# Merging Similar Control-Flow Regions in LLVM for Performance and Code Size Benefits

**Charitha Saumya (Presenter),** Kirshanthan Sundararajah, Milind Kulkarni

# Similar Code inside Conditional Branches is Plentiful

```
if ( fixed_mult_quo(bx, ay, by) < ax ) {
  left.end.x = px + bx, left.end.y = py + by;
  right.end.x = px + ax, right.end.y = py + ay;
  code = (*fill_trapezoid)(dev, &left, &right, py, ym, false, pdevc, lop);
  right.start = right.end;
  right.end = left.end;
} else {
  left.end.x = px + ax, left.end.y = py + ay;
  right.end.x = px + bx, right.end.y = py + by;
  code = (*fill_trapezoid)(dev, &left, &right, py, ym, false, pdevc, lop);
  left.start = left.end;
  left.end = right.end;
}
```

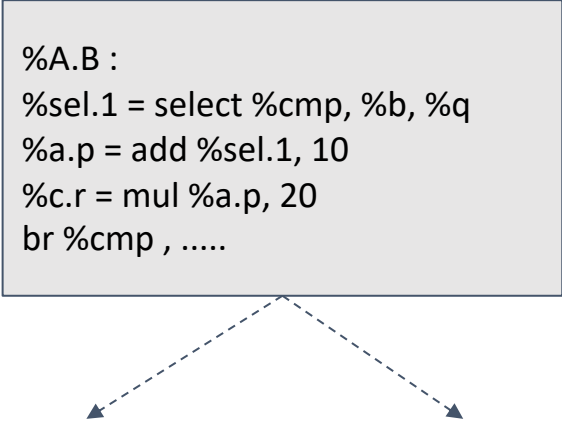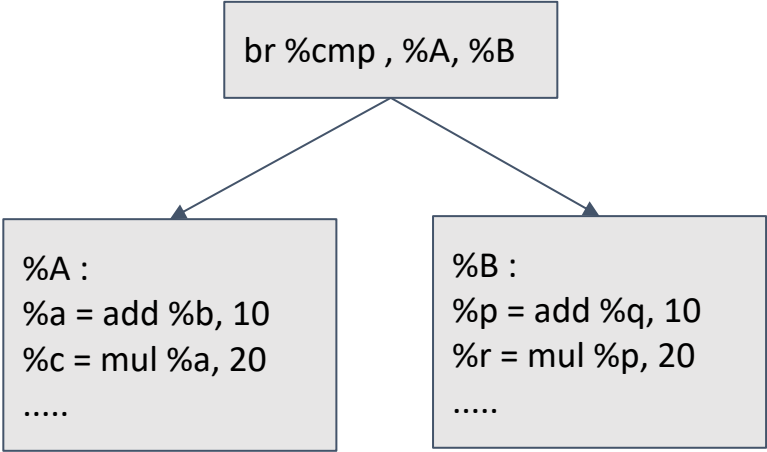Source : `mibench/ghostscript` benchmark

```
if ((tid & k) == 0)
{
    if (shared[tid] > shared[ixj])
        swap(shared[tid], shared[ixj]);
}
else
{
    if (shared[tid] < shared[ixj])
        swap(shared[tid], shared[ixj]);
}
```
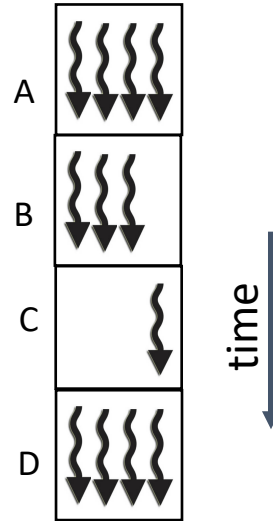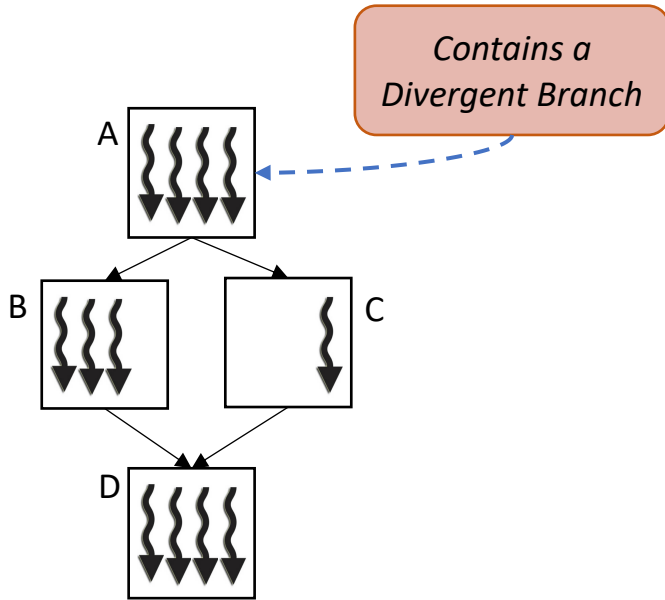
Source : *CUDA* implementation of bitonic sort

*Identical/similar operation sequences inside if and else sections*

*Similar control–flow regions inside if and else sections*

# Code Size Reduction

br %cmp , %A, %B

%A :
%a = add %b, 10
%c = mul %a, 20
.....

