



Regular Expressions

Hana Dusíková

Who am I?

- I live in Brno, Czechia.
- I'm a researcher at Avast.
- Organizer of **Avast C++ Prague Meetup**.
- (soon-ish) Czech national body to WG21.
- Amateur photographer and occasional hiker.

Let's talk about Regular Expressions...

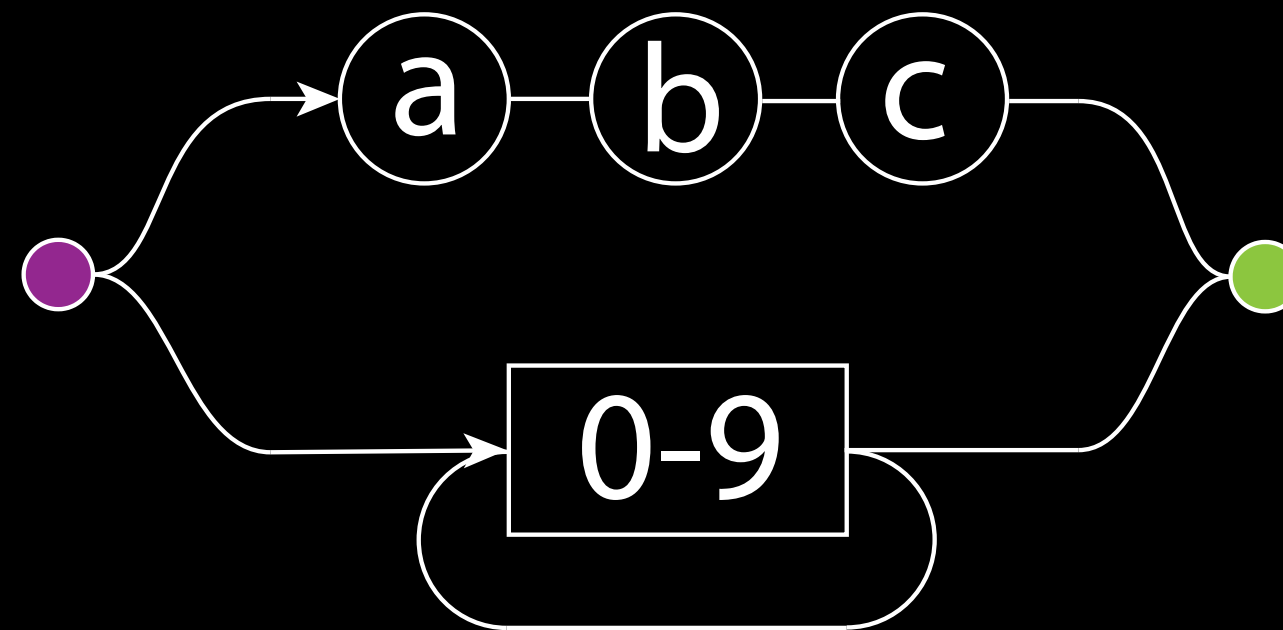
What can I do with a regular expression?*

- **Match** the subject against a pattern.
- **Extract** the content of the subject.
- **Modify** parts of the subject.

*Depending on algorithm used internally.

What does a regular expression look like*

abc|[0-9]+



*There are many slightly different dialects (Perl, ECMA, Posix).

How can I use regular expressions in C++?

```
// compile the RE only once
std::regex re{"abc|[0-9]+"};

bool match(std::string_view sv) {
    return std::regex_match(sv.begin(), sv.end(), re);
}
```


How can I use regular expressions in C?

```
const unsigned char * p = "[0-9]{4,16}?[aA]";
int err = 0;
size_t errOff = 0;

pcre2_code * re = pcre2_compile(p, PCRE2_ZERO_TERMINATED,
    0, &err, &errOff, NULL);

int match(const char * str, size_t sz) {
    return pcre2_match(re, str, sz, 0, 0, NULL, NULL);
}

// TODO: don't forget to call pcre2_code_free at the end ;)
```

What's going on in the constructor?

- Syntax parsing of the pattern.
- Building a data-structure (tree or table) representing the pattern.
- The calculations **are very expensive** (matching is not).

