

Compiler Optimization with Retrofitting Transformations: Is there a Semantic Mismatch?

Santosh Nagarakatte

Joint work with my PhD student Jay Lim and collaborator Vinod Ganapathy
Department of Computer Science,
Rutgers University



RUTGERS

Retrofitting Transformation



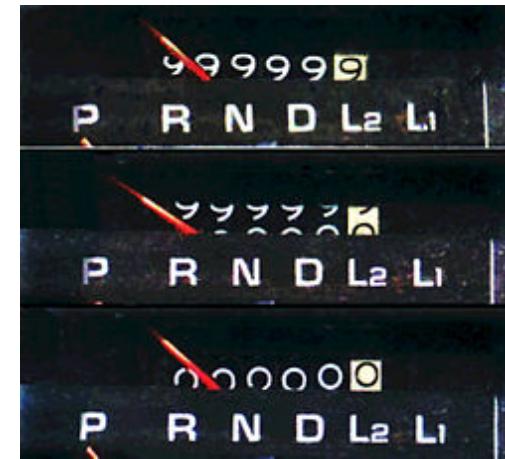
Source

Property of interest

Retrofitting Transformation



Source



Undefined Behavior

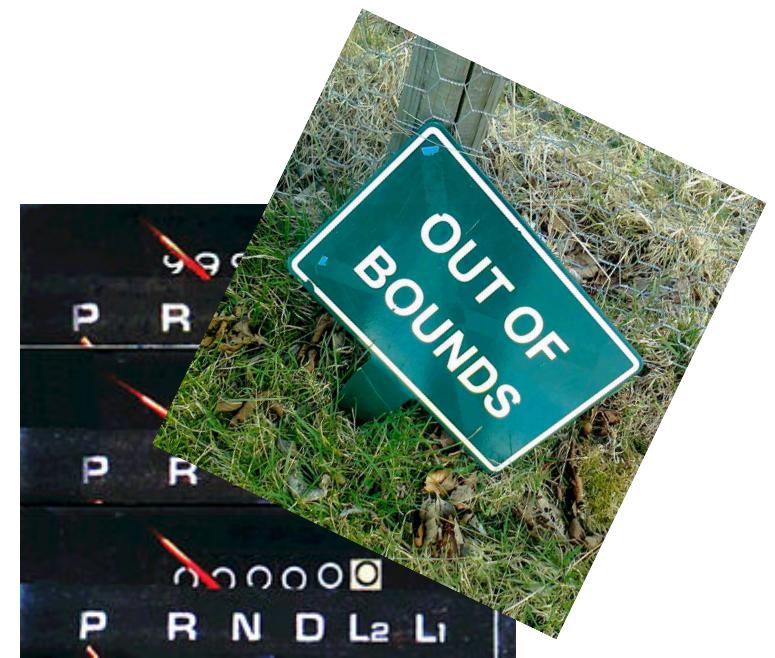


RUTGERS

Retrofitting Transformation



Source



Bounds Safety



RUTGERS

Retrofitting Transformation



Source

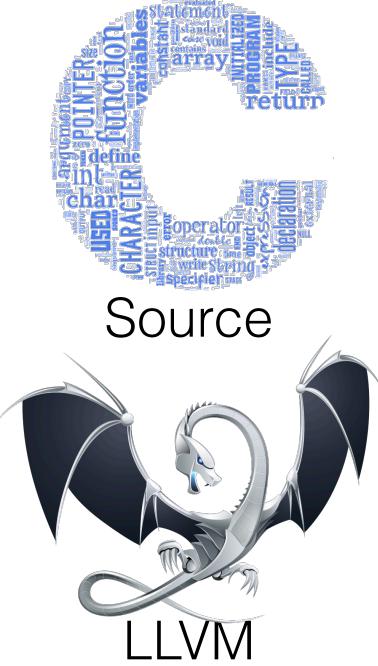


Control Flow Integrity



RUTGERS

Retrofitting Transformation



Source



```
.globl __add_forty_two<Int32>:Int32
.align 4, 0x90
__add_forty_two<Int32>:Int32:
.cfi_startproc
pushq %rbp
Ltmp1992:
.cfi_offset %rbp, 16
Ltmp1993:
.cfi_offset %rbp, -16
movq %rsp, %rbp
Ltmp1994:
.cfi_def_cfa_register %rbp
addl $42, %edi
movl %edi, %eax
popq %rbp
retq
.cfi_endproc
```

Assembly



Instrumentation



RUTGERS

Retrofitting Transformation



LLVM

:

Optimize



Instrument

Optimize



RUTGERS

Motivating Example

```
int foo(int a, int b) {  
    int c = a + b;  
    return c;  
}
```



LLVM : Optimize → Instrument → Optimize



RUTGERS

Motivating Example

```
int foo(int a, int b) {  
    int c = a + b;  
    return c;  
}
```

signed integer overflow
= undefined behavior



LLVM : Optimize → Instrument → Optimize



RUTGERS

Motivating Example

```
int foo(int a, int b) {  
    int c = a + b;  
    return c;  
}
```

signed integer overflow
= undefined behavior

Integer overflow checker

```
Instrumentation  
↓  
int foo(int a, int b) {  
    if (a > 0 && b > 0 && a > a + b)  
        exit(1);  
    if (a < 0 && b < 0 && a < a + b)  
        exit(1);  
    int c = a + b;  
    return c;  
}
```



LLVM : Optimize → Instrument → Optimize



RUTGERS

Motivating Example

```
int foo(int a, int b) {  
    int c = a + b;  
    return c;  
}
```

signed integer overflow
= undefined behavior

Integer overflow checker

Optimize (Linux LLVM)

Instrumentation removed

```
int foo(int a, int b) {  
    if (a > 0 && b > 0 && a > a + b)  
        exit(1);  
    if (a < 0 && b < 0 && a < a + b)  
        exit(1);  
    int c = a + b;  
    return c;  
}
```

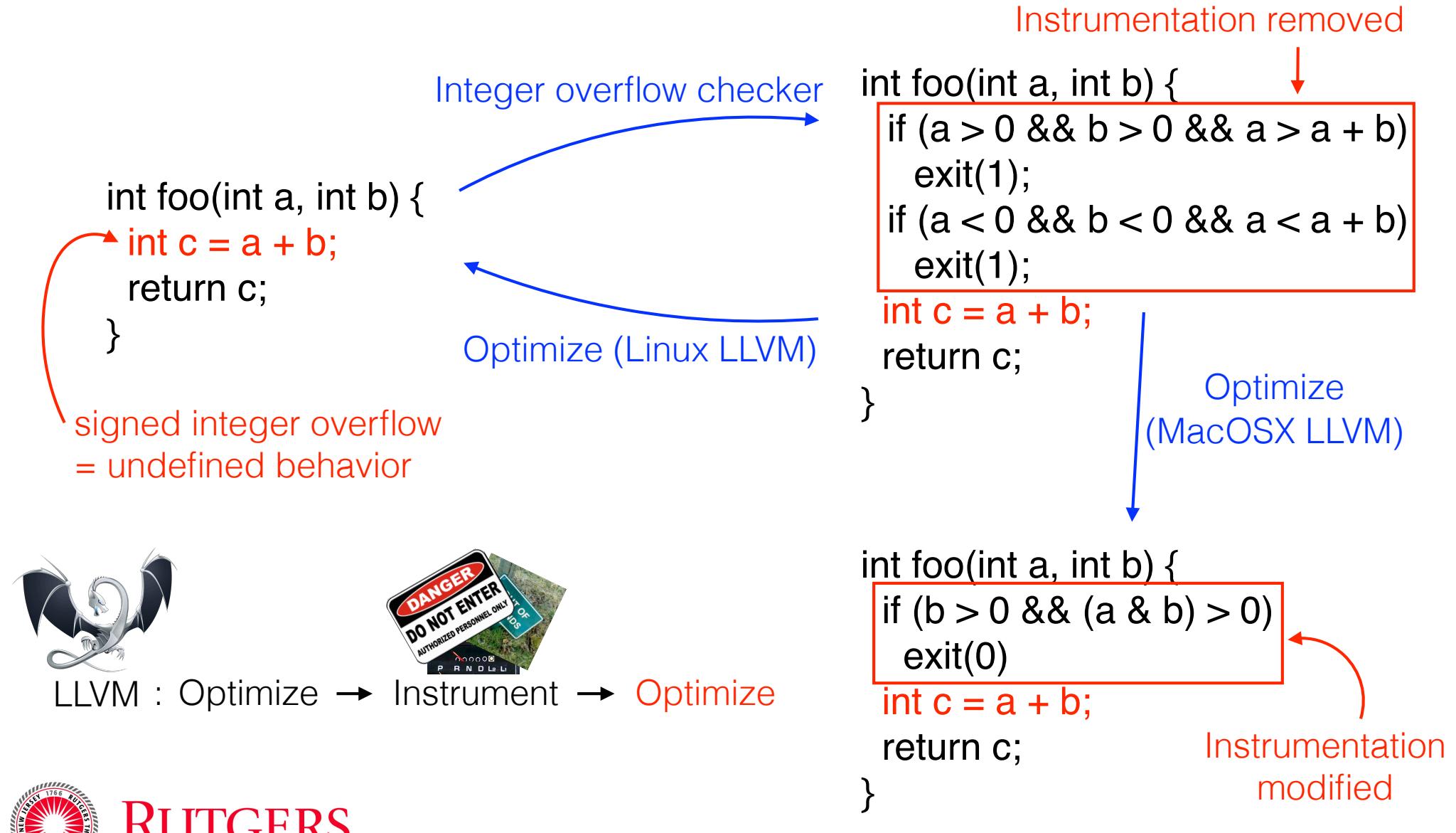


LLVM : Optimize → Instrument → Optimize



RUTGERS

Motivating Example



RUTGERS

Problem Statement

Can we detect erroneously removed/modified instrumentation due to compiler optimizations?

Challenges:

1. Checks may be completely removed.
2. Checks may be partially removed.
3. Checks may be moved.
4. Some checks are indeed redundant.

Solution:

Can we frame this as a reachability problem?



RUTGERS

P_{retro} : retrofitted program

```
int foo(int a, int b) {  
    if (a > 0 && b > 0 && a > a + b)  
        exit(1);  
    if (a < 0 && b < 0 && a < a + b)  
        exit(1);  
    int c = a + b;  
    return c;  
}
```

P_{opt} : optimized P_{retro}

```
int foo(int a, int b) {  
    if (b > 0 && (a & b) > 0)  
        exit(0);  
    int c = a + b;  
    return c;  
}
```



RUTGERS

Reachability

P_{retro} : retrofitted program

```
int foo(int a, int b) {  
    if (a > 0 && b > 0 && a > a + b)  
        exit(1);  
    if (a < 0 && b < 0 && a < a + b)  
        exit(1);  
    int c = a + b;  
    return c;  
}
```

P_{opt} : optimized P_{retro}

```
int foo(int a, int b) {  
    if (b > 0 && (a & b) > 0)  
        exit(0);  
    int c = a + b;  
    return c;  
}
```

Event of Interest



RUTGERS

Reachability

P_{retro} : retrofitted program

```
int foo(int a, int b) {  
    if (a > 0 && b > 0 && a > a + b)  
        exit(1);  
    if (a < 0 && b < 0 && a < a + b)  
        exit(1);  
    int c = a + b;  
    return c;  
}
```

P_{opt} : optimized P_{retro}

```
int foo(int a, int b) {  
    if (b > 0 && (a & b) > 0)  
        exit(0);  
    int c = a + b;  
    return c;  
}
```

Direct unsafe execution to exit path.

Reachability

P_{retro} : retrofitted program

```
int foo(int a, int b) {  
    if (a > 0 && b > 0 && a > a + b)  
        exit(1);  
    if (a < 0 && b < 0 && a < a + b)  
        exit(1);  
    int c = a + b;  
    return c;  
}
```

Direct unsafe execution to exit path.

P_{opt} : optimized P_{retro}

```
int foo(int a, int b) {  
    if (b > 0 && (a & b) > 0)  
        exit(0);  
    int c = a + b;  
    return c;  
}
```

Semantics same as P_{retro}



Reachability

P_{retro} : retrofitted program

```
int foo(int a, int b) {  
    if (a > 0 && b > 0 && a > a + b)  
        exit(1);  
    if (a < 0 && b < 0 && a < a + b)  
        exit(1);  
    int c = a + b;  
    return c;  
}
```

Direct unsafe execution to exit path.

P_{opt} : optimized P_{retro}

```
int foo(int a, int b) {  
    if (b > 0 && (a & b) > 0)  
        exit(0);  
    int c = a + b;  
    return c;  
}
```

Semantics same as P_{retro}

Given same inputs,



RUTGERS

Reachability

P_{retro} : retrofitted program

```
int foo(int a, int b) {  
    if (a > 0 && b > 0 && a > a + b)  
        exit(1);  
    if (a < 0 && b < 0 && a < a + b)  
        exit(1);  
    int c = a + b;  
    return c;  
}
```

Direct unsafe execution to exit path.

