

Performance analysis of a stencil code in modern C++

Victor Eijkhout, Yojan Chitkara
Texas Advanced Computing Center
`eijkhout@tacc.utexas.edu`

lxPUG2024@SC24

Scientific computing parallelism

- ▶ Large amounts of data:
often cartesian multi-dimensional arrays, sometimes unstructured data
- ▶ 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

Why modern C++

- ▶ Support for multi-D arrays
- ▶ ‘Range algorithms’: iteration without indices or bound checking
- ▶ Native support for parallelism. Ish.
- ▶ Libraries / language extensions for heterogeneous CPU/GPU/FPGA computing.

Basic algorithm

Power method

Let A a matrix of interest

Let x be a random vector

For iterations until convergence

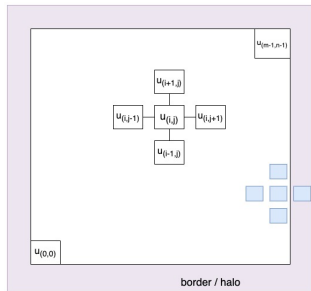
compute the product $y \leftarrow Ax$

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, eigenvalues

Stencil operations



.	-1	.
-1	4	-1
.	-1	.

- ▶ This rectangular $m \times n$ thing is the vector
- ▶ The $4, -1, \dots$ stencil is / stands for the matrix.
- ▶ Goes by: difference stencil, convolution, Toeplitz matrix

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) ] * factor;  
  
1  // seq.cpp  
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 = in[ 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 atomic operation

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

Implementations

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 = in[ 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, 'clps' for both loops collapsed
- ▶ Can be formulated as range algorithm.

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,2>
6             > cartesian_data;

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

1 // base.cpp
2 template< typename real >
3 bordered_array_base<real>::bordered_array_base
4     ( idxint m,idxint n,int border )
5     : _m(m),_n(n),_border(border)
6     , _data( new real[ (m+2*border)*(n+2*border) ] )
7     , data_owned(true)
8     , cartesian_data
9     ( md::mdspan
10     ( _data,md::dextents{m+2*border,n+2*border} )
11     )
```

No performance loss [?].

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 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);
```

- ▶ Range over indices, not over data
(Indices are a subset of the full data!)
- ▶ OpenMP can handle iterators.

Stencil operation

Most complicated operation of the bunch:

```
1 // span.cpp
2 auto out = this->data2d();
3 const auto in = other.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 }
```

- ▶ Hard to formulate as range algorithm
- ▶ Performance not necessarily determined by floating point operations.

Implementation 3: 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      });
```

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

Implementation 4: Sycl

Open standard, but mostly pushed by Intel

```
1  q.submit([&] (handler &h) {  
2      accessor D_a(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               D_a[row][col] = 1.;  
9           });  
10 }) .wait();
```

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

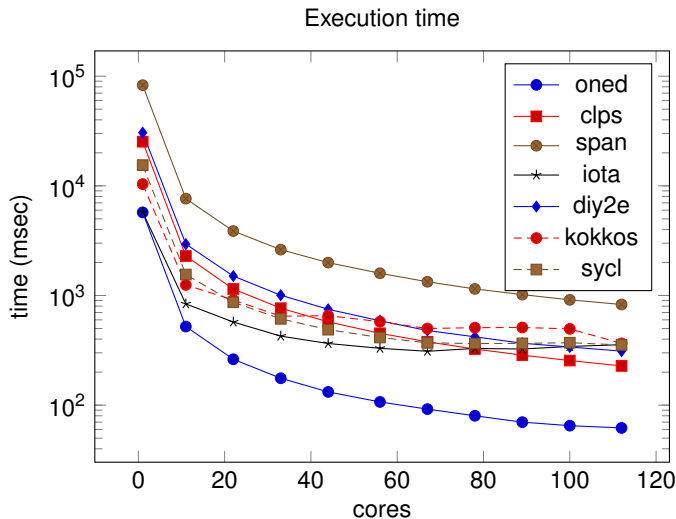
Implementation X: ranges

