

How to detect machine regions where threads diverge

1

Late Divergence Analysis

How to detect machine regions where threads diverge

Are all threads enabled at this point?

Is only one threads enabled ?

Where did threads diverge?

Where will they converge?

Late Divergence Analysis

Outline

A motivating example

Existing LLVM infrastructure

Custom pipeline and data model

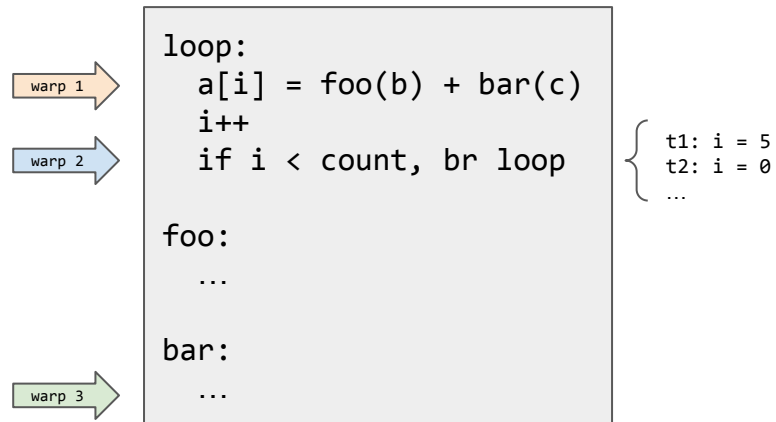
Extensions and applications

Summary

A motivating example

A motivating example

SIMT machine (eg GPU)



Predicate flags control enabled threads

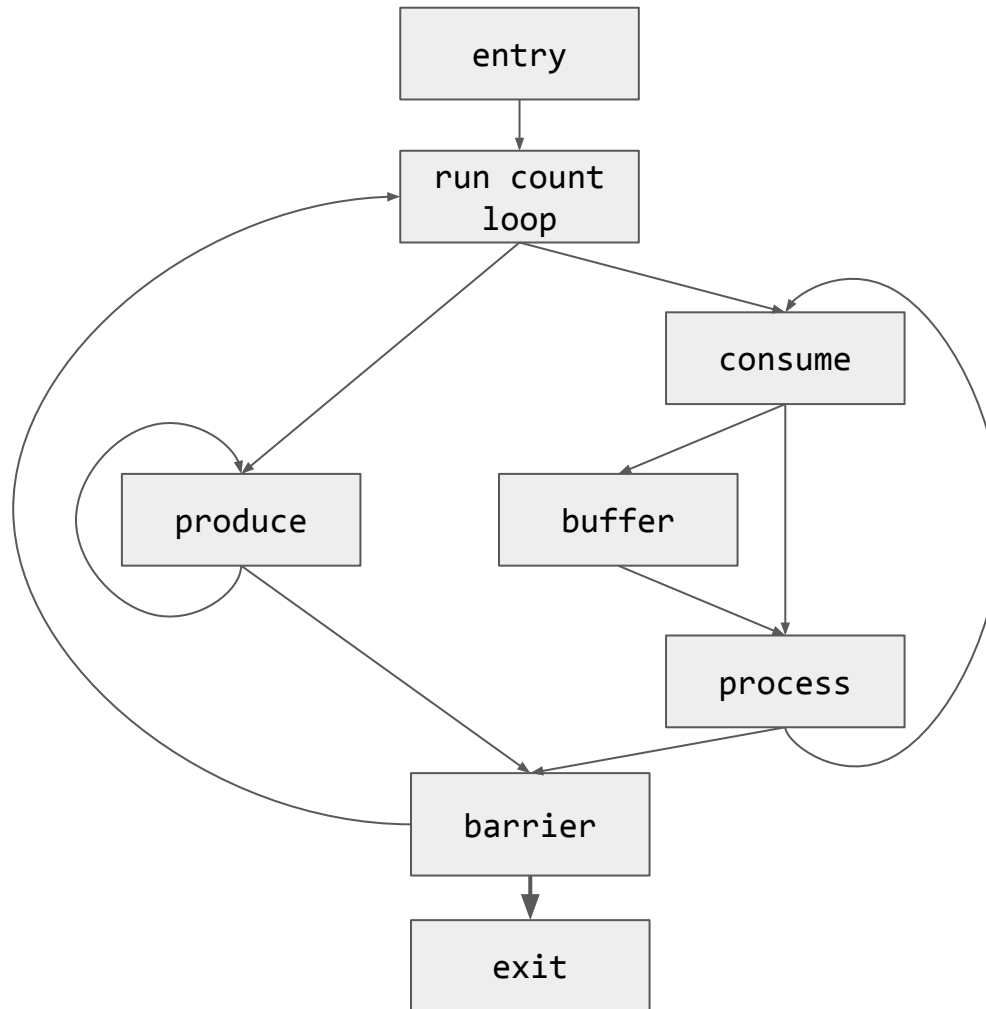
A motivating example

Producer / consumer algorithm on a SIMT machine

```
1.  void run() {
2.    for (unsigned i = 0; i < RUN_COUNT; i++) {
3.
4.        if (is_producer_warp()) {
5.            while(more_to_produce()) produce();
6.
7.        } else {
8.            while(more_to_consume()) {
9.                if (empty_buffer()) buffer_data();
10.                process_data();
11.            }
12.        }
13.
14.        barrier();
15.    }
16. }
```

