

Lecture 27

# **CIS 341: COMPILERS**

# Announcements

- HW 7: Optimization & Experiments
  - Post your benchmark programs early (i.e. tonight!)
  - Due: Tomorrow April 29<sup>th</sup>
- Final Exam:
  - Thursday, May 7<sup>th</sup>
  - 9:00AM
  - Moore 216

Vellvm

# VERIFYING COMPILER TRANSFORMATIONS

# LLVM<sub>ND</sub> Operational Semantics

- Define a transition relation:

$$f \vdash \sigma_1 \mapsto \sigma_2$$

- $f$  is the program
- $\sigma$  is the program state: pc, locals( $\delta$ ), stack, heap
- Nondeterministic
  - $\delta$  maps local `%uids` to sets.
  - Step relation is nondeterministic
- Mostly straightforward (given the heap model)
  - Another wrinkle: phi-nodes executed atomically

# Operational Semantics

	Small Step	Big Step
Nondeterministic	$\text{LLVM}_{ND}$	
Deterministic		

# Deterministic Refinement

	Small Step	Big Step
Nondeterministic	$\text{LLVM}_{ND}$	
	$\Downarrow$	
Deterministic	$\text{LLVM}_D$	

Instantiate 'undef' with default value (0 or null)  $\Rightarrow$  deterministic.

# Big-step Deterministic Refinements

	Small Step	Big Step
Nondeterministic	$\text{LLVM}_{ND}$	
Deterministic	$\text{LLVM}_{Interp} \approx \text{LLVM}_D$	

$\cup$

Bisimulation up to “observable events”:

- external function calls

# Big-step Deterministic Refinements

	Small Step	Big Step
Nondeterministic	<div style="border: 1px solid black; padding: 5px; display: inline-block;"> <math>LLVM_{ND}</math> </div>	
Deterministic	<div style="display: flex; align-items: center; gap: 10px;"> <div style="border: 1px solid black; padding: 5px; display: inline-block;"> <math>LLVM_{Interp}</math> </div> <math>\approx</math> <div style="border: 1px solid black; padding: 5px; display: inline-block;"> <math>LLVM_D</math> </div> </div>	<div style="display: flex; align-items: center; gap: 10px;"> <div style="border: 1px solid black; padding: 5px; display: inline-block;"> <math>LLVM^*_{DFn}</math> </div> <math>\approx</math> <div style="border: 1px solid black; padding: 5px; display: inline-block;"> <math>LLVM^*_{DB}</math> </div> </div>

Simulation up to “observable events”:

- useful for encapsulating behavior of function calls
- large step evaluation of basic blocks

[Tristan, et al. *POPL* '08, Tristan, et al. *PLDI* '09]



# Strategy for Proving Optimizations

- Decompose the program transformation into a sequence of “micro” transformations
  - e.g. code motion =
    1. insert “redundant” instruction
    2. substitute equivalent definitions
    3. remove the “dead” instruction
- Use the backward simulations to show each “micro” transformation correct.
  - Often uses a *safety property*
  - Safety: establish an invariant of the execution of the program
- Compose the individual proofs of correctness

# Safety Properties

- A well-formed program never accesses undefined variables.

If  $\vdash f$  and  $f \vdash \sigma_0 \mapsto^* \sigma$  then  $\sigma$  is not stuck.

$\vdash f$  program  $f$  is well formed  
 $\sigma$  program state  
 $f \vdash \sigma \mapsto^* \sigma$  evaluation of  $f$

- *Initialization:*

If  $\vdash f$  then  $\text{wf}(f, \sigma_0)$ .

- *Preservation:*

If  $\vdash f$  and  $f \vdash \sigma \mapsto \sigma'$  and  $\text{wf}(f, \sigma)$  then  $\text{wf}(f, \sigma')$

- *Progress:*

If  $\vdash f$  and  $\text{wf}(f, \sigma)$  then  $f \vdash \sigma \mapsto \sigma'$

# Safety Properties

- A well-formed program never accesses undefined variables.

If  $\vdash f$  and  $f \vdash \sigma_0 \mapsto^* \sigma$  then  $\sigma$  is not stuck.

$\vdash f$  program  $f$  is well formed  
 $\sigma$  program state  
 $f \vdash \sigma \mapsto^* \sigma$  evaluation of  $f$

- *Initialization:*

If  $\vdash f$  then  $\text{wf}(f, \sigma_0)$

- *Preservation:*

If  $\vdash f$  and  $f \vdash \sigma \mapsto \sigma'$  and  $\text{wf}(f, \sigma)$  then  $\text{wf}(f, \sigma')$

- *Progress:*

If  $\vdash f$  and  $\text{wf}(f, \sigma)$  then  $f \vdash \sigma \mapsto \sigma'$