How can I avoid calling the constructor?

You can write a **finite state machine** by yourself.

But that's hard to write properly.

You can easily **introduce an error**.

You can use an **external utility** to generate the parser.

Really? There is no other way?

Compile Time Regular Expressions

Hana Dusíková

What do I want?

- Work with REs as with any other code.
- Don't pay any **runtime cost** for parsing the pattern.
- Have REs as quick as possible with **zero-overhead**.

How do I think REs should be used in C++?

```
bool match(std::string_view sv) noexcept {  
    return std::regex_match<"abc|[0-9]+">(sv);  
}
```

This is possible to implement with C++20's language support for Non Type Template Parameters.

The **C**ompile **T**ime **R**egular **E**xpression library with **C++20**

```
bool match(std::string_view sv) noexcept {  
    return ctcre::match<"abc|[0-9]+>(sv);  
}
```

Disabled for compilers without **class NTTP** support.

The **C**ompile **T**ime **R**egular **E**xpression library with **C++17**

```
static constexpr ctcre::fixed_string  
    ptn = "abc|[0-9]+";  
  
bool match(std::string_view sv) noexcept {  
    return ctcre::match<ptn>(sv);  
}
```

The **C**ompile **T**ime **R**egular **E**xpression library with **C++17** and **N3599** extension*

```
bool match(std::string_view sv) noexcept {  
    using namespace ctre::literals;  
    return "abc|[0-9]+"_ctre.match(sv);  
}
```

Example: Constexpr

```
static_assert (  
    ctcre::match<"[0-9]+\\. [0-9]+"> ("123.456"sv)  
);
```

Example: Constexpr

```
static_assert (  
    ctcre::match<"[0-9]+\\. ([0-9]+)"> ("123.456"sv)  
    .get<1>() == "456"sv  
);
```

Example: Extracting a date from a string

```
struct date { unsigned year, month, day; };

std::optional<date> extract(std::string_view sv) noexcept {
    if (auto re = ctcre::match<"(\\n+) / (\\n{1,2}) / (\\n{1,2})">(sv)) {
        return date{
            conv(re.get<0>()),
            conv(re.get<1>()),
            conv(re.get<2>())
        };
    }
    return std::nullopt;
}
```

Example: Extracting a name from a CSV-like input

```
struct name { std::string_view first, family; };

std::optional<name> extract(std::string_view sv) noexcept {
    if (auto [re, f, l] = ctre::match<"([A-Za-z]+), ([A-Za-z]+)">(sv); re) {
        return name{f, l};
    } else {
        return std::nullopt;
    }
}
```

**But this must generate
horrible assembly!**

Not really :)

Optimized assembly is just **68 LoC** long!

```
A Save/Load + Add new... C++ x86-64 gcc 7.3 -std=c++17 -O3
1 #include <ctre.hpp>
2 #include <string_view>
3 #include <optional>
4
5 struct name { std::string_view first, family; };
6
7 static constexpr ctre::basic_fixed_string pattern = "[A-Za-z][A-Za-z]";
8
9 std::optional<name> extract(std::string_view sv) noexcept {
10     if (auto [re,f,l] = ctre::re<pattern>().match(sv); re) {
11         return name{f,l};
12     } else {
13         return std::nullopt;
14     }
15 }
16
```

```
A 11010 .LX0: .text // \s+ Intel Demangle Libraries + Add new... Add tool...
1 extract(std::basic_string_view<char, std::char_traits<char> >):
2     add    rsi, rdx
3     mov    rax, rdi
4     cmp    rdx, rsi
5     je     .L2
6     movzx ecx, BYTE PTR [rdx]
7     and    ecx, -33
8     sub    ecx, 65
9     cmp    cl, 25
10    jbe    .L15
11 .L2:
12    xor    edx, edx
13    mov    BYTE PTR [rax+32], dl
14    ret
15 .L15:
16    lea    rcx, [rdx+1]
17    cmp    rsi, rcx
18    je     .L2
19 .L5:
20    movzx r9d, BYTE PTR [rcx]
21    mov    r8d, r9d
22    and    r8d, -33
23    sub    r8d, 65
24    cmp    r8b, 25
25    jbe    .L3
26    cmp    rsi, rcx
27    je     .L2
28    cmp    r9b, 44
29    jne    .L2
30    lea    r9, [rcx+1]
31    cmp    rsi, r9
32    je     .L2
33    movzx edi, BYTE PTR [rcx+1]
34    and    edi, -33
35    sub    edi, 65
36    cmp    dil, 25
37    ja     .L2
38    lea    r8, [rcx+2]
39    cmp    rsi, r8
40    je     .L8
41 .L9:
42    movzx edi, BYTE PTR [r8]
43    and    edi, -33
44    sub    edi, 65
45    cmp    dil, 25
```


