# Unlocking 15% More Performance: A Case Study in LLVM Optimization for RISC-V

Phd Mikhail R. Gadelha

Igalia

**16%**

**Unlocking ~~15%~~ More Performance: A Case Study in LLVM Optimization for RISC-V**

Phd Mikhail R. Gadelha

Igalia

# Up to 16% faster SPEC CPU® 2017 on SpacemiT-X60

- This RISE project was a focused, ten-month effort to optimize the LLVM compiler for RISC-V.

- This RISE project was a focused, ten-month effort to optimize the LLVM compiler for RISC-V.

- Our target board was the Banana Pi BPI-F3 with a SpacemiT-X60 8-core RISC-V processor:
  - In-order processor.
  - Supports the RVA22U64 Profile and 256-bit RVV 1.0 standard.

- This RISE project was a focused, ten-month effort to optimize the LLVM compiler for RISC-V.

- Our target board was the Banana Pi BPI-F3 with a SpacemiT-X60 8-core RISC-V processor:
  - In-order processor.
  - Supports the RVA22U64 Profile and 256-bit RVV 1.0 standard.

- Our goal: to close the performance gap between LLVM and the GCC compiler.

- This RISE project was a focused, ten-month effort to optimize the LLVM compiler for RISC-V.

- Our target board was the Banana Pi BPI-F3 with a SpacemiT-X60 8-core RISC-V processor:
  - In-order processor.
  - Supports the RVA22U64 Profile and 256-bit RVV 1.0 standard.

- Our goal: to close the performance gap between LLVM and the GCC compiler.

- Our result: individual contributions boosted performance by up to 16% on SPEC CPU® 2017 benchmarks.

# The Project

- Prior to this work, a clear performance gap existed between code generated by LLVM and GCC for RISC-V.

# The Project

- Prior to this work, a clear performance gap existed between code generated by LLVM and GCC for RISC-V.

- There is no single solution to close the gap, as improvements and regressions occur daily within the codebase.

- Prior to this work, a clear performance gap existed between code generated by LLVM and GCC for RISC-V.

- There is no single solution to close the gap, as improvements and regressions occur daily within the codebase.

- This presentation will focus on our three main contributions to help close the gap:
  - Introducing a scheduling model for the SpacemiT-X60.
  - Improvements to vectorization across calls.
  - Register Allocation with IPRA Support for RISC-V.

# Our Contributions

(major contributions first)

# SpacemiT-X60 Scheduling Model

# SpacemiT-X60 Scheduling Model

- Instruction Scheduling == Performance.
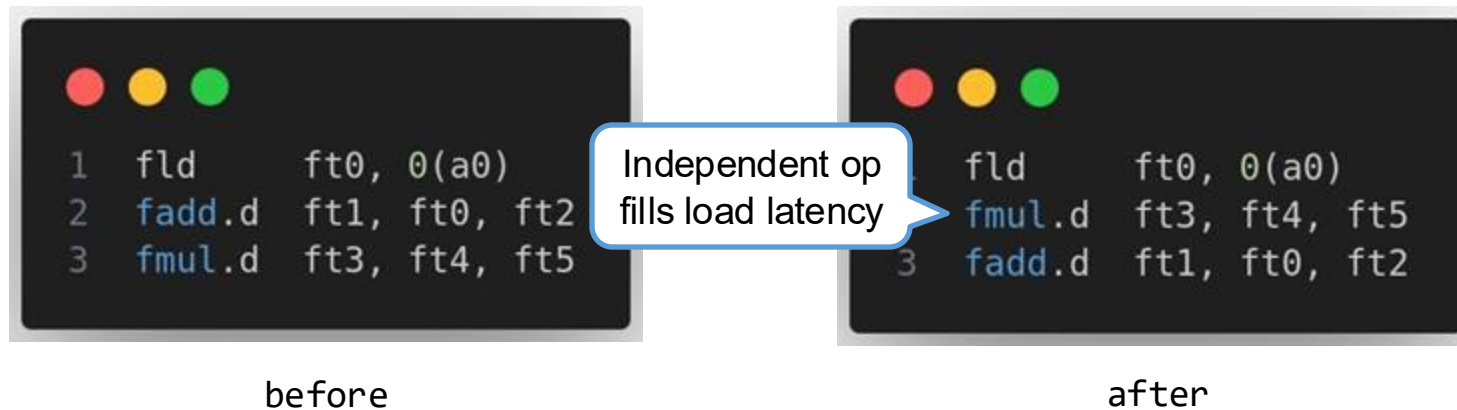- Wrong latencies/resources → compiler makes poor choices.



before



after

# SpacemiT-X60 Scheduling Model

- Instruction Scheduling == Performance.
- Wrong latencies/resources → compiler makes poor choices.