# Well-formed States

**entry:**

$r_0 = \dots$

$r_1 = \dots$

$r_2 = \dots$

**br**  $r_0$  loop exit

**loop:**

$r_3 = \phi [0; \text{entry}] [r_5; \text{loop}]$

$r_4 = r_1 \times r_2$

$r_5 = r_3 + r_4$

$r_6 = r_5 \geq 100$

**br**  $r_6$  loop exit

**exit:**

$r_7 = \phi [0; \text{entry}] [r_5; \text{loop}]$

$r_8 = r_1 \times r_2$

$r_9 = r_7 + r_8$

**ret**  $r_9$

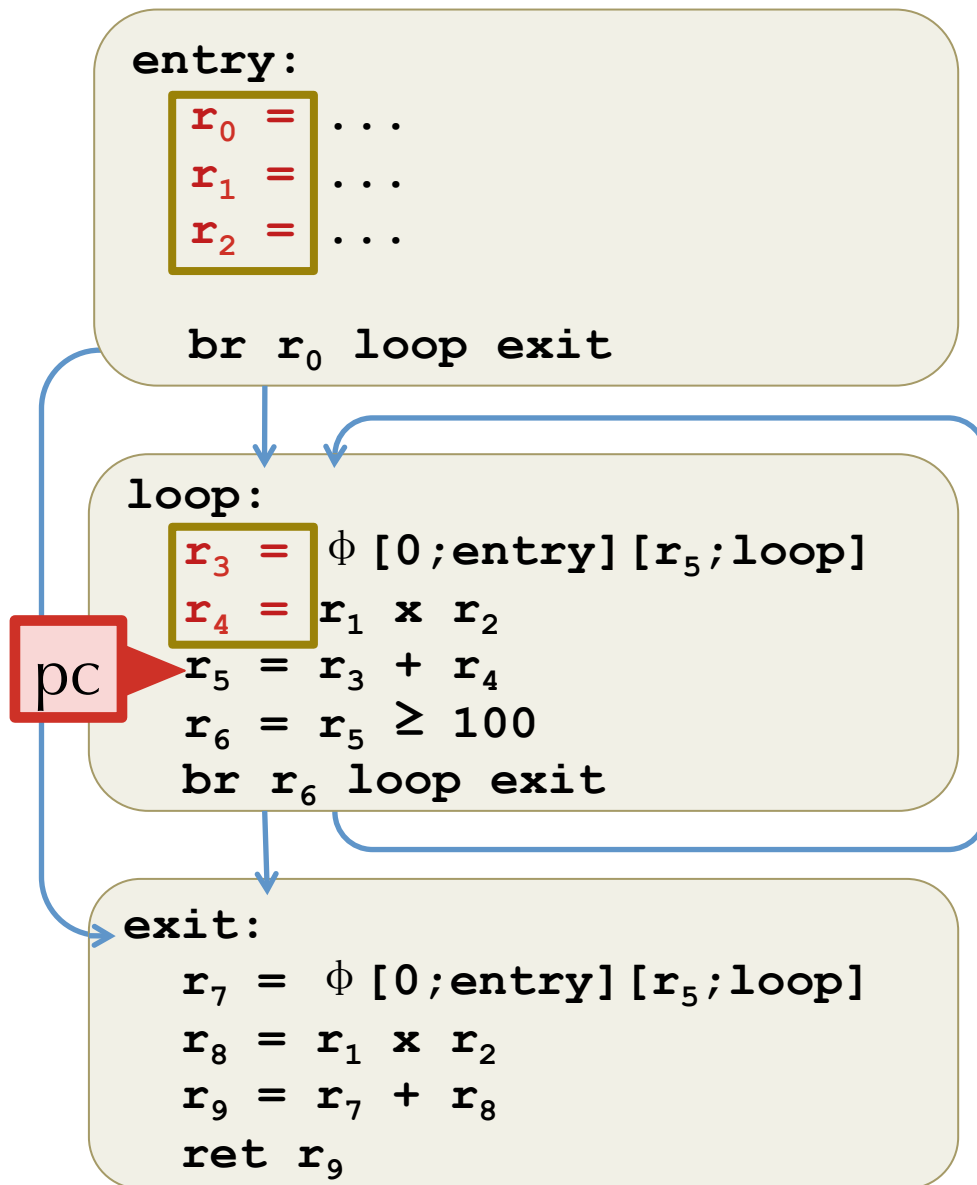
State  $\sigma$  is:

pc = program counter

$\delta$  = local values

pc

# Well-formed States (Roughly)



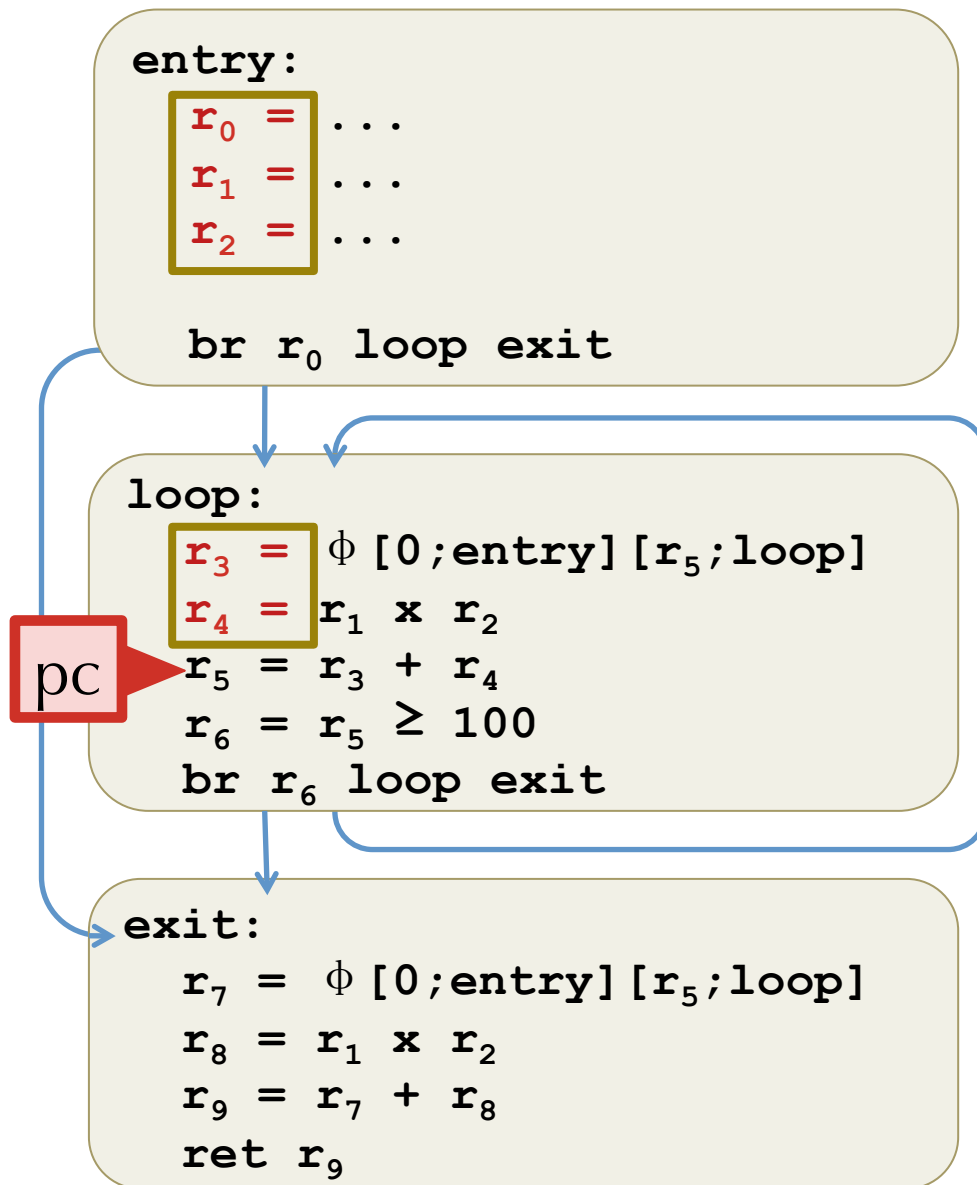
State  $\sigma$  is:

pc = program counter

$\delta$  = local values

$\text{sdom}(f, \text{pc})$  = variable  
defns. that *strictly*  
*dominate* pc.

