October 10-12, 2023

# ALCF Hands-on HPC Workshop
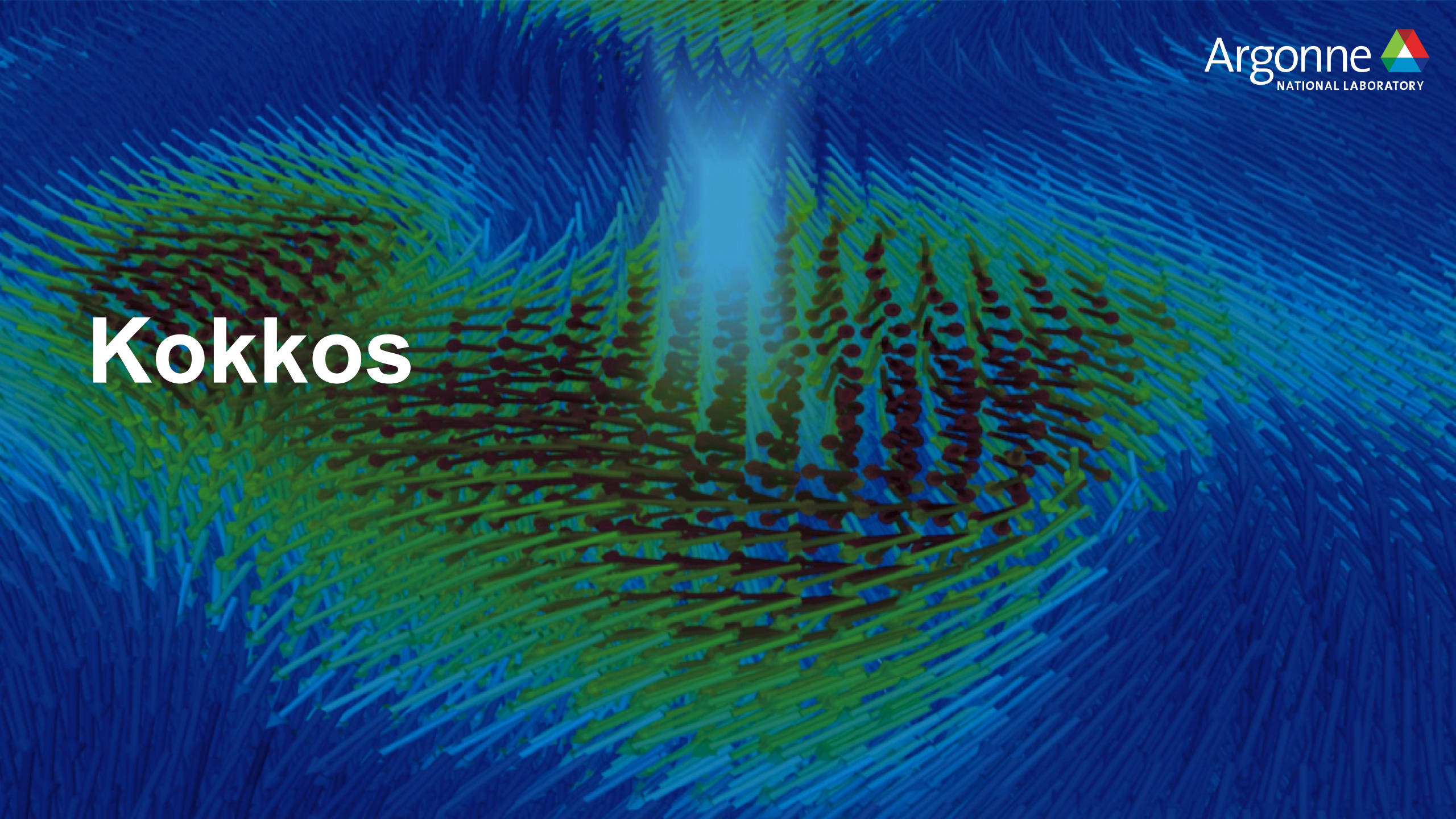
Kokkos

# CG Solve: The AXPBY

- Simple data parallel loop: Kokkos::parallel_for

- Easy to express in most programming models

- Bandwidth bound

- Serial Implementation:

```cpp
void axpby(int n, double* z, double alpha, const double* x,
                             double beta,  const double* y) {
   for(int i=0; i<n; i++)
     z[i] = alpha*x[i] + beta*y[i];
}
```

