

# Recent Advances in Batch Linear Algebra Computation

Mawussi ZOUNON

joint work with

Samuel RELTON

# Overview

**Aim:** Design of efficient batch linear algebra kernels

**Target:** batch of very small matrices ( $2 \times 2$  to  $32 \times 32$ )

**Example:** solve  $C^{(i)} \leftarrow \alpha^{(i)}A^{(i)}B^{(i)} + \beta^{(i)}C^{(i)}$ ,  $i = 1:\text{batch\_count}$

**Common Approach:** OpenMP for loop

**Challenge:** keep CPU cores running at full efficiency

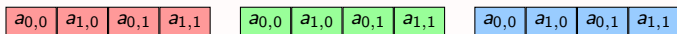
1. The OpenMP for loop approach
2. The interleaved memory layout approach
3. Experimental results
4. Concluding remarks

# The OpenMP for loop approach

## Example of batch DGEMM

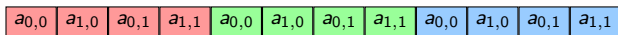
- 1: **for**  $i \leftarrow 1$  to *batch\_count* **do**
  - 2: | `#pragma omp parallel for`
  - 3: | `cblas_dgemm(...)`
  - 4: **end for**
- **Principle:** Each core solve one problem at a time
- **Advantages:**
  - ▶ Easy to implement
  - ▶ Use existing vendor optimized kernels
- **Disadvantage:**
  - ▶ Ineffective CPU cache exploitation
  - ▶ Fails to use vector units at full efficiency

## Pointer to pointer layout:



- **High memory access:** thousands of matrices spread out in memory
- **Ineffective cache use:** matrices smaller than cache size

## Stride layout:



- **Possible overhead:** may require layout conversion
- **Vectorization issue:** may fail to use wide vector units at full efficiency

# Batch linear algebra and wide vector units

## Issues of very small matrices

- Intel AVX-512: 8 double precision
- $2 \times 2$  matrix: 4 double precision
- Alternative: cross-matrix vectorization



matrix



Intel AVX-512

## Issues of very small matrices

- Intel AVX-512: 8 double precision
- $2 \times 2$  matrix: 4 double precision
- Alternative: cross-matrix vectorization



matrix

empty

# Batch linear algebra and wide vector units

## Issues of very small matrices

- Intel AVX-512: 8 double precision
- $2 \times 2$  matrix: 4 double precision
- Alternative: cross-matrix vectorization

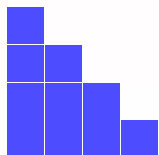


matrix

empty

## Issues of sequential algorithms

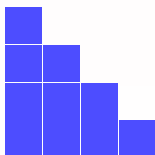
- Synchronization point
- One division: Only 1 DP flop over 8



$A^{(1)}$



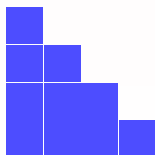
$B^{(1)}$



$A^{(2)}$




$B^{(2)}$



$A^{(3)}$



$B^{(3)}$

 Division

 Update



# Batch linear algebra and wide vector units

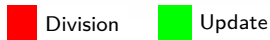
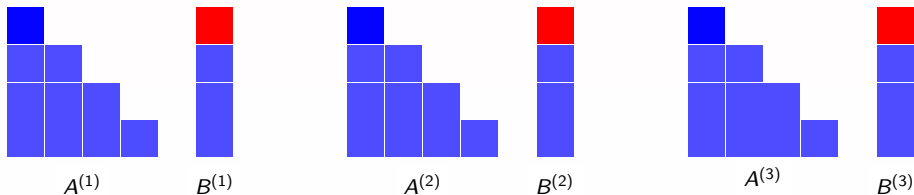
## Issues of very small matrices

- Intel AVX-512: 8 double precision
- $2 \times 2$  matrix: 4 double precision
- Alternative: cross-matrix vectorization



## Issues of sequential algorithms

- Synchronization point
- One division: Only 1 DP flop over 8



# Batch linear algebra and wide vector units

## Issues of very small matrices

- Intel AVX-512: 8 double precision
- $2 \times 2$  matrix: 4 double precision
- Alternative: cross-matrix vectorization

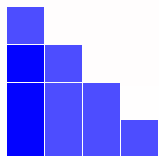


matrix

empty

## Issues of sequential algorithms

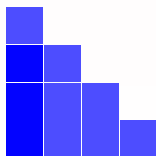
- Synchronization point
- One division: Only 1 DP flop over 8



$A^{(1)}$



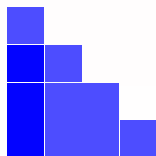
$B^{(1)}$



$A^{(2)}$



$B^{(2)}$



$A^{(3)}$



$B^{(3)}$

 Division

 Update

# Batch linear algebra and wide vector units

## Issues of very small matrices

- Intel AVX-512: 8 double precision
- $2 \times 2$  matrix: 4 double precision
- **Alternative:** cross-matrix vectorization

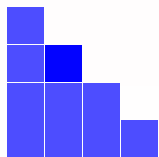


matrix

empty

## Issues of sequential algorithms

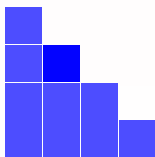
- Synchronization point
- **One division:** Only 1 DP flop over 8
- **Alternative:** cross-matrix vectorization
- Investigate novel memory layouts



