
Verifying the LLVM

Steve Zdancewic

DeepSpec Summer School 2017



Substitution.v

SUBSTITUTION

Generalizing Safety

- Definition of wf:

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

- Generalize like this:

$$\text{wf}(f, (\text{pc}, \delta)) = \mathbf{P} f (\delta \upharpoonright_{\text{sdom}(f, \text{pc})})$$

where $\mathbf{P} : \text{Program} \rightarrow \text{Local} \rightarrow \text{Prop}$

- Methodology: for a given \mathbf{P} prove that

Initialization(\mathbf{P})
Preservation(\mathbf{P})
Progress(\mathbf{P})

Consider only variables in scope $\Rightarrow \mathbf{P}$ defined relative to the dominator tree of the CFG.

Instantiating

- For usual safety:

$$P_{\text{safety}} f \delta = \forall r \in \text{dom}(\delta). \exists v. \delta(r) = \llbracket v \rrbracket$$

- For semantic properties:

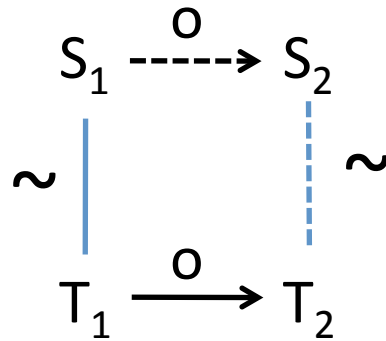
$$P_{\text{sem}} f \delta = \forall r. f[r] = \llbracket \text{rhs} \rrbracket \Rightarrow \delta(r) = \llbracket \text{rhs} \rrbracket_{\delta}$$

- Useful for creating the simulation relation for correctness of:
 - code motion, dead variable elimination, common expression elimination, etc.

Is Backward Simulation Hopeless?

- Suppose the source & target languages are the same.
 - So they share the same definition of program state.
- Further suppose that the steps are very “small”.
 - Abstract machine (i.e. no “complex” instructions).
- Further suppose that “compilation” is only a very minor change.
 - add or remove a single instruction
 - substitute a value for a variable
- Then: backward simulation is more achievable
 - it’s easier to invent the “decompilation” function because the “compilation” function is close to trivial
- Happily: This is the situation for some LLVM transformations

Lock-Step Backward Simulation

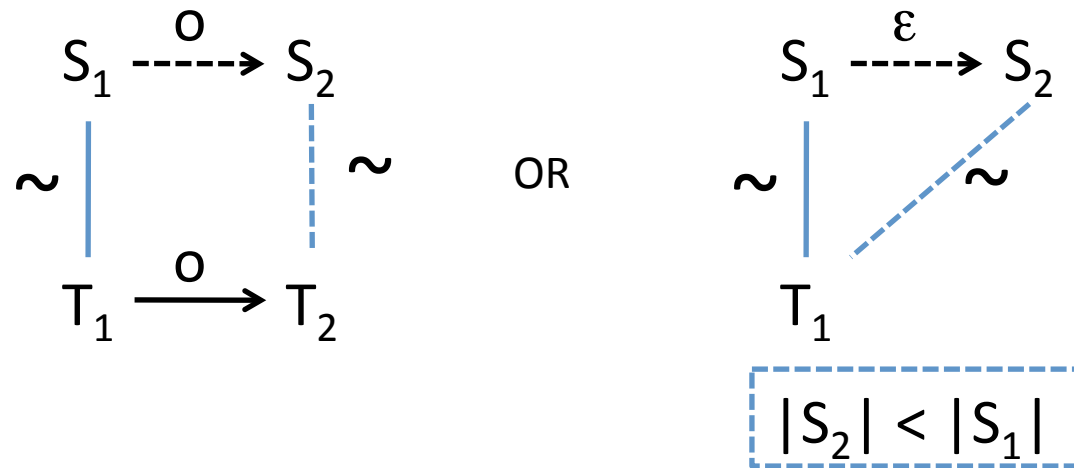


o is either an “observable event” or a “silent event”

$o ::= e \mid \varepsilon$

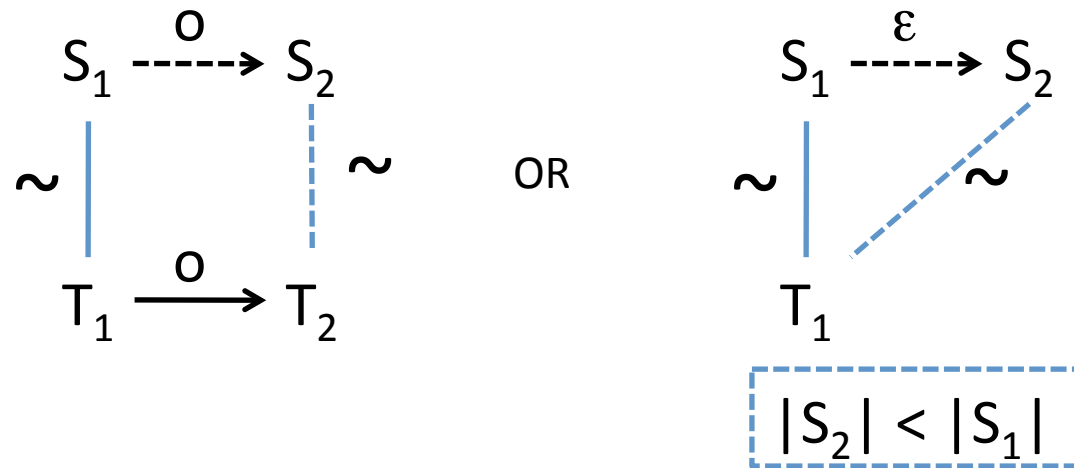
Example use: proving variable substitution correct.

Right-Option Backward Simulation



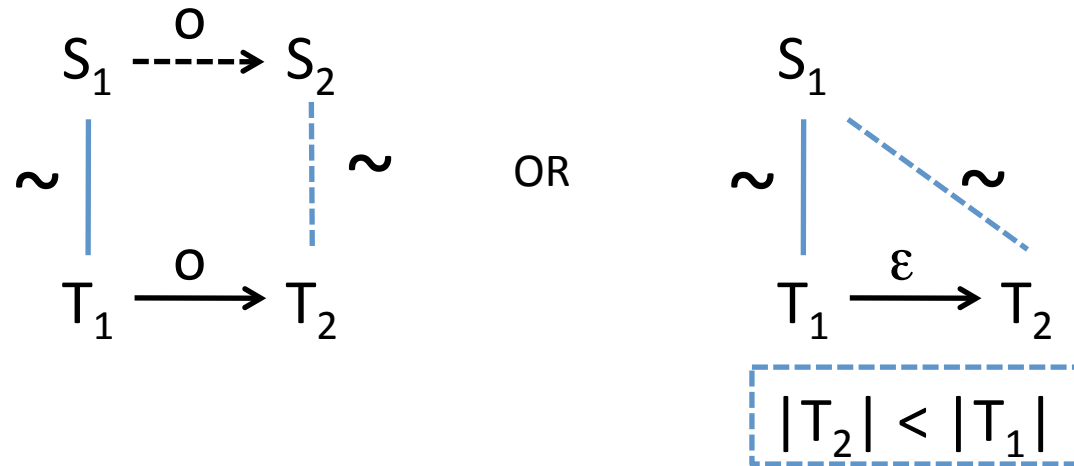
- Either:
 - the source and target are in lock-step simulation.
 - Or
 - the source takes a silent transition to a smaller state
- Example use: removing an instruction in the target.

Right-Option Backward Simulation



- Either:
 - the source and target are in lock-step simulation.
 - Or
 - the source takes a silent transition to a smaller state
- Example use: removing an instruction in the target.

Left-Option Backward Simulation



- Either:
 - the source and target are in lock-step simulation.
 - Or
 - the target takes a silent transition to a smaller state
- Example use: adding an instruction to the target.

Strategy for Proving Transformations

- 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 generalization of the Vminus safety property
- Compose the individual proofs of correctness

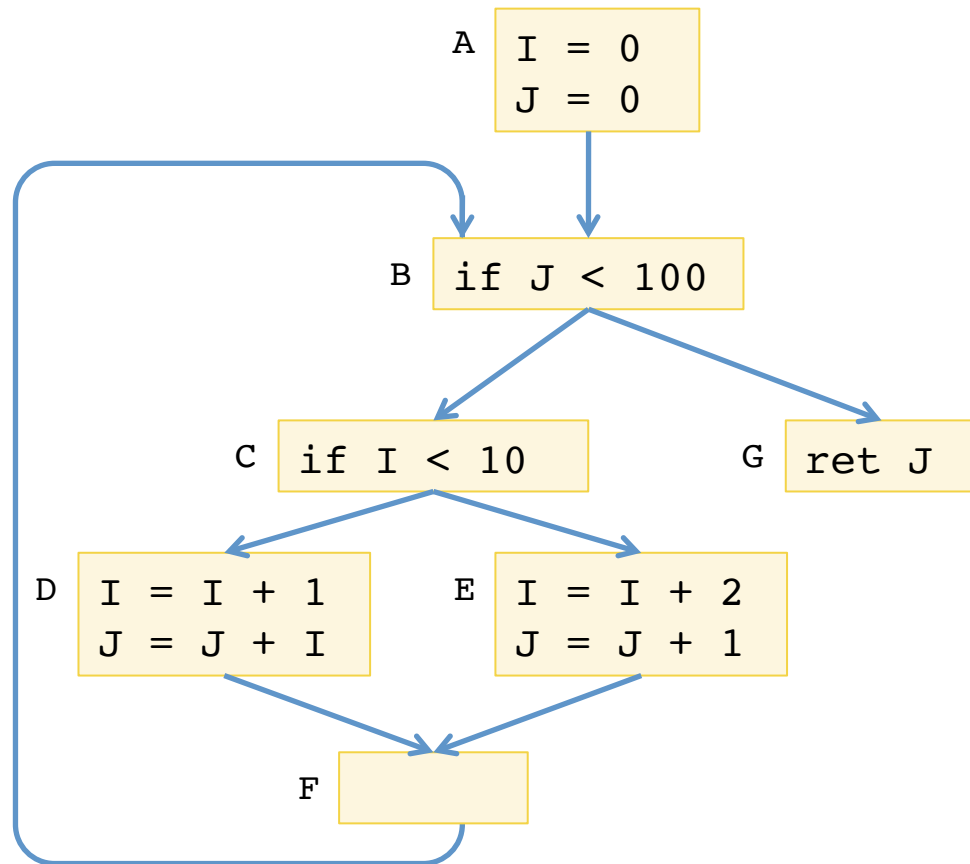
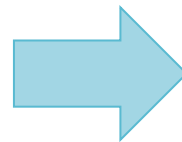
SSA CONSTRUCTION

SSA Construction by Example

```
I := 0;;
J := 0;;
WHILE J < 100 DO
  IF I < 10 THEN
    I := I + 1;;
    J := J + I
  ELSE
    I := I + 2;;
    J := J + 1;
  FI
END;;
RETURN J
```

SSA Construction by Example

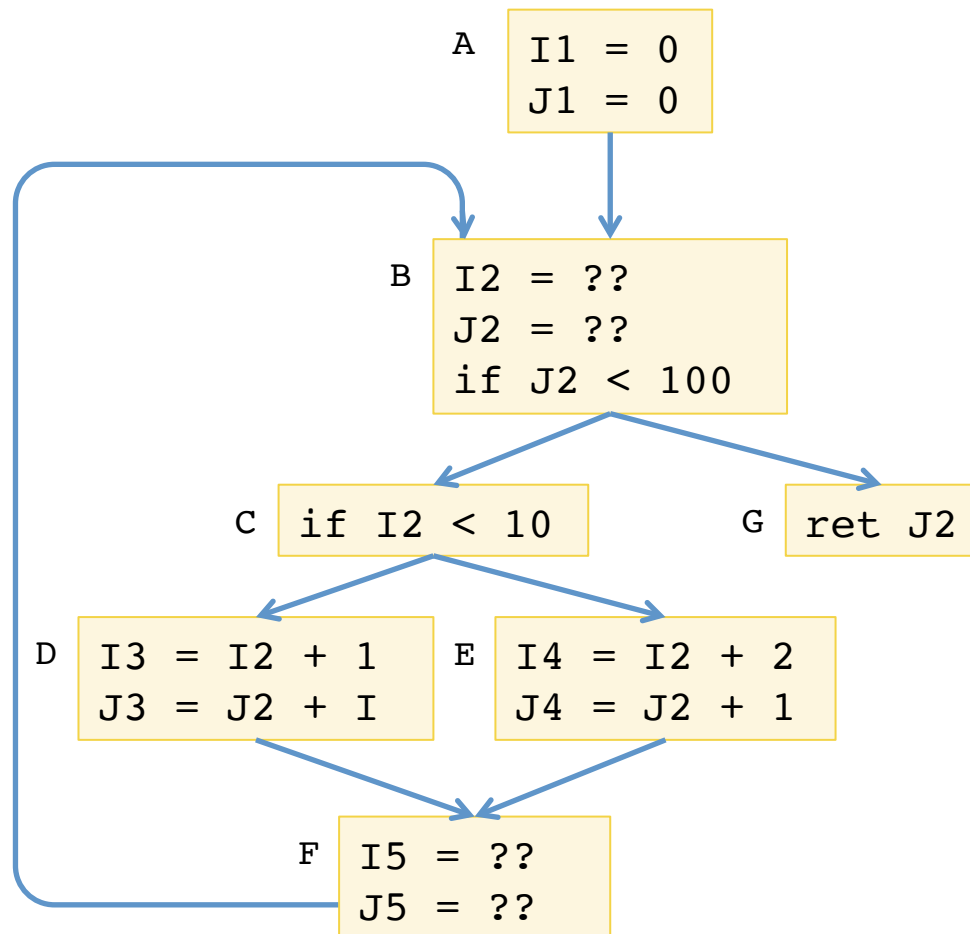
```
I := 0;;  
J := 0;;  
WHILE J < 100 DO  
  IF I < 10 THEN  
    I := I + 1;;  
    J := J + I  
  ELSE  
    I := I + 2;;  
    J := J + 1;  
  FI  
END;;  
RETURN J
```



Step 1: Convert to a control-flow graph.

SSA Construction by Example

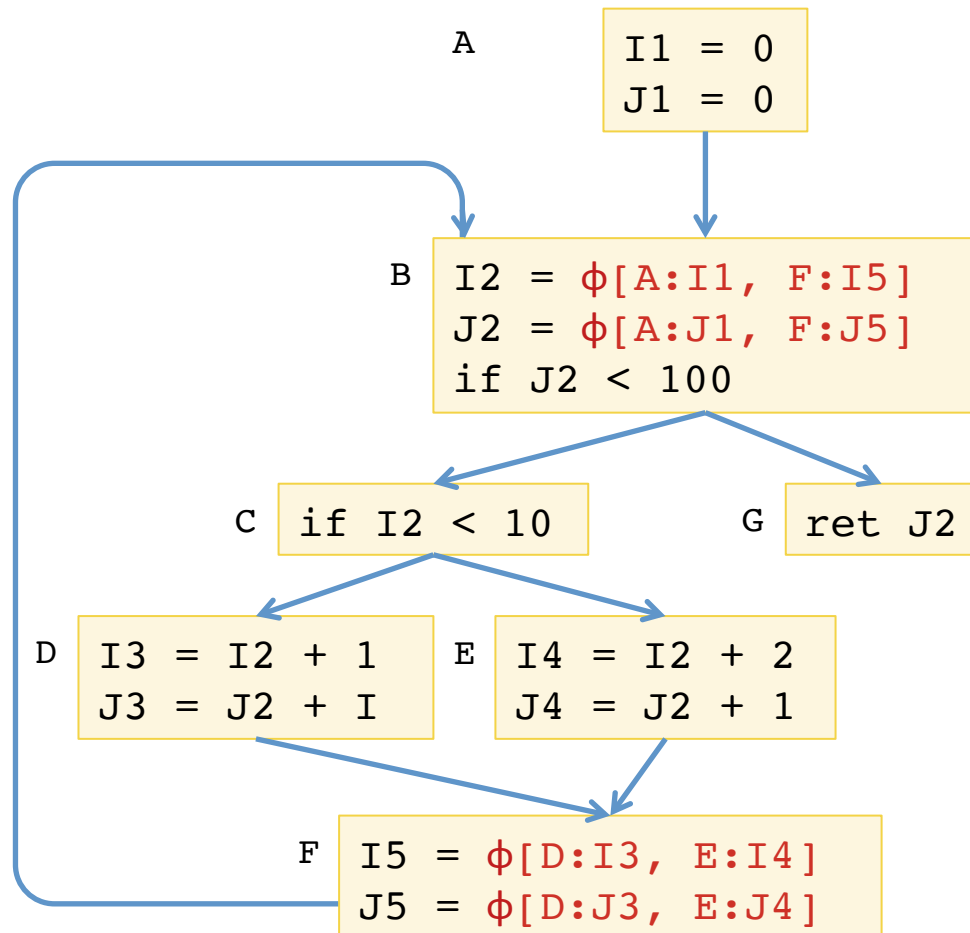
```
I := 0;;  
J := 0;;  
WHILE J < 100 DO  
  IF I < 10 THEN  
    I := I + 1;;  
    J := J + I  
  ELSE  
    I := I + 2;;  
    J := J + 1;  
  FI  
END;;  
RETURN J
```



Step 2: Rename variables to satisfy single assignment.