P_{opt} : optimized P_{retro}

```
int foo(int a, int b) {  
    if (b > 0 && (a & b) > 0)  
        exit(0);  
    int c = a + b;  
    return c;  
}
```

Semantics same as P_{retro}

Given same inputs,

1. P_{retro} and P_{opt} reaches the Event of Interest



RUTGERS

Reachability

P_{retro} : retrofitted program

```
int foo(int a, int b) {  
    if (a > 0 && b > 0 && a > a + b)  
        exit(1);  
    if (a < 0 && b < 0 && a < a + b)  
        exit(1);  
    int c = a + b;  
    return c;  
}
```

Direct unsafe execution to exit path.

P_{opt} : optimized P_{retro}

```
int foo(int a, int b) {  
    if (b > 0 && (a & b) > 0)  
        exit(0);  
    int c = a + b;  
    return c;  
}
```

Semantics same as P_{retro}

Given same inputs,

1. P_{retro} and P_{opt} reaches the Event of Interest
2. P_{retro} and P_{opt} reaches exit(0)



Reachability

P_{retro}: retrofitted program

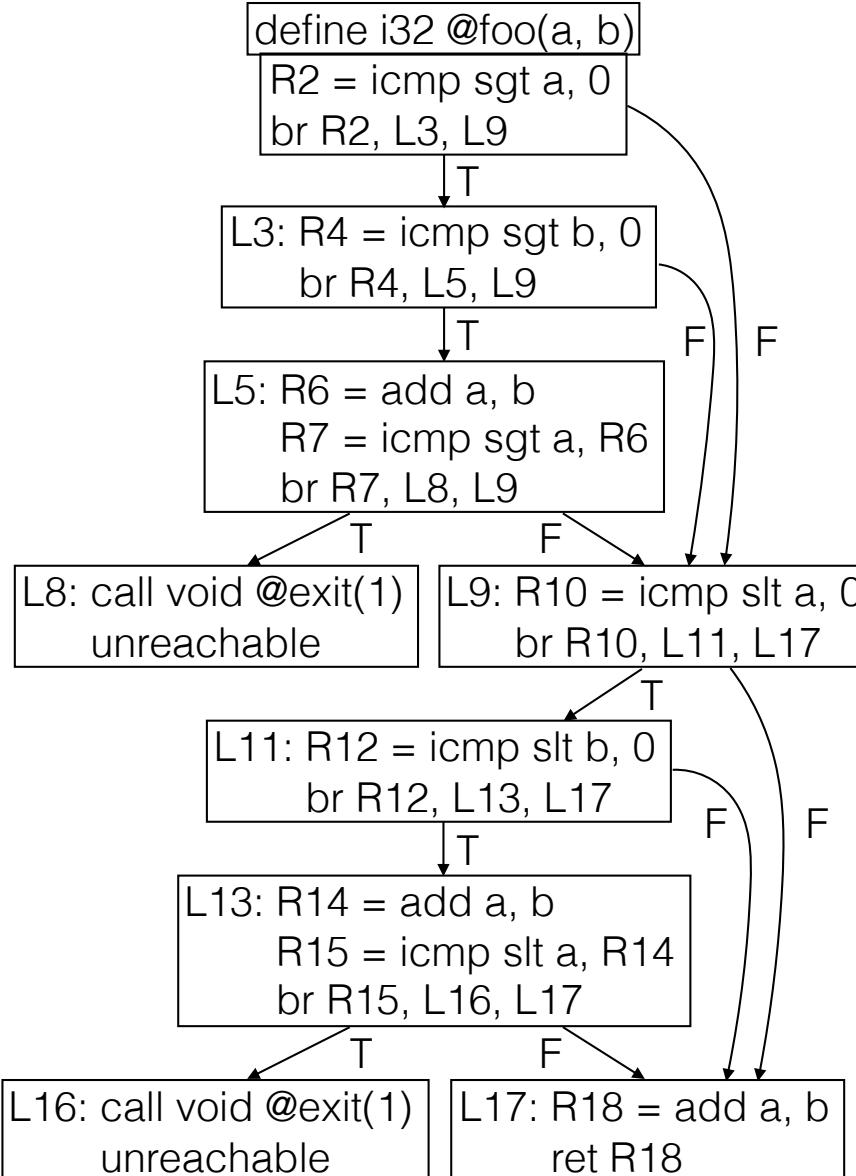
```
int foo(int a, int b) {  
    if (a > 0 && b > 0 && a > a + b)  
        exit(1);  
    if (a < 0 && b < 0 && a < a + b)  
        exit(1);  
    int c = a + b;  
    return c;  
}
```



RUTGERS

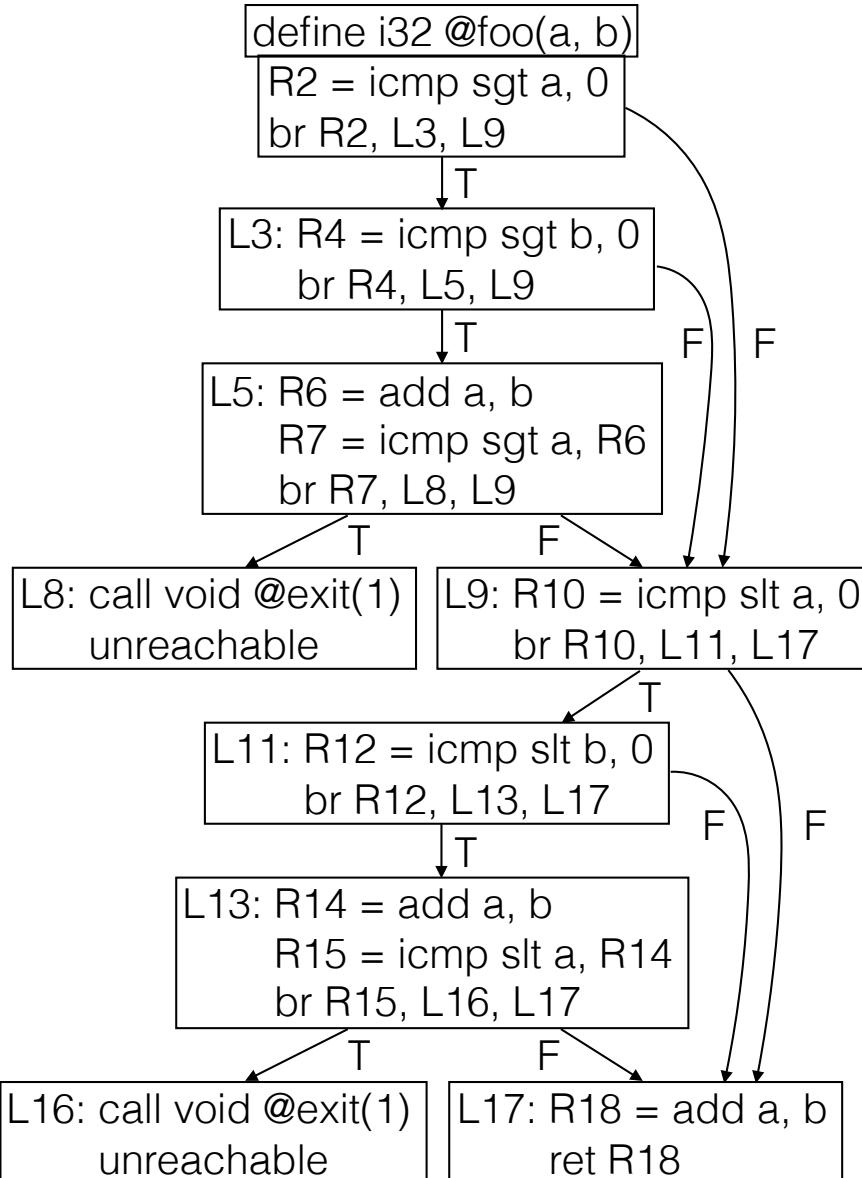
Reachability

P_{retro}: retrofitted program



Reachability

P_{retro} : retrofitted program

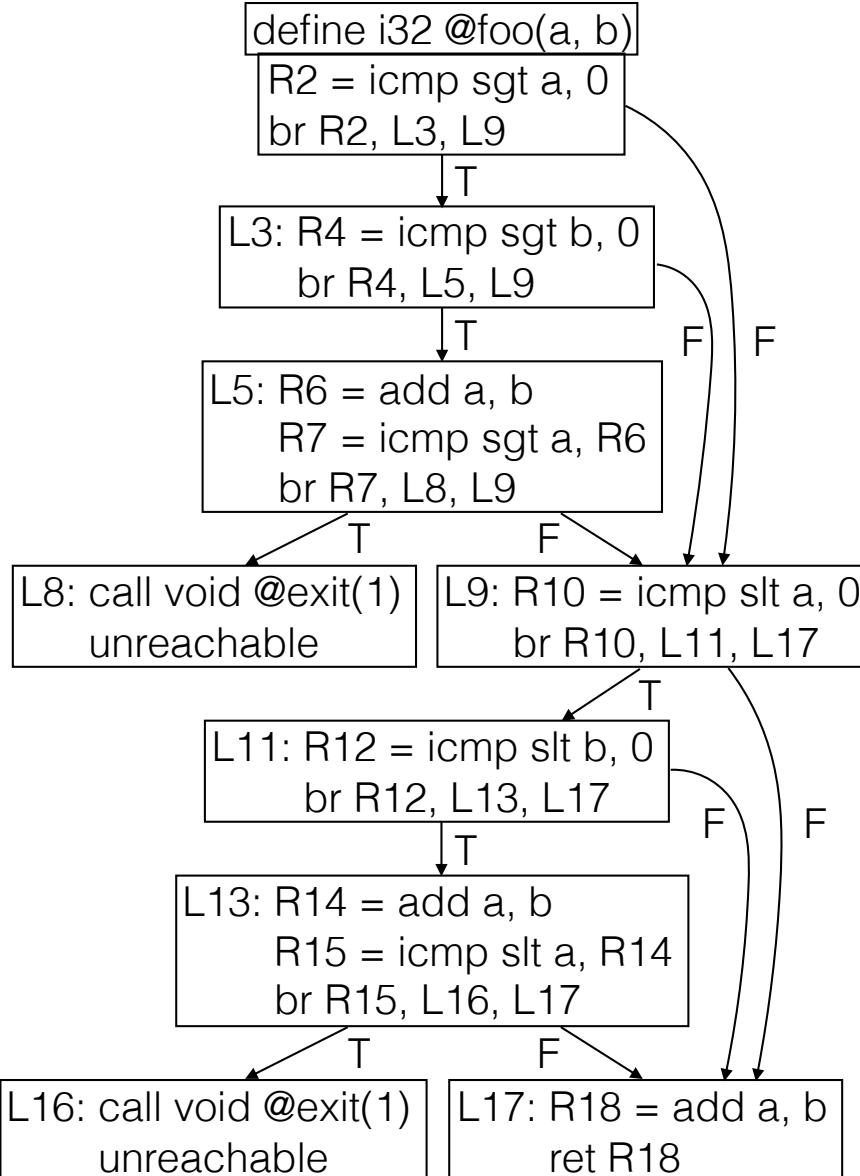


P_{opt} : optimized P_{retro}

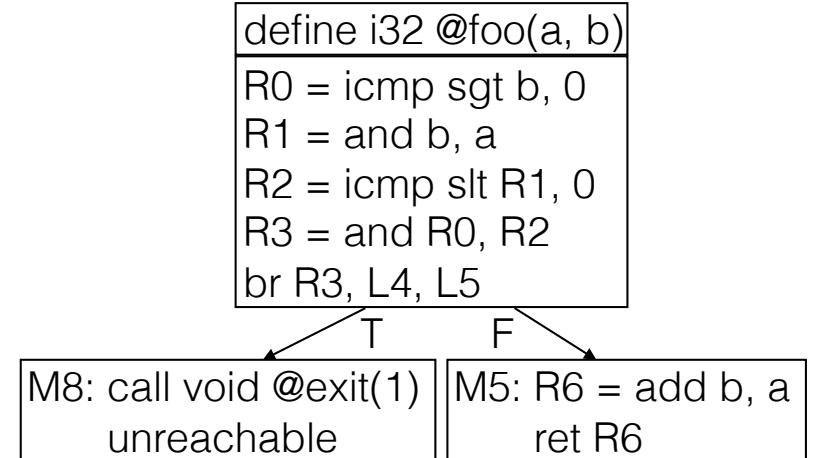
```
int foo(int a, int b) {  
    if (b > 0 && (a & b) > 0)  
        exit(0);  
    int c = a + b;  
    return c;  
}
```