> Parallel Pattern: for loop

> String Label: Profiling/Debugging

> Execution Policy: do n iterations

> Loop Body

> Iteration handle: integer index

- Kokkos Implementation:

```cpp
void axpby(int n, View<double*> z, double alpha, View<const double*> x,
                                   double beta,  View<const double*> y) {
   parallel_for("AXpBY", n, KOKKOS_LAMBDA ( const int i) {
     z(i) = alpha*x(i) + beta*y(i);
   });
}
```

# CG Solve: The Dot Product

- Simple data parallel loop with reduction: Kokkos::parallel_reduce
- Non trivial in CUDA due to lack of built-in reduction support
- Bandwidth bound
- Serial Implementation:

```cpp
double dot(int n, const double* x, const double* y) {
    double sum = 0.0;
    for(int i=0; i<n; i++)
        sum += x[i]*y[i];
    return sum;
}
```
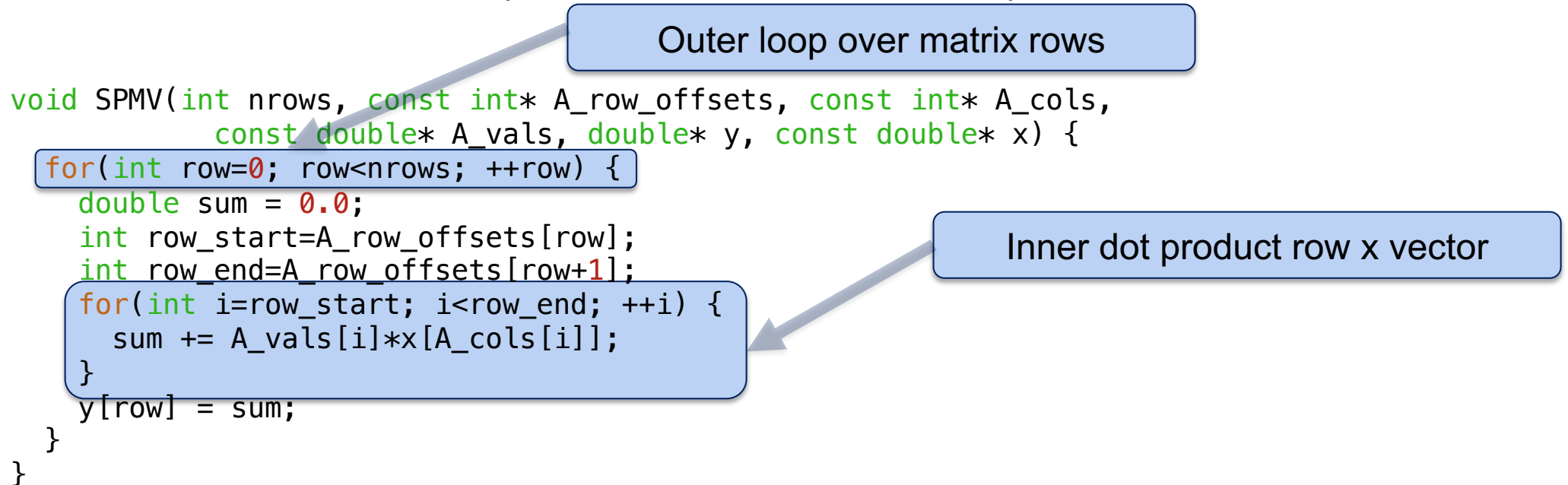
- Kokkos Implementation:

```cpp
double dot(int n, View<const double*> x, View<const double*> y) {
    double x_dot_y = 0.0;
    parallel_reduce("Dot",n, KOKKOS_LAMBDA (const int i,double& sum) {
        sum += x[i]*y[i];
    }, x_dot_y);
    return x_dot_y;
}
```

