

A Stride Towards Generating Segment Accesses in RVV

Athanasios Kastoras | Roger Ferrer Ibañez | Lorenzo Albano
 {athanasios.kastoras, roger.ferrer, lorenzo.albano}@bsc.es

The Segment Access Lowering Pass

An LLVM IR Pass in the RISC-V Target Backend

Analysis

Set of rules applied to group of instructions, split to obey them.

We are left with a list of lists of strided vector instructions that can be correctly replaced by a segment instruction.

Transformation

Insertion point: Before the first strided instruction.

Segment Access Intrinsics operate on target-specific data types: tuples of vectors strided instructions replaced by tuple inserts (stores) or extracts (loads).

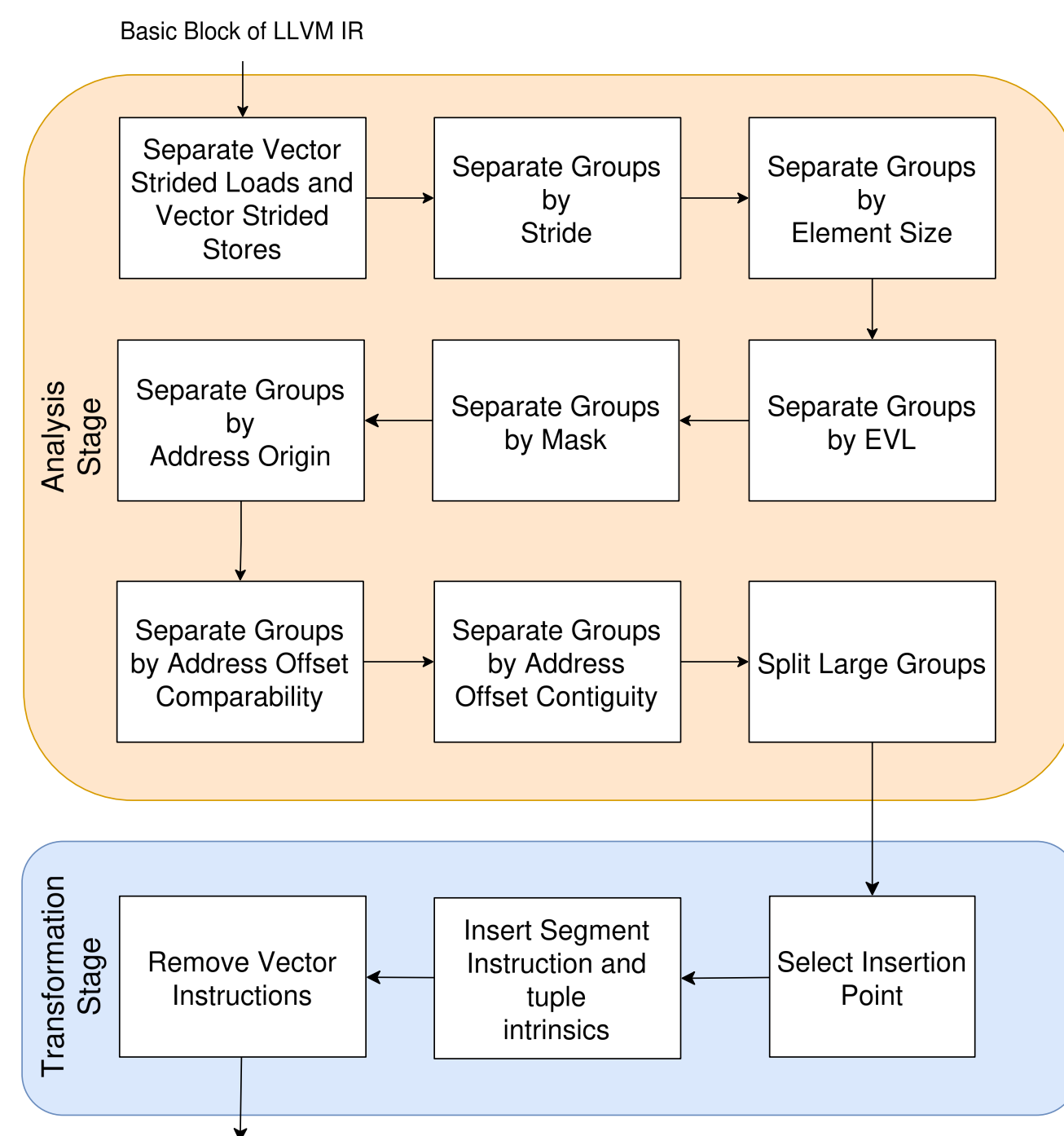
Why in the Backend?

No target-agnostic Segment intrinsics in LLVM IR.

No representation of contiguous vector groups in target-agnostic LLVM IR.

llvm.vector(de)interleave[2-8] would work only for unit-stride instructions.

A VPlan pass would need non-trivially extending LLVM IR.

