

Bridge Over Troubled Water: Offloading OpenMP Regions to AI Compiler via StableHLO/XLA

Problem Formulation: The Gap

To meet the intense computational demands of AI, manufacturers have developed specialized hardware, such as NVIDIA's Tensor Cores, to accelerate matrix operations.

However, despite sharing workload characteristics with AI, scientific computing has yet to fully leverage these advancements (or requires painful manual work and expertise).

This research aims to bridge the gap between Fortran-based scientific computing and AI-specialized hardware by introducing a novel compilation path that utilizes the XLA compiler.

Research Scope: The `workdistribute` Directive

Fortran natively supports array operations as first-class language constructs. For example, given arrays X , Y , Z of the same size n , and a scalar a . AXPY operation could be expressed in:

```
1 Z = a * X + Y
```

Listing 1. array like syntax

Traditionally, achieving parallelism required explicit do-loops rewriting within OpenMP:

```
1 !$omp target teams distribute parallel do
2 do i = 1, n
3   Z(i) = a * X(i) + Y(i)
4 end do
```

Listing 2. OpenMP syntax

However, OpenMP 6.0 standard introduced the `workdistribute` directive, enabling automatic parallelism for array syntax[1]. Notice the following code is exactly the same with Listing 1 except for the OpenMP directives:

```
1 !$omp target teams workdistribute
2 Z = a * X + Y
3 !$omp end target teams workdistribute
```

Listing 3. workdistribute syntax

Taking Advantage of AI Compiler: The Era of Configurable Compiler Infrastructure

The modular design of the LLVM/MLIR framework [2], which powers the flang compiler, provides the flexibility to decouple the frontend from the backend.