# Well-formed States (Roughly)



State  $\sigma$  contains:  
pc = program counter  
 $\delta$  = local values

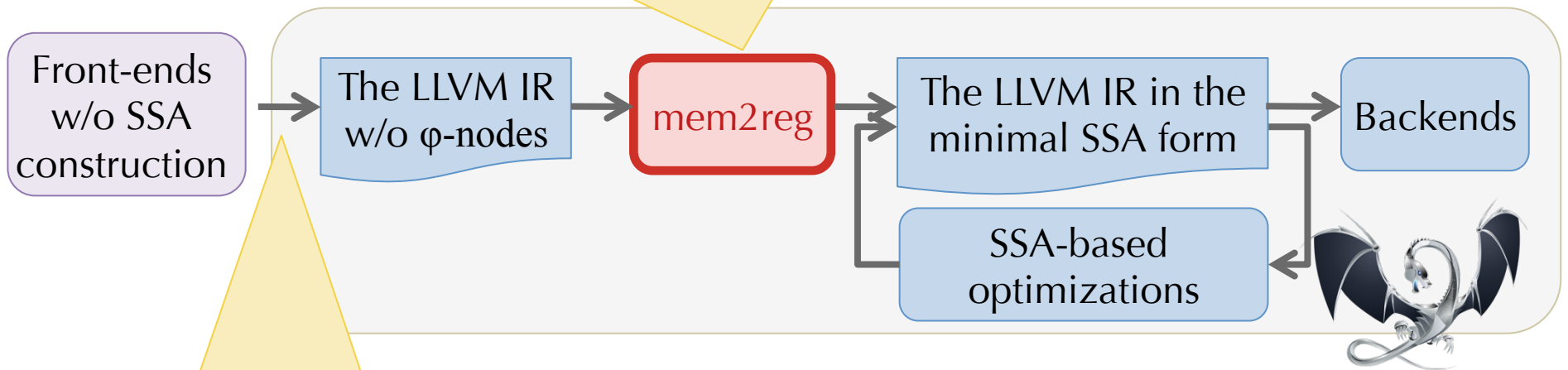
$\text{sdom}(f, \text{pc})$  = variable defns. that *strictly dominate* pc.

$\text{wf}(f, \sigma) = \forall r \in \text{sdom}(f, \text{pc}). \exists v. \delta(r) = \lfloor v \rfloor$

“All variables in scope are initialized.”

# mem2reg in LLVM (part of SROA)

- Promote stack allocas to temporaries
- Insert minimal  $\phi$ -nodes



- imperative variables  $\Rightarrow$  stack allocas
- no  $\phi$ -nodes
- trivially in SSA form

# mem2reg Example

```
int x = 0;  
if (y > 0)  
    x = 1;  
return x;
```

```
l1: %p = alloca i32  
      store 0, %p  
      %b = %y > 0  
      br %b, %l2, %l3
```

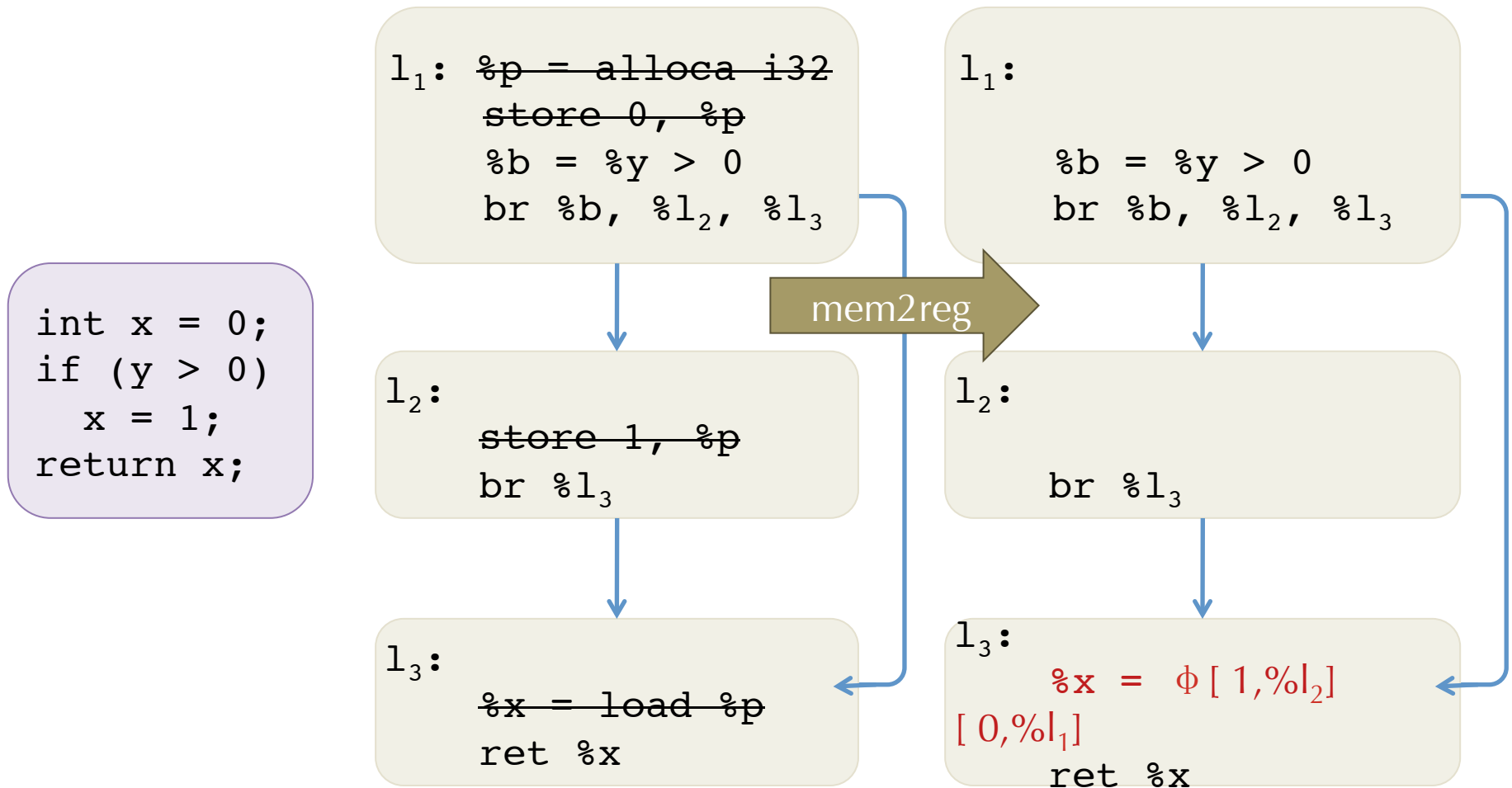
```
l2:  
      store 1, %p  
      br %l3
```

```
l3:  
      %x = load %p  
      ret %x
```

