



# Understanding, Using, and Improving `std::generator`

Johannes Kalmbach

2025

## About Myself (Johannes Kalmbach)

- C++ programmer since 2017
- PhD student at University of Freiburg (Germany)
- Co-founder of QLeverize
- First time speaker + attendee
- [johannes.kalmbach@gmail.com](mailto:johannes.kalmbach@gmail.com)
- <https://github.com/joka921>

# std::generator - A Simple Example

```
1 import std; // Save slide space.  
2  
3 std::generator<size_t> fibonacci_gen() {  
4     size_t i = 0, j = 1;  
5     while (true) {  
6         co_yield j;  
7         i = std::exchange(j, i + j);  
8     }  
9 }  
10  
11 int main() {  
12     for (auto fib : fibonacci_gen() | std::views::take(10)) {  
13         std::println("yielded {}", fib);  
14     }  
15 }
```

*A std::generator generates a sequence of elements by repeatedly resuming the coroutine from which it was returned. Each time a co\_yield statement is evaluated, the coroutine produces one element of the sequence.*

*C++ Standard, section 26.8.1*

## **std::generator is...**

- ... a convenient way to produce values one by one
- ... a `std::ranges::input_range`
- ... a coroutine mechanism (using C++20 `<coroutine>`s)
- ... a C++23 library type (can be implemented using C++20)
- ... inherently type-erased (like all coroutines)

## `std::generator` is type-erased

```
1 // A coroutine of type `generator<size_t>`  
2 generator<size_t> fibonacci_gen() {  
3     /* as before, including `co_yield` */  
4 }  
5  
6 // Another coroutine of type `generator<size_t>`  
7 generator<size_t> forty_two_gen() {  
8     co_yield 42;  
9 }  
10  
11 // An ordinary function that returns `generator<size_t>`  
12 generator<size_t> choose_gen(bool b) {  
13     return b ? fibonacci_gen() : forty_two_gen();  
14 }
```

# The Template Interface

```
1 template <typename Ref, typename V = void, typename Allocator = void>
2 class generator;
```

- Typical case: Ref specified, V = void
- Infer reference\_type and value\_type from Ref

Ref	reference_type	value_type
string	string&&	string
string&&	string&&	string
string&	string&	string
const string&	const string&	string

# The Template Interface

```
1 template <typename Ref, typename V = void, typename Allocator = void>
2 class generator;
```

- If both Ref and V are specified, they are reference\_type and value\_type of the generator
- Needed for proxy reference types, e.g.

```
generator<tuple<const string&, const vector<int>&>, tuple<string, vector<int>>>
```

## Caveat

```
1 std::generator<std::string> gen() {  
2     std::string s;  
3     // This line silently creates a copy.  
4     co_yield s;  
5 }
```

- `reference_type` is an rvalue-reference
- An lvalue-reference is `co_yielded`
- Then the yielded value is copied
- Mandated by the standard, currently not found on `cppreference.com`

## Which of Those Have Issues?

```
1 generator<string> g1() {  
2     std::string s;  
3     co_yield s;  
4 }  
5  
6 generator<string&&> g2() {  
7     std::string s;  
8     co_yield s;  
9 }  
10  
11 generator<const string&> g3() {  
12     std::string s;  
13     co_yield s;  
14 }
```

```
1 generator<string&> g4() {  
2     const std::string s;  
3     co_yield s;  
4 }  
5  
6 generator<string&> g5() {  
7     std::string s;  
8     co_yield s;  
9 }
```

## Which of Those Have Issues?

```
1 generator<string> g1() {  
2     std::string s;  
3     // Silent copy  
4     co_yield s;  
5 }  
6  
7 generator<string&&> g2() {  
8     std::string s;  
9     // Silent copy  
10    co_yield s;  
11 }  
12  
13 generator<const string&> g3() {  
14     std::string s;  
15     co_yield s;  
16 }
```

```
1 generator<string&> g4() {  
2     const std::string s;  
3     // Fails to compile.  
4     co_yield s;  
5 }  
6  
7 generator<string&> g5() {  
8     std::string s;  
9     co_yield s;  
10 }
```