Parallel Pattern: loop with reduction

Iteration Index + Thread-Local Red. Varible

- Loop over rows

- Dot product of matrix row with a vector

- Example of Non-Tightly nested loops

- Random access on the vector (Texture fetch on GPUs)

Outer loop over matrix rows

```
void SPMV(int nrows, const int* A_row_offsets, const int* A_cols,
          const double* A_vals, double* y, const double* x) {
  for(int row=0; row<nrows; ++row) {
    double sum = 0.0;
    int row_start=A_row_offsets[row];
    int row_end=A_row_offsets[row+1];
    for(int i=row_start; i<row_end; ++i) {
      sum += A_vals[i]*x[A_cols[i]];
    }
    y[row] = sum;
  }
}
```

Inner dot product row x vector

# CG Solve: Sparse Matrix Vector Multiply

```cpp
void SPMV(int nrows, View<const int*> A_row_offsets,
          View<const int*> A_cols, View<const double*> A_vals,
          View<double*> y,
          View<const double*, MemoryTraits< RandomAccess>> x) {

  // Performance heuristic to figure out how many rows to give to a team
  int rows_per_team = get_row_chunking(A_row_offsets);

  parallel_for("SPMV:Hierarchy", TeamPolicy< Schedule< Static > >
      ((nrows+rows_per_team-1)/rows_per_team,AUTO,8),
    KOKKOS_LAMBDA (const TeamPolicy<>::member_type& team) {

    const int first_row = team.league_rank()*rows_per_team;
    const int last_row = first_row+rows_per_team<nrows? first_row+rows_per_team : nrows;

    parallel_for(TeamThreadRange(team,first_row,last_row),[&] (const int row) {
      const int row_start=A_row_offsets[row];
      const int row_length=A_row_offsets[row+1]-row_start;

      double y_row;
      parallel_reduce(ThreadVectorRange(team,row_length),[&] (const int i, double& sum) {
        sum += A_vals(i+row_start)*x(A_cols(i+row_start));
      } , y_row);
      y(row) = y_row;
    });
  });
}
```

Enable Texture Fetch on x

Row x Vector dot product

Distribute rows in workset over team-threads

Team Parallelism over Row Worksets

# Kokkos

```
> git clone https://github.com/kokkos/kokkos.git
> git clone https://github.com/kokkos/kokkos-tutorials.git
> cd kokkos-tutorials/Exercises/01/Begin
```

# Kokkos

## exercise_1_begin.cpp

```
51 #include <sys/time.h>
52
53 // EXERCISE: Include Kokkos_Core.hpp.
54 //             cmath library unnecessary after.
55 // #include <Kokkos_Core.hpp>
56 #include <cmath>
57
58 void checkSizes( int &N, int &M, int &S, int &nrepeat );
```

Argonne
NATIONAL LABORATORY

# Kokkos

## exercise_1_begin.cpp

```
95    // Check sizes.
96    checkSizes( N, M, S, nrepeat );
97
98    // EXERCISE: Initialize Kokkos runtime.
99    // Include braces to encapsulate code between initialize and finalize calls
100   // Kokkos::initialize( argc, argv );
101   // {
102
103   // Allocate y, x vectors and Matrix A:
104   double * const y = new double[ N ];
```

# Kokkos

## exercise_1_begin.cpp

```
108    // Initialize y vector.
109    // EXERCISE: Convert outer loop to Kokkos::parallel_for.
110    for ( int i = 0; i < N; ++i ) {
111      y[ i ] = 1;
112    }
113
114    // Initialize x vector.
115    // EXERCISE: Convert outer loop to Kokkos::parallel_for.
116    for ( int i = 0; i < M; ++i ) {
117      x[ i ] = 1;
118    }
119
120    // Initialize A matrix, note 2D indexing computation.
121    // EXERCISE: Convert outer loop to Kokkos::parallel_for.
122    for ( int j = 0; j < N; ++j ) {
123      for ( int i = 0; i < M; ++i ) {
124        A[ j * M + i ] = 1;
```

