Document Number: D0350R5

Date: 2023-02-13

Reply-to: Matthias Kretz < m.kretz@gsi.de>

Audience: LEWG

Target: C++23 (pending simd merge)

INTEGRATING SIMD WITH PARALLEL ALGORITHMS

Λ.	D^{0}	T	\Box	Λ.	\frown	Ī
А	D	ווכ	K/	4	CI	

This paper discusses a new execution policy for integrating simd with parallel algorithms.

CONTENTS

D0350R5

1 CHANGELOG

1.1 CHANGES FROM REVISION 0

Previous revision: P0350R0

• Update to apply against C++17 wording.

- Removed executors discussion because the executors design has not left SG1 yet.
- Updated example code to reflect changes in PO214.

1.2 CHANGES FROM REVISION 1

Previous revision: P0350R1

- Updated code to match [N4744].
- Fixed a bug in the for_each example implementation.
- Improved iota and for_each example implementations with constexpr-if.
- Discuss impact on all algorithms.

1.3 CHANGES FROM REVISION 2

Previous revision: P0350R2

- Discuss ABI tag of std::generate callables.
- Add a Tony Table.
- Note that remove and remove_copy are implicitly vectorizable.

1.4 CHANGES FROM REVISION 3

Previous revision: P0350R3

- Add heap algorithms to list of algorithms that may benefit from new execution policy.
- Discuss the 3 options for generate.
- Settle on simd_mask predicates.

CHANGES FROM REVISION 4

1.5

Previous revision: P0350R4

Discuss algorithms in <numeric>. (transform_reduce is a bit more involved because of different value_types)

- Mention <memory> algorithms (they don't call user-defined code).
- Add pointer to PO214 for motivation.
- Fix mention of swap_ranges, which was swap before.
- Add C++20's shift_left and shift_right.
- make_heap, push_heap, pop_heap, and sort_heap have no execution policy overload. Do fix.
- Discuss permission to fall back to sequential execution (i.e. T vs. simd<T, scalar>).

2

STRAW POLLS

2.1 sgl at oulu

Poll: Ship it to LEWG?

SF	F	N	А	SA
6	6	2	0	0

Poll: Forward the paper to LWG?

SF

→ Paper needs a revision: LEWG wants a list of affected algorithms and an update to concept requirements.

2.3

2.2

LEWG AT COLOGNE 2019

LEWG AT ALBUQUERQUE

Poll: Prefer expressing simd execution via policy tag (not newly named algorithms, std::simd_transform. etc).

SF	F	N	А	SA
4	9	3	0	0

D0350R5 3 Introduction

Poll: What do we do with predicates? (Roughly option 2 from the paper, be mindful of the amount of wording that is thrown out of conformance.)

→ unanimous consent

Poll: Require that this goes first into std::ranges algorithms (which haven't been parallelized yet).

Revise and come back.

3 INTRODUCTION

Parallel Algorithms enable implementations of the existing STL algorithms to use non-sequential semantics when executing the user-supplied code (explicit callable or implicit operator call). The first argument to the algorithm function determines this change in execution semantics via an *execution policy*. This paper introduces a new execution policy, called execution::simd¹. execution::simd requires user-provided function objects to be callable with simd<T, Abi> arguments instead of the T arguments the execution::seq variant would use. The algorithm therefore processes chunks of simd<T, Abi>::size() objects concurrently. The execution order of the chunks retains the sequential semantics of the non-parallel algorithms.

As a consequence, the applicability of the execution policy is limited to iterators where Iterator::value_-type is a vectorizable type [N4744]. A future extension of simd may lift this restriction by allowing certain (or all) user-defined types as first template argument to simd. A different conceivable extensions is a recursive destructuring applied inside the algorithm, subsequent creation of a corresponding number of simd objects, and a call to the function object with a corresponding number of arguments. (E.g. application of an algorithm on std::vector<std::pair<float, float>> calls the function object with simd<float>, simd<float> instead of simd<std::pair<float, float>>.)

4

PARALLEL ALGORITHMS

4.1 EXAMPLE

Consider the example in Listing ??. The iota and for_each functions each could create an internal simd iterator adaptor, depending on the iterator category. Being able to determine whether the storage, the iterator points to, is contiguous, is most important in this context as it enables vector loads and stores. Since the std::vector iterators are contiguous iterators, the example implementations shown in Listing ?? and Listing ?? could be used for the example.

