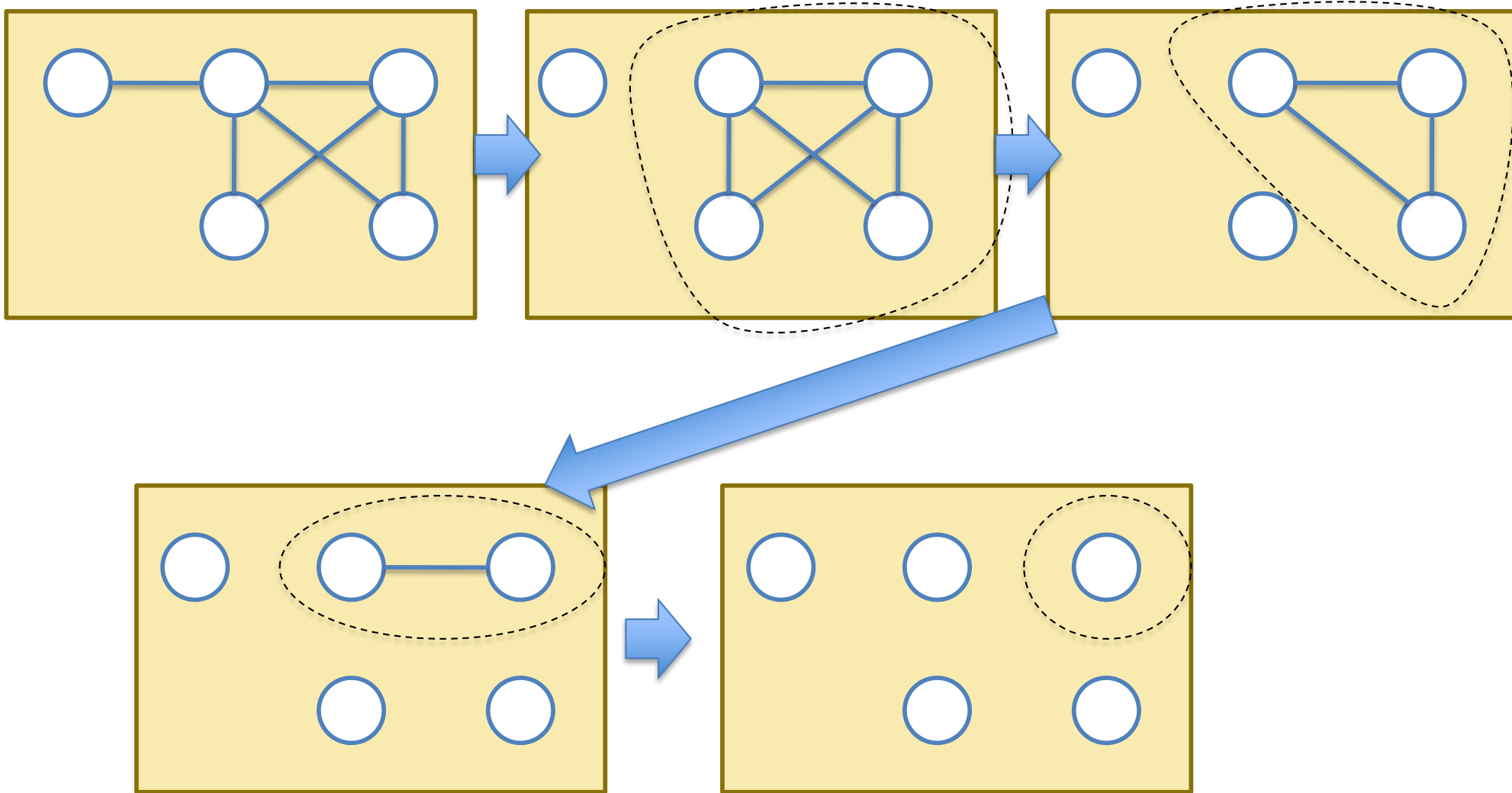
# Coloring a Graph: Kempe's Algorithm

Kempe [1879] provides this algorithm for K-coloring a graph.

It's a recursive algorithm that works in three steps:

1. Find a node with degree  $< K$  and cut it out of the graph.
  - Remove the nodes and edges.
  - This is called *simplifying* the graph
2. Recursively K-color the remaining subgraph
3. When remaining graph is colored, there must be at least one free color available for the deleted node (since its degree was  $< K$ ). Pick such a color.