%B :
%p = add %q, 10
%r = mul %p, 20
.....

%A.B :
%sel.1 = select %cmp, %b, %q
%a.p = add %sel.1, 10
%c.r = mul %a.p, 20
br %cmp , .....

Size reduction =  cost(add) + cost(mul) - cost(select)

# Reducing Control-flow Divergence



Contains a
Divergent Branch

- Thread group (warp/wavefront)

time

# Reducing Control-flow Divergence



%B:
  ...
  %a = load ..
  %a = mul ....
  %b = div ....
  ....

%C:
  ...
  %p = load ..
  %q = mul ....
  %r = div ....
  ....

A
B
C
D

time

# Reducing Control-flow Divergence



%B:
...
%a = load ..
%a = mul ....
%b = div ....
....

%C:
...
%p = load ..
%q = mul ....
%r = div ....
....

A
B
C
D

time

A
B'
C'
BC
B''
C''
D

*Convergent execution of common instructions*

# Code Sinking and Code Hoisting



*Code Sinking*

**Identical instruction sequences are moved to common successor**

# Code Sinking and Code Hoisting

```
....
br %c, label %A, label %B
```

```
%A :
add ....
store ...
....
```

```
%B :
add ....
store ....
....
```

*Code Hoisting*

```
select ..
select ..
add ..
select ..
select ..
store ..
br %c, label %A, label %B
```

```
%A:
....
```

```
%B:
.....
```

**Identical instruction sequences are moved to common predecessor**

# Branch Fusion



(a)

$l_1$ load($t_0$, tid)
$p_0 = t_0 \neq 0.0$
branch($p_0$, $l_{13}$)

T

$l_4$ load($t_1$, tid)
$t_2 = t_1 * t_1$
$t_3 = t_2 * t_1$
$t_4 = t_3 * 3.14$
$t_5 = t_4 / t_0$
$t_6 = t_5 / t_0$
$t_7 = t_1 * 2.71$
$t_8 = t_6 + t_7$
store($t_8$, tid)

F

$l_{13}$ load($t_9$, tid)
$t_{10} = t_9$
$t_{11} = t_{10}$
$t_{12} = t_{11}$
$t_{13} = t_9$
$t_{14} = t_{13}$
$t_{15} = t_{12}$
store($t_{15}$, tid)
jump($l_{22}$)

$l_{22}$ sync
stop

*Branch Fusion using*

*These techniques are limited to exploiting similarity at basic block level (e.g., diamond shaped control-flow)*

(c)

$l_1$ load($t_0$, tid)
$p_0 = t_0 \neq 0.0$
load($t_{1\text{-}9}$, tid)
$t_{2\text{-}10} = t_{1\text{-}9} * t_{1\text{-}9}$
$s_1 = \text{sel}(p_0, t_{2\text{-}10}, 3.14)$
$t_{3\text{-}11} = t_{2\text{-}10} * s_1$
branch($p_0$, $l_{10}$)

$l_8$ $t_4 = t_{3\text{-}11} * 3.14$

$\phi(t_5, t_{3\text{-}11})$
$\text{sel}(p_0, t_0, 2.0)$
$s_{5\text{-}3\text{-}11} / s_3$
$t_{1\text{-}9} * 2.71$
ch($p_0$, $l_{16}$)

$l_{15}$ $t_{14} = t_{7\text{-}13} * t_{1\text{-}9}$

$l_{16}$ $t_{7\text{-}14} = \phi(t_{7\text{-}13}, t_{14})$
sync
$t_{8\text{-}15} = t_{6\text{-}12} + t_{7\text{-}14}$
store($t_{8\text{-}15}$, tid)

**Convergent execution of common instructions**

9

# Control-Flow Melding (CFM)



Control–Flow Melding
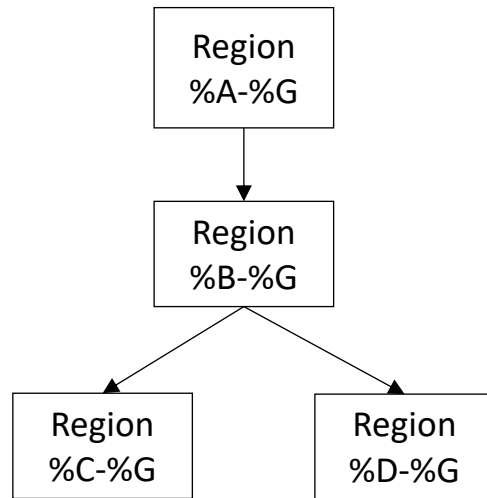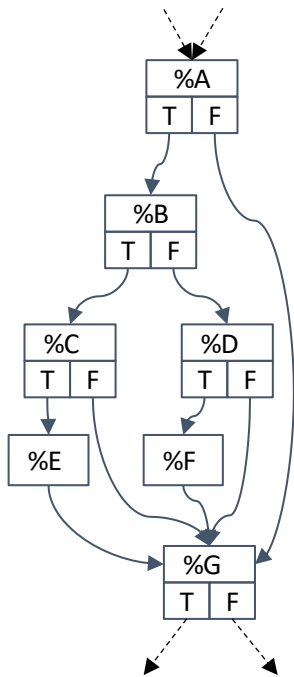
*Control-flow Regions with Similar computations*
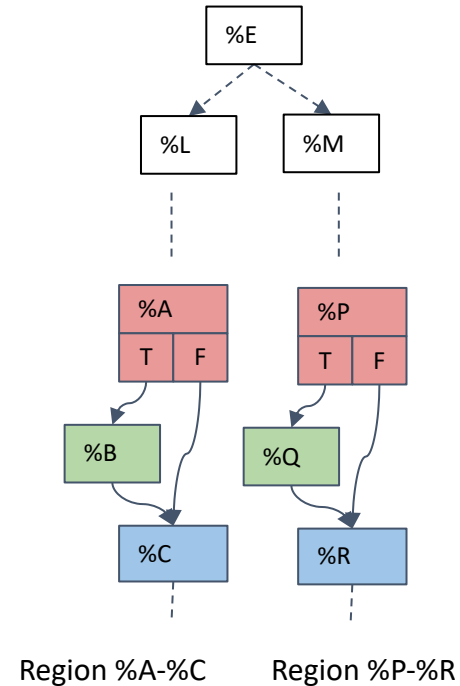
*Melded control-flow*