¹ An alternative suggestion for the name is execution::simd_type.

```
before
                                                                           with D0350R5
using V = stdx::native_simd<float>;
                                                                           using V = stdx::native_simd<float>;
constexpr int N = 60;
                                                                           constexpr int N = 60;
template <class T> T something(T);
                                                                           template <class T> T something(T);
auto f(const std::array<float, N>& data)
                                                                           auto f(const std::array<float, N>& data)
  std::array<float, N> output;
                                                                             std::array<float, N> output;
 size_t i = 0;
for (; i + V::size() <= N; i += V::size()) {
   V x(&data[i], stdx::element_aligned);
   x = something(x + 1);</pre>
                                                                             \verb|stdx::transform(std::execution::simd|,
                                                                               data.begin(), data.end(), output.begin(),
                                                                               [](auto x) {
                                                                                 return something(x + 1);
   x.copy_to(&output[i], stdx::element_aligned);
 for (; i < N; ++i) {
  output[i] = something(data[i] + 1);
}</pre>
 return output;
                                                                             return output;
```

Tony Table 1: Transformation of an array. c.f. https://godbolt.org/z/mEL3CK

```
std::vector<float> data;
data.resize(99);
iota(execution::simd, data.begin(), data.end(), 0.f);
for_each(execution::simd, data.begin(), data.end(), [](auto &x) {
    x *= x;
});
```

Listing 1: Example using execution::simd with iota and for_each.

```
template <size_t N, class ContiguousIterator>
inline void epilogue(ContiguousIterator first, ContiguousIterator last,
                        typename ContiguousIterator::value_type first_value) {
    if constexpr (N > 0) {
      if (distance(first, last) >= N) {
         using T = ContiguousIterator::value_type;
        using V = simd<T, simd_abi::deduce_t<T, N>>;
        const V init = V([&](auto i) { return T(i); }) + first value;
8
        store(init, std::addressof(*first), element_aligned);
        first += V::size();
10
       epilogue<V::size() / 2>(first, last, init[V::size() - 1] + 1);
12
    }
13
14 }
15
  template <class ContiguousIterator>
16
void iota(execution::simd_policy, ContiguousIterator first, ContiguousIterator last,
             typename ContiguousIterator::value_type first_value) {
18
    using T = ContiguousIterator::value_type;
    using V = native simd<T>;
    V init = V([&](auto i) { return T(i); }) + first_value;
    const V stride = T(V::size());
    for (; distance(first, last) >= V::size(); first += V::size(), init += stride) {
      store(init, std::addressof(*first), element_aligned);
24
25
    epilogue < V::size() / 2>(first, last, init[V::size() - 1] + 1);
```

Listing 2: Implementation idea for the iota function used in Listing ??.

```
template <size_t N, class ContiguousIterator, class UnaryFunction>
   inline void epilogue(ContiguousIterator first, ContiguousIterator last,
                         UnaryFunction f) {
     if constexpr (N > 0) {
       using T = ContiguousIterator::value_type;
       using V = simd<T, simd_abi::deduce_t<T, N>>;
       if (distance(first, last) >= V::size()) {
         V tmp(std::addressof(*first), element_aligned);
         f(tmp);
         if constexpr (is_functor_argument_mutable_v < UnaryFunction, V >) {
10
           store(tmp, std::addressof(*first), element_aligned);
11
13
       epilogue < V::size() / 2>(first, last, f);
14
15
  }
16
17
  template <class ContiguousIterator, class UnaryFunction>
18
   void for_each(execution::simd_policy, ContiguousIterator first,
                 ContiguousIterator last, UnaryFunction f) {
20
21
     using V = native_simd<ContiguousIterator::value_type>;
     for (; distance(first, last) >= V::size(); first += V::size()) {
22
       V tmp(std::addressof(*first), element_aligned);
23
       f(tmp);
       if constexpr (is_functor_argument_mutable_v < UnaryFunction, V >) {
25
         store(tmp, std::addressof(*first), element_aligned);
26
27
     }
     epilogue < V::size() / 2>(first, last, f);
```

Listing 3: Implementation idea for the for_each function used in Listing ??.