SSA Construction by Example

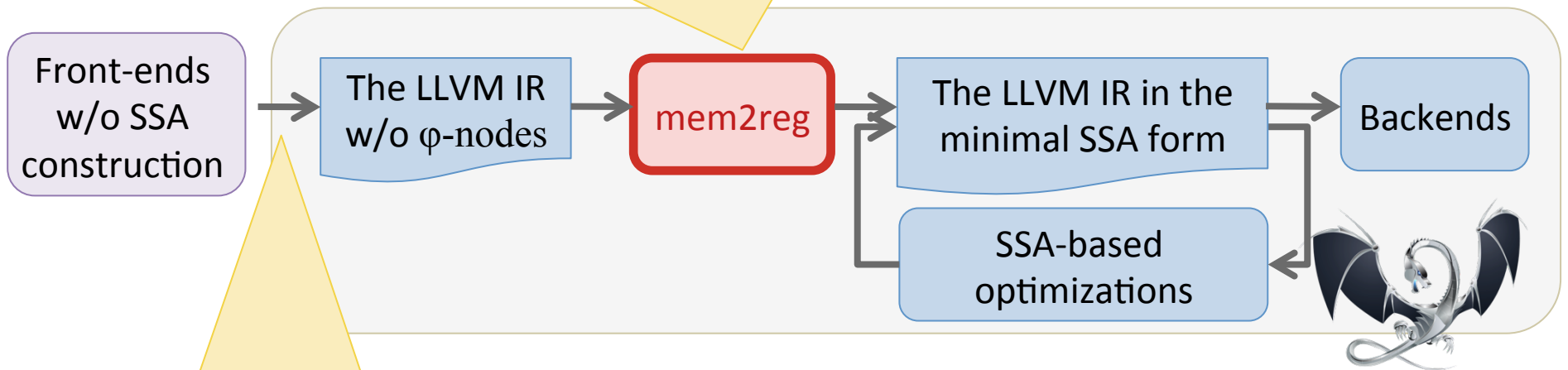
```
I := 0;;  
J := 0;;  
WHILE J < 100 DO  
  IF I < 10 THEN  
    I := I + 1;;  
    J := J + I  
  ELSE  
    I := I + 2;;  
    J := J + 1;  
  FI  
END;;  
RETURN J
```



Step 3: Insert “ ϕ ” functions that capture control dependence.

mem2reg in LLVM

- Promote stack allocas to temporaries
- Insert minimal ϕ -nodes



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

mem2reg Example

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

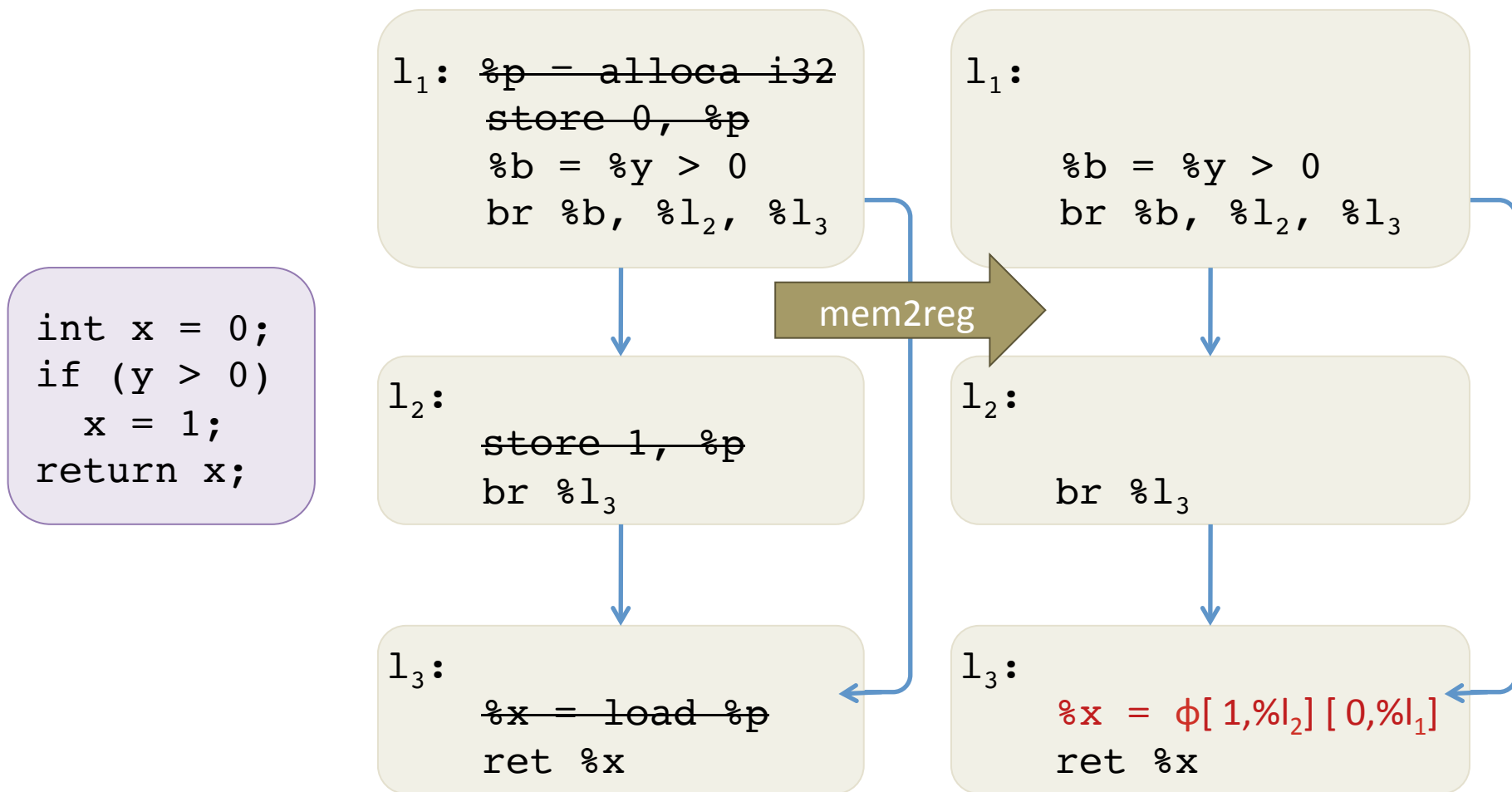
```
l1: %p = alloca i32  
      store 0, %p  
      %b = %y > 0  
      br %b, %l2, %l3
```

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

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

The LLVM IR in the trivial SSA form

mem2reg Example



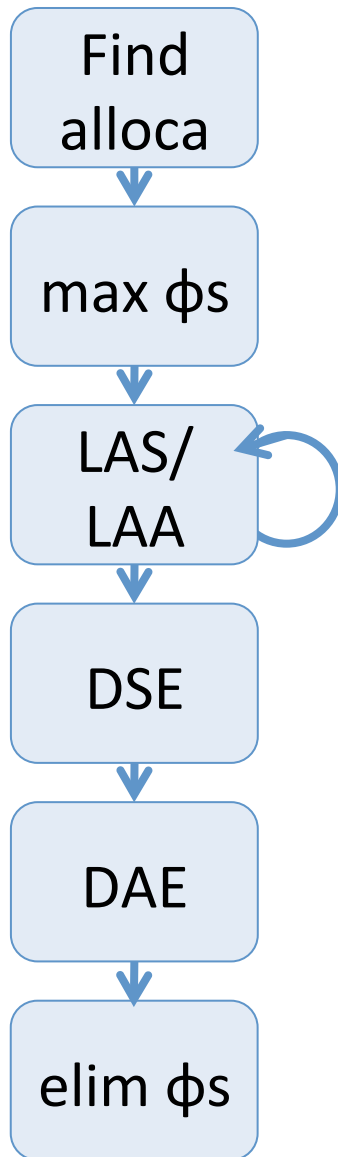
The LLVM IR in the trivial SSA form

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

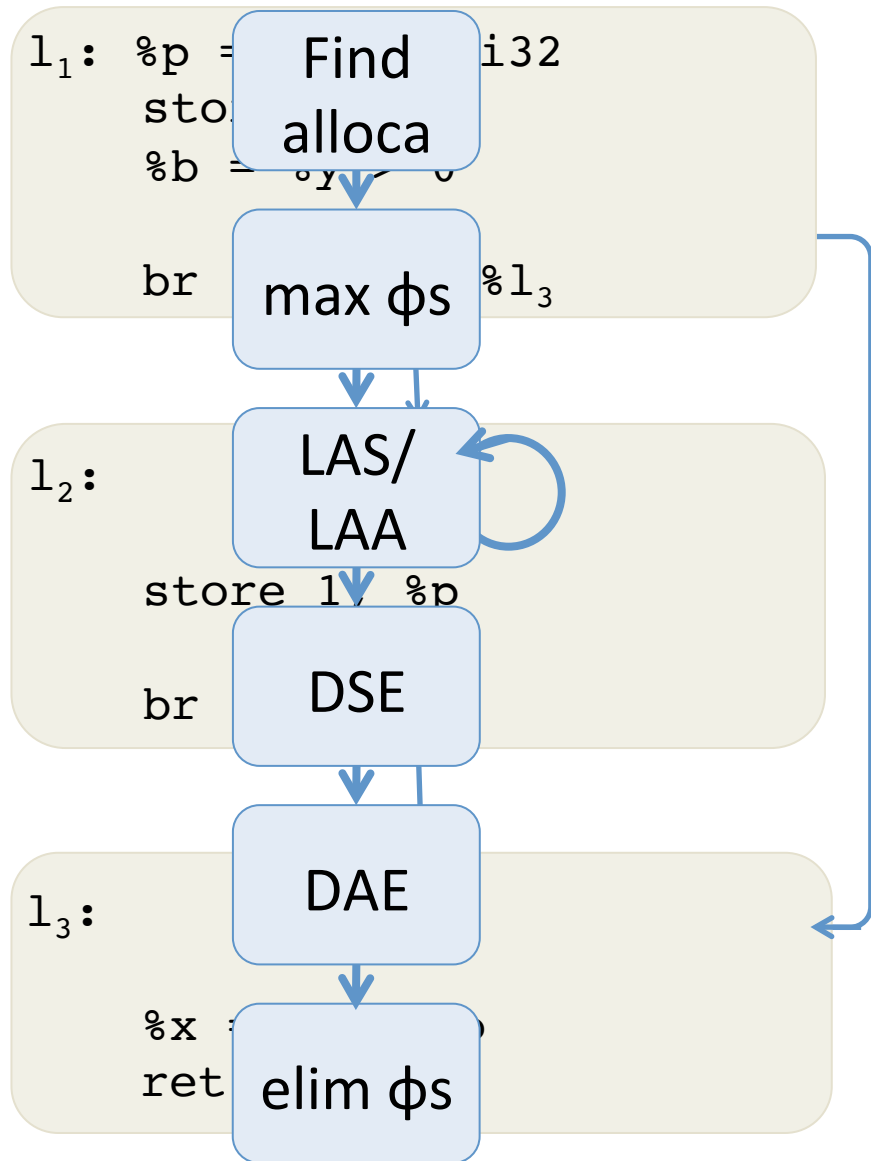
vmem2reg Algorithm



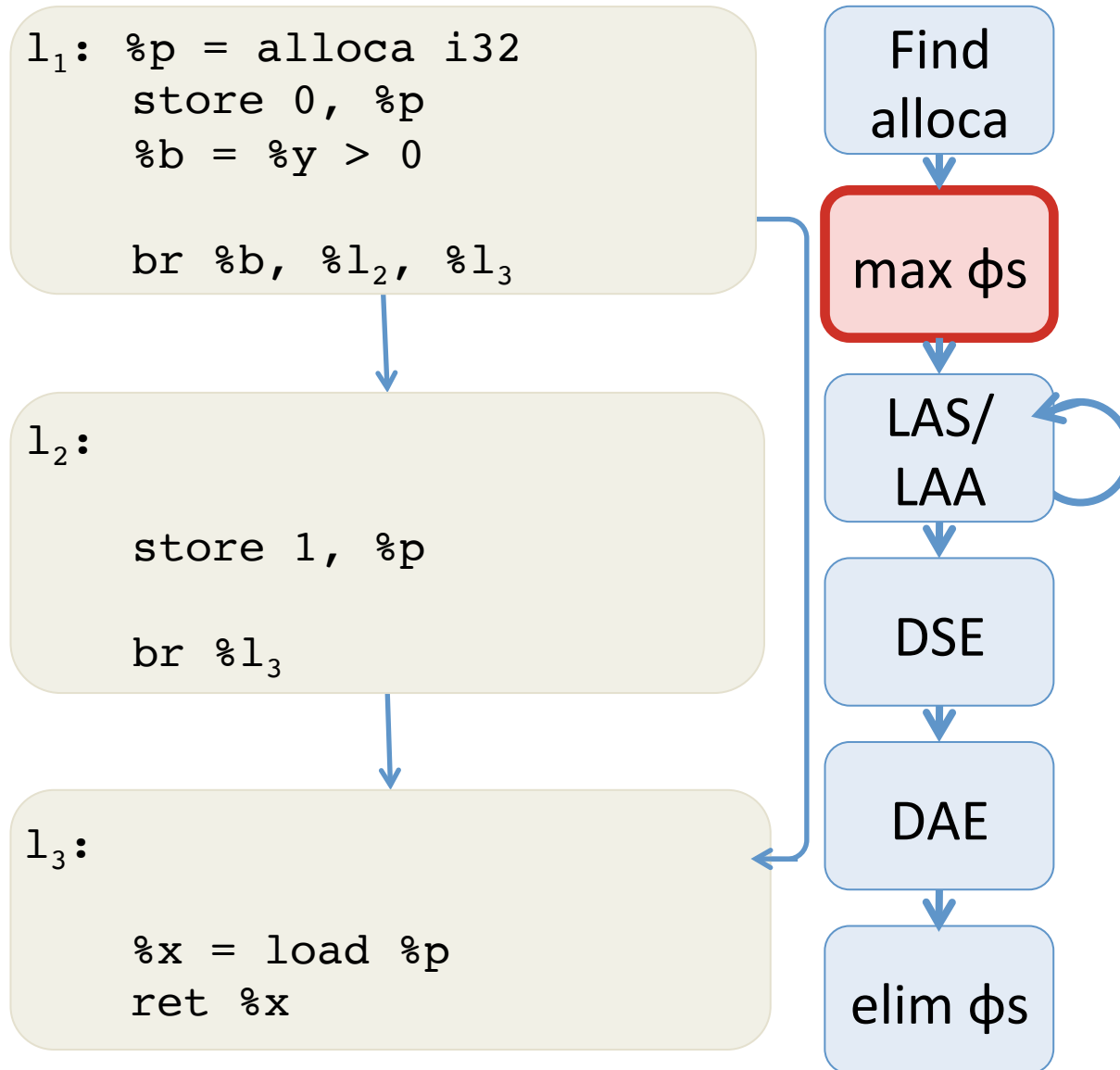
- Incremental algorithm
- Pipeline of micro-transformations
 - Preserves SSA semantics
 - Preserves well-formedness

- Inspired by Aycock & Horspool 2002.