# Program Structure Tree

- Represents all Single-Entry Single-Exit (SESE) regions in a CFG
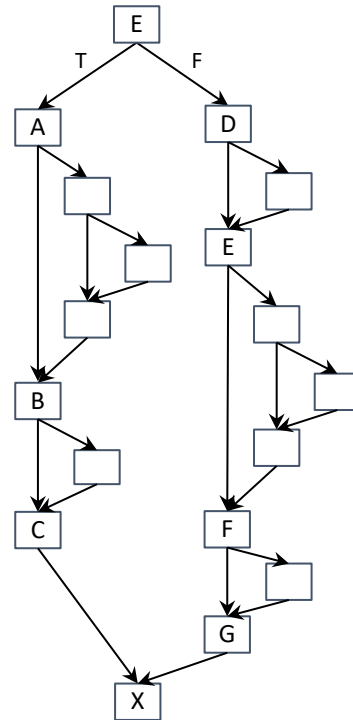- Can be obtained using the **RegionInfo** interface in LLVM

# Meldable Regions

- Two SESE regions can be **melded** if,

  - Dominated by a conditional branch

  - No path exists that goes through both the SESE regions

  - Entry blocks of the regions must post-dominate either the left or right successor of the conditional branch
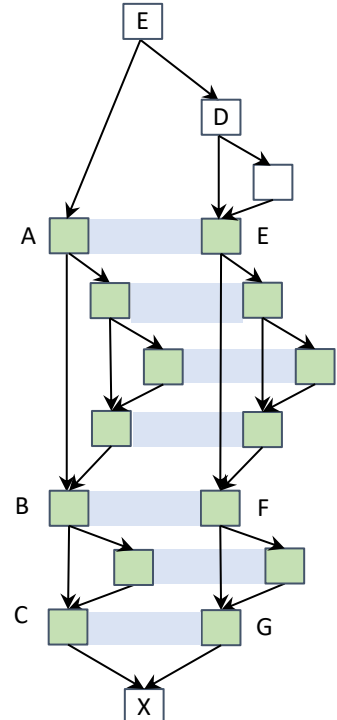
  - They are **isomorphic** (have same control-flow signature)



Region %A-%C    Region %P-%R

# Region Alignment

- Multiple isomorphic regions in if and else paths?

- Regions are aligned based on **Melding Profitability**

- **Melding Profitability** : metric that measures the similarity of two regions base on instruction frequencies
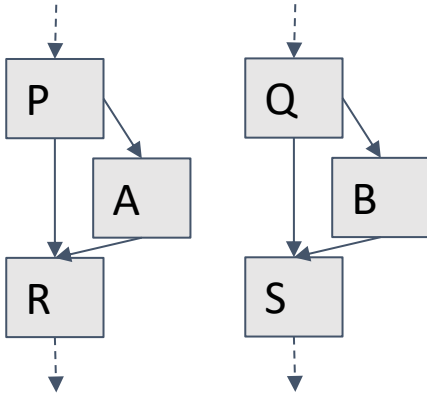


