# Kokkos

## A Parallel, Portable Programming Model for CPUs and GPUs

Rohit Kakodkar

Research Software Engineer

- Kokkos is a C++* performance portability library
  - Write a single source implementation
  - Descriptive programming model
  - Compile for CPUs or GPUs
  - Kokkos is a shared memory programming model (works in conjunction with MPI)

- Major buy-in by DOE and national labs
  - LAAMPS, Trilinos, PETSc
  - Over 100 projects using kokkos
  - Contributing to the C++ standard
  - Active kokkos developers community via slack (invite only)

Python and Fortran bindings are available

The goal here is to motivate a use case for Kokkos. Given time limit this is not an in-depth tutorial of Kokkos (probably sometime in the future).

- Outline of the talk:
  - Need for performance portability
  - Understanding performance
    - Caching (CPUs) vs Coalescing(GPUs)
  - Kokkos views
  - Kokkos parallel construct
    - Hands-on example
    - Parallel policies
  - Packaging Kokkos projects
  - Why use kokkos?

Modern high-performance computing applications need to run efficiently across multiple architectures:

|  | GPU | | | CPU |
|---|---|---|---|---|
| **Architecture** |  |  |  |  |
| **Programming model** | CUDA | HIP | DPC++ (SYCL) | OpenMP pthreads |

# Modern high-performance computing applications need to run efficiently across multiple architectures:

Pre-Exascale

Exascale

LANL Trinity
Intel Haswell/ Intel KNL
OpenMP 3

ORNL Summit
NVIDIA Volta100
OpenMP/CUDA

ORNL Frontier
AMD GPUs
OpenMP/HIP

ANL Aurora
Intel GPUs
DPC++

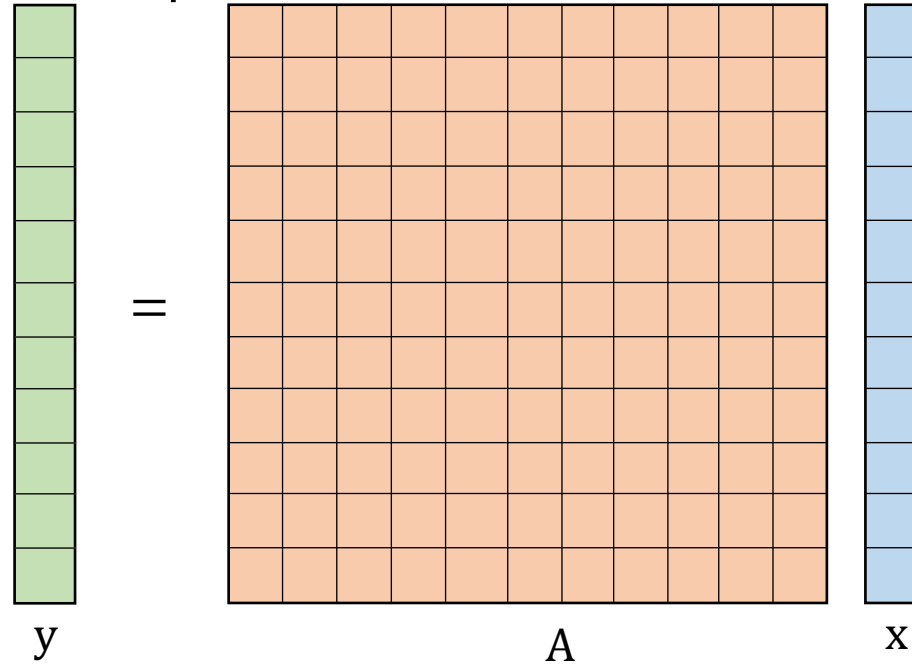Moving into the exascale era will require adapting your applications to utilize modern architectures

- Problem: Porting applications for various architectures is time consuming
  - Typical HPC application : 300k – 600k lines of code
    - Smaller scale applications : SpecFEM2D + SpecFEM3D ~ 100k lines of code (conservative estimate)
    - Porting requires ~10% rewrite of the application
    - Typical software engineer writes about 20k LOC/year


