# A Quest Toward the Perfect Optimizing Compiler

Theodoros Theodoridis

AST.ethz.ch
ADVANCED SOFTWARE TECHNOLOGIES

ETH zürich

Source



Compiler

Target

1010
1010

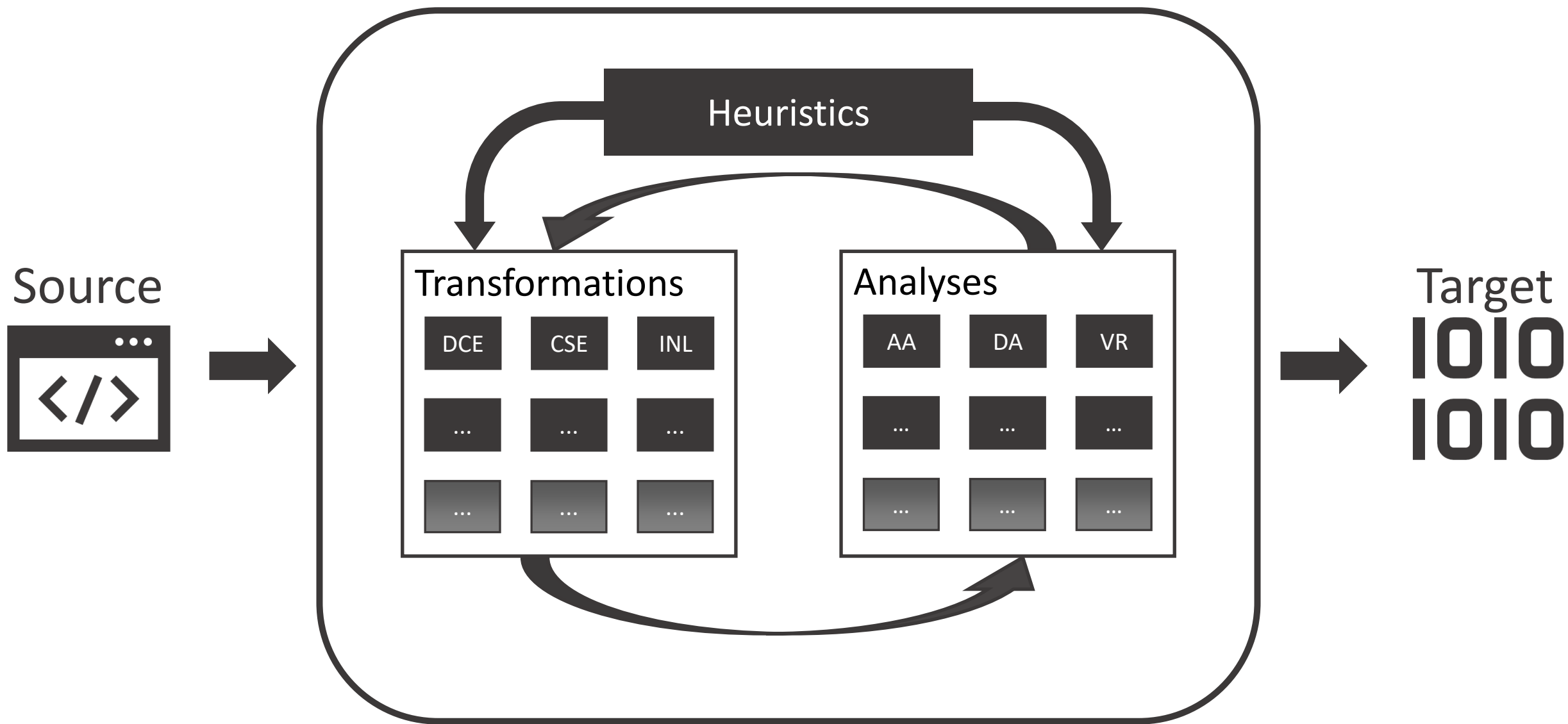Source → Heuristics [ Transformations (DCE, CSE, INL, ...) ← → Analyses (AA, DA, VR, ...) ] → Target

# How far are we from the optimum? Pretty far...



**BOLT: A Practical Binary Optimizer for Data Centers and Beyond**

Maksim Panchenko, Rafael Auler, Bill Nell, Guilherme Ottoni
Facebook, Inc.
Menlo Park, CA, USA
{maks,rafaelauler,bnell,ottoni}@fb.com

**Amir**
@disruptnhandlr · Follow

I've been working on CMake magic to fully automatically apply BOLT to Clang, similar to how PGO two-stage build is automated with CMake cache file. It appears that fully automatically BOLT-optimized Clang is ~30% faster than the Release (-O3) Clang: **reviews.llvm.org/D132975**

3:08 AM · Sep 2, 2022

Read the full conversation on Twitter

♡ 166    ◯ Reply    ⬆ Share

**1.3x Speedup**

timizations (FDO), also called *profile-guided optimizations* (PGO), particularly code layout. At the same time, due to their large sizes, applying FDO to these applications poses scalability challenges. Instrumentation-based profilers incur significant memory and computational performance costs, often making it impractical to gather accurate profiles from a production system. To simplify deployment and increase adoption, it is desirable to have a system that can obtain profile data for FDO from unmodified binaries running in their normal production environments. This is possible through the use of *sample-based profiling*, which enables high-quality profiles to be gathered with minimal operational complexity. This is the approach taken by tools such as Ispike [1], AutoFDO [2], and

code layout.

