

Automatically Generating High-performance Matrix Multiplication Kernels on the Latest Sunway Processor

Xiaohan Tao **Yu Zhu**
Boyang Wang Jinlong Xu Jianmin Pang Jie Zhao

State Key Laboratory of Mathematical Engineering and Advanced Computing
i.zhuyu@126.com

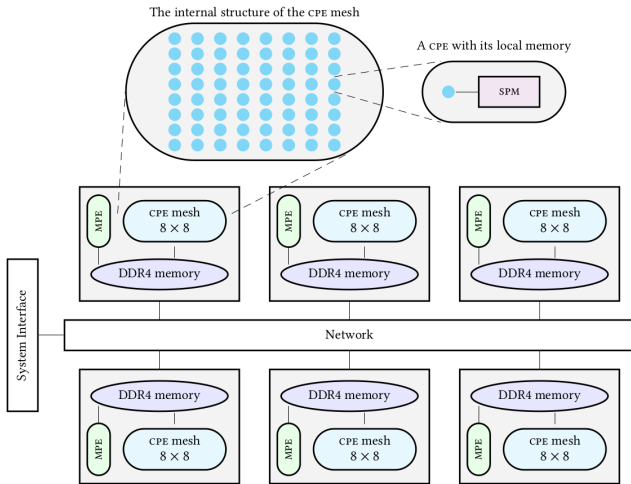
51st International Conference on Parallel Processing
August 31,2022

Background and Motivation

- The Sunway TaiHuLight supercomputer
 - Run high-performance computing applications, which involve massive linear algebra operations
- GEMM: Essential BLAS kernel for AI/ML
- Problem and Motivation
 - Manual efforts incurring high cost
 - Compilation approach is still missing
- Existing methods
 - Have to expose the architecture information
 - Only generated the innermost loop
- Our Method
 - An **Automatic Code Generation Approach** for GEMM on SW26010Pro

Background and Motivation

The Sunway Architecture



Background and Motivation

Polyhedral Compilation

Polyhedral Model

A mathematical abstraction to reason about loop transformations and memory optimizations using integer sets and affine relations.

Background and Motivation

Polyhedral Compilation

Polyhedral Model

A mathematical abstraction to reason about loop transformations and memory optimizations using integer sets and affine relations.

for i in $[0, M)$ and j in $[0, N)$ and k in $[0, K)$
 $C[i, j] = C[i, j] + A[i, k] * B[k, j] /* S_1 */$

DOMAIN: $\{S_1(i, j, k) : 0 \leq i < M \wedge 0 \leq j < N \wedge 0 \leq k < K\}$
BAND: $[S_1(i, j, k) \rightarrow (i, j, k)]$

(a) A 3D loop nest of GEMM code.

(b) The initial schedule tree.

DOMAIN: $\{S_1(i, j, k) : 0 \leq i < M \wedge 0 \leq j < N \wedge 0 \leq k < K\}$
BAND: $[S_1(i, j, k) \rightarrow (blockIdx.y, blockIdx.x, \lfloor \frac{k}{32} \rfloor)]$
EXTENSION: $[(d_0, d_1, d_2) \rightarrow readA[d_3, d_4]; (d_0, d_1, d_2) \rightarrow readB[d_3, d_4]]$
BAND: $[S_1(i, j, k) \rightarrow (threadIdx.y, threadIdx.x, k - 32 \lfloor \frac{k}{32} \rfloor)]$
SEQUENCE:
FILTER: $\{readA[d_3, d_4]\}$
FILTER: $\{readB[d_3, d_4]\}$
FILTER: $\{S_1(i, j, k)\}$

(c) The schedule tree with shared memory promotion statements.

Compute Decomposition

Tiling All Dimensions

tile size selection issue has not yet been modeled by isl or other polyhedral tools.

Our Objective

match the shape configuration of the **assembly micro kernel**.

Compute Decomposition

Tiling All Dimensions

tile size selection issue has not yet been modeled by isl or other polyhedral tools.