- Potential portability options: OpenMP 5 or OpenACC
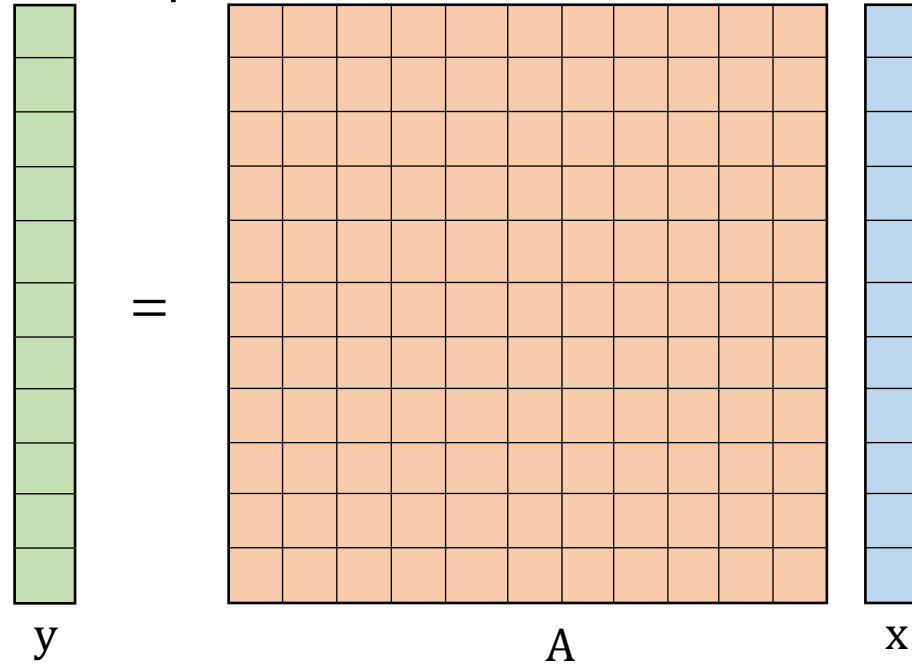
What about performance??

- Example: Matrix – vector multiplication



y            A        x

Serial Implementation

```
for (int j = 0; j < N; j++){
        for (int i = 0; i < M; i++){
                y(j) += A(j,i)*x(i);
        }
}
```

• Example: Matrix – vector multiplication



y           A           x

Parallel Implementation

```
parallel_for (int j = 0; j < N; j++){
// distribute different rows to different threads
    for (int i = 0; i < M; i++){
        y(j) += A(j,i)*x(i);
    }
}
```

- Performance is dependent on matrix layout

Row major format

Column major format

- Good performance

i →

j

Row major format

```
parallel_for j: N
// every thread gets a specific j to execute
    for i: M
        A(j,i)
```

CPU 1    CPU 2

L1       L1

CPU 3    CPU 4

L1       L1

• Bad performance

i →

j ↓

Column major format



```
parallel_for j: N
// every thread gets a specific j to execute
    for i: M
        A(j,i)
```

# • GPU execution model

### Software view

Host

Device

Grid 1

Kernel 1 →

Block (0, 0)  Block (1, 0)  Block (2, 0)

Block (0, 1)  Block (1, 1)  Block (2, 1)

Grid 2

Kernel 2 →

Block (1, 1)

(0,0,1) (1,0,1) (2,0,1) (3,0,1)

Thread (0,0,0)  Thread (1,0,0)  Thread (2,0,0)  Thread (3,0,0)

Thread (0,1,0)  Thread (1,1,0)  Thread (2,1,0)  Thread (3,1,0)