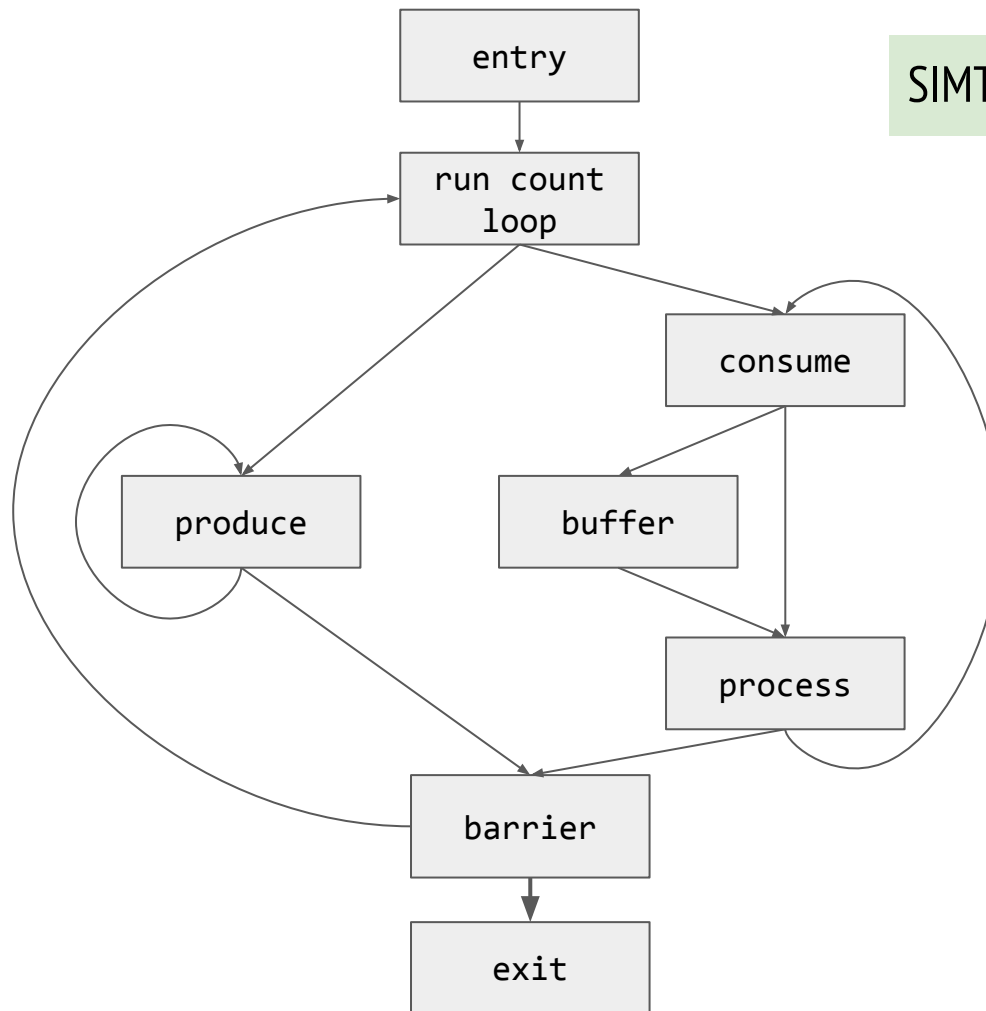
A motivating example

Program control flow



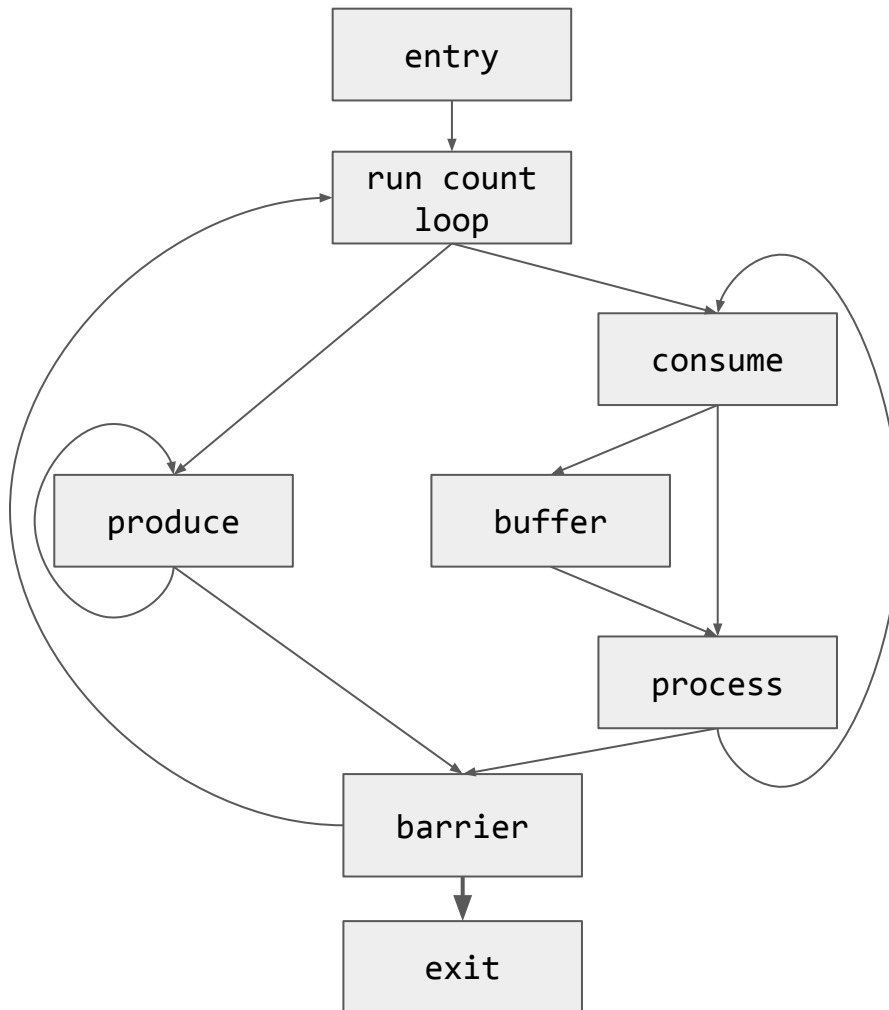
A motivating example

Program control flow



A motivating example

SIMT compatible code



```
entry:
...
run_count_loop:
...
    enabled = is_producer

produce:
...
    enabled = more_to_produce
    if any thread enabled, br produce

consume:
    enabled = is_consumer
...

buffer:
    enabled = empty_buffer
...

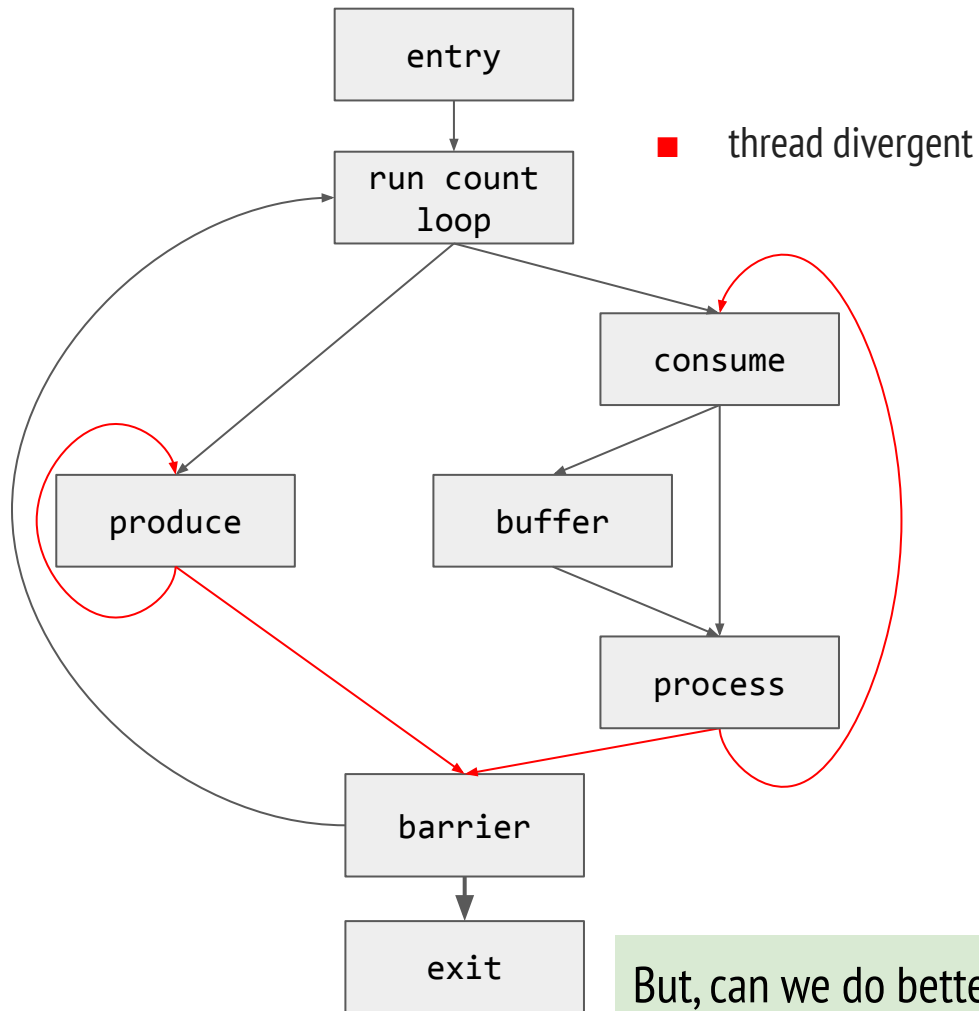
process:
    enabled = is_consumer
...
    enabled = more_to_consume
    if any thread enabled, br consume

barrier:
    enabled = i < RUN_COUNT
...
    i += 1
    enabled = i < RUN_COUNT
    if any thread enabled, br run_count_loop

exit:
```

A motivating example

SIMT compatible code



```
entry:
...
run_count_loop:
...
    enabled = is_producer

produce:
...
    enabled = more_to_produce
    if any thread enabled, br produce

consume:
    enabled = is_consumer
...

buffer:
    enabled = empty_buffer
...

process:
    enabled = is_consumer
...
    enabled = more_to_consume
    if any thread enabled, br consume

barrier:
    enabled = i < RUN_COUNT
...
    i += 1
    enabled = i < RUN_COUNT
    if any thread enabled, br run_count_loop

exit:
```

LLVM infrastructure

Existing infrastructure

LLVM IR divergence

Divergent values pass

Targets specify divergence sources

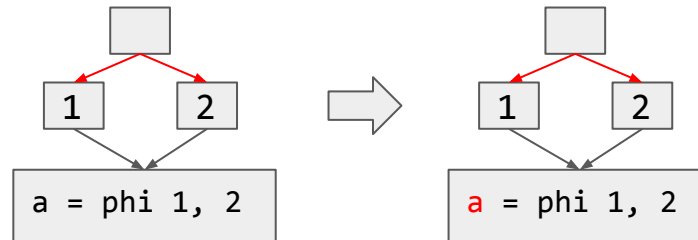
- thread id, iteration variable, memory, etc.

Pass propagates divergence to dependent values

`a = b + c`



`a = b + c`



Existing infrastructure

Draw backs

IR only

Instruction level analysis

- what about block or region level?