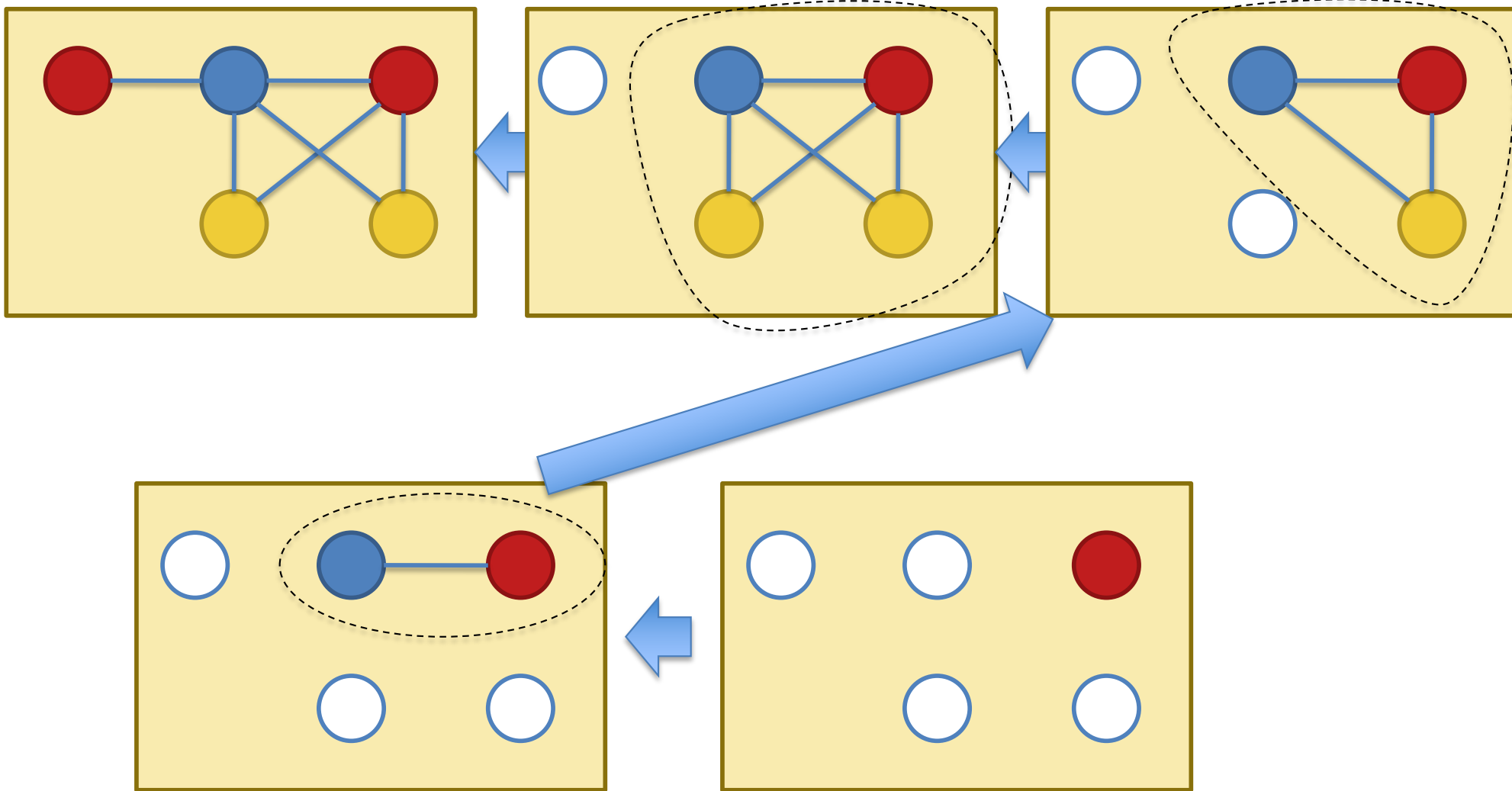
# Example: 3-color this Graph



Recurring Down the Simplified Graphs



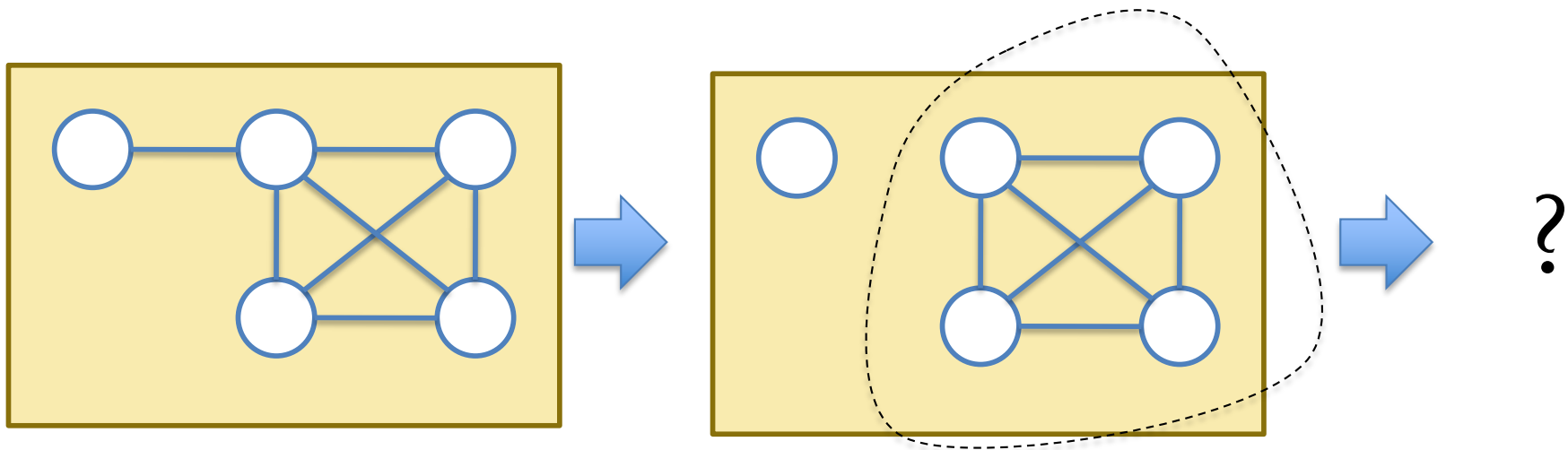
# Example: 3-color this Graph



Assigning Colors on the way back up.

# Failure of the Algorithm

- If the graph cannot be colored, it will simplify to a graph where every node has at least  $K$  neighbors.
  - This can happen even when the graph is  $K$ -colorable!
  - This is a symptom of NP-hardness (it requires search)
- Example: When trying to 3-color this graph:

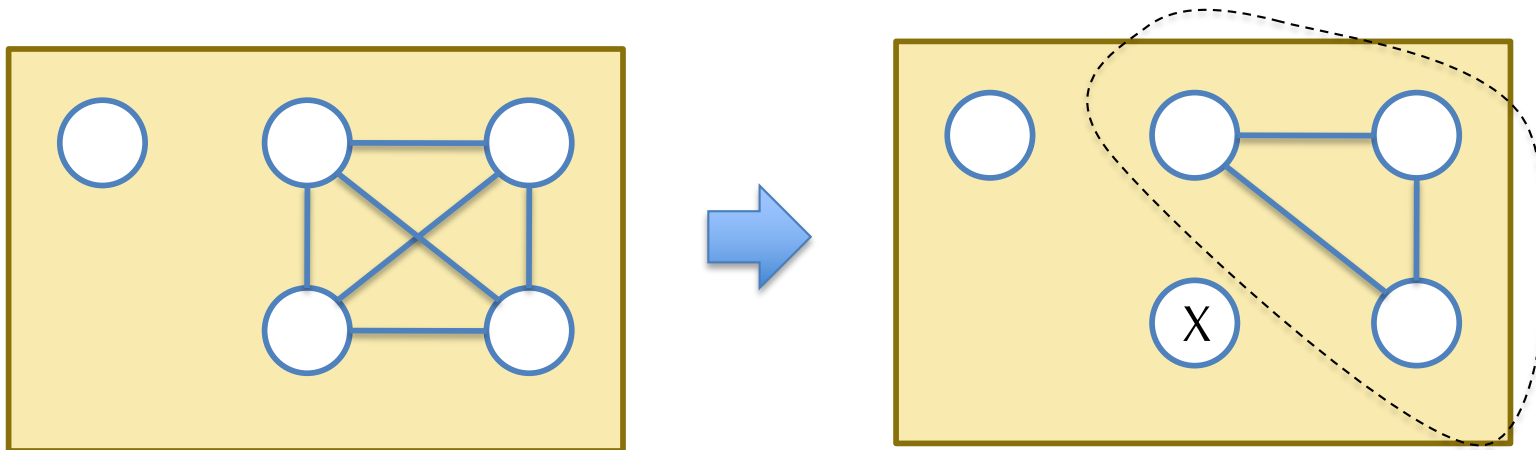


# Spilling

- Idea: If we can't K-color the graph, we need to store one temporary variable on the stack.
- Which variable to spill?
  - Pick one that isn't used very frequently
  - Pick one that isn't used in a (deeply nested) loop
  - Pick one that has high interference  
(since removing it will make the graph easier to color)
- In practice: some weighted combination of these criteria
- When coloring:
  - Mark the node as spilled
  - Remove it from the graph
  - Keep recursively coloring

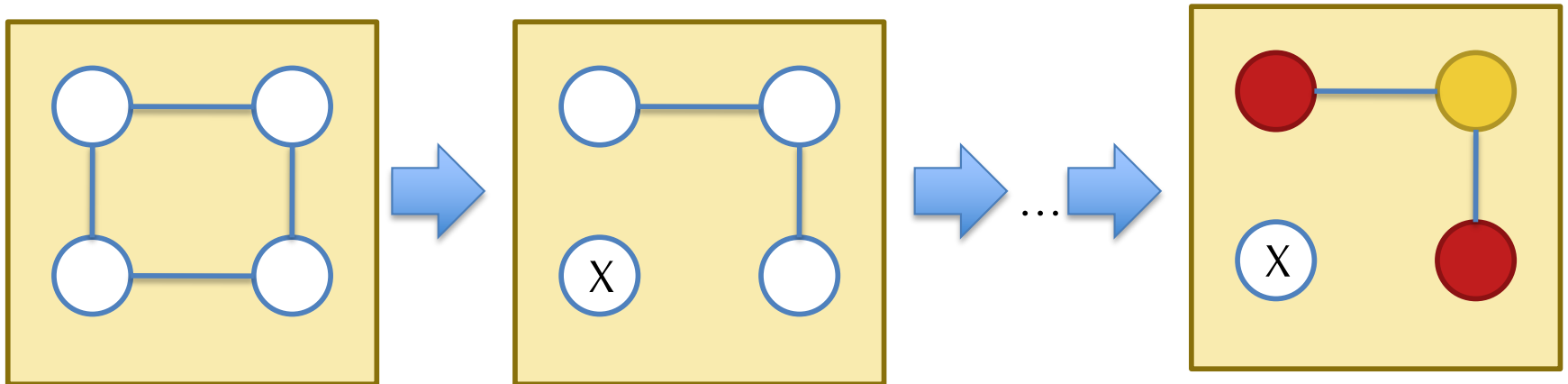
# Spilling, Pictorially

- Select a node to spill
- Mark it and remove it from the graph
- Continue coloring



# Optimistic Coloring

- Sometimes it is possible to color a node marked for spilling.
  - If we get “lucky” with the choices of colors made earlier.
- Example: When 2-coloring this graph:



- Even though the node was marked for spilling, we can color it.
- So: on the way down, mark for spilling, but don't actually spill...

# Accessing Spilled Registers

- If optimistic coloring fails, we need to generate code to move the spilled temporary to & from memory.
- Option 1: Reserve registers specifically for moving to/from memory.
  - Con: Need at least two registers (one for each source operand of an instruction), so decreases total # of available registers by 2.
  - Pro: Only need to color the graph once.
  - Not good on 32bit x86 because there are too few registers & too many constraints on how they can be used.
  - OK on 64bit x86 and other processors. (We use this for HW6)
- Option 2: Rewrite the program to use a new temporary variable, with explicit moves to/from memory.
  - Pro: Need to reserve fewer registers.
  - Con: Introducing temporaries changes live ranges, so must recompute liveness & recolor graph

# Example Spill Code

- Suppose temporary  $t$  is marked for spilling to stack slot located at  $[rbp+offs]$

- Rewrite the program like this:

```
t = a op b;
```

```
...
```

```
x = t op c
```

```
...
```

```
y = d op t
```



```
t = a op b // defn. of t
```

```
Mov [rbp+offs], t
```

```
...
```

```
Mov t37, [rbp+offs] // use 1 of t
```

```
x = t37 op c
```

```
...
```

```
Mov t38, [rbp+offs] // use 2 of t
```

```
y = d op t38
```

- Here,  $t37$  and  $t38$  are freshly generated temporaries that replace  $t$  for different uses of  $t$ .
- Rewriting the code in this way breaks  $t$ 's live range up:  
 $t$ ,  $t37$ ,  $t38$  are only live across one edge

# Precolored Nodes

- Some variables must be pre-assigned to registers.
  - e.g., on x86 the shift instructions must use %rcx
- To properly allocate temporaries, we treat registers as nodes in the interference graph with pre-assigned colors.
  - Pre-colored nodes can't be removed during simplification.
  - When the graph is empty except the pre-colored nodes, we have reached the point where we start coloring the rest of the nodes.