The LLVM IR in the trivial SSA form



# mem2reg Example



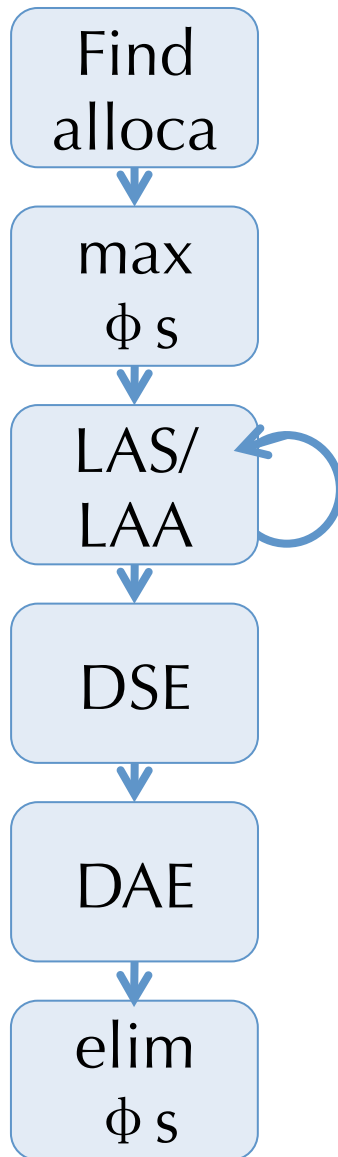
The LLVM IR in the trivial SSA form

Minimal SSA after mem2reg

# mem2reg Algorithm

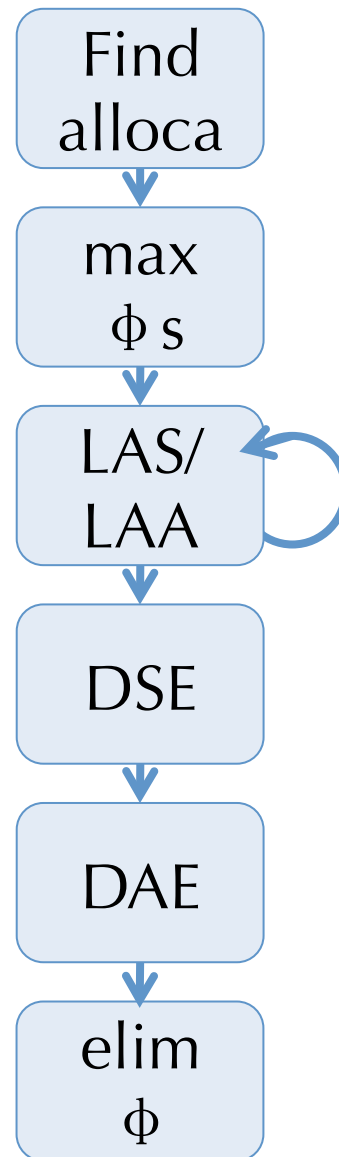
- Two main operations
  - Phi placement (Lengauer-Tarjan algorithm)
  - Renaming of the variables
- Intermediate stage breaks SSA invariant
  - Defining semantics & well formedness non-trivial