We demonstrate the finding above in the context of a static binary optimizer we built, called BOLT. BOLT is a modern, re-targetable binary optimizer built on top of the LLVM compiler infrastructure [8]. Our experimental evaluation on large real-world applications shows that BOLT can improve performance by up to 20.4% on top of FDO and LTO. Furthermore, our analysis demonstrates that this improvement is mostly due to the improved code layout that is enabled by the more accurate usage of sample-based profile data at the binary level.
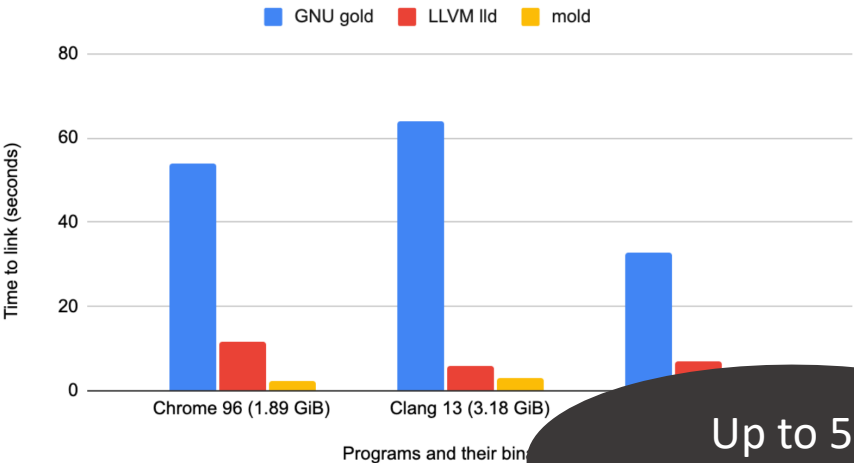Overall, this paper makes the following contributions:

---

☰ **README.md**

## mold: A Modern Linker

mold is a faster drop-in replacement for existing Unix linkers. It is several times faster than the LLVM lld linker, the second-fastest open-source linker which I originally created a few years ago. mold is designed to increase developer productivity by reducing build time, especially in rapid debug-edit-rebuild cycles.

Here is a performance comparison of GNU gold, LLVM lld, and mold for linking final debuginfo-enabled executables of major large programs on a simulated 8-core 16-threads machine.



**Up to 5x Speedup**

| Program (linker output size) | GNU gold | LLVM lld | mold |
|---|---|---|---|
| Chrome 96 (1.89 GiB) | 53.86s | 11.74s | 2.21s |
| Clang 13 (3.18 GiB) | 64.12s | 5.82s | 2.90s |
| Firefox 89 libxul (1.64 GiB) | 32.95s | 6.80s | 1.42s |

mold is so fast that it is only 2x *slower* than `cp` on the same machine. Feel free to file a bug if you find mold is not faster than other linkers.