# Custom Allocators

- Coroutines need to allocate
- Default is `std::allocator`

```
1 template<typename T>
2 class LoggingAllocator {
3     constexpr T* allocate(size_t n) {
4         std::println("Allocating {} bytes", n * sizeof(T));
5         // Actual allocation...
6     }
7
8     constexpr void deallocate(T* p size_t n) { /* Analogously */}
9 }
10 };
11
12 std::generator<int, void, LoggingAllocator<int>> gen() {
13     co_yield 23;
14     co_yield 64;
15 }
16
17 int main() {
18     { gen(); }
19     // Prints "Allocating 96 bytes" on GCC-14, and "...64 bytes" on Clang-19
20 }
```

## std::ranges::elements\_of

```
1 // Example from `cppreference.com`
2 #include <generator>
3
4 template<typename T>
5 struct Tree
6 {
7     T value;
8     Tree *left{}, *right{};
9
10    std::generator<const T&> traverse_inorder() const
11    {
12        if (left)
13            co_yield std::ranges::elements_of(left->traverse_inorder());
14
15        co_yield value;
16
17        if (right)
18            co_yield std::ranges::elements_of(right->traverse_inorder());
19    }
20};
```

- More efficient than for (auto&& el : range) {co\_yield std::forward<decltype(el)>(el);}

# Performance

# Performance / Codegen

## According to Trusted Community Members

*Coroutines are simple and elegant, but the codegen is atrocious – definitely unsuitable for hot paths.*

[Vittorio Romeo on his blog](#)

*... Well, that's relative ...*

[Phil Nash on CppCast, starting at 16:00](#)

*... I mean, that's relative and it depends on your code and you have to measure and blablabla ...*

Timur Doumler, same place

*... so let's measure ...*

## A Simple Benchmark

Compare various implementations of fibonacci\_gen:

```
1 // Generator, body visible in header
2 std::generator<size_t> fib_gen() { /* as before */}
3
4 // Can be completely inlined
5 struct Inlineable {
6     size_t i = 0, j = 1;
7     size_t next() {
8         auto res = j;
9         i = std::exchange(j, i + j);
10        return res;
11    }
12
13 // Indirection via function call
14 struct NonInlineable {
15     size_t i = 0, j = 1;
16     size_t next(); // Defined in different TU
17 }
18
19 // Indirection via virtual function
20 struct RangeOfSizeTBase {
21     virtual size_t next() = 0;
22     virtual ~RangeOfSizeTBase = default;
23 }
24 std::unique_ptr<RangeOfSizeTBase> getVirtualFibRange();
```

## Initial Benchmark Results

- Using `google::benchmark`
- AMD Ryzen 9 7950x
- time in nanoseconds/element

	GCC-14	Clang-19
<code>std::generator</code>	2.7	1.7
<code>inlineable class</code>	0.6	0.6
<code>non-inlineable class</code>	2.6	1.7
<code>virtual function</code>	2.6	1.7

# Performance Analysis

- `std::generator` has to allocate
- Allocation implies indirection
- ... and hinders inlining
- The overhead introduced can be
  - Negligible / Acceptable
  - Significant / Inacceptable
  - Unavoidable
- HALO is allowed for coroutines...
- ... but neither GCC nor Clang currently apply it ([Clang used to](#))
- Everything above holds for all type-erasure mechanisms.

## **Summary: When to Consider std::generator**

- When you can afford the overhead
- When you need type-erasure anyway
  - See excellent talks by Klaus Iglberger about type erasure
- When you need to define a complex input range
  - Writing iterators manually is tedious
  - Where possible, consider std::ranges/std::views

# Mitigating Performance Issues

Idea: Let the generator yield a batch of values at once

```
1 generator<const vector<size_t>&> batched_fib() {
2     std::vector<size_t> buffer;
3     buffer.reserve(BUF_SIZE);
4     size_t i = 0, j = 1;
5     while (true) {
6         size_t next = j;
7         i = std::exchange(j, i + j);
8         buffer.push_back(next);
9         if (buffer.size() == BUF_SIZE) {
10             co_yield buffer;
11             buffer.clear();
12         }
13     }
14 }
```

```
1 void consume_using_join_view() {
2     for (size_t s : batched_iota()
3          | std::views::join) {
4         // Do something.
5     }
6 }
7
8 void consume_using_nested_loop() {
9     for (const auto& batch : batched_iota()) {
10        for (size_t s : batch) {
11            // Do something
12        }
13    }
14 }
```