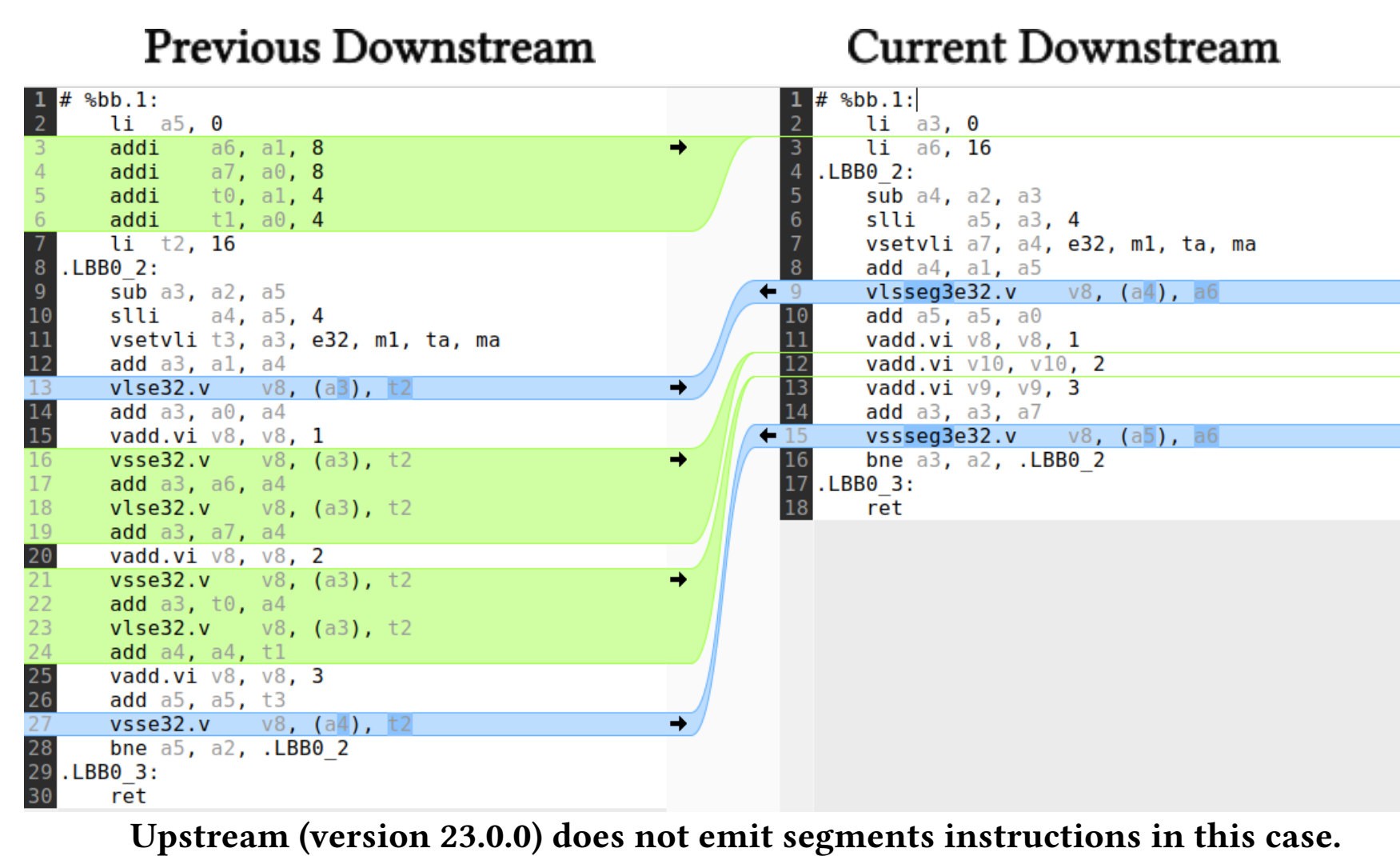


Emitting Strided Instructions

Example Code: rgba.c

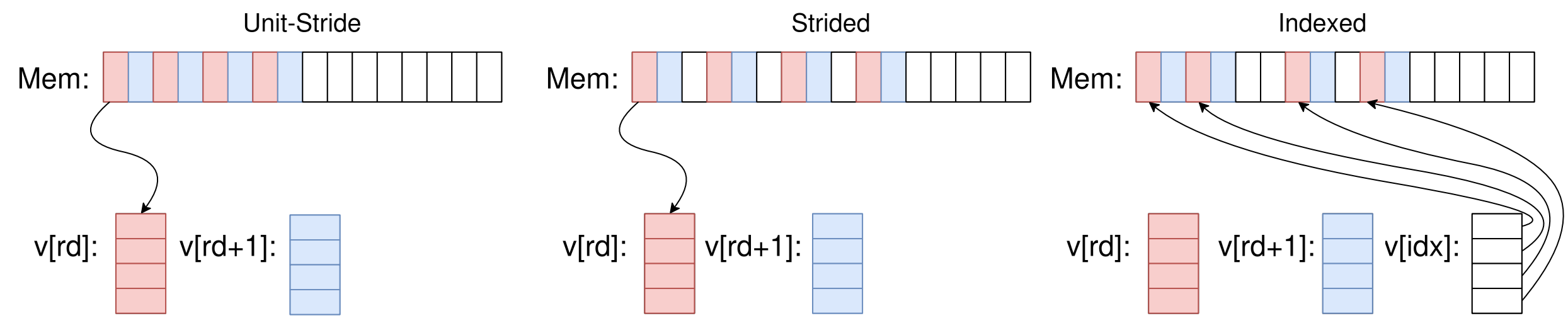
```
struct pixel { int r, b, g, a; };

void vecops(struct pixel * __restrict_ outimg,
            struct pixel * __restrict_ img, int n) {
    for (int i = 0; i < n; i++) {
        outimg[i].r = img[i].r + 1;
        outimg[i].g = img[i].g + 2;
        outimg[i].b = img[i].b + 3;
    }
}
```