Reachability

P_{retro} : retrofitted program

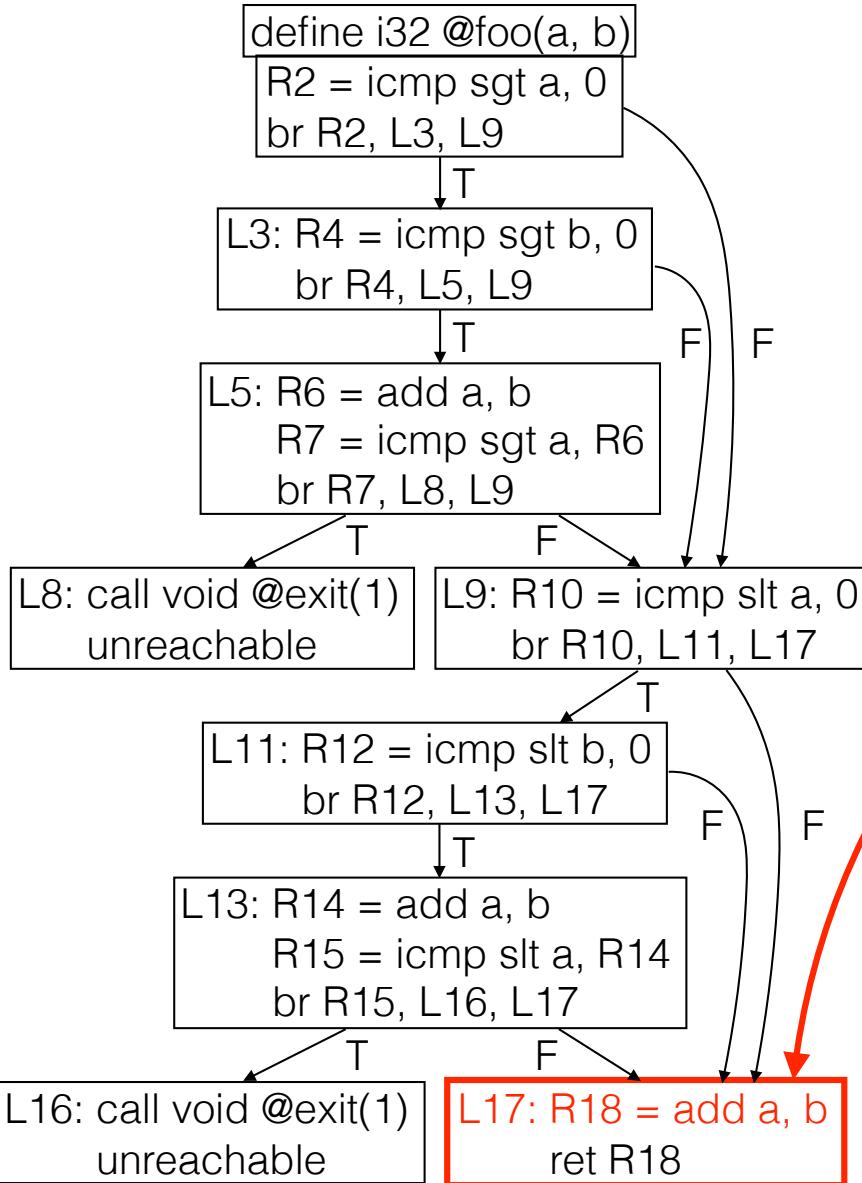


P_{opt} : optimized P_{retro}

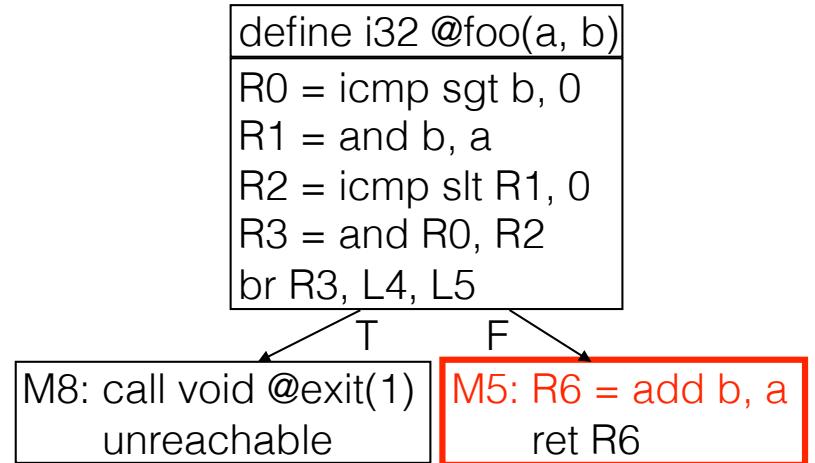


Reachability

P_{retro} : retrofitted program



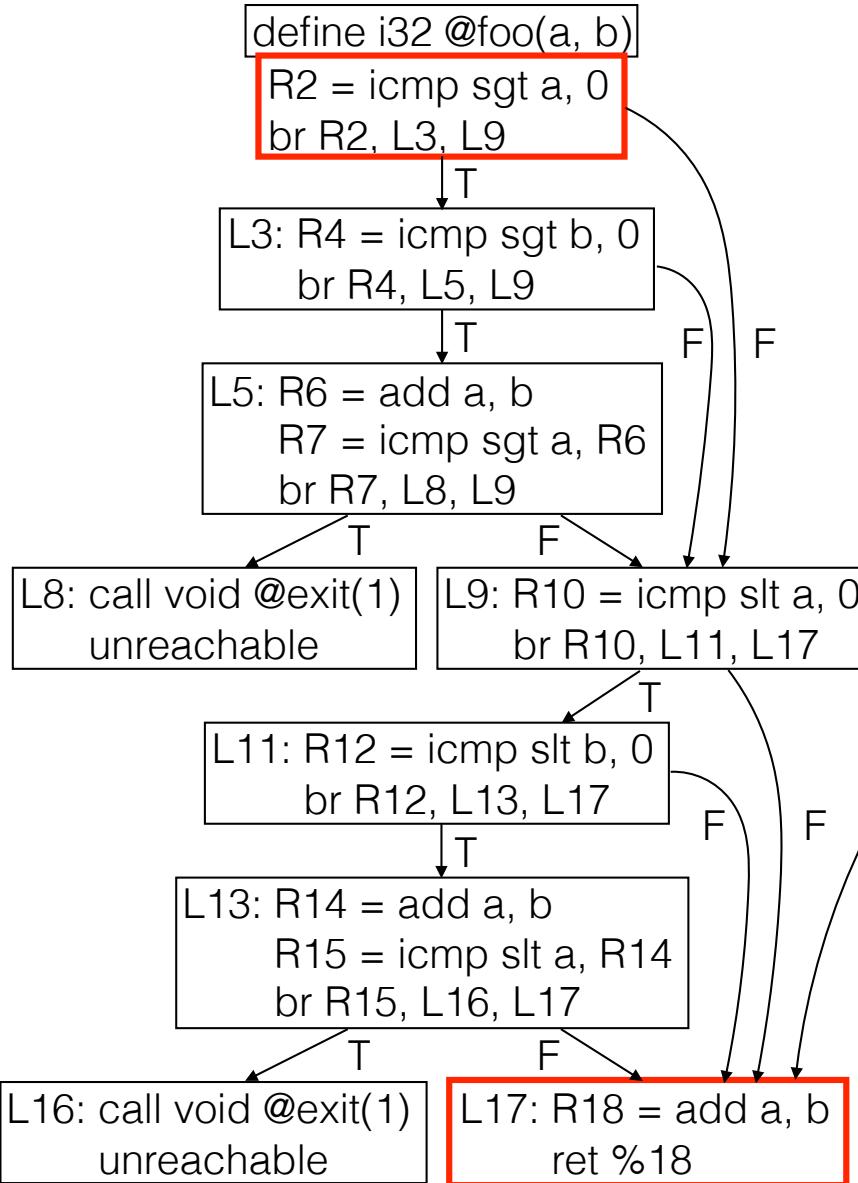
P_{opt} : optimized P_{retro}



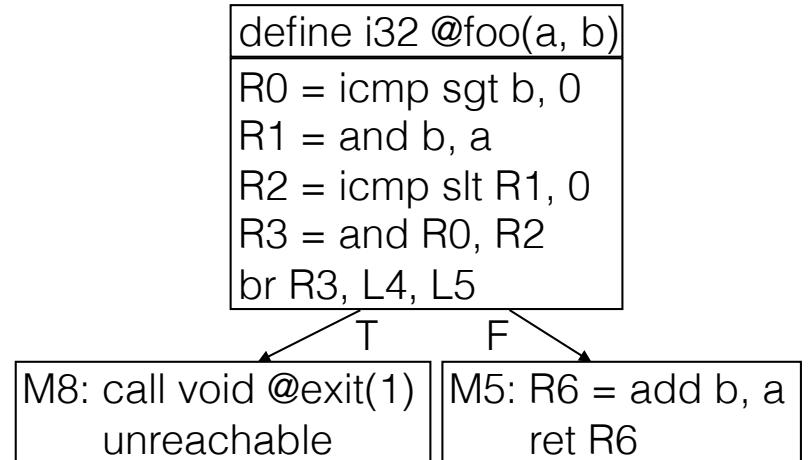
Event of interest

Reachability

P_{retro} : retrofitted program



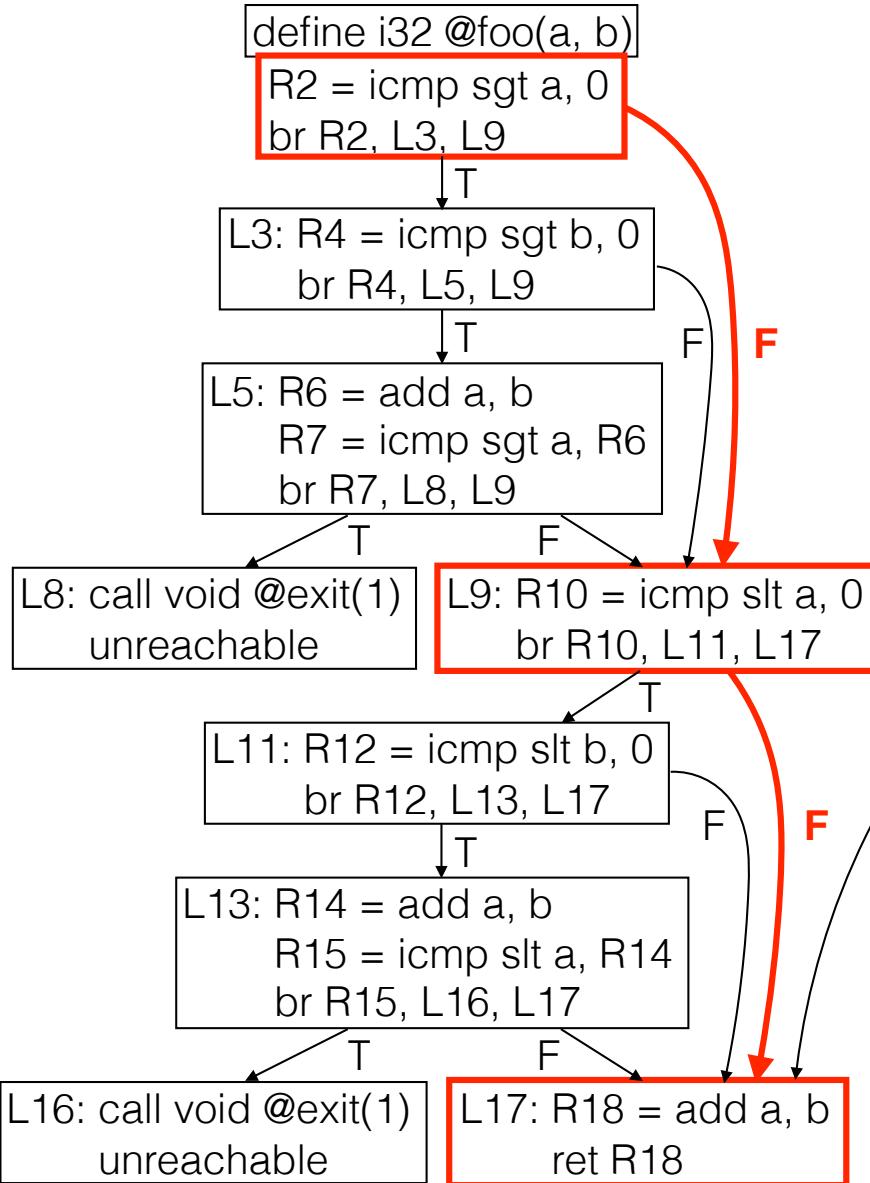
P_{opt} : optimized P_{retro}



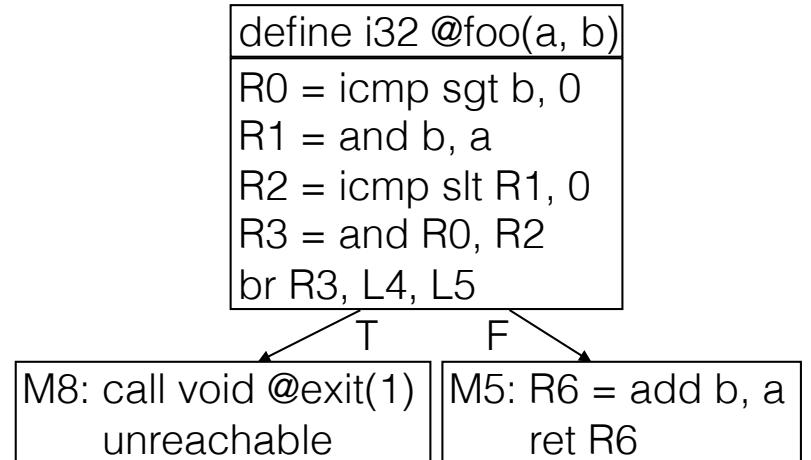
Event of interest

Reachability

P_{retro} : retrofitted program