Example of vmem2reg Algorithm

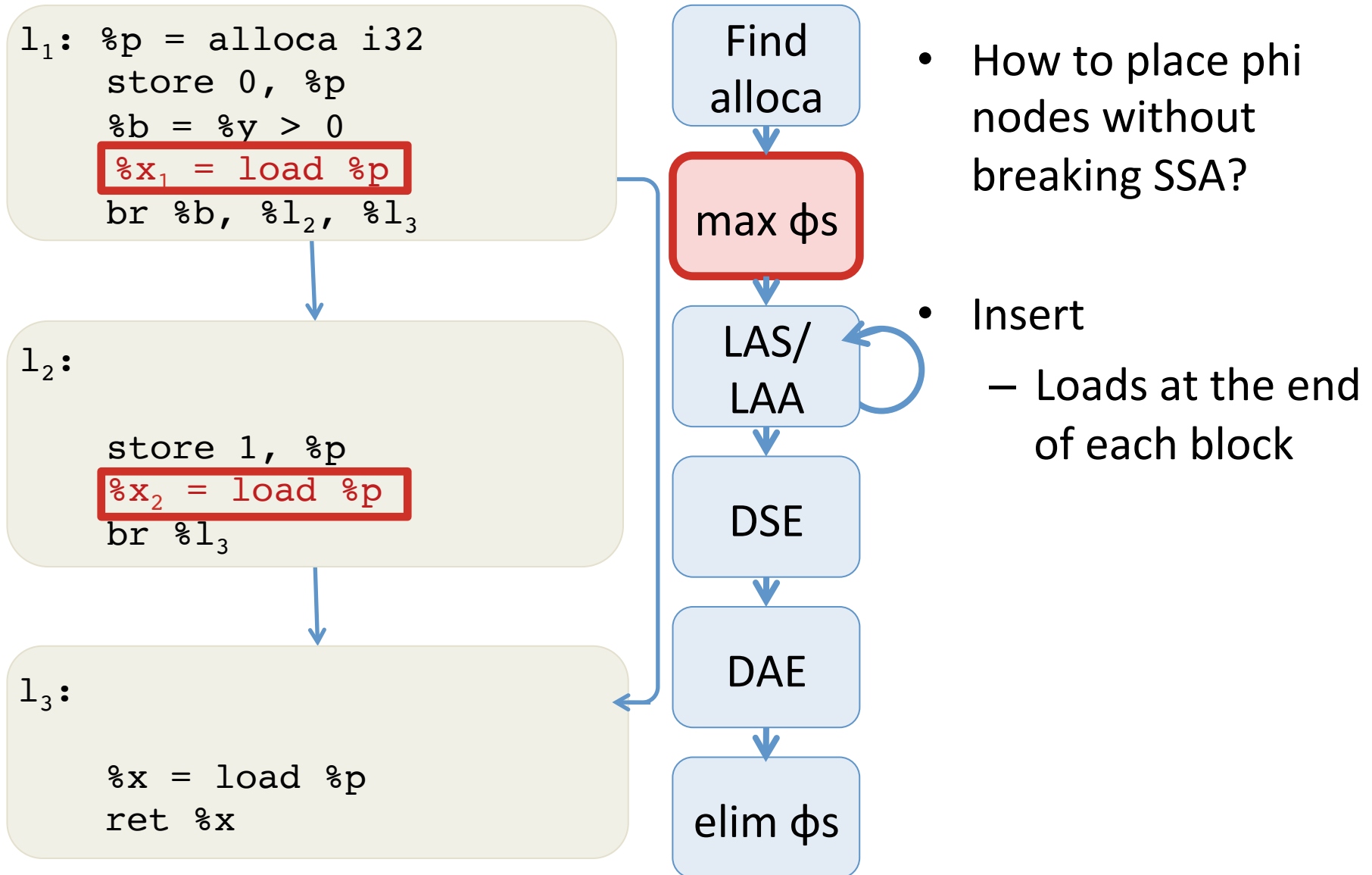


Example of vmem2reg Algorithm

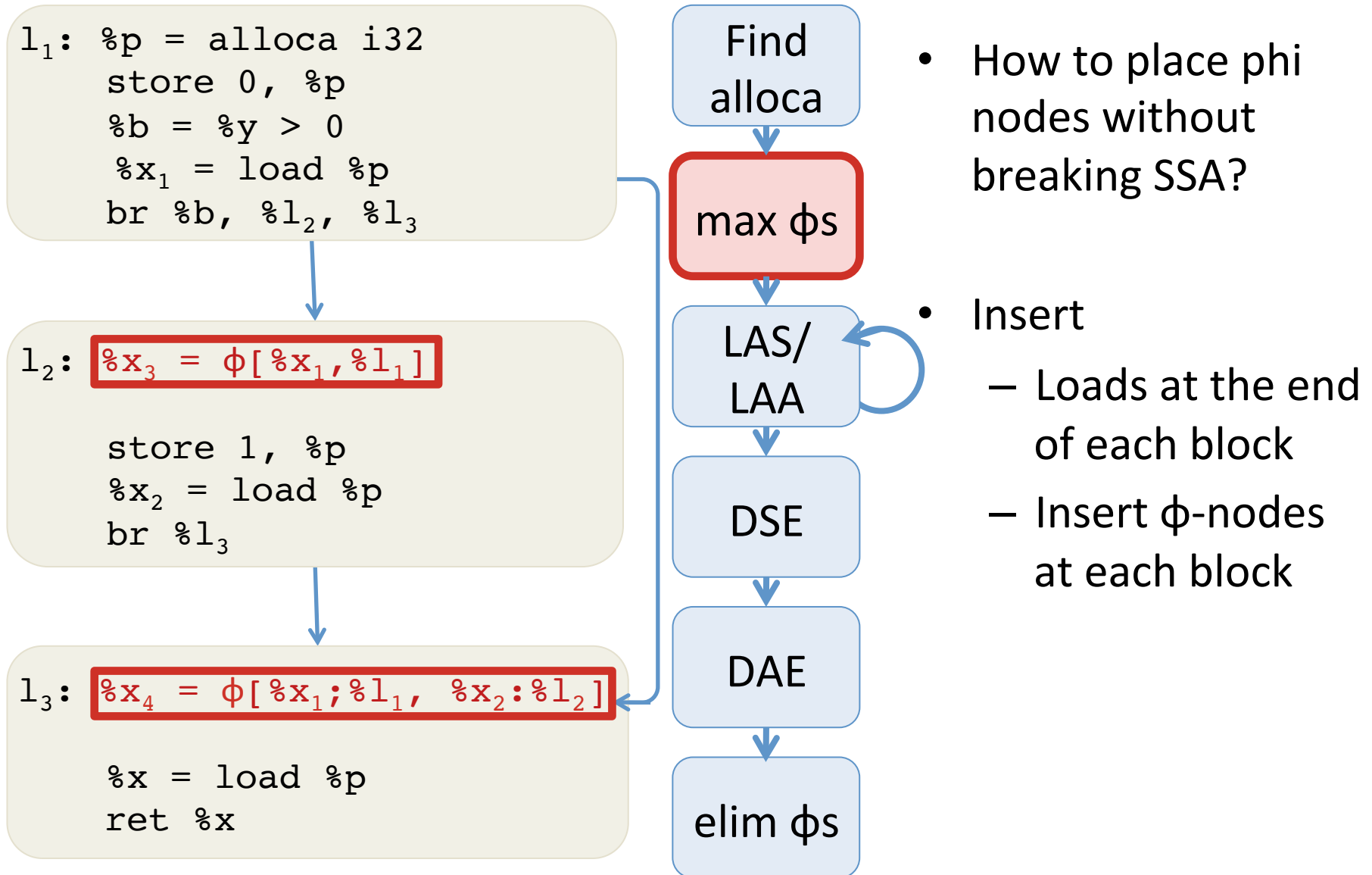


- How to place phi nodes without breaking SSA?