Background

The RVV segment instructions support accessing memory in a AoS-to-SoA fashion, allowing for a more efficient memory access pattern. There are three kinds of segment accesses: Unit-Strided, Strided, and Indexed. Here is an illustration of segment accesses with a segment size of 2.



Upstream, supports emitting Unit-Strided instructions. We present an alternative way to emit Unit-Strided and Strided Segment instructions that could be extended to emitting Indexed instructions as well.

Downstream LLVM-based EPI compiler's different approach to loop vectorization:

Extra analysis on gathers and scatters to emit vp.experimental.strided intrinsics already in the VPlan when possible

Smarter VPlan costing can be achieved by more precisely costing memory accesses (e.g. cost of gather > cost of a strided load) leading to more precise vectorization decisions.

This allows us to simply analyze strided accesses in the Loop Vectorizer (instead of the backend) and make decisions early.

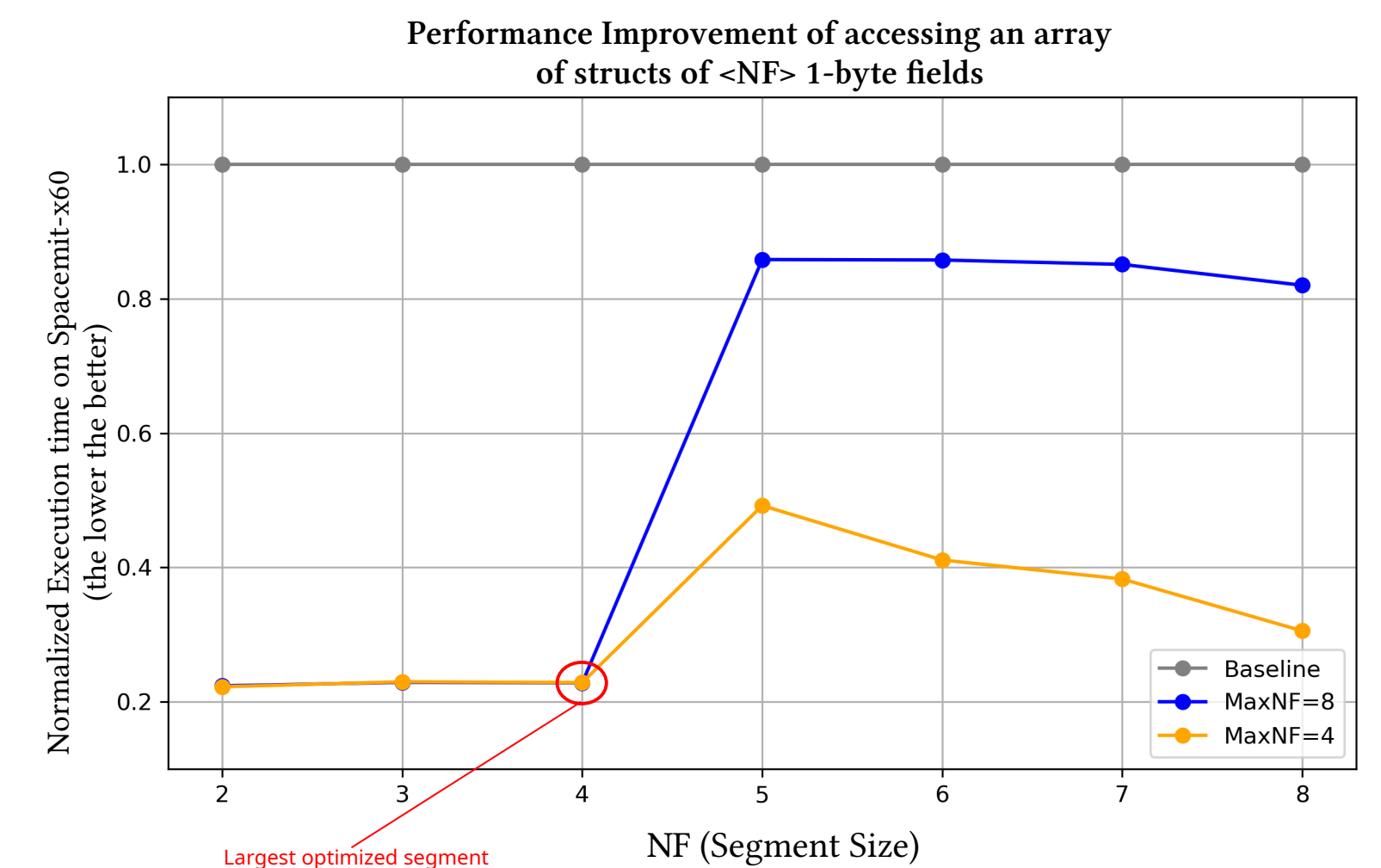
The larger the segment the merrier?