### Hardware view

| SM1 | SM2 | SM3 | SM4 |

Memory

| SM5 | SM6 | SM7 | SM8 |

### Logical view

### Streaming Multiprocessor

Core

### Hardware view

32 threads
32 threads
32 threads
32 threads
32 threads

warp

Shared memory

• Takeaway: Threads in a warp execute the same instruction
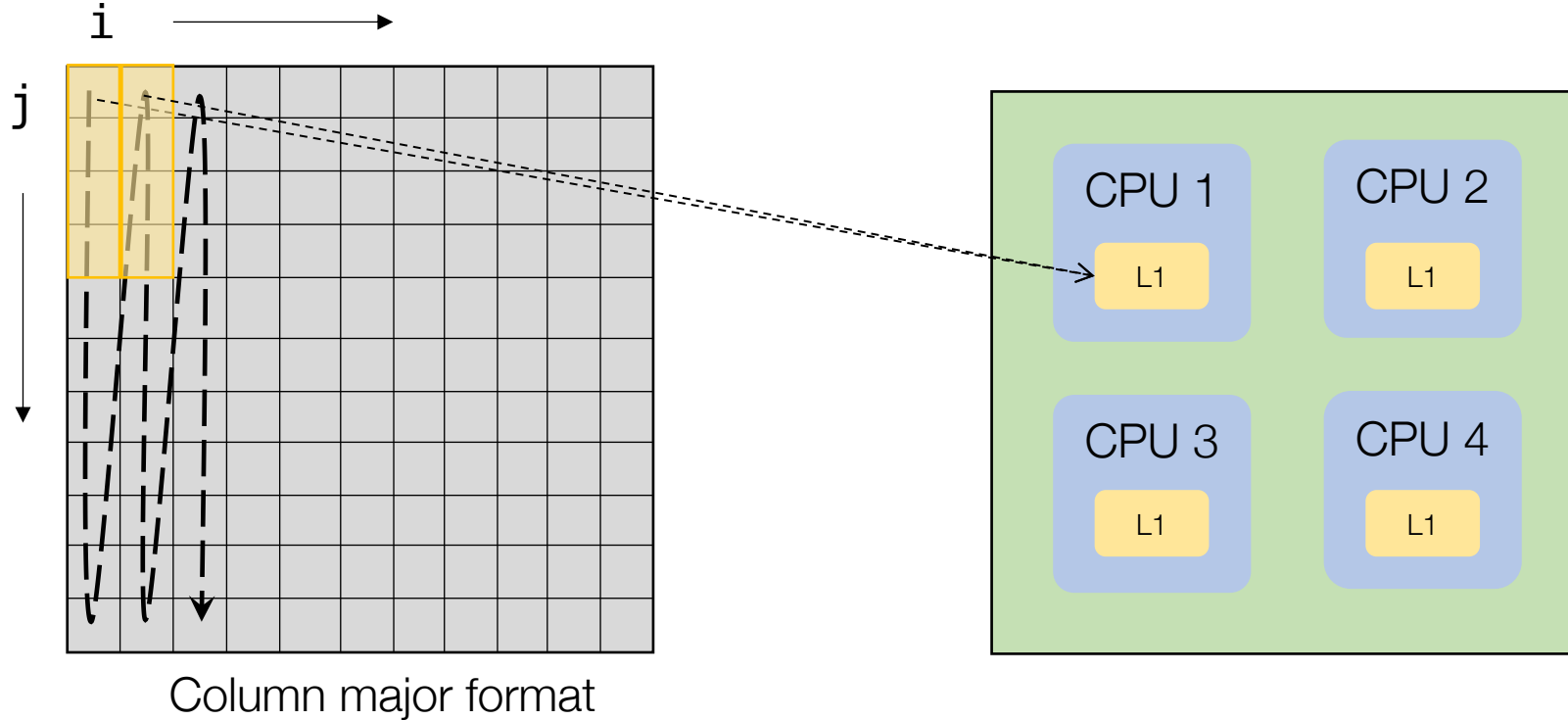
• Bad performance



Warp

```
parallel_for j: N
// every thread gets a specific j to execute
    for i: M
        A(j,i)
```

- Bad performance

i →

j ↓

32 threads
32 threads
32 threads
32 threads
32 threads

Warp

```
parallel_for j: N
// every thread gets a specific j to execute
    for i: M
        A(j,i)
```