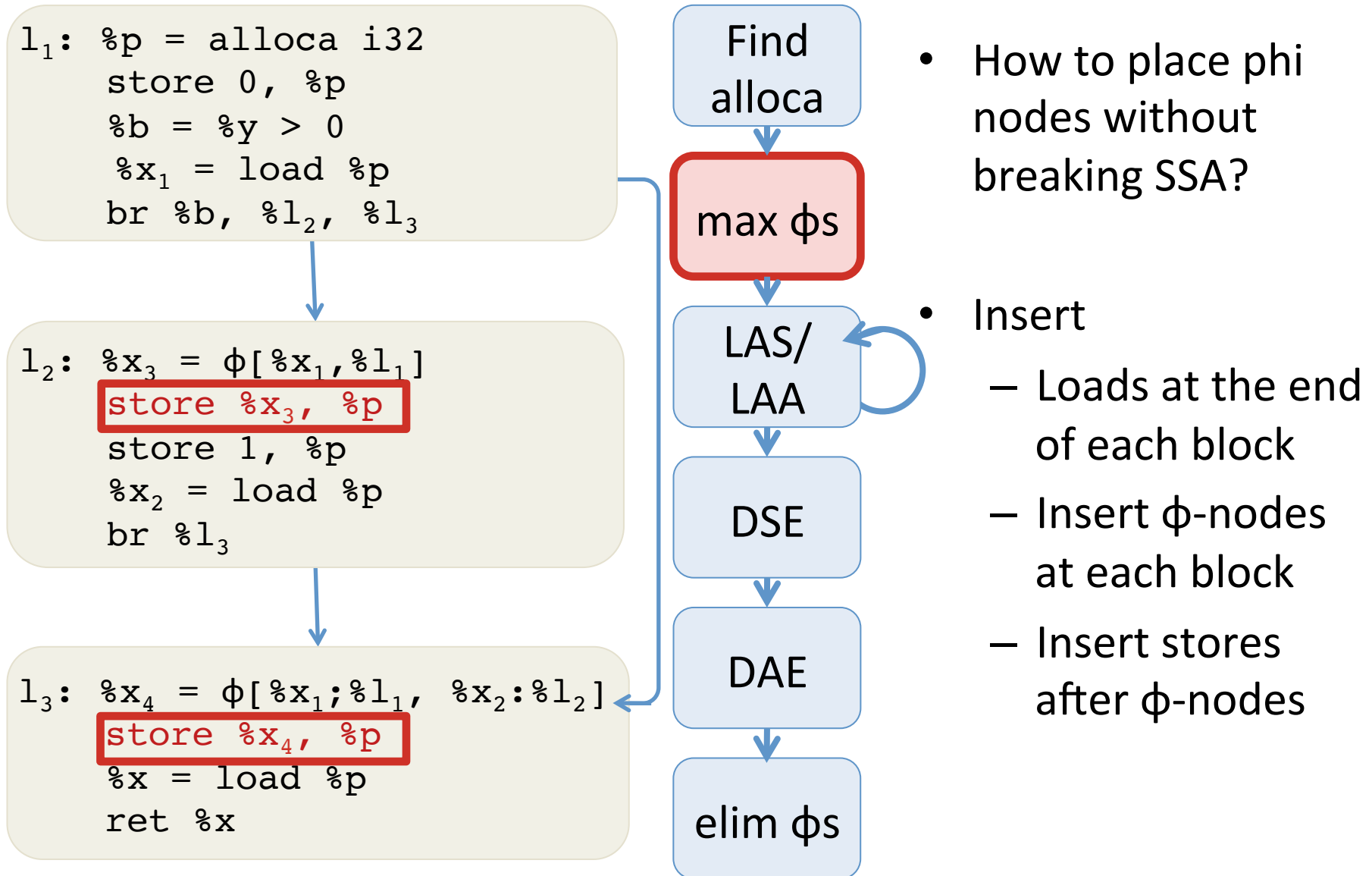
Example of vmem2reg Algorithm



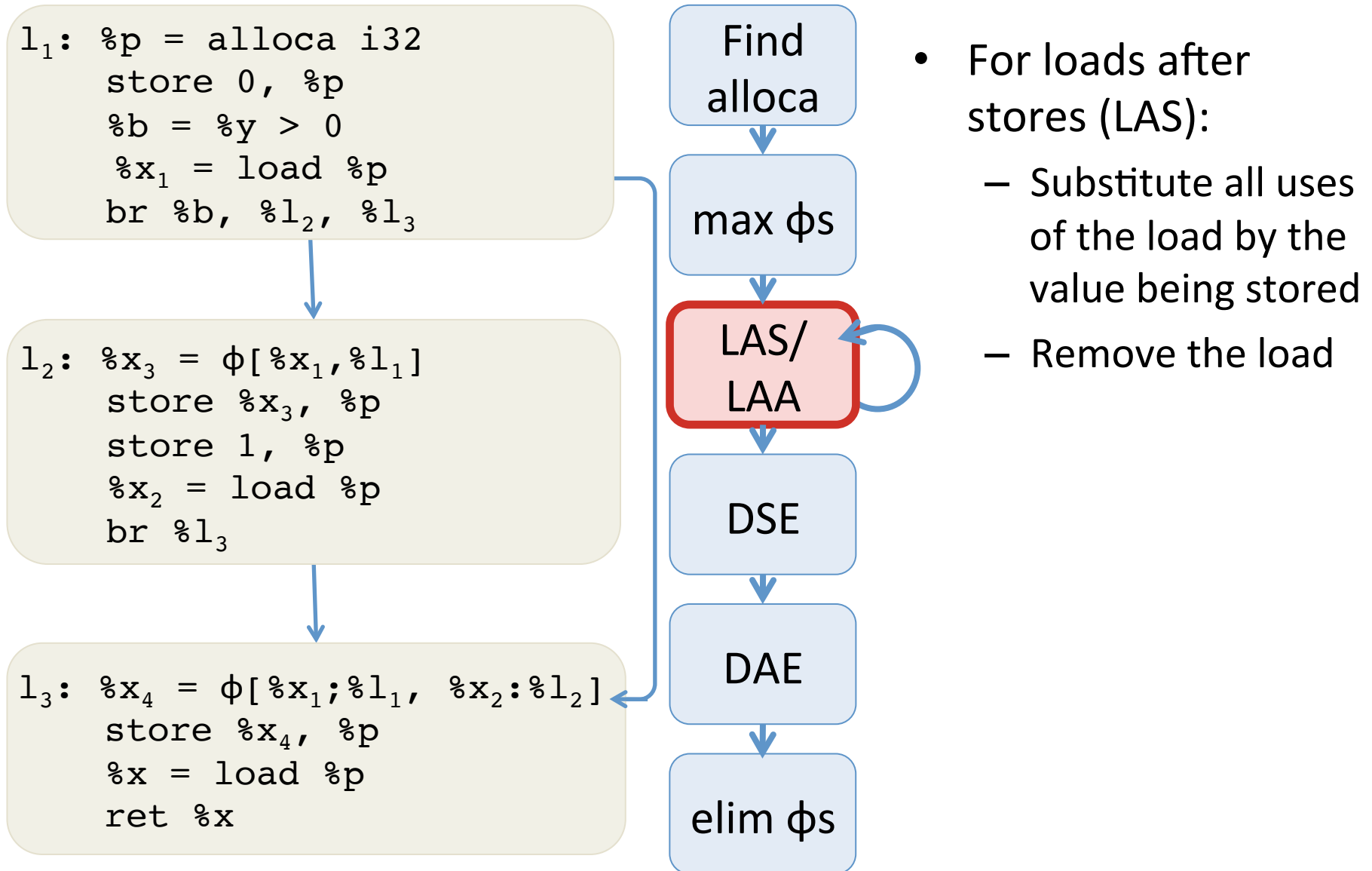
Example of vmem2reg Algorithm



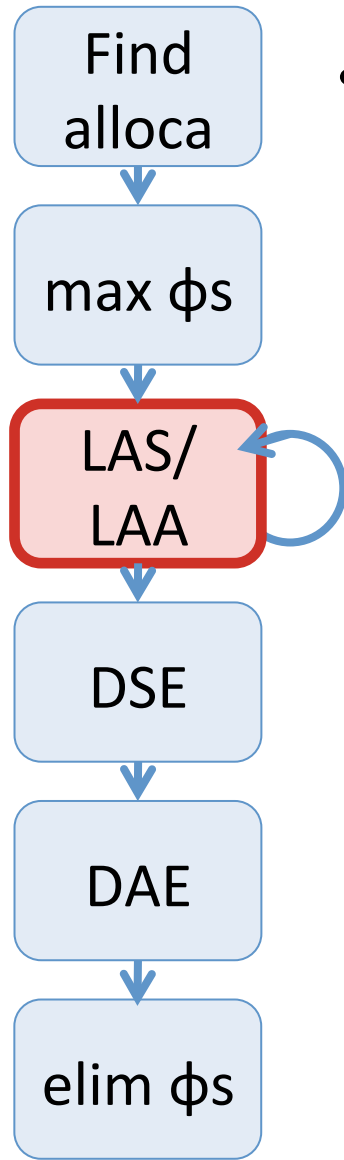
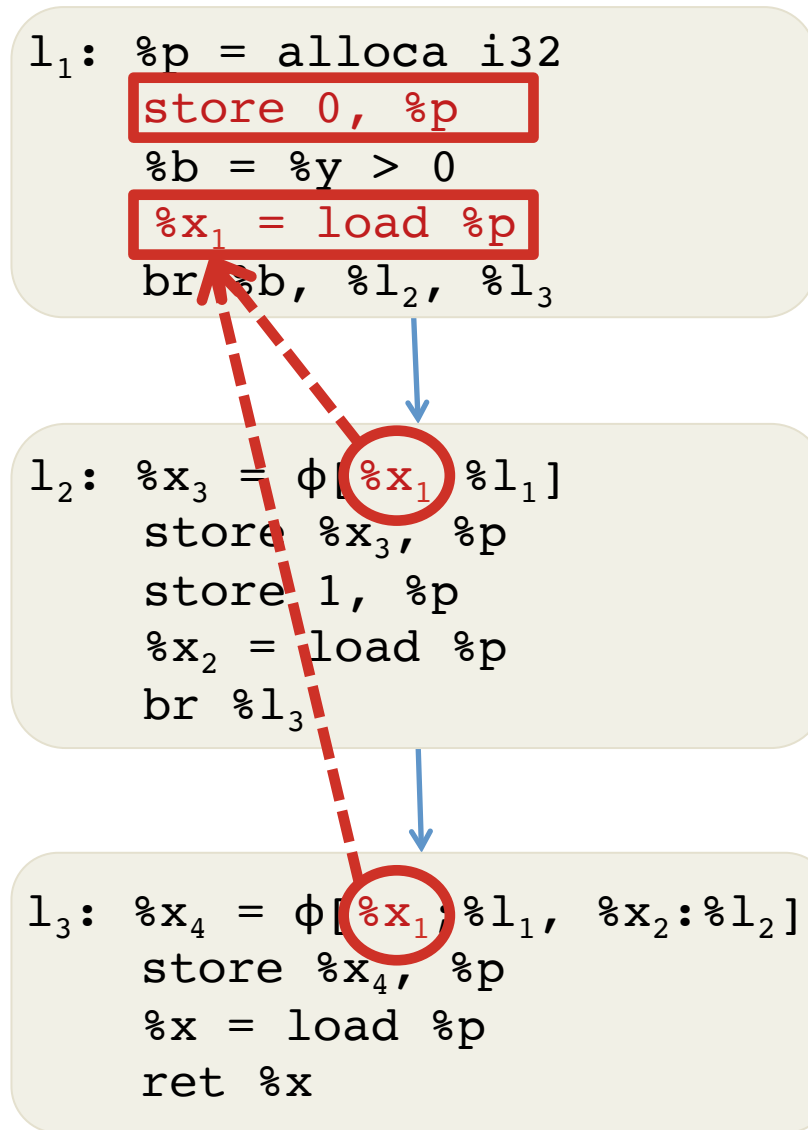
Example of vmem2reg Algorithm



Example of vmem2reg Algorithm

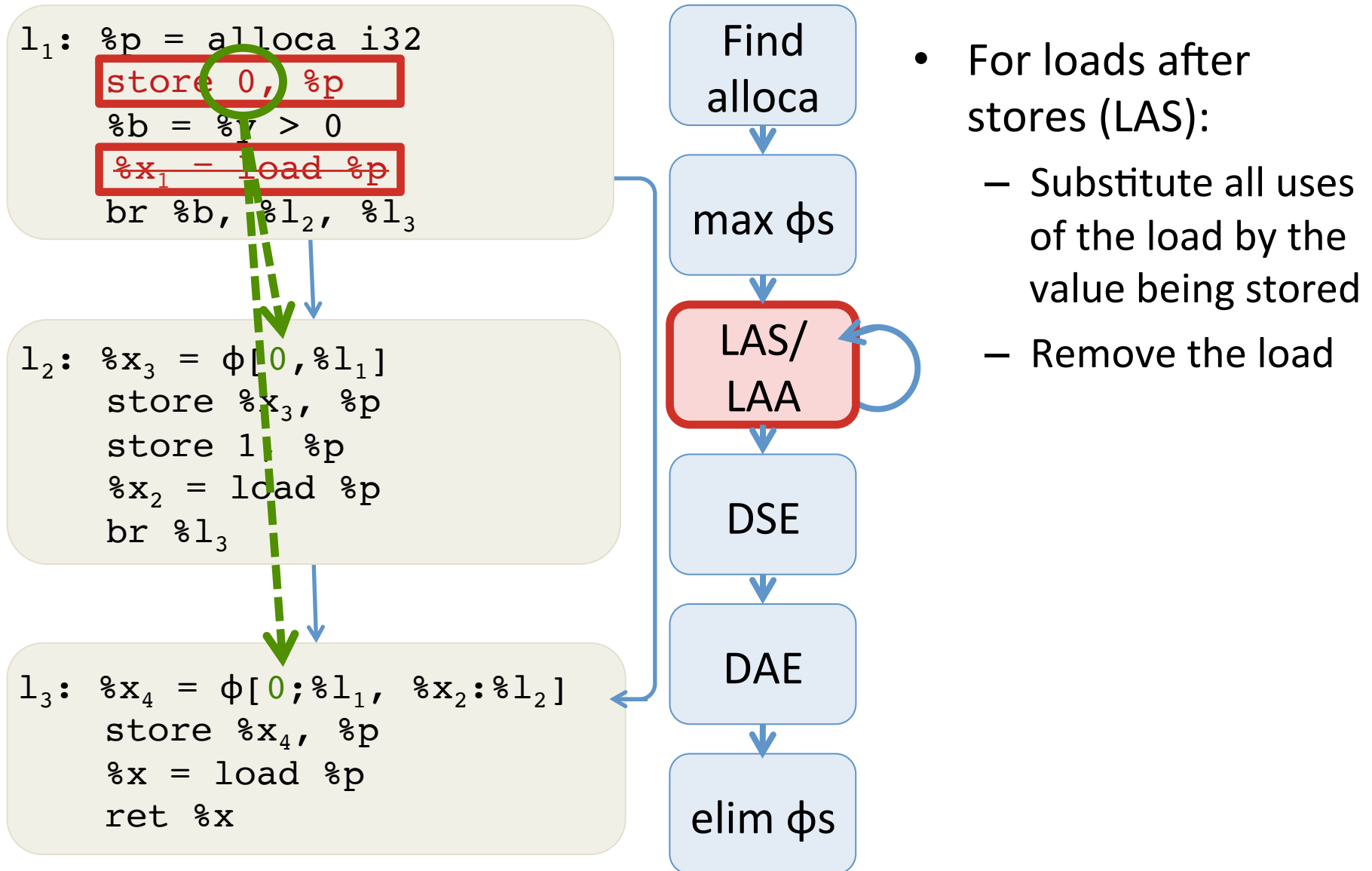


Example of vmem2reg Algorithm

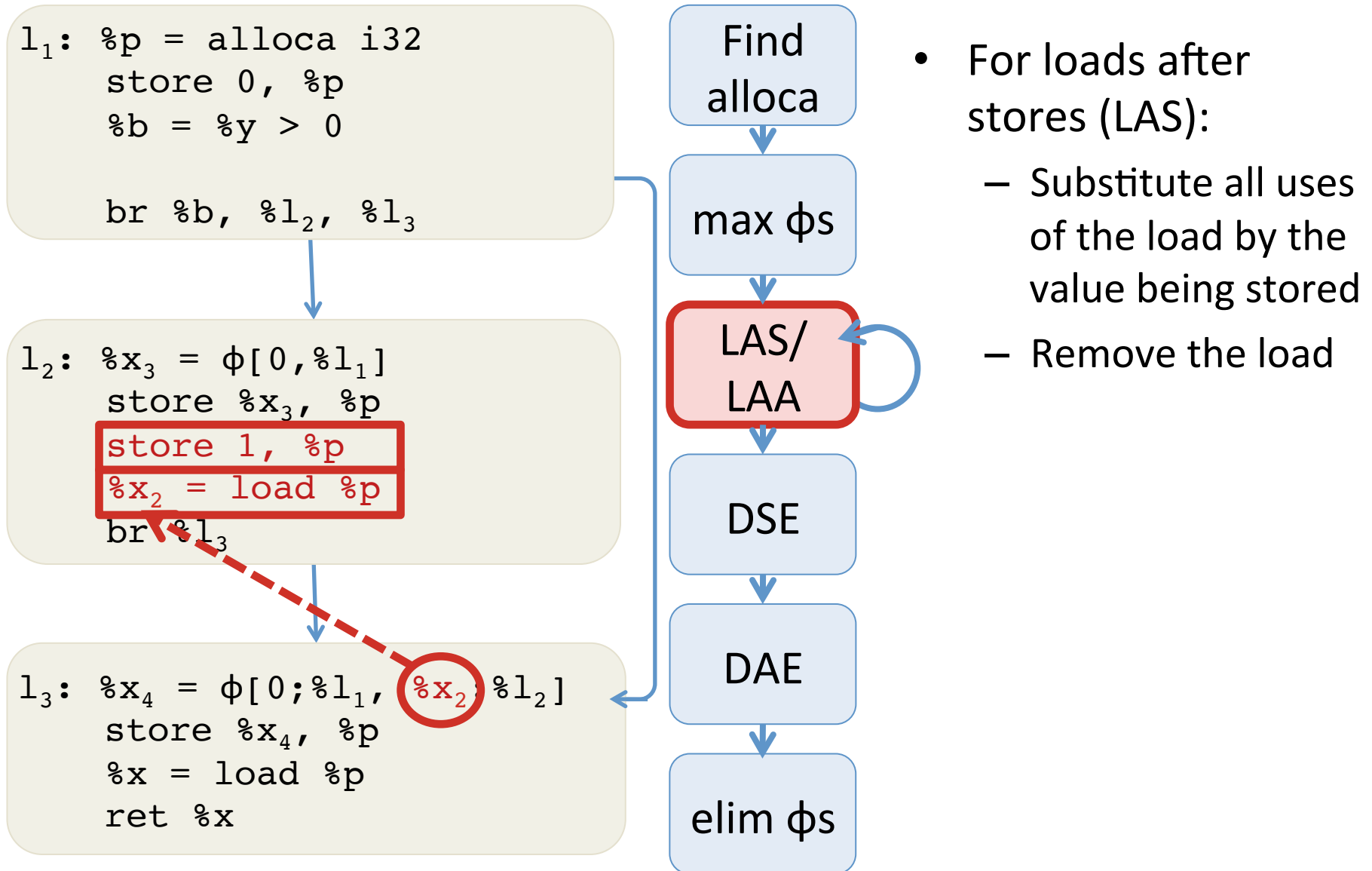


- For loads after stores (LAS):
 - Substitute all uses of the load by the value being stored
 - Remove the load