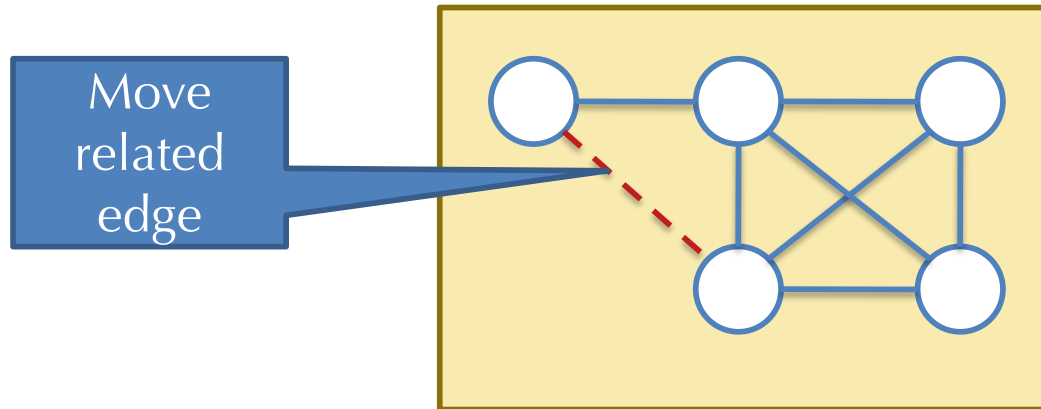


# Picking Good Colors

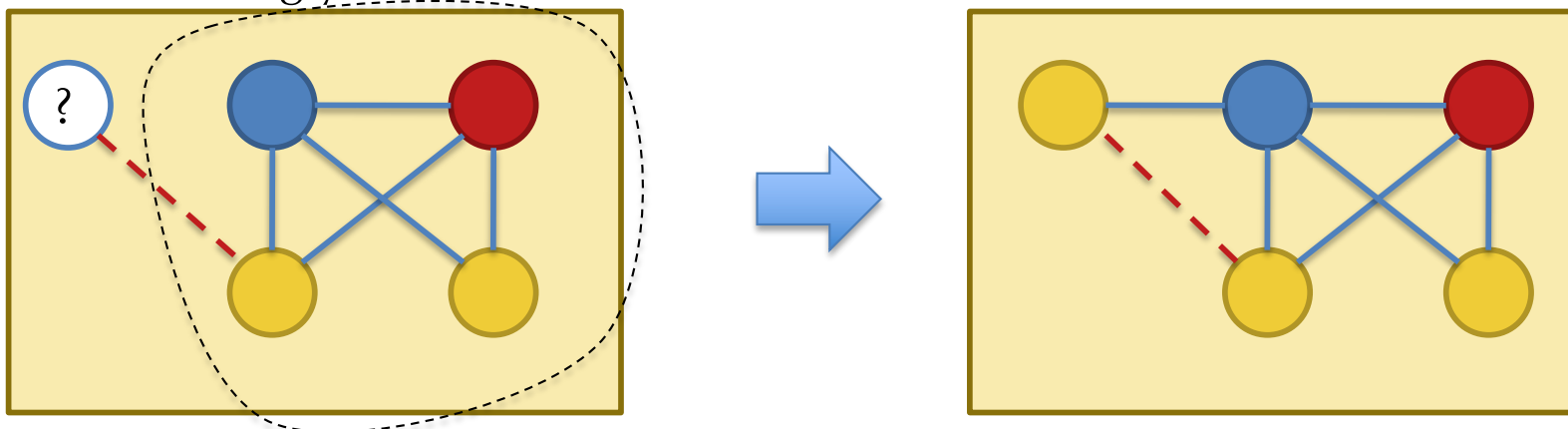
- When choosing colors during the coloring phase, *any* choice is semantically correct, but some choices are better for performance.
- Example:  
`movq t1, t2`
  - If t1 and t2 can be assigned the same register (color) then this move is redundant and can be eliminated.
- A simple color choosing strategy that helps eliminate such moves:
  - Add a new kind of “move related” edge between the nodes for t1 and t2 in the interference graph.
  - When choosing a color for t1 (or t2), if possible, pick a color of an already colored node reachable by a move-related edge.

# Example Color Choice

- Consider 3-coloring this graph, where the dashed edge indicates that there is a Move from one temporary to another.

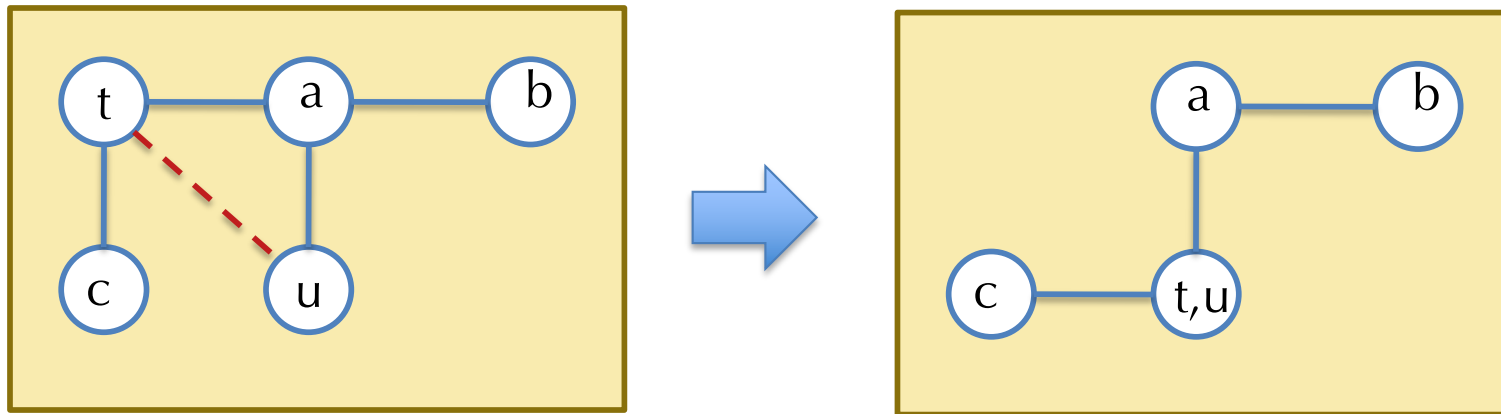


- After coloring the rest, we have a choice:
  - Picking yellow is better than red because it will eliminate a move.