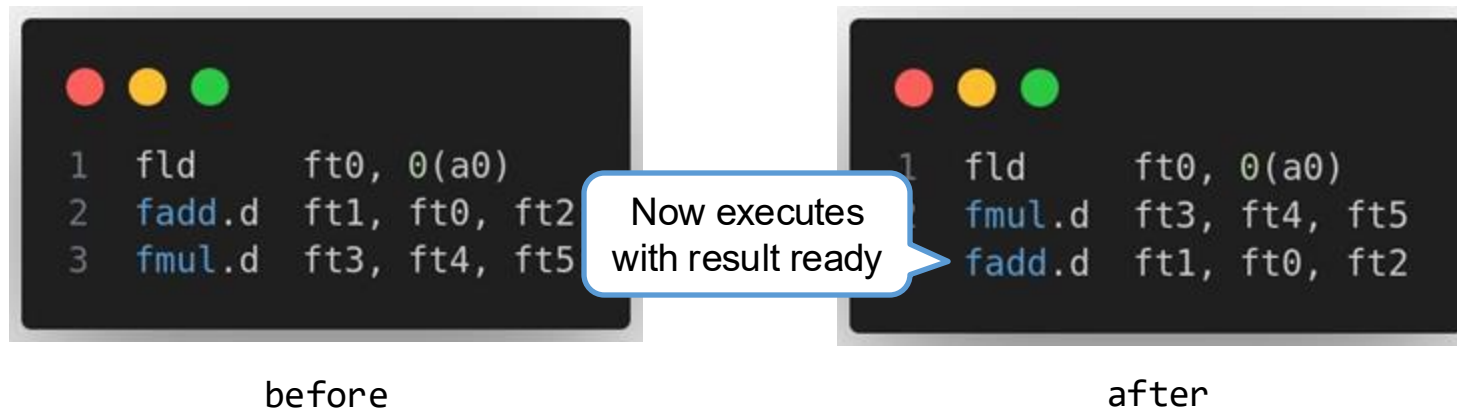before



after

# SpacemiT-X60 Scheduling Model

- Instruction Scheduling == Performance.
- Wrong latencies/resources → compiler makes poor choices.



before



after

# SpacemiT-X60 Scheduling Model

- Instruction Scheduling == Performance.
- Wrong latencies/resources → compiler makes poor choices.



before



after

# SpacemiT-X60 Scheduling Model

- Instruction Scheduling == Performance.
- Wrong latencies/resources → compiler makes poor choices.



before

after

# SpacemiT-X60 Scheduling Model

- Instruction Scheduling == Performance.
- Wrong latencies/resources → compiler makes poor choices.



before

after

# SpacemiT-X60 Scheduling Model

- Instruction Scheduling == Performance.
- Wrong latencies/resources → compiler makes poor choices.



before

after

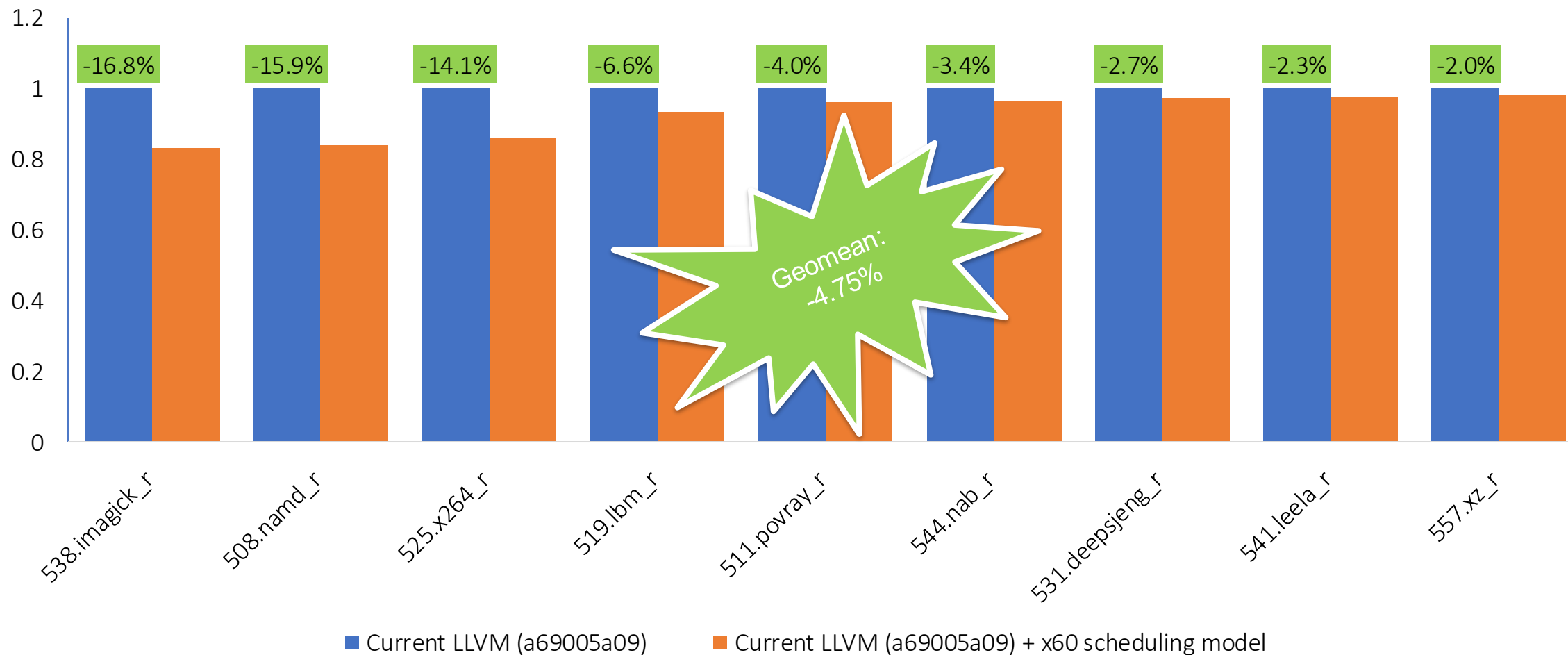- We built custom microbenchmarks to measure instruction latencies.

