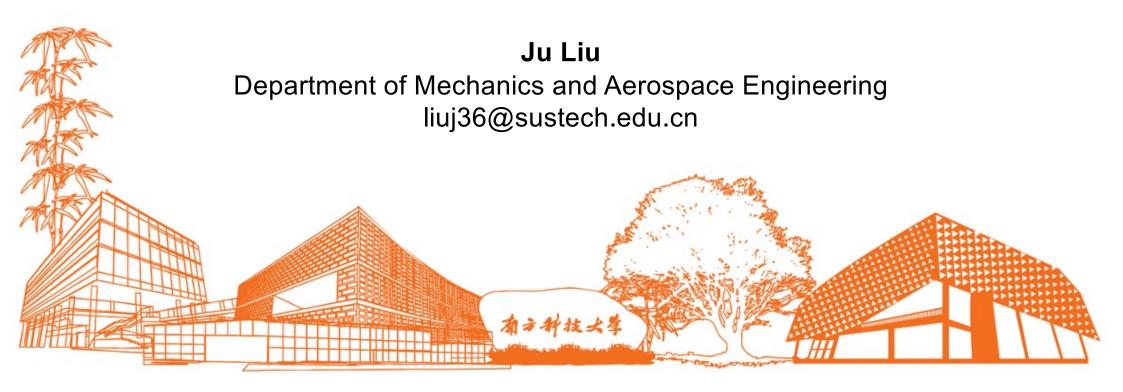
# MAE 5032 High Performance Computing: Methods and Applications

Lab 3: SIMD & BLAS



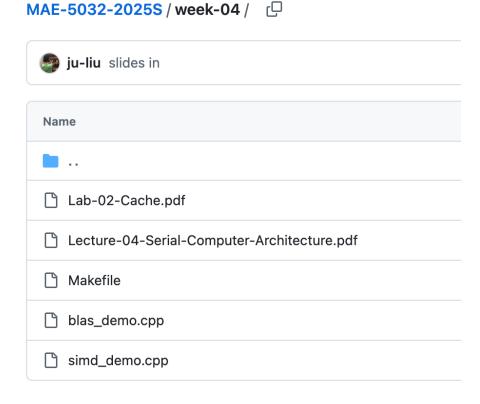
# **Objective**

- You will experiment two codes
  - experience the invocation of SIMD in an explicit manner
  - explore the block matrix multiplication
  - learn the use of BLAS functions

# Task 1: SIMD

Go to <a href="https://github.com/ju-liu/MAE-5032-2025S/tree/main/week-04">https://github.com/ju-liu/MAE-5032-2025S/tree/main/week-04</a>

Download the code



## Task 1: SIMD

```
#include <iostream>
#include <vector>
#include <chrono>
#include <immintrin.h> // SIMD Intrinsics for AVX
#include <cstdlib> // For aligned memory allocation

using namespace std;
using namespace std::chrono;

const int N = 128;
const int REPEATS = 10000000;

// Allocate aligned memory (32-byte aligned for AVX)
float* aligned_alloc(int size) {
   void* ptr = nullptr;
   posix_memalign(&ptr, 32, size * sizeof(float)); // POSIX aligned allocation
   return (float*)ptr;
}
```

Allocate the memory with size x 4 bytes whose starting address is guaranteed to be a multiple of 32 bytes = 256 bits.

This is necessary for AVX.

### Task 1: SIMD

```
// Scalar (regular) vector addition
void scalar_add(const float* A, const float* B, float* C, int n) {
  for(int r =0; r<REPEATS; r++)</pre>
    for (int i = 0; i < n; i++)
      C[i] = A[i] + B[i];
// SIMD vector addition using AVX with aligned memory
void simd_add(const float* A, const float* B, float* C, int n) {
  for(int r =0; r<REPEATS; r++)</pre>
    for (int i = 0; i < n; i += 8) { // AVX processes 8 floats at a time
                                         // Aligned load
      m256 a = mm256 load ps(&A[i]);
        m256 b = _mm256_load_ps(&B[i]);
       m256 c = mm256 add ps(a, b);
                                         // SIMD parallel addition
      _mm256_store_ps(&C[i], c);
                                         // Aligned store
```

\_mm256\_load\_ps is part of the AVX instruction set.

It load 8 single-precision floatingpoint values from memory into AVX 256-bit register

\_mm256\_add\_ps: addition in parallel

\_mm256\_store\_ps: writes 8 floating point values from an AVX 256-bit register into memory

```
int main() {
  // Allocate 32-byte aligned memory
  float* A = aligned_alloc(N);
  float* B = aligned alloc(N);
  float* C = aligned_alloc(N);
  // Initialize vectors
  for (int i = 0; i < N; i++) {
    A[i] = i * 1.0f;
    B[i] = (N - i) * 1.0f;
  // Measure scalar addition time
  auto start = high resolution clock::now();
  scalar_add(A, B, C, N);
  auto end = high_resolution_clock::now();
  cout << "Scalar time: "</pre>
    << duration<double>(end - start).count() << " sec" << endl:
  // Measure SIMD addition time
  start = high_resolution_clock::now();
  simd add(A, B, C, N);
  end = high_resolution_clock::now();
  cout << "SIMD time: "</pre>
    << duration<double>(end - start).count() << " sec" << endl;
 // Free aligned memory
  free(A); free(B); free(C);
  return 0:
```