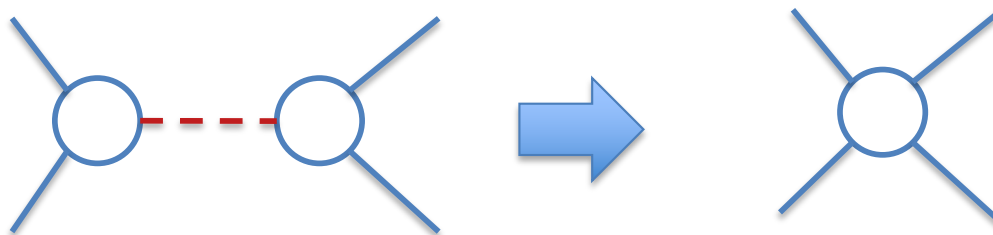


# Coalescing Interference Graphs

- A more aggressive strategy is to *coalesce* nodes of the interference graph if they are connected by move-related edges.
  - Coalescing the nodes *forces* the two temporaries to be assigned the same register.



- Idea: interleave simplification and coalescing to maximize the number of moves that can be eliminated.
- Problem: coalescing can sometimes increase the degree of a node.



# Conservative Coalescing

- Two strategies are guaranteed to preserve the  $k$ -colorability of the interference graph.
  1. *Brigg's strategy*: It's safe to coalesce  $x$  &  $y$  if the resulting node will have fewer than  $k$  neighbors (with degree  $\geq k$ ).
  2. *George's strategy*: We can safely coalesce  $x$  &  $y$  if for every neighbor  $t$  of  $x$ , either  $t$  already interferes with  $y$  or  $t$  has degree  $< k$ .

# Complete Register Allocation Algorithm

1. Build interference graph (precolor nodes as necessary).
  - Add move related edges
2. Reduce the graph (building a stack of nodes to color).
  1. Simplify the graph as much as possible without removing nodes that are move related (i.e. have a move-related neighbor). Remaining nodes are high degree or move-related.
  2. Coalesce move-related nodes using Brigg's or George's strategy.
  3. Coalescing can reveal more nodes that can be simplified, so repeat 2.1 and 2.2 until no node can be simplified or coalesced.
  4. If no nodes can be coalesced *freeze* (remove) a move-related edge and keep trying to simplify/coalesce.
3. If there are non-precolored nodes left, mark one for spilling, remove it from the graph and continue doing step 2.
4. When only pre-colored node remain, start coloring (popping simplified nodes off the top of the stack).
  1. If a node must be spilled, insert spill code as on slide 14 and rerun the whole register allocation algorithm starting at step 1.

# Details for HW6

- New backend code generator uses a PMov instruction for
  - Function Declaration:  
moving arguments from calling convention registers to their allocated slots
  - Compiling a Call instruction:  
moving arguments from their allocated slot to their locations according to the calling conventions
- Therefore, to do well, your allocation strategy should prefer to put `%uid` values in locations that work well for the calling conventions.

# PMov Instruction

- The PMov instruction:
  - takes a list of "assignments"
 
$$(dst_1 \leftarrow src_1) (dst_2 \leftarrow src_2) (dst_3 \leftarrow src_3) \dots (dst_n \leftarrow src_n)$$

Algorithm:

- filter out
  - srcs already in the right location
  - moves to destinations that aren't live
- emit "ready" moves
  - dst is not also a src (so no collision)
- if none are ready:
  - push src of first move
  - recursively process results
  - pop src into its dst

- generates relatively efficient code to shuffle all of the  $src_i$  to the corresponding  $dst_i$

ol		ol'		2		2		3		2
x <- y		x <- y		w <- x		MOV x, w		MOV x, w		MOV x, w
z <- z	==>	<del>z &lt;- z</del>	==>	-----	==>	-----	==>	PUSH y	==>	PUSH y
w <- x		w <- x		x <- y		x <- y		y <- z		MOV z, y
y <- z		y <- z		y <- z		y <- z		POP x		POP x