P_{opt} : optimized P_{retro}



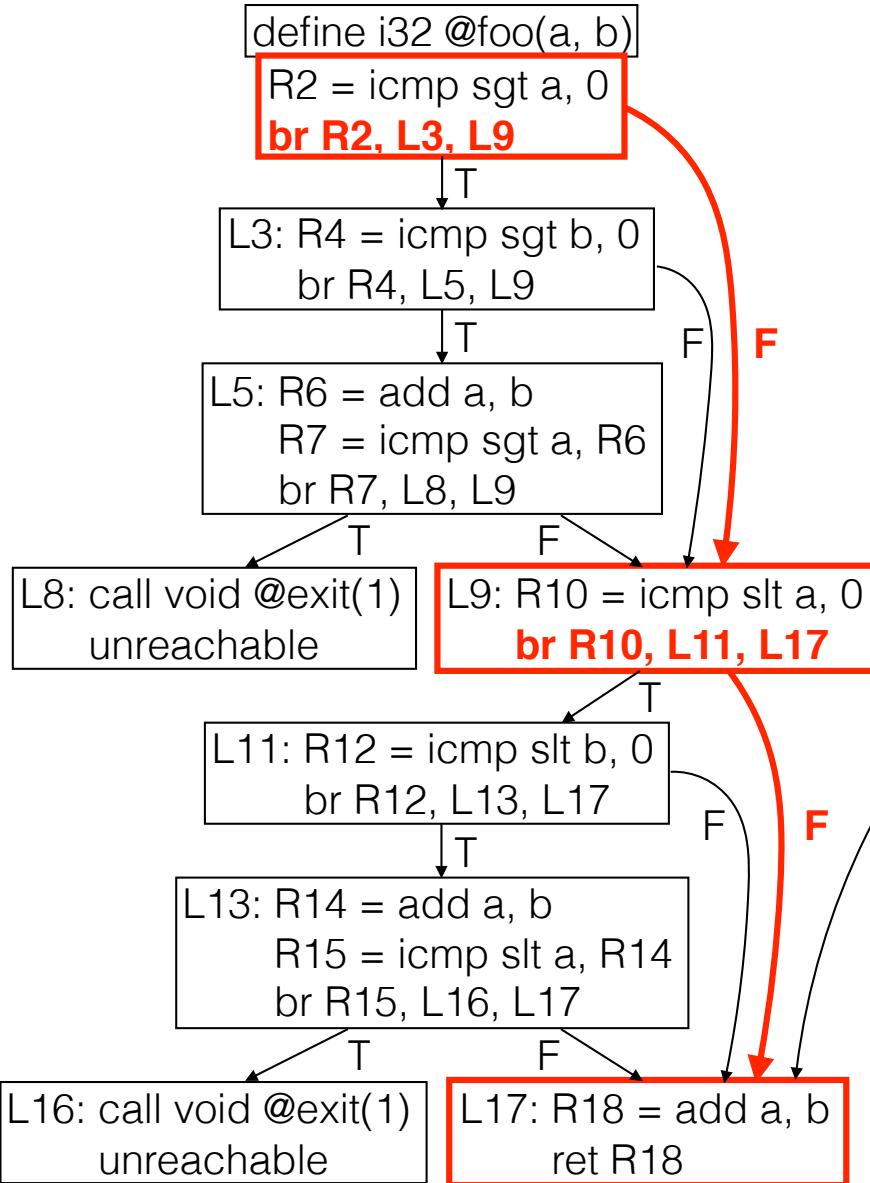
Event of interest



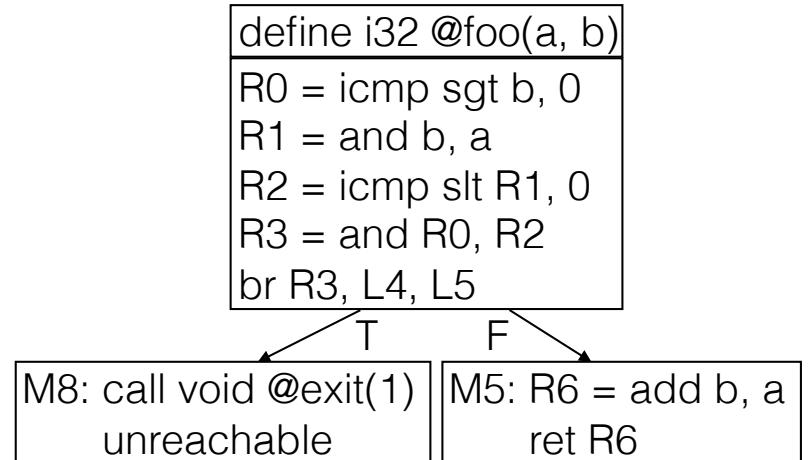
RUTGERS

Reachability

P_{retro} : retrofitted program



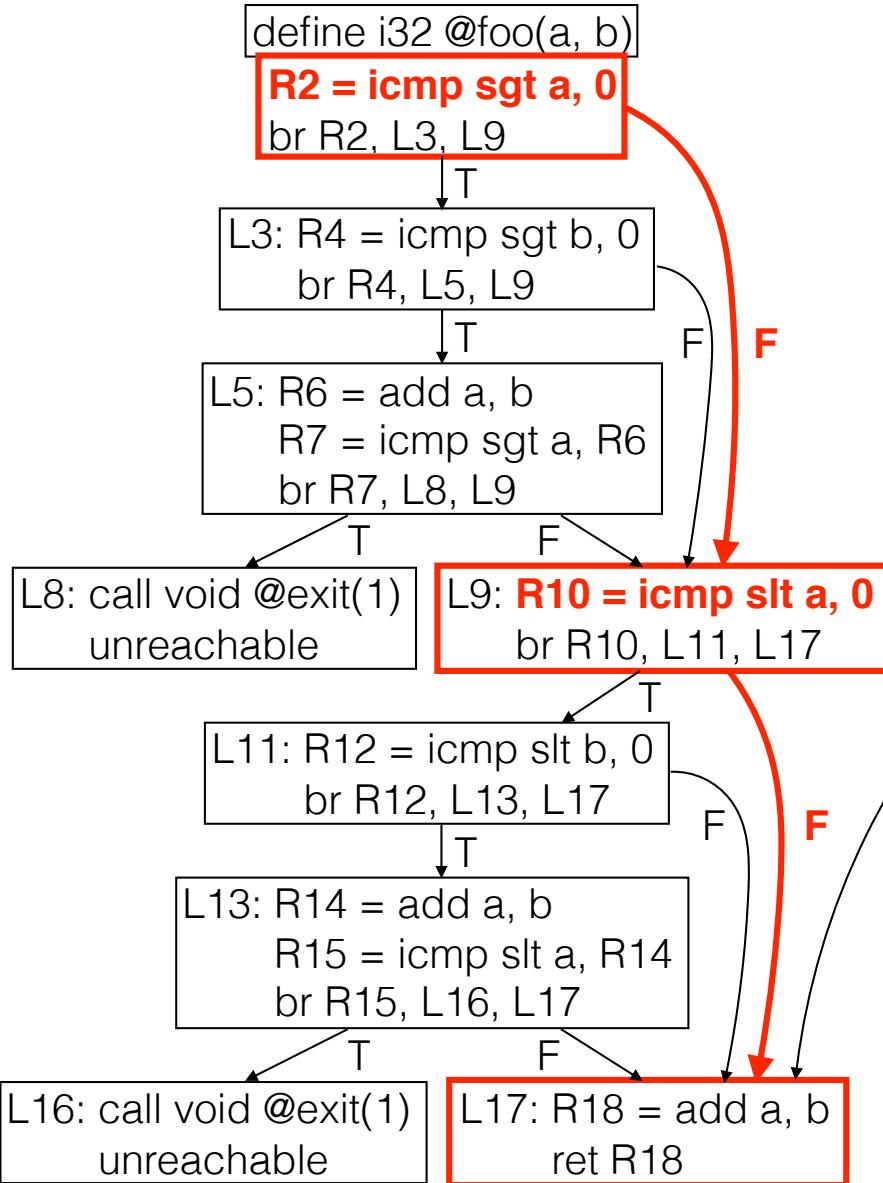
P_{opt} : optimized P_{retro}



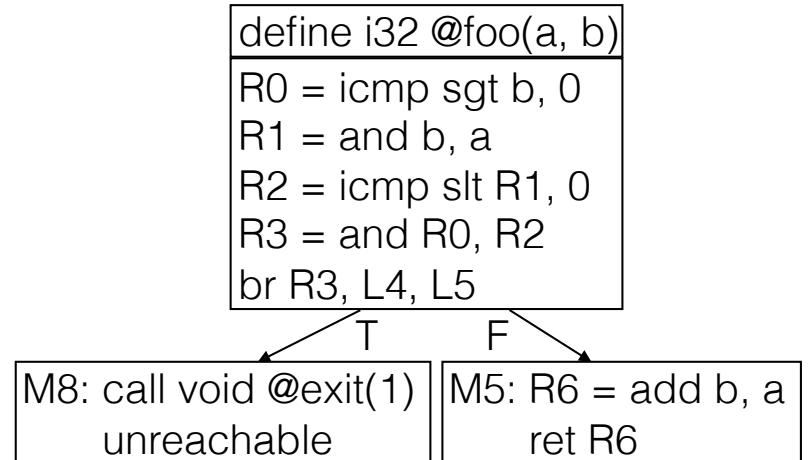
$$P_{\text{0retro}}: \neg(R2) \wedge \neg(R10)$$

Reachability

P_{retro} : retrofitted program



P_{opt} : optimized P_{retro}



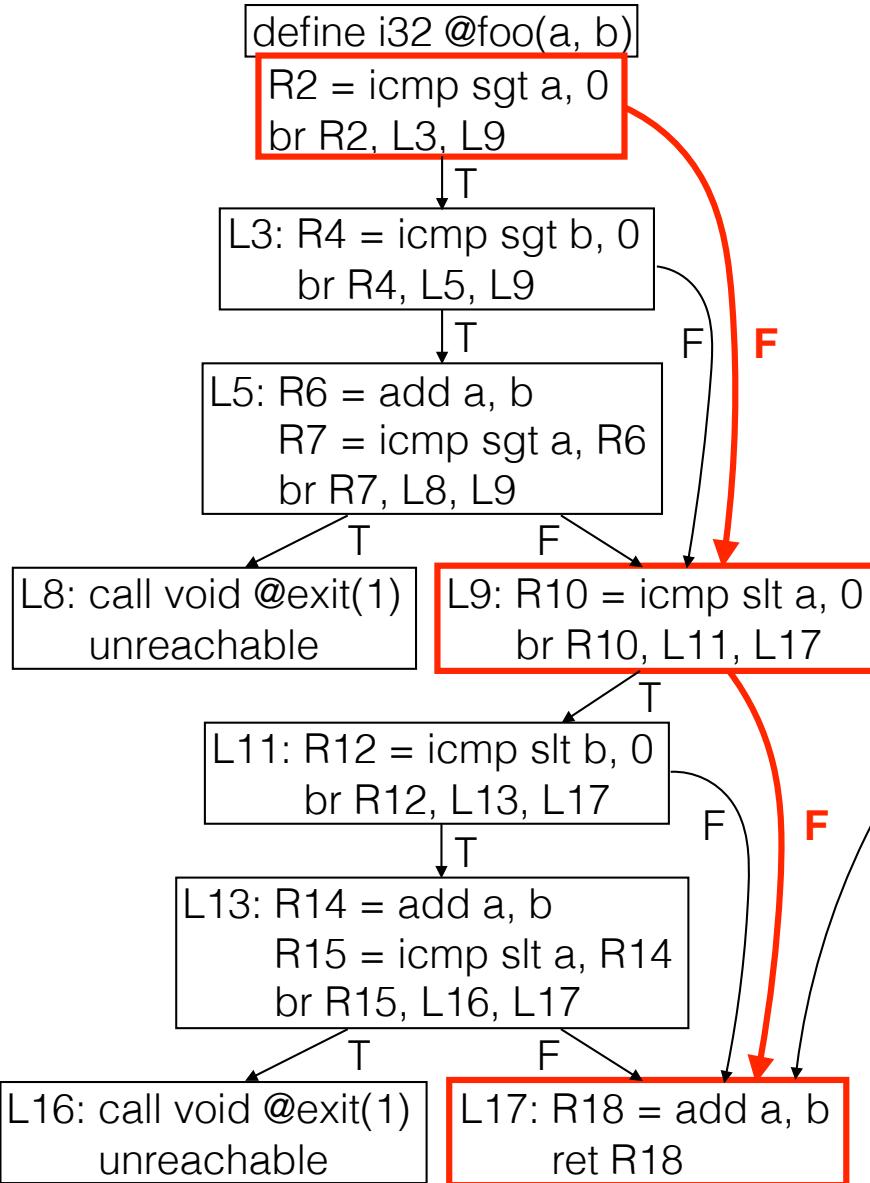
Event of interest

$$P_{0\text{retro}}: \neg(R2) \wedge \neg(R10) \wedge \\ (R2 = (a > 0)) \wedge (R10 = (a < 0))$$