Choosing the Segment Size
 Not all segment sizes are efficient on all CPUs.

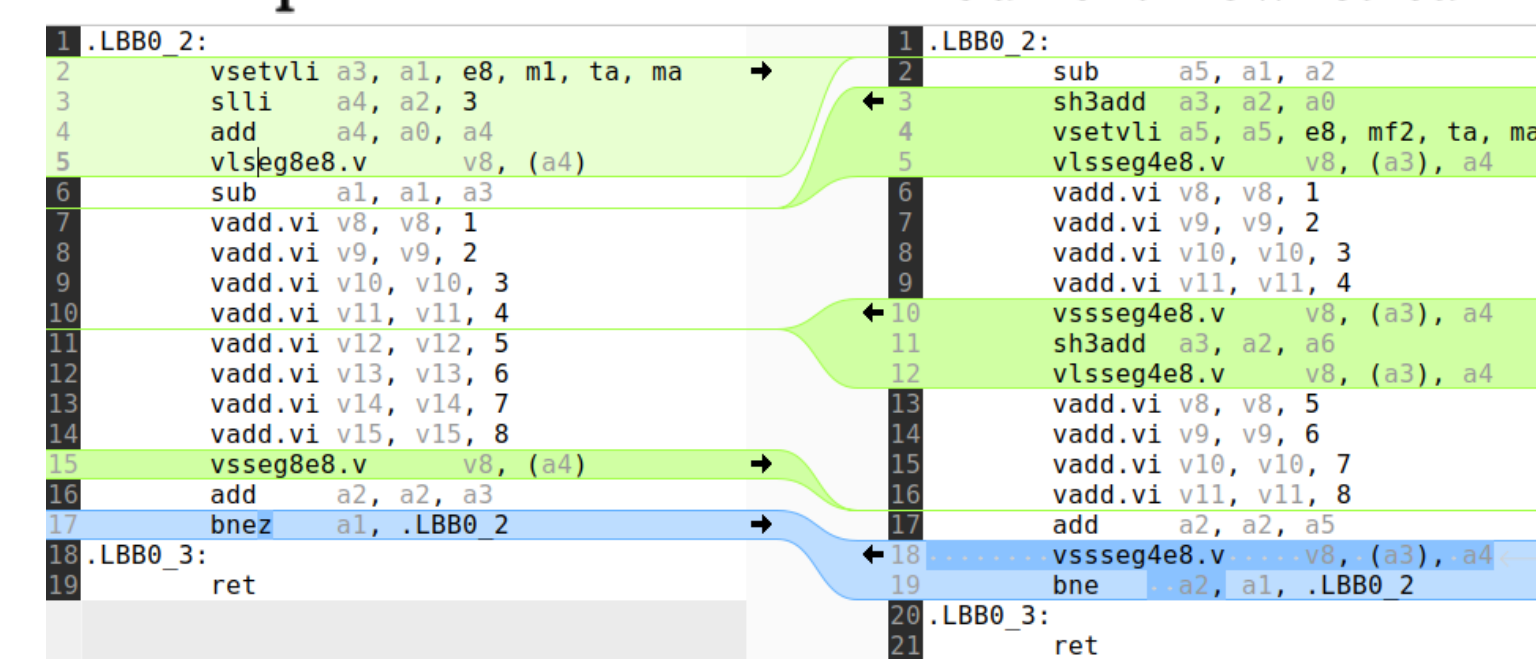
Largest efficient segment size already in RISCVProcessors.td.

Upper limit in the 'Split Large Groups' rule of the analysis part of the pass.

Larger segment accesses are split in smaller segment accesses.