Both implementations might be improved with a prologue that enables aligned loads and stores. Also note that for_each allows the Function parameter to mutate the argument if the iterator is a mutable iterator. The implementation uses a compile-time trait to determine whether the function f uses a reference parameter, in which case it stores the temporary simd object back. Otherwise, the store is optimized away.

Figure ?? shows a visualization how the iota implementation works. The init simd object is stored via vector stores to 4 (assuming native simd::size() == 4) elements in the std::vector. In each iteration the init object is incremented by simd::size() and stored to the following elements in the std::vector. Since the std::vector has 99 elements, the last three elements cannot be initialized with a vector store of four elements. Instead the epilogue recursion generates a new init simd object for size 2 and subsequently for size 1.

Figure ?? visualizes the end of the for_each implementation. The main for loop processes four elements of the std::vector in parallel. It executes a vector load, calls the user-provided function with the temporary simd object, and executes a vector store back to the same memory location. The remaining three elements are again handled by an epilogue recursion which divides the number of processed elements by 2 with every step.

For both algorithms it would be perfectly valid to implement the epilogue as a sequential loop using simd objects with size 1.

4.2 DISCUSSION OF ALGORITHMS

COPIES In general, the execution::simd policy requires algorithms to make a copy from the input sequence. For now, since simd only supports arithmetic types and simd does not return Ivalue references to its values, it is not observable whether a copy was made. With two exceptions:

- Modification of the input sequence via different means than the function parameter(s), which is UB anyway, will not modify the value of the function parameter(s).
- Using mutable iterators, assignment to the simd (Ivalue reference) parameter of the user-supplied function object will not modify the output sequence until after the function has returned (cf. Listing ??).

Note that most non-modifying sequence operations allow modification of the sequence by using a non-const lvalue reference parameter for the user-supplied function object.

PREDICATES Algorithms that take a predicate returning a bool have two possible vectorization strategies:

1. The predicate still returns bool. In this case, every predicate must execute a simd_mask reduction. This makes it simple to short-circuit in the algorithm implementation but may unnecessarily restrict the achievable parallelization.

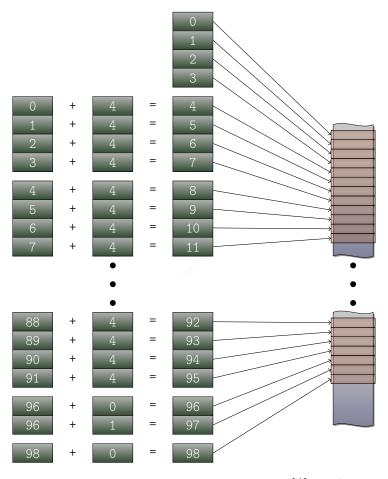


Figure 1: Visualization of chunking the iota call with $\mathcal{W}_{\mathrm{T}}=4$ in Listing ??.

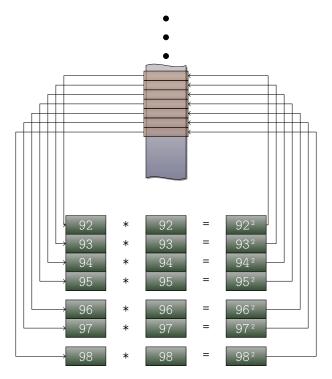


Figure 2: Visualization of chunking the for_each call with $\mathcal{W}_{\scriptscriptstyle T}=4$ in Listing ??.

2. The predicate returns $simd_mask$. In this case $\mathcal{W}_{ForwardIterator::value_type}$ reductions can happen in parallel. Short-circuiting is still possible, but requires a $simd_mask$ reduction on each step (QoI question).

It would also be possible to allow both and let the algorithm switch the strategy depending on the return type of the predicate.

In Cologne 2019, LEWG unanimously recommended to exclusively go with predicates returning simd_mask: The use of bool would effectively change the algorithms.

COMPLEXITY REQUIREMENTS For many algorithms, the complexity requirement states "Applies f exactly last - first times". In the execution::simd case, the number of applications of f is reduced by an unspecified factor.

Sorting The Compare function object type is required to return a value that is contextually convertible to bool. For sorting, it is important that overloads using the execution::simd policy work with simd_mask instead of bool. It is not useful for the sort algorithm to know whether all/any/some/none of the compared values are "less than". It requires a mask object to know the "less than" relation for each individual value.

