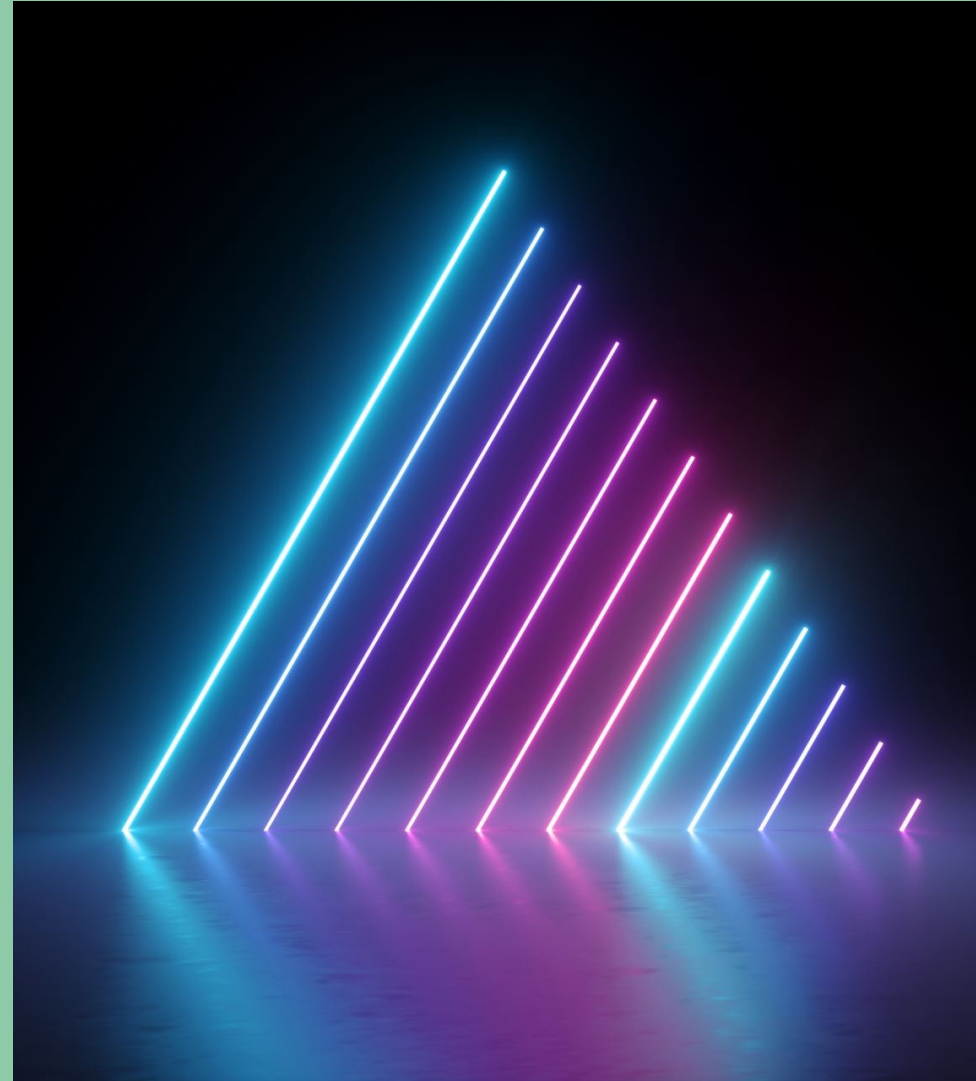


OpenMP-- just a few things about C++

CSC Summer Institute
October 9-11, 2023

- Kent Milfeld (TACC)
- Emanuele Vitali (CSC)



Using OpenMP in C++

- Using OpenMP directives and features in C++
 - Some basics
 - Range-base Iterations
 - Vector and Array Differences
 - Sophisticated Reduction Schemes

Thanks to Victor Eijkhout for creating this course and allowing me to present his work:

https://github.com/VictorEijkhout/TheArtOfHPC_vol2_parallelprogramming/tree/main/examples/omp/cxx

See also his HPC book at:

<https://theartofhpc.com/pcse.html>

Output streams in parallel

- *cout* is thread safe.
- However, a line (<< x << y << z) can be broken apart at each <<.

```
void main(){
    #pragma omp parallel
    {
        int t = omp_get_thread_num();
        int s = t+1;

        cout << s << " from " << t << endl;
    }
}
```

Output streams in parallel

- Solution, form a single string with *stringstream*.

```
#include <sstream>
int main(){
    #pragma omp parallel
    {
        int t = omp_get_thread_num();
        int s = t+1;
        stringstream onestring;
        onestring << s << " from " << t << endl;
        cout << onestring.str();
    }
}
```

Parallel Construct in lambda

- OpenMP parallel directives can create parallel region in lambdas.

```
int main(){  
    const int s = [] () {  
        int s;  
        #pragma omp parallel  
        #pragma omp master  
        s = 2 * omp_get_num_threads();  
        return s;  
    } ();  
  
    cout << s << endl;  
    return 0;  
}
```

Class Methods: dynamic scope

- Dynamic scope holds for class methods as for any other function:

```
class C {  
public:  
    void class_func() {  
        #pragma omp for  
        for(int i=0; i<8; i++) {  
            cout<< omp_get_thread_num() ;}  
        }  
    };  
  
int main() {  
    C c;  
    #pragma omp parallel num_threads(4)  
    c.class_func();  
}
```

OUTPUT:

23320110

(order varies)

Privatizing class members

Class members can only be privatized from (non-static) class methods.

In this example `f` can not be static:

```
1  // private.cxx
2  class foo {
3  private:
4      int x;
5  public:
6      void f() {
7          #pragma omp parallel private(x)
8              somefunction(x);
9      };
10 };
```

You can not privatize just a member:

```
1  // privateno.cxx
2  class foo { public: int x; };
3  int main() {
4      foo thing;
5      #pragma omp parallel private(thing.x) // NOPE
```

Parallel loops

Range based loops

Library

Questions

1. Do regular OpenMP loops look different in C++?
2. Is there a relation between OpenMP parallel loops and iterators?
3. OpenMP parallel loops vs parallel execution (algorithm) policies.

Range syntax

Parallel loops in C++ can use range-based syntax as of OpenMP-5.0

```
1  // vecdata.cxx
2  vector<float> values(100);
3
4  #pragma omp parallel for
5  for ( auto& elt : values ) {
6      elt = 5.f;
7  }
8
9  float sum{0.f};
10 #pragma omp parallel for reduction(+:sum)
11 for ( auto elt : values ) {
12     sum += elt;
13 }
```

Performance (not shown here) is the same as C code.

C dynamic (“array”) storage

- C dynamic storage use pointers:
pointer manipulation by thread requires privatization.

```
#define Nthreads 4
int main() {
    int *array = (int*) malloc(Nthreads*sizeof(int));

    for (int i=0; i<Nthreads; i++) array[i] = 0;

    #pragma omp parallel firstprivate(array) num_threads(4)
    { int t = omp_get_thread_num();
      array += t;
      array[0] = t;
    }

    for(int i=0; i<Nthreads; i++)
        cout <<i<< " "<<array[i]<<endl;
}
```

OUTPUT:

```
0 0
1 1
2 2
3 3
```

C++ vectors

- C++ vector: copy constructor copies data
(the same applies to an array: `int array[Nthreads];`)

```
#include <vector>
#define Nthreads 4

int main() {
    vector<int> array(Nthreads, 9);

    #pragma omp parallel firstprivate(array) num_threads(4)
    { int t = omp_get_thread_num(); array[t] = t; }
    for (auto i: array) cout << i;  cout<<endl;

    #pragma omp parallel num_threads(4)
    { int t = omp_get_thread_num(); array[t] = t; }
    for (auto i: array) cout << i;  cout<<endl;
}
```

OUTPUT:

9999
0123

OpenMP can parallelize any loop over
a C++ construct that has a 'random-access' iterator.

Ranges

Ranges: An abstraction of "something iterable", requiring `begin()` and `end()` on the range (iterable).