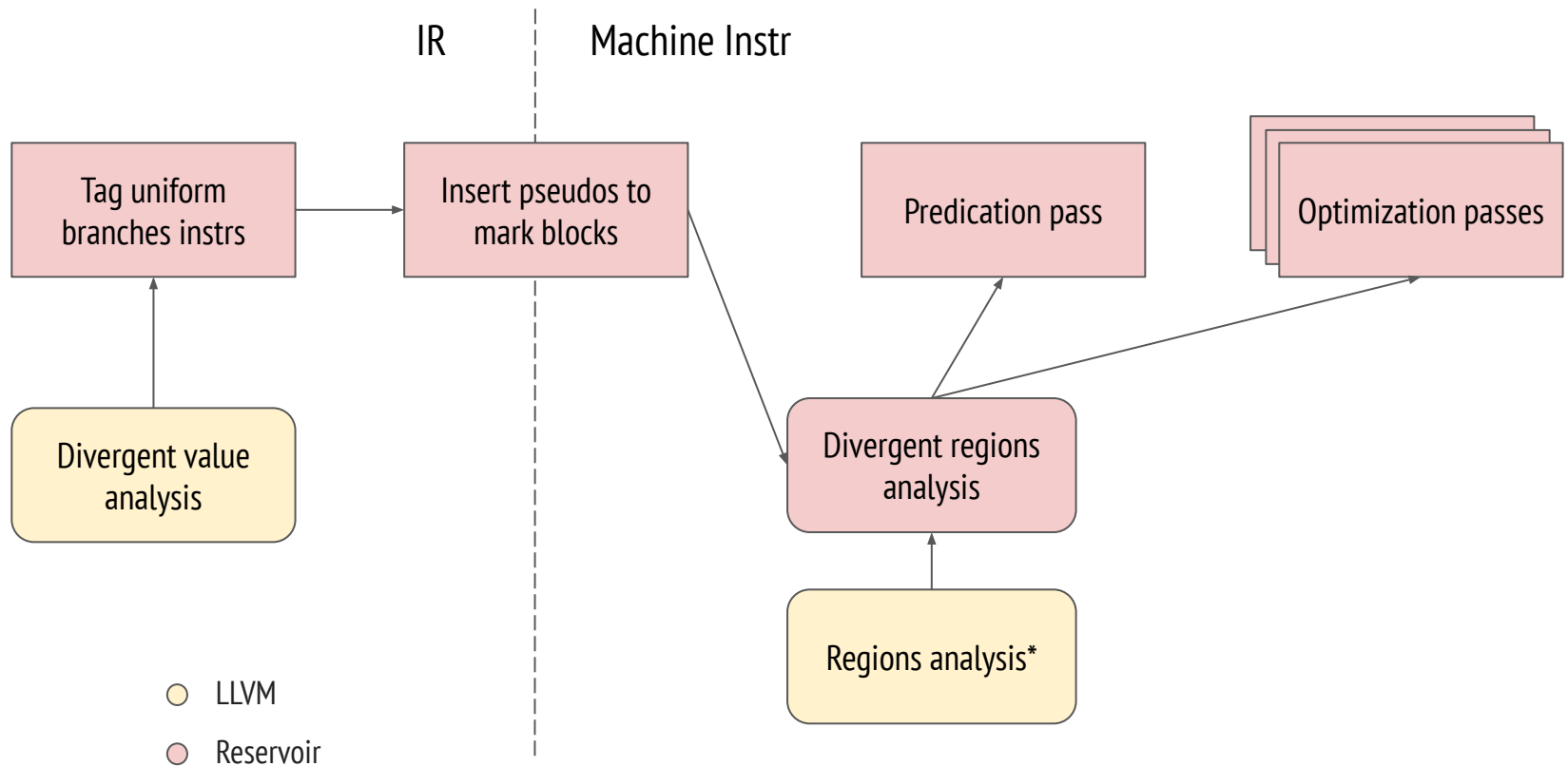
Can't answer the questions

- are all threads enabled at this point?
- is only one thread enabled ?

Custom pipeline

Custom pipeline

Overview



Custom pipeline

Region analysis

LLVM regions

- Program Tree Structure [Johnson, Pearson and Pingali 1994]

Bounded by entry and exit blocks

Custom pipeline

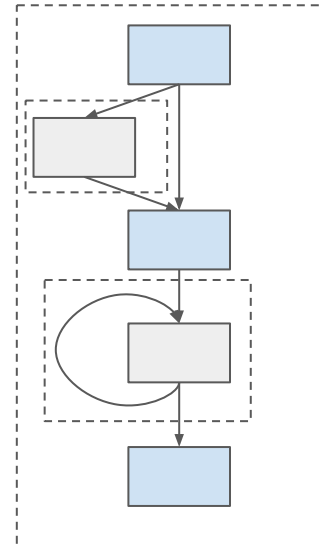
Region analysis

LLVM regions

- Program Tree Structure [Johnson, Pearson and Pingali 1994]

Bounded by entry and exit blocks

Control dependent equivalence



Custom pipeline

Region analysis

LLVM regions

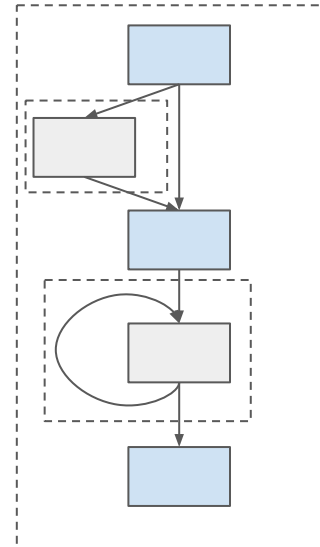
- Program Tree Structure [Johnson, Pearson and Pingali 1994]

Bounded by entry and exit blocks

Control dependent equivalence

Needed custom regions for increase
granularity, modified CFGs

Detail appendix - A1



Custom pipeline and data model

Divergent regions

Threads diverge from immediate dominator region

Custom pipeline and data model

Divergent regions

Threads diverge from immediate dominator region

Propagation:

- Incoming divergent edge(s) -> divergent

Custom pipeline and data model

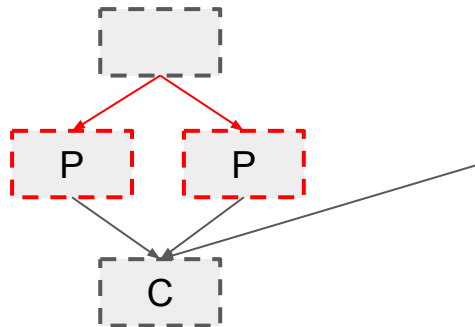
Divergent regions

Threads diverge from immediate dominator region

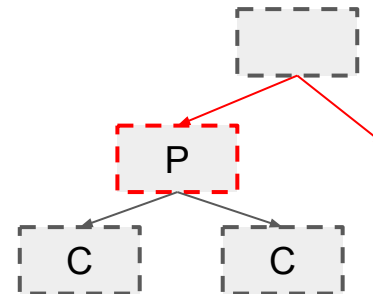
Propagation:

- Incoming divergent edge(s) -> divergent
- Parent divergent -> child divergent, unless

child post-dominates divergent source(s)



parent dominates child



Custom pipeline and data model

Divergent regions

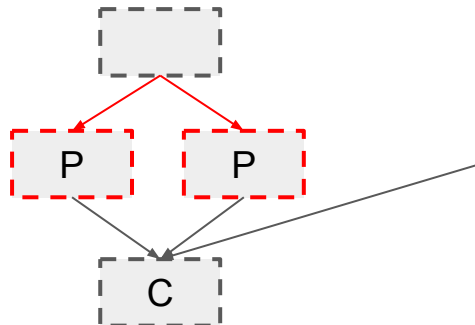
Threads diverge from immediate dominator region

Propagation:

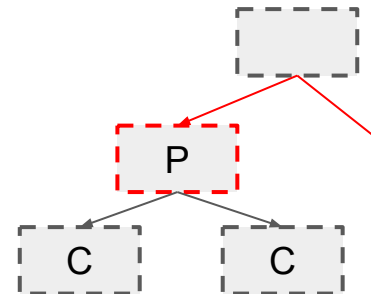
- Incoming divergent edge(s) -> divergent
- Parent divergent -> child divergent, unless

Propagation example in appendix - A2

child post-dominates divergent source(s)