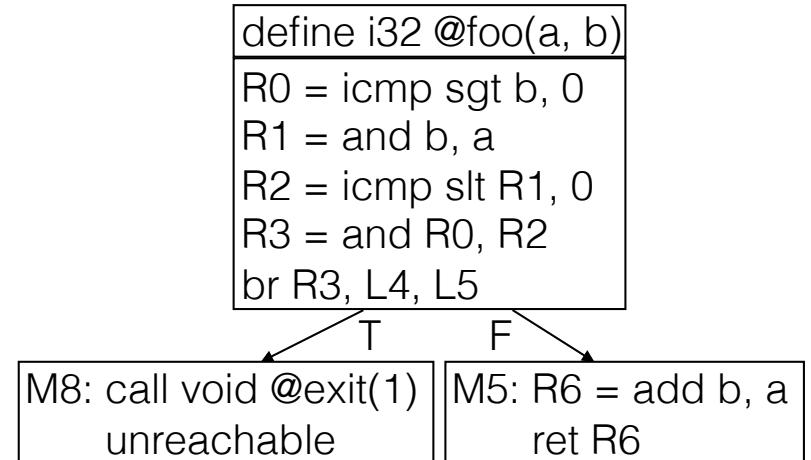


Reachability

P_{retro} : retrofitted program



P_{opt} : optimized P_{retro}



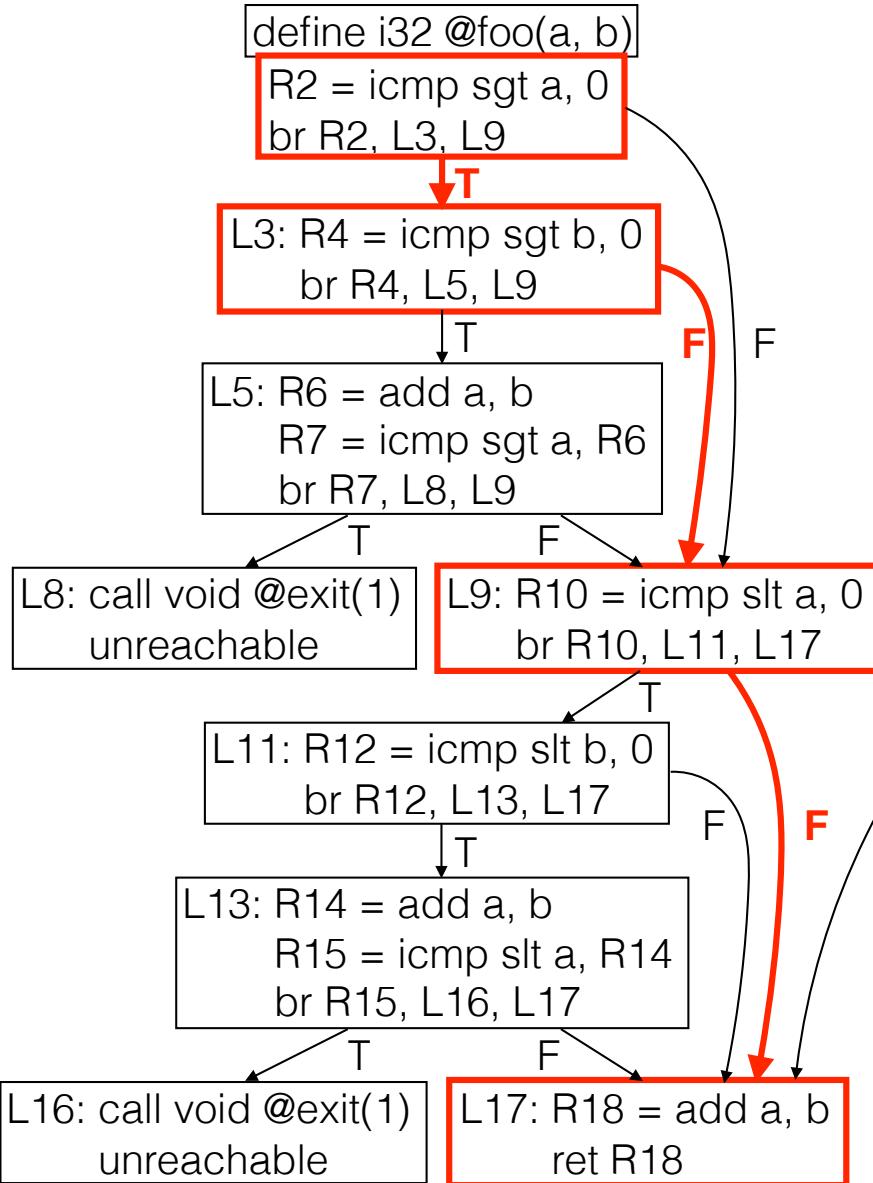
Event of interest

$$P_{0\text{retro}}: !(R2) \wedge !(R10) \wedge (R2 = (a > 0)) \wedge (R10 = (a < 0))$$

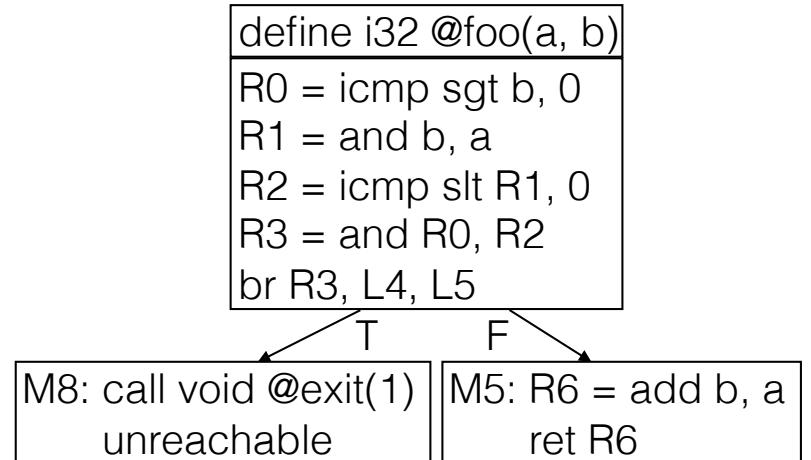
$$E_{\text{retro}} = P_{0\text{retro}}$$

Reachability

P_{retro} : retrofitted program



P_{opt} : optimized P_{retro}

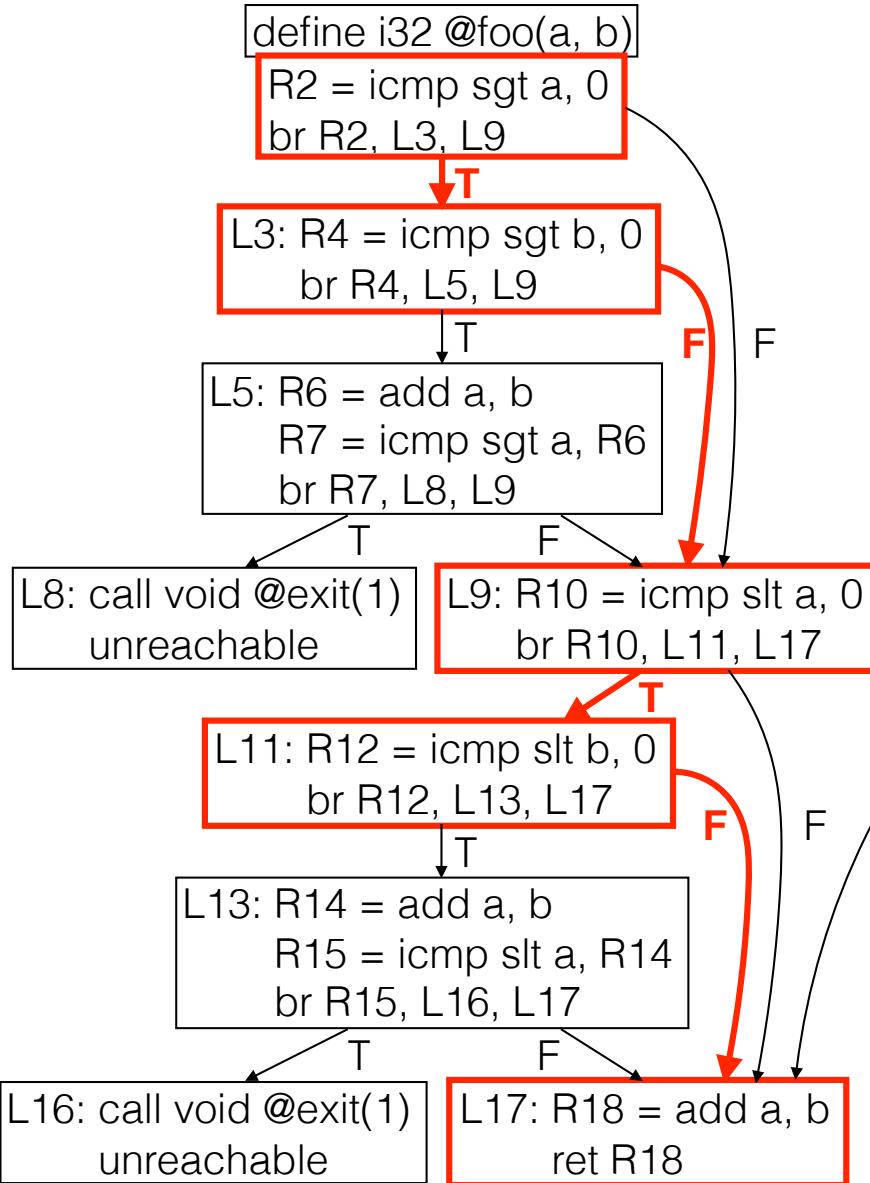


$$P_{\text{0retro}}: !(R2) \wedge !(R10) \wedge (R2 = (a > 0)) \wedge (R10 = (a < 0))$$

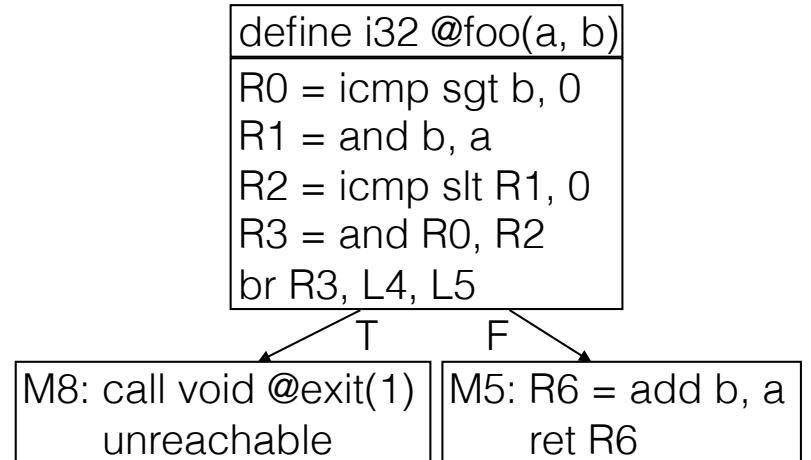
$$E_{\text{retro}} = P_{\text{0retro}} \vee \dots \vee P_{\text{iretro}}$$