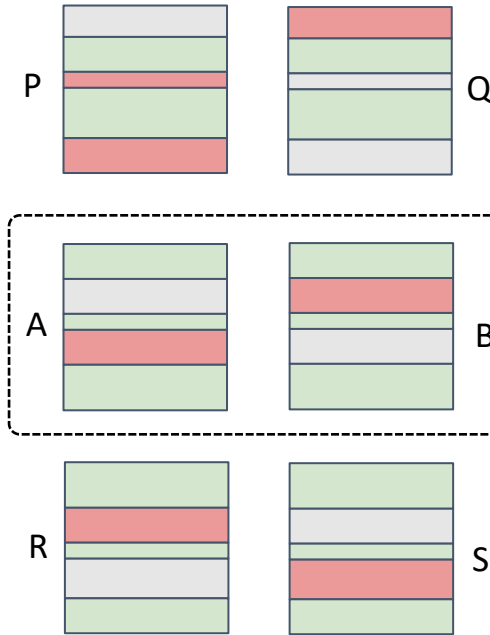Left regions : A-B, B-C
Right Regions : D-E, E-F, F-G
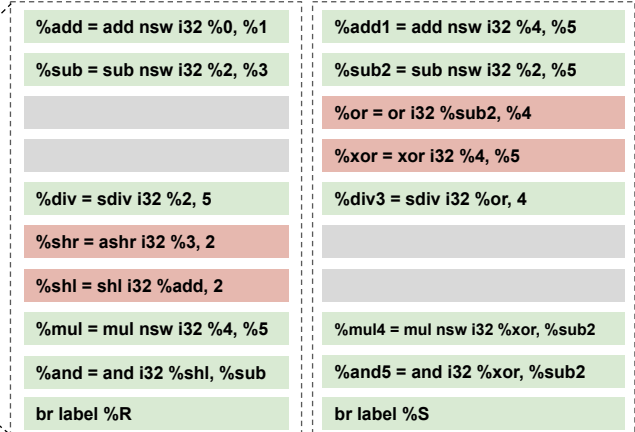
Region Alignment :
A-B with E-F
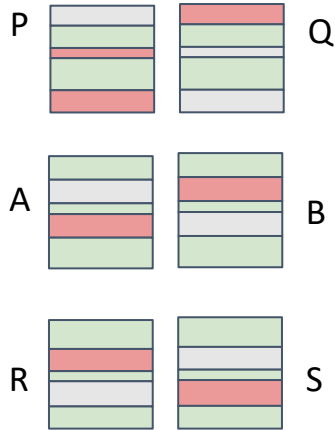B-C with F-G

# Instruction Alignment
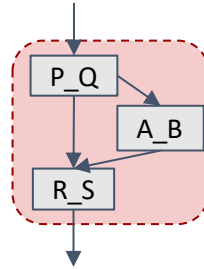


Aligned region pair

Instruction alignment for A and B computed using **instruction compatibility** and **cost**

```
%add = add nsw i32 %0, %1          %add1 = add nsw i32 %4, %5
%sub = sub nsw i32 %2, %3          %sub2 = sub nsw i32 %2, %5
                                   %or = or i32 %sub2, %4
                                   %xor = xor i32 %4, %5
%div = sdiv i32 %2, 5              %div3 = sdiv i32 %or, 4
%shr = ashr i32 %3, 2
%shl = shl i32 %add, 2
%mul = mul nsw i32 %4, %5          %mul4 = mul nsw i32 %xor, %sub2
%and = and i32 %shl, %sub          %and5 = and i32 %xor, %sub2
br label %R                        br label %S
```

Aligned          Unaligned

# Code Generation



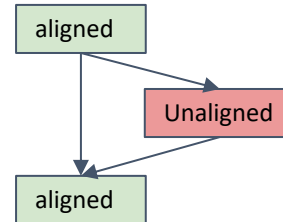Generated melded control-flow



Generated melded instructions

Aligned instruction pair

```
%add = add nsw i32 %0, %1          %add1 = add nsw i32 %4, %5
```
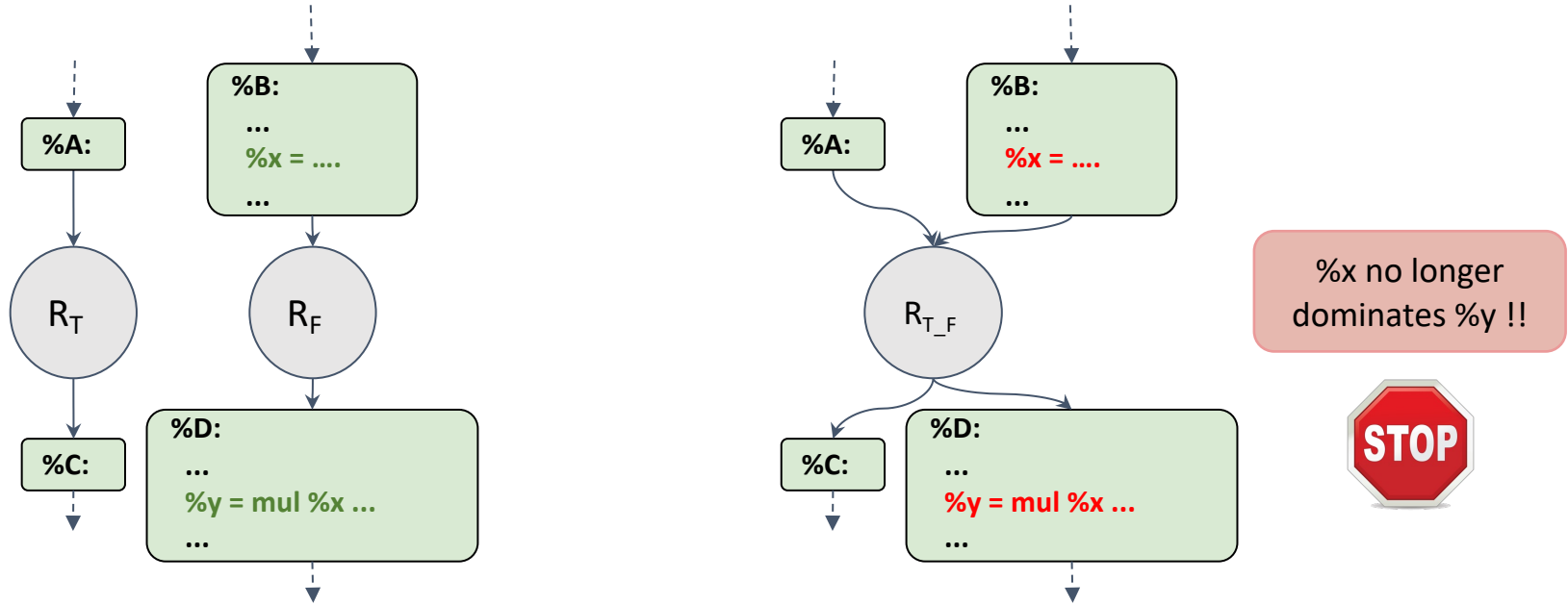
Generated code

```
%sel1 = select i1 %cmp, i32 %0, i32 %4
%sel2 = select i1 %cmp, i32 %1, i32 %5
 %6 = add nsw i32 %sel1, %sel2
```