D0350R5 4 PARALLEL ALGORITHMS

4.3 DESIGN ALTERNATIVE

In Cologne 2019, LEWG recommended to not pursue this design alternative. It is still provided in this paper for completeness.

There are subtle differences in how the execution::simd specializations need to be used (e.g. std::generate currently requires the generator function to return objects that can be assigned to a dereferenced ForwardIt; the execution::simd specialization requires the generator function to return objects of type simd<ForwardIt::value_type>). An attempt to fit execution::simd_policy into the existing wording results in some special-casing in the algorithm specifications. This observation leads to the question whether a new execution policy is really the best approach. The alternative would be a duplication of algorithms to variants with a simd_ prefix in their name. Example:

```
simd_for_each(data.begin(), data.end(), [](auto &x) {
  x *= x;
});
```

This alternative would not reduce the amount of wording/complexity though, since now a lot of the algorithm wording would need to be duplicated. However, this would allow a very simple reduction of the number of algorithms that support simd execution.

4.4 AFFECTED ALGORITHMS

The following algorithms have an ExecutionPolicy overload and can work with a execution::simd_policy specialization:

- all_of, any_of, none_of
- for_each, for_each_n
- find, find_if, find_if_not
- find_end
- find_first_of
- adjacent_find
- count, count_if
- mismatch
- equal
- search, search_n

- copy, copy_n (no real need; can be implicitly vectorized)
- copy_if
- swap (no real need; can be implicitly vectorized)
- transform
- replace, replace_if, replace_copy, replace_copy_if
- fill, fill_n (no real need; can be implicitly vectorized)
- generate_n (see Section ??)
- remove, remove_copy (no real need; can be implicitly vectorized)
- remove_if, remove_copy_if
- unique, unique_copy
- reverse, reverse_copy (no real need; can be implicitly vectorized)
- rotate, rotate_copy (no real need; can be implicitly vectorized)
- is_partitioned, partition, stable_partition, partition_copy, partition_point
- sort,stable_sort, partial_sort, partial_sort_copy, is_sorted, is_sorted_until
- nth_element
- merge, inplace_merge
- includes, set_union, set_intersection, set_difference, set_symmetric_difference
- min_element, max_element, minmax_element
- lexicographical_compare
- is_heap, is_heap_until, make_heap, push_heap, pop_heap, and sort_heap (The comparison function object can use simd and simd_mask.)

The remaining algorithms have no obvious use for the specialization:

• move makes no sense until we can create simd<T> types for pointers (likely) and class types (less likely).

lower_bound, upper_bound, equal_range, and binary_search may benefit from simd usage, but currently do not provide ExecutionPolicy overloads.

D0350R5 4 PARALLEL ALGORITHMS

```
std::vector<float> v;
3 // O. Existing (scalar) interface
4 | std::generate(v.begin(), v.end(), [&]() { return 0.f; })
6 // 1. Choose ABI via template parameter
7 // template <class Gen>
8 // void generate(/*...*/, Gen g) { g.template operator()<native_simd<float>>(); }
std::generate(v.begin(), v.end(), []<typename T>() { return T(); });
11 // 2. Choose ABI via return type
12 // template <class Gen>
// void generate(/*...*/, Gen g) {
14 //
      auto x = g(); // discard if x has more values than needed
15 // }
std::generate(v.begin(), v.end(), []() { return native_simd<float>(); });
18 // 3. Choose ABI via unused parameter
19 // template <class Gen>
20 // void generate(/*...*/, Gen g) {
21 // constexpr native_simd<float>* tag = nullptr;
22 //
       g(tag);
23 // }
std::generate(v.begin(), v.end(), []<typename T>(T*) { return T(); });
```

Listing 4: Use std::generate to fill a container with random numbers.

4.5 THE GENERATE ALGORITHM

The generator function passed to generate/generate_n does not expect any arguments and thus has no interface for the algorithm to request a certain ABI tag from the function (template). Consequently, there are three ideas how to make it work for simd:

- 1. Require the generator function object to take a template argument (no function arguments).
- 2. Let the algorithm implementation cope with the return type defined by the generator function object. It may have to discard values if the number of values in the range is not a multiple of the number of values in the return type.
- 3. Pass a parameter that is not used other than for deducing the expected return type of the generator function object.