**(In theory) The optimizer should be
able to see the same intent
with `std::regex` too...**

Unfortunately it doesn't.

The equivalent code is ~19.7 kLoC.

```
A Save/Load + Add new... C++ x86-64 gcc 7.3 -std=c++17 -O3
1 #include <regex>
2 #include <string_view>
3 #include <optional>
4
5 struct name { std::string_view first, family; };
6
7 std::optional<name> extract(std::string_view sv) noexcept {
8     static std::regex re{"^([a-z]++),([a-z]++)$"};
9
10    std::match_results<std::string_view::iterator> results;
11
12    if (std::regex_match(sv.begin(), sv.end(), results, re)) {
13        //[[assert: results.size() == 2]];
14        return name{results[0].str(), results[1].str()};
15    } else {
16        return std::nullopt;
17    }
18 }
19
A 11010 .LX0: .text // \s+ Intel Demangle Libraries + Add new... Add tool...
1 std::ctype<char>::do_widen(char) const:
2     mov     eax, esi
3     ret
4 std::ctype<char>::do_narrow(char, char) const:
5     mov     eax, esi
6     ret
7 std::_Function_base::_Base_manager<std::__detail::_AnyMatcher<std::__cxx11::regex_traits<char>, fa
8     test    edx, edx
9     je     .L6
10    cmp     edx, 1
11    jne    .L5
12    mov     QWORD PTR [rdi], rsi
13 .L5:
14    xor     eax, eax
15    ret
16 .L6:
17    mov     QWORD PTR [rdi], OFFSET FLAT:typeinfo for std::__detail::_AnyMatcher<std::__cxx11:
18    xor     eax, eax
19    ret
20 std::_Function_base::_Base_manager<std::__detail::_AnyMatcher<std::__cxx11::regex_traits<char>, fa
21    cmp     edx, 1
22    je     .L10
23    jb     .L11
24    cmp     edx, 2
25    jne    .L9
26    mov     rax, QWORD PTR [rsi]
27    mov     QWORD PTR [rdi], rax
28 .L9:
29    xor     eax, eax
30    ret
31 .L11:
32    mov     QWORD PTR [rdi], OFFSET FLAT:typeinfo for std::__detail::_AnyMatcher<std::__cxx11:
33    xor     eax, eax
34    ret
35 .L10:
36    mov     QWORD PTR [rdi], rsi
37    xor     eax, eax
38    ret
39 std::_Function_base::_Base_manager<std::__detail::_AnyMatcher<std::__cxx11::regex_traits<char>, fa
40    cmp     edx, 1
41    je     .L15
42    jb     .L16
43    cmp     edx, 2
44    jne    .L14
45    mov     rax, QWORD PTR [rsi]
```

What about errors?

```
static inline constexpr ctre::fixed_string pattern = ")hello";

bool fnc(std::string_view view) {
    return ctre::match<pattern>(view);
}
```

In file included from test.cpp:1:

In file included from include/ctre.hpp:5:

```
include/ctre/functions.hpp:48:2: error: static_assert failed due to requirement 'correct'
    static_assert(correct, "Regular Expression contains syntax error.");
    ^              ~~~~~
```

```
test.cpp:21:15: note: in instantiation of function template specialization 'ctre::match'
    return ctre::match<pattern>(view);
           ^
```

Compile Time Regular Expression library

- is easy to use,
- uses existing and well known RE syntax,
- is perfect for checking inputs of your program,
- emits nice assembly,
- works on Clang and GCC and **MSVC 15.8.8+**
- and we will talk later about performance...