Reachability

P_{retro} : retrofitted program



P_{opt} : optimized P_{retro}

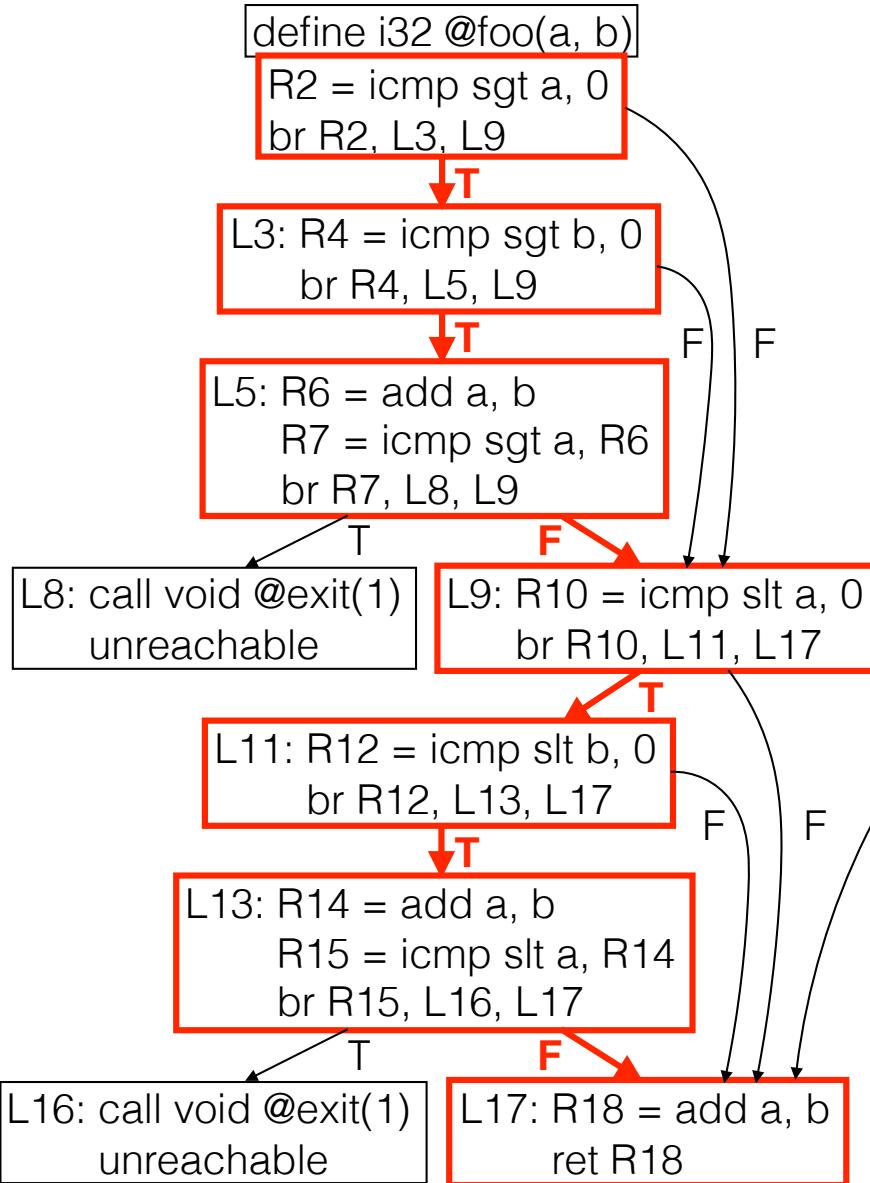


$$P_{\text{0retro}}: !(R2) \wedge !(R10) \wedge (R2 = (a > 0)) \wedge (R10 = (a < 0))$$

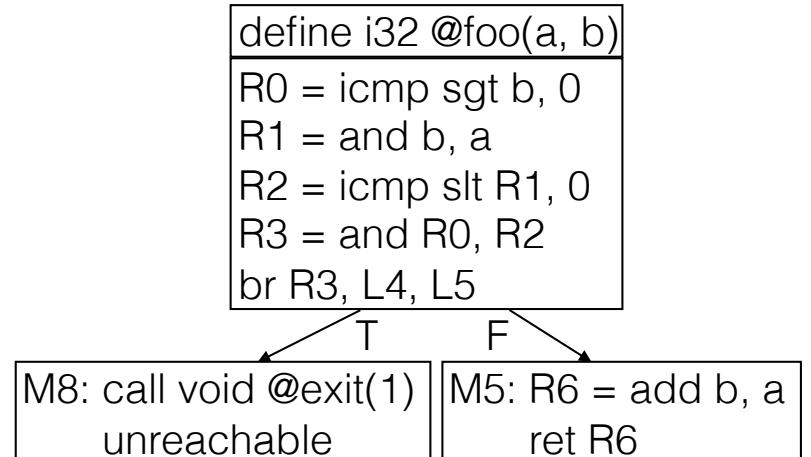
$$E_{\text{retro}} = P_{\text{0retro}} \vee \dots \vee P_{\text{Iretro}}$$

Reachability

P_{retro} : retrofitted program



P_{opt} : optimized P_{retro}



Event of interest

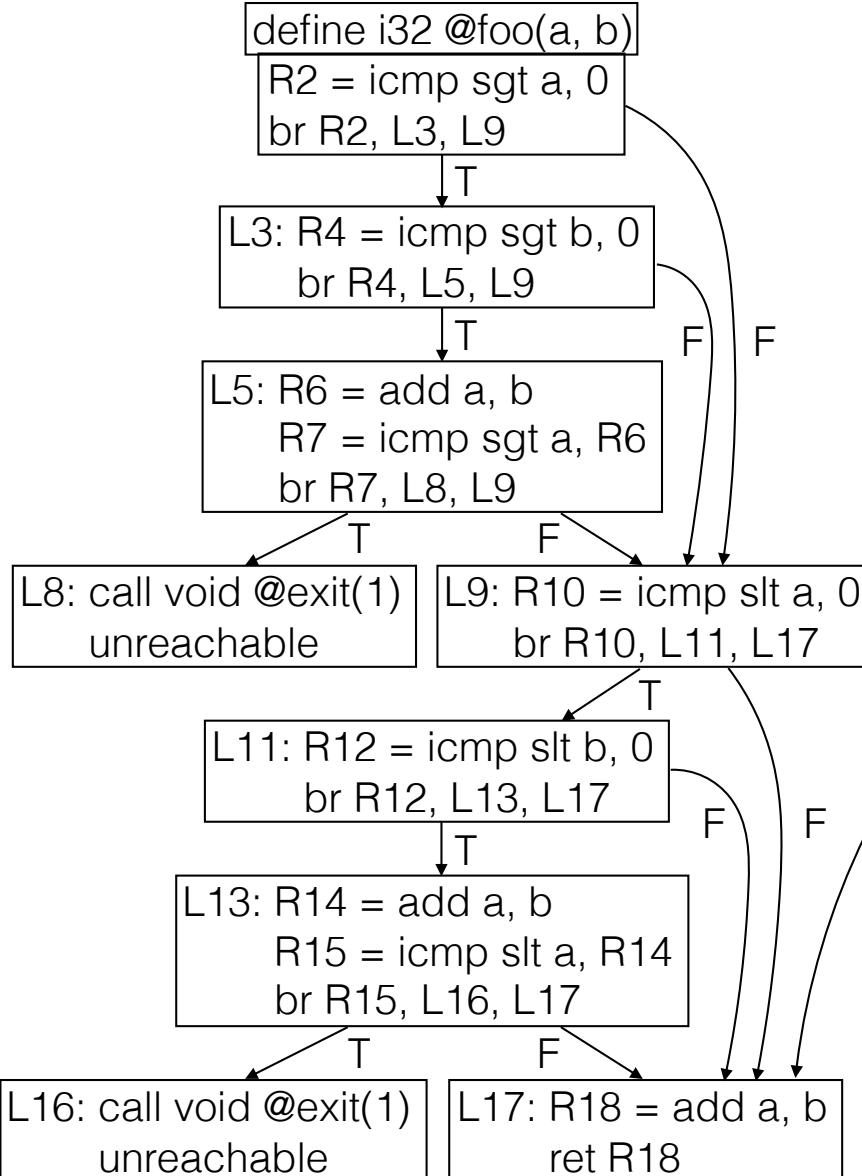
$$P_{0\text{retro}}: !(R2) \wedge !(R10) \wedge (R2 = (a > 0)) \wedge (R10 = (a < 0))$$

$$E_{\text{retro}} = P_{0\text{retro}} \vee \dots \vee P_{i\text{retro}}$$

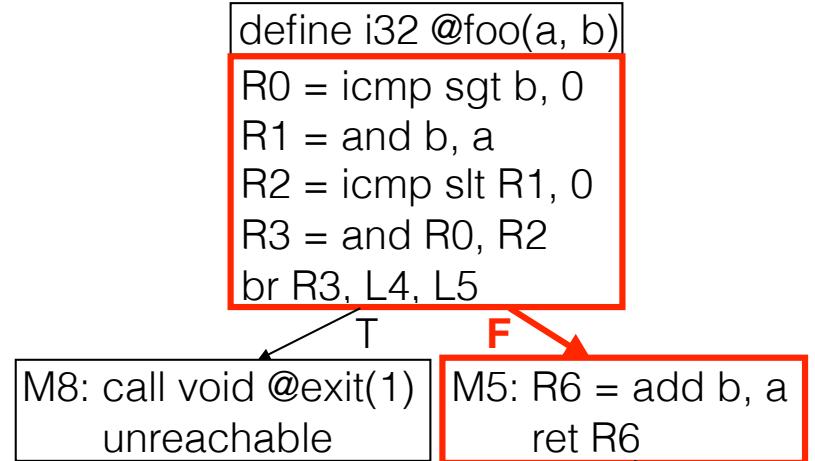


Reachability

P_{retro} : retrofitted program



P_{opt} : optimized P_{retro}



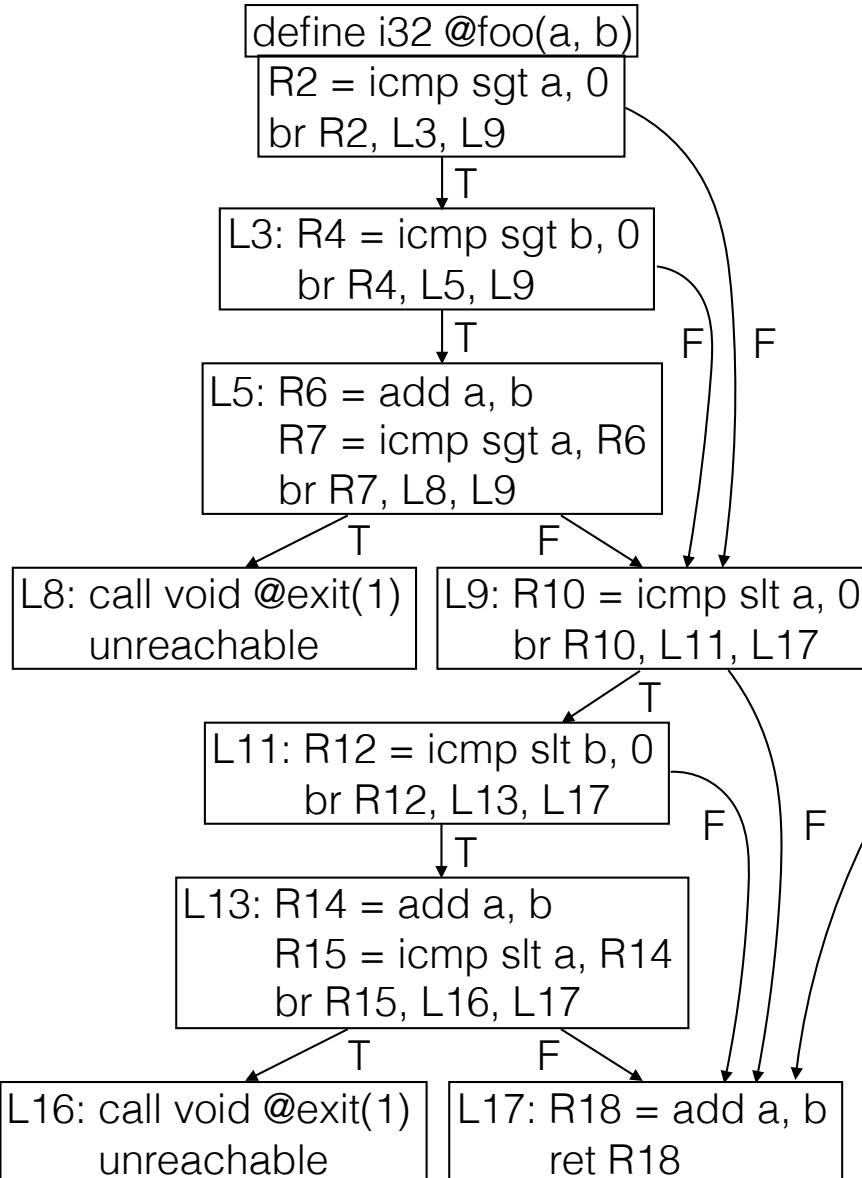
$$P_{\text{retro}}: !(R2) \wedge !(R10) \wedge (R2 = (a > 0)) \wedge (R10 = (a < 0))$$

$$E_{\text{retro}} = P_{\text{retro}} \vee \dots \vee P_i_{\text{retro}}$$

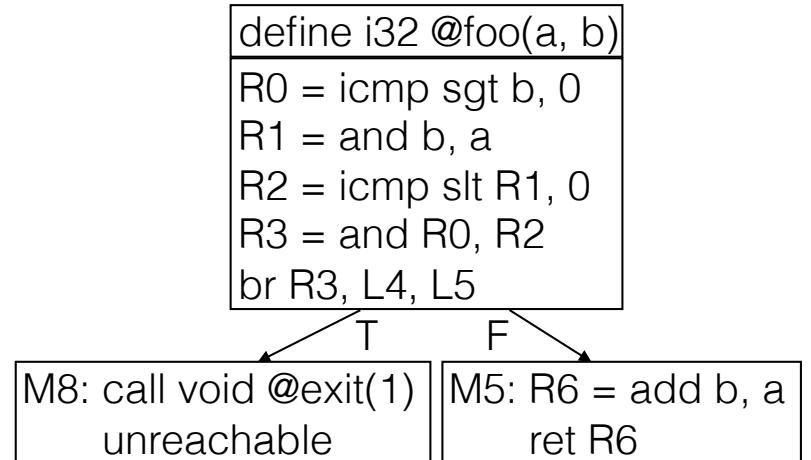
$$\mathbf{E}_{\text{opt}} = \mathbf{P}_{\text{opt}} \vee \dots \vee \mathbf{P}_j_{\text{opt}}$$