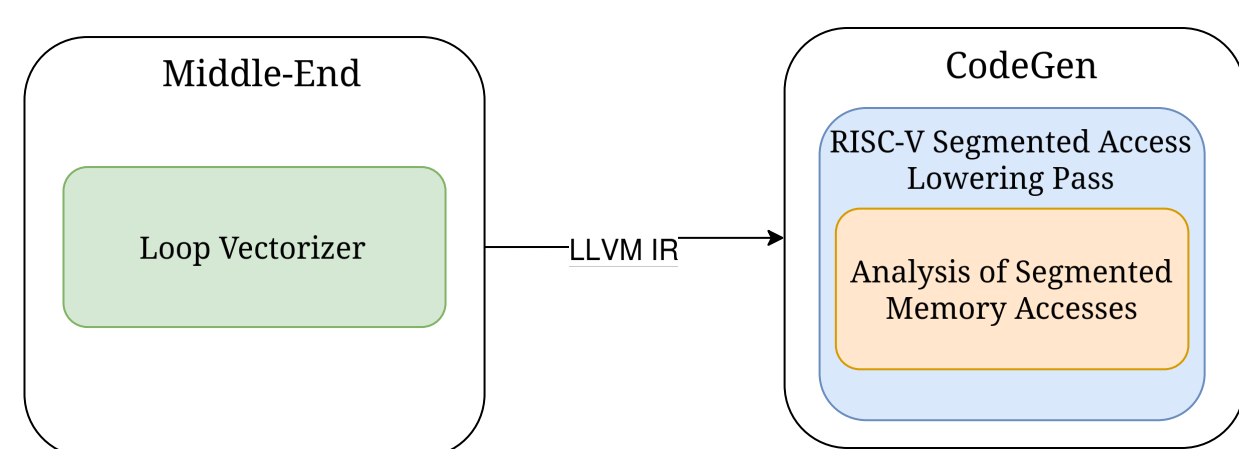
Upstream



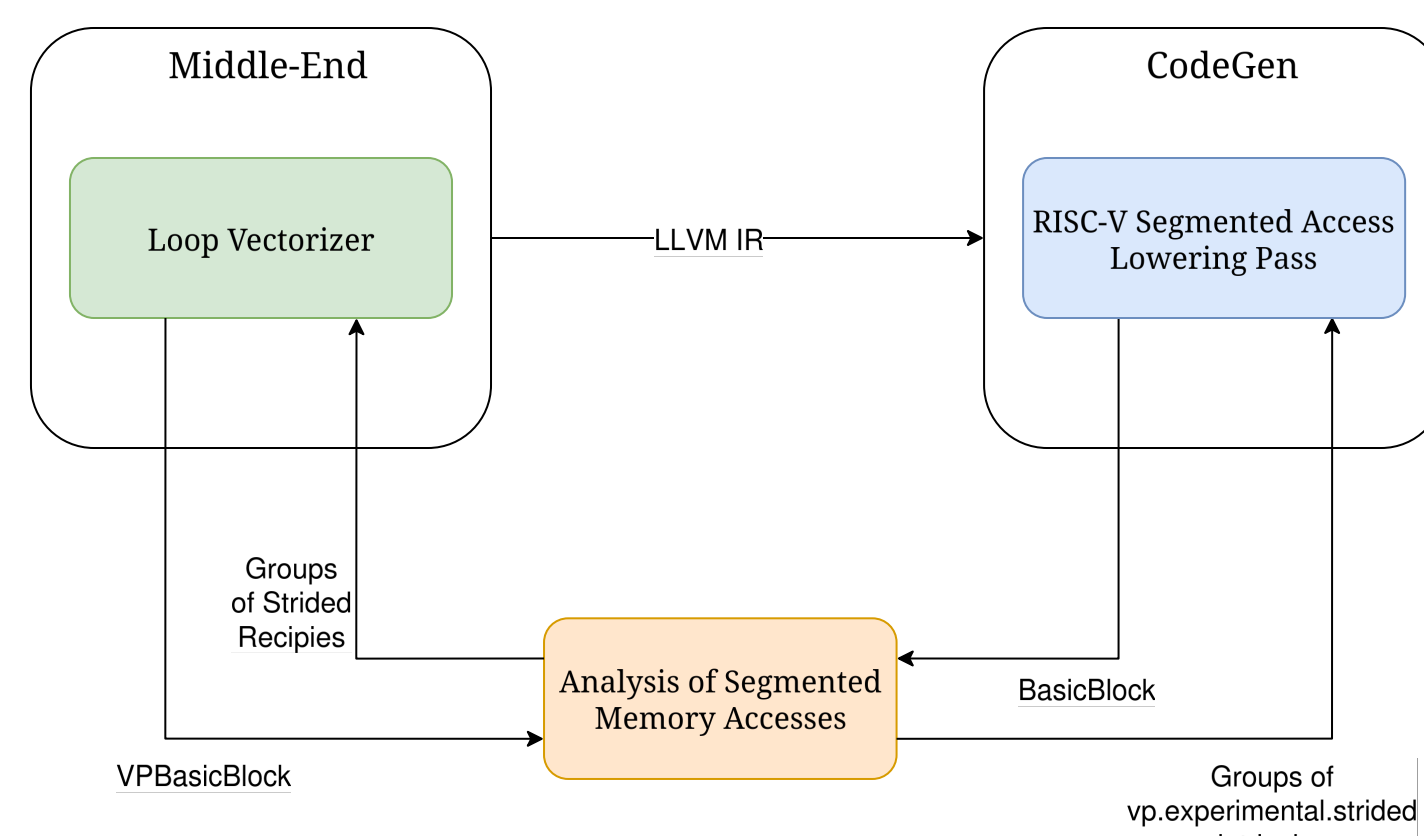
The microbenchmark code

```
void foo(char a[][NF], int N) {
    for (int i = 0; i < N; i++) {
        // inner-loop will be unrolled
        for (int j = 0; j < NF; j++) {
            a[i][j] += j+1;
        }
    }
}
```