Our Objective

match the shape configuration of the **assembly micro kernel**.

$$\text{DOMAIN: } \{S_1(i, j, k) : 0 \leq i < M \wedge 0 \leq j < N \wedge 0 \leq k < K\}$$

$$\text{BAND: } [S_1(i, j, k) \rightarrow (\lfloor \frac{i}{64} \rfloor, \lfloor \frac{j}{64} \rfloor, \lfloor \frac{k}{32} \rfloor)]$$

$$\text{BAND: } [S_1(i, j, k) \rightarrow (i - 64 \lfloor \frac{i}{64} \rfloor, j - 64 \lfloor \frac{j}{64} \rfloor, k - 32 \lfloor \frac{k}{32} \rfloor)]$$

(a) The schedule tree after tiling.

$$\text{DOMAIN: } \{S_1(i, j, k) : 0 \leq i < M \wedge 0 \leq j < N \wedge 0 \leq k < K\}$$

$$\text{BAND: } [S_1(i, j, k) \rightarrow (Rid, Cid, \lfloor \frac{k}{32} \rfloor)]$$

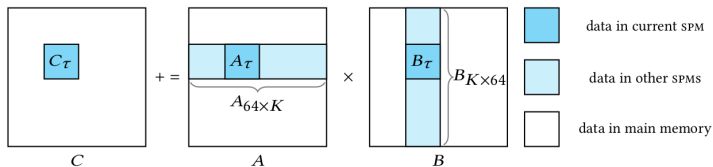
$$\text{BAND: } [S_1(i, j, k) \rightarrow (i - 64 \lfloor \frac{i}{64} \rfloor, j - 64 \lfloor \frac{j}{64} \rfloor, k - 32 \lfloor \frac{k}{32} \rfloor)]$$

(b) The schedule tree with CPE mesh parameters.

Compute Decomposition

Strip-mining the Reduced Dimension

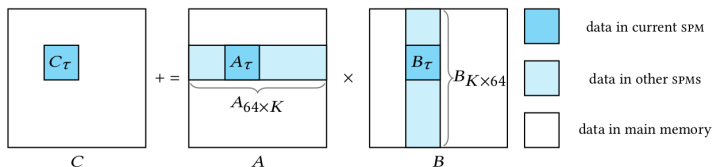
Each CPE can buffer a size of 64×64 tile of the output matrix and 64×32 tiles of the input matrices on its own SPM (allowed by RMA).



Compute Decomposition

Strip-mining the Reduced Dimension

Each CPE can buffer a size of 64×64 tile of the output matrix and 64×32 tiles of the input matrices on its own SPM (allowed by RMA).



DOMAIN: $\{S_1(i, j, k) : 0 \leq i < M \wedge 0 \leq j < N \wedge 0 \leq k < K\}$

BAND: $[S_1(i, j, k) \rightarrow (Rid, Cid)]$ /* This band is mapped to the 2D CPE mesh. */

BAND: $[S_1(i, j, k) \rightarrow (\lfloor \frac{k}{256} \rfloor)]$

BAND: $[S_1(i, j, k) \rightarrow (\lfloor \frac{k}{32} \rfloor - 8 \lfloor \frac{k}{256} \rfloor)]$

BAND: $[S_1(i, j, k) \rightarrow (i - 64 \lfloor \frac{i}{64} \rfloor, j - 64 \lfloor \frac{j}{64} \rfloor, k - 32 \lfloor \frac{k}{32} \rfloor)]$

Athread programming model for SW26010Pro

- *dma_iget* (void *dst, void *src, int size, int len, int strip, int *reply)
- *dma_iput* (void *dst, void *src, int size, int len, int strip, int *reply)

Inserting extension nodes in schedule tree

- *the schedule tree with shared memory promotion statements*

DOMAIN: $\{S_1(i, j, k) : 0 \leq i < M \wedge 0 \leq j < N \wedge 0 \leq k < K\}$

BAND: $[S_1(i, j, k) \rightarrow (\text{blockIdx.y}, \text{blockIdx.x}, \lfloor \frac{k}{32} \rfloor)]$

EXTENSION: $[(d_0, d_1, d_2) \rightarrow \text{readA}[d_3, d_4]; (d_0, d_1, d_2) \rightarrow \text{readB}[d_3, d_4]]$

BAND: $[S_1(i, j, k) \rightarrow (\text{threadIdx.y}, \text{threadIdx.x}, k - 32 \lfloor \frac{k}{32} \rfloor)]$

SEQUENCE:

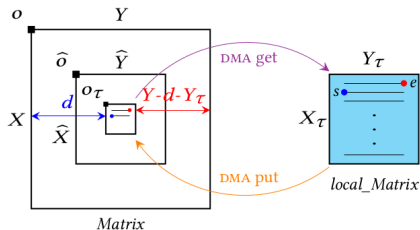
FILTER: $\{\text{readA}[d_3, d_4]\}$

FILTER: $\{\text{readB}[d_3, d_4]\}$

FILTER: $\{S_1(i, j, k)\}$

Automating DMA Communication

The DMA Mechanism



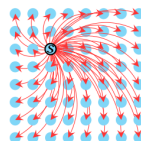
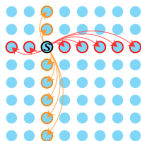
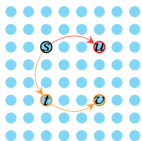
$$\begin{aligned} \text{Matrix=C} & \begin{cases} X = M \wedge Y = N \\ \widehat{X} = 512 \wedge \widehat{Y} = 512 \\ X_{\tau} = 64 \wedge Y_{\tau} = 64 \end{cases} \\ \text{Matrix=A} & \begin{cases} X = M \wedge Y = K \\ \widehat{X} = 512 \wedge \widehat{Y} = 256 \\ X_{\tau} = 64 \wedge Y_{\tau} = 32 \end{cases} \\ \text{Matrix=B} & \begin{cases} \widehat{X} = K \wedge \widehat{Y} = N \\ X = 256 \wedge Y = 512 \\ X_{\tau} = 32 \wedge Y_{\tau} = 64 \end{cases} \end{aligned}$$

■ data in SPM
□ data in main memory

- each CPE mesh executes a GEMM kernel of 512x512x256.
- `dma_iget (&Matrix[0][0], &Matrix[r][c], $X_{\tau} \times Y_{\tau}$, Y_{τ} , $Y - Y_{\tau}$, &reply)`
- **Automatically implementing data movements** from the main memory to the SPMs

Implementing RMA Broadcast

The communication manners between CPEs



(a) Point to Point. (b) Row/column broadcast. (c) All broadcast.

Implementing RMA Broadcast

Insert extension nodes for DMA and RMA

DOMAIN: $\{S_1(i, j, k) : 0 \leq i < M \wedge 0 \leq j < N \wedge 0 \leq k < K\}$
BAND: $[S_1(i, j, k) \rightarrow (Rid, Cid)]$
EXTENSION: $[(d_0, d_1, d_2) \rightarrow getC(d_3, d_4) / get_replyC()]$
EXTENSION: $[(d_0, d_1, d_2) \rightarrow putC(d_3, d_4) / put_replyC()]$
SEQUENCE:
FILTER: $\{ getC(d_3, d_4) \} \otimes$ **FILTER:** $\{ get_replyC() \}$
FILTER: $\{ S_1(i, j, k) \}$
BAND: $[S_1(i, j, k) \rightarrow (\lfloor \frac{k}{256} \rfloor)]$
EXTENSION: $[(d_0, d_1, d_2) \rightarrow getA(d_3, d_4) / get_reply_A()]$
EXTENSION: $[(d_0, d_1, d_2) \rightarrow getB(d_3, d_4) / get_reply_B()]$
SEQUENCE:
FILTER: $\{ getA(d_3, d_4) \} \oplus$ **FILTER:** $\{ get_replyA() \}$
FILTER: $\{ getB(d_3, d_4) \} \oplus$ **FILTER:** $\{ get_replyB() \}$
FILTER: $\{ S_1(i, j, k) \}$
BAND: $[S_1(i, j, k) \rightarrow (\lfloor \frac{k}{32} \rfloor - 8 \lfloor \frac{k}{256} \rfloor)]$
EXTENSION: $[(d_0, d_1, d_2, d_3) \rightarrow rbcstA(d_4, d_5) / rbcst_replyA()]$
EXTENSION: $[(d_0, d_1, d_2, d_3) \rightarrow rbcstB(d_4, d_5) / rbcst_replyB()]$
SEQUENCE:
FILTER: $\{ rbcstA(d_4, d_5) \} \oplus$ **FILTER:** $\{ rbcst_replyA() \}$
FILTER: $\{ rbcstB(d_4, d_5) \} \oplus$ **FILTER:** $\{ rbcst_replyB() \}$
FILTER: $\{ S_1(i, j, k) \}$
BAND: $[S_1(i, j, k) \rightarrow (i - 64 \lfloor \frac{i}{64} \rfloor, j - 64 \lfloor \frac{j}{64} \rfloor, k - 32 \lfloor \frac{k}{32} \rfloor)]$
FILTER: $\{ putC(d_3, d_4) \} \otimes$ **FILTER:** $\{ put_replyC() \}$