## Updated Benchmark Results

	GCC-14	Clang-19
std::generator	2.7	1.7
inlineable class	0.6	0.6
non-inlineable class	2.6	1.7
virtual function	2.6	1.7
batched, join_view	1.1	1.1
batched, nested loop	1.1	1.1

# Removing The Boilerplate, Version 1

- We can use a macro...
- ... or maybe reflection in the future?

```
1 #define CO_YIELD_BATCHED(buffer, el) \
2     buffer.push_back(el); \
3     if (buffer.size() == BUF_SIZE) { \
4         co_yield buffer; \
5         buffer.clear(); \
6     } \
7     void(0)
8
9 generator<const vector<size_t>&> batched_fib() {
10    std::vector<size_t> buffer;
11    buffer.reserve(BUF_SIZE);
12    size_t i = 0, j = 1;
13    while (true) {
14        size_t next = j;
15        i = std::exchange(j, i + j);
16        CO_YIELD_BATCHED(buffer, i);
17    }
18 }
```

# Removing the Boilerplate, Version 2

- Previously: Manual batching of values (hidden by macro)
- Can we move the boilerplate into the generator type?

```
1 #include "./batched_generator.hpp"
2
3 // Only choose a different return type, everything else stays the same.
4 // Second template argument is the batch size
5 batched_generator<size_t, 100> batched_fib() {
6     size_t i = 0, j = 1;
7     while (true) {
8         // Yield values one by one
9         size_t next = j;
10        i = std::exchange(j, i + j);
11        // Implicit batching by the machinery.
12        co_yield next;
13    }
14 }
15
16 int main() {
17     for (size_t s : batched_fib()) {
18         // Do something useful.
19         std::println("{}", s);
20     }
21 }
```

## The promise\_type (Simplified)

```
1 // `Yielded` is the corresponding generators yielded type
2 // (e.g. `string&&` for `generator<string>`, or
3 // `const string&` for `generator<const string&>`).
4 template<typename Yielded>
5 class PromiseType {
6     // Store a pointer to the most recently yielded value
7     add_pointer_t<Yielded> value_ = nullptr;
8
9     // Called for each `co_yield` statement. Suspend the coroutine,
10    // s.t. the value can be used by the `generator`.
11    suspend_always yield_value(Yielded val) noexcept {
12        value_ = std::addressof(val);
13        return {};
14    }
15
16    /* Omitted: Special overload for `co_yield` that creates a copy if necessary */
17
18    // Suspend the coroutine at the beginning and the end.
19    suspend_always initial_suspend() const noexcept { return {}; }
20    std::suspend_always final_suspend() noexcept { return {}; }
21
22    /* Omitted: boilerplate code for exception handling and `co_return` */
23};
```

## The Batched promise\_type (Simplified)

```
1 template<typename T, size_t BatchSize = 100>
2 class BatchedPromiseType {
3     static_assert(is_object_v<T>);
4
5     std::vector<T> buffer_;
6     // Called for each `co_yield` statement.
7     // Suspend the coroutine only if the `buffer_` is full.
8     template <typename U>
9     SuspendIf yield_value(U&& val) noexcept {
10         buffer_.emplace_back(std::forward<U>(val));
11         return SuspendIf{buffer_.size() >= BatchSize};
12     }
13
14     // Has to be called when advancing the iterator
15     void resetBuffer() {buffer_.clear(); buffer_.reserve(BatchSize)}
16
17     // Suspend the coroutine at the beginning and the end.
18     suspend_always initial_suspend() const noexcept { return {}; }
19     std::suspend_always final_suspend() noexcept { return {}; }
20
21     /* Omitted: boilerplate code for exception handling and `co_return` */
22 };
```

## SuspendIf

- Consists purely of special functions that the <coroutine> framework expects

```
1 struct SuspendIf {
2     bool suspend_;
3
4     // If `true`, then the coroutine is NOT suspended
5     constexpr bool await_ready() noexcept { return !suspend_; }
6
7     // When suspending, return control to the caller
8     void await_suspend(std::coroutine_handle<>) noexcept {
9         return;
10    }
11    // Nothing to do when resuming the coroutine
12    constexpr void await_resume() const noexcept {}
13};
```