# How We Got the Numbers

- We built custom microbenchmarks to measure instruction latencies.

- Most of instruction throughput data available at https://camel-cdr.github.io/rvv-bench-results/bpi_f3/index.html.

# How We Got the Numbers

- We built custom microbenchmarks to measure instruction latencies.

- Most of instruction throughput data available at [https://camel-cdr.github.io/rvv-bench-results/bpi_f3/index.html](https://camel-cdr.github.io/rvv-bench-results/bpi_f3/index.html).

- It's RISC but:
  - 201 scalar instructions.
  - 82 floating-point instructions.
  - 9185 RVV instructions (because of the combination of different LMULs and SEWs).
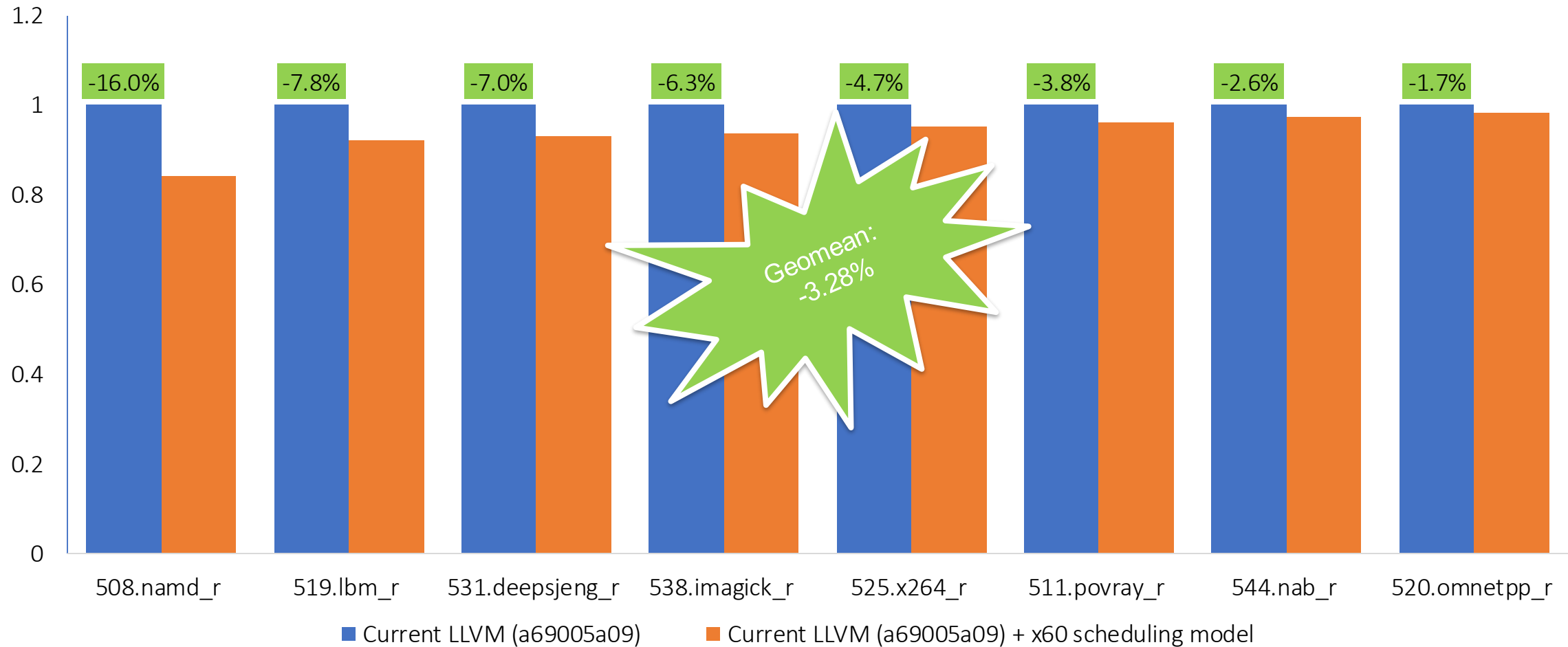
# RVA22U64 SPEC exec time, O3+LTO+mcpu=spacemit-x60



Bar chart of SPEC benchmark execution times comparing "Current LLVM (a69005a09)" (blue) vs "Current LLVM (a69005a09) + x60 scheduling model" (orange), normalized to 1.0.

| Benchmark | Improvement |
|-----------|-------------|
| 538.imagick_r | -16.8% |
| 508.namd_r | -15.9% |
| 525.x264_r | -14.1% |
| 519.lbm_r | -6.6% |
| 511.povray_r | -4.0% |
| 544.nab_r | -3.4% |
| 531.deepsjeng_r | -2.7% |
| 541.leela_r | -2.3% |
| 557.xz_r | -2.0% |

Legend: ■ Current LLVM (a69005a09)　■ Current LLVM (a69005a09) + x60 scheduling model

# RVA22U64_V SPEC exec time, O3+LTO+mcpu=spacemit-x60



| | 508.namd_r | 519.lbm_r | 531.deepsjeng_r | 538.imagick_r | 525.x264_r | 511.povray_r | 544.nab_r | 520.omnetpp_r |
|---|---|---|---|---|---|---|---|---|
| | -16.0% | -7.8% | -7.0% | -6.3% | -4.7% | -3.8% | -2.6% | -1.7% |

■ Current LLVM (a69005a09)   ■ Current LLVM (a69005a09) + x60 scheduling model

# RVA22U64_V SPEC exec time, O3+LTO+mcpu=spacemit-x60



Geomean: -3.28%

| Benchmark | Current LLVM (a69005a09) | Current LLVM (a69005a09) + x60 scheduling model |
|---|---|---|
| 508.namd_r | | -16.0% |
| 519.lbm_r | | -7.8% |
| 531.deepsjeng_r | | -7.0% |
| 538.imagick_r | | -6.3% |
| 525.x264_r | | -4.7% |
| 511.povray_r | | -3.8% |
| 544.nab_r | | -2.6% |
| 520.omnetpp_r | | -1.7% |

- Scheduling nearly eliminated the gap between scalar and vector configs.