# Kokkos

## exercise_1_begin.cpp

```
108    // Initialize y vector.
109    Kokkos::parallel_for( "y_init", N, KOKKOS_LAMBDA ( int i ) {
110      y[ i ] = 1;
111    });
112
113    // Initialize x vector.
114    Kokkos::parallel_for( "x_init", M, KOKKOS_LAMBDA ( int i ) {
115      x[ i ] = 1;
116    });
117
118    // Initialize A matrix, note 2D indexing computation.
119    Kokkos::parallel_for( "matrix_init", N, KOKKOS_LAMBDA ( int j ) {
120      for ( int i = 0; i < M; ++i ) {
121        A[ j * M + i ] = 1;
122      }
123    });
```

Argonne
NATIONAL LABORATORY

# Kokkos

## exercise_1_begin.cpp

```
138        // EXERCISE: Convert outer loop to Kokkos::parallel_reduce.
139     for ( int j = 0; j < N; ++j ) {
140       double temp2 = 0;
141
142       for ( int i = 0; i < M; ++i ) {
143         temp2 += A[ j * M + i ] * x[ i ];
144       }
145
146       result += y[ j ] * temp2;
147     }
148
```

Argonne
NATIONAL LABORATORY

# Kokkos

## exercise_1_begin.cpp

```
138     Kokkos::parallel_reduce( "yAx", N, KOKKOS_LAMBDA ( int j, double &update ) {
139       double temp2 = 0;
140
141       for ( int i = 0; i < M; ++i ) {
142         temp2 += A[ j * M + i ] * x[ i ];
143       }
144
145       update += y[ j ] * temp2;
146     }, result );
```