## The generator Class

```
1 template<typename Ref, typename Val>
2 class generator : public ranges::view_interface<generator<Ref, Val>> {
3     /* Omitted typedefs for `Value`, `Reference`, `Yielded`. */
4     using promise_type = PromiseType<Yielded>;
5
6 public:
7     /* Omitted: move operations (implemented) and copy operations (deleted) */
8
9     ~generator() { if (coro_) { coro_.destroy(); }
10
11    // Iterator `resumes` the coroutine via the handle and accesses
12    // yielded values via the `promise`.
13    struct Iterator;
14    Iterator begin();
15    std::default_sentinel_t end() const noexcept;
16 private:
17    coroutine_handle<promise_type> coro_;
18
19    generator(coroutine_handle<promise_type> coro) noexcept
20        : coro_{move(coro)} {}
21};
```

## The batched\_generator Class

```
1 template<typename T>
2 class batched_generator : public ranges::view_interface<generator<T> {
3     using promise_type = BatchedPromiseType<T>;
4 public:
5     /* Omitted: move operations, constructor, destructor (same as before) */
6
7     struct Iterator;
8     Iterator beginInternal() { return Iterator{coro_}; }
9
10    void emplaceJoin() {
11        if (!join_) {
12            join_.emplace(std::views::join(std::subrange{beginInternal(), std::default_sentinel}));
13        }
14    }
15
16    auto begin() {
17        emplaceJoin();
18        return join_->begin();
19    }
20    auto end() {
21        emplaceJoin();
22        return join_->end();
23    }
24 private:
25     std::optional<std::join_view<std::subrange<Iterator, std::default_sentinel_t>>> join_;
26     coroutine_handle<promise_type> coro_;
27 }
```

## The Valley of Tears

- `batched_generator` showed no performance improvements...
- Codegen of generator is very vulnerable to small changes
- Chicken vs. Egg

# Disassembly of Hot Path, iota\_gen

```
1 .L20:  
2     leaq    1(%rax), %rcx      ; i++  
3     movq    %rax, 80(%rbx)    ; Store old `i` to RAM  
4     movq    %rcx, 64(%rbx)    ; Store `i` to RAM  
5     cmpq    %rdx, %rbp        ; buf.size() < buf.capacity()  
6     je     .L10              ; Jump to reallocation  
7     movq    %rax, 0(%rbp)      ; push_back (*it = i)  
8     addq    $8, %rbp          ; it++  
9     movq    %rbp, 24(%rbx)    ; Store `it` to RAM  
10    .L11:  
11     movq    %rbp, %rax  
12     subq    %r13, %rax  
13     cmpq    $792, %rax       ; buf.size() >= 100  
14     seta    %al  
15     movb    %al, 72(%rbx)    ; Store `SuspendIf::suspend`  
16     testb   %al, %al  
17     jne     .L21              ; Jump to suspension  
18     movq    64(%rbx), %rax    ; Load `i` from RAM  
19     jmp     .L20              ; Back to beginning of loop
```

# Disassembly of Hot Path, iota\_gen, Manually Improved

```
1 .L20:  
2  
3     movq    %rax, 0(%rbp)    ; push_back()  
4     addq    $8,  %rbp        ; ++it;  
5  
6     movq    %rbp, %rcx  
7     subq    %r13, %rcx  
8     cmpq    $792, %rcx       ; if (buf.size() >= 100  
9     ja     .L21_pre         ; Jump to suspension  
10  
11    leaq    1(%rax), %rax   ; ++i;  
12    jmp     .L20            ; Back to beginning of loop  
13  
14 /* dumping of registers before suspend ... */
```

- Performance as expected/desired
- Compilers need to learn this + HALO

## Final Thoughts

- `std::generator` is useful and convenient
- `std::generator` has its performance issues
  - Be nice to your compiler vendors, maybe we get consistent HALO one day
  - ... and better codegen for the coroutine state machines
- `(std::)generator` is a good starting point for your `<coroutine>` journey
  - Available in the standard library
  - Not too complex (understandable for mere mortals)

mod