# Our Approach



1. We obtain the optimum.

2. We compare with the compiler and find the gap.

# The Benefits of Inlining

```
int bar(int a, int b) {
   if ((a * b) % 2)
      return a + b;
   else
      return a - b;
}


int foo(int x) {
   return bar(x,2) + 2;
}
```
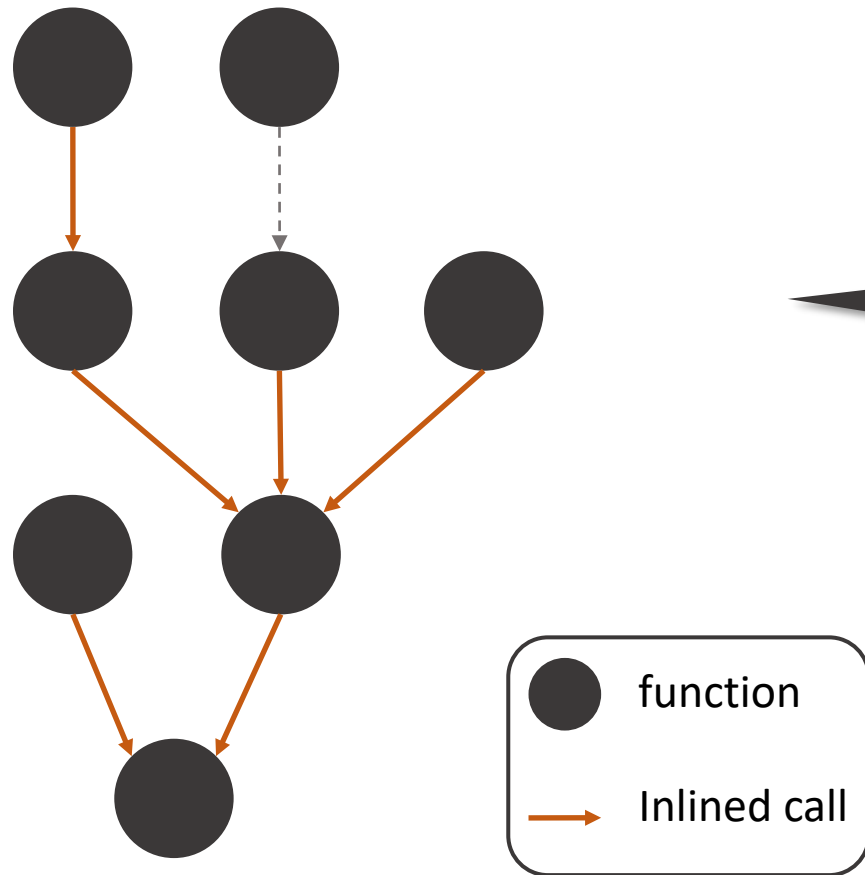
after inlining

```
int foo(int x) {
   return x;
}
```
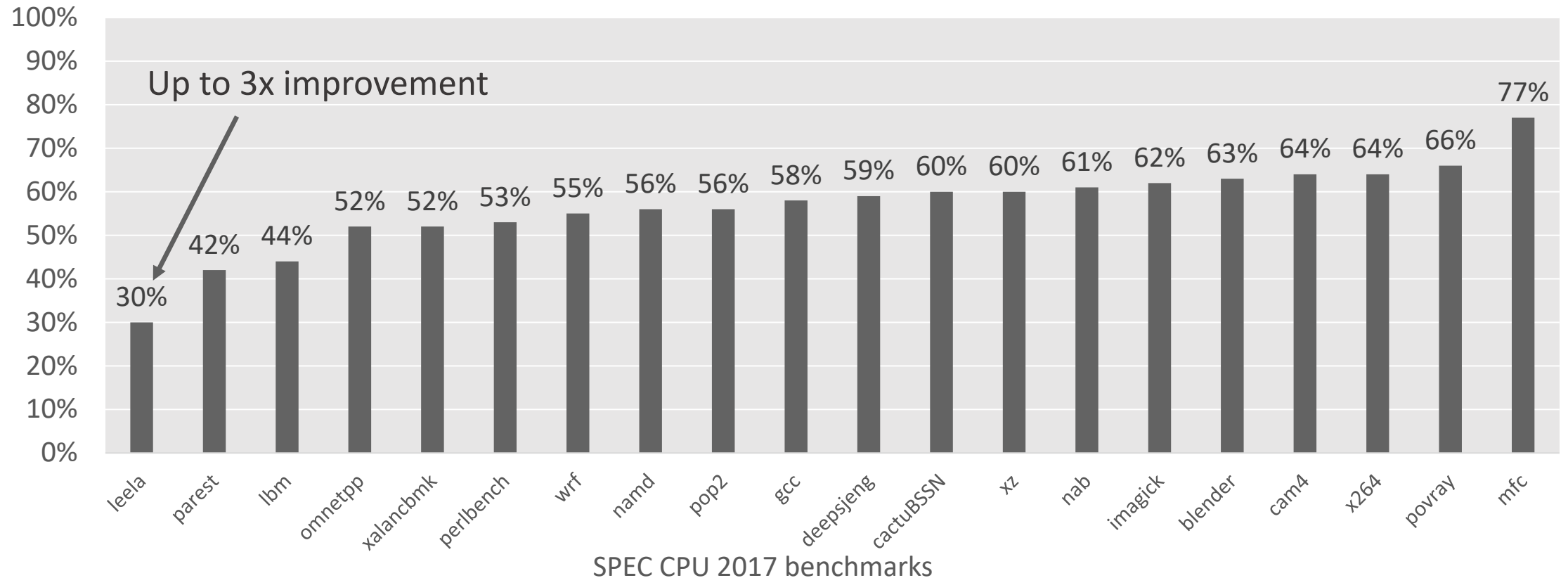
# Too much Inlining is Bad
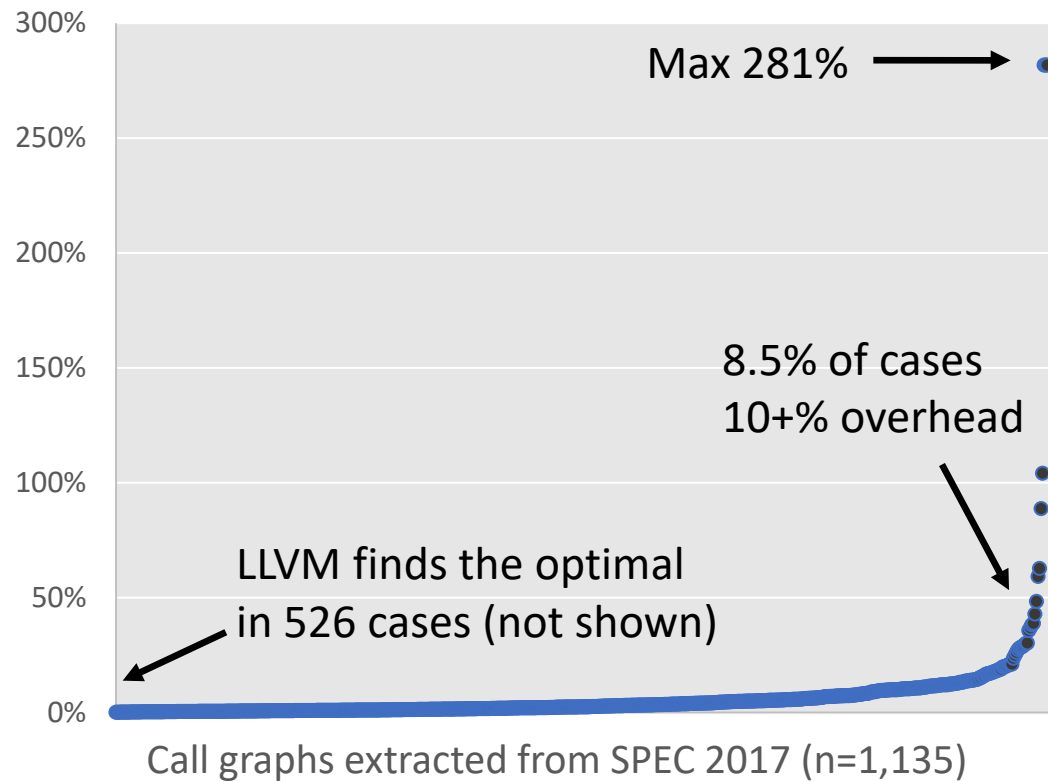