# Kokkos

## exercise_1_begin.cpp

```
182
183   // EXERCISE: finalize Kokkos runtime
184   // }
185   // Kokkos::finalize();
186
187   return 0;
188 }
```
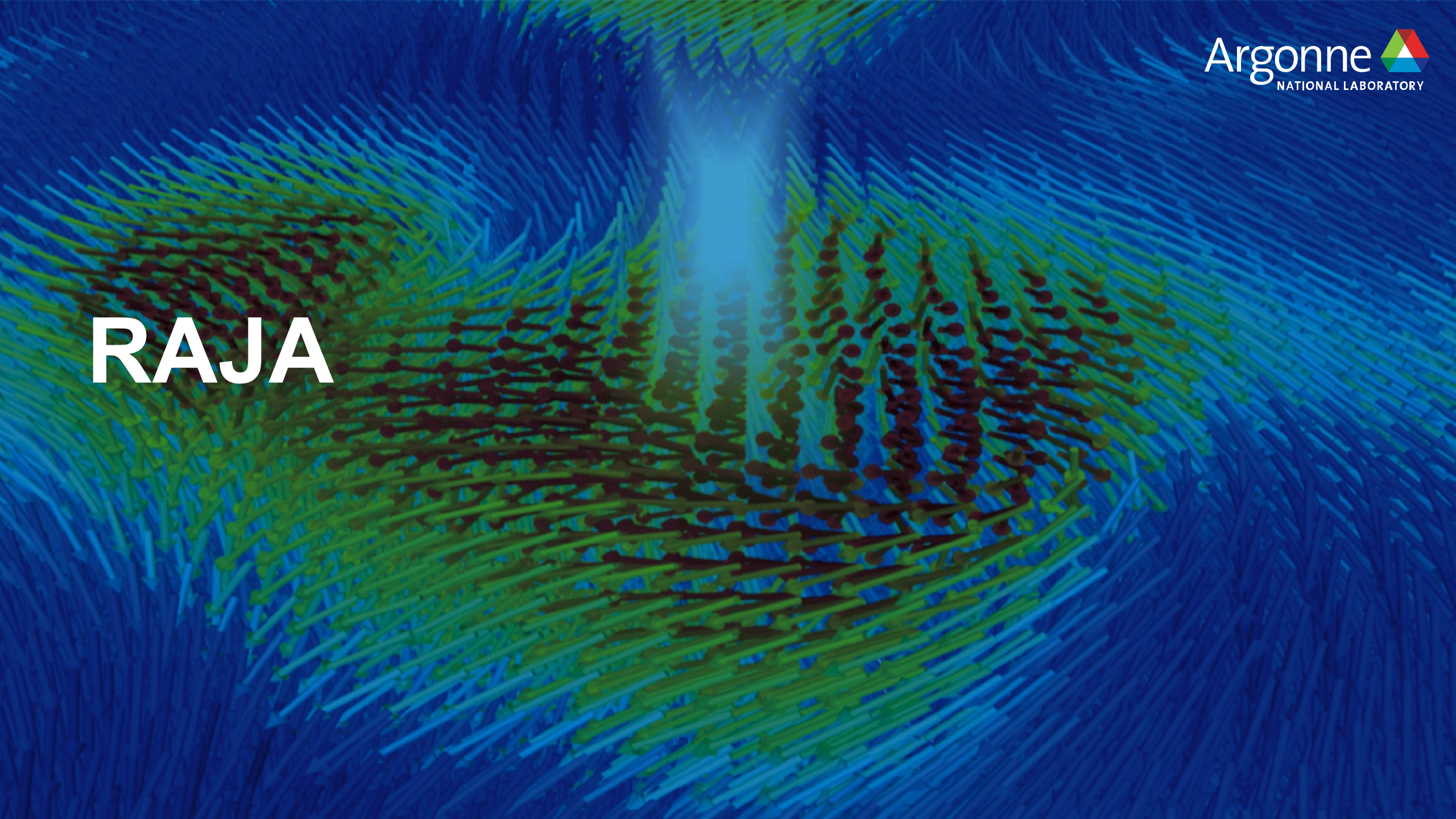
# Kokkos

## Build and Environment

| POLARIS |
| --- |

```
> module swap PrgEnv-nvhpc/8.3.3 PrgEnv-gnu

> module load nvhpc-mixed

> cd /path/to/kokkos-tutorials/02/Begin


> make KOKKOS_PATH=/path/to/kokkos KOKKOS_DEVICES="Cuda" KOKKOS_ARCH="Ampere80"
```

Argonne
NATIONAL LABORATORY

RAJA

Argonne
NATIONAL LABORATORY

# RAJA loop execution has four core concepts

```
using EXEC_POLICY = ...;
RAJA::RangeSegment range(0, N);

RAJA::forall< EXEC_POLICY >( range, [=] (int i)
{
  // loop body...
} );
```

1.  Loop **execution template** (e.g., 'forall')

2.  Loop **execution policy type** (EXEC_POLICY)

3.  Loop **iteration space** (e.g., 'RangeSegment')

4.  Loop **body** (C++ lambda expression)

# Reduction is a common and important parallel pattern

dot product: $dot = \sum_{i=0}^{N-1} a_i \, b_i$, where a and b are vectors, dot is a scalar
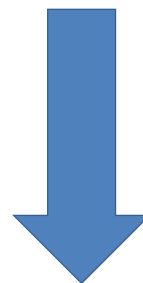
C-style

```
double dot = 0.0;
for (int i = 0; i < N; ++i) {
  dot += a[i] * b[i];
}
```

# A RAJA reduction object hides the complexity of a parallel reduction operation

C-style

```
double dot = 0.0;
for (int i = 0; i < N; ++i) {
    dot += a[i] * b[i];
}
```