# LOOPS AND DOMINATORS



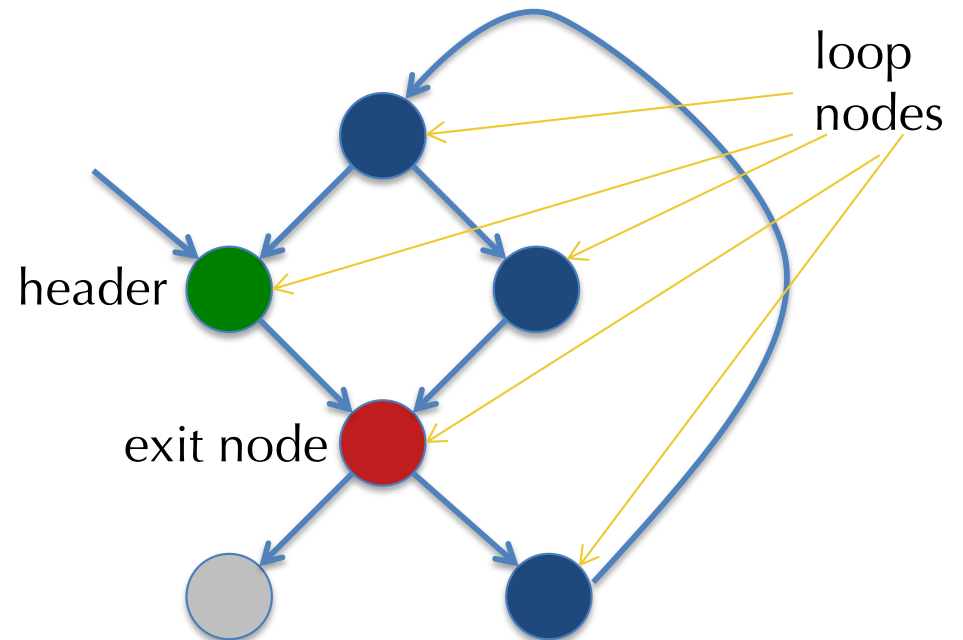
# Loops in Control-flow Graphs

- Taking into account loops is important for optimizations.
  - The 90/10 rule applies, so optimizing loop bodies is important
- Should we apply loop optimizations at the AST level or at a lower representation?
  - Loop optimizations benefit from other IR-level optimizations and vice-versa, so it is good to interleave them.
- Loops may be hard to recognize at the quadruple / LLVM IR level.
  - Many kinds of loops: while, do/while, for, continue, goto...
- Problem: *How do we identify loops in the control-flow graph?*

# Definition of a Loop

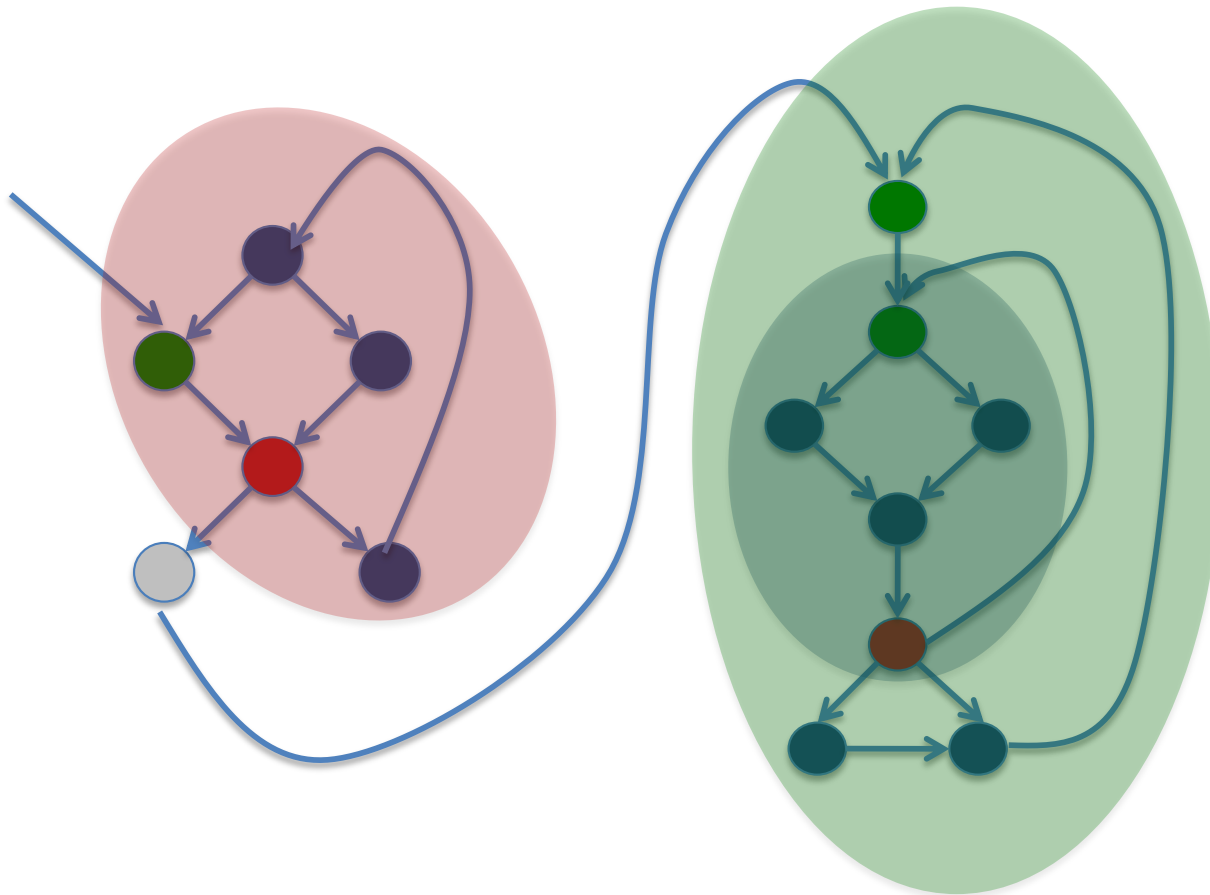
- A *loop* is a set of nodes in the control flow graph.
  - One distinguished entry point called the *header*

- Every node is reachable from the header & the header is reachable from every node.
  - A loop is a *strongly connected component*
- No edges enter the loop except to the header
- Nodes with outgoing edges are called loop exit nodes

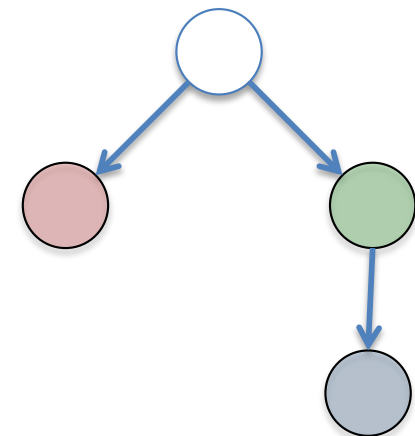


# Nested Loops

- Control-flow graphs may contain many loops
- Loops may contain other loops:



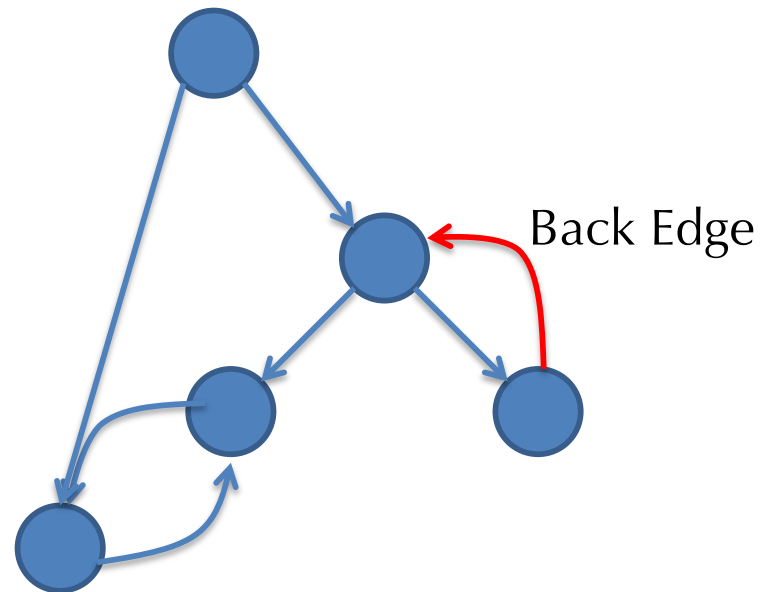
Control Tree:



The control tree depicts the nesting structure of the program.

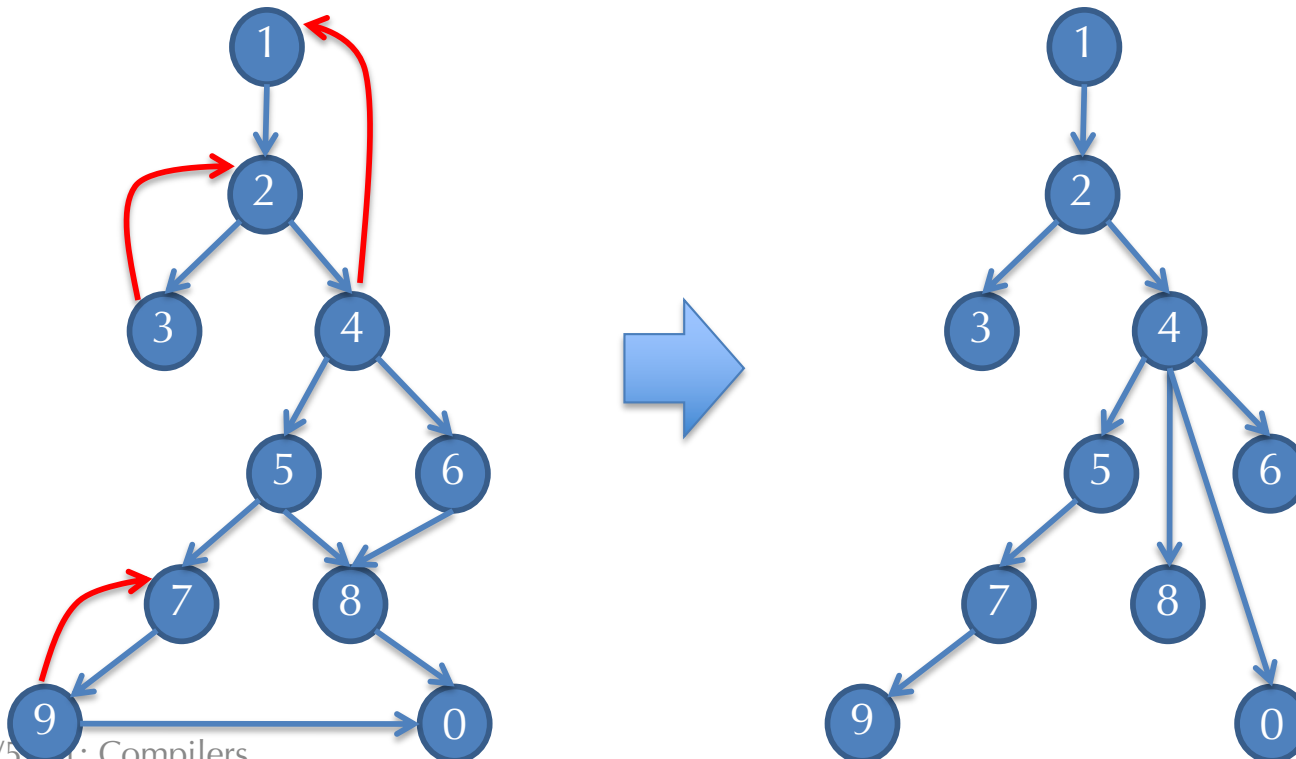
# Control-flow Analysis

- Goal: Identify the loops and nesting structure of the CFG.
- Control flow analysis is based on the idea of *dominators*:
- Node A *dominates* node B if the only way to reach B from the start node is through node A.
- An edge in the graph is a *back edge* if the target node dominates the source node.
- A loop contains at least one back edge.



# Dominator Trees

- Domination is transitive:
  - if A dominates B and B dominates C then A dominates C
- Domination is anti-symmetric:
  - if A dominates B and B dominates A then A = B
- Every flow graph has a dominator tree
  - The Hasse diagram of the dominates relation