The Parser

Disclaimer: all examples were simplified.

How does it work?

- It uses a generic **LL(1)** parser for converting a pattern into a type.
- The parser uses a provided PCRE compatible grammar.
- Output type encodes the pattern as an **expression template**.
- Evaluation of the expression template matches the subject.

How does the LL(1) parser work?

- Starts with a **start symbol** on the stack.
- On every **step** it **pops** one symbol from the stack and **checks the current character** at the input.
- Based on the pair of symbol and character it decides to:
 - **push** a string of symbols to the stack,
 - **pop** a character from the input,
 - or **reject**.
- Repeat until the stack and input are empty then **accept**.

Let's create RE grammar...

What does the grammar look like?

$f(\text{symbol}, \text{char}) \rightarrow$

	()	*	+	?	\	a	d		other	ϵ
$\rightarrow S$	(alt0) mod seq alt					\ esc mod seq alt	a mod seq alt	d mod seq alt		other mod seq alt	ϵ
alt0	(alt0) mod seq alt					\ esc mod seq alt	a mod seq alt	d mod seq alt		other mod seq alt	
alt		ϵ							seq0 alt		ϵ
esc							a	d			
mod	ϵ	ϵ	*	+	?	ϵ	ϵ	ϵ	ϵ	ϵ	ϵ
seq0	(alt0) mod seq					\ esc mod seq	a mod seq	d mod seq		other mod seq	
seq	(alt0) mod seq	ϵ				\ esc mod seq	a mod seq	d mod seq	ϵ	other mod seq	ϵ
(<u>pop</u>										
)		<u>pop</u>									
*			<u>pop</u>								
+				<u>pop</u>							
?					<u>pop</u>						
\						<u>pop</u>					
a							<u>pop</u>				
d								<u>pop</u>			
									<u>pop</u>		
other										<u>pop</u>	
Z ₀											<u>accept</u>

How can I write it in C++?

How can I represent symbols in C++?

```
struct S {};  
struct alt0 {};  
struct alt {};  
struct esc {};  
struct mod {};  
struct seq0 {};  
struct seq {};  
  
using start_symbol = S;
```

How can I represent an LL(1) table in C++?

```
//  $f(\text{symbol}, \text{char}) \rightarrow (\dots)$   
auto f(symbol, term<'c'>) -> list{...};  
  
//  $f(\text{symbol}, \text{symbol}) \rightarrow \text{pop input}$   
template <auto S> auto f(term<S>, term<S>) -> pop_input;  
  
//  $f(\text{symbol}, \text{char}) \rightarrow \text{reject}$   
auto f(...) -> reject;  
  
//  $f(Z_0, \epsilon) \rightarrow \text{accept}$   
auto f(empty_stack, epsilon) -> accept;
```

How can I pass the grammar?

```
struct my_grammar {  
    struct S {};  
    struct alt0 {};  
    struct alt {};  
    // ...  
  
    using start_symbol = S;  
  
    auto f(...) -> reject;  
    // ...  
}
```

Let's use it...

```
constexpr bool ok = parser<my_grammar, "^hello">::correct;
```

What about parser?

```
template <typename Grammar, ...> struct parser {  
    //...  
    auto next_move = Grammar::f(top_of_stack, current_char);  
    //...  
}
```

How do I get current character?

```
template <..., fixed_string Str> struct parser {  
    // ...  
    template <size_t Pos> constexpr auto get_character() const {  
        if constexpr (Pos < Str.size()) return term<Str[Pos]>{};  
        else return epsilon{};  
    }  
    // ...  
};
```


How do I implement the input string?

```
template <typename CharT, size_t N> struct fixed_string {
    std::array<CharT, N> data;
    // constexpr constructor from const char[N]

    constexpr auto operator[](size_t i) const noexcept { return data[i]; }
    constexpr size_t size() const noexcept { return N; }
    constexpr auto operator<=>(const fixed_string &) = default;
};

template <typename CharT, size_t N>
    fixed_string(const CharT[N]) -> fixed_string<CharT, N>;

// more info about class NTPP in p0732 by Jeff Snyder and Louis Dionne
```