parent dominates child



Custom pipeline and data model

Region properties

Divergent

- threads **can** diverge from immediate dominator region

Uniform

- threads **never** diverge from immediate dominator region

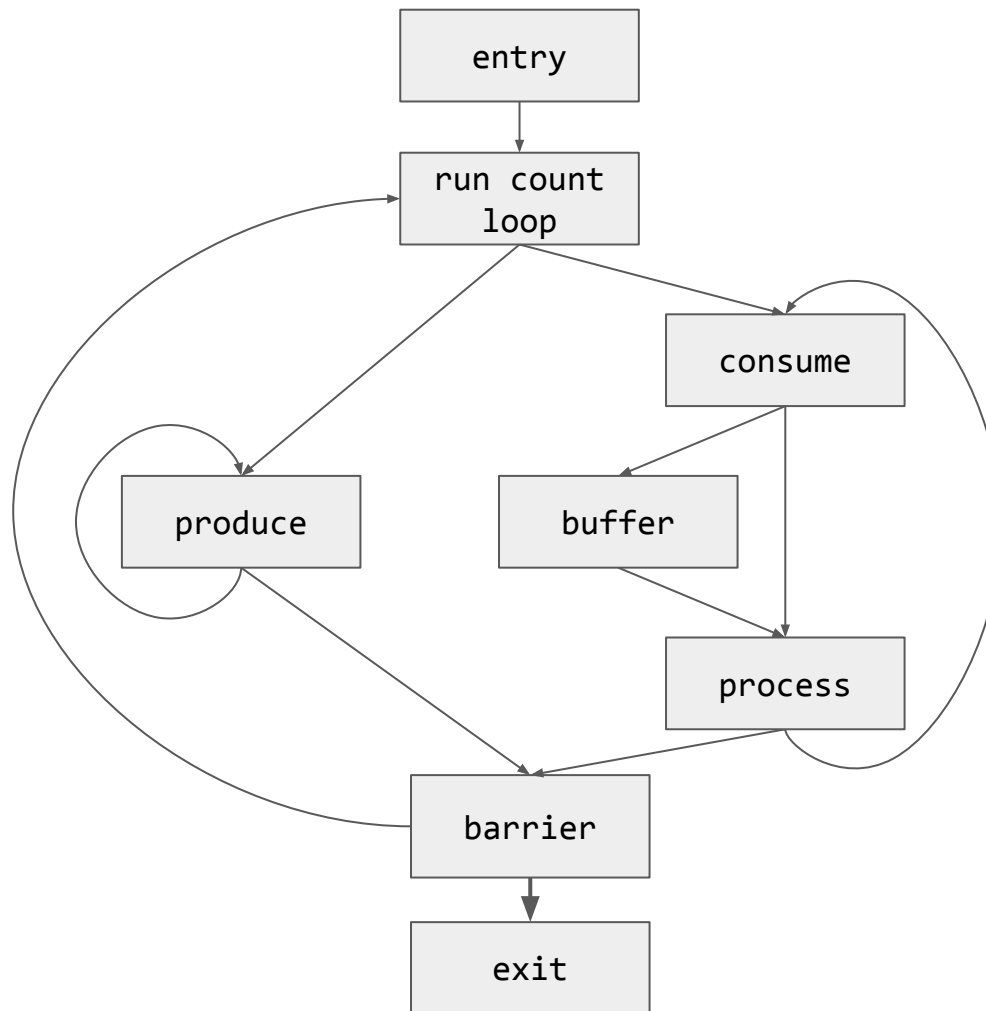
Function uniform

- every dominator is uniform
- threads **never** diverge from function entry

Example revisited

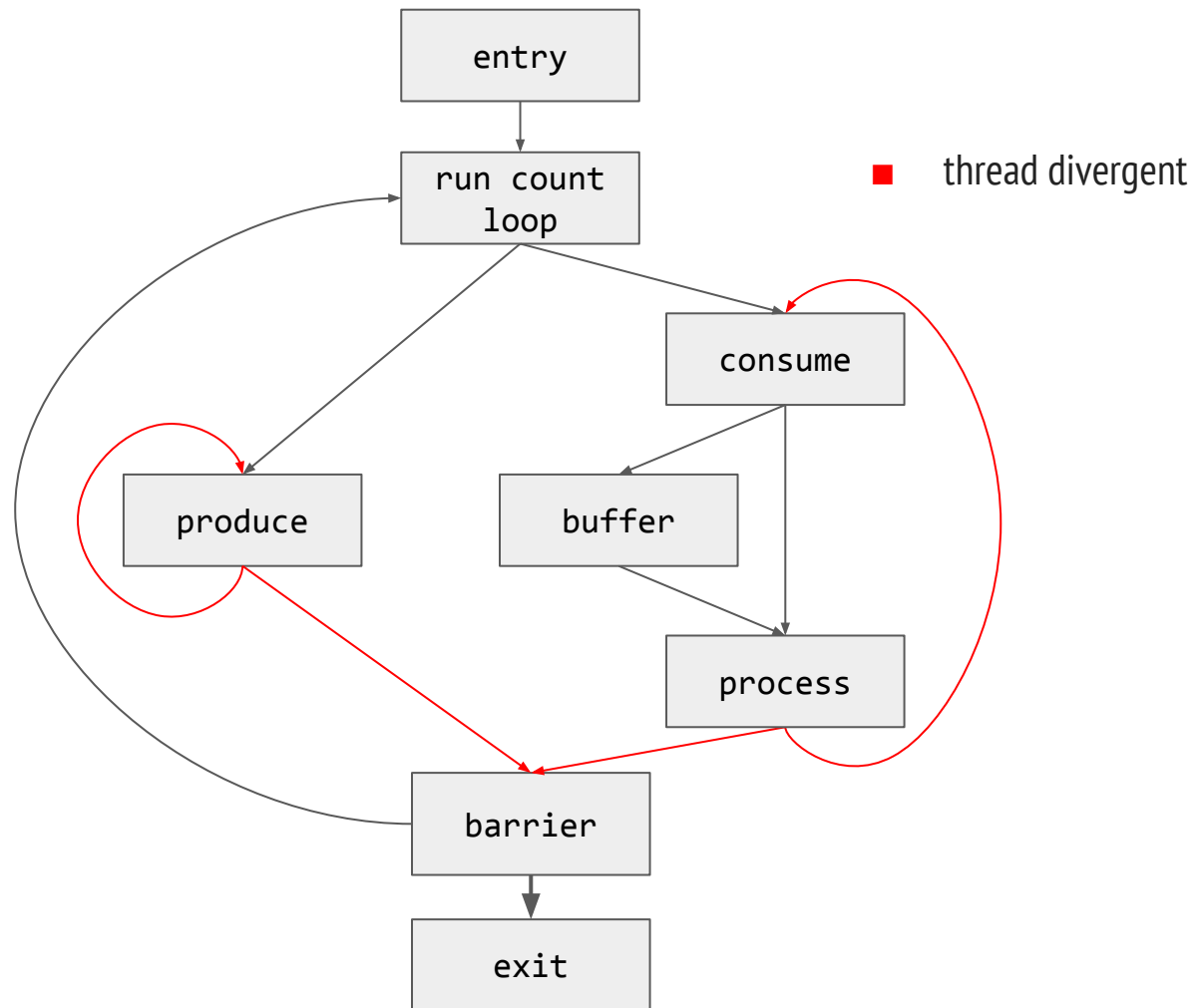
Example revisited

Program control flow



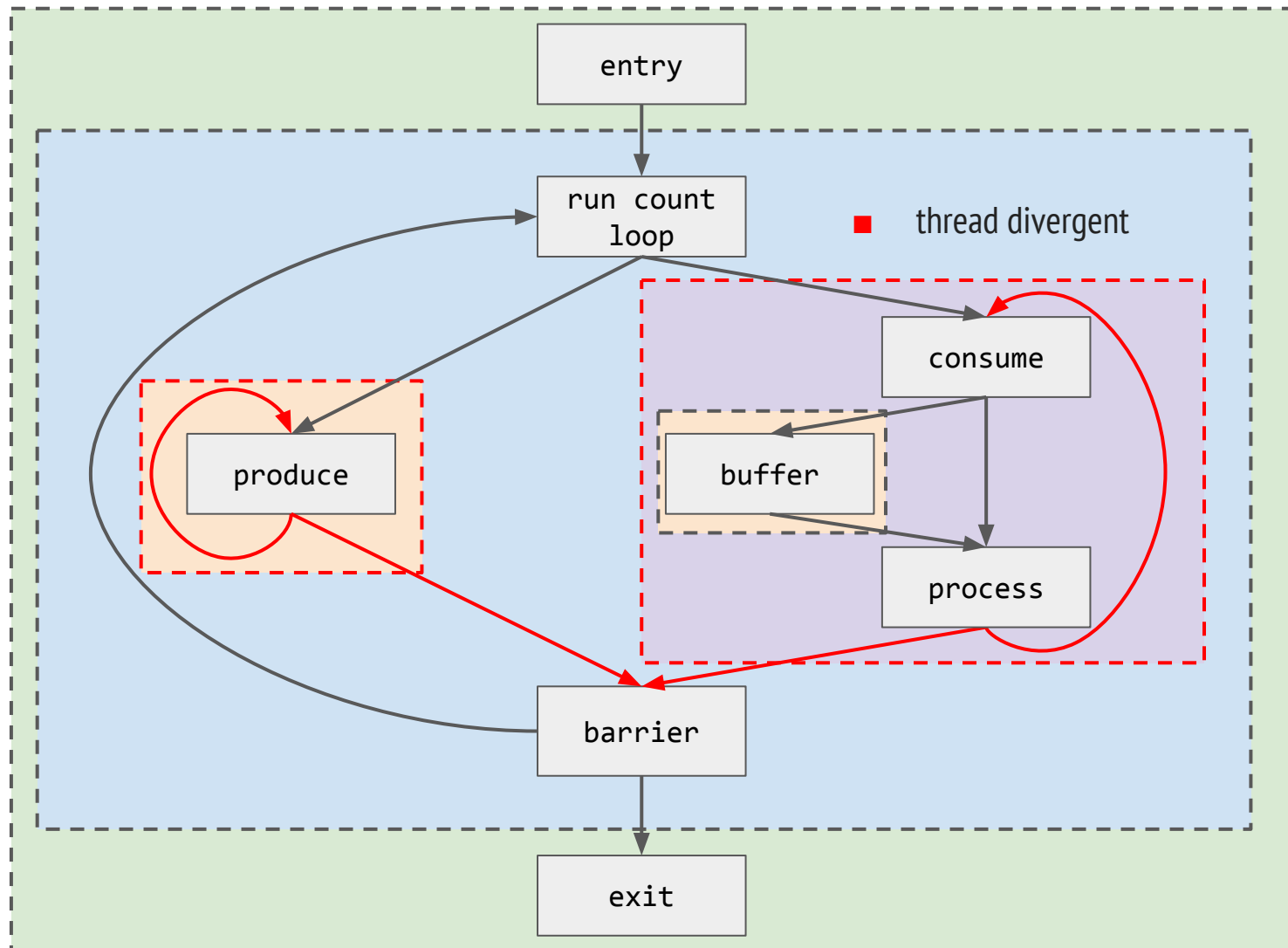
Example revisited

Program control flow



Example revisited

Divergent regions



Example revisited

Emitting SIMD compatible code

Only divergent regions need predicates

Uniform branches can be left in

Example revisited

Emitting SMT compatible code