- Bad performance

i →

j



32 threads
32 threads
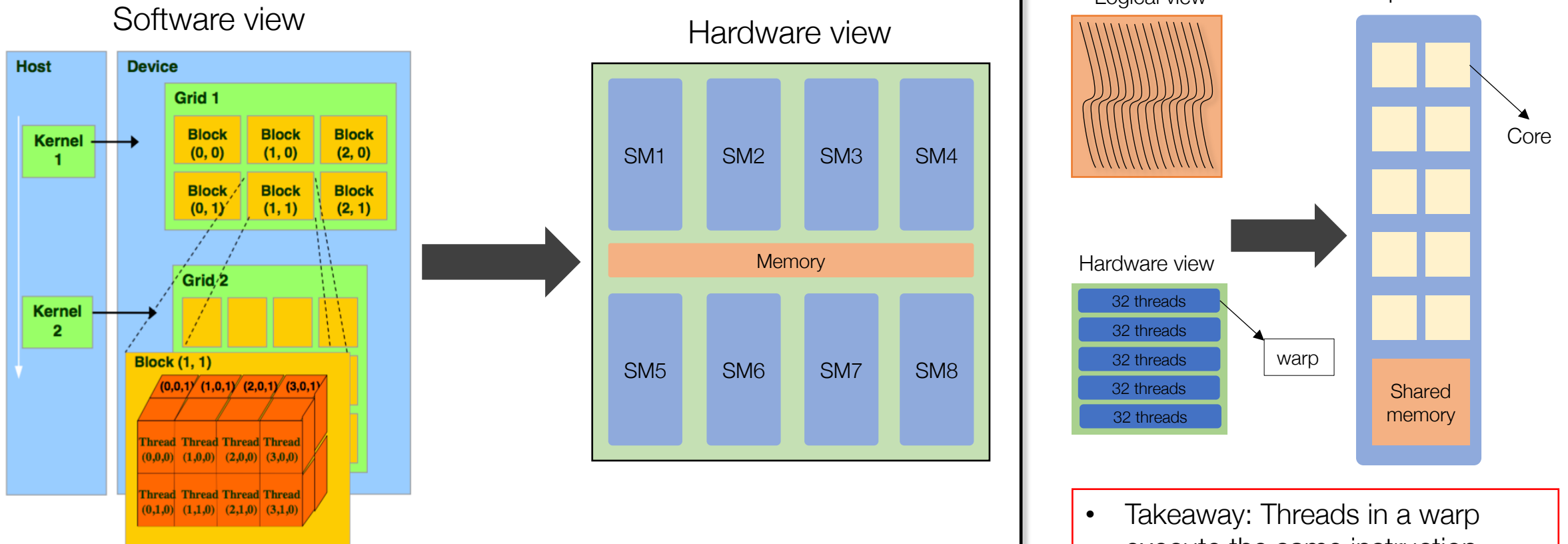32 threads
32 threads
32 threads

Warp

```
parallel_for j: N
// every thread gets a specific j to execute
    for i: M
        A(j,i)
```