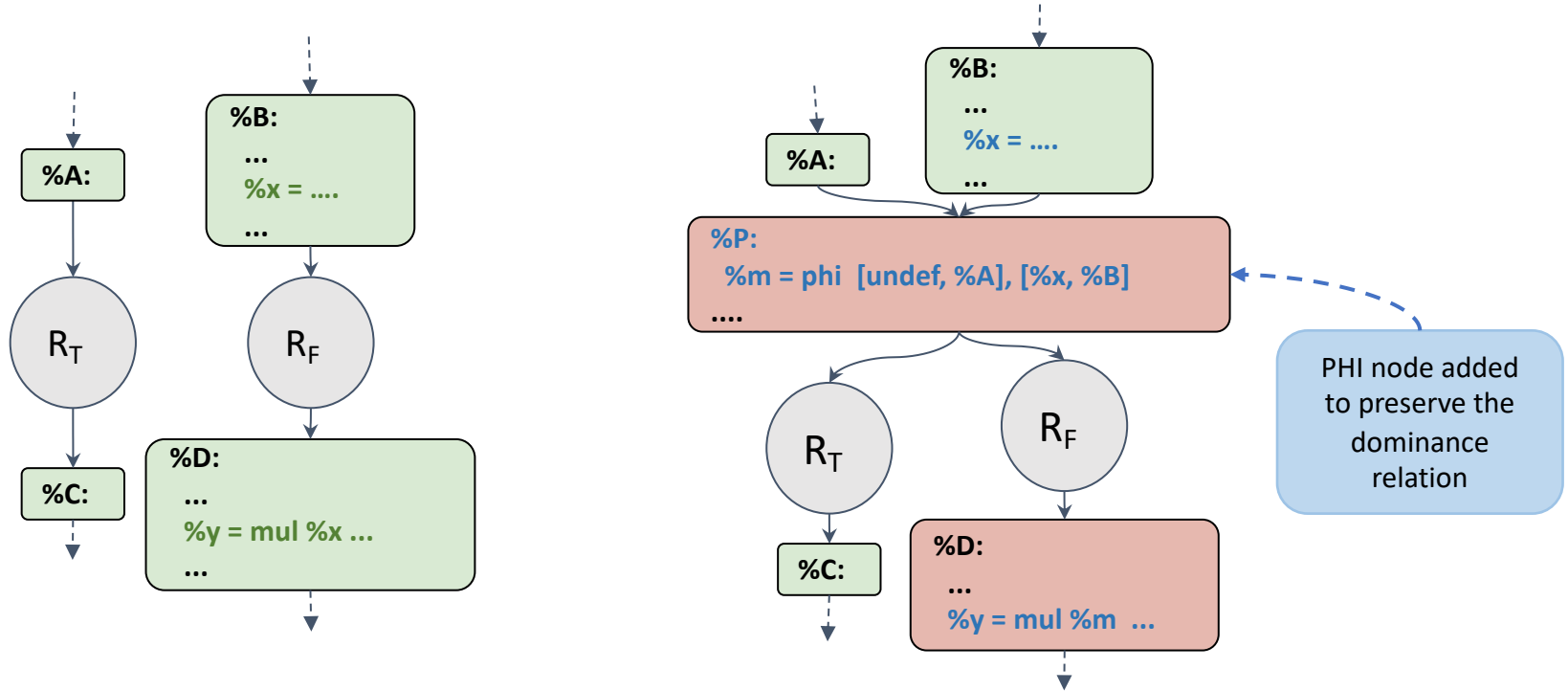Unaligned instructions are executed conditionally

