

Modern C++ for Parallelism in Scientific Computing

Victor Eijkhout
Texas Advanced Computing Center
eijkhout@tacc.utexas.edu

CppCon 2024

Scientific computing parallelism

- Large amounts of data: often cartesian multi-dimensional arrays, sometimes unstructured d
- Large amounts of parallelism: each element of output array independent.
- No explicit threading parallelism created by some runtime
- Range algorithm notion: do some operation on each element of a dataset

Power method

Let A a matrix of interest
Let x be a random vector
For iterations until convergence
compute the product $y \leftarrow Ax$
compute the norm $\gamma = \|y\|$
normalize $x \leftarrow y/\gamma$

- Method for computing largest eigenvalue of a matrix
- Also Google PageRank
- Stands for many scientific codes: Krylov methods, eigenvalue:

Stencil operations



- This rectangular $m \times n$ thing is the vector
- The 4, ..., stencil is / stands for the matrix.
- Goes by: difference stencil, convolution, Toeplitz matrix



TACC
TEXAS ADVANCED COMPUTING CENTER

TEXAS
The University of Texas at Austin



Array parallelism

Traditional C implementation:

```
1 for ( idxint i=0; i<m; i++)
2   for ( idxint j=0; j<n; j++)
3     out[ IINDEX(i,j,m,n,b) ] = in[ IINDEX(i,j,m,n,b) ]
// seq.ccp
2 #define IINDEX( i,j,m,n,b ) ((i)*b)*(n+2*b) + (j)*b
```

- Two / three-dimensional loop
- all dimensions large
- every output element independent

Reductions

ℓ_2 reduction:

```
1 for ( idxint i=0; i<m; i++)
2   for ( idxint j=0; j<n; j++) {
3     auto v = out[ IINDEX(i,j,m,n,b) ];
4     sum_of_squares += v*v;
5   }
6 return std::sqrt(sum_of_squares);
```

- Parallel except for the accumulation
- Obviously should not be done through at

Stencil computation

Apply stencil to each (i,j) index:

```
1 for ( idxint i=0; i<m; i++) {
2   for ( idxint j=0; j<n; j++) {
3     out[ IINDEX(i,j,m,n,b) ] = 4*in[ IINDEX(i,j,m,n,b) ]
4     - in[ IINDEX(i-1,j,m,n,b) ] - in[ IINDEX(i+1,j,m,n,b) ]
5     - in[ IINDEX(i,j-1,m,n,b) ] - in[ IINDEX(i,j+1,m,n,b) ]
6   }
7 }
```

- Differential operator / image convolution
- Structure can be more complicated in scientific codes

Tools: mdspan and cartesian_product

Data is declared as `mdspan`:

```
1 private:
2   real *_data(nullptr);
3   md::mdspan<
4     real,
5     md::dextents<idxint,>
6     > cartesian_data;
```

```
1 // pointer to the data as 2D array
2 auto& data2d() {
3   return cartesian_data;
4 }
5 const auto& data2d() const {
6   return cartesian_data;
7 }
```

mdspan and cartesian_product

Index range is declared as `range::views::cartesian_product`:

```
1 const auto& s = data2d();
2 int b = this->border();
3 idxint
4 lo_m = static_cast<idxint>(b),
5 hi_m = static_cast<idxint>(s.extent(0)-b),
6 lo_n = static_cast<idxint>(b),
7 hi_n = static_cast<idxint>(s.extent(1)-b);
8 range2d = rng::views::cartesian_product
9   ( rng::views::iota(lo_m,hi_m),rng::views::iota(lo_n,hi_n) );
```

- Vector allocated with size $(m+2b) \times (n+2b)$ to include border
- for handling of boundary conditions / halo regions in PDEs.

Implementation 1: OpenMP parallelism

Annotate loops as parallel and/or reduction:

```
1 #pragma omp parallel for reduction(+:sum_of_squares)
2 for ( idxint i=0; i<m; i++)
3   for ( idxint j=0; j<n; j++) {
4     auto v = out[ IINDEX(i,j,m,n,b) ];
5     sum_of_squares += v*v;
6   }
7 return std::sqrt(sum_of_squares);
8 ;
```

- Static assignment of iterations to threads by default
- Highly controlled affinity
- 'oned' as above, 'clips' for both loops collapsed
- Can be formulated as range algorithm.