Aggressive Inlining:
69% binary size increase

function

Inlined call

# Proper Inlining Reduces Program Size

Relative size: clang -Os vs clang -Os -fno-inline



SPEC CPU 2017 benchmarks

# Gap between LLVM and Optimal

## Heuristic size overhead

Max 281% →

8.5% of cases
10+% overhead

LLVM finds the optimal
in 526 cases (not shown)

Call graphs extracted from SPEC 2017 (n=1,135)

## Common inlining choices

LLVM's heuristic
is too aggressive

None inline | Only LLVM inlines | Only optimal inlines | Both inline

```c
static int a = 0;

int main () {
  if (a != 0) {
    return 1;
  }
  a = 1;
  return 0;
}
```

```asm
main:
  xorl %eax, %eax
  retq
```

C source #1 ✕                                    ☐ ✕

A▾   🖫 Save/Load   ＋ Add new... ▾   ⅴ Vim          C                    ▾

```c
 3
 4
 5
 6
 7   static int a = 0;
 8
 9   int main () {
10     if (a != 0) {
11       return 1;
12     }
13     a = 1;
14     return 0;
15   }
16
17
18
19
20
21
```

clang 14.0.0 ✎ ✕                                 ☐ ✕

x86-64 clang 14.0.0        ▾   ✅   -O3                    ▾

A▾   ⚙ Output... ▾   ▼ Filter... ▾   ▤ Libraries   ＋ Add new... ▾   ✎ Add tool... ▾

```asm
1   main:                                    # @main
2           movl    $1, %eax
3           cmpb    $0, a(%rip)
4           jne     .LBB0_2
5           movb    $1, a(%rip)
6           xorl    %eax, %eax
7   .LBB0_2:
8           retq
```