Reachability

P_{retro} : retrofitted program



P_{opt} : optimized P_{retro}



$$P_{\text{retro}}: !(R2) \wedge !(R10) \wedge (R2 = (a > 0)) \wedge (R10 = (a < 0))$$

$$E_{\text{retro}} = P_{\text{retro}} \vee \dots \vee P_i_{\text{retro}}$$

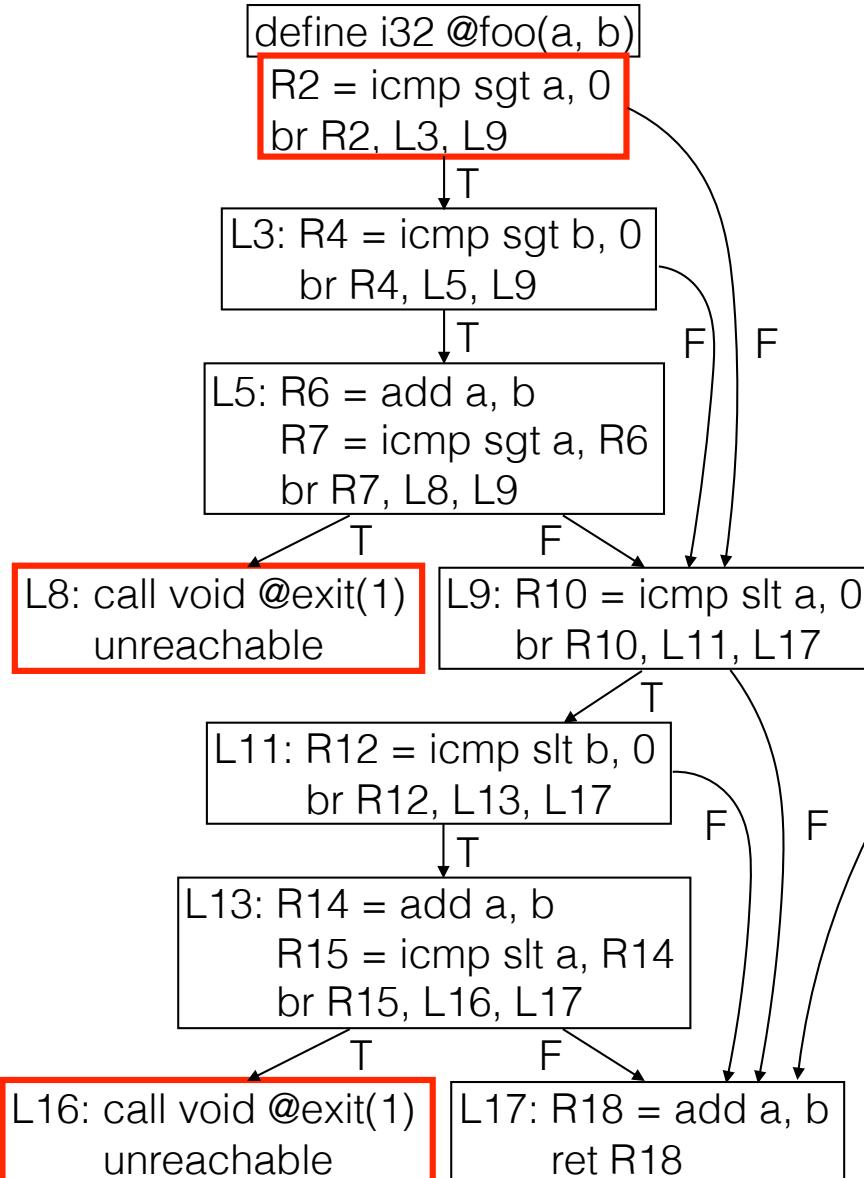
$$E_{\text{opt}} = P_{\text{opt}} \vee \dots \vee P_j_{\text{opt}}$$

If P_{retro} reaches event of interest, then P_{opt} reaches event of interest
($E_{\text{retro}} \Rightarrow E_{\text{opt}}$)

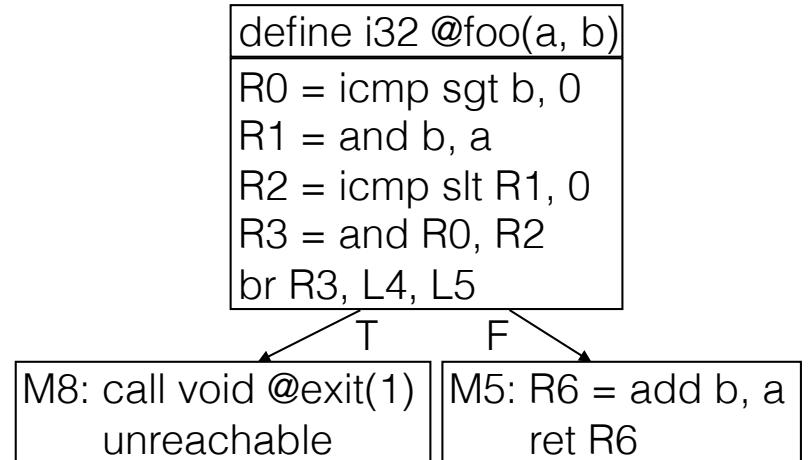


Reachability

P_{retro} : retrofitted program



P_{opt} : optimized P_{retro}



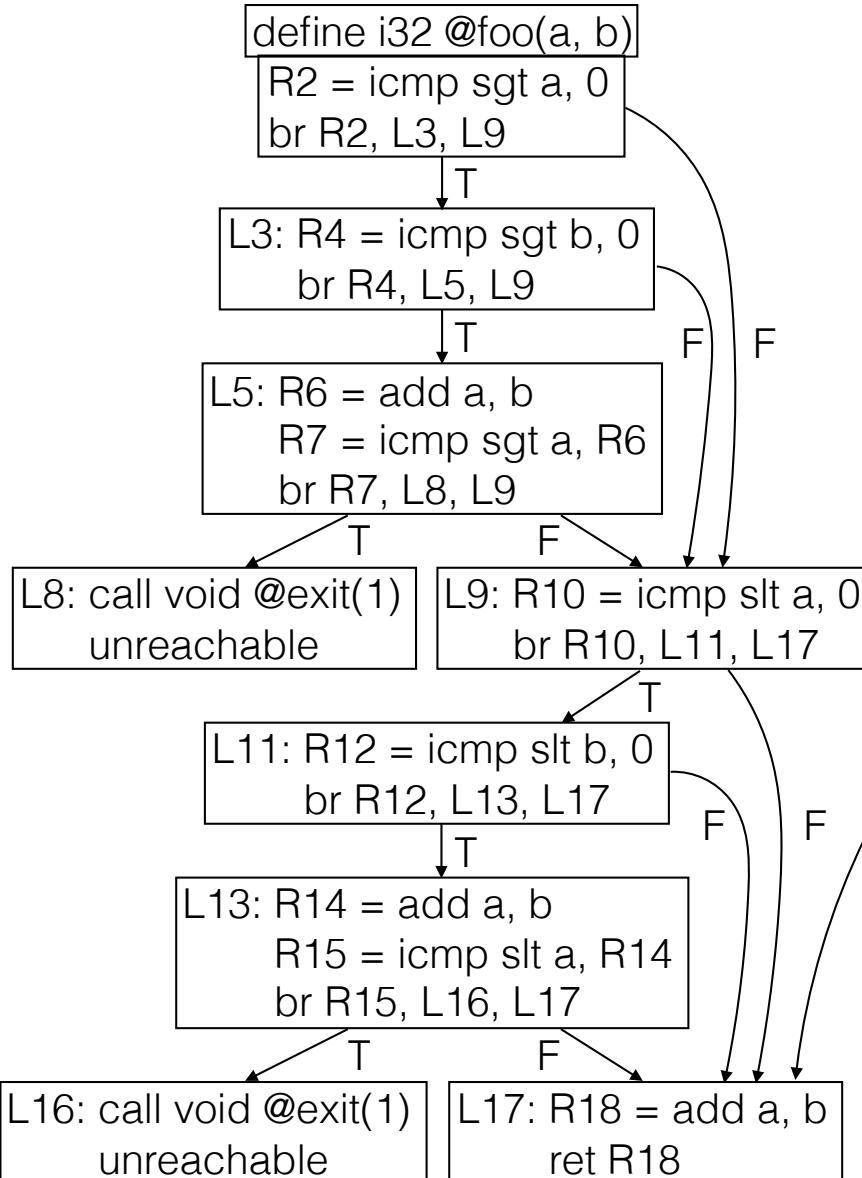
Event of interest

$P_{\text{retro}}: !(R2) \wedge !(R10) \wedge$
 $(R2 = (a > 0)) \wedge (R10 = (a < 0))$
 $E_{\text{retro}} = P_{\text{retro}} \vee \dots \vee P_i_{\text{retro}}$
 $E_{\text{opt}} = P_{\text{opt}} \vee \dots \vee P_j_{\text{opt}}$
 If P_{retro} reaches event of interest, then P_{opt} reaches event of interest ($E_{\text{retro}} \Rightarrow E_{\text{opt}}$)

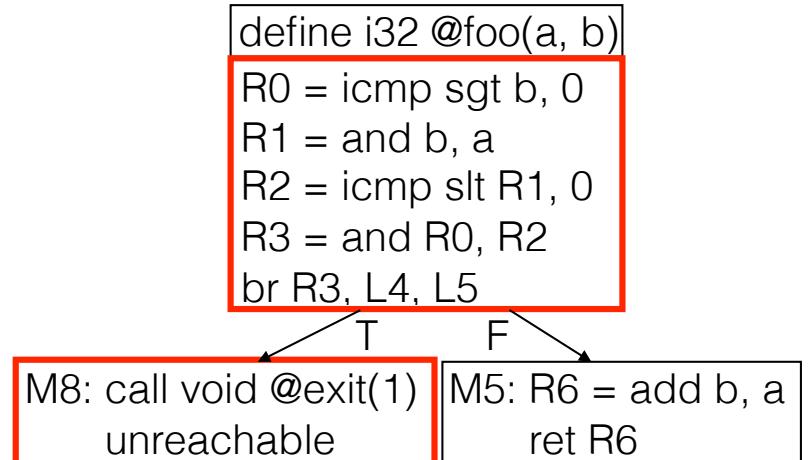
$\sim E_{\text{retro}}$

Reachability

P_{retro} : retrofitted program



P_{opt} : optimized P_{retro}



Event of interest

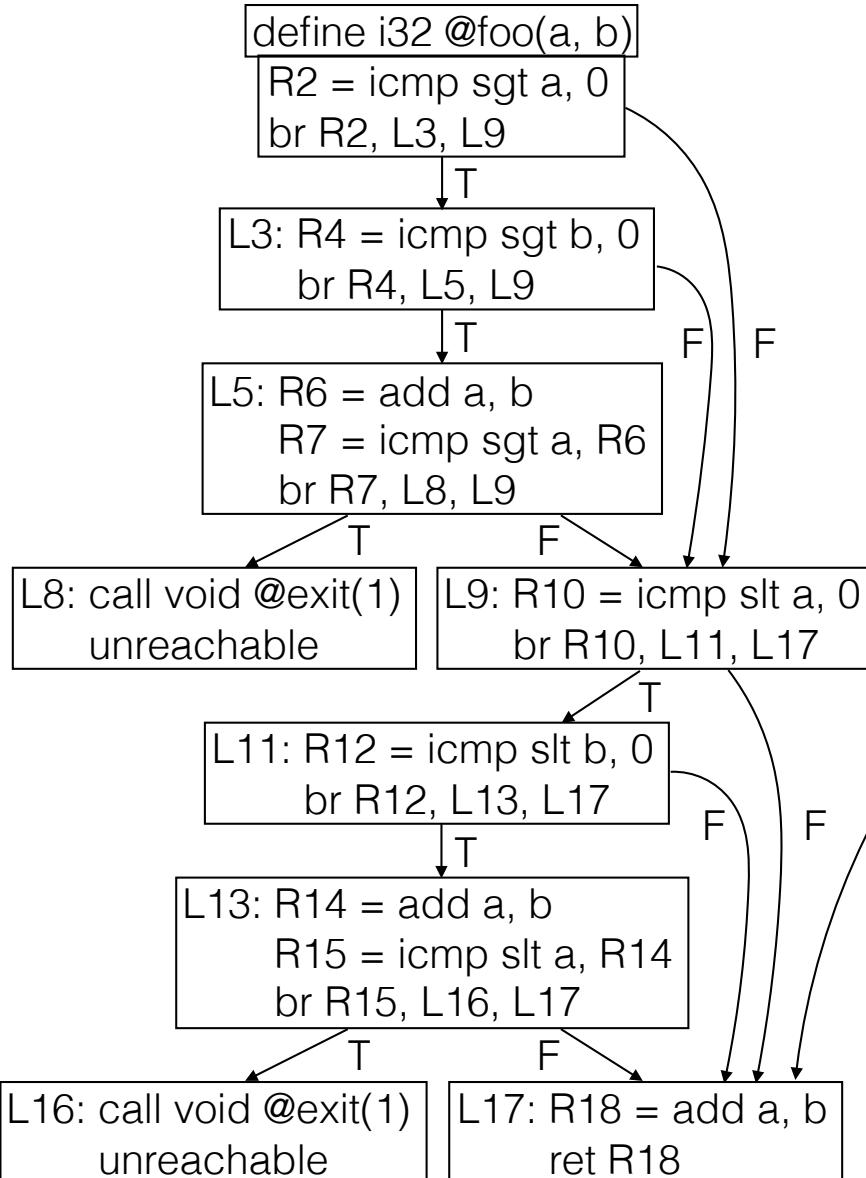
$P_{\text{retro}}: !(R2) \wedge !(R10) \wedge$
 $(R2 = (a > 0)) \wedge (R10 = (a < 0))$
 $E_{\text{retro}} = P_{\text{retro}} \vee \dots \vee P_i_{\text{retro}}$
 $E_{\text{opt}} = P_{\text{opt}} \vee \dots \vee P_j_{\text{opt}}$
 If P_{retro} reaches event of interest, then P_{opt} reaches event of interest ($E_{\text{retro}} \Rightarrow E_{\text{opt}}$)