Implementation 2: range over indices

Range-based for loop:

```
1 auto array = this->data2d();
2 #pragma omp parallel for reduction(+:sum_of_squares)
3 for ( auto ij : this->inner() ) {
4   auto [i,j] = ij;
5   auto v = array[i,j];
6   sum_of_squares += v*v;
7 }
8 return std::sqrt(sum_of_squares);
9 ;
```

- Range over indices, not over data
- Indices are a subset of the full data!

Stencil operation

Most complicated operation of the bunch:

```
1 // span.ccp
2 auto out = this->data2d();
3 for ( auto ij : this->data2d() ) {
4   #pragma omp parallel for
5   for ( auto ij : this->inner() ) {
6     auto [i,j] = ij;
7     out[ i,j ] = 4*in[ i,j ]
8     - in[ i-1,j ] - in[ i+1,j ] - in[ i,j-1 ] - in[ i,j+1 ]
9   }
10 }
```

- Hard to formulate as range algorithm
- Performance not necessarily determined by floating point op

Implementation 3: Sycl

Open standard, but mostly pushed by Intel

```
1 q.submit([&](handler& h) {
2   accessor Da(buf_a,h,write_only);
3   h.parallel_for
4     (range<2>(msize-2,nsize-2),
5      [=](auto index){
6        auto row = index.get_id(0) + 1;
7        auto col = index.get_id(1) + 1;
8        Da_a[row][col] = 1;
9      });
10  }).wait();
```

- Heterogeneous CPU/GPU code, transparent data movement
- Range algorithm-like syntax, but explicit task queue

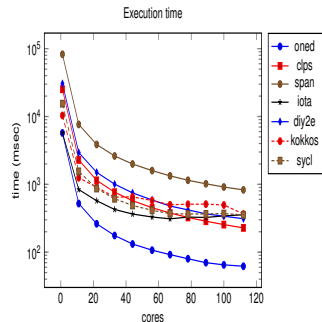
Implementation 4: Kokkos

Open Source heterogeneous execution layer

```
1 Kokkos::parallel_for
2   ("Update x",
3    Kokkos::MDRangePolicy(Kokkos::Rank<2>
4      ((1,1), (msize-1, nsize-1)),
5      Kokkos::LAMBDA(int i, int j) {
6        x(i, j) = Ax(i, j) / norm;
7      })
8   );
```

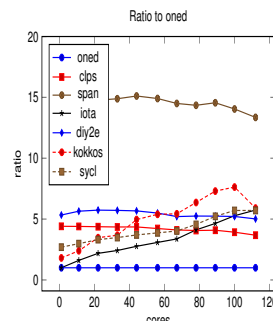
- Same code for CPU and GPU
- Implicit task queue
- Two-dimensional indexing
- Range algorithm-like philosophy

Comparing models (Intel)

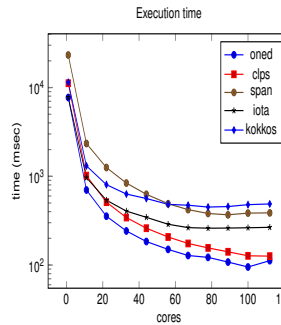


Intel compiler. C-style variant fastest.

Ratio to fastest (Intel)

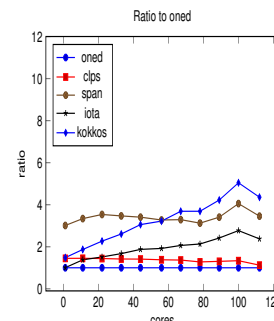


Comparing models (Gcc)



Gcc compiler. less variance between variants

Ratio to fastest (Gcc)



Where do we lose performance?

Hint: perf output on the 'span' variant:

```
1 55.68% [.]
2 --std::ranges::cartesian_product_view::std::ranges::iota_view::long,
3 --loop, std::ranges::iota_view::long
4 --:::iterator::true::operator+=
5 18.73% [.] __div13
6 11.33% [.]
7 --linalg::bordered_array_span::float::central_difference_from
8 5.37% [.]
9 --linalg::bordered_array_span::float::scale_interior
10 5.03% [.] linalg::bordered_array_span::float::l2norm
11 2.69% [.] __div13gpt
```

Index calculations take lots of time.

Conclusion and Acknowledgement

- 'Fancy' schemes suffer from indexing overhead strongly implementation and compiler dependent.
- This work was supported by the Intel oneAPI Center of Excellence, and the TACC STAR Scholars program, funded by generous gifts from TACC industry partners, including Intel, Shell, Exxon, and Chevron.

TACC intel STAR
TEXAS ADVANCED COMPUTING CENTER PARTNERSHIP PROGRAM

- Sycl code contributed by Yojan Chitkara