Only divergent regions need predicates

Uniform branches can be left in*

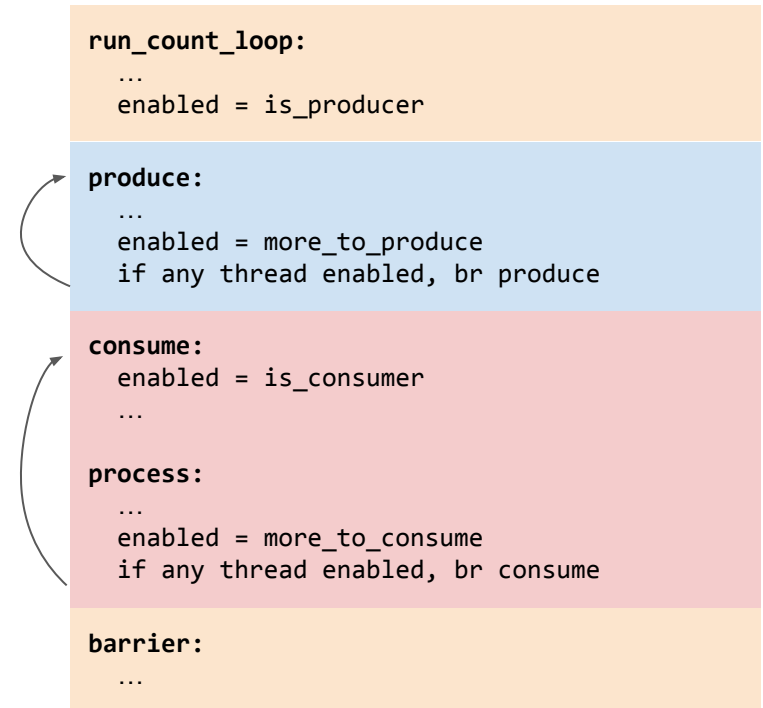
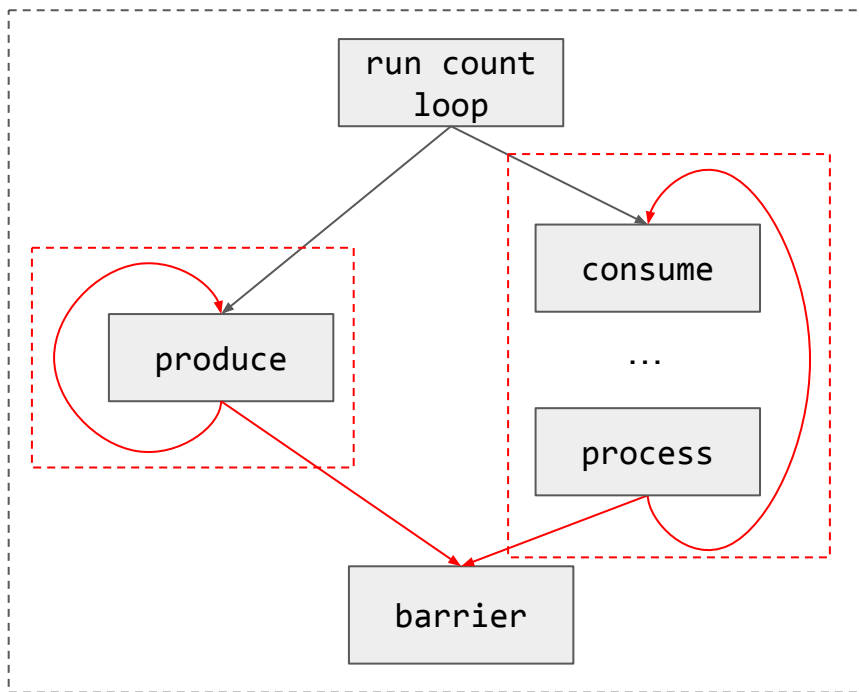
Example in appendix - A3

Example revisited

Emitting SMT compatible code

Only divergent regions need predicates

Uniform branches can be left in

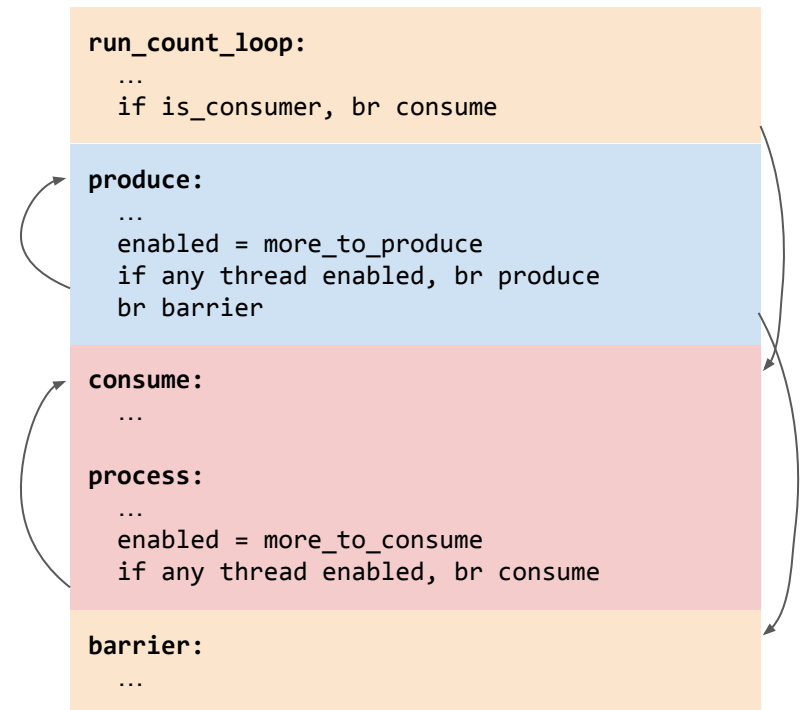
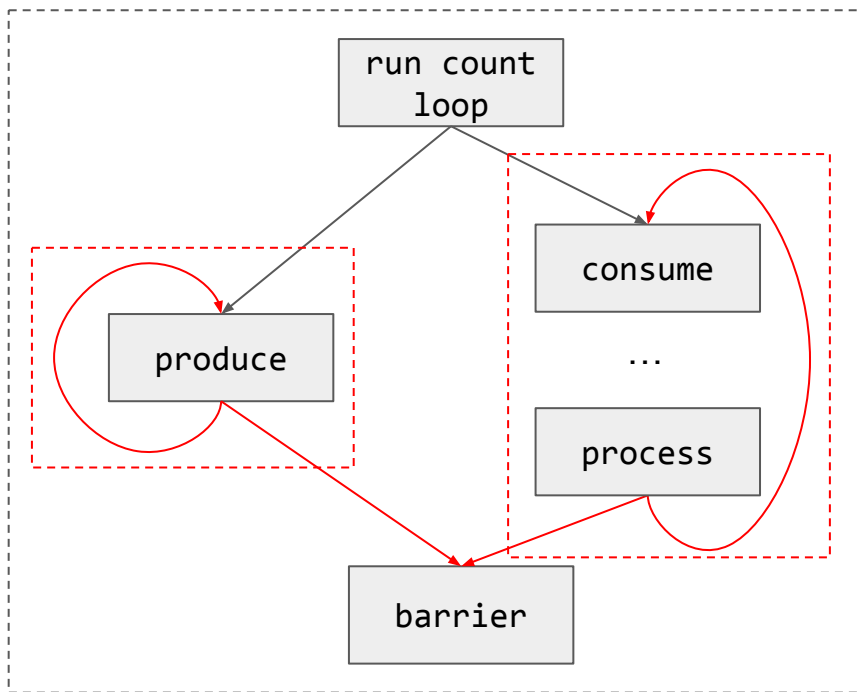


Example revisited

Emitting SMT compatible code

Only divergent regions need predicates

Uniform branches can be left in

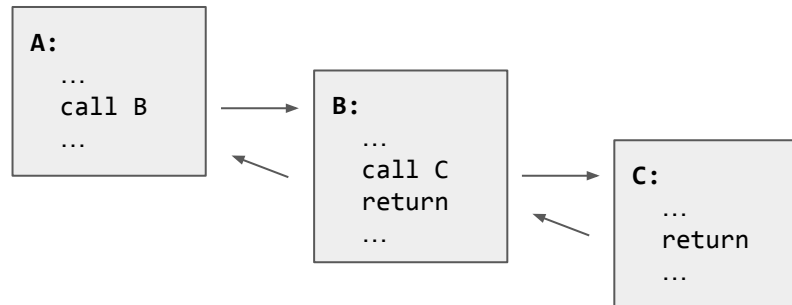


Extensions / applications

Extensions and applications

Applications

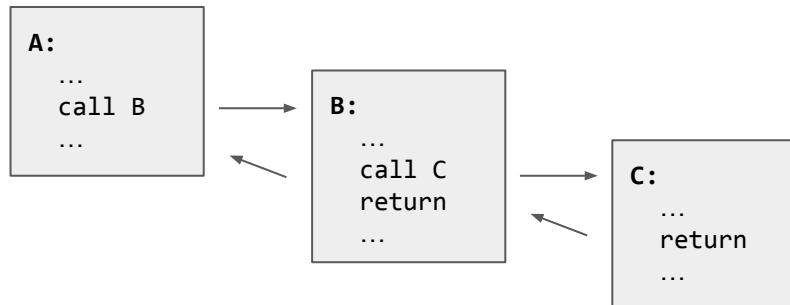
- hardware loop candidates
- tail call optimization candidates