# RVA22U64 vs RVA22U64_V, O3+LTO+mcpu=spacemit-x60



- Scheduling nearly eliminated the gap between scalar and vector configs.

- On in-order processors like X60, scheduling is critical; on out-of-order, impact would be smaller and vectorization more decisive.

# Improvements to Vectorization Across Calls

# Improvements to Vectorization Across Calls

- Initial surprise: vectorized code sometimes underperformed scalar.

# Improvements to Vectorization Across Calls

- Initial surprise: vectorized code sometimes underperformed scalar.

- Root cause: poor cost modeling and suboptimal spill behavior.

# Improvements to Vectorization Across Calls

- Initial surprise: vectorized code sometimes underperformed scalar.

- Root cause: poor cost modeling and suboptimal spill behavior.

- The extra cycles were due to register spilling, particularly around function call boundaries.
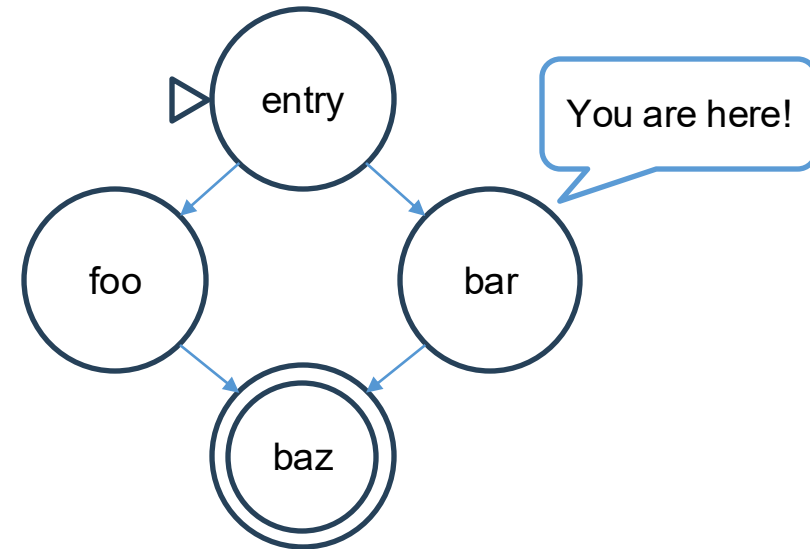
```
 1  define void @f(i1 %c, ptr %p, ptr %q) {
 2  entry:
 3    %x0 = load i64, ptr %p
 4    %p1 = getelementptr i64, ptr %p, i64 1
 5    %x1 = load i64, ptr %p1
 6    br i1 %c, label %foo, label %bar
 7  foo:
 8    call void @g()
 9    br label %baz
10  bar:
11    call void @g()
12    br label %baz
13  baz:
14    store i64 %x0, ptr %q
15    %q1 = getelementptr i64, ptr %q, i64 1
16    store i64 %x1, ptr %q1
17    ret void
18  }
```
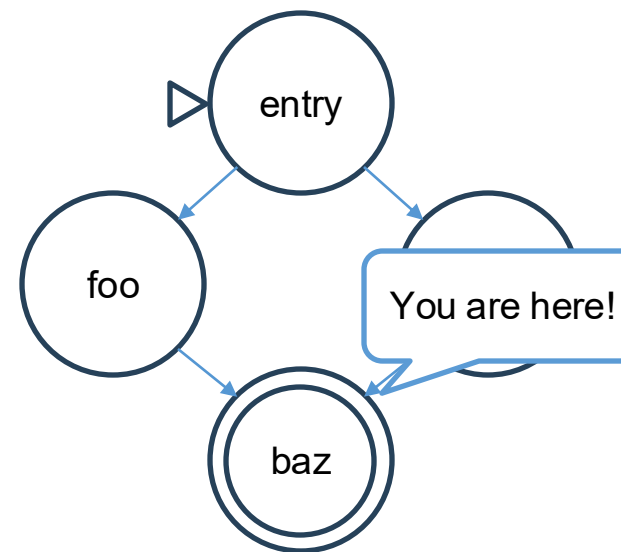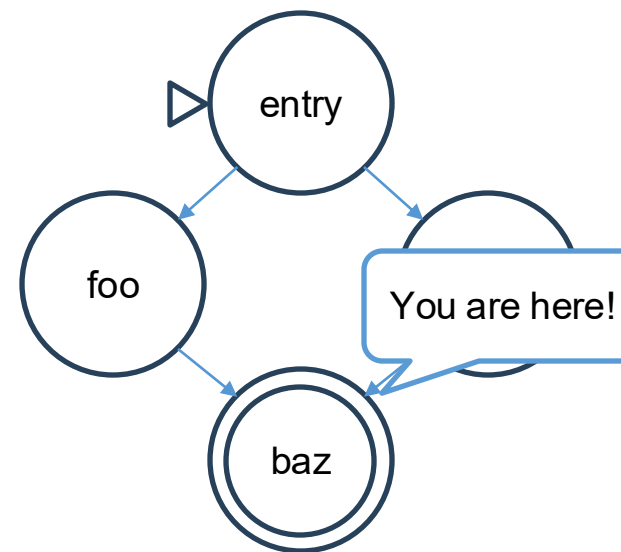
# The 544.nab_r case

# Fixing Real Bugs

- Since we are storing and loading from continuous regions, the accesses can be vectorized.