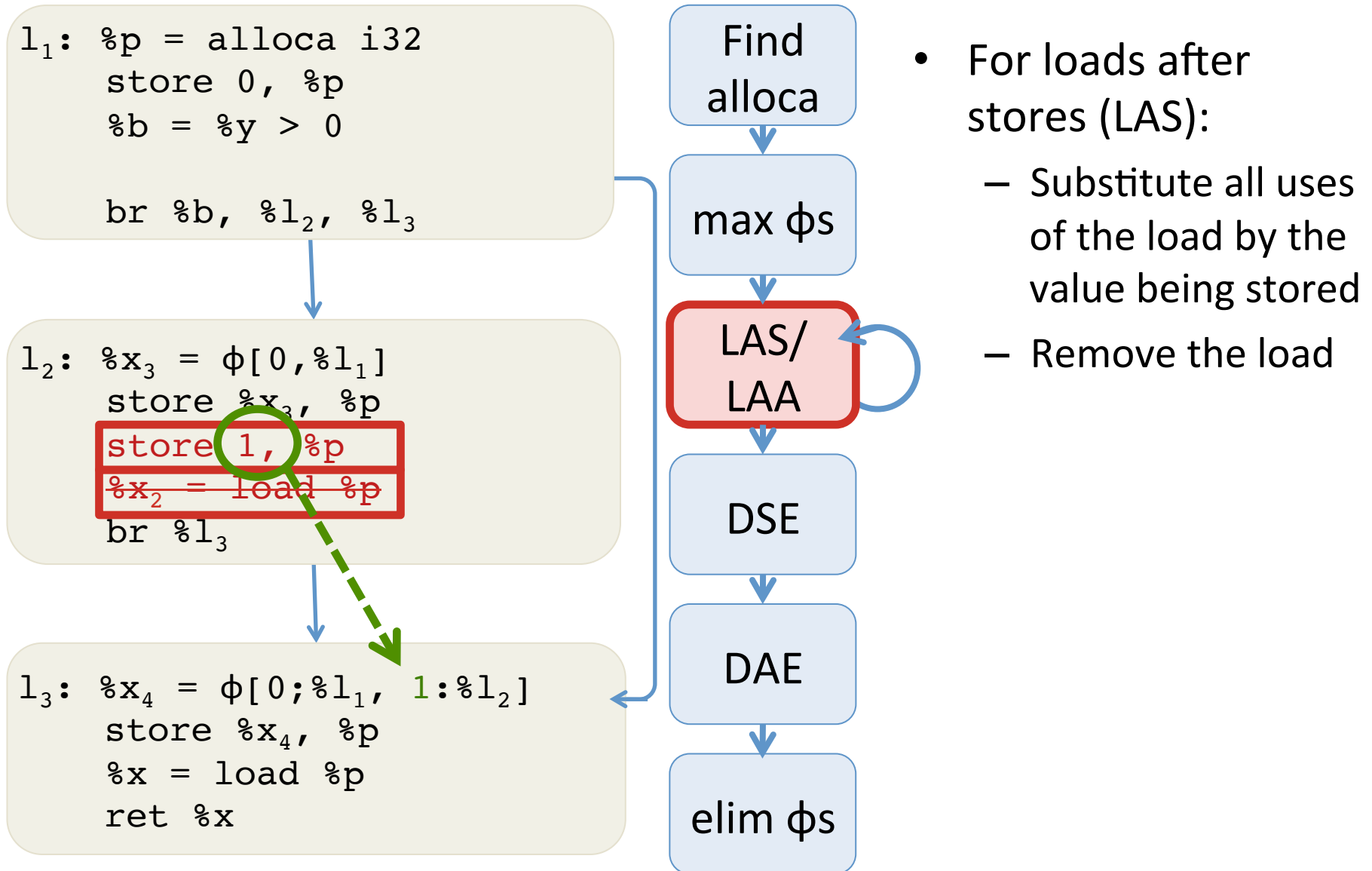
Example of vmem2reg Algorithm



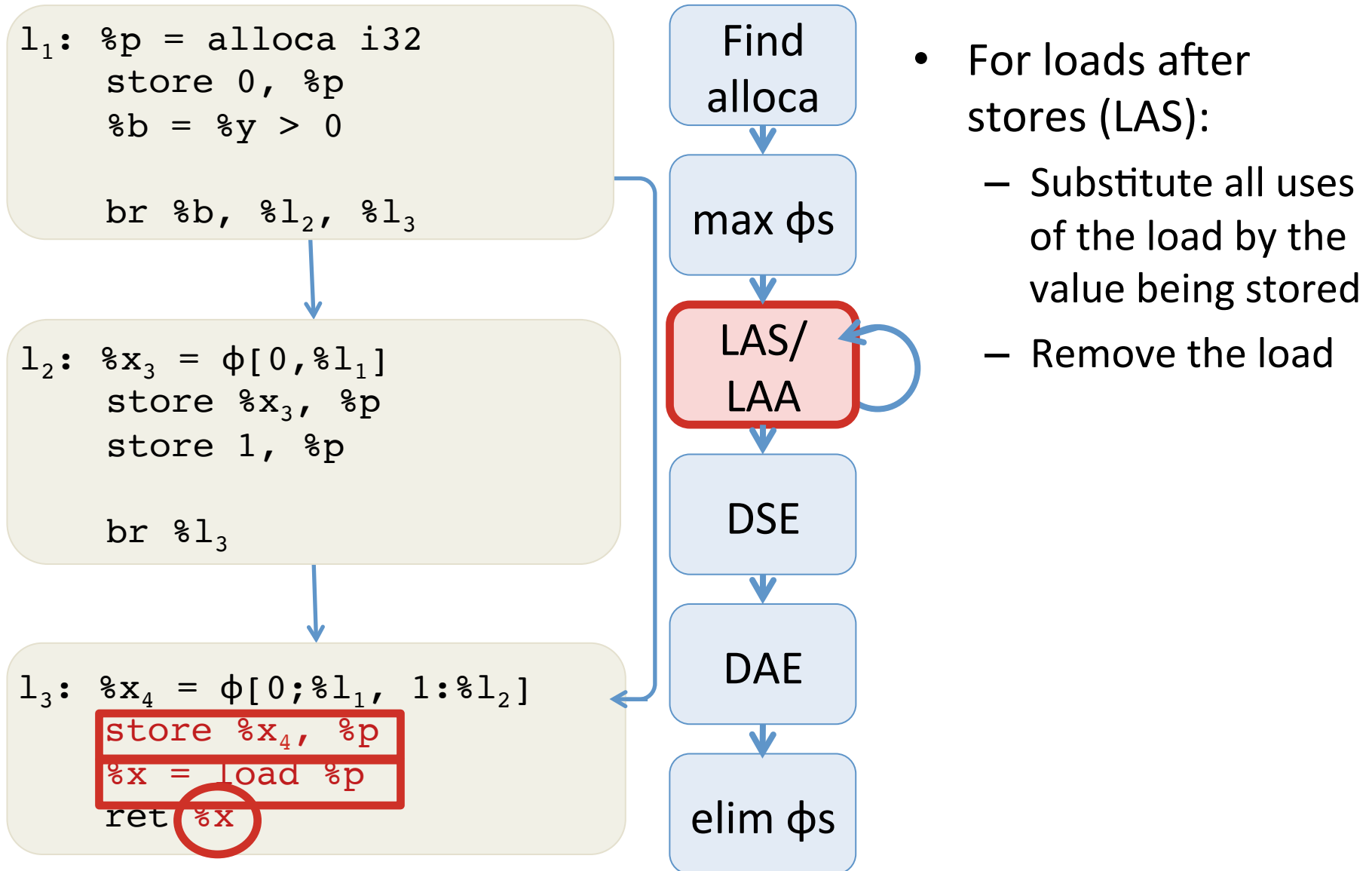
Example of vmem2reg Algorithm



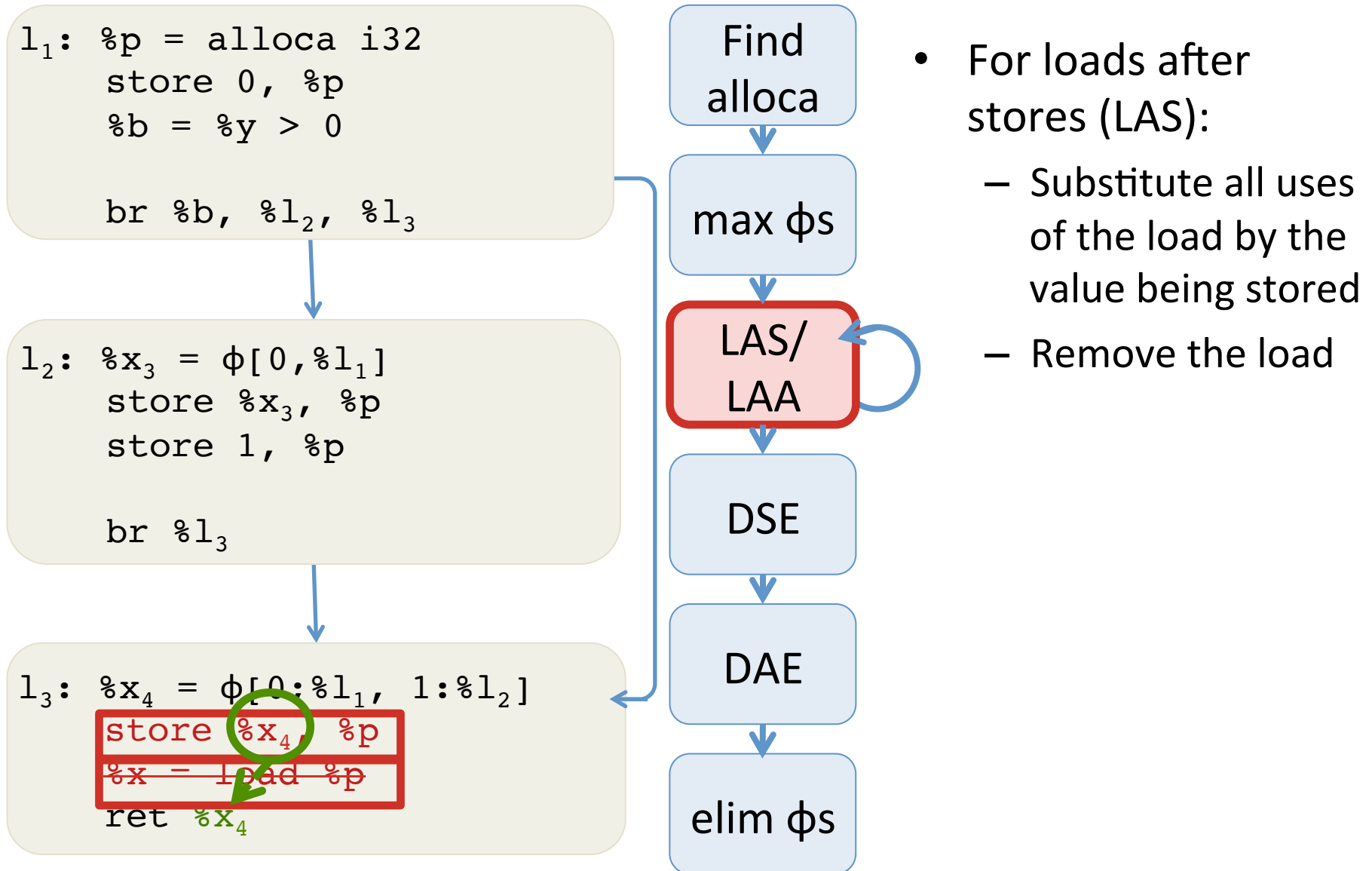
Example of vmem2reg Algorithm



Example of vmem2reg Algorithm

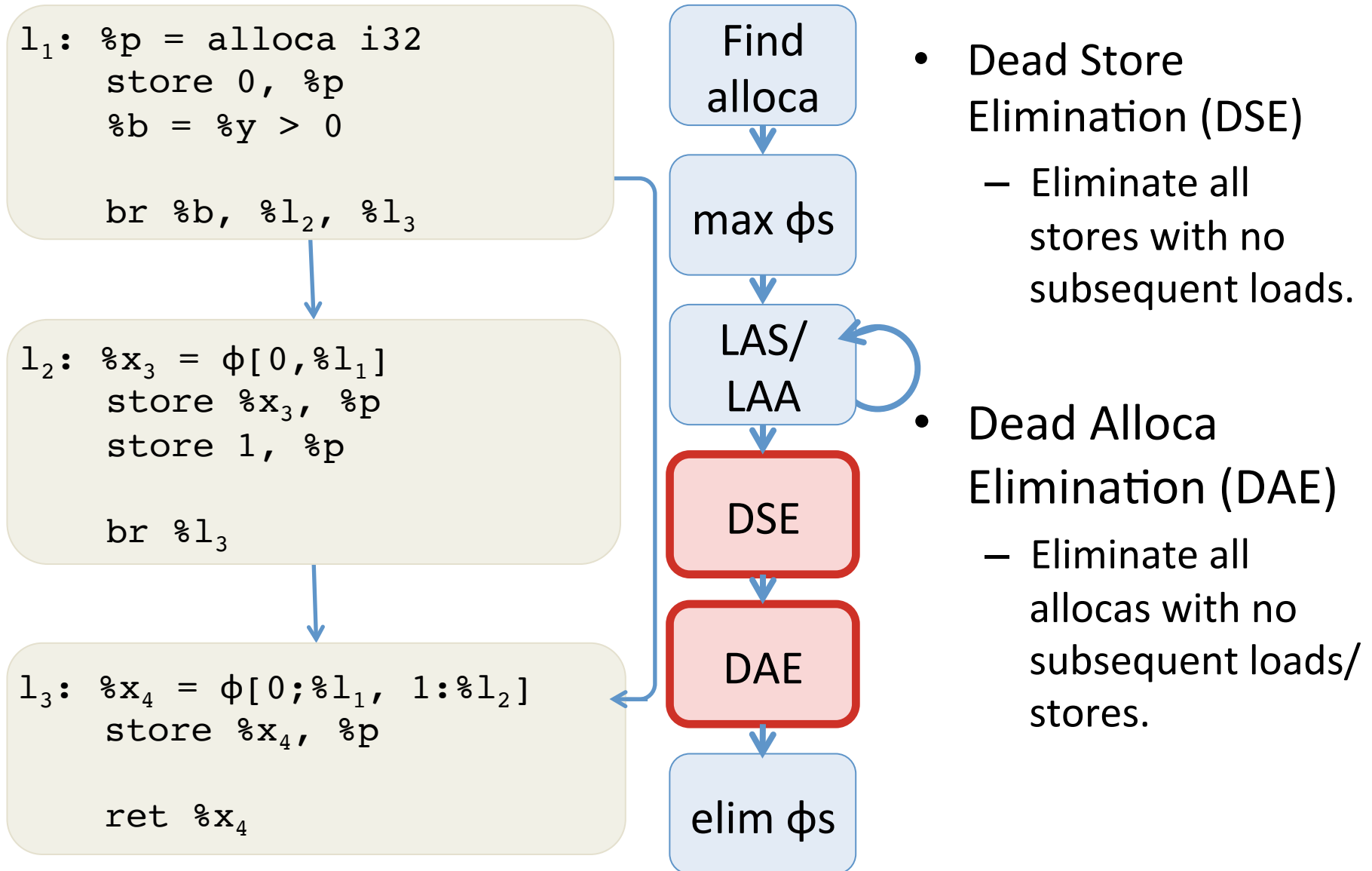


Example of vmem2reg Algorithm

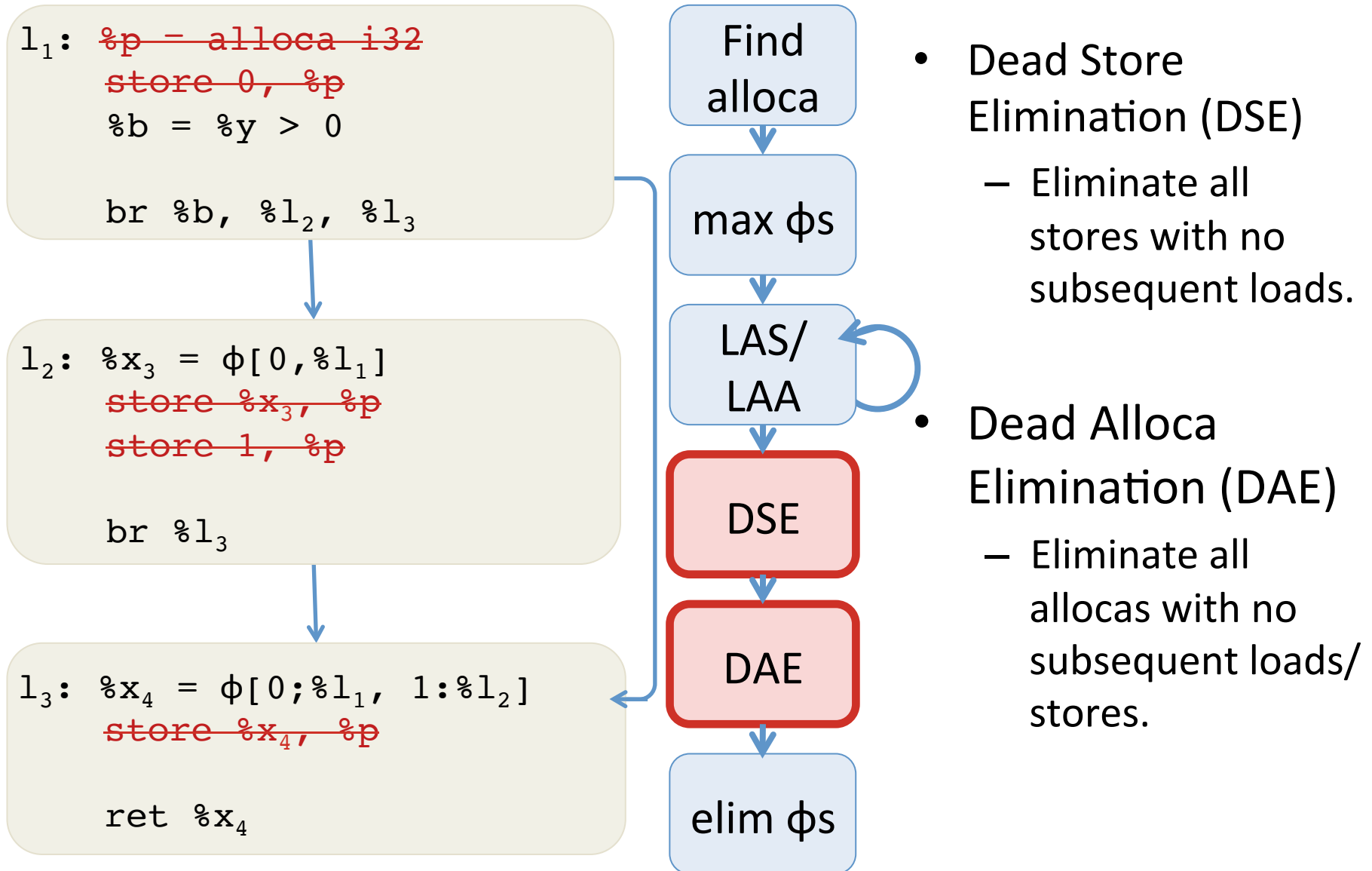


- For loads after stores (LAS):
 - Substitute all uses of the load by the value being stored
 - Remove the load