```
1  #include <range>
2  std::ranges::for_each
3    ( std::execution::par_unseq,
4      some_container,
5      [] ( /* ... */ ) { /* ... */ }
6    );
```

- ▶ ‘Apply this point function to each element of this range’
- ▶ Elegant, trivially parallelized
- ▶ Eh ... 1. relies too much on TBB and such: no affinity control
- ▶ Eh ... 2. hard to account for the halo region
- ▶ Eh ... 3. hard to express stencil operations
(although see <https://www.youtube.com/watch?v=ImM7f5IQ0aw>)
- ▶ Eh ... 4. compiler bug prevented parallel execution of the cartesian product iterator.

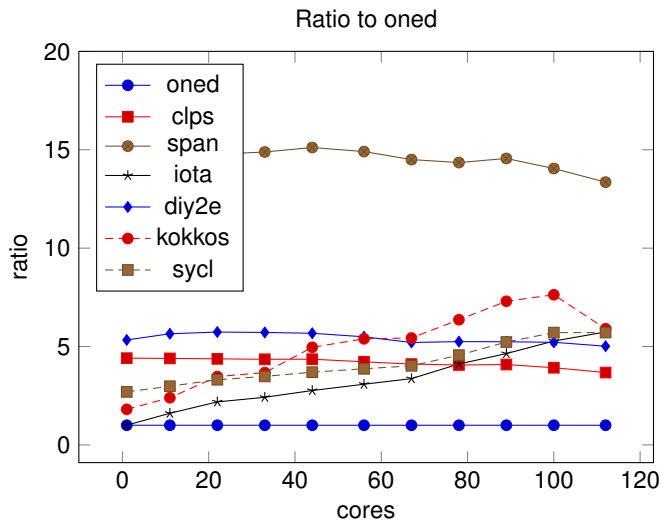
Tests

Comparing models (Intel)

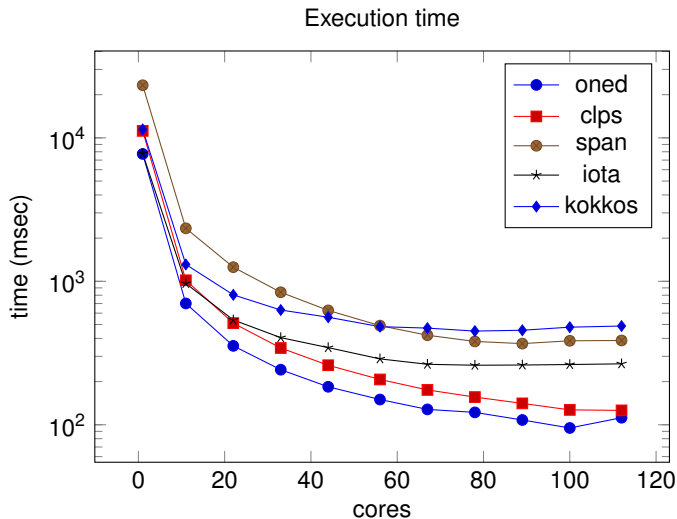


Intel compiler. C-style variant fastest.

Ratio to fastest (Intel)

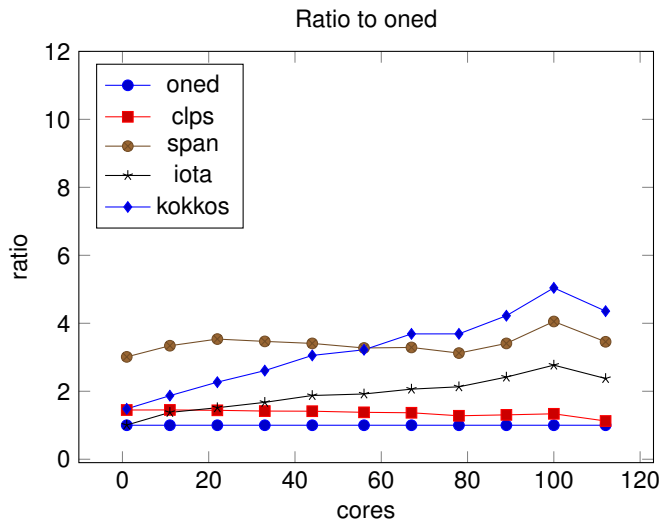


Comparing models (Gcc)



Gcc compiler. less variance between variants

Ratio to fastest (Gcc)



Where do we lose performance?

Hint: perf output on the 'span' variant:

```
1      55.60%  [.]  
        ↪ std::ranges::cartesian_product_view<std::ranges::iota_view<long,  
        ↪ long>, std::ranges::iota_view<long, long>  
        ↪ >::_Iterator<true>::operator+=  
2      18.73%  [.] __divti3  
3      11.33%  [.]  
        ↪ linalg::bordered_array_span<float>::central_difference_from  
4       5.37%  [.]  
        ↪ linalg::bordered_array_span<float>::scale_interior  
5       5.01%  [.] linalg::bordered_array_span<float>::l2norm  
6       2.69%  [.] __divti3@plt
```

Index calculations take lots of time.

GCC vs Intel

Observation:

GCC performance seems better than Intel on the 2D indexing schemes

I'm working on that.

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.



Bibliography



David S. Hollman, Bryce Adelstein Lelbach, H. Carter Edwards, Mark Hoemmen, Daniel Sunderland, and Christian R. Trott.

mdspan in c++: A case study in the integration of performance portable features into international language standards.

In 2019 IEEE/ACM International Workshop on Performance, Portability and Productivity in HPC (P3HPC), volume 74, pages 60–70. IEEE, November 2019.