- Since we are storing and loading from continuous regions, the accesses can be vectorized.

- We found that the SLP Vectorizer was aggressively vectorizing regions without properly accounting for the cost of spilling vector registers across calls.
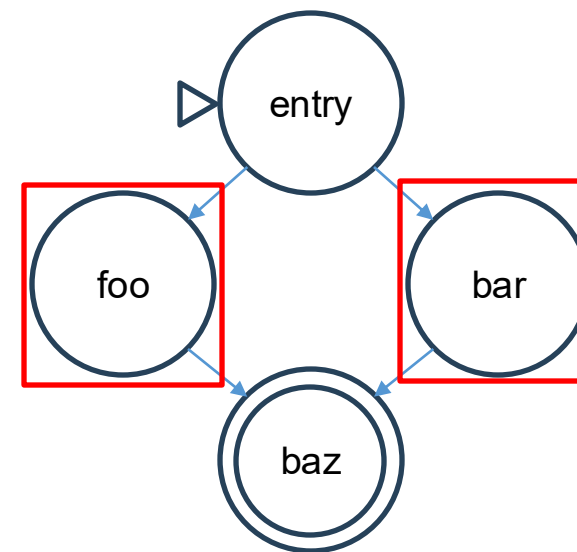
- Since we are storing and loading from continuous regions, the accesses can be vectorized.

- We found that the SLP Vectorizer was aggressively vectorizing regions without properly accounting for the cost of spilling vector registers across calls.

- Previously, the SLP vectorizer only analyzed the entry and baz blocks, ignoring foo and bar entirely.

# The 544.nab_r case

- To address the issue, we first proposed a patch which modified the SLP vectorizer to properly walk through all basic blocks when analyzing cost.

# Addressing the Issue

- To address the issue, we first proposed a patch which modified the SLP vectorizer to properly walk through all basic blocks when analyzing cost.

- Promising results: execution time dropped by **9.92%** in 544.nab_r.

# Addressing the Issue

- To address the issue, we first proposed a patch which modified the SLP vectorizer to properly walk through all basic blocks when analyzing cost.

- Promising results: execution time dropped by **9.92%** in 544.nab_r.