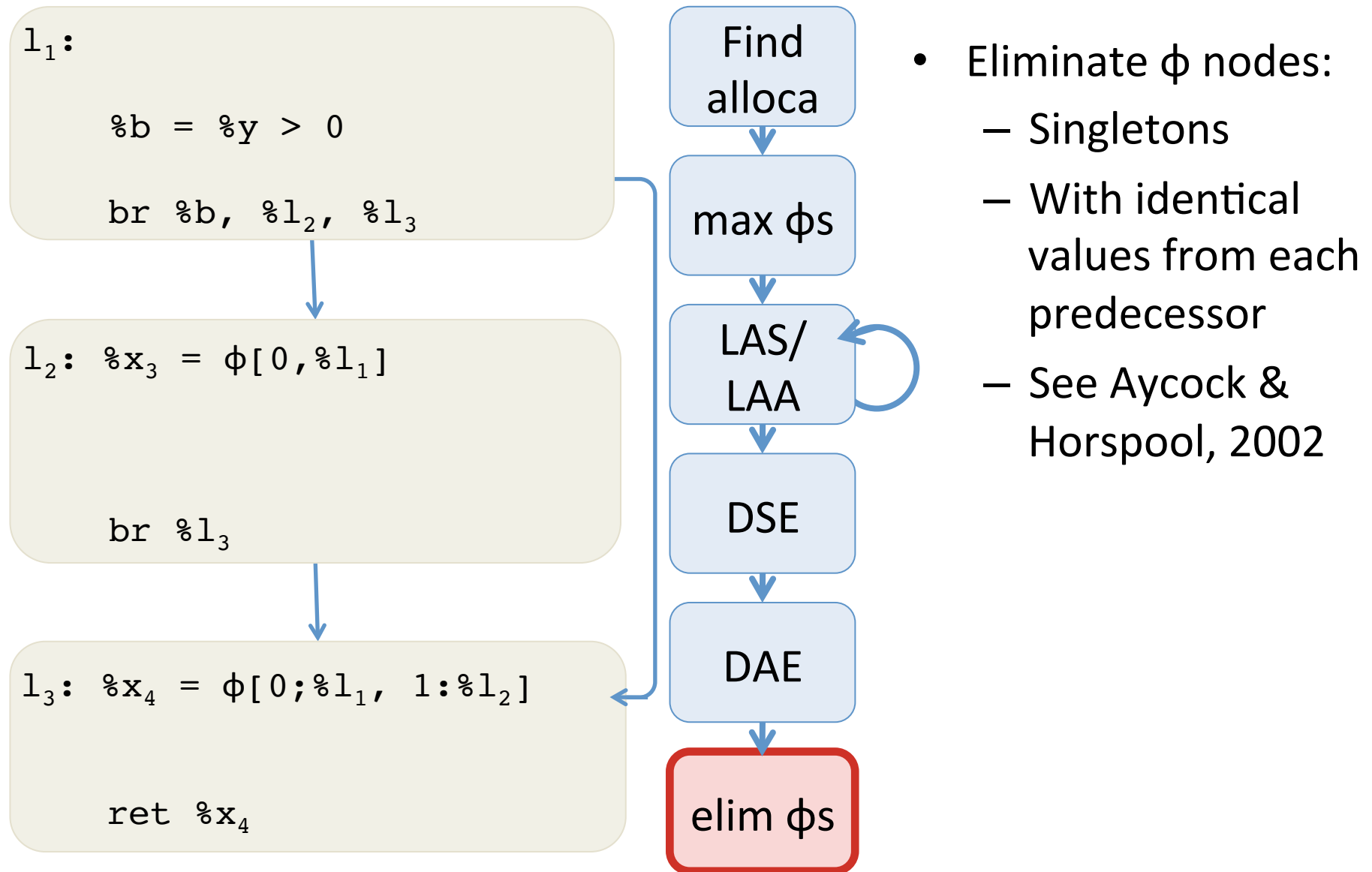
Example of vmem2reg Algorithm



Example of vmem2reg Algorithm



Example of vmem2reg Algorithm

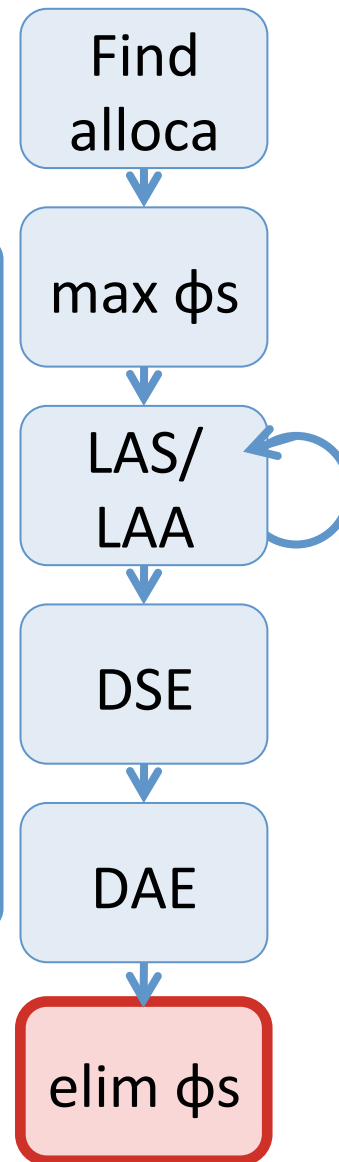


Example of vmem2reg Algorithm

```
l1:  
%b = %y > 0  
br %b, %l2, %l3
```

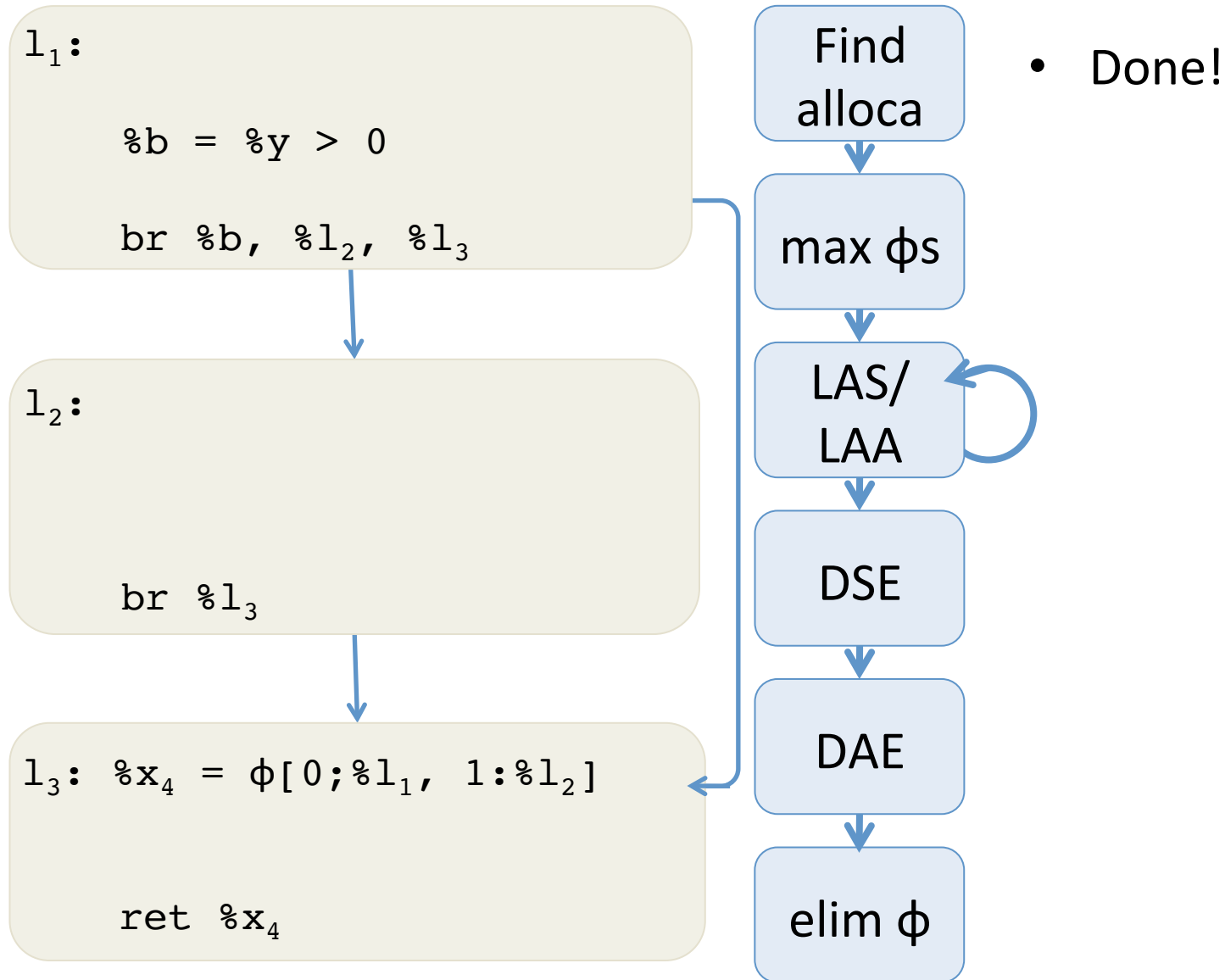
```
l2: %x3 = φ[0, %l1]  
br %l3
```

```
l3: %x4 = φ[0; %l1, 1: %l2]  
ret %x4
```

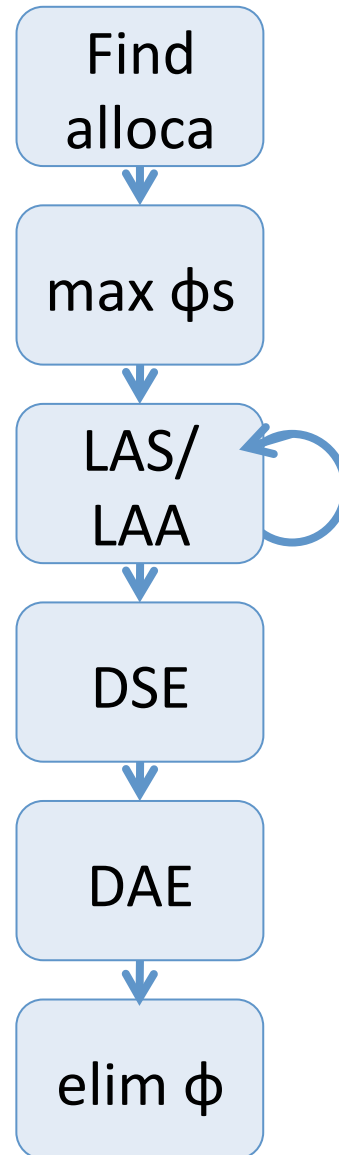


- Eliminate ϕ nodes:
 - Singletons
 - With identical values from each predecessor
 - See Aycock & Horspool, 2002

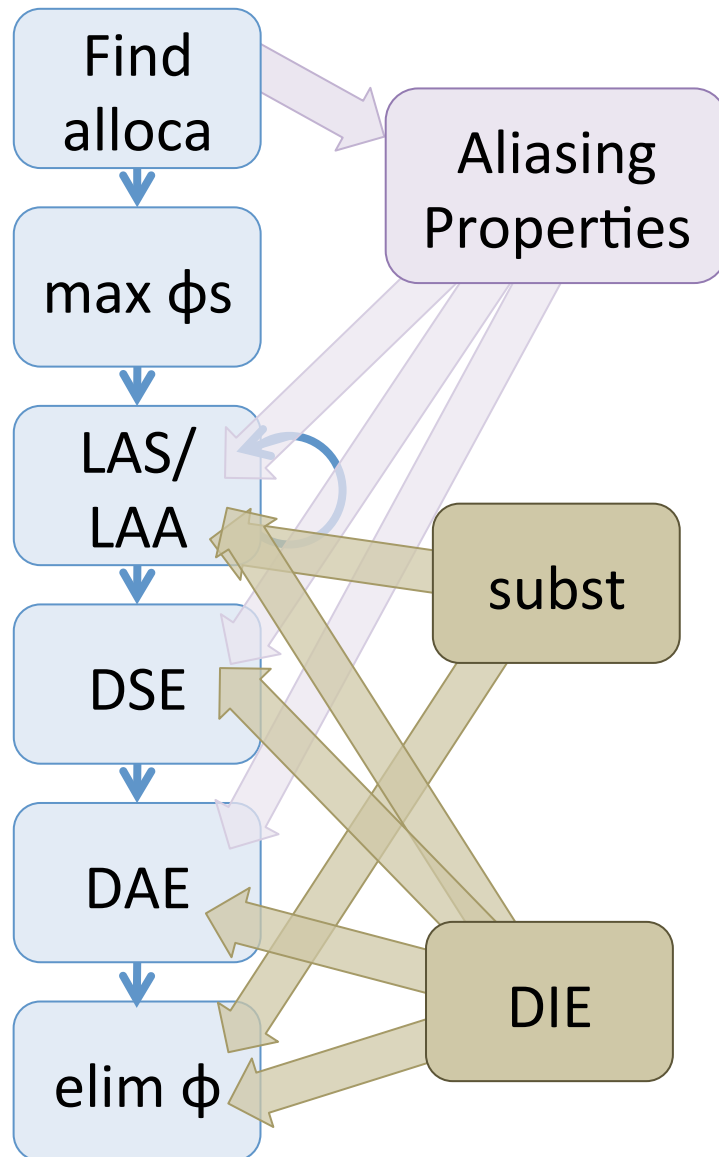
Example of vmem2reg Algorithm



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.) □

SCALING UP: LLVM

Other Parts of the LLVM IR

```
op      ::= %uid | constant | undef
bop     ::= add | sub | mul | shl | ...
cmpop  ::= eq | ne | slt | sle | ...
```

Operands
Operations
Comparison

```
insn ::=
| %uid = alloca ty
| %uid = load ty op1
| store ty op1, op2
| %uid = getelementptr ty op1 ...
| %uid = call rt fun(...args...)
| ...
```

Stack Allocation
Load
Store
Address Calculation
Function Calls

```
phi ::=
| φ[op1;lbl1]...[opn;lbln]
```

```
terminator ::=
| ret %ty op
| br op label %lbl1, label %lbl2
| br label %lbl
```

Structured Data in LLVM

- LLVM's IR uses types to describe the structure of data.

```
ty ::=
| i1 | i8 | i32 | ...           N-bit integers
| [<#elts> x t]                 arrays
| r (ty1, ty2, ... , tyn)    function types
| {ty1, ty2, ... , tyn}     structures
| ty*                           pointers
| %Tident                        named (identified) type

r ::=           Return Types
  ty           first-class type
  void        no return value
```

- <#elts> is an integer constant ≥ 0
- (Recursive) Structure types can be named at the top level:

```
%T1 = type {ty1, ty2, ... , tyn}
```

Distilling the LLVM

Documentation for the LLVM System January 2012 Archives by thread

- [LLVM Design](#)
- [LLVM Publications](#)
- [LLVM User Guides](#)
- [General LLVM Program](#)
- [LLVM Subsystem Docur](#)
- [LLVM Mailing Lists](#)

Written by [The LLVM Team](#)

- Messages sorted by: [\[subject \]](#) [\[author \]](#) [\[date \]](#)
- [More info on this list...](#)

Starting: Sun Jan 1 12:44:27 CST 2012

Ending: Thu Jan 19 18:21:55 CST 2012

Messages: 348

LLVM

- [LLVM Language Reference](#)
- [Introduction to the LLVM C](#)
- [The LLVM Compiler Frame](#)
exploring the system.
- [LLVM: A Compilation Fra](#)
overview.
- [LLVM: An Infrastructure](#)
- [GetElementPtr FAQ](#) - An
misunderstood instructor

- [The LLVM Getting Sta](#)

infrastructure. Everything from unpack...