↻ ▤ Output (0/0)   x86-64 clang 14.0.0  ℹ - cached (7549B) ~158 lines filtered 📊

gcc 11.3 ✎ ✕                                     ☐ ✕

x86-64 gcc 11.3        ▾   ✅   -O3                        ▾

A▾   ⚙ Output... ▾   ▼ Filter... ▾   ▤ Libraries   ＋ Add new... ▾   ✎ Add tool... ▾

```asm
1   main:
2           movl    a(%rip), %eax
3           testl   %eax, %eax
4           jne     .L3
5           movl    $1, a(%rip)
6           ret
7   .L3:
8           movl    $1, %eax
9           ret
```
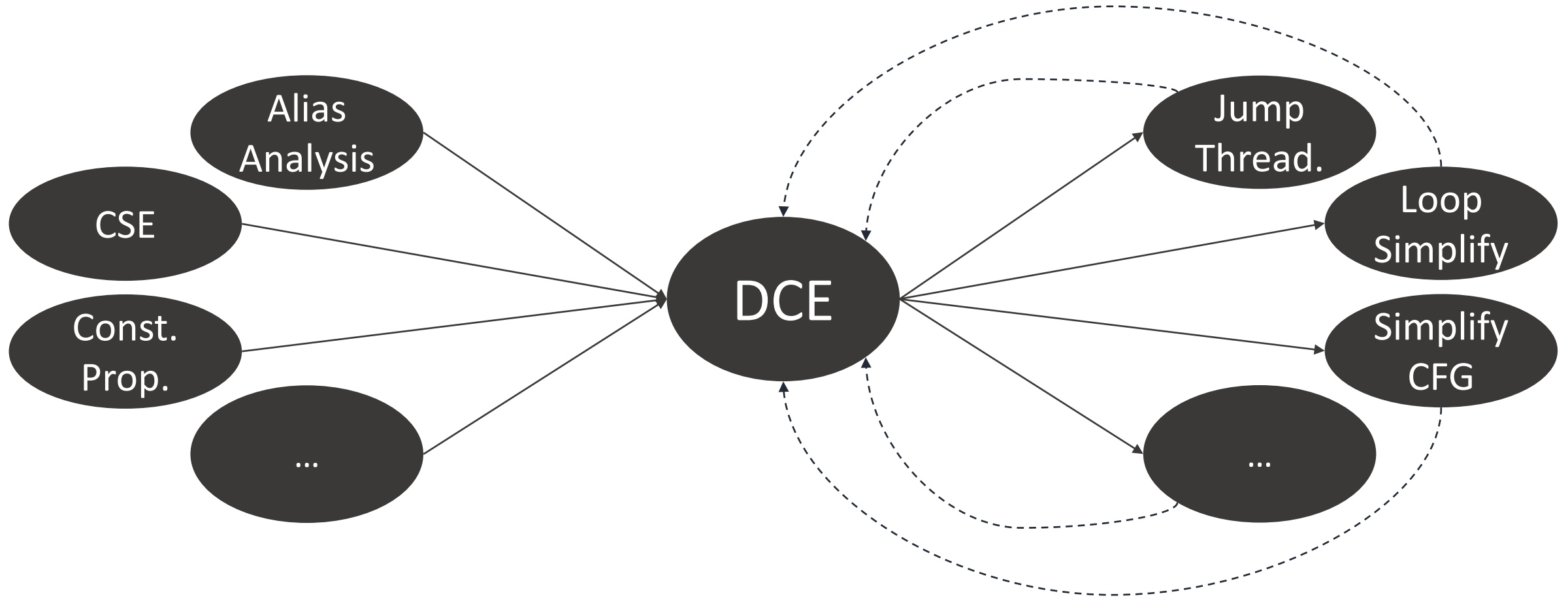14

↻ ▤ Output (0/0)   x86-64 gcc 11.3  ℹ - cached (3047B) ~182 lines filtered 📊

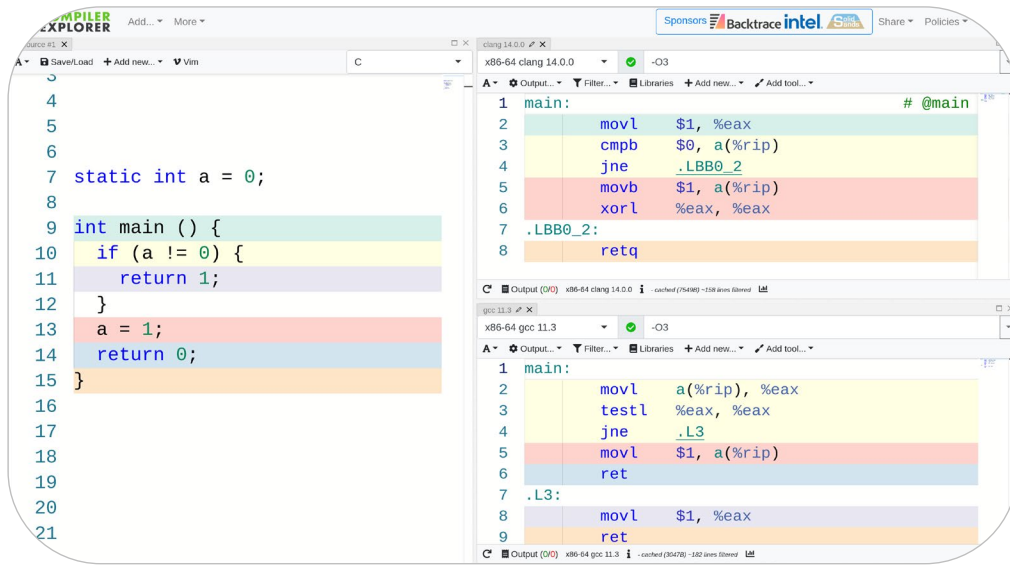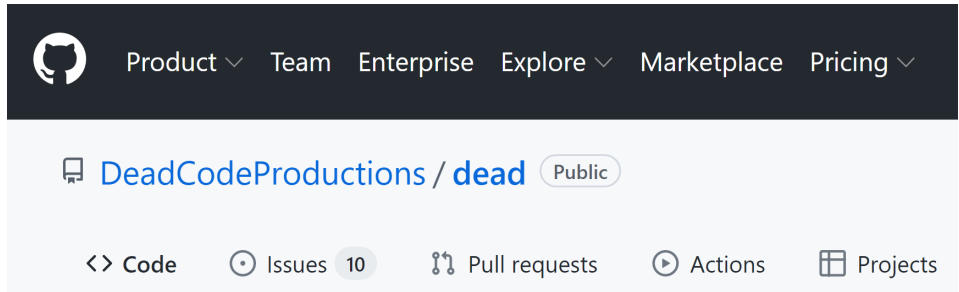# Dead Code Elimination: An Optimization Sink