- Good performance



```
parallel_for j: N
// every thread gets a specific j to execute
    for i: M
        A(j,i)
```
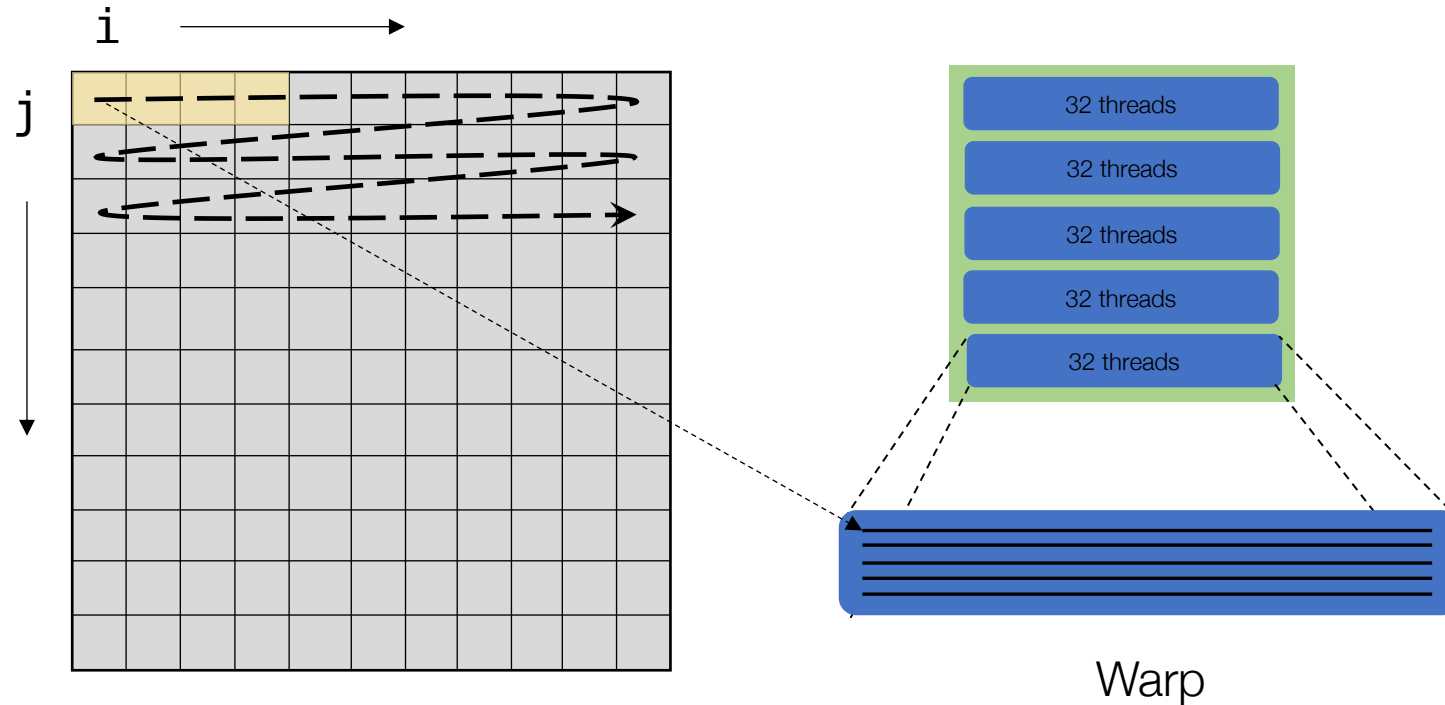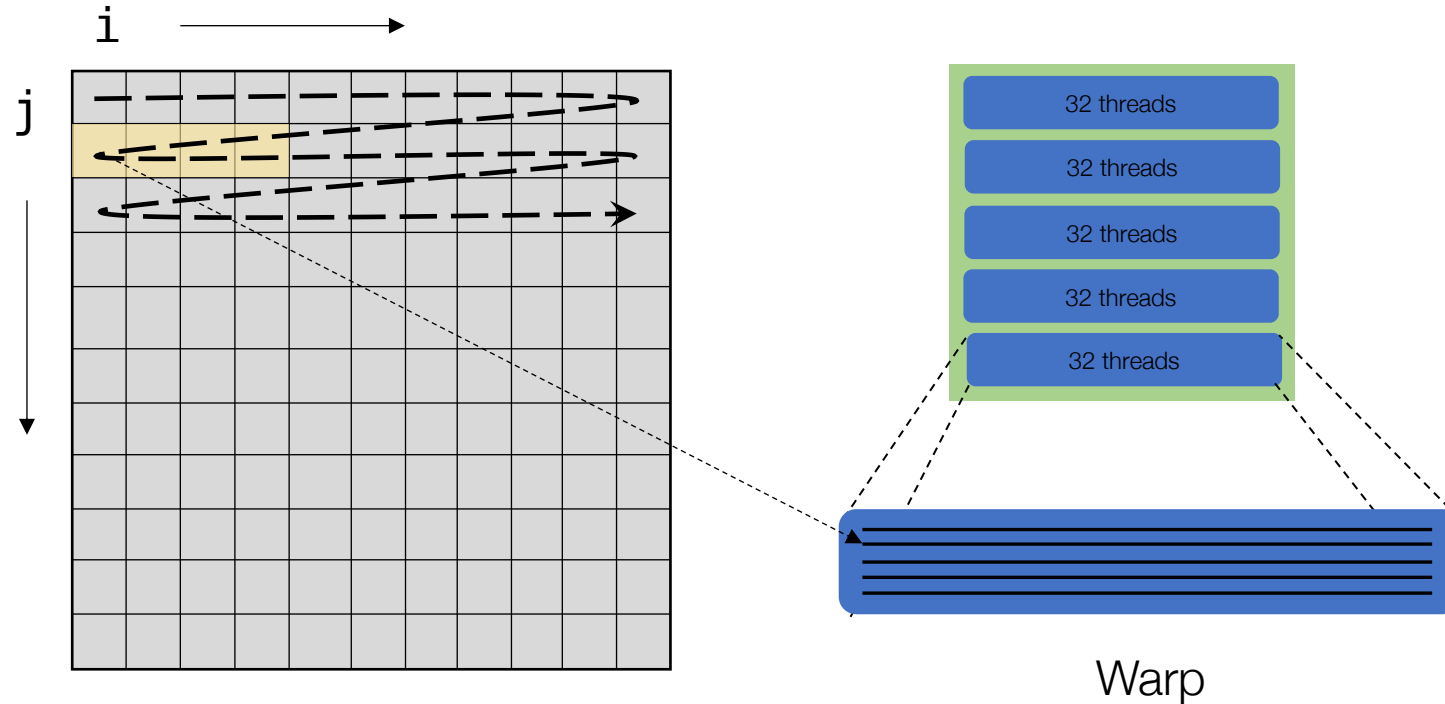