Segment Access Analysis in the BackEnd



Segment Analysis Informed VPlan Costing

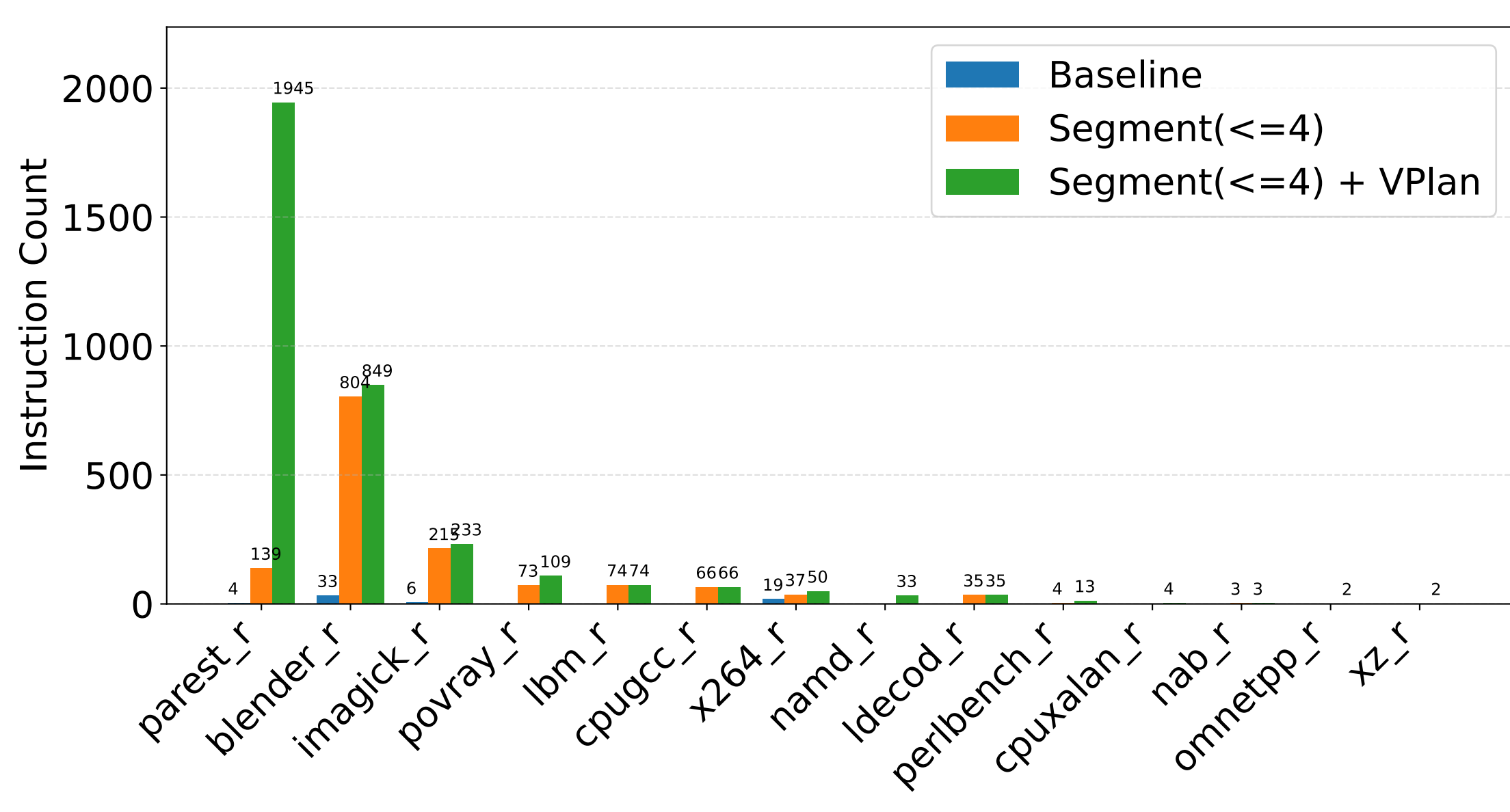


Real Life Usecase: SPEC2017 - 557.xz_r/sha-2/sha-512.c

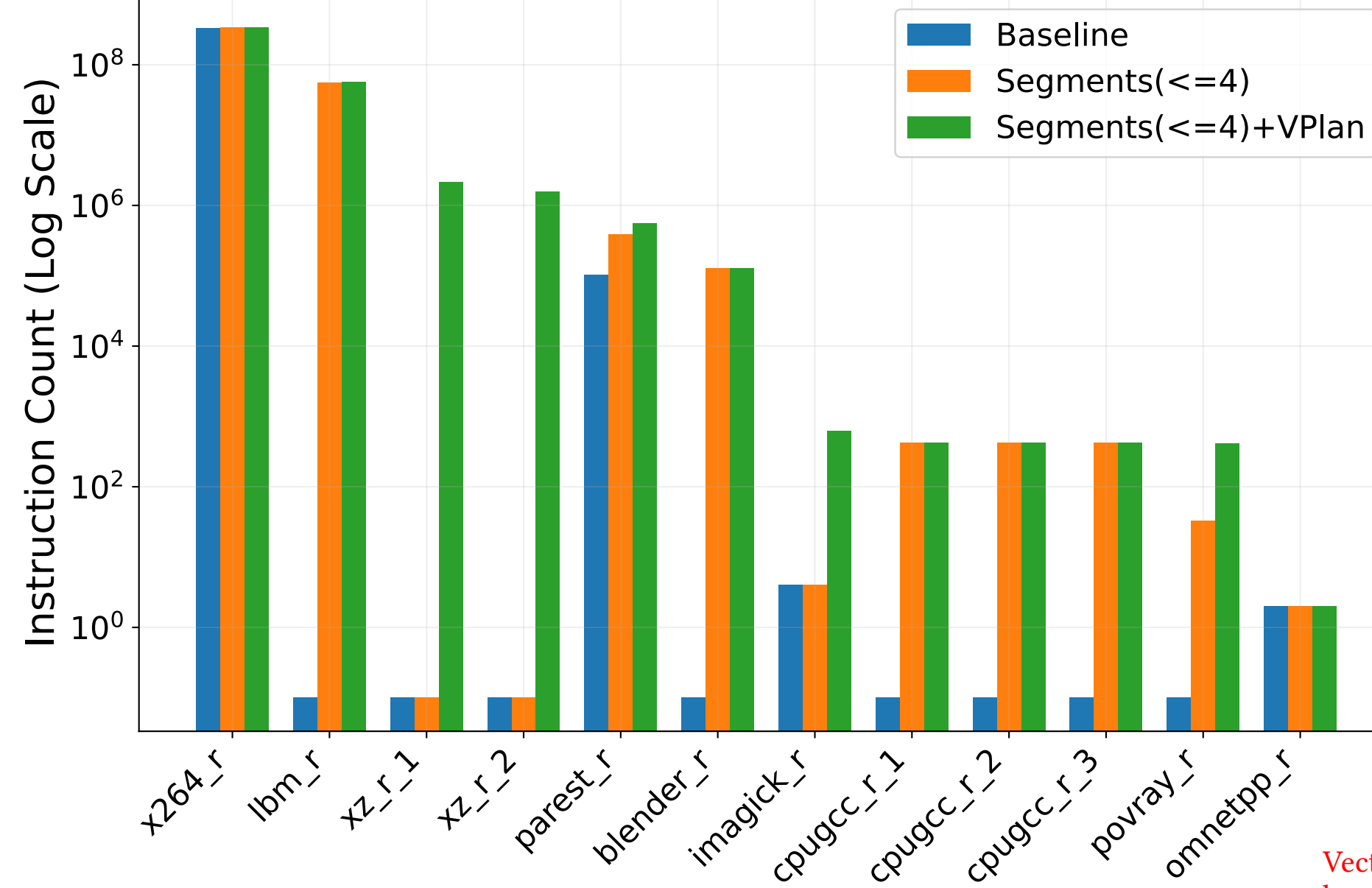
VPlan Before: Is not vectorized	VPlan After: Is vectorized
VPLAN: Computing cost for VF vscale x 1. VPLAN: computing cost for VPBasicBlock vector.body with cost kind RecipThroughput.	VPLAN: Computing cost for VF vscale x 1. VPLAN: computing cost for VPBasicBlock vector.body with cost kind RecipThroughput.
VPLAN: cost of 2 for recipe: EPI-WIDEN ir<4> = vp.strided.load ir<arrayptr>, vp<4>, ir<...>	VPLAN: cost of 0 for recipe: EPI-WIDEN ir<4> = vp.strided.load ir<...>
VPLAN: cost of 2 for recipe: EPI-WIDEN ir<5> = vp.strided.load ir<arrayidx.1>, vp<4>, ir<...>	VPLAN: cost of 0 for recipe: EPI-WIDEN ir<5> = vp.strided.load ir<...>
VPLAN: cost of 2 for recipe: EPI-WIDEN ir<6> = vp.strided.load ir<arrayidx.2>, vp<4>, ir<...>	VPLAN: cost of 0 for recipe: EPI-WIDEN ir<6> = vp.strided.load ir<...>