# Ensuring Correctness

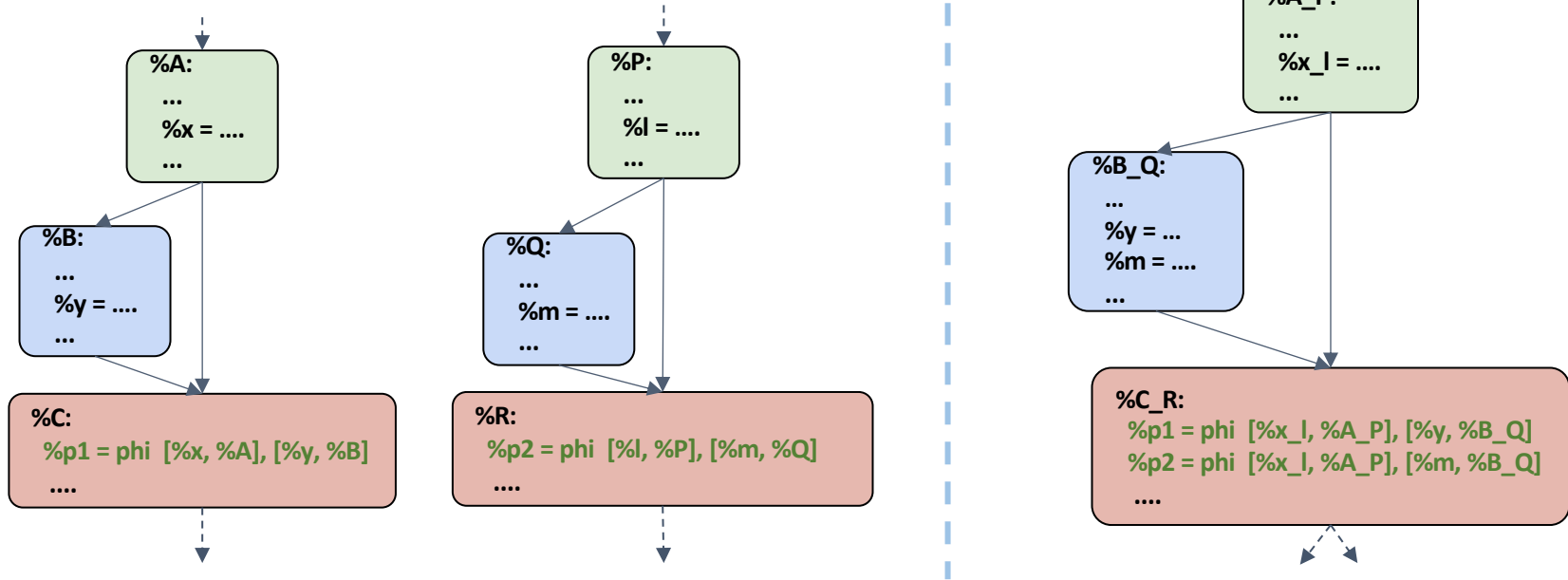- Melding can break the def-use chains outside the melded regions



%x no longer dominates %y !!

# Ensuring Correctness

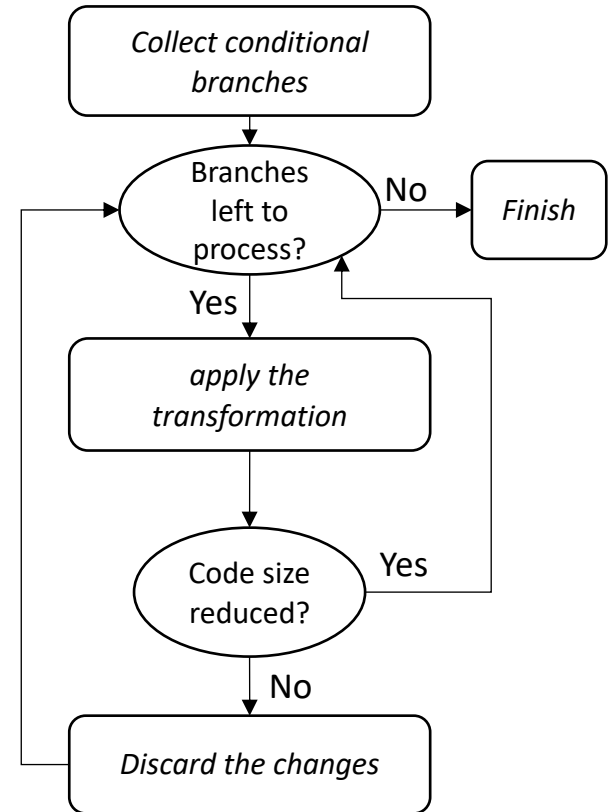- Melding can break the def-use chains outside the melded regions