Implementing RMA Broadcast

Insert extension nodes for DMA and RMA

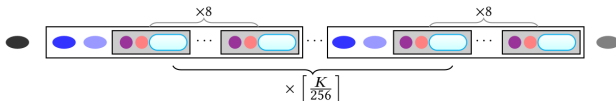
```
DOMAIN: {S1(i, j, k) : 0 ≤ i < M ∧ 0 ≤ j < N ∧ 0 ≤ k < K}
BAND: [S1(i, j, k) → (Rid, Cid)]
EXTENSION: [(d0, d1, d2) → getC(d3, d4) / get_replyC()]
EXTENSION: [(d0, d1, d2) → putC(d3, d4) / put_replyC()]
SEQUENCE:
  FILTER: { getC(d3, d4) } ⊗ FILTER: { get_replyC() }
  FILTER: { S1(i, j, k) }
  BAND: [S1(i, j, k) → (⌊  $\frac{k}{256}$  ⌋)]
  EXTENSION: [(d0, d1, d2) → getA(d3, d4) / get_reply_A()]
  EXTENSION: [(d0, d1, d2) → getB(d3, d4) / get_reply_B()]
  SEQUENCE:
    FILTER: { getA(d3, d4) } ⊕ FILTER: { get_replyA() }
    FILTER: { getB(d3, d4) } ⊕ FILTER: { get_replyB() }
    FILTER: { S1(i, j, k) }
    BAND: [S1(i, j, k) → (⌊  $\frac{k}{32}$  ⌋ - 8 ⌊  $\frac{k}{256}$  ⌋)]
    EXTENSION: [(d0, d1, d2, d3) → rbcstA(d4, d5) / rbcst_replyA()]
    EXTENSION: [(d0, d1, d2, d3) → rbcstB(d4, d5) / rbcst_replyB()]
    SEQUENCE:
      FILTER: { rbcstA(d4, d5) } ⊕ FILTER: { rbcst_replyA() }
      FILTER: { rbcstB(d4, d5) } ⊕ FILTER: { rbcst_replyB() }
      FILTER: { S1(i, j, k) }
      BAND: [S1(i, j, k) → (i - 64 ⌊  $\frac{i}{64}$  ⌋, j - 64 ⌊  $\frac{j}{64}$  ⌋, k - 32 ⌊  $\frac{k}{32}$  ⌋)]
  FILTER: { putC(d3, d4) } ⊗ FILTER: { put_replyC() }
```

- Automatically implementing data communication within CPE mesh

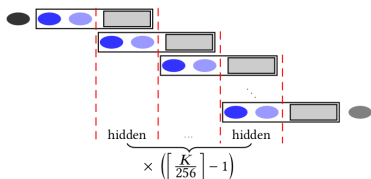
Memory Latency Hiding

Software Pipelining

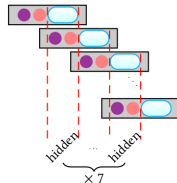
The two-level memory latency hiding strategy



(a) Sequential execution along the k loop dimension.



(b) Hiding DMA.



(c) Hiding RMA.



Code Generation

- Code generation phases
 - Reusing the AST generator of isl
 - Printing Athread syntax

- Code generation phases
 - Reusing the AST generator of isl
 - Printing Athread syntax
- Inline Assembly Routine
 - The assembly micro kernel is provided as a compiled object, which has been highly optimized by the Sunway architects
 - We use a mark node in schedule tree to instruct the code generator to print an assembly function call.