# vmem2reg Algorithm

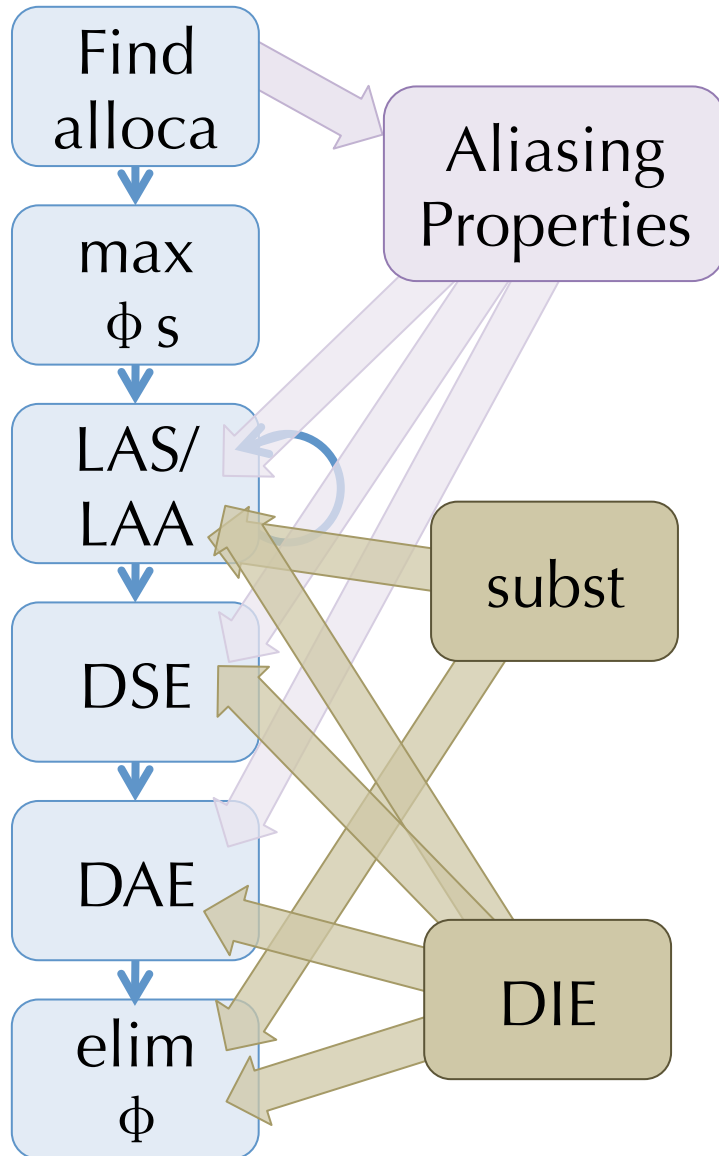


- Incremental algorithm
- Pipeline of micro-transformations
  - Preserves SSA semantics
  - Preserves well-formedness
- Inspired by Aycok & Horspool 2002.

# How to Establish Correctness?



# How to Establish Correctness?



1. Simple aliasing properties (e.g. to determine promotability)
2. Instantiate proof technique for
  - Substitution
  - Dead Instruction Elimination

$P_{DIE} = \dots$   
Initialize( $P_{DIE}$ )  
Preservation( $P_{DIE}$ )  
Progress( $P_{DIE}$ )
4. Put it all together to prove composition of “pipeline” correct.

# vmem2reg is Correct

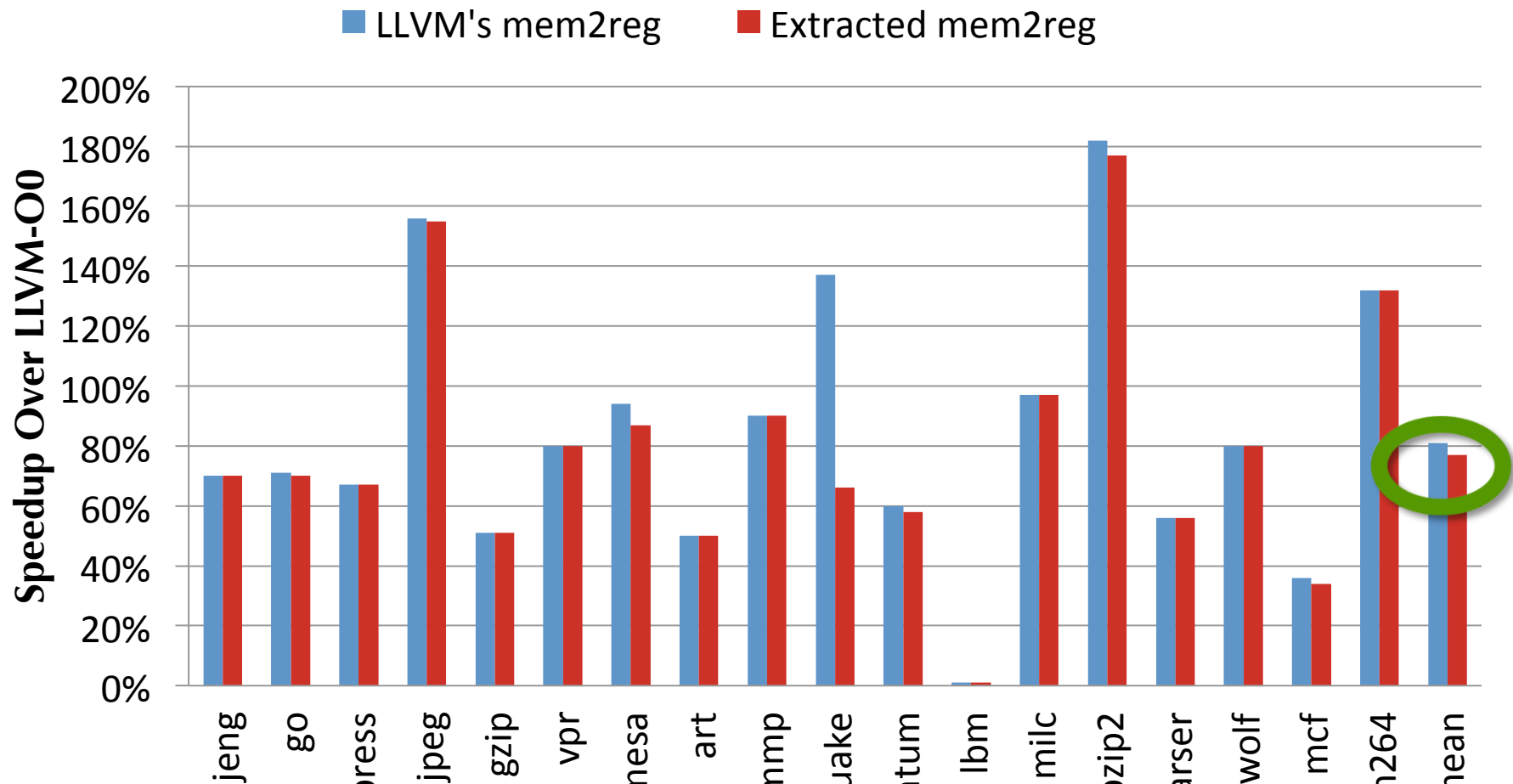
Theorem: The vmem2reg algorithm preserves the semantics of the source program.

Proof:

Composition of simulation relations from the “mini” transformations, each built using instances of the sdom proof technique.