$\sim E_{\text{retro}}, \sim E_{\text{opt}}$

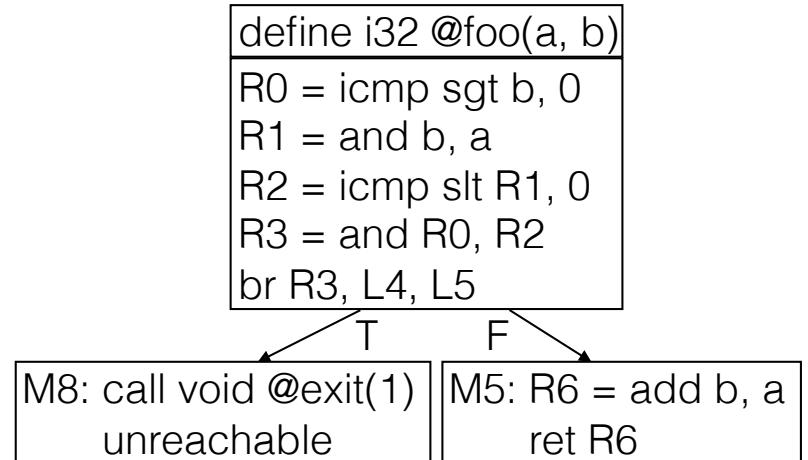


Reachability

P_{retro} : retrofitted program



P_{opt} : optimized P_{retro}



$P_{\text{retro}}: !(R2) \wedge !(R10) \wedge$
 $(R2 = (a > 0)) \wedge (R10 = (a < 0))$
 $E_{\text{retro}} = P_{\text{retro}} \vee \dots \vee P_i_{\text{retro}}$
 $E_{\text{opt}} = P_{\text{opt}} \vee \dots \vee P_j_{\text{opt}}$
 If P_{retro} reaches event of interest, then P_{opt} reaches event of interest ($E_{\text{retro}} \Rightarrow E_{\text{opt}}$)

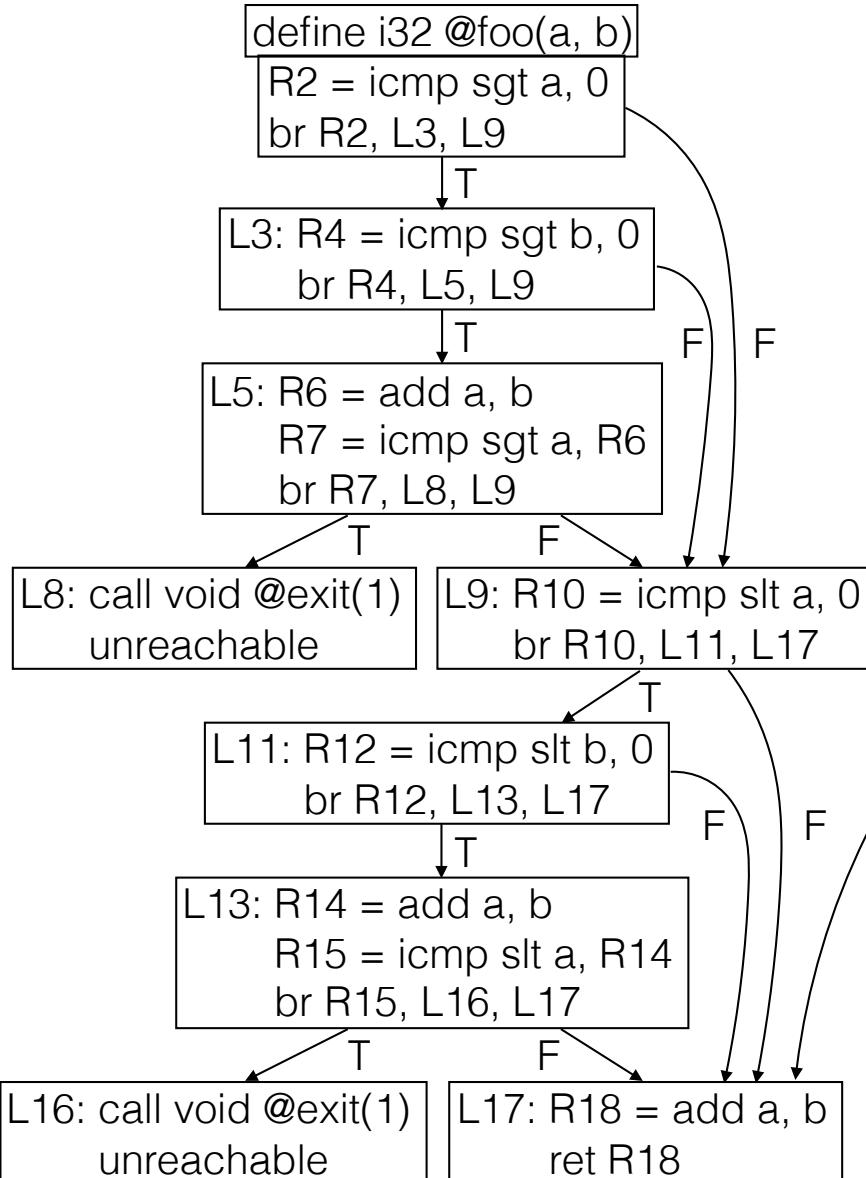
$\sim E_{\text{retro}}, \sim E_{\text{opt}}$

If P_{retro} does not reach event of interest, then P_{opt} does not reach event of interest ($\sim E_{\text{retro}} \Rightarrow \sim E_{\text{opt}}$)

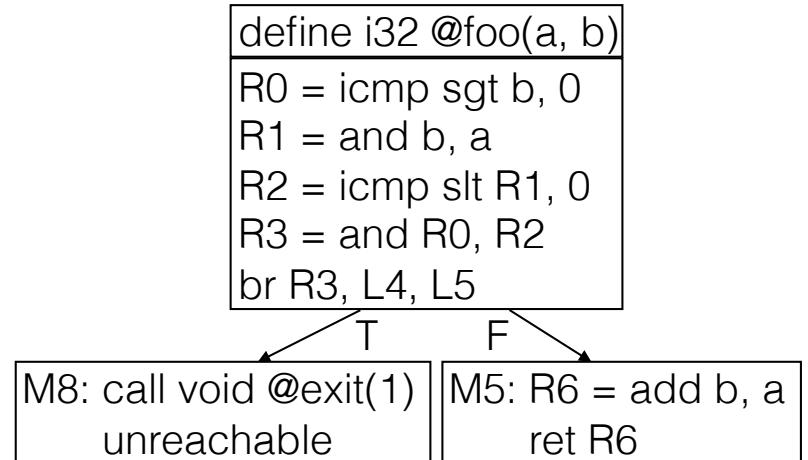


Reachability

P_{retro} : retrofitted program



P_{opt} : optimized P_{retro}



$P_{\text{retro}}: !(R2) \wedge !(R10) \wedge$
 $(R2 = (a > 0)) \wedge (R10 = (a < 0))$
 $E_{\text{retro}} = P_{\text{retro}} \vee \dots \vee P_i_{\text{retro}}$
 $E_{\text{opt}} = P_{\text{opt}} \vee \dots \vee P_j_{\text{opt}}$
 If P_{retro} reaches event of interest, then P_{opt} reaches event of interest ($E_{\text{retro}} \Rightarrow E_{\text{opt}}$)

$\sim E_{\text{retro}}, \sim E_{\text{opt}}$

If P_{opt} reaches event of interest, then
 P_{retro} reaches event of interest
 $(E_{\text{opt}} \Rightarrow E_{\text{retro}})$



RUTGERS

Initial Prototype

- Built prototype for LLVM IR programs.
- Modified retrofitting transformation to mark event of interests.
- Z3 for query.
- Naive integer overflow checker, Address Sanitizer, SoftboundCETS
- Discovered a bug in SoftBoundCETS

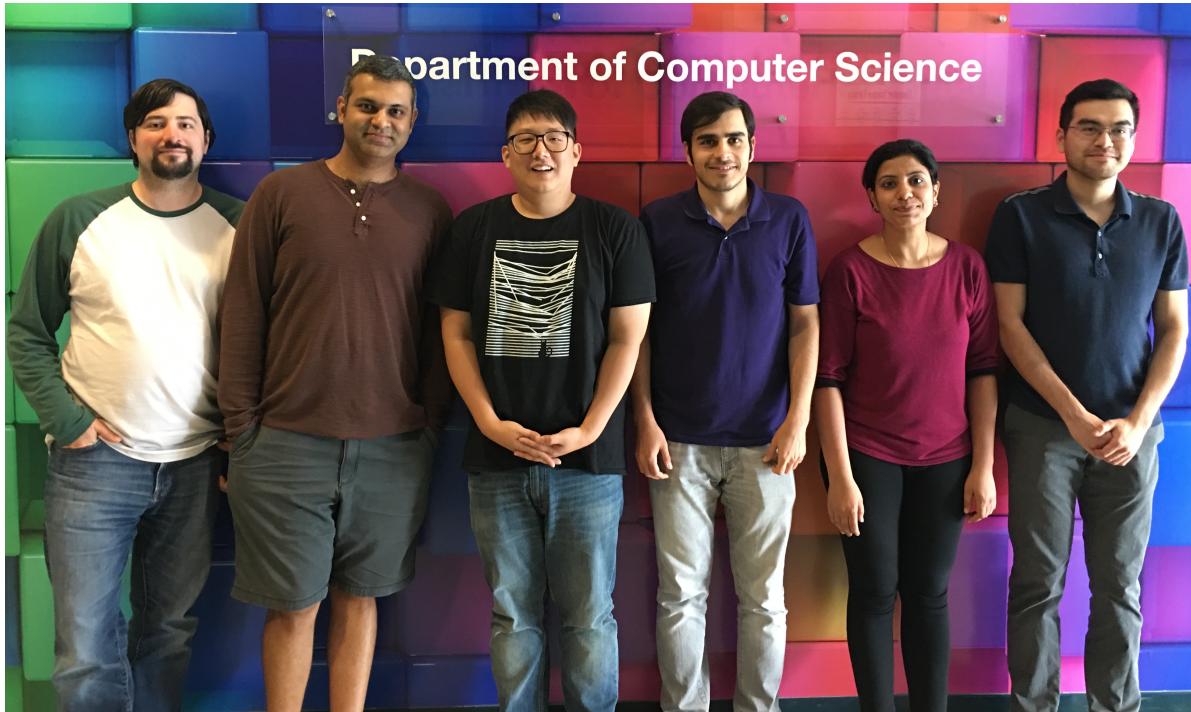


Summary

- Optimization do erroneously remove checks.
- Formulating as Reachability detects errors.
- Leverage translation validation approaches to identify cut-points
 - Memory equivalence axioms?
 - More precise encoding for loops?



Push formal methods into the tool chain of mainstream developers



Alive-NJ:

<https://github.com/rutgers-apl/alive-nj/>

TaskProf:

<https://github.com/rutgers-apl/TaskProf>

PTRacer:

<https://github.com/rutgers-apl/PTRacer>

DON'T STOP...



Other software prototypes from the Rutgers
Programming Languages Group:
<https://github.com/rutgers-apl/>



RUTGERS