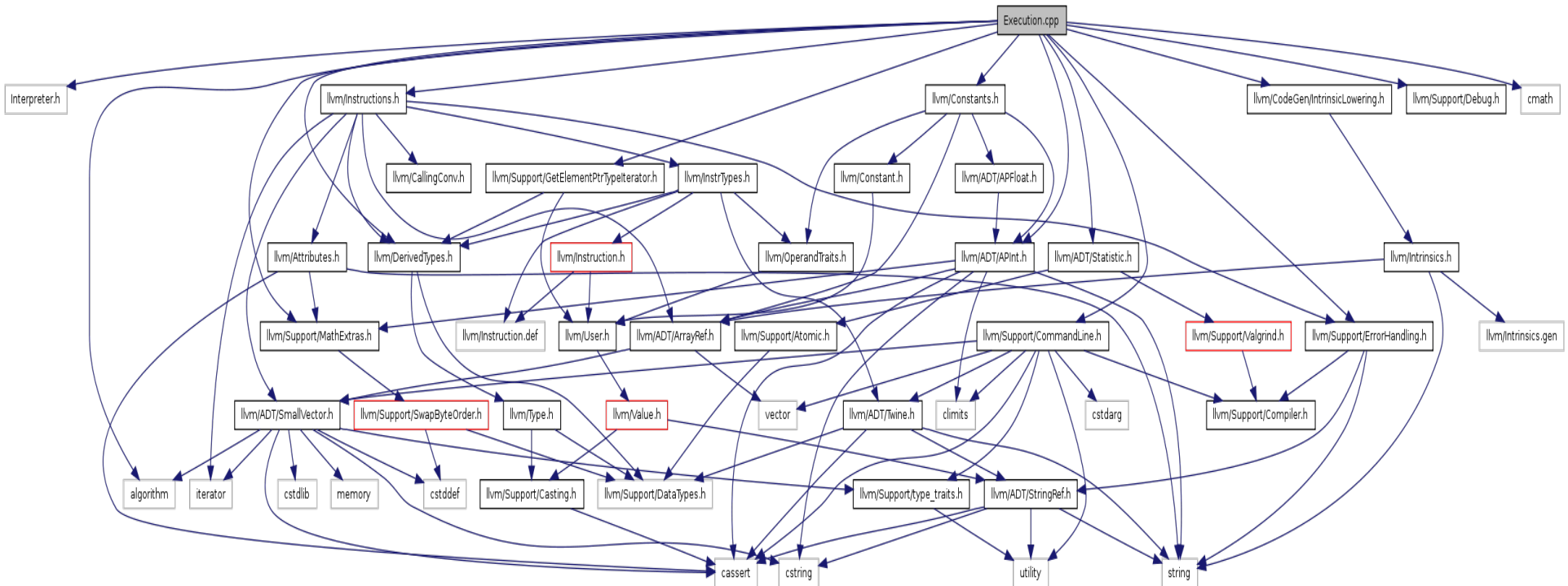
- [\[LLVMdev\] \[PATCH\] TLS support for Windows 32+64bit](#) Kai
- [\[LLVMdev\] \[PATCH\] TLS support for Windows 32+64bit](#) Eli Friedman
- [\[LLVMdev\] \[PATCH\] TLS support for Windows 32+64bit](#) Kai
- [\[LLVMdev\] \[PATCH\] TLS support for Windows 32+64bit](#) Kai
- [\[LLVMdev\] tbaa](#) Jianzhou Zhao
- [\[LLVMdev\] Checking validity of metadata in an .ll file](#) Seb
- [\[LLVMdev\] Checking validity of metadata in an .ll file](#) Devang Patel
- [\[LLVMdev\] Using llvm command line functions from within a plugin?](#) Talin
- [\[LLVMdev\] Using llvm command line functions from within a plugin?](#) Duncan Sands
- [\[LLVMdev\] Using llvm command line functions from within a plugin?](#) Talin
- [\[LLVMdev\] Comparison of Alias Analysis in LLVM](#) Jianzhou Zhao
- [\[LLVMdev\] Comparison of Alias Analysis in LLVM](#) Chris Lattner
- [\[LLVMdev\] Comparison of Alias Analysis in LLVM](#) Jianzhou Zhao
- [\[LLVMdev\] Comparison of Alias Analysis in LLVM](#) Chris Lattner
- [\[LLVMdev\] Comparison of Alias Analysis in LLVM](#) Jianzhou Zhao

Distilling the LLVM

Documentation for the LLVM System
January 2012 Archives by thread

- [LLVM Design](#)
- [LLVM Publications](#)
- [LLVM User Guides](#)
- [General LLVM Program](#)

- Messages sorted by: [\[subject \]](#) [\[author \]](#) [\[date \]](#)
- [More info on this list...](#)



infrastructure. Everything... [Comparison of Alias Analysis in LLVM](#) Chris Lattner
Jianzhou Zhao

Example LLVM Types

- An array of 341 integers: `[341 x i32]`
- A 2D array of integers: `[3 x [4 x i32]]`
- C-style linked lists:
`%Node = type { i32, %Node* }`
- Structs:
`%Rect = { %Point, %Point,
 %Point, %Point }
%Point = { i32, i32 }`

GetElementPtr

- LLVM provides the `getelementptr` instruction to compute pointer values
 - Given a pointer and a “path” through the structured data pointed to by that pointer, `getelementptr` computes an address
 - This is the abstract analog of the X86 LEA (load effective address). It does not access memory.
 - It is a “type indexed” operation, since the size computations involved depend on the type

```
insn ::= ...  
      | %uid = getelementptr t*, %val, t1 idx1, t2 idx2 ,...
```


Example

```

struct RT {
    int A;
    int B[10][20];
    int C;
}
struct ST {
    struct RT X;
    int Y;
    struct RT Z;
}
int *foo(struct ST *s) {
    return &s[1].Z.B[5][13];
}
    
```

1. %s is a pointer to an (array of) ST structs, suppose the pointer value is ADDR

2. Compute the index of the 1st element by adding sizeof(struct ST).

3. Compute the index of the Z field by adding sizeof(struct RT) + sizeof(int) to skip past X and Y.

4. Compute the index of the B field by adding sizeof(int) to skip past A.

5. Index into the 2d array.

```

%RT = type { i32, [10 x [20 x i32]], i32 }
%ST = type { %RT, i32, %RT }
define i32* @foo(%ST* %s) {
entry:
    %arrayidx = getelementptr %ST* %s, i32 1, i32 2, i32 1, i32 5, i32 13
    ret i32* %arrayidx
}
    
```

Final answer: ADDR + sizeof(struct ST) + sizeof(struct RT) + sizeof(int) + sizeof(int) + 5*20*sizeof(int) + 13*sizeof(int)

LLVM's memory model

```
%ST = type {i10,[10 x i8*]}
```

High-level
Representation

i10
i8*
i8*
i8*
i8*
i8*
i8*
i8*
i8*
i8*
i8*

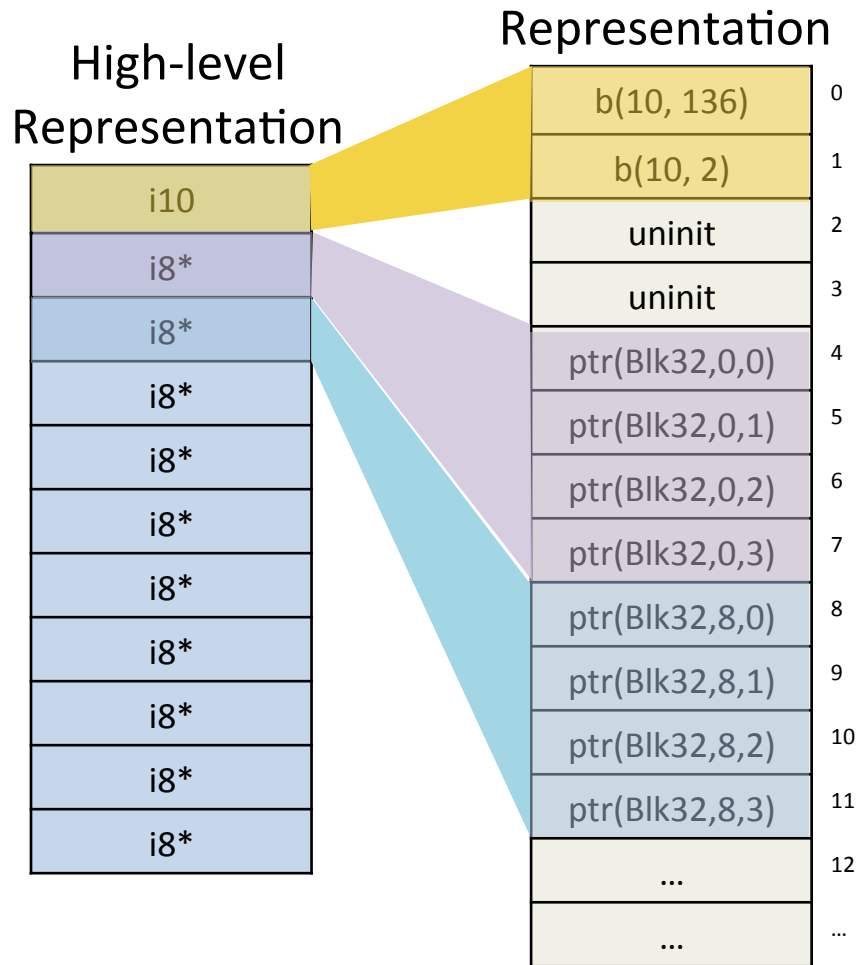
- Manipulate structured types.

```
%val = load %ST* %ptr  
...  
store %ST* %ptr, %new
```

LLVM's memory model

```
%ST = type {i10,[10 x i8*]}
```

Low-level

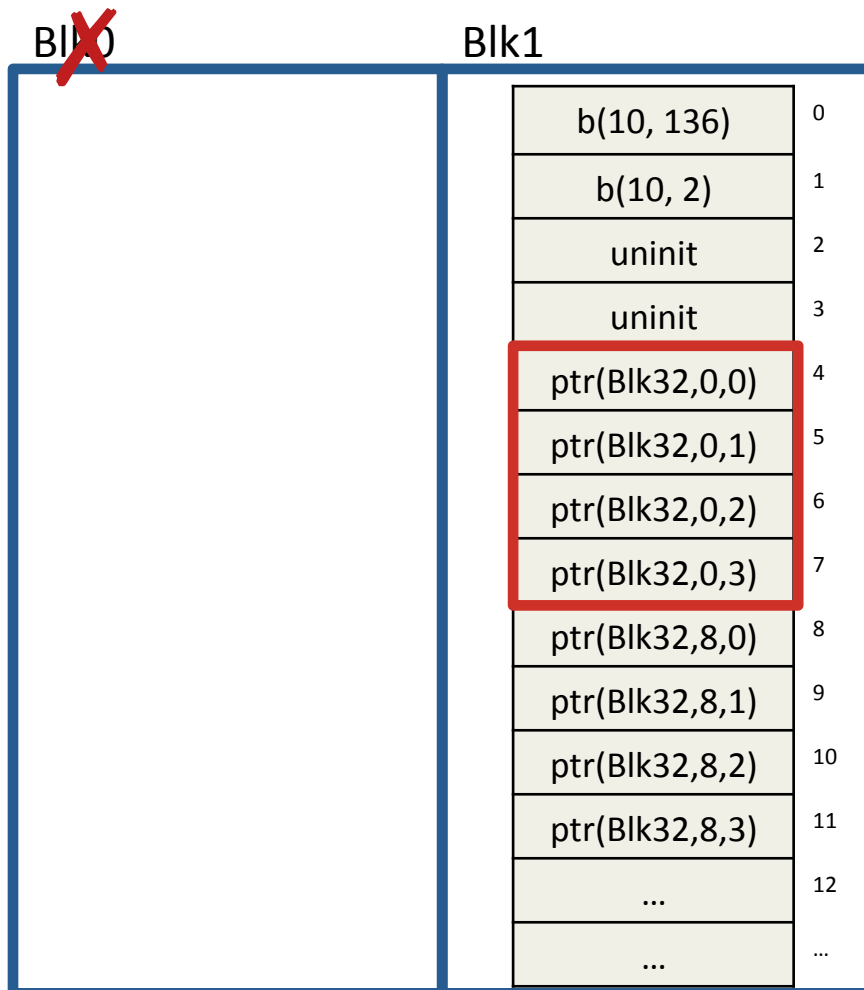


- Manipulate structured types.

```
%val = load %ST* %ptr  
...  
store %ST* %ptr, %new
```

- Semantics is given in terms of byte-oriented low-level memory.
 - padding & alignment
 - physical subtyping

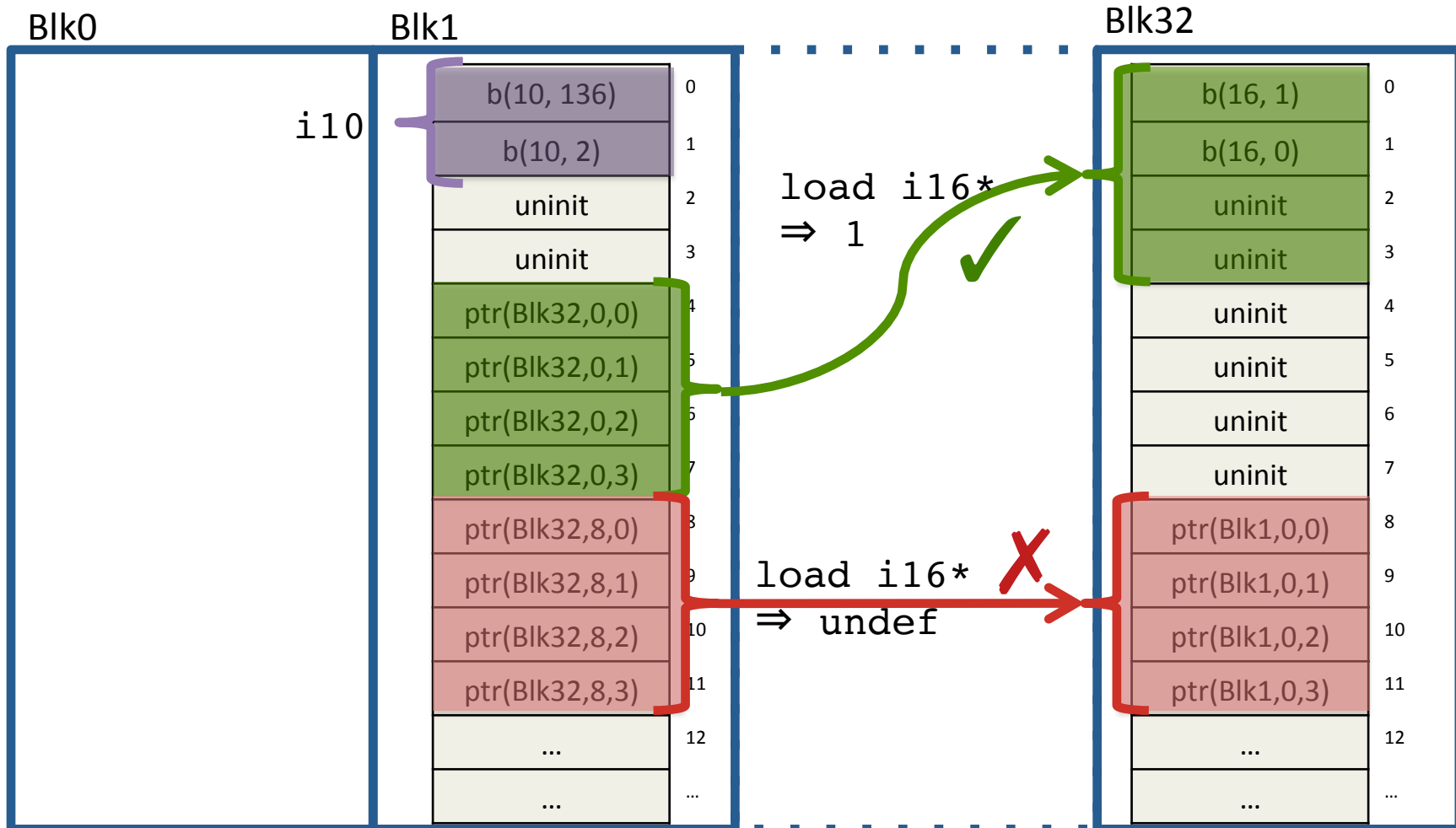
Adapting CompCert's Memory Model



- Data lives in blocks
- Represent pointers abstractly
 - block + offset
- Deallocate by invalidating blocks
- Allocate by creating new blocks
 - infinite memory available

Dynamic Physical Subtyping

[Nita, et al. *POPL '08*]



Sources of Undefined Behavior

Target-dependent Results

- Uninitialized variables:

```
%v = add i32 %x, undef
```

- Uninitialized memory:

```
%ptr = alloca i32  
%v = load (i32*) %ptr
```

- Ill-typed memory usage

Nondeterminism

Fatal Errors

- Out-of-bounds accesses
- Access dangling pointers
- Free invalid pointers
- Invalid indirect calls

Stuck States

Sources of Undefined Behavior

Target-dependent Results

- Uninitialized variables:

```
%v = add i32 %x, undef
```

- Uninitialized memory:

```
%ptr = alloca i32  
%v = load (i32*) %ptr
```

- Ill-typed memory usage

Nondeterminism

Defined by a predicate on the program configuration.

```
Stuck(f,  $\sigma$ ) = BadFree(f,  $\sigma$ )  
                   $\vee$  BadLoad(f,  $\sigma$ )  
                   $\vee$  BadStore(f,  $\sigma$ )  
                   $\vee$  ...  
                   $\vee$  ...
```

Stuck States

undef

- What is the value of %y after running the following?

```
%x = or i8 undef, 1
%y = xor i8 %x %x
```

- One plausible answer: 0
- Not LLVM's semantics!
(LLVM is more liberal to permit more aggressive optimizations)

undef

- Partially defined values are interpreted *nondeterministically* as sets of possible values:

```
%x = or i8 undef, 1
%y = xor i8 %x %x
```