(See Coq Vellvm development.)  $\square$

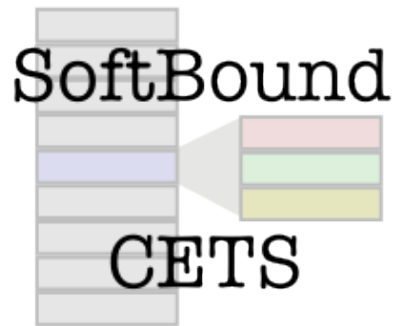
# Runtime overhead of verified mem2reg



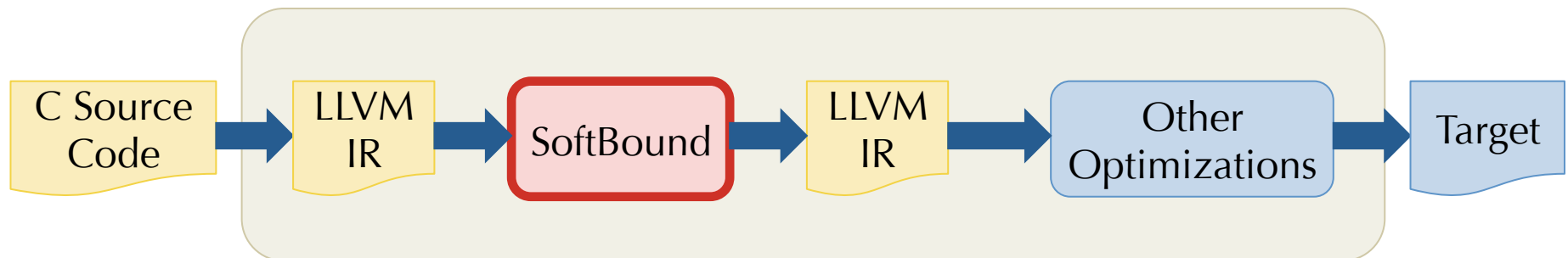
**Vmem2reg: 77% LLVM's mem2reg: 81%**

(LLVM's mem2reg promotes allocas used by  
intrinsic)

# SoftBound



- Implemented as an LLVM pass.
- Detect spatial/temporal memory safety violations in legacy C code.
- Good test case:
  - Safety Critical  $\Rightarrow$  Proof cost warranted
  - Non-trivial Memory transformation





# SoftBound

```
%p = call malloc [10 x i8]
```

Maintain base and bound for all pointers

```
%q = gep %p, i32 0, i32 255
```

Propagate metadata on assignment

Check that a pointer is within its bounds when being accessed

```
store i8 0, %q
```

```
%p = call malloc [10 x i8]
```

```
%p_base = gep %p, i32 0
```

```
%p_bound = gep %p, i32 0, i32 10
```

```
%q = gep %p, i32 0, i32 255
```

```
%q_base = %p_base
```

```
%q_bound = %p_bound
```

```
assert %q_base <= %q
```

```
    /\ %q+1 < %q_bound
```

```
store i8 0, %q
```

C Source Code

LLVM IR

SoftBound

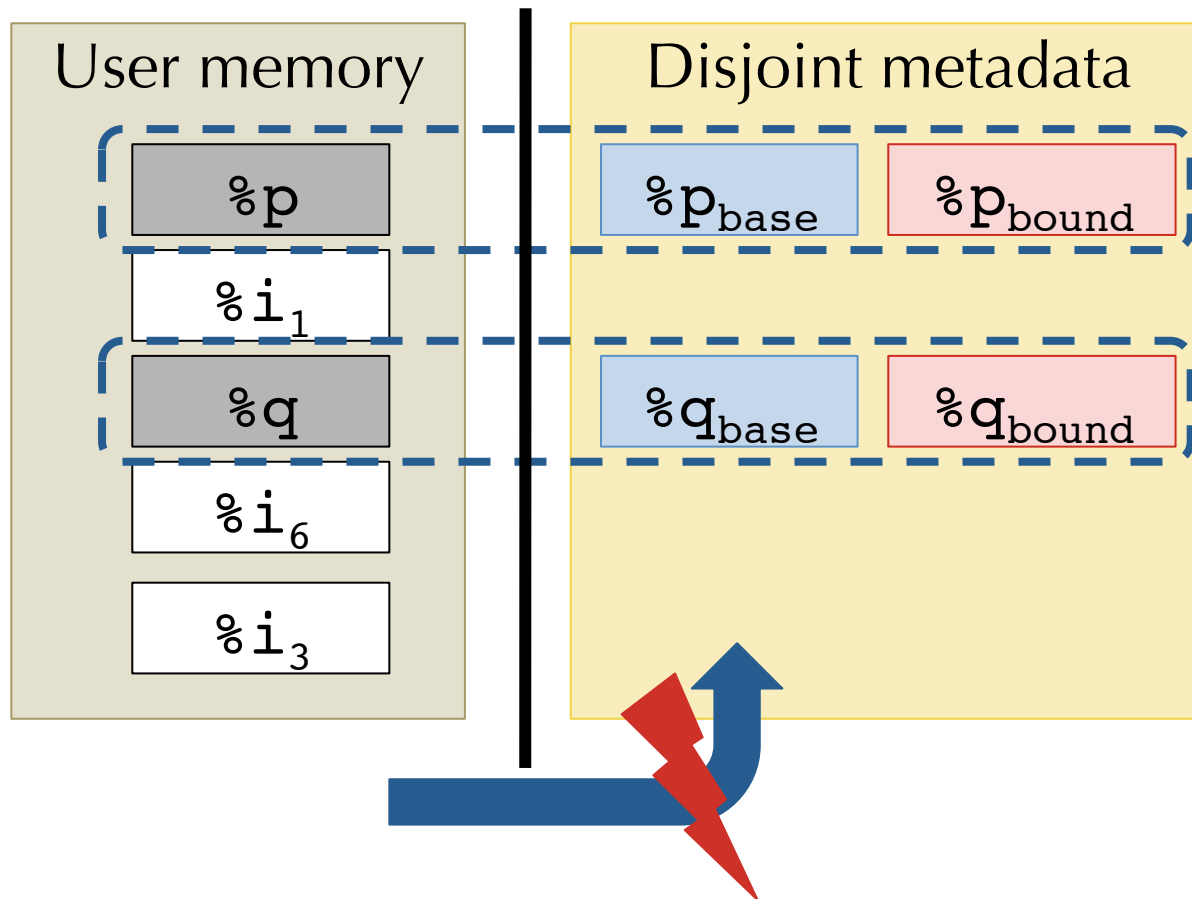
LLVM IR

Other Optimizations

Target

# Disjoint Metadata

- Maintain pointer bounds in a separate memory space.
- Key Invariant: Metadata cannot be corrupted by bounds violation.