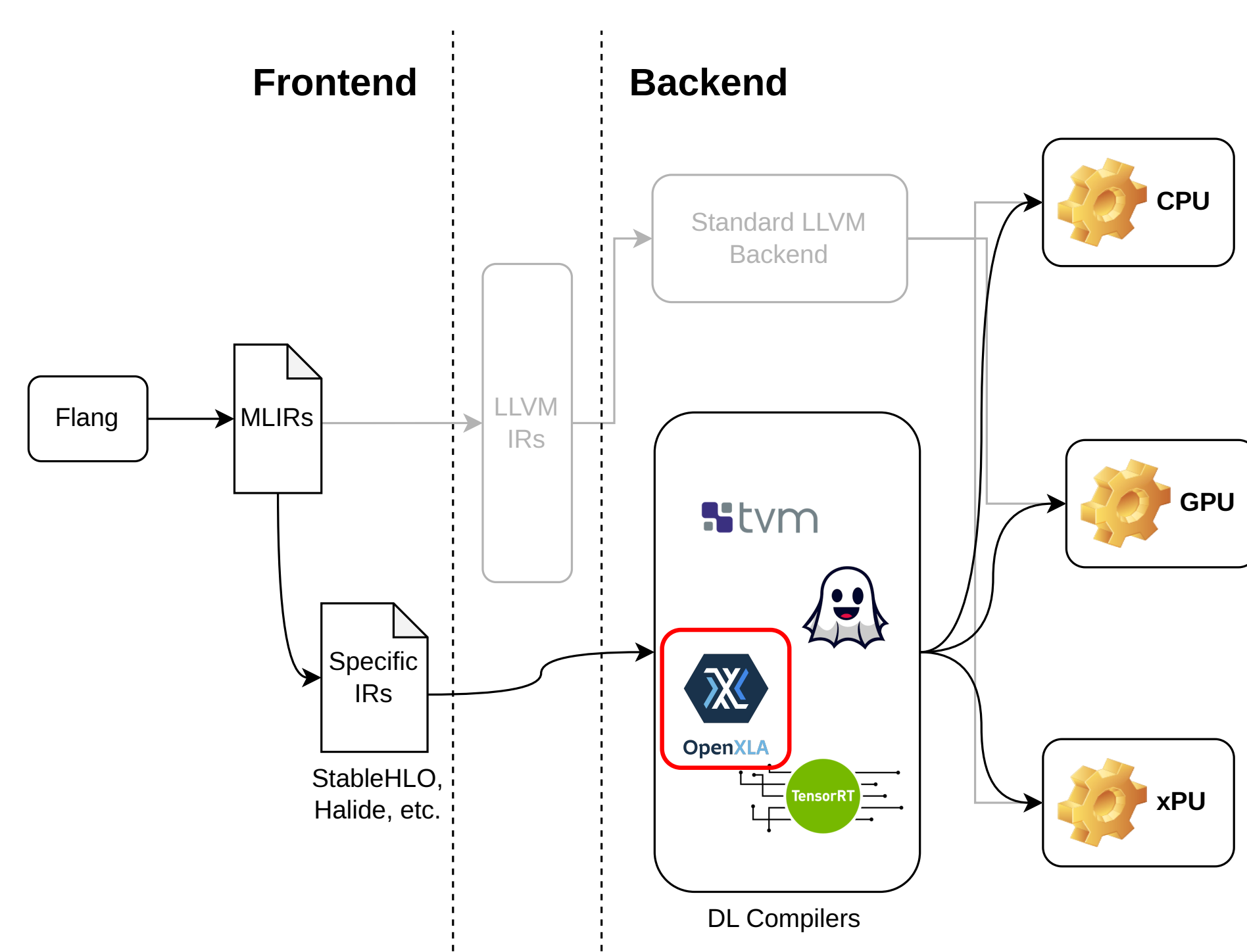


Figure 1. The Big Picture

This allows us to replace traditional backends with specialized AI compilers that already offer highly optimized code generation for GPUs and other accelerators.

System Design: Intercept Modern Flang Compilation and OpenMP Runtime

Our custom LLVM/MLIR pass, `WorkdistributeJITPass`, extracts `workdistribute` regions IR. They are then transformed into StableHLO and compiled JIT into a kernel binary for the target accelerator:

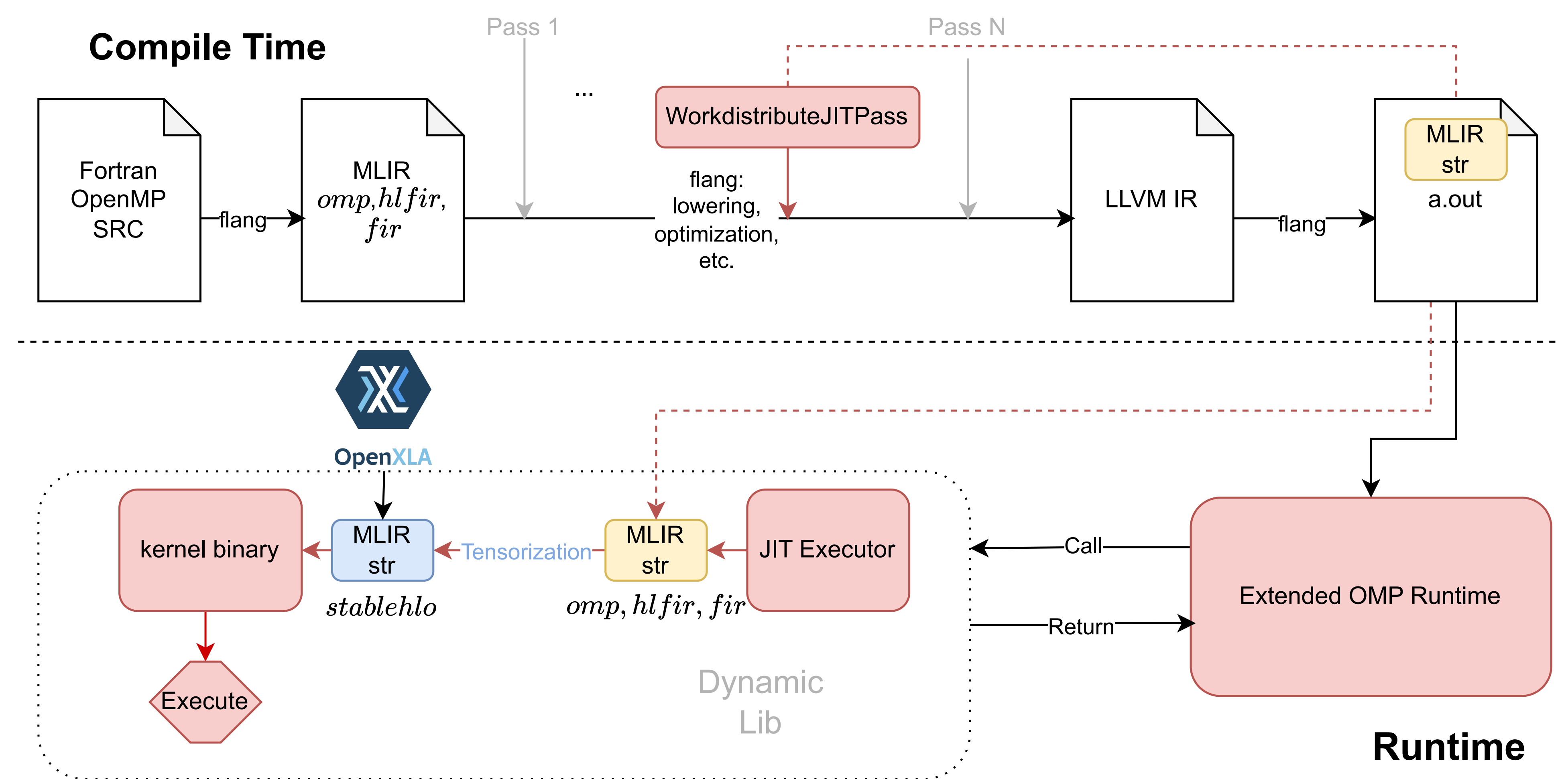


Figure 2. The Compilation Architecture

Tensorization: Translate from Fortran IRs to StableHLO

Memory related operations in Fortran IRs like `load`, `designate`, `apply` are transformed to inputs and outputs in StableHLO, leaving only computations of tensors.

```
1 omp.workdistribute {
2   %a = fir.load %148#0: f32
3   %tmp = hlfir.elemental %149: %Xshape->expr<10xf32> {
4     ^bb0(%i: index):
5       %154 = hlfir.designate %X[%i]: ref<f32>
6       %155 = fir.load %154: f32
7       %156 = arith.mulf %151, %a: f32
8       hlfir.yield_element %156: f32
9     }
10  %153 = hlfir.elemental %149: %Yshape->expr<10xf32> {
11    ...
12    %157 = arith.addf %154, %156: f32
13    ...
14  }
15  hlfir.assign %153 to %Z: expr<10xf32>, ref<10xf32>
16  ...
17 }
```

Listing 4. AXPY function before *Tensorization*: omp hlfir fir

```
1 func.func @kernel(%Z: tensor<f32>, %a:
2   tensor<f32>, %X: tensor<10xf32>, %Y:
3   tensor<10xf32>) {
4   %0 = stablehlo.broadcast_in_dim %a:
5     (tensor<f32>) -> tensor<10xf32>
6   %1 = stablehlo.multiply %X, %0:
7     tensor<10xf32>
8   %2 = stablehlo.add %1, %Y: tensor<10
9     xf32>
10  %return = %2, %a, %X, %Y: tensor<10xf32
11    >, tensor<f32>, tensor<10xf32>, tensor
12    <10xf32>
13 }
```

Listing 5. AXPY function after *Tensorization*: stableHLO

Preliminary Benchmark

We evaluated the above AXPY operations kernel execution time comparing **traditional Fortran OpenMP kernels** and **XLA kernels compiled from handwritten StableHLO** on an Nvidia A100 GPU:

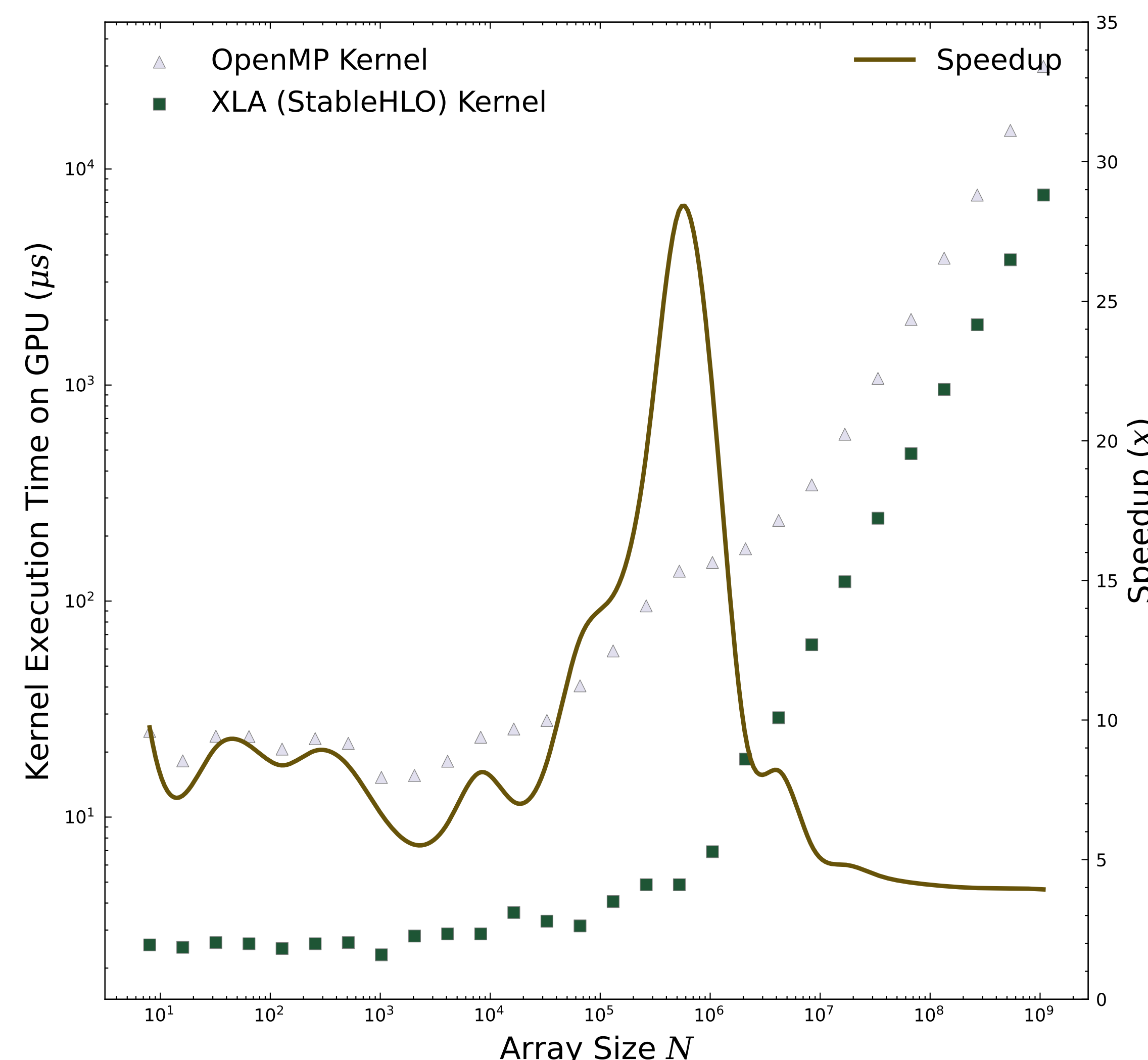


Figure 3. OpenMP Kernel vs. XLA kernel Performance Comparison

Summary and Future Research

XLA kernels demonstrate significant speedup over the native OpenMP kernels, though total performance includes runtime overheads not fully captured by GPU profiling.

Our other Matrix multiplication tests also confirm the automatic generation of NVIDIA CUTLASS tensor instructions.

Future work includes further system integration and evaluation, mitigating JIT overhead and broadening syntax support.

- [1] Ivanov et al. Automatic Parallelization and OpenMP Offloading of Fortran Array Notation. Cham: Springer Nature Switzerland, 2024. isbn: 978-3-031-72567-8.
- [2] Lattner and Adve. LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation. CGO '04. Palo Alto, California: IEEE Computer Society, 2004. isbn: 0769521029.