- But with a major drawback: **+6.9%** increase compilation time in 502.gcc_r.
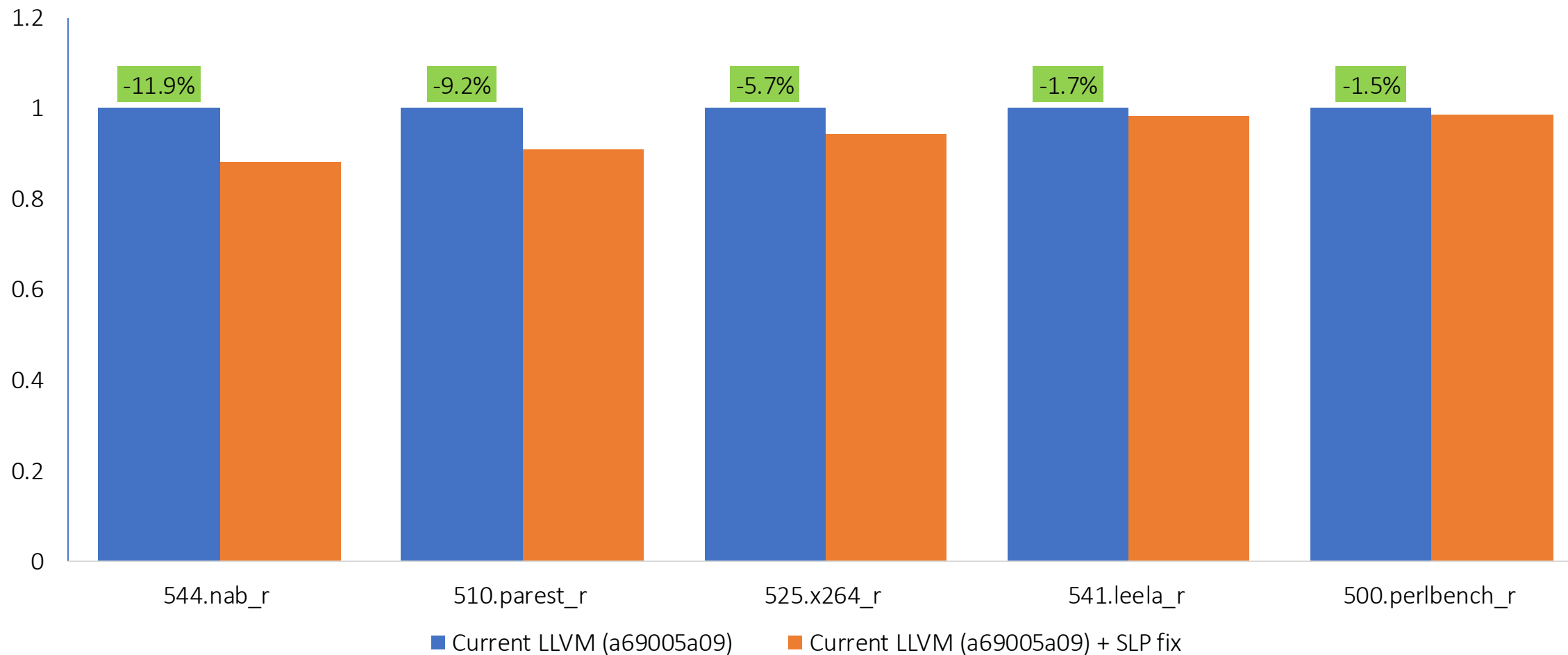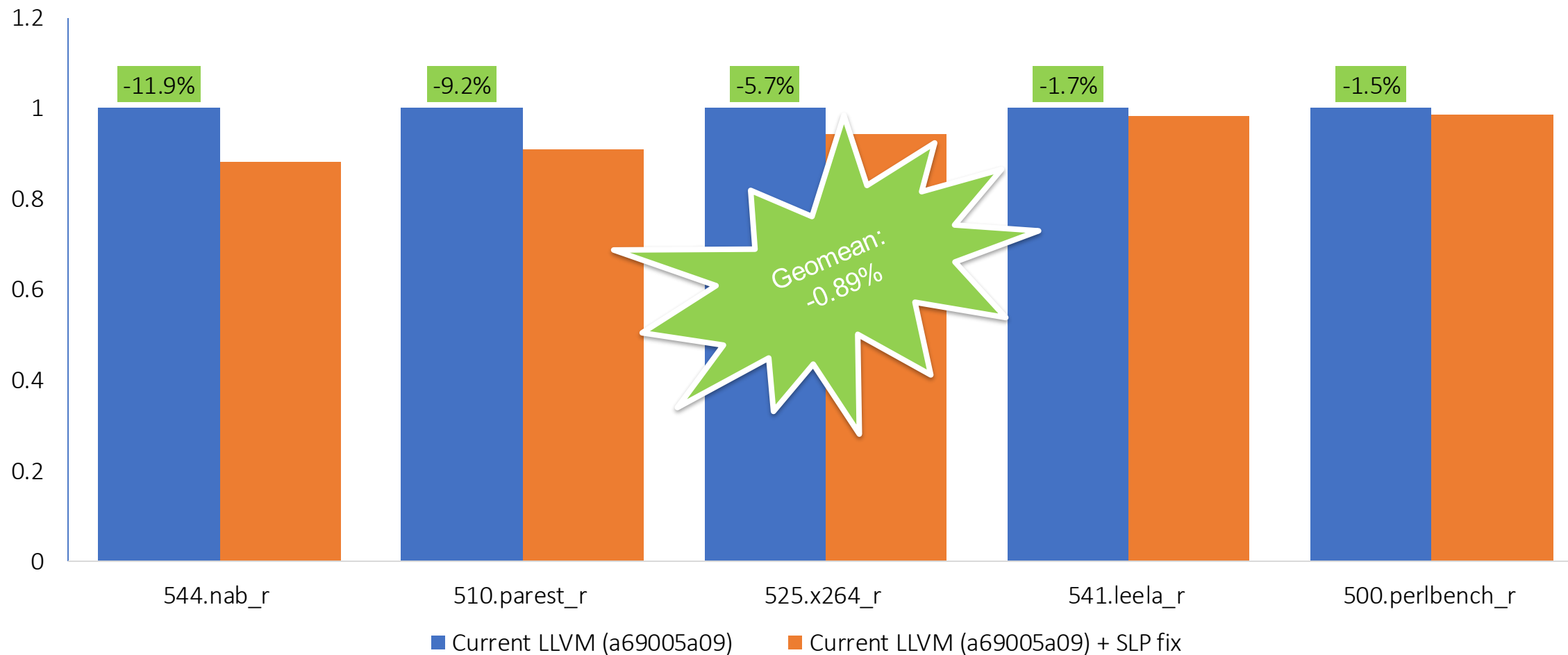
# Addressing the Issue

- To address the issue, we first proposed a patch which modified the SLP vectorizer to properly walk through all basic blocks when analyzing cost.

- Promising results: execution time dropped by **9.92%** in 544.nab_r.

- But with a major drawback: **+6.9%** increase compilation time in 502.gcc_r.

- Following discussions with the community, Alexey Bataev (SLP Vectorizer code owner) proposed and landed refined solution.

# RVA22U64_V SPEC exec time, O3+LTO, SLP fix

# RVA22U64_V SPEC exec time, O3+LTO, SLP fix

# IPRA (Inter-Procedural Register Allocation) Support

- Function calls can often spill (when you store a register to the stack) more registers than necessary → wasted cycles on every call.
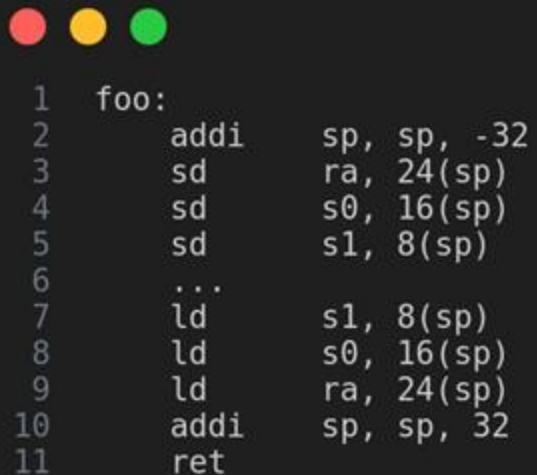
# IPRA (Inter-Procedural Register Allocation) Support

- Function calls can often spill (when you store a register to the stack) more registers than necessary → wasted cycles on every call.