How do I implement the stack?

```
template <typename... Ts> struct list { };

template <typename... Ts, typename... As>
constexpr auto push(list<Ts...>, As...) -> list<As..., Ts...>;

template <typename T, typename... As>
constexpr auto pop(list<T, Ts...>) -> list<Ts...>;

template <typename T, typename... Ts>
constexpr auto top(list<T, Ts...>) -> T;

struct empty { };
constexpr auto top(list<>) -> empty;
```

How does it fit together?

```
constexpr bool ok = parser<my_grammar, "there$">::correct;
```

```

template <typename Grammar, fixed_string Str> struct parser {
    static constexpr bool correct = parse(list<Grammar::start_symbol>());

    // prepare each step and move to next
    template <size_t Pos = 0, typename S>
        static constexpr bool parse(S stack) {

            using next_stack = decltype( pop(stack) );
            using symbol     = decltype( top(stack) );
            using current    = decltype( get_character<Pos>() );
            using next_move  = decltype( Grammar::f(symbol{}, current{}) );
            return move<Pos>(next_move{}, next_stack{});
        }
    //...

```

```

// pushing something to stack (epsilon case included)
template <size_t Pos, typename... Push, typename Stack>
    static constexpr bool move(list<Push...>, Stack stack) {
        using next_stack = decltype( push(stack, Push{}...) )
        return parse<Pos>(next_stack{});
    }

// move to next character
template <size_t Pos, typename Stack>
    static constexpr bool move(pop_input, Stack stack) {
        return parse<Pos+1>(stack);
    }

```

```
template <size_t Pos, typename Stack>
    static constexpr bool move(reject, Stack) {
        return false;
    }

template <size_t Pos, typename Stack>
    static constexpr bool move(accept, Stack) {
        return true;
    }

}; // end of Parser struct
```

**But this is returning a boolean value!
Is there something missing?**

Building The Expression Template

How can I **build a type** from a string?

Where are the semantic actions placed?

	()	*	+	?	\	a	d		other	ϵ
$\rightarrow S$	(<i>alt0</i>) <i>mod seq alt</i>					\ <i>esc mod seq alt</i>	a char <i>mod seq alt</i>	d char <i>mod seq alt</i>		other char <i>mod seq alt</i>	ϵ
alt0	(<i>alt0</i>) <i>mod seq alt</i>					\ <i>esc mod seq alt</i>	a char <i>mod seq alt</i>	d char <i>mod seq alt</i>		other char <i>mod seq alt</i>	
alt		ϵ							<i>seq0 alt alt</i>		ϵ
esc							a alpha	d digit			
mod	ϵ	ϵ	* star	+ plus	? opt	ϵ	ϵ	ϵ	ϵ	ϵ	ϵ
seq0	(<i>alt0</i>) <i>mod seq</i>					\ <i>esc mod seq</i>	a char <i>mod seq</i>	d char <i>mod seq</i>		other char <i>mod seq</i>	
seq	(<i>alt0</i>) <i>mod seq seq</i>	ϵ				\ <i>esc mod seq seq</i>	a char <i>mod seq seq</i>	d char <i>mod seq seq</i>	ϵ	other char <i>mod seq seq</i>	ϵ
(<u>pop</u>										
)		<u>pop</u>									
*			<u>pop</u>								
+				<u>pop</u>							
?					<u>pop</u>						
\						<u>pop</u>					
a							<u>pop</u>				
d								<u>pop</u>			
									<u>pop</u>		
other										<u>pop</u>	
Z_0											<u>accept</u>

What do the **SemanticAction** symbols look like?

```
struct my_grammar {  
    struct _char: action {};  
    struct alpha: action {};  
    struct digit: action {};  
    struct seq: action {};  
    struct star: action {};  
    struct plus: action {};  
    struct opt: action {};  
    // ...  
}
```