# How good are compilers at DCE?

Corpus of 10,000 test programs:

- Generated with Csmith

- 3,109,167 dead blocks

| Optimization Level | % of dead blocks that are missed | |
|---|---|---|
| | GCC | LLVM |
| O0 | 85.2% | 83.2% |
| O1 | 8.2% | 5.2% |
| Os | 6.0% | 4.8% |
| O2 | 5.7% | 4.4% |
| O3 | 5.6% | 4.3% |

# Finding Missed Optimization Opportunities Automatically
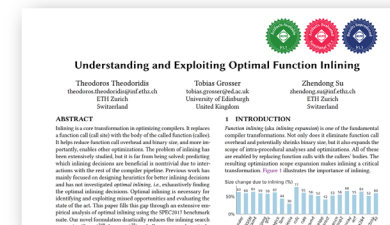




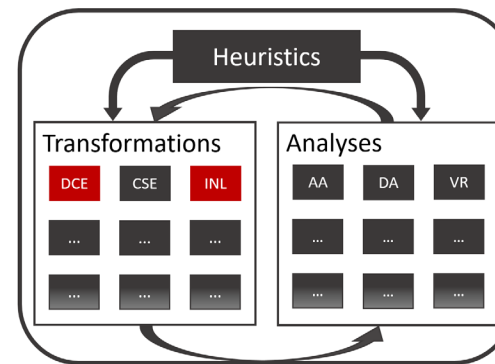| | LLVM | GCC |
|---|---|---|
| Reported | 47 | 55 |
| Confirmed | 35 | 46 |
| Fixed | 15 | 15 |