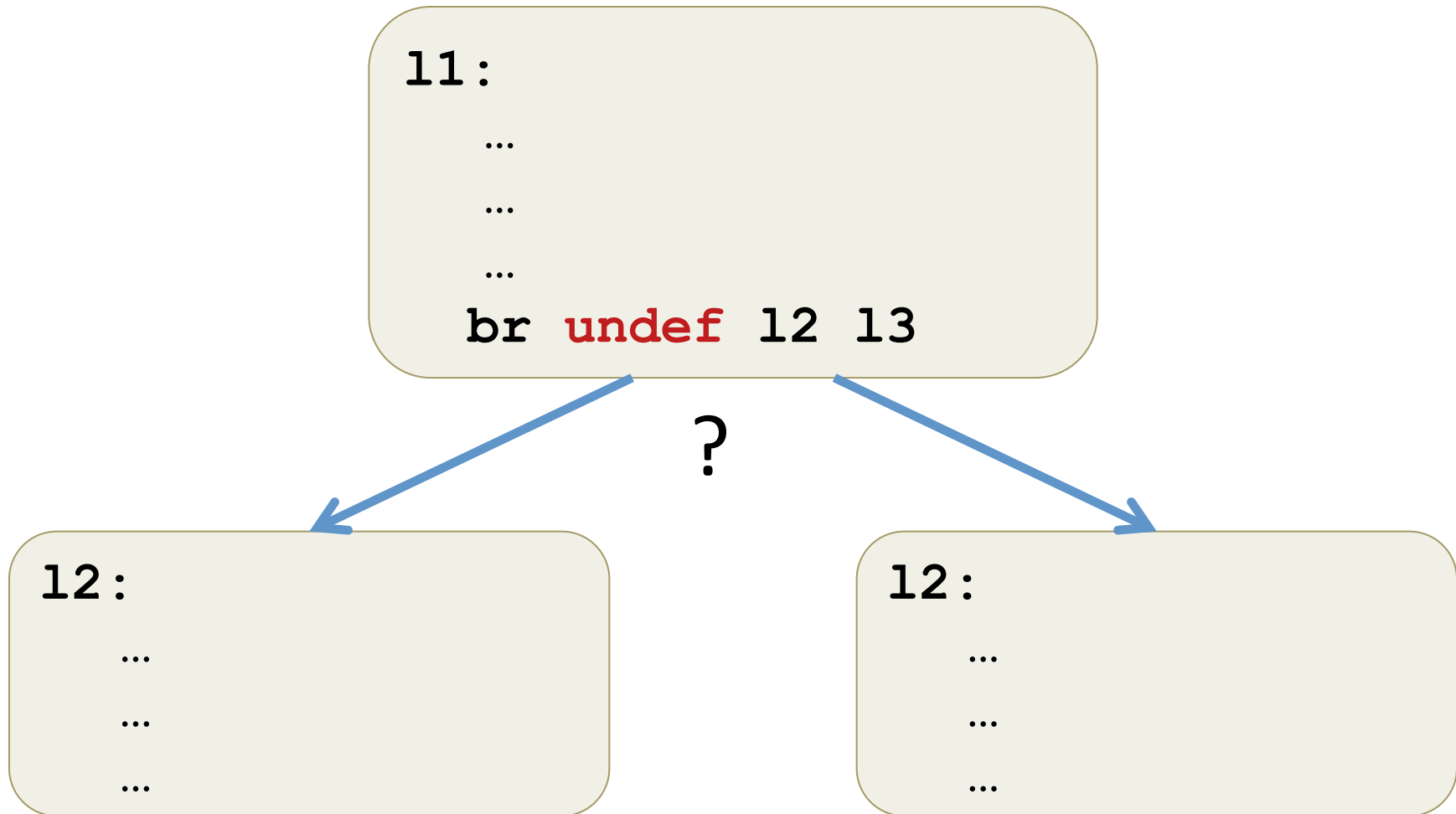
$[[i8\ undef]] = \{0, \dots, 255\}$

$[[i8\ 1]] = \{1\}$

$[[\%x]] = \{a\ or\ b\ |\ a \in [[i8\ undef]],\ b \in [[1]]\}$
 $= \{1, 3, 5, \dots, 255\}$

$[[\%y]] = \{a\ xor\ b\ |\ a \in [[\%x]],\ b \in [[\%x]]\}$
 $= \{0, 2, 4, \dots, 254\}$

Nondeterministic Branches



LLVM_{ND} Operational Semantics

- Define a transition relation:

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

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

Operational Semantics

	Small Step	Big Step
Nondeterministic	LLVM_{ND}	
Deterministic		

Deterministic Refinement

	Small Step	Big Step
Nondeterministic	LLVM_{ND}	
	\Downarrow	
Deterministic	LLVM_D	

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

Big-step Deterministic Refinements

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

\cup

Bisimulation up to “observable events”:

- external function calls

Big-step Deterministic Refinements

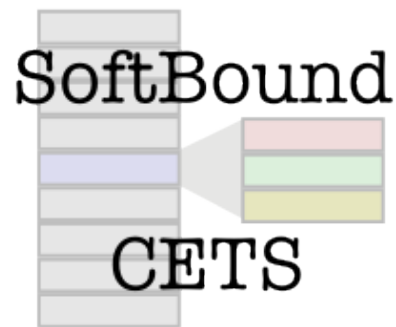
	Small Step	Big Step
Nondeterministic	<div style="border: 1px solid black; padding: 5px; display: inline-block;"> $LLVM_{ND}$ </div>	
Deterministic	<div style="display: flex; align-items: center; gap: 10px;"> <div style="border: 1px solid black; padding: 5px; display: inline-block;"> $LLVM_{Interp}$ 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

Simulation up to “observable events”:

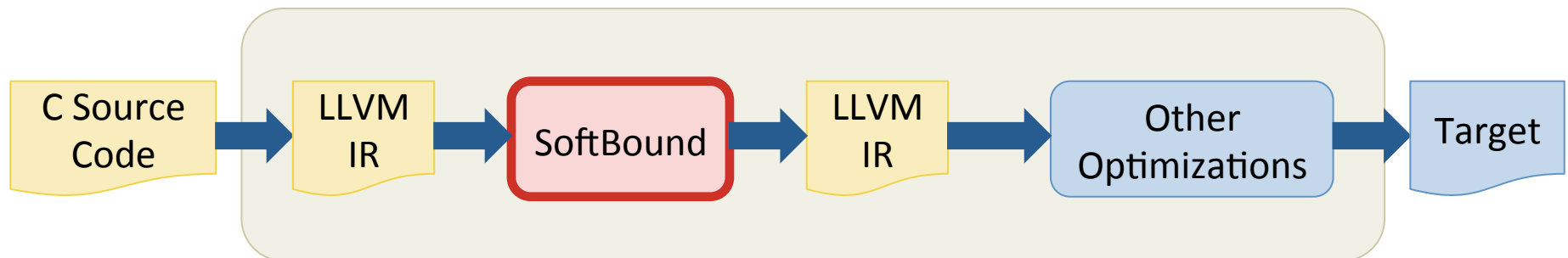
- useful for encapsulating behavior of function calls
- large step evaluation of basic blocks

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

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

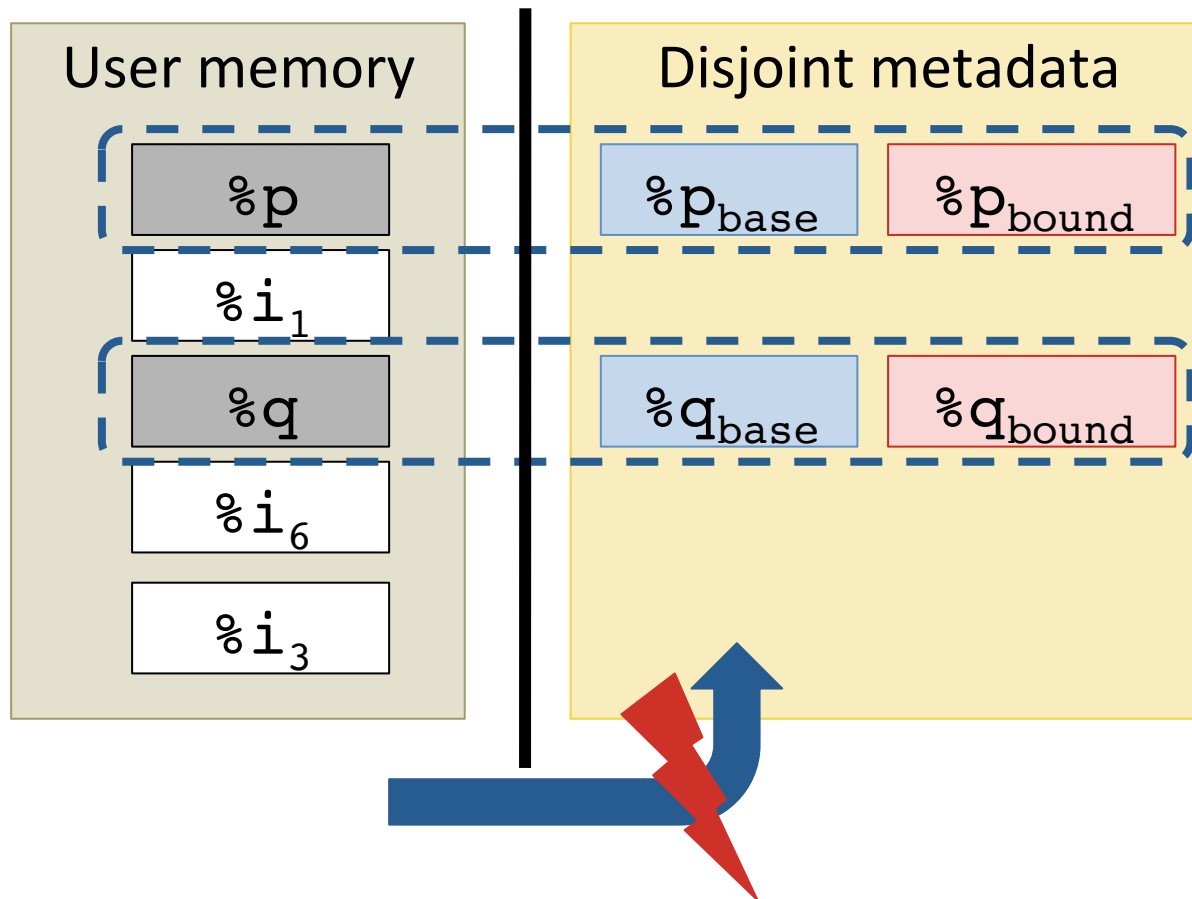
LLVM IR

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:

$$f \vdash \sigma \xrightarrow{\text{SB}} \sigma'$$

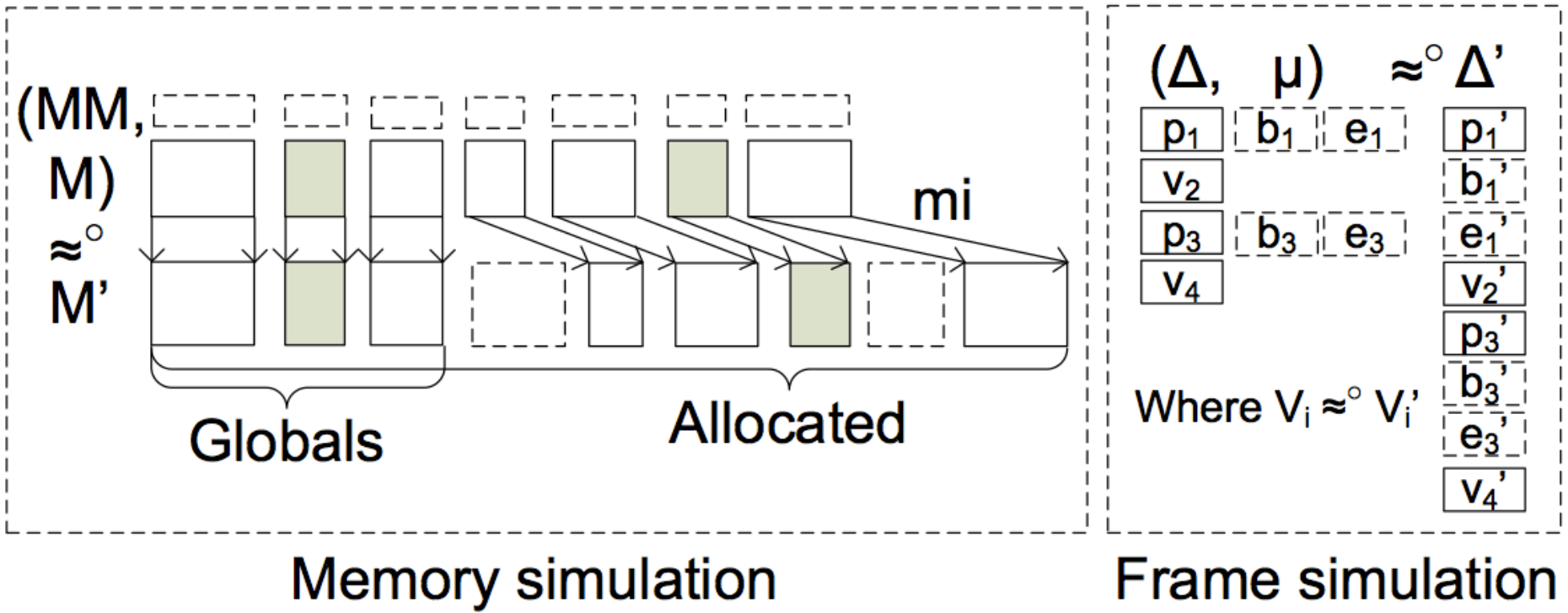
- Builds in safety invariants “by construction”

$$f \vdash \sigma \xrightarrow{\text{SB}*} \sigma' \Rightarrow \neg \text{MemoryViolation}(f, \sigma')$$

4. Show that the instrumented code simulates the “correct” code:

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

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)

VELLM \Rightarrow VELLM II

The Bad

- Large, monolithic code base
- Clunky proofs

Development	LOC (Defns. + Proofs)
syntax + semantics	~45K
mem2reg + optimizations	~60K
SoftBound	~15K

⇒ hard for others to adopt/adapt

The Bad

- Representation & semantics very syntactic
 - cfg \approx lists of blocks
 - block \approx lists of instructions
 - operational semantics uses this syntax
 - \Rightarrow lots of boiler plate everywhere
- Tightly coupled to the memory model, design non-modular
- LLVM transformations not designed for verification
 - translating "informal" to "formal" proofs is difficult
- Limited use of proof automation

\Rightarrow Hard to deal with change:

LLVM 2.6 \Rightarrow LLVM 3.0 \Rightarrow LLVM 3.6 \Rightarrow ...

Coq 8.2 \Rightarrow Coq 8.3 \Rightarrow Coq 8.4 \Rightarrow Coq 8.5 \Rightarrow ...

CompCert's memory model evolution

Ongoing Work

- "Legacy Vellvm"
 - defunct
- Vellvm II
 - <https://github.com/vellvm/vellvm>
- Modernize & Refactor the development
 - more streamlined
 - experiment with LTS operational semantics

Partial Bibliography

- CompCert [Leroy et al.]
- CompCertSSA [Barthe, Demange et al. ESOP 2012]
 - Translation validate the SSA construction
- Verified Software Toolchain [Appel et. al]
- Verifiable SSA Representation [Menon et al. POPL 2006]
 - Identify the well-formedness safety predicate for SSA
- Specification of SSA
 - Temporal checking & model checking for proving SSA transforms [Mansky et al, ITP 2010]
 - Matrix representation of ϕ nodes [Yakobowski, INRIA]
 - Type system equivalent to SSA [Matsuno et al]
- LLVM Semantics
 - Taming Undefined Behavior in LLVM [Lee et al., PLDI17]

Conclusions

- Proof techniques for verifying SSA transformations
 - Generalize the SSA scoping predicate
 - Preservation/progress + simulations.
 - Simulation proofs
- Verified:
 - Softbound & vmem2reg
 - Similar performance to native implementations
- See the papers/coq sources for details!
- Future:
 - Clean up + make more accessible
 - Alias analysis? Concurrency?
 - Applications to more LLVM-SSA optimizations



<http://www.cis.upenn.edu/~stevez/vellvm/>