What are the changes in the parser?

```
template <typename Grammar, fixed_string Str> struct parser {  
  
    template <size_t Pos = 0, typename S, typename T = list<>>  
        static constexpr auto parse(S stack, T subject = {}) {  
            using next_stack = decltype( pop(stack) );  
            using symbol      = decltype( top(stack) );  
            if constexpr (SemanticAction<symbol>) {  
                using previous = decltype( prev_character<Pos>() );  
                using next_subject = decltype( modify(symbol{}, previous{}, subject) );  
                return parse<Pos>(next_stack(), next_subject{});  
            } else {  
                using current      = decltype( get_character<Pos>() );  
                using next_move    = decltype( Grammar::f(symbol{}, current{}) );  
                return move<Pos>(next_move{}, next_stack{}, subject);  
            }  
        }  
};
```

```

// pushing something to stack (epsilon case included)
template <size_t Pos, typename... Push, typename Stack, typename T>
static constexpr auto move(list<Push...>, Stack stack, T subject) {
    using next_stack = decltype( push(stack, Push{}...) )
    return parse<Pos>(next_stack{}, subject);
}

// move to next character
template <size_t Pos, typename Stack, typename T>
static constexpr auto move(pop_input, Stack stack, T subject) {
    return parse<Pos+1>(stack, subject);
}

```

```
template <size_t Pos, typename Stack, typename T>
    static constexpr auto move(reject, Stack, T subject) {
        return pair{false, subject};
    }

template <size_t Pos, typename Stack, typename T>
    static constexpr auto move(accept, Stack, T subject) {
        return pair{true, subject};
    }

}; // end of Parser struct
```

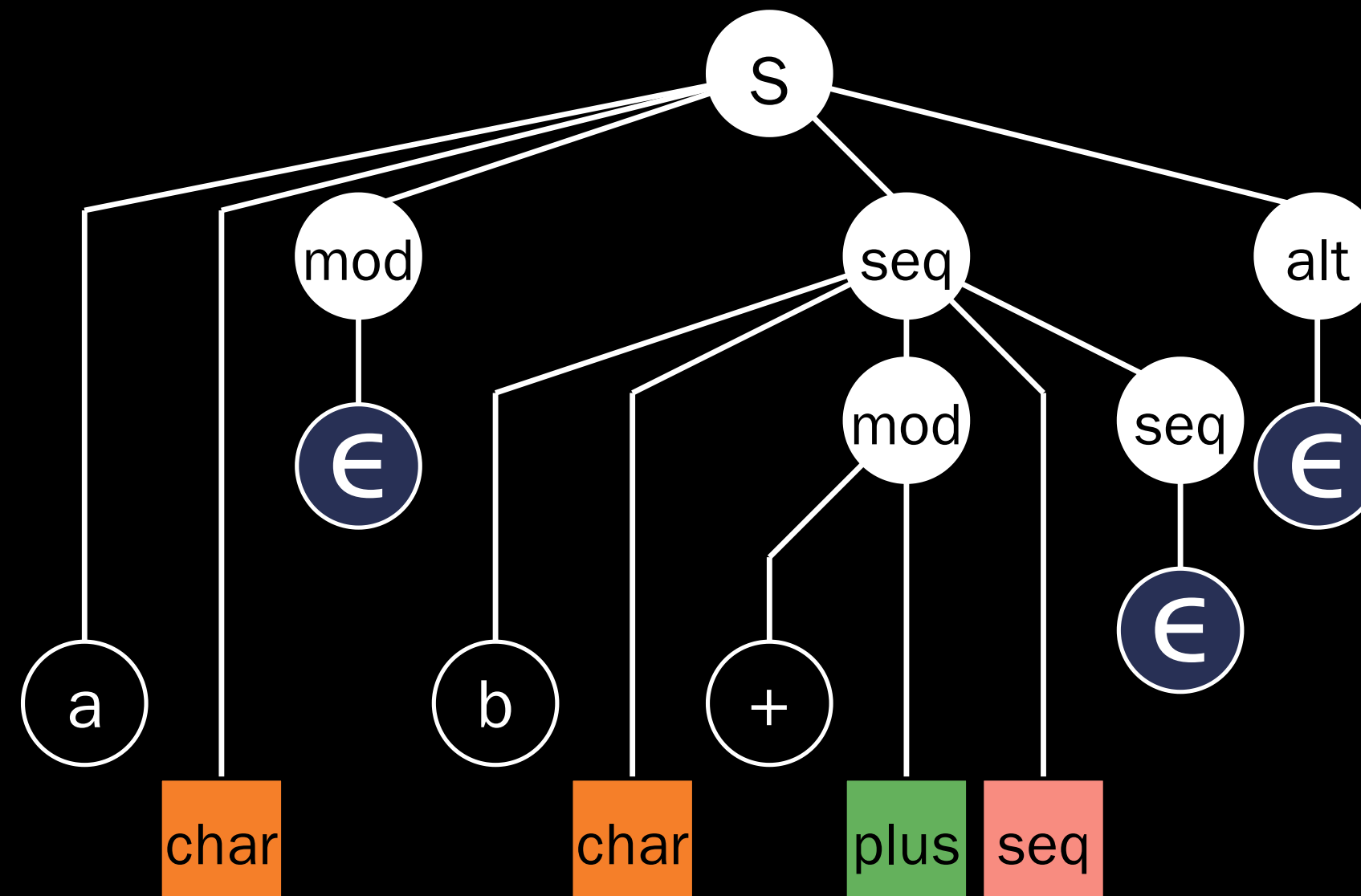
Where is my subject?

```
template <typename Grammar, fixed_string Str> struct parser {
    static constexpr pair result
        = parse(list<Grammar::start_symbol>(), list<>());
    static constexpr auto correct = result.first;
    static_assert(correct, "RE syntax error");
    using output_type = decltype(front(result.second));
    // ...
};

using T = parser<my_grammar, "^wow$" >::output_type;
```