Views: ranges that are defined by another range, that are used to transform the underlying range.

```
#include <ranges>
using namespace std;

int main() {
    vector<int> vec = {1, 2, 3, 4, 5, 6, 7, 8};
    auto vw = vec | views::reverse;
    cout << *vw.begin() << std::endl;
}
```

vw is a view: it does not change vec and it does not store any elements
-- the creation time and size of vw is independent of the size of vec.

Ranges library / `std::ranges::zip_view`

Range access

<code>begin</code>	<code>rbegin</code>	<code>size</code>	<code>empty</code>
<code>cbegin</code>	<code>crbegin</code>	<code>ssize</code>	
<code>end</code>	<code>rend</code>	<code>data</code>	
<code>cend</code>	<code>crend</code>	<code>cdata</code>	

Range conversions

to (C++23)

Range primitives

<code>range_size_t</code>	<code>iterator_t</code>	<code>range_const_reference_t</code> (C++23)	<code>elements_of</code> (C++23)
<code>range_difference_t</code>	<code>const_iterator_t</code> (C++23)	<code>range_rvalue_reference_t</code>	
<code>range_value_t</code>	<code>sentinel_t</code>	<code>range_common_reference_t</code>	
<code>range_reference_t</code>	<code>const_sentinel_t</code> (C++23)		

Dangling iterator handling

<code>dangling</code>	<code>borrowed_iterator_t</code>	<code>borrowed_subrange_t</code>
-----------------------	----------------------------------	----------------------------------

Range concepts

<code>range</code>	<code>common_range</code>	<code>input_range</code>	<code>bidirectional_range</code>
<code>borrowed_range</code>	<code>view</code>	<code>output_range</code>	<code>random_access_range</code>
<code>sized_range</code>	<code>viewable_range</code>	<code>forward_range</code>	<code>contiguous_range</code>
<code>constant_range</code> (C++23)			

Views

<code>view_interface</code>	<code>subrange</code>
-----------------------------	-----------------------

Range factories

<code>empty_view</code>	<code>iota_view</code>	<code>basic_istream_view</code>
<code>views::empty</code>	<code>views::iota</code>	<code>views::istream</code>
<code>single_view</code>	<code>repeat_view</code> (C++23)	
<code>views::single</code>	<code>views::repeat</code> (C++23)	

Range adaptors

<code>views::all_t</code>	<code>drop_view</code>	<code>reverse_view</code>	<code>adjacent_view</code> (C++23)
<code>views::all</code>	<code>views::drop</code>	<code>views::reverse</code>	<code>views::adjacent</code> (C++23)
<code>ref_view</code>	<code>drop_while_view</code>	<code>as_const_view</code> (C++23)	<code>views::pairwise</code> (C++23)
<code>owning_view</code>	<code>views::drop_while</code>	<code>views::as_const</code> (C++23)	<code>adjacent_transform_view</code> (C++23)
<code>as_rvalue_view</code> (C++23)	<code>join_view</code>	<code>elements_view</code>	<code>views::adjacent_transform</code> (C++23)
<code>views::as_rvalue</code> (C++23)	<code>views::join</code>	<code>views::elements</code>	<code>views::pairwise_transform</code> (C++23)
<code>filter_view</code>	<code>join_with_view</code> (C++23)	<code>keys_view</code>	<code>chunk_view</code> (C++23)
<code>views::filter</code>	<code>views::join_with</code> (C++23)	<code>views::keys</code>	<code>views::chunk</code> (C++23)
<code>transform_view</code>	<code>lazy_split_view</code>	<code>values_view</code>	<code>slide_view</code> (C++23)
<code>views::transform</code>	<code>views::lazy_split</code>	<code>views::values</code>	<code>views::slide</code> (C++23)
<code>take_view</code>	<code>split_view</code>	<code>enumerate_view</code> (C++23)	<code>chunk_by_view</code> (C++23)
<code>views::take</code>	<code>views::split</code>	<code>views::enumerate</code> (C++23)	<code>views::chunk_by</code> (C++23)
<code>take_while_view</code>	<code>views::counted</code>	<code>zip_view</code> (C++23)	<code>stride_view</code> (C++23)
<code>views::take_while</code>	<code>common_view</code>	<code>views::zip</code> (C++23)	<code>views::stride</code> (C++23)
	<code>views::common</code>	<code>zip_transform_view</code> (C++23)	<code>cartesian_product_view</code> (C++23)
		<code>views::zip_transform</code> (C++23)	<code>views::cartesian_product</code> (C++23)