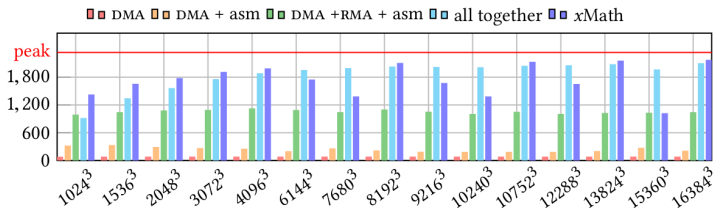
Code Generation

- Code generation phases
 - Reusing the AST generator of isl
 - Printing Athread syntax
- Inline Assembly Routine
 - The assembly micro kernel is provided as a compiled object, which has been highly optimized by the Sunway architects
 - We use a mark node in schedule tree to instruct the code generator to print an assembly function call.
- Fusion Patterns
 - Fusion patterns with a prologue/epilogue operation

- Experiment Environment.
 - PPCG (Polyhedral Parallel Code Generation)
 - SW26010Pro Processor
 - SWGCC Compiler
- Comparison
 - xMath version 2.0 (tuned BLAS library of the SW26010 processor).

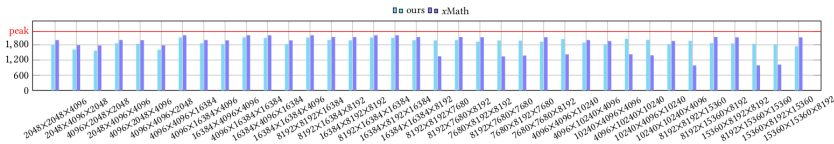
Experiments

- Performance Breakdown



- DMA (Baseline version): 84.89 Gflops
- DMA+asm: 240.39 Gflops
- DMA+RMA+asm: 1052.94 Gflops
- all together: 1849.06 Gflops

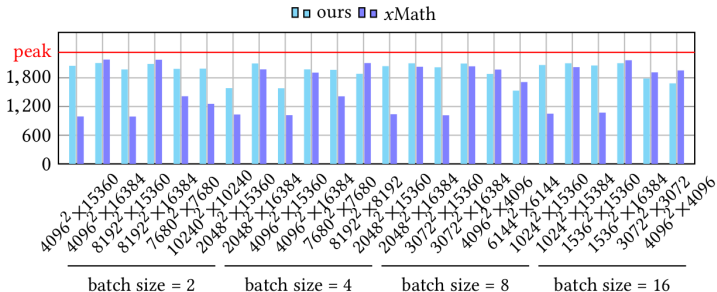
- Performance Comparison of GEMM and xMath



- xMath achieves a mean 1746.97 Gflops
- Our approach outperforms it by 9.62%

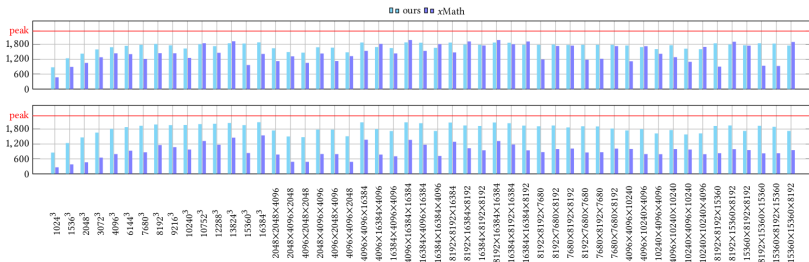
Experiments

- Performance Comparison of Batched GEMM



- Batch sizes: 2, 4, 8, 16
- Average performance: 1949.92 Glops
- Our approach outperforms xMath by 1.30x

- Performance Comparison of Fusion



- 2.11x speedup when fused with epilogue
- 1.26x speedup when fused with prologue

Conclusion

- We present a method to automatically generate matrix multiplication kernels for the latest Sunway processor.

Conclusion

- We present a method to automatically generate matrix multiplication kernels for the latest Sunway processor.
- Polyhedral Transformations
 - Compute decomposition
 - DMA/RMA communication
 - Memory latency hiding
- Low-level optimizations
 - inline assembly kernel

Conclusion

- We present a method to automatically generate matrix multiplication kernels for the latest Sunway processor.
- Polyhedral Transformations
 - Compute decomposition
 - DMA/RMA communication
 - Memory latency hiding
- Low-level optimizations
 - inline assembly kernel
- One can obtain up to more than **90% of the theoretical performance** within few lines of C code

The End

Questions? Comments?