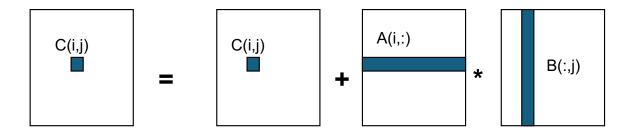
#### Compile the code by

```
g++ -00 -std=c++11 -mavx
-o simd simd_demo.cpp
```

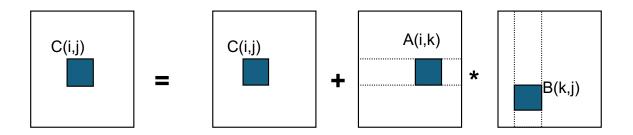
# You will get an executable named simd

./simd

# Task 2: BLAS



# Task 2: BLAS



## Task 2: BLAS

cblas\_dgemm is part of the BLAS library and is used for matrix-matrix multiplication

C = alpha AB + beta C

A is m-by-k, B is k-by-n, and C is m-by-n. alpha and beta are scalars.

#### Linux:

#### Install openblas

```
sudo apt update
Sudo apt install openblas
```

#### Check to see if the lib is installed

```
ls /usr/lib/x86_64-linux-gnu | grep openblas
ls /usr/include/x86_64-linux-gnu | grep openblas
```

#### Compile the code

```
g++ -std=c++11 blas_demo.cpp -I/usr/include/x86_64-linux-gnu -L/usr/lib/x86_64-linux-gnu -lopenblas
```

#### Mac:

#### Install homebrew if not already installed

```
/bin/bash -c "$(curl -fsSL
https://raw.githubusercontent.com/Homebrew/install/HEAD/install.s
h)"
Or go to https://brew.sh/
```

#### Install openblas

brew install openblas

#### Mac:

#### Check install info

brew info openblas

```
juliu::Riemann {~ }
-> brew info openblas
==> openblas: stable 0.3.29 (bottled), HEAD [keg-only]
Optimized BLAS library
https://www.openblas.net/
Installed
/usr/local/Cellar/openblas/0.3.29 (24 files, 134.8MB)
  Poured from bottle using the formulae.brew.sh API on 2025-03-11 at 16:37:00
From: https://github.com/Homebrew/homebrew-core/blob/HEAD/Formula/o/openblas.rb
License: BSD-3-Clause AND BSD-2-Clause-Views AND BSD-3-Clause-Open-MPI AND BSD-2-Clau
se
==> Dependencies
Required: gcc <
==> Options
--HEAD
        Install HEAD version
==> Caveats
openblas is keg-only, which means it was not symlinked into /usr/local,
because macOS provides BLAS in Accelerate.framework.
For compilers to find openblas you may need to set:
  export LDFLAGS="-L/usr/local/opt/openblas/lib"
  export CPPFLAGS="-I/usr/local/opt/openblas/include"
For pkg-config to find openblas you may need to set:
  export PKG CONFIG PATH="/usr/local/opt/openblas/lib/pkgconfig"
==> Analytics
install: 31,948 (30 days), 89,333 (90 days), 503,982 (365 days)
install-on-request: 7,486 (30 days), 20,601 (90 days), 125,029 (365 days)
build-error: 23 (30 days)
```

```
int main() {
   // Allocate matrices on the heap to prevent stack overflow
   std::vector<double> A(N * N, 1.0);
   std::vector<double> B(N * N, 1.0);
   std::vector<double> C(N * N, 0.0);
   // Naive multiplication
   auto start = std::chrono::high_resolution_clock::now();
   matrix_mult_naive(A, B, C);
   auto end = std::chrono::high_resolution_clock::now();
   std::cout << "Naive Matrix Multiplication Time: "</pre>
              << std::chrono::duration_cast<std::chrono::milliseconds>(end - start).count()
              << " ms" << std::endl;
   // Reset C
   std::fill(C.begin(), C.end(), 0.0);
   // Blocked multiplication
   start = std::chrono::high_resolution_clock::now();
   matrix_mult_blocked(A, B, C, Bsize);
   end = std::chrono::high_resolution_clock::now();
   std::cout << "Blocked Matrix Multiplication Time: "</pre>
              << std::chrono::duration_cast<std::chrono::milliseconds>(end - start).count()
              << " ms" << std::endl;
   // Reset C
   std::fill(C.begin(), C.end(), 0.0);
   // BLAS DGEMM multiplication
   start = std::chrono::high_resolution_clock::now();
   matrix_mult_blas(A, B, C);
   end = std::chrono::high_resolution_clock::now();
   std::cout << "BLAS DGEMM Matrix Multiplication Time: "</pre>
              << std::chrono::duration cast<std::chrono::milliseconds>(end - start).count()
              << " ms" << std::endl:
   return 0:
```

#### Compile the code by

g++ -std=c++11 I/usr/local/opt/openblas/include
-L/usr/local/opt/openblas/lib lopenblas -O3 blas\_demo.cpp

# You will get an executable named simd

./a.out