# Ensuring Correctness

- Handling PHI instructions

# Implementation (CPU Code Size)

- LLVM-IR transformation pass

- Process all if-then-else branches

- Uses **TargetTransformInfo** interface
  - For instruction alignment
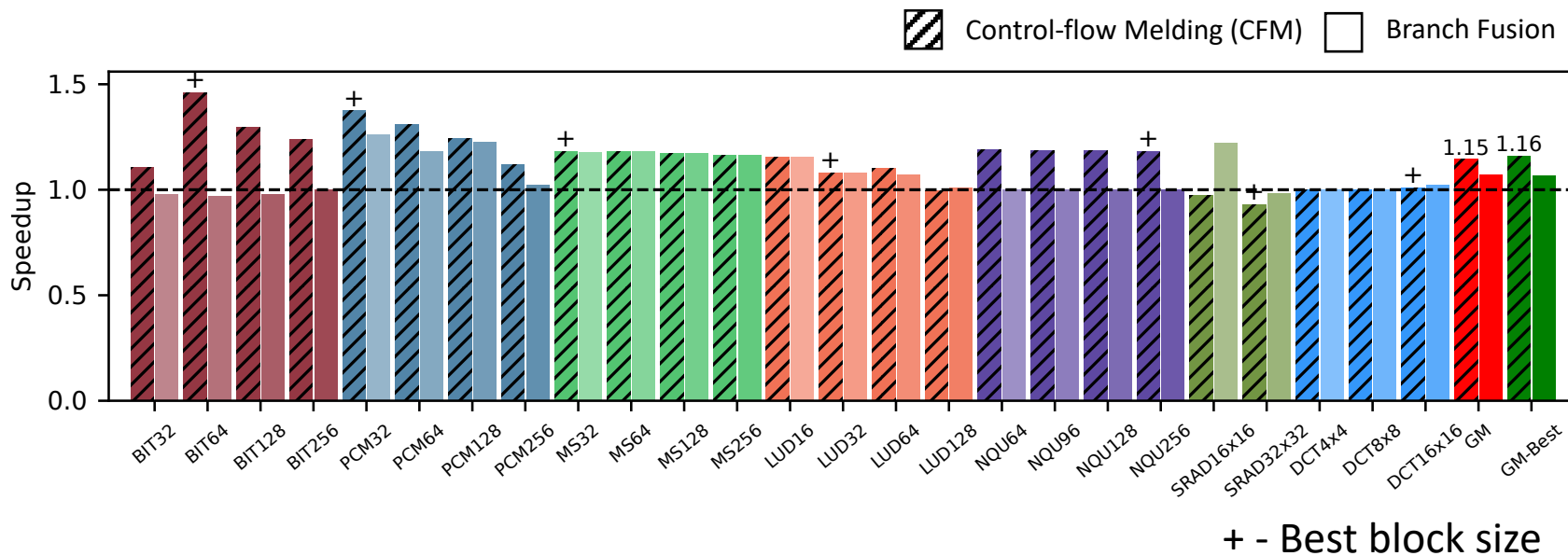  - For evaluating the code size before and after

```
Collect conditional
branches
        │
        ▼
   Branches          No
   left to  ──────────────▶  Finish
   process?
        │                      ▲
        │ Yes                  │
        ▼                      │
   apply the                   │
   transformation              │
        │                      │
        ▼                      │
   Code size        Yes        │
   reduced?  ──────────────────┘
        │
        │ No
        ▼
   Discard the changes
```

# Implementation (GPU Divergence Reduction)

- Process only divergent if-the-else branches
  - Uses LLVM ***DivergenceAnalysis***

- Uses a custom instruction latency cost model to evaluate the profitability

# Evaluation (GPU Divergence Reduction)

● Speedups on AMD Radeon Pro Vega 20 GPU + AMD Ryzen CPU

● Comparison against both *-O3 (full opts)* and *Branch Fusion with full opts*

● Benchmarks

1. Bitonic Sort (BIT)
2. Parallel and Concurrent Merge (PCM)
3. Mergesort (MS)
4. LU-decomposition (LUD)
5. N-Queens (NQU)
6. Speckle Reducing Anisotropic Diffusion (SRAD)
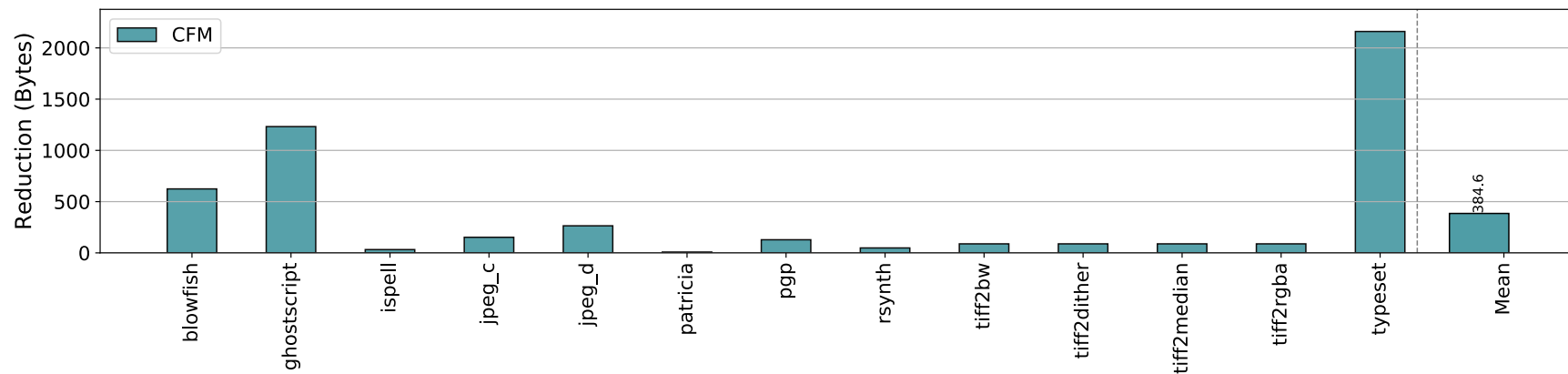7. DCT Quantization (DCT)