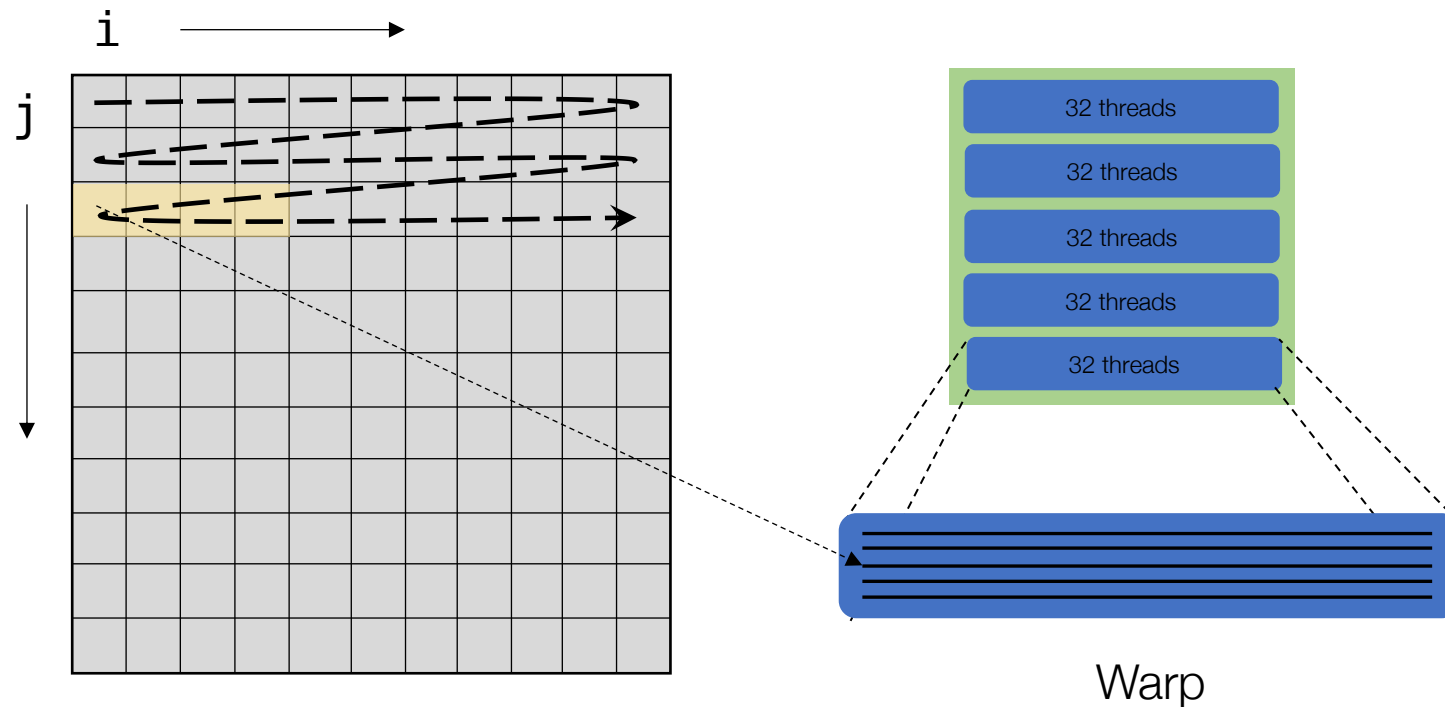
- CPUs require row major format
  - $A(i_1, i_2, …, i_n) \rightarrow i_n$ needs to be continuous in memory (Layout Right)
  - Results in data caching