Ranges (requires header)

The C++20 ranges library is also supported:

```
1  #      pragma omp parallel for reduction(+:count)
2      for ( auto e : data
3          | std::ranges::views::drop(1) )
4          count += e;
5  #      pragma omp parallel for reduction(+:count)
6      for ( auto e : data
7          | std::ranges::views::transform
8          ( []( auto e ) { return 2*e; } ) )
9          count += e;
```

C++ ranges speedup

```
1  ==== Run range on 1 threads ====
2  sum of vector: 50000005000000 in 6.148
3  sum w/ drop 1: 50000004999999 in 6.017
4  sum times 2 : 100000010000000 in 6.012
5  ==== Run range on 25 threads ====
6  sum of vector: 50000005000000 in 0.494
7  sum w/ drop 1: 50000004999999 in 0.477
8  sum times 2 : 100000010000000 in 0.489
9  ==== Run range on 51 threads ====
10 sum of vector: 50000005000000 in 0.257
11 sum w/ drop 1: 50000004999999 in 0.248
12 sum times 2 : 100000010000000 in 0.245
13 ==== Run range on 76 threads ====
14 sum of vector: 50000005000000 in 0.182
15 sum w/ drop 1: 50000004999999 in 0.184
16 sum times 2 : 100000010000000 in 0.185
17 ==== Run range on 102 threads ====
18 sum of vector: 50000005000000 in 0.143
19 sum w/ drop 1: 50000004999999 in 0.139
20 sum times 2 : 100000010000000 in 0.134
21 ==== Run range on 128 threads ====
22 sum of vector: 50000005000000 in 0.122
23 sum w/ drop 1: 50000004999999 in 0.11
```


Ranges and indices

Use `iota_view` to obtain indices– be careful with `auto`:

```
std::vector<int> data(N,2);

#pragma omp parallel for
for(auto i : std::ranges::iota_view{0UZ, data.size()})
    data[i]=f(i);
std::cout<< data[0] << " " << data[N-1] << std::endl;
```

Note C++23 suffix `UZ` is used to type `0` as an unsigned `size_t`.
For older version use:

```
iota_view{static_cast<size_t>(0),data.size()} )
```

Custom iterators

Recall that

Short hand:

```
1  vector<float> v;  
2  for ( auto e : v )  
3      ... e ...
```

for:

```
1  for ( vector<float>::iter  
2      e=v.begin();  
3      e!=v.end(); e++ )  
4      ... *e ...
```

If we want

```
1  for ( auto e : my_object )  
2      ... e ...
```

we need a sub-class for the iterator with methods such as *begin*, *end*, *** and *++*.

Probably also *+=* and *-*

Custom iterators

OpenMP can parallelize any range-based loop with a random-access iterator.

Class:

```
1  // iterator.cxx
2  template<typename T>
3  class NewVector {
4  protected:
5      T *storage;
6      int s;
7  public:
8      // iterator stuff
9      class iter;
10     iter begin();
11     iter end();
12 };
```

Main:

```
1  NewVector<float> v(s);
2  #pragma omp parallel for
3  for ( auto e : v )
4      cout << e << " ";
```

Custom iterators, exercise

Required iterator methods:

(1)

```
1 NewVector<T>::iter& operator++();
2 T& operator*();
3 bool operator==( const NewVector::iter &other ) const;
4 bool operator!=( const NewVector::iter &other ) const;
5 // needed to OpenMP
6 int operator-
7     ( const NewVector::iter& other ) const;
8 NewVector<T>::iter& operator+=( int add );
```

This is a little short of a full random-access iterator; the difference depends on the OpenMP implementation

Write the missing iterator methods.

Here's something to get you started.