- IPRA: caller/callee register use is tracked across the functions, by eliminating unnecessary save/restore sequences.

# What IPRA Brings



before



after

# What IPRA Brings



before

after

# What IPRA Brings



before

after

- Reduction in register pressure, shorter prologue/epilogue code.

- Reduction in register pressure, shorter prologue/epilogue code.
- SPEC benchmarks showed measurable improvements (small but consistent).

# Impact of IPRA

- Reduction in register pressure, shorter prologue/epilogue code.

- SPEC benchmarks showed measurable improvements (small but consistent).

- Unfortunately, it can't be enabled by default: IPRA is not enabled by default due to a bug (described in issue 119556), however, it does not affect the SPEC benchmarks.

# RVA22U64 SPEC exec time, O3+LTO+IPRA

# RVA22U64_V SPEC exec time, O3+LTO+IPRA

# Conclusions

# Faster RISC-V Today

- We contributed with:
  - Scheduling: largest wins, especially for scalar-heavy code.
  - Vectorization: enabled smarter spilling cost calculations.
  - IPRA: smaller but consistent improvements across workloads.

- We contributed with:
  - Scheduling: largest wins, especially for scalar-heavy code.
  - Vectorization: enabled smarter spilling cost calculations.
  - IPRA: smaller but consistent improvements across workloads.
- Almost all changes are upstream benefiting everyone. Under review:
  - https://github.com/llvm/llvm-project/pull/150618
  - https://github.com/llvm/llvm-project/pull/150644
  - https://github.com/llvm/llvm-project/pull/152557
  - https://github.com/llvm/llvm-project/pull/152737
  - https://github.com/llvm/llvm-project/pull/152738

# What We Learned Along the Way

- Scheduling is critical for performance.
  - No scheduling model → LLVM pessimises the final code.
  - We should likely adopt some scheduling model as default, like other targets do.
  - Should we make the X60 scheduling model default for in-order RISC-V processors?

- Scheduling is critical for performance.
  - No scheduling model → LLVM pessimises the final code.
  - We should likely adopt some scheduling model as default, like other targets do.
  - Should we make the X60 scheduling model default for in-order RISC-V processors?
- Many contributions don't have immediate benchmark impact.

- Scheduling is critical for performance.
  - No scheduling model → LLVM pessimizes the final code.
  - We should likely adopt some scheduling model as default, like other targets do.
  - Should we make the X60 scheduling model default for in-order RISC-V processors?
- Many contributions don't have immediate benchmark impact.
- Vectorization needs careful tuning to avoid regressions.

# Did we close the performance gap between LLVM and the GCC compiler?

- **Note it's not a direct apples-to-apples comparison.**

- **Note it's not a direct apples-to-apples comparison**.
- The code was compiled with the same RISC-V extensions enabled.

- **Note it's not a direct apples-to-apples comparison**.

- The code was compiled with the same RISC-V extensions enabled.