- GPUs require a column major format
  - $A(i_1, i_2, …, i_n) \rightarrow i_1$ needs to be continuous in memory (Layout Left)
  - Results in data coalescing

- vector - matrix – vector multiplication
- Total number of elements in matrix is a constant (M*N = constant)
- Bandwidth ~ amount of data processed per second

- CPUs require row major format
  - $A(i_1, i_2, ..., i_n) \rightarrow i_n$ needs to be continuous in memory (Layout Right)
  - Results in data caching

- GPUs require a column major format
  - $A(i_1, i_2, ..., i_n) \rightarrow i_1$ needs to be continuous in memory (Layout Left)
  - Results in data coalescing

Data layout needs to be determined at compile time based on target architecture

- Kokkos views are multidimensional arrays
- Can reside on CPUs or GPUs
  - Use "mirror" of your arrays to sync between CPUs and GPUs
  - deep_copy can be initialized between original and mirror views
- View layout can be specified at compile time
  - Default behavior is to optimized layout based on execution space
    - GPUs = Layout Left
    - CPUs = Layout Right
- Views are reference counted
  - Easy memory management
- Copy construction leads to both views pointing to the same data
  - deep copy requires use of "deep_copy" routine rather than a copy construction

```
Kokkos::View<float **, Kokkos::LayoutLeft, Kokkos::CudaSpace>
A("A_matrix", N, M);
```

Policy

```
Pattern  for (int element = 0; element < nelements; element++){
             total = 0;
             for (int qp = 0; qp < numQP; qp++){
                     // reduction over every qp in element
                     total += foo(variable, element, qp);
             }

             elementValue[element] = total;
         }
```

Body

Policy

```
        #pragma omp parallel for schedule(<static, dynamic>)
Pattern for (int element = 0; element < nelements; element++){
            total = 0;
            for (int qp = 0; qp < numQP; qp++){
                    // reduction over every qp in element
                    total += foo(variable, element, qp);
            }

            elementValue[element] = total;
        }
```