(2)

```
1 template<typename T>
2 class NewVector<T>::iter {
3 private: T *searcher;
4 };
5 template<typename T>
6 NewVector<T>::iter::iter( T* searcher )
7     : searcher(searcher) {};
8 template<typename T>
9 NewVector<T>::iter NewVector<T>::begin() {
10     return NewVector<T>::iter(storage); };
11 template<typename T>
12 NewVector<T>::iter NewVector<T>::end() {
13     return NewVector<T>::iter(storage+NewVector<T>::s); };
```

(3)

```
1 template<typename T>
2 bool NewVector<T>::iter::operator==
3     ( const NewVector<T>::iter &other ) const {
4     return searcher==other.searcher; };
5 template<typename T>
6 bool NewVector<T>::iter::operator!=
7     ( const NewVector<T>::iter &other ) const {
8     return searcher!=other.searcher; };
9 template<typename T>
10 NewVector<T>::iter& NewVector<T>::iter::operator++() {
11     searcher++; return *this; };
12 template<typename T>
13 NewVector<T>::iter& NewVector<T>::iter::operator+=( int add )
14     searcher += add; return *this; };
```

(4)

```
1 template<typename T>
2 T& NewVector<T>::iter::operator*() {
3     return *searcher; };
4 // needed for OpenMP
5 template<typename T>
6 int NewVector<T>::iter::operator-
7     ( const NewVector<T>::iter& other ) const {
8     return searcher-other.searcher; };
```


OpenMP vs standard parallelism

Application: prime number marking (load unbalanced)

```
1  #pragma omp parallel for schedule(guided,8)
2  for ( int i=0; i<nsize; i++) {
3      results[i] = one_if_prime( number(i) );
4  }
```

```
1  // primepolicy.cxx
2  transform( std::execution::par,
3             numbers.begin(),numbers.end(),
4             results.begin(),
5             [] (int n ) -> int {
6                 return one_if_prime(n); }
7             );
```

Standard parallelism uses Thread Building Blocks (TBB) as backend

Timings

```
1  ==== Run primepolicy on 1 threads ====
2  OMP: found 0 primes; Time:      390 msec (threads= 1)
3  TBB: found 0 primes; Time:      392 msec
4  ==== Run primepolicy on 25 threads ====
5  OMP: found 0 primes; Time:       17 msec (threads=25)
6  TBB: found 0 primes; Time:       19 msec
7  ==== Run primepolicy on 51 threads ====
8  OMP: found 0 primes; Time:        9 msec (threads=51)
9  TBB: found 0 primes; Time:       13 msec
10 ==== Run primepolicy on 76 threads ====
11 OMP: found 0 primes; Time:         6 msec (threads=76)
12 TBB: found 0 primes; Time:        15 msec
13 ==== Run primepolicy on 102 threads ====
14 OMP: found 0 primes; Time:         5 msec (threads=102)
15 TBB: found 0 primes; Time:        71 msec
16 ==== Run primepolicy on 128 threads ====
17 OMP: found 0 primes; Time:         4 msec (threads=128)
18 TBB: found 0 primes; Time:        55 msec
```

Reductions vs standard parallelism

Application: prime number counting (load unbalanced)

```
1  #pragma omp parallel for schedule(guided,8) reduction(+:prime_count)
2  for ( auto n : numbers ) {
3      prime_count += one_if_prime( n );
4  }
```

```
1  // reducepolicy.cxx
2  prime_count = transform_reduce
3      ( std::execution::par,
4        numbers.begin(), numbers.end(),
5        0,
6        std::plus<>{},
7        [] ( int n ) -> int {
8            return one_if_prime(n); }
9        );
```