RAJA

```
RAJA::ReduceSum< REDUCE_POLICY, double> dot(0.0);

RAJA::forall< EXEC_POLICY >( range, [=] (int i) {
    dot += a[i] * b[i];
} );
```

# Elements of RAJA reductions…

```
RAJA::ReduceSum< REDUCE_POLICY, DTYPE > sum(init_val);

RAJA::forall< EXEC_POLICY >(... {
  sum += func(i);
});

DTYPE reduced_sum = sum.get();
```

- A **reduction type** requires:
  - A reduction policy
  - A reduction value type
  - An initial value

# Elements of RAJA reductions…

```
RAJA::ReduceSum< REDUCE_POLICY, DTYPE > sum(init_val);

RAJA::forall< EXEC_POLICY >(... {
    sum += func(i);
});

DTYPE reduced_sum = sum.get();
```

Note that you cannot access the reduced value inside a kernel.

- A reduction type requires:
  - A reduction policy
  - A reduction value type
  - An initial value

- **Updating reduction value is what you expect (+=, min, max)**

# Elements of RAJA reductions…

```
RAJA::ReduceSum< REDUCE_POLICY, DTYPE > sum(init_val);

RAJA::forall< EXEC_POLICY >(... {
  sum += func(i);
});
```

`DTYPE reduced_sum = sum.get();`

- A reduction type requires:
  - A reduction policy
  - A reduction value type
  - An initial value

- Updating reduction value is what you expect (+=, min, max)

- **After loop runs, get reduced value via 'get' method**

# The reduction policy must be <u>compatible</u> with the loop execution policy

```
RAJA::ReduceSum< REDUCE_POLICY, DTYPE > sum(init_val);

RAJA::forall< EXEC_POLICY >(... {
    sum += func(i);
});


DTYPE reduced_sum = sum.get();
```

An OpenMP execution policy requires an OpenMP reduction policy, similarly for CUDA, etc.

# RAJA provides reduction policies for all supported programming model back-ends

```
RAJA::ReduceSum< REDUCE_POLICY, int > sum(0);

          RAJA::seq_reduce;

          RAJA::omp_reduce;

          RAJA::cuda_reduce;

          RAJA::tbb_reduce;

          RAJA::omp_target_reduce;
```

A sample of RAJA reduction policy types.

Note: OpenMP target and HIP reductions are works-in-progress.

# RAJA supports five common reductions types

```
RAJA::ReduceSum<    REDUCE_POLICY, DTYPE > r(in_val);

RAJA::ReduceMin<    REDUCE_POLICY, DTYPE > r(in_val);

RAJA::ReduceMax<    REDUCE_POLICY, DTYPE > r(in_val);

RAJA::ReduceMinLoc< REDUCE_POLICY, DTYPE > r(in_val,
                                             in_loc);

RAJA::ReduceMaxLoc< REDUCE_POLICY, DTYPE > r(in_val,
                                             in_loc);
```
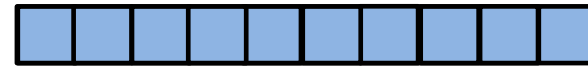
Initial "loc" values

**"Loc" reductions give a loop index where reduced value was found.**

# A RAJA "Segment" defines a set of loop iterates

- A **Segment** is a set of loop indices to <u>run for a kernel</u>

**Contiguous range** [beg, end)

**Strided range** [beg, end, stride)

**List of indices** (indirection)

# A RangeSegment defines a contiguous sequence of indices (stride-1)

```
RAJA::RangeSegment range( 0, N );
```

```
RAJA::forall< RAJA::seq_exec >( range , [=] (int i)
{
  // ...
} );
```

Runs loop indices: 0, 1, 2, … , N-1

# A RangeStrideSegment defines a strided sequence of indices

```
RAJA::RangeStrideSegment srange1( 0, N, 2 );

RAJA::forall< RAJA::seq_exec >( srange1 , [=] (int i)
{
  // ...
} );
```