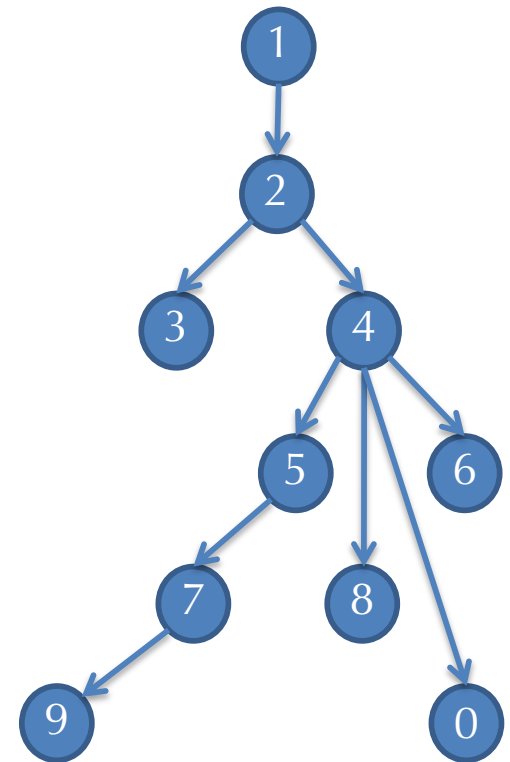


# Dominator Dataflow Analysis

- We can define  $\text{Dom}[n]$  as a forward dataflow analysis.
  - Using the framework we saw earlier:  $\text{Dom}[n] = \text{out}[n]$  where:
- “A node B is dominated by another node A if A dominates *all* of the predecessors of B.”
  - $\text{in}[n] := \bigcap_{n' \in \text{pred}[n]} \text{out}[n']$
- “Every node dominates itself.”
  - $\text{out}[n] := \text{in}[n] \cup \{n\}$
- Formally:  $\mathcal{L}$  = set of nodes ordered by  $\subseteq$ 
  - $T = \{\text{all nodes}\}$
  - $F_n(x) = x \cup \{n\}$
  - $\sqcap$  is  $\cap$
- Easy to show monotonicity and that  $F_n$  distributes over meet.
  - So algorithm terminates and is MOP

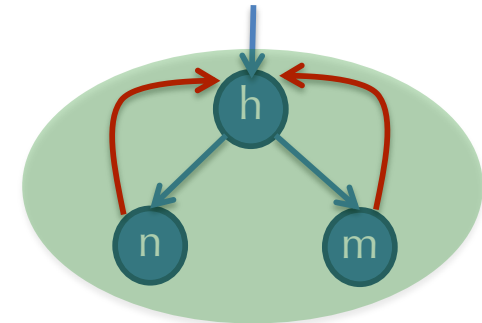
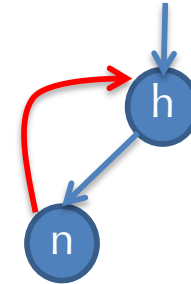
# Improving the Algorithm

- Dom[b] contains just those nodes along the path in the dominator tree from the root to b:
  - e.g., Dom[8] = {1,2,4,8}, Dom[7] = {1,2,4,5,7}
  - There is a lot of sharing among the nodes
- More efficient way to represent Dom sets is to store the dominator *tree*.
  - doms[b] = immediate dominator of b
  - doms[8] = 4, doms[7] = 5
- To compute Dom[b] walk through doms[b]
- Need to efficiently compute intersections of Dom[a] and Dom[b]
  - Traverse up tree, looking for least common ancestor:
  - Dom[8]  $\cap$  Dom[7] = Dom[4]
- See: “A Simple, Fast Dominance Algorithm” Cooper, Harvey, and Kennedy



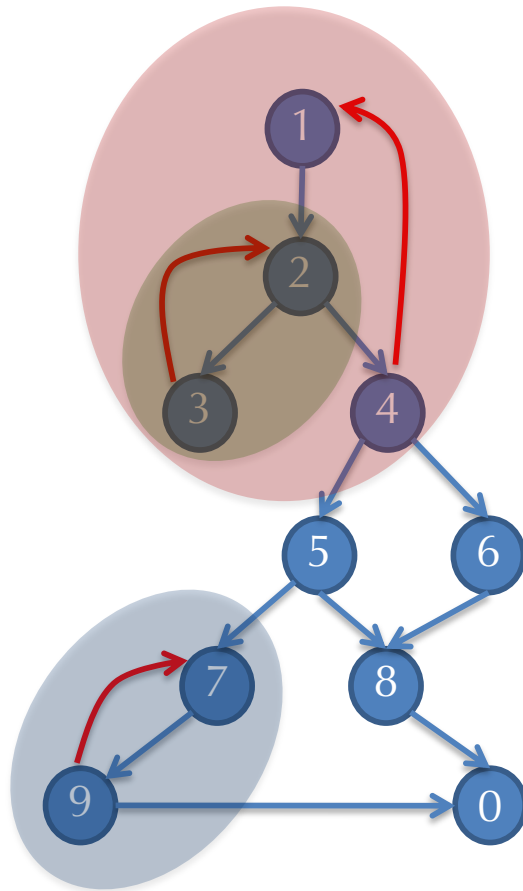
# Completing Control-flow Analysis

- Dominator analysis identifies *back edges*:
  - Edge  $n \rightarrow h$  where  $h$  dominates  $n$
- Each back edge has a *natural loop*:
  - $h$  is the header
  - All nodes reachable from  $h$  that also reach  $n$  without going through  $h$
- For each back edge  $n \rightarrow h$ , find the natural loop:
  - $\{n' \mid n \text{ is reachable from } n' \text{ in } G - \{h\}\} \cup \{h\}$
- Two loops may share the same header: merge them
- Nesting structure of loops is determined by set inclusion
  - Can be used to build the control tree





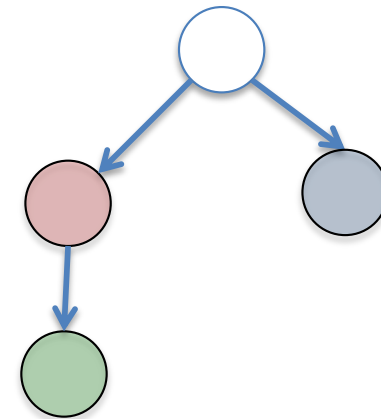
# Example Natural Loops



Natural Loops



Control Tree:



The control tree depicts the nesting structure of the program.

# Uses of Control-flow Information

- Loop nesting depth plays an important role in optimization heuristics.
  - Deeply nested loops pay off the most for optimization.
- Need to know loop headers / back edges for doing
  - loop invariant code motion
  - loop unrolling
- Dominance information also plays a role in converting to SSA form
  - Used internally by LLVM to do register allocation.