# Performance



+ - Best block size

**BIT, PCM, NQU :** *Performance improvement over branch fusion comes from CFM's ability to meld at region level*
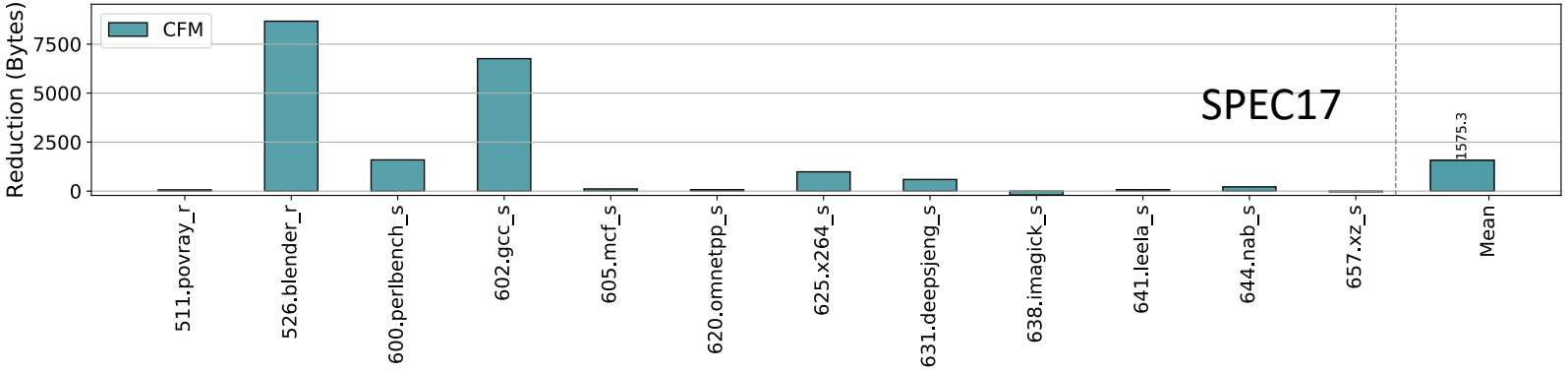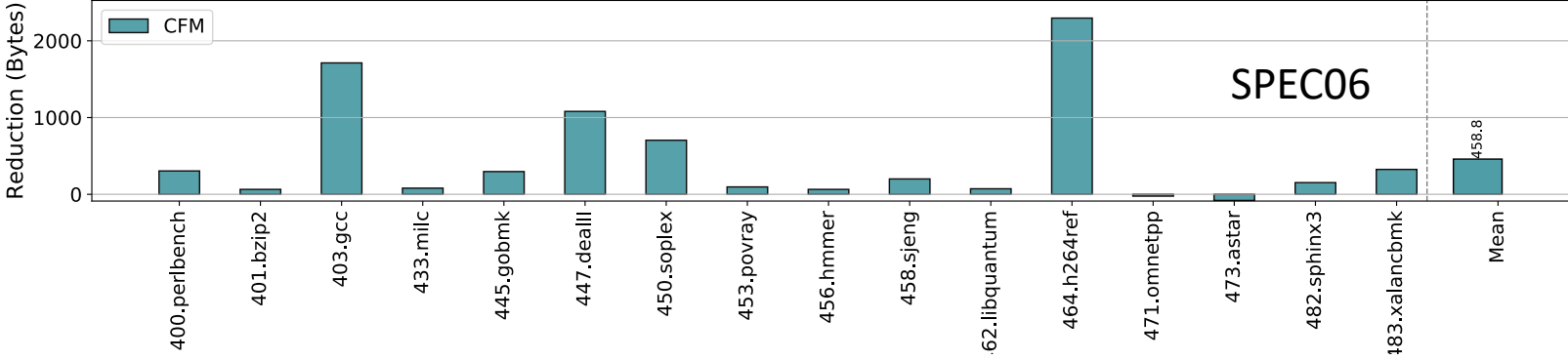
22

# Evaluation (Code Size Reduction)

- Reduction in the text section of the final binary in x86

- Comparison against *–Oz*

- Benchmarks Suites
  - MiBench
  - SPEC 2006
  - SPEC 2017

# Code Size Reductions – MiBench

# Code Size Reductions – SPEC

# Summary

- Similar code inside conditional branches is quite common

- Traditional techniques like code hoisting/sinking are not sufficient to fully exploit the code similarity

- We propose *Control-Flow Melding (CFM)* that merge similar code at region level

- Our LLVM implementation of CFM shows its utility as a general compiler transformation

**Code :**

https://github.com/charitha22/llvm-project-codesize/tree/cfm-dev

**Contact Us :**

Charitha Saumya (cgusthin@purdue.edu)

https://charitha22.github.io/

# Merging Similar Control-Flow Regions in LLVM for Performance and Code Size Benefits

**Charitha Saumya (Presenter),** Kirshanthan Sundararajah, Milind Kulkarni