Runs loop indices: 0, 2, 4, …

# Segments also allow negative indices and strides

```
RAJA::RangeStrideSegment srange2( N-1, -1, -1 );

RAJA::forall< RAJA::seq_exec >( srange2 , [=] (int i)
{
  // ...
} );
```

Runs loop in reverse: N-1, N-2, … , 1, 0

# The RAJA::kernel API is designed for composing and transforming complex kernels

Nested loops

```cpp
using namespace RAJA;
using KERNEL_POL = KernelPolicy<
                      statement::For<1, exec_policy_row,
                        statement::For<0, exec_policy_col,
                          statement::Lambda<0>
                        >
                      >
                    >;
```

```cpp
RAJA::kernel<KERNEL_POL>( RAJA::make_tuple(col_range, row_range),
                          [=](int col, int row ) {
    double dot = 0.0;
    for (int k = 0; k < N; ++k) {
      dot += A(row, k) * B(k, col);
    }
    C(row, col) = dot;
} );
```

Note: lambda expression for inner loop body is the same as C-style variant.

# The RAJA::kernel interface uses four basic concepts, analogous to RAJA::forall

1. Kernel **execution template** ('RAJA::kernel')

2. Kernel **execution policies** (in 'KERNEL_POL')

3. Kernel **iteration spaces** (e.g., 'RangeSegments')

4. Kernel **body** (lambda expressions)

# Each loop level has an iteration space and loop variable

```cpp
using namespace RAJA;
using KERNEL_POL = KernelPolicy<
                      statement::For<1, exec_policy_row,
                        statement::For<0, exec_policy_col,
                          statement::Lambda<0>
                        >
                      >
                    >;



RAJA::kernel<KERNEL_POL>( RAJA::make_tuple(col_range, row_range),
                          [=](int col, int row ) {
    // ...
} );
```

The order (and types) of tuple items and lambda arguments <u>must match</u>.

# Each loop level has an execution policy

```cpp
using namespace RAJA;
using KERNEL_POL = KernelPolicy<
                    statement::For<1, exec_policy_row,
                      statement::For<0, exec_policy_col,
                        statement::Lambda<0>
                      >
                    >
                  >;


RAJA::kernel<KERNEL_POL>( RAJA::make_tuple(col_range, row_range),
                          [=](int col, int row ) {
    // ...
  } );
```

'For' statement integer parameter indicates tuple item it applies to: '0' → col, 1' → row.

# To transform the loop order, change the execution policy, not the kernel code

```
using KERNEL_POL = KernelPolicy<
                        statement::For<1, exec_policy_row,
                            statement::For<0, exec_policy_col,
                            ...
                    >;
```

Outer row loop (1), inner col loop (0)

'For' statements are swapped.

```
using KERNEL_POL = KernelPolicy<
                        statement::For<0, exec_policy_col,
                            statement::For<1, exec_policy_row,
                            ...
                    >;
```

Outer col loop (0), inner row loop (1)

This is analogous to swapping for-loops in a C-style implementation.

# RAJA::KernelPolicy constructs comprise a simple DSL that relies only on standard C++11 support

- A KernelPolicy is built from "Statements" and "StatementLists"

  — A **Statement is an action**: execute a loop, invoke a lambda, synchronize threads, etc. ,

  `For<0, exec_pol, ...>`          `Lambda<0>`          `CudaSyncThreads`

  — A **StatementList is an ordered list of Statements** processed as a sequence; e.g.,

  ```
  For<0, exec_policy0,
      Lambda<0>,
      For<2, exec_policy2,
          Lambda<1>
      >
  >
  ```

**A RAJA::KernelPolicy type is a StatementList.**

# RAJA

```
> git clone --recursive https://github.com/llnl/raja.git
> cd raja/exercises/tutorial_halfday
```

# RAJA

## ex1_vector-addition.cpp

```
105    /// EXERCISE: Implement the vector addition kernel using a RAJA::forall
106    ///           method and RAJA::seq_exec execution policy type.
107    ///
108    /// NOTE: We've done this one for you to help you get started...
109    ///
110
111    using EXEC_POL1 = RAJA::seq_exec;
112
113    RAJA::forall< EXEC_POL1 >(RAJA::RangeSegment(0, N), [=] (int i) {
114      c[i] = a[i] + b[i];
115    });
116
117    checkResult(c, c_ref, N);
```