Body

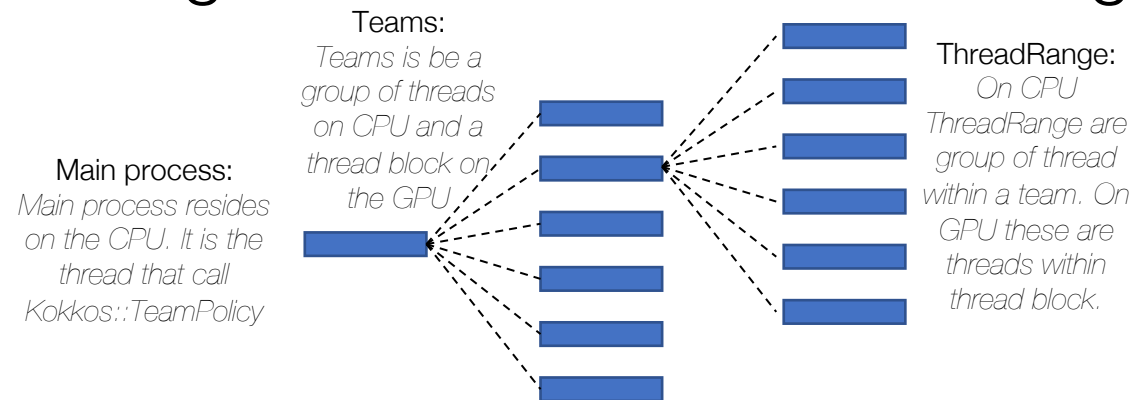Pattern

Policy

```cpp
Kokkos::parallel_for("LoopName", Kokkos::RangePolicy(nelements),
        [=](const int element){
        total = 0;
        for (int qp = 0; qp < numQP; qp++){
                // reduction over every qp in element
                total += foo(variable, element, qp);
        }

        elementValue[element] = total;
});
```

Body

PRINCETON UNIVERSITY

# Example

- **`Kokkos::RangePolicy –`** Each iterate is an integer in a contiguous range

- **`Kokkos::MDRangePolicy –`** Each iterate for each rank is an integer in a contiguous range

- **`Kokkos::TeamPolicy –`** Assigns to each iterate in a contiguous range a team of threads

Teams:
*Teams is be a group of threads on CPU and a thread block on the GPU*

Main process:
*Main process resides on the CPU. It is the thread that call Kokkos::TeamPolicy*

ThreadRange:
*On CPU ThreadRange are group of thread within a team. On GPU these are threads within thread block.*

# Packaging Kokkos projects

- CMAKE - Recommended
  - Provide Kokkos as a dependency package

```
include(FetchContent)
FetchContent_Declare(
kokkos
URL
https://github.com/kokkos/kokkos/archive/refs/tags/3.6.01.zip
)
FetchContent_MakeAvailable(kokkos)

add_executable(question question.cpp)
target_link_libraries(question Kokkos::kokkos)
```

- CMAKE – Recommended
  - In house installations

```
set(KOKKOS_PATH <PATH TO KOKKOS ROOT>)

add_subdirectory(KOKKOS_PATH)
add_executable(question question.cpp)
target_link_libraries(question Kokkos::kokkos)
```

- Makefiles – maybe suitable for individual private projects
  - Check this example out:
    https://github.com/kokkos/kokkos/blob/master/benchmarks/bytes_and_flops/Makefile

PRINCETON UNIVERSITY

# Documentation Overview

| | CUDA/HIP/DPC++ | OpenMP 5/ OpenACC | Kokkos |
|---|---|---|---|
| Portability across architectures | No. Need to write separate kernels for every architecture | Yes. Single source code with pragma-based approach | Yes. Single source code implemented using Kokkos functions |
| Performance | Optimized performance | Tough to optimize for performance | Very good performance |
| Cost of portability | Very high | | High |
| Cost of maintenance | Very high. Newer architectures might require tuning of kernels | Low. Assuming compilers do a good job of implementing the standard | Low. Assuming Kokkos backend is always optimized |
| Compiler dependence | N/A | These are standards, vendors have a flexibility on implementation | N/A |
| Fortran support | No. Could use bindings | Yes | No. Could use bindings |

Why use Kokkos?

- Kokkos is a ***performance portability*** library
- Performance is dependent on data layout
  - Kokkos Views are easy way to manage data layout
- Compile time definition of data using C++ template meta-programming
- Low maintenance overhead for future architectures

Useful resources:
- [Documentation](#)
- [Tutorials](#)
- [Slack](#) (invite only) : ctrott [a] sandia.gov
- Local help - cses@princeton.edu

# Thank you