Timings

```
1  ==== Run reducepolicy on 1 threads ====
2  OMP: found 9592 primes; Time:      390 msec (threads= 1)
3  TBB: found 9592 primes; Time:      392 msec
4  ==== Run reducepolicy on 25 threads ====
5  OMP: found 9592 primes; Time:       17 msec (threads=25)
6  TBB: found 9592 primes; Time:       20 msec
7  ==== Run reducepolicy on 51 threads ====
8  OMP: found 9592 primes; Time:        8 msec (threads=51)
9  TBB: found 9592 primes; Time:       13 msec
10 ==== Run reducepolicy on 76 threads ====
11 OMP: found 9592 primes; Time:         6 msec (threads=76)
12 TBB: found 9592 primes; Time:        23 msec
13 ==== Run reducepolicy on 102 threads ====
14 OMP: found 9592 primes; Time:         5 msec (threads=102)
15 TBB: found 9592 primes; Time:       105 msec
16 ==== Run reducepolicy on 128 threads ====
17 OMP: found 9592 primes; Time:         4 msec (threads=128)
18 TBB: found 9592 primes; Time:       54 msec
```

Reductions

Questions

1. Are simple reductions the same as in C?
2. Can you reduce `std::vector` like an array?
3. Precisely what can you reduce?
4. Any interesting examples?
5. Compare reduction to native C++ mechanisms.

reductions on scalars

Same as in C,
you can now use range syntax for the loop.

```
1  // range.cxx
2  #      pragma omp parallel for reduction(+:count)
3      for ( auto e : data )
4          count += e;
```

Output streams in parallel

- Arrays (`a[N]`) can be used in reductions: specified as an array section (`a[:N]`).
- Vectors (`vector<T> v(N)`) can also be used, but vectors array section (`v[:N]`)
- are not allowed— must use a pointer to the vector

```
int aums[N];
#pragma omp parallel for reduction(+ : aums[0:N])    ALLOWED

vector<int> vsums(N,0);
#pragma omp parallel for reduction(+ : vsums[0:N]) NOT ALLOWED
```

```
vector<int>      vsums(N,0);    //vector of sums
int *ptr_vsums = vsums.data(); //pointer to sums

#pragma omp parallel for reduction(+ : ptr_vsums[:N]) \
                        schedule (static, 1)
for (    int it=0; it<n; it++){
    for (int iv=0; iv<N; iv++)
        ptr_vsums[iv]+=iv*data[it];
}
```

reduction on class objects

Reduction can be applied to any class for which the reduction operator is defined as `operator+` or whichever operator the case may be.

```
1  // reductclass.cxx
2  class Thing {
3  private:
4      float x{0.f};
5  public:
6      Thing() = default;
7      Thing( float x ) : x(x) {};
8      Thing operator+
9          ( const Thing& other ) {
10         return Thing( x + other.x );
11     };
12 };
```

```
1  vector< Thing >
2      things(500,Thing(1.f) );
3  Thing result(0.f);
4  #pragma omp parallel for \
5      reduction( +:result )
6  for ( const auto& t : things )
7      result = result + t;
```

A default constructor is required for the internally used init value.

Uninitialized containers

Multi-socket systems:

parallel initialization instantiates pages on sockets:
'first touch'

```
1      double *x = (double*)malloc( N*sizeof(double));
2      #pragma omp parallel for
3      for (int i=0; i<N; i++)
4          x[i] = f(i);
5
```

This does not work with

```
1      std::vector<double> x(N);
2      #pragma omp parallel for
3      for (int i=0; i<N; i++)
4          x[i] = f(i);
5
```

because of value initialization in the `vector` container.

Uninitialized containers (cont. 1)

Trick to create a vector of uninitialized data:

```
1 // heatalloc.cxx
2 template<typename T>
3 struct uninitialized {
4     uninitialized() {}
5     T val;
6     constexpr operator T() const {return val;};
7     T operator=( const T&& v ) { val = v; return val; };
8 };
```

so that we can create vectors that behave normally:

```
1 vector<uninitialized<double>> x(N),y(N);
2
3 #pragma omp parallel for
4 for (int i=0; i<N; i++)
5     y[i] = x[i] = 0.;
6 x[0] = 0; x[N-1] = 1.;
```

(Question: why not use *reserve*?)

Uninitialized containers (cont. 2)

Easy way of dealing with this:

```
1  template<typename T>
2  class ompvector : public vector<uninitialized<T>> {
3  public:
4      ompvector( size_t s )
5          : vector<uninitialized<T>>::vector<uninitialized<T>>(s) {};
6  };
```