Argonne
NATIONAL LABORATORY

# RAJA

## ex1_vector-addition.cpp

```
128   std::cout << "\n Running RAJA SIMD vector addition...\n";

129

130   ///

131   /// TODO...

132   ///

133   /// EXERCISE: Implement the vector addition kernel using a RAJA::forall

134   ///           method and RAJA::simd_exec execution policy type.

135   ///

136

137   checkResult(c, c_ref, N);
```

# RAJA

## ex1_vector-addition.cpp

```
128   std::cout << "\n Running RAJA SIMD vector addition...\n";

129

130   using EXEC_POL2 = RAJA::simd_exec;

131

132   RAJA::forall< EXEC_POL2 >(RAJA::RangeSegment(0, N), [=] (int i) {

133     c[i] = a[i] + b[i];

134   });

135

136   checkResult(c, c_ref, N);
```

# RAJA

## ex1_vector-addition.cpp

```
191   std::cout << "\n Running RAJA OpenMP multithreaded vector addition...\n";
192
193   ///
194   /// TODO...
195   ///
196   /// EXERCISE: Implement the vector addition kernel using a RAJA::forall
197   ///           method and RAJA::omp_parallel_for_exec execution policy type.
198   ///
199
200   checkResult(c, c_ref, N);
```

# RAJA

## ex1_vector-addition.cpp

```cpp
191   std::cout << "\n Running RAJA OpenMP multithreaded vector addition...\n";
192
193   using EXEC_POL4 = RAJA::omp_parallel_for_exec;
194
195   RAJA::forall< EXEC_POL4 >(RAJA::RangeSegment(0, N), [=] (int i) {
196     c[i] = a[i] + b[i];
197   });
198
199   checkResult(c, c_ref, N);
```

# RAJA

## ex1_vector-addition.cpp

```
213   std::cout << "\n Running RAJA CUDA vector addition...\n";
214
215   ///
216   /// TODO...
217   ///
218   /// EXERCISE: Implement the vector addition kernel using a RAJA::forall
219   ///           method and RAJA::cuda_exec execution policy type.
220   ///
221
222   checkResult(c, c_ref, N);
```

# RAJA

## ex1_vector-addition.cpp

```
 43 #if defined(RAJA_ENABLE_CUDA)
 44 const int CUDA_BLOCK_SIZE = 256;
 45 #endif
{…}
213   std::cout << "\n Running RAJA CUDA vector addition...\n";
214
215   using EXEC_POL5 = RAJA::cuda_exec<CUDA_BLOCK_SIZE>;
216
217   RAJA::forall< EXEC_POL5 >(RAJA::RangeSegment(0, N),
218                             [=] RAJA_DEVICE (int i) {
219     c[i] = a[i] + b[i];
220   });
221
222   checkResult(c, c_ref, N);
```

Argonne
NATIONAL LABORATORY

# RAJA

## Build and Environment

### POLARIS

```
> module swap PrgEnv-nvhpc/8.3.3 PrgEnv-gnu

> module load nvhpc-mixed cmake

> cd /path/to/raja

> mkdir build

> cd build

> cmake -DCMAKE_BUILD_TYPE=Release -DCMAKE_CXX_COMPILER=g++ -DCMAKE_C_COMPILER=gcc
        -DCUDA_COMMON_OPT_FLAGS="-restrict -arch sm_80 --expt-extended-lambda"
         -C ../host-configs/ubuntu-builds/nvcc_gcc_X.cmake
        -DENABLE_CUDA=On -DRAJA_ENABLE_EXAMPLES=On ..

> Make -j 16
```

Argonne
NATIONAL LABORATORY