See Listing ?? for examples.

Ideas 1 and 3 require the user to supply a type with a template call operator. Idea 2 may lead the user to unintentionally drop generated values. For the user, idea 3 requires to write boilerplate with no apparent use: it's an implementation detail shining through. Idea 1 was preferred in the discussion in LEWG in Cologne 2019. However, it is still unclear whether the need for explicit

instantiation of the call operator in the standard library implementation is acceptable. There is no apparent technical reason not to use this variant.

Add a new execution policy to [N4842]:

// 20.18.7, unsequenced execution policy
class unsequenced_policy;

// 20.18.8, simd execution policy
class simd_policy;

// 20.18.89, execution policy objects:
inline constexpr sequenced_policy seq{ unspecified };
inline constexpr parallel_policy par{ unspecified };
inline constexpr parallel_unsequenced_policy par_unseq{ unspecified };
inline constexpr unsequenced_policy unseq{ unspecified };
inline constexpr unsequenced_policy unseq{ unspecified };
inline constexpr simd_policy simd{ unspecified };

Renumber §20.18.8 to §20.18.9 and add §20.18.8 [execpol.simd]:

class simd_policy { unspecified };

The class simd_policy is an execution policy type used as a unique type to disambiguate parallel algorithm overloading and indicate that a parallel algorithm's execution may be vectorized using simd for interfacing with user-provided functionality.

During the execution of a parallel algorithm with the execution::simd_policy policy, if the invocation of an element access function exits via an uncaught exception, terminate() shall be called.

Add to §20.18.9 [execpol.objects]:

inline constexpr execution::simd_policy execution::simd{ unspecified };

[N4842] defines requirements on user-provided function objects. This might be the right place to add:

_ §25.3.2 [algorithms.parallel.user]

- ² Function objects passed into parallel algorithms instantiated with the execution::simd execution policy shall:
 - be callable with arguments of type simd<Iterator::value_type, Abi>, for any ABI tag Abi, for all arguments that otherwise would be of type Iterator::value_type;

- return objects of type simd<Iterator::value_type, Abi>, if the function object is otherwise expected to return objects assignable to a dereferenced Iterator object;
- return objects of type simd_mask<Iterator::value_type, Abi> or bool, if the function object is otherwise expected to return bool.

The following subsection in [N4842] defines the semantics of the execution policies. A new paragraph for execution::simd is needed. The intent is to

- 1. constrain execution to the calling thread,
- 2. allow implementations to assume unordered access for all internal element access functions (most importantly loads and stores),
- 3. apply user-provided function objects in the order the simd chunks are created from sequential iteration over the iterator(s).

§25.3.3 [algorithms.parallel.exec] The invocations of element access functions in parallel algorithms invoked with an execution policy object of

type execution::simd_policy are permitted to execute in an unordered fashion in the calling thread, except for the application of user-provided function objects. User-provided function objects are called with an unspecified number of sequence elements combined into a simd<T, Abi> object. The type for Abi is chosen by the implementation. It may be different for subsequent applications of the user-provided function in the same parallel algorithm invocation. The type for T is the decayed type of the sequence elements. The order of elements in the simd object is equal to the order of the corresponding elements in the sequence argument. The invocation order of user-provided function objects is sequential.

It is my understanding that we do not want to add anything to [N4842] at this point. The situation

is simpler for the execution::simd policy. It is almost equivalent to the seg policy.

4.7 WORDING FOR INDIVIDUAL ALGORITHMS §25.7 [alg.sorting]

² Compare is a function object type. The return value of the function call operation applied to an object of type Compare, when contextually converted to bool, yields true if the first argument of the call is less than the second, and false otherwise. If the ExecutionPolicy is execution::simd_policy, the return type of the function call operation applied to an object of type Compare is a specialization of simd_mask. Its i-th element in the simd_mask yields true if the value of the i-th element of the first argument of the call is less than the corresponding element of the second, and false otherwise. Compare comp is used throughout for algorithms assuming an ordering relation. It is assumed that comp will not apply any non-constant function through the dereferenced iterator.

Further wording work is necessary where individual algorithms refer to boolean results from predicates. (E.g. all of returns false if any of (E) is false ...)