What does building from a string look like?

ab+



seq<ch<'a'>,plus<ch<'b'>>>

What about the modify function?

Pushing character

```
template <auto C, typename... Ts>  
  auto modify(my_grammar::_char, term<C>, list<Ts...>)  
    -> list<ch<C>, Ts...>;
```

(Notice the term and ch types.)

Making something optional

```
template <auto C, typename Opt, typename... Ts>  
  auto modify(my_grammar::opt, term<C>, list<Opt, Ts...>)  
    -> list<opt<Opt>, Ts...>;
```

Making cycles

```
template <auto C, typename Plus, typename... Ts>  
  auto modify(my_grammar::plus, term<C>, list<Plus, Ts...>)  
  -> list<plus<Plus>, Ts...>;
```

```
template <auto C, typename Star, typename... Ts>  
  auto modify(my_grammar::star, term<C>, list<Star, Ts...>)  
  -> list<star<Star>, Ts...>;
```

Creating a sequence

```
template <auto C, typename A, typename B, typename... Ts>  
  auto modify(my_grammar::seq, term<C>, list<B, A, Ts...>)  
  -> list<seq<A, B>, Ts...>;
```

(Notice the switched order of A & B on the stack.)

Adding to a sequence

```
template <auto C, typename... As, typename B, typename... Ts>  
  auto modify(my_grammar::seq, term<C>, list<B, seq<As...>, Ts...>)  
    -> list<seq<As..., B>, Ts...>;
```

Making an alternate

```
template <auto C, typename A, typename B, typename... Ts>  
  auto modify(my_grammar::alt, term<C>, list<B, A, Ts...>)  
  -> list<alt<A, B>, Ts...>;
```

(Notice the switched order of A & B on the stack.)

Adding to an alternate

```
template <auto C, typename... As, typename B, typename... Ts>  
  auto modify(my_grammar::alt, term<C>, list<B, alt<As...>, Ts...>)  
    -> list<alt<As..., B>, Ts...>;
```


Adding a digit or alpha class

```
template <auto C, typename... Ts>  
  auto modify(my_grammar::alpha, term<C>, list<Ts...>)  
  -> list<alpha, Ts...>;
```

```
template <auto C, typename... Ts>  
  auto modify(my_grammar::digit, term<C>, list<Ts...>)  
  -> list<digit, Ts...>;
```

We have the expression template.

```
static_assert (std::is_same_v<  
    seq<ch<'a'>, plus<ch<'b'>>>,  
    parser<my_grammar, "ab+">::output_type  
>);
```

Evaluating The Expression Template

Let's make this work...

```
match<" \\a \\d+"> ("a42"sv) ;
```

How are REs matched against a string?

```
template <fixed_string re> bool