Our Approach

1. We obtain the optimum.

2. We compare with the compiler and find the gap.
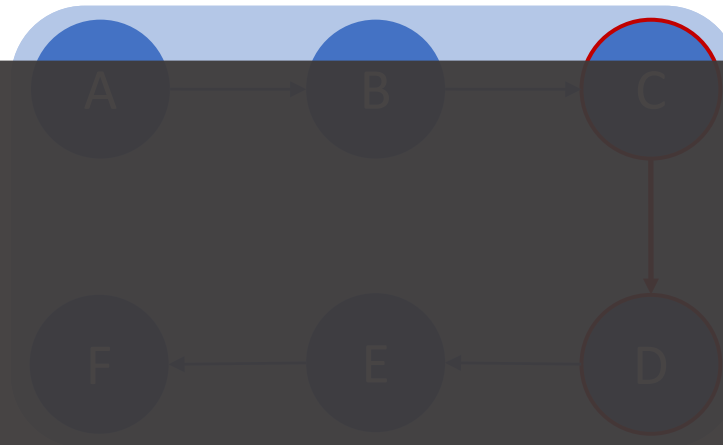
"What if we had optimal ...?"

Ongoing Work

• Optimal Alias Analysis Information

• Optimal Pass Pipelines

• Learning Heuristics based on Optimal Inlining Choices

# Backup Slides

# This can be done recursively!

C → D inlined

C → D not inlined

$2^2$ combinations

$2^2$ combinations

$2^4$ combinations

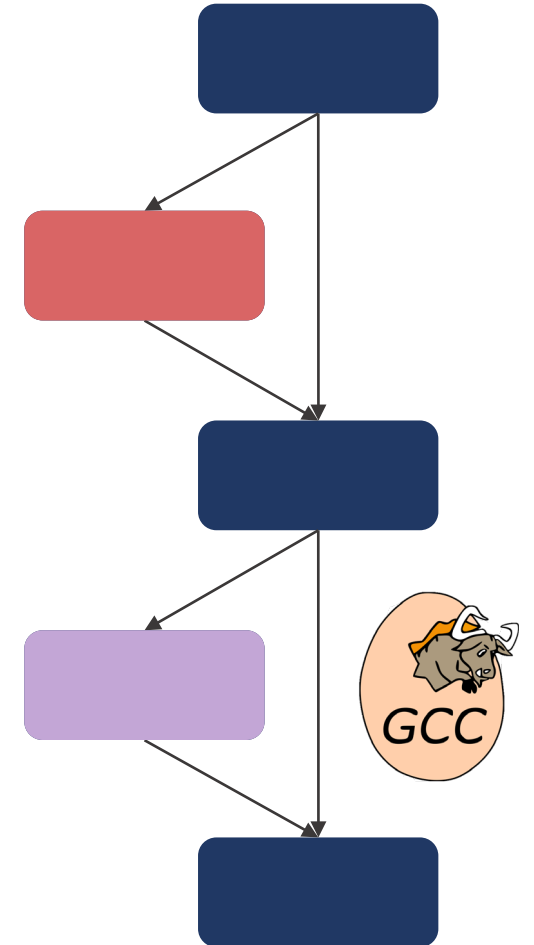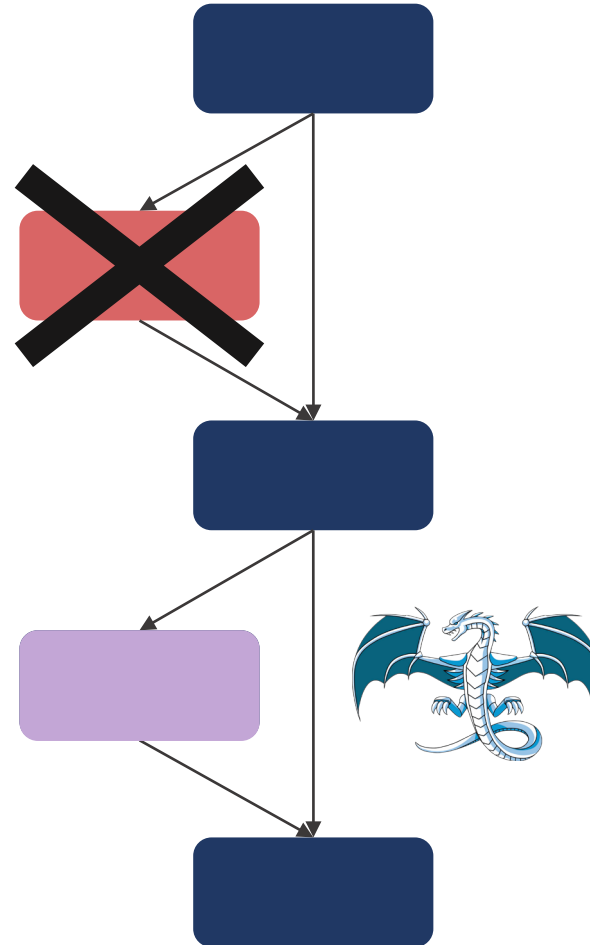Total:$(2^2 + 2^2 + 1) + 2^4 = 25 < 32$ naïve

# Lens of Dead Code Elimination

# Different Compilers Eliminate Different Parts

```
int a = 0;
static int b[2] = {0,0}, c = 0;

int main() {
  if (b[a]) {
    return 1;
  }
  if (c) {
    return 2;
  }
  c = 1;
  return 0;
}
```

# Missed Dead Code Elimination Detection

```
int a = 0;
static int b[2] = {0,0}, c = 0;

int main() {
  if (b[a]) {
    return 1;
  }
  if (c) {
    return 2;
  }
  c = 1;
  return 0;
}
```

```
main:
  movl   $2, %eax
  cmpb   $0, c(%rip)
  jne    .LBB0_2
  movb   $1, c(%rip)
  xorl   %eax, %eax
.LBB0_2:
  retq
```

```
main:
  movslq a(%rip), %rdx
  movl   $1, %eax
  movl   b(,%rdx,4),%edx
  testl  %edx, %edx
  jne    .L1
  movl   c(%rip), %eax
  testl  %eax, %eax
  jne    .L4
  movl   $1, c(%rip)
  ret
.L4:
  movl   $2, %eax
.L1:
  ret
```