VPLAN: cost of 2 for recipe: EPI-WIDEN ir<7> = vp.strided.load ir<arrayidx.3>, vp<4>, ir<...>	VPLAN: cost of 0 for recipe: EPI-WIDEN ir<7> = vp.strided.load ir<...>
VPLAN: cost of 2 for recipe: EPI-WIDEN ir<8> = vp.strided.load ir<arrayidx.4>, vp<4>, ir<...>	VPLAN: cost of 0 for recipe: EPI-WIDEN ir<8> = vp.strided.load ir<...>
VPLAN: cost of 2 for recipe: EPI-WIDEN ir<9> = vp.strided.load ir<arrayidx.5>, vp<4>, ir<...>	VPLAN: cost of 0 for recipe: EPI-WIDEN ir<9> = vp.strided.load ir<...>
VPLAN: cost of 2 for recipe: EPI-WIDEN ir<10> = vp.strided.load ir<arrayidx.6>, vp<4>, ir<...>	VPLAN: cost of 0 for recipe: EPI-WIDEN ir<10> = vp.strided.load ir<...>
VPLAN: cost of 2 for recipe: EPI-WIDEN ir<11> = vp.strided.load ir<arrayidx.7>, vp<4>, ir<...>	VPLAN: cost of 0 for recipe: EPI-WIDEN ir<11> = vp.strided.load ir<...>
VPLAN: cost of 2 for recipe: EPI-WIDEN ir<12> = vp.strided.load ir<arrayidx.8>, vp<4>, ir<...>	VPLAN: cost of 0 for recipe: EPI-WIDEN ir<12> = vp.strided.load ir<...>
VPLAN: cost of 33 for VF vscale x 1 (assuming a minimum vscale of 2).	VPLAN: cost of 26 for VF vscale x 1 (assuming a minimum vscale of 2).