Extensions and applications

Applications

- hardware loop candidates
- tail call optimization candidates



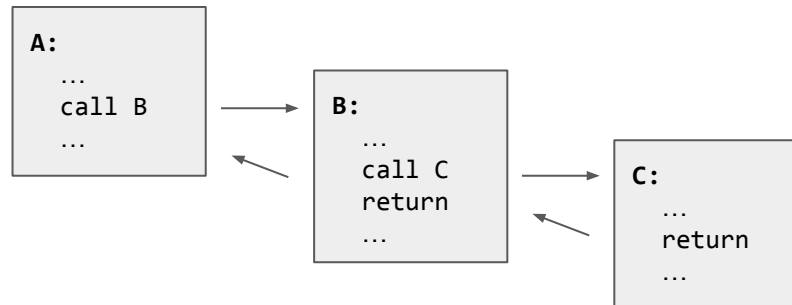
Extension

- thread uniqueness
1. `while(more_to_produce())`
 2. `produce();`
 - 3.
 4. `if (lead_produce(warp_id, thread_id))`
 5. `notify_completion();`

Extensions and applications

Applications

- hardware loop candidates
- tail call optimization candidates



Extension

- thread uniqueness

Details in appendix - A5

```
1. while(more_to_produce())
2.   produce();
3.
4. if (lead_produce(warp_id, thread_id))
5.   notify_completion();
```

Summary

Identifies SESE regions that are

- divergent, uniform, function uniform, thread unique

Analysis combines

- IR divergence, region detection, divergence / uniqueness propagation

Use cases

- predication, tail call optimizations, hardware loops, more?

Appendix

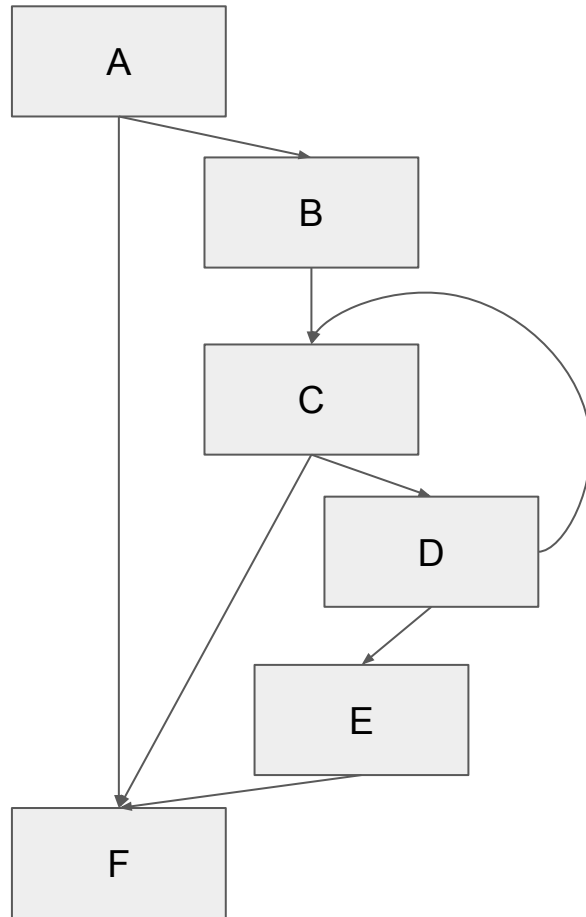
A1

A1 - LLVM regions vs custom region

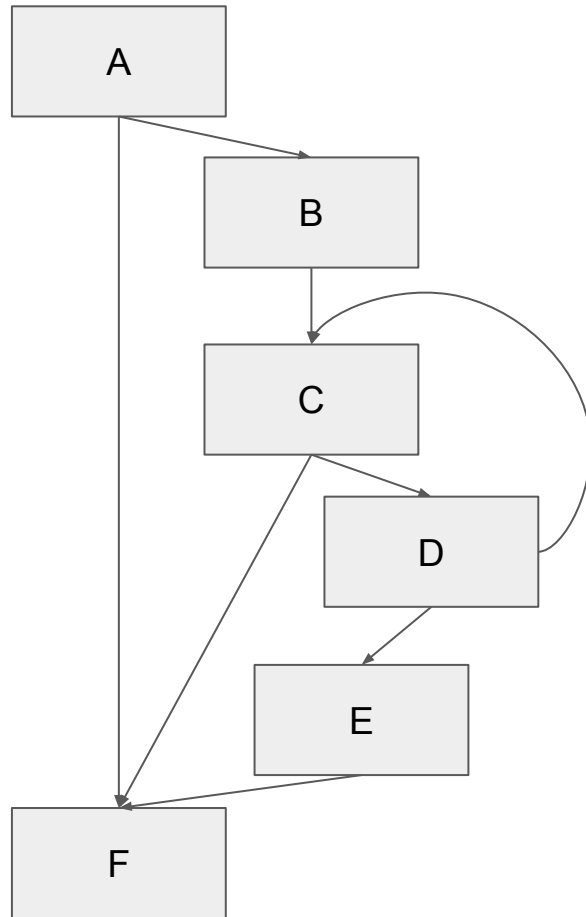
Example from the gcc testsuite: gcc.c-torture/execute/pr54985.c

```
1.  int foo(ST *s, int c) {
2.    int first = 1;
3.    int count = c;
4.    ST *item = s;
5.    int a = s->a;
6.    int x;
7.
8.    while (count--)
9.    {
10.        x = item->a;
11.        if (first)
12.            first = 0;
13.        else if (x >= a) {
14.            return 1;
15.        }
16.        a = x;
17.        item++;
18.    }
19.    return 0;
20. }
```

A1 - LLVM regions vs custom region



A1 - LLVM regions vs custom region



LLVM regions

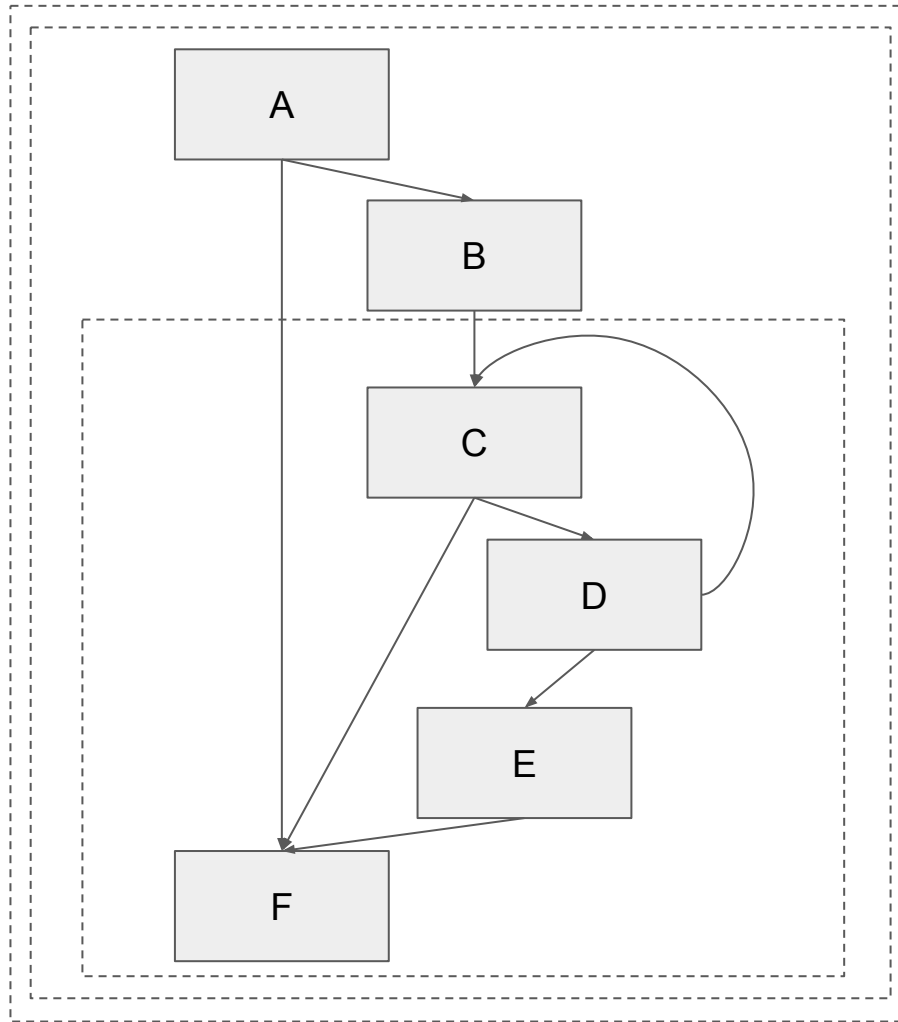
Region tree:

[0] A => <Function Return>

[1] A => F

[2] C => F

A1 - LLVM regions vs custom region



LLVM regions

Region tree:

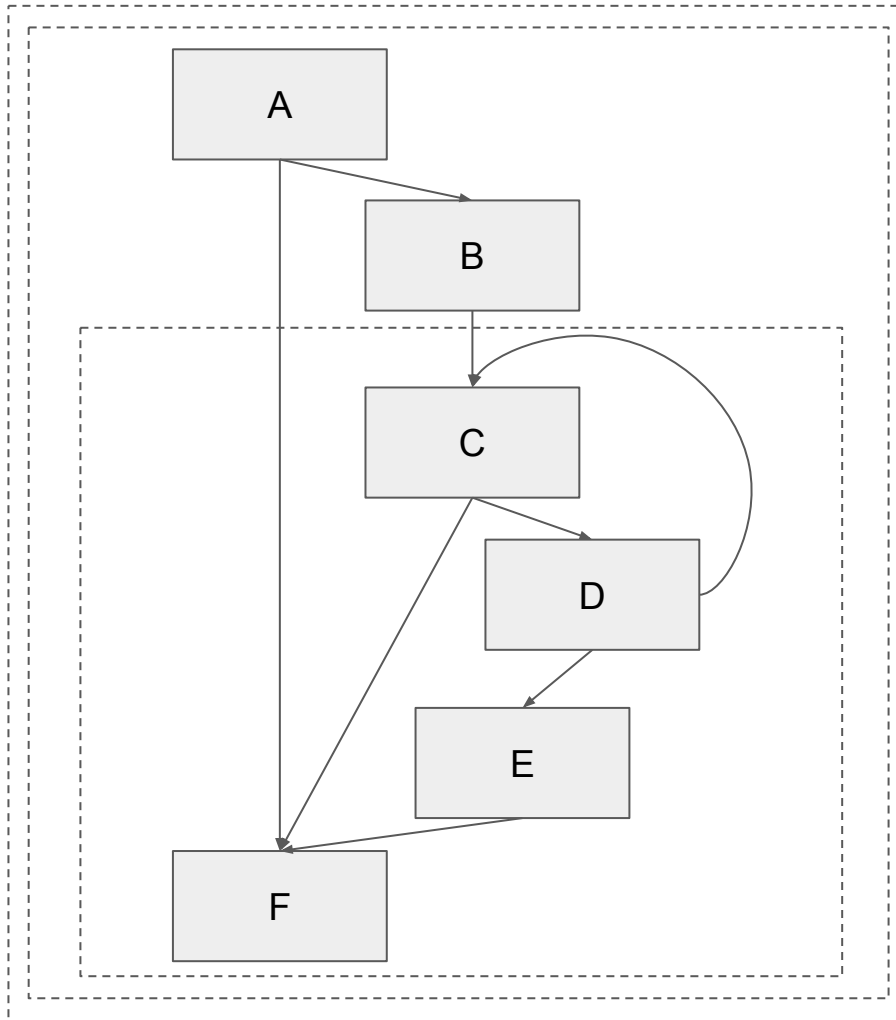
[0] A => <Function Return>

[1] A => F

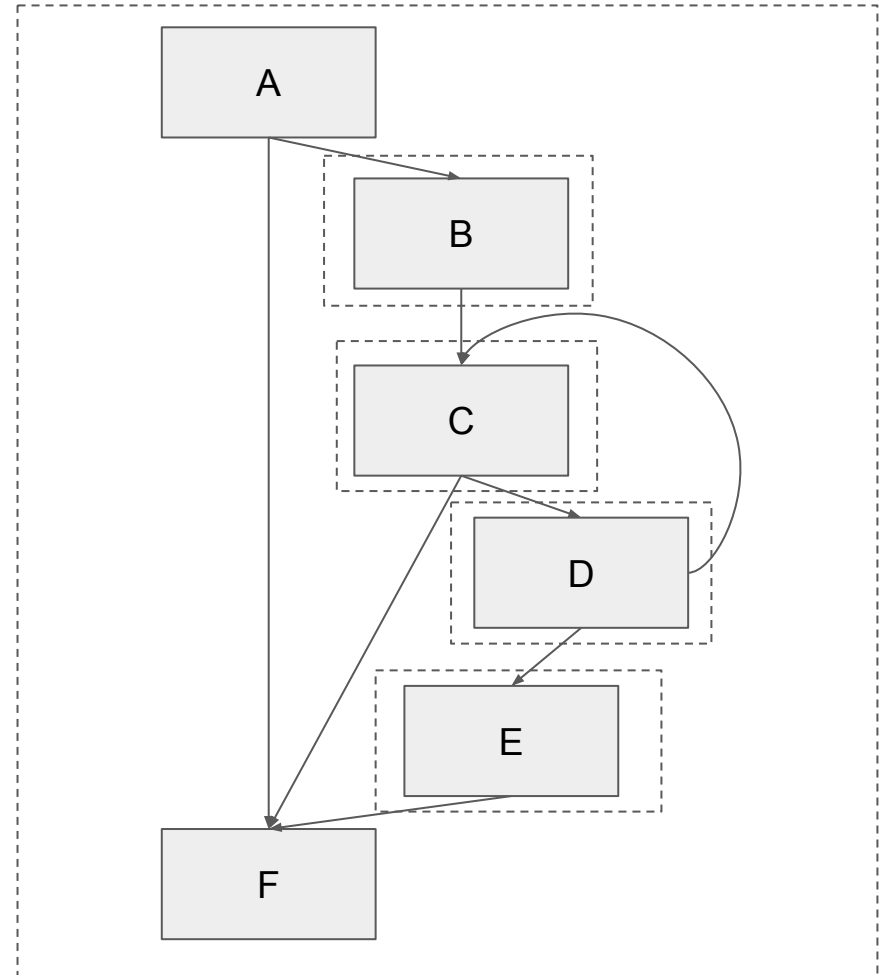
[2] C => F

A1 - LLVM regions vs custom region

LLVM regions



Custom regions



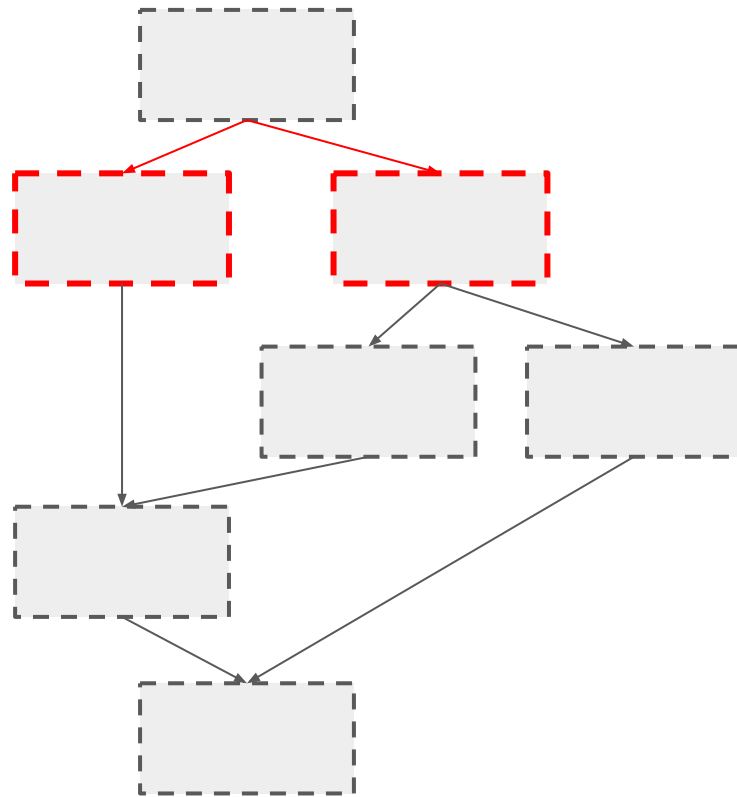
A2

Example 1



A2 - Propagating divergence to SESE regions

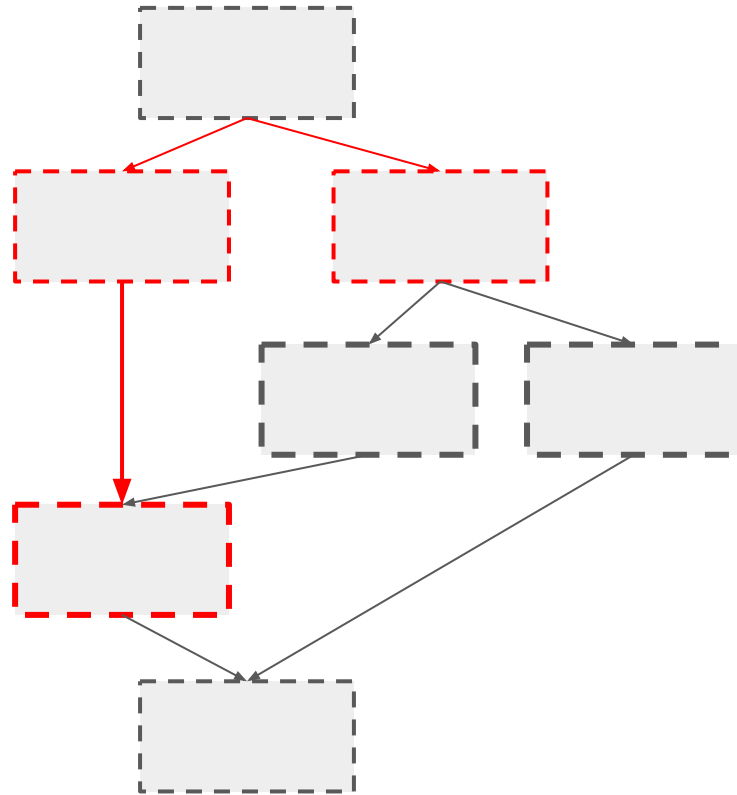
Example 1



Regions with incoming
divergent edge(s) -> divergent

A2 - Propagating divergence to SESE regions

Example 1



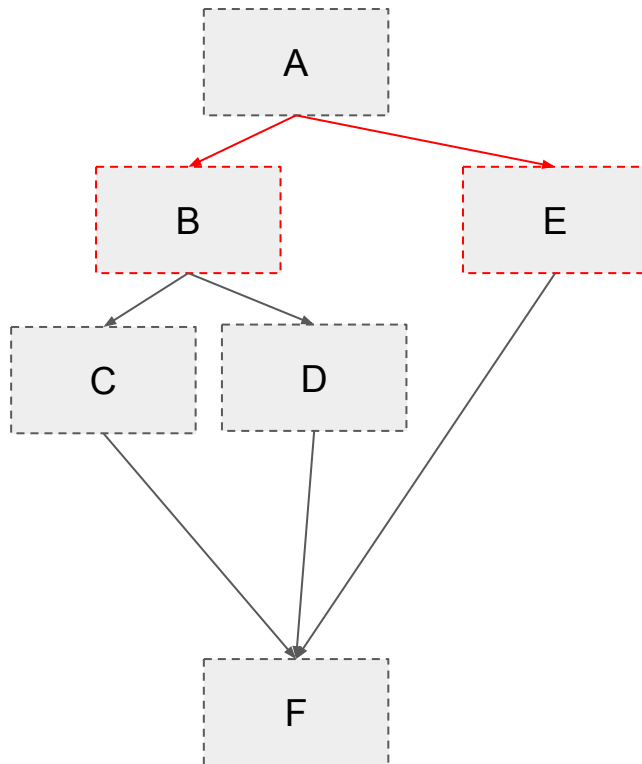
Divergence propagates from regions to child regions, unless

- child dominates divergent source(s)
- parent dominates child

A3

A3 - Generating SIMT code

Restricting branch aheads



Uniform branches can be left in

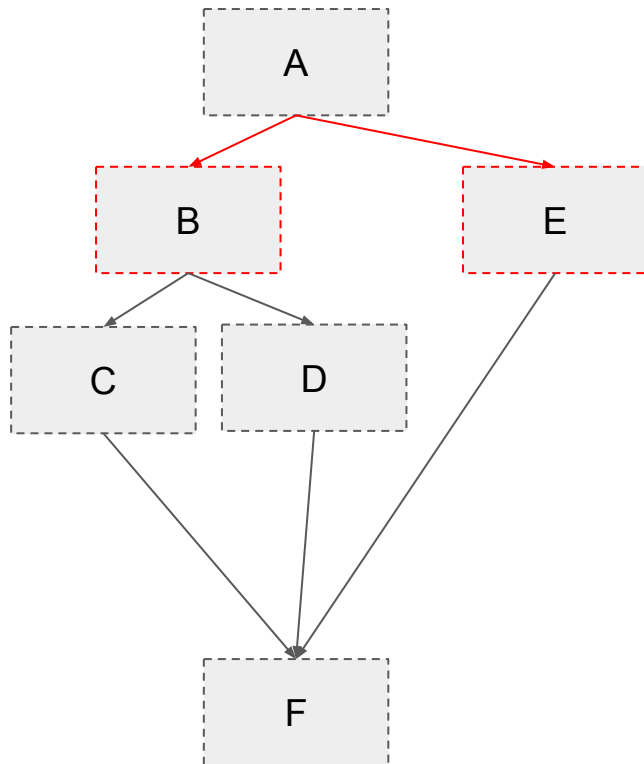
```
A:
...
B:
    enable = do_b
...
    if do_d, br D
C:
...
    br F
D:
...
    br F
E:
    enable = do_e
...
F:
...
```



Will skip over E for
threads that need to
execute it

A3 - Generating SIMT code

Restricting branch aheads



Uniform branches can be left in

- but can only branch as far ahead as they dominate

```
A:
...
B:
    enable = do_b
...
    if do_d, br D
C:
...
    br E
D:
...
E:
    enable = do_e
...
F:
...
```

A4

A4 - Custom pipeline

IR pass

Tag uniform branches

```
1.  DA = &getAnalysis<DivergenceAnalysis>();
2.
3.  for (BasicBlock &B : F) {
4.      Instruction *Term = B.getTerminator();
5.
6.      if ((isa<BranchInst>(Term) || isa<SwitchInst>(Term))
7.          && DA->isUniform(Term)) {
8.
9.          Term->setMetadata("my-target-divergent-branch",
10.                           MDNode::get(Term->getContext(), {}));
11.      }
12. }
```

A4 - Custom pipeline

Add pseudo instructions

Before selection

1. run_count_loop:
2. ...