$A^{(1)}$



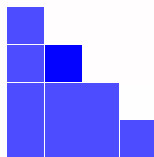
$B^{(1)}$



$A^{(2)}$



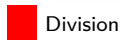
$B^{(2)}$



$A^{(3)}$



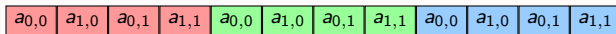
$B^{(3)}$



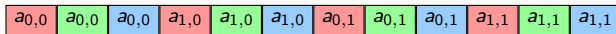
1. The OpenMP for loop approach
2. The interleaved memory layout approach
3. Experimental results
4. Concluding remarks

# Interleaved memory layout

Stride layout:



Interleaved memory layout:



Less prone to cache misses  
Maximizes use of vector units

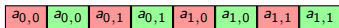
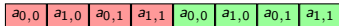
Requires layout conversion  
Requires appropriate algorithms

# Algorithm design for interleaved layout

Full interleave of a large batch: **high cache misses**



Block interleaved



- Divide the batch in blocks
- Interleave each block
- Block big enough to fill vector units
- Block small enough to fit in cache

# Algorithm design for interleaved layout

## Basic GEMM algorithm

```
for  $i \leftarrow 1$  to  $m$  do
  for  $j \leftarrow 1$  to  $n$  do
     $c_{i,j} \leftarrow 0$ 
    for  $k \leftarrow 1$  to  $K$  do
       $c_{i,j} \leftarrow c_{i,j} + a_{i,k} \times b_{k,j}$ 
    end for
  end for
end for
```

# Algorithm design for interleaved layout

## Basic GEMM algorithm

```
for i ← 1 to m do
  for j ← 1 to n do
    ci,j ← 0
    for k ← 1 to K do
      ci,j ← ci,j + ai,k × bk,j
    end for
  end for
end for
```

- Extra inner for loop
- Vectorization transformations

## Interleaved basic GEMM algorithm

```
for i ← 1 to m do
  for j ← 1 to n do
    for k ← 1 to K do
      #pragma simd
      for idx ← 1 to BlockSize do
        if k == 1 then
          ci,j(idx) ← 0
        end if
        ci,j(idx) ← ci,j(idx) + ai,k(idx) × bk,j(idx)
      end for
    end for
  end for
end for
```

- Kernel for a single block
- Parallelize blocks over cores

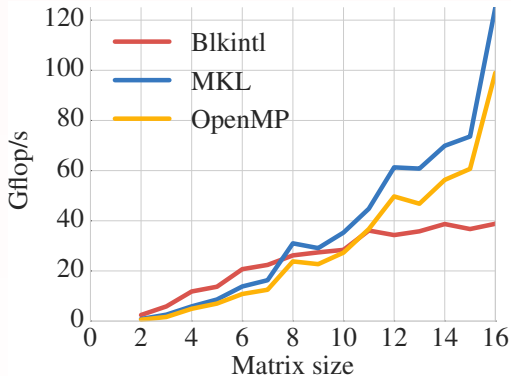


# Algorithm design for interleaved layout

- **Auto-tuning**: set the optimal block size
- **Conversion**: translate each block to the interleaved layout
- **Batch computation**: call block interleaved kernels
- **Conversion**: translate results back to the user layout

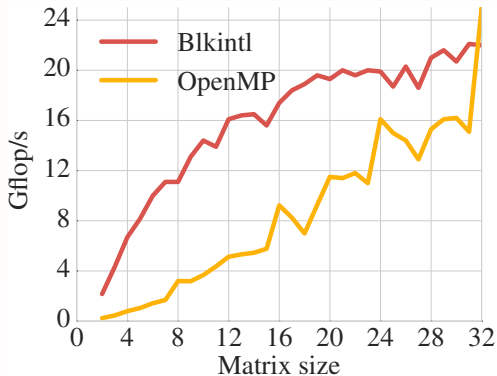
1. The OpenMP for loop approach
2. The interleaved memory layout approach
3. Experimental results
4. Concluding remarks

# Interleaved memory layout - GEMM (Intel KNL)



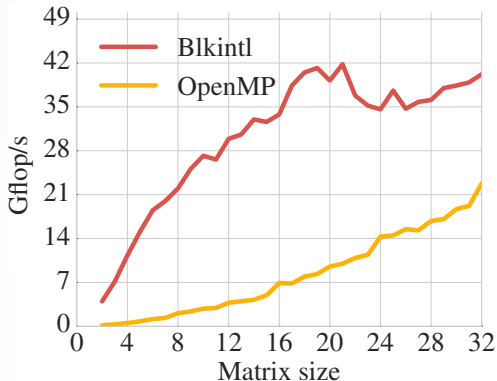
- Batch of 10,000 DGEMM
- 68 threads on 68-core Intel KNL
- Flush the cache
- Include conversion time

# Interleaved memory layout - TRSM (Intel KNL)



- Batch of 10,000 DTRSM (4rhs)
- 68 threads on 68-core Intel KNL
- Flush the cache
- Include conversion time

# Interleaved memory layout - DPOSV (Intel KNL)



- Batch of 10,000 POSV (4rhs)
- 68 threads on 68-core Intel KNL
- Flush the cache
- Include conversion time

1. The OpenMP for loop approach
2. The interleaved memory layout approach
3. Experimental results
4. Concluding remarks

## Remarks

- Keep OpenMP for loop for medium size matrices
- Block interleaved for very small matrices

## Perspective

- Extend interleave to other kernels
- Interleaved LU factorization

A Batch Of Thanks