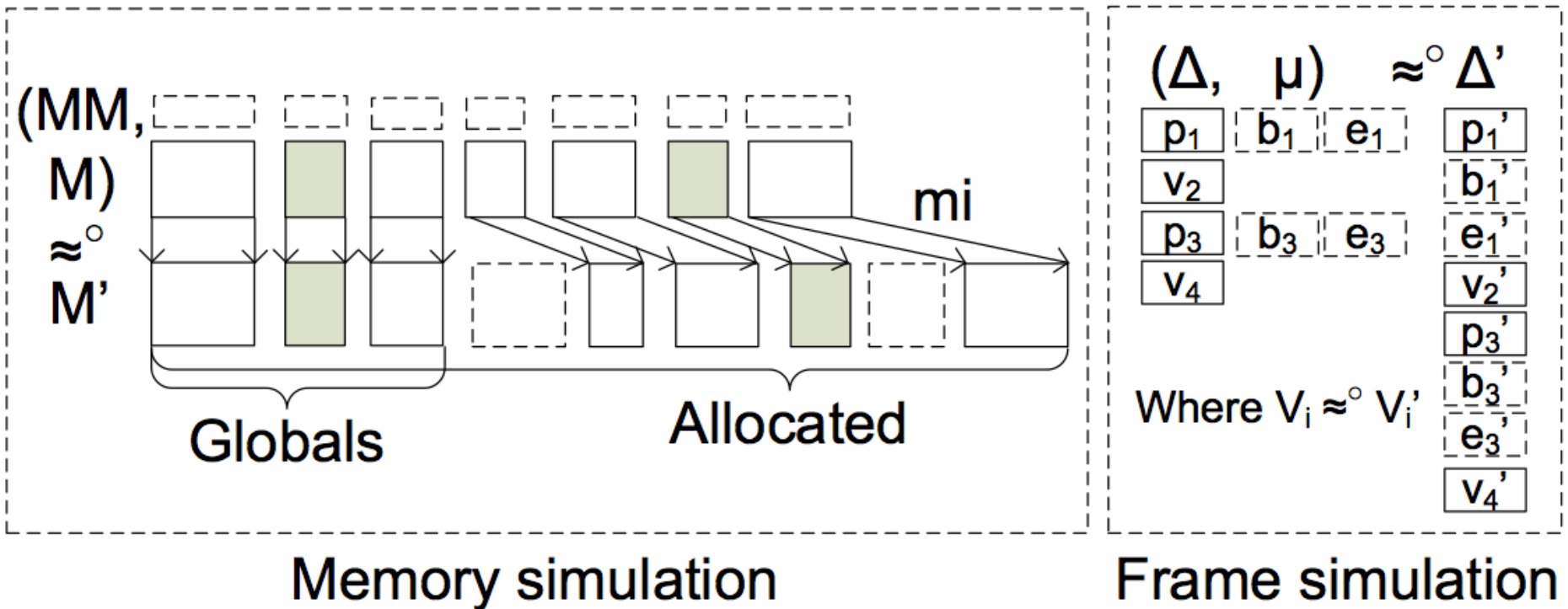


# Proving SoftBound Correct

1. Define  $\text{SoftBound}(f, \sigma) = (f_s, \sigma_s)$ 
  - Transformation pass implemented in Coq.
2. Define predicate:  $\text{MemoryViolation}(f, \sigma)$
3. Construct a *non-standard* operational semantics:
$$f \vdash \sigma \xrightarrow{\text{SB}} \sigma'$$
  - Builds in safety invariants “by construction”
$$f \vdash \sigma \xrightarrow{\text{SB}}^* \sigma' \Rightarrow \neg \text{MemoryViolation}(f, \sigma')$$
4. Show that the instrumented code simulates the “correct” code:

$$\text{SoftBound}(f, \sigma) = (f_s, \sigma_s) \Rightarrow [f \vdash_B^S \sigma \xrightarrow{*} \sigma'] \approx [f_s \vdash \sigma_s \xrightarrow{*} \sigma'_s]$$