Results on SPEC2017 on Spacemix-x60

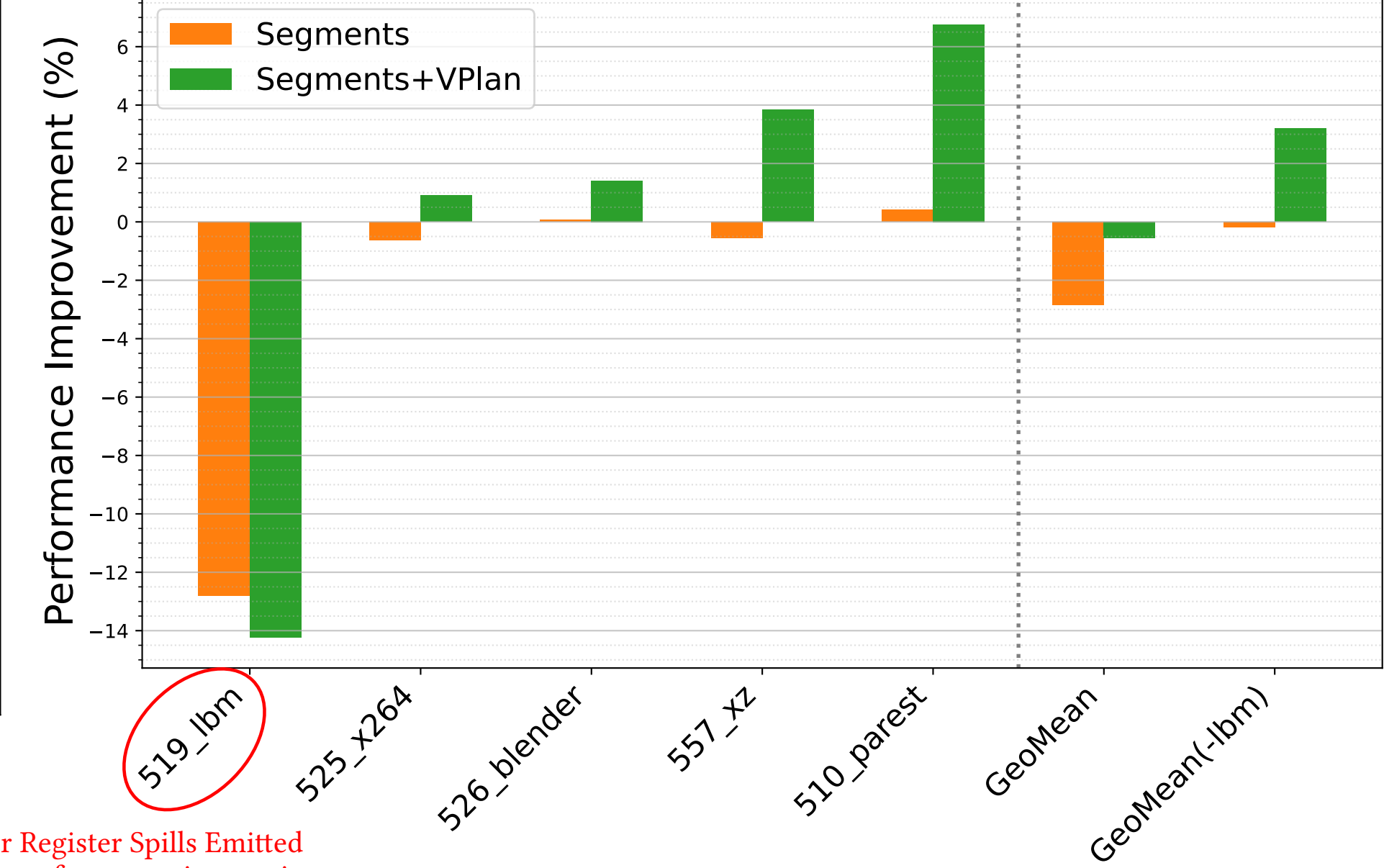
Static Instruction Count: Segment Instructions by Benchmark



Dynamic Instruction Count: Segment Instructions



Performance Improvements



Vector Register Spills Emitted because of segment instructions

Next Steps

- 1) Compute the vector register pressure with and without segment instructions. So, before we emit any segment instruction, we can approximate if forming segments will increase the register pressure above a threshold that is likely to cause additional vector spilling.
- 2) Extend the analysis to forming groups of gathers and scatters that can be replaced by indexed segment instructions.
- 3) Explore more smart decision making regarding when to emit segment instructions and of what size.

This project is co-funded by the Ministerio para la Transformación Digital y de la Función Pública, within the framework of the Plan de Recuperación Transformación y Resiliencia, and by the European Union - NextGenerationEU. The views and opinions expressed are solely those of the author(s) and do not necessarily reflect those of the European Union. Neither the European Union nor the European Commission can be held responsible for them.