3. br i1 %cond, label %consume, label %produce_entry, !my.target.uniform !2

After selection

1. bb.0.run_count_loop:
2. ...
3. P_UNIFORM_MBB
4. I_BR %2:cond, %bb.1,consume, %bb.3.produce_entry

A5

A5 - An extension

Thread unique property

Only one thread enabled

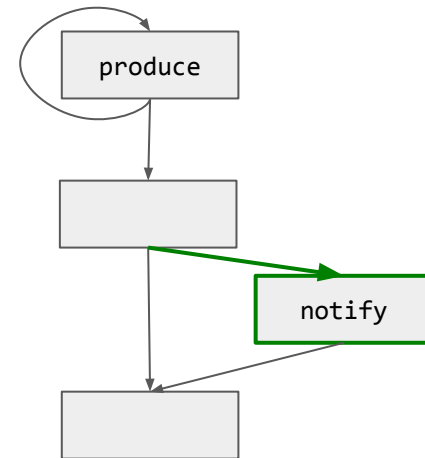
1. `while(more_to_produce())`
2. `produce();`
- 3.
4. `if (lead_produce(warp_id, thread_id))`
5. `notify_completion();`

A5 - An extension

Thread unique property

Only one thread enabled

1. `while(more_to_produce())`
2. `produce();`
- 3.
4. `if (lead_produce(warp_id, thread_id))`
5. `notify_completion();`

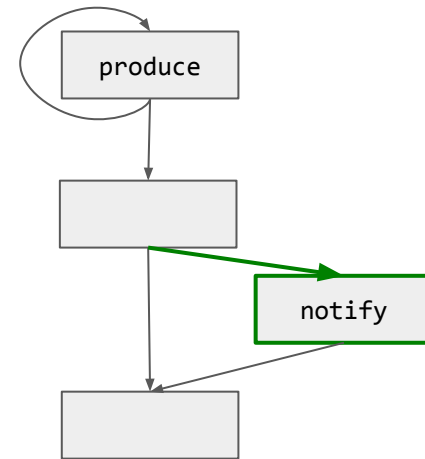


A5 - An extension

Thread unique property

Only one thread enabled

```
1. while(more_to_produce())
2.   produce();
3.
4. if (lead_produce(warp_id, thread_id))
5.   notify_completion();
```



Value propagation:

- Thread unique sources
- Operations that persevere uniqueness

Region propagation:

- One thread unique entry edge-> thread unique
- IDom unique -> child unique

A5 - An extension

Thread unique property

Use case: predication

Regions dominated by thread unique regions

- don't need a predicate
- all branches can be left in*

A6

A6 - Other applications

Uniform regions

Use case: hardware loops

Some hardware loops may require uniform loop conditions

Loops in uniform regions are candidates

A6 - Other applications

Function uniform regions

Use case: tail call optimizations

Standard tail call optimizations can be applied when return is in a function uniform region

References

References

- [1] Johnson, R., Pearson, D., and Pingali, K. “The program structure tree: computing control regions in linear time.” Programming Language Design and Implementation, 1994
- [2] Vanhatalo, J., Völzer, H., Koehler, J. “The refined process structure tree.” Data Knowl. Eng. 2009, 68, 793–818