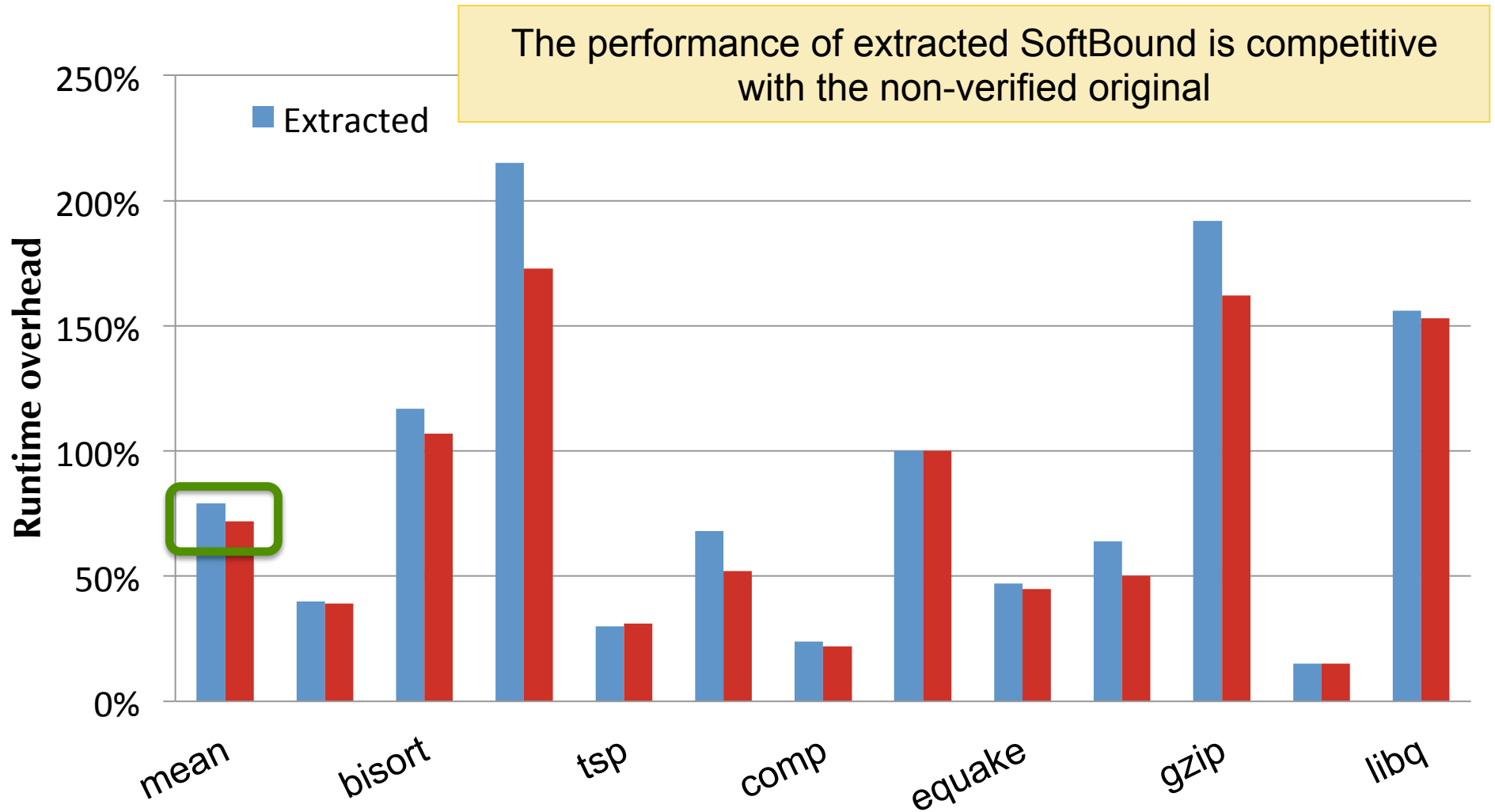
# Memory Simulation Relation



# Lessons About SoftBound

- Found several bugs in our C++ implementation
  - Interaction of undef, 'null', and metadata initialization.
- Simulation proofs suggested a redesign of SoftBound's handling of stack pointers.
  - Use a "shadow stack"
  - Simplify the design/implementation
  - Significantly more robust (e.g. varargs)

# Competitive Runtime Overhead



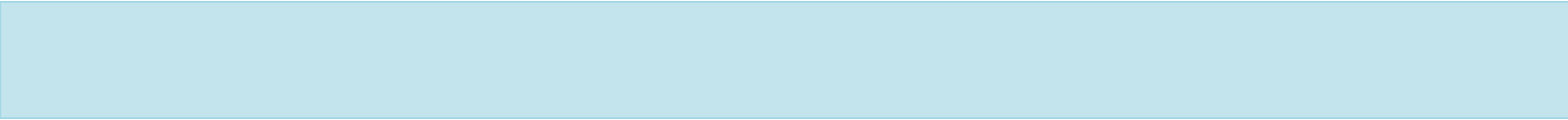


# FINAL EXAM

# Final Exam

- Will cover material since the midterm almost exclusively
  - Starting from Lecture 14
  - Objects, inheritance, types, implementation of dynamic dispatch
  - Basic optimizations
  - Dataflow analysis (forward vs. backward, fixpoint computations, etc.)
    - Liveness
  - Control flow analysis
    - Loops, dominator trees
  - SSA
  - Graph-coloring Register Allocation
- Will focus more on the theory side of things
- Format will be similar to the midterm
  - Simple answer, computation, multiple choice, etc.
  - Sample exam from last time is on the web





What have we learned?  
Where else is it applicable?  
What next?

# COURSE WRAP-UP

# Why CIS 341?

- You will learn:
  - Practical applications of theory
  - Parsing
  - How high-level languages are implemented in machine language
  - (A subset of) Intel x86 architecture
  - A deeper understanding of code
  - A little about programming language semantics
  - Functional programming in OCaml
  - How to manipulate complex data structures
  - How to be a better programmer
- Did we meet these goals?

# Stuff we didn't Cover

- We skipped stuff at every level...
- Concrete syntax/parsing:
  - Much more to the theory of parsing...
  - Good syntax is art not science!
- Source language features:
  - Exceptions, recursive data types (easy!), advanced type systems, type inference, concurrency
- Intermediate languages:
  - Intermediate language design, bytecode, bytecode interpreters, just-in-time compilation (JIT)
- Compilation:
  - Continuation-passing transformation, efficient representations, scalability
- Optimization:
  - Scientific computing, cache optimization, instruction selection/optimization

# Course Work

- 72% Projects: *The Quaker OAT Compiler*
- 12% Midterm
- 16% Final exam
- Expect this to be a challenging, implementation-oriented course.



I think we met this goal...

# Related Courses: Fall 2013

- CIS 500: Software Foundations
  - Dr. Pierce
  - Theoretical course about functional programming, proving program properties, type systems, lambda calculus. Uses the theorem prover Coq.
- CIS 501: Computer Architecture
  - Dr. Devietti
  - 371++: pipelining, caches, VM, superscalar, multicore,...
- CIS 552: Advanced Programming
  - Dr. Weirich
  - Advanced functional programming in Haskell, including generic programming, metaprogramming, embedded languages, cool tricks with fancy type systems
- CIS 670: Special topics in programming languages
  - TBA

# Where to go from here?

- Conferences (proceedings available on the web):
  - Programming Language Design and Implementation (PLDI)
  - Principles of Programming Languages (POPL)
  - Object Oriented Programming Systems, Languages & Applications (OOPSLA)
  - International Conference on Functional Programming (ICFP)
  - European Symposium on Programming (ESOP)
  - ...
- Technologies / Open Source Projects
  - Yacc, lex, bison, flex, ...
  - LLVM – low level virtual machine
  - Java virtual machine (JVM), Microsoft's Common Language Runtime (CLR)
  - Languages: OCaml, F#, Haskell, Scala, Go, Rust, ...?

# Where else is this stuff applicable?

- General programming
  - In C/C++, better understanding of how the compiler works can help you generate better code.
  - Ability to read assembly output from compiler
  - Experience with functional programming can give you different ways to think about how to solve a problem
- Writing domain specific languages
  - lex/yacc very useful for little utilities
  - understanding abstract syntax and interpretation
- Understanding hardware/software interface
  - Different devices have different instruction sets, programming models

# Thanks!

- To the TAs: Dmitri, Rohan, and Mitchell
  - for doing an amazing job putting together the projects for the course.
- To *you* for taking the class!
  
- How can I improve the course?
  - Feedback survey posted to Piazza