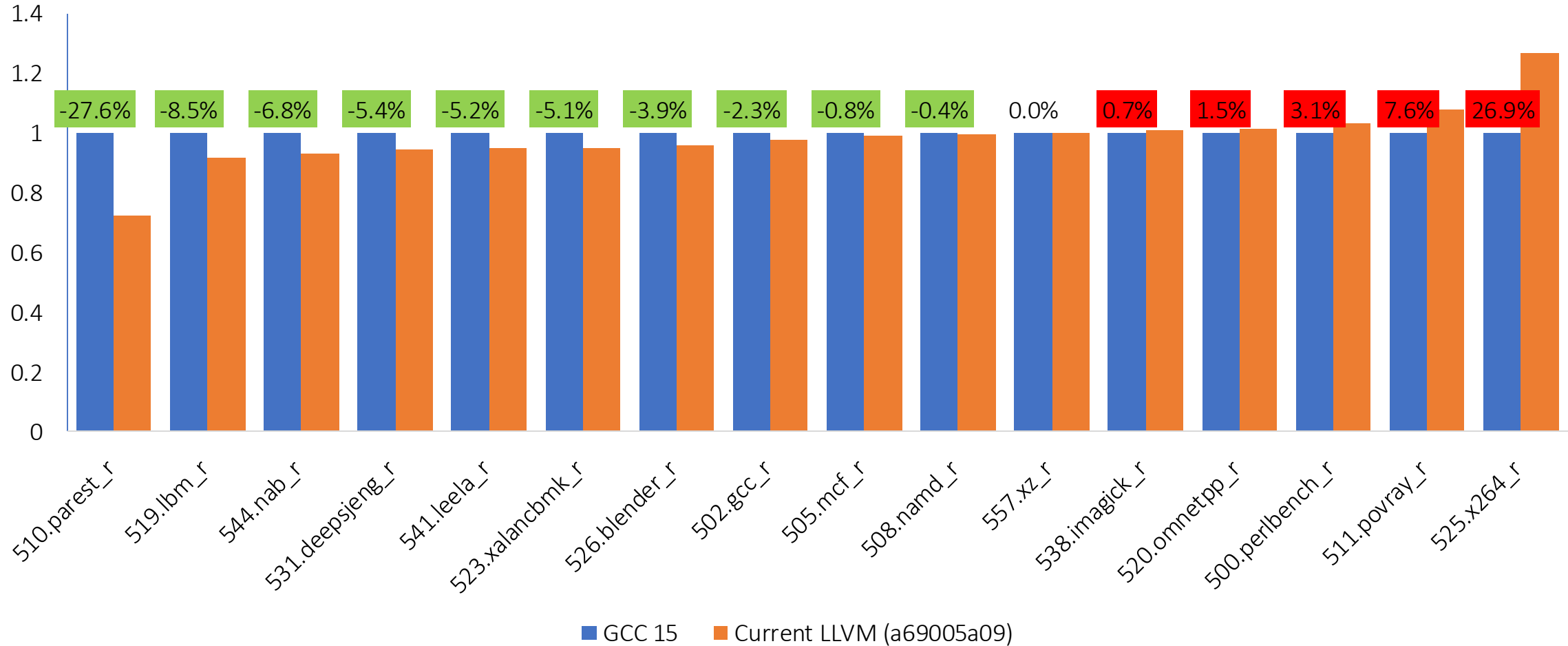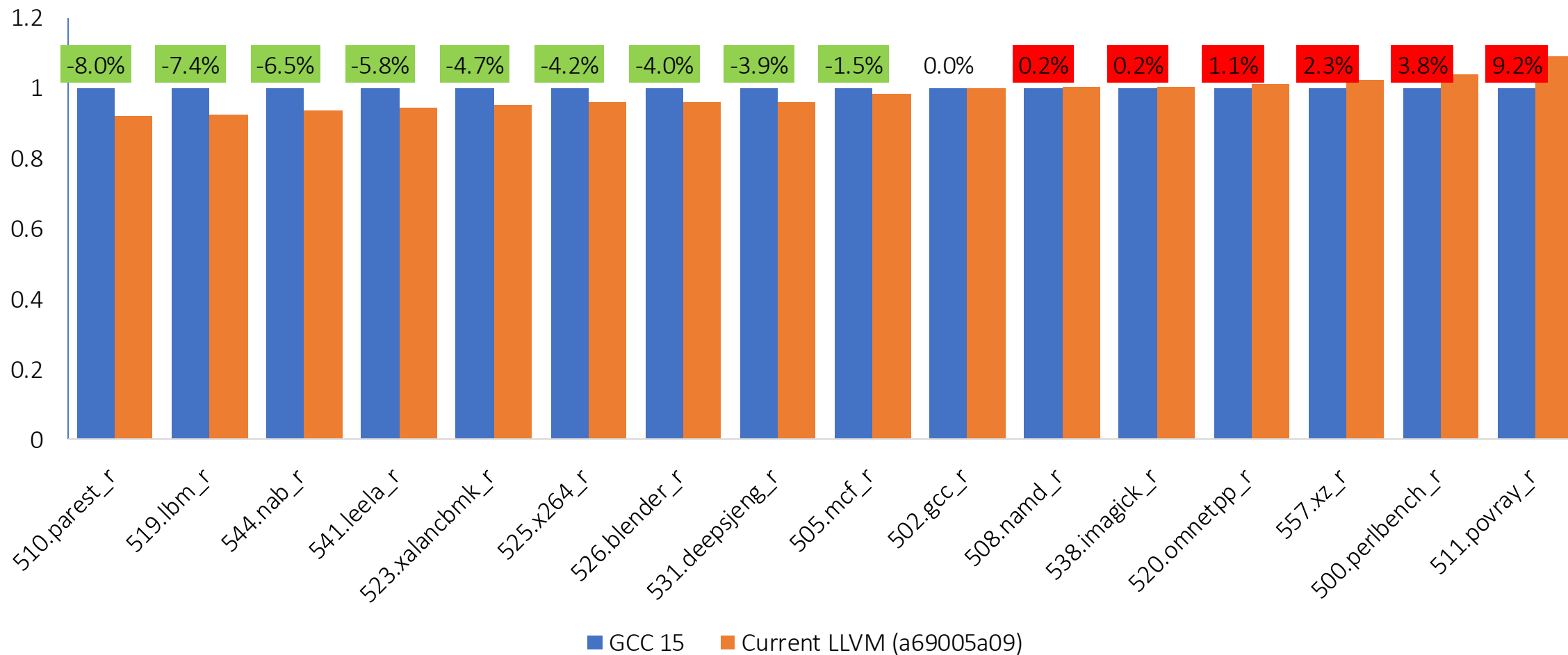- But GCC doesn't have X60-specific scheduling latencies, while LLVM does.

- **Note it's not a direct apples-to-apples comparison**.

- The code was compiled with the same RISC-V extensions enabled.

- But GCC doesn't have X60-specific scheduling latencies, while LLVM does.

- Still useful to show relative progress and identify where LLVM has caught up.

# RVA22U64_V SPEC execution time, GCC vs LLVM

# Thanks!

- This work at Igalia was made possible thanks to support from RISE, under Project RP009.

- Thanks to my Igalia colleagues for discussions, and feedback.

- And to the LLVM RISC-V community for reviews and for getting patches upstream quickly.

- Accidental Dataflow Analysis: Extending the RISC-V VL Optimizer @ 2025 EuroLLVM by Luke Lau: https://www.youtube.com/watch?v=bkOwPr36SrQ

- Improvements to RISC-V Vector code generation in LLVM @ 2025 RISC-V Summit Europe by Luke Lau and Alex Bradbury: https://www.youtube.com/watch?v=0NjugW7FF48

- RISC-V nightly performance testing of top-of-tree GCC and clang:

  https://cc-perf.igalia.com/