*GCC*

# Missed Dead Code Elimination: Markers

```c
int a = 0;
static int b[2] = {0,0}, c = 0;

int main() {
 if (b[a]) {
   return 1;
 }
 if (c) {
   return 2;
 }
 c = 1;
 return 0;
}
```
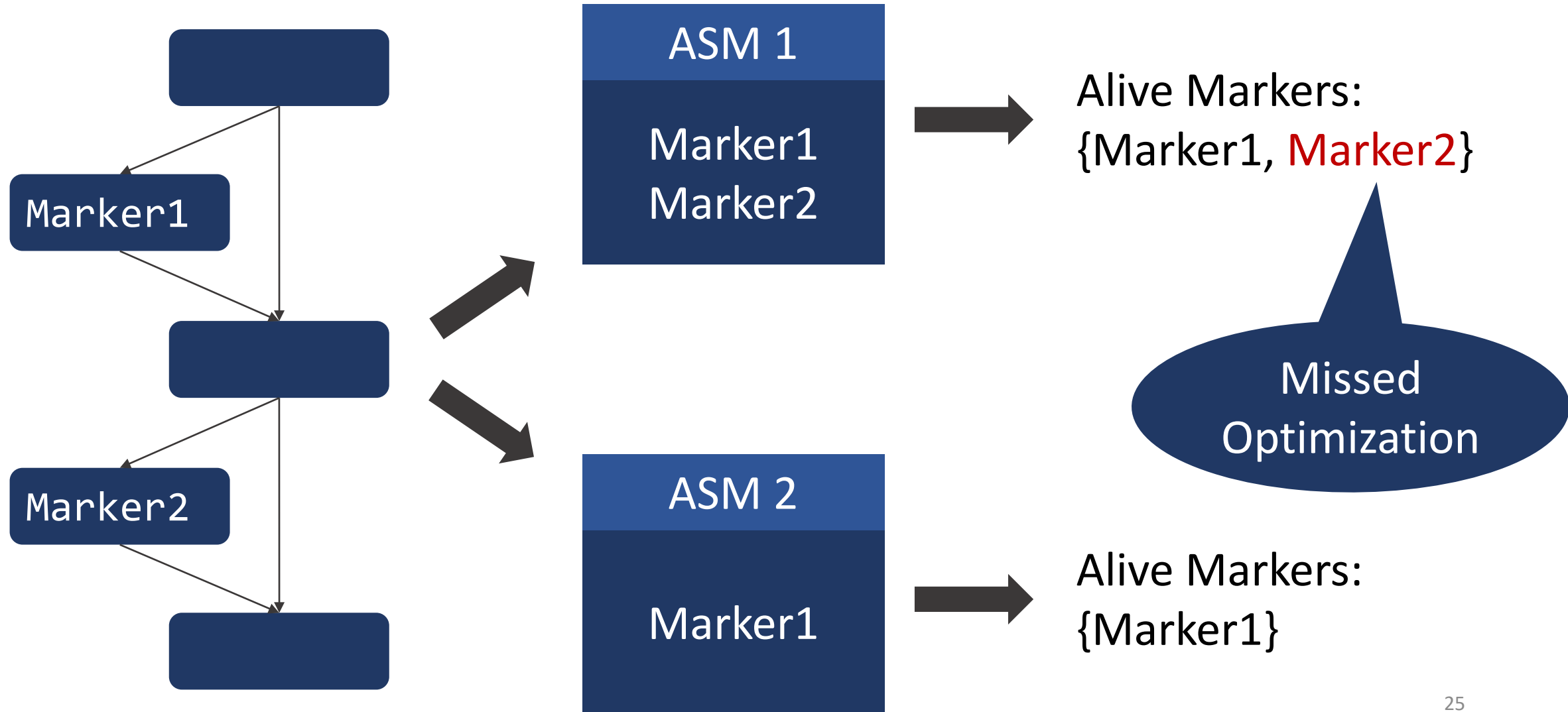
```asm
main:
 pushq    %rax
 cmpb     $1, c(%rip)
 jne      .LBB0_2
 callq    DCEMarker2
 movl     $2, %eax
 popq     %rcx
 retq
.LBB0_2:
 movb     $1, c(%rip)
 xorl     %eax, %eax
 popq     %rcx
 retq
```

```asm
main:
 subq     $8, %rsp
 movslq a(%rip), %rax
 movl     b(,%rax,4),%eax
 testl   %eax, %eax
 jne      .L7
 …
.L7:
 call     DCEMarker1
 movl     $1, %eax
 jmp      .L1
.L8:
 call     DCEMarker2
 movl     $2, %eax
 jmp      .L1
```

# The Lens of Dead Code Elimination

# DCE Examples

```
static long a = 78240;
static int b, d;
static short e;
static short c(short f, short h) {
  return h == 0 ||
    (f && h == 1) ? 0 : f % h; }
int main() {
  short g = a;
  for (b = 0; b < 1; b++) {
    e = a;
    d = c((e == a) ^ g, a);
  }

  if (d) {
    DCEMarker();
    for (; a; a++);
  }
}
```

LLVM 13 -O1

Modulo on constant ranges: [X,X+1) % [X,X+1) not simplified

LLVM 13 -O3

```
static int b = -1, e = 1;
static short c = 0, d = 0;
short a(unsigned short f, int g) {
  return f >> g;
}

int main() {
  c++;
  d  = a(4294967295 + (c > 0),1);
  e ^= (short)(d * 3) /(unsigned)b;
  if (!e)
    DCEMarker();
}
```
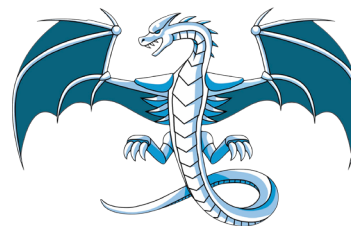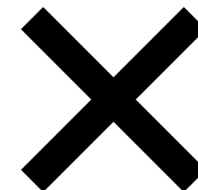
Regression on shift peephole optimization

e != 0

LLVM 13 -O3

LLVM dev -O3

[SimplifyCFG] don't sink common insts too soon (PR34603)

This should solve:
https://bugs.llvm.org/show_bug.cgi?id=34603
...by preventing SimplifyCFG from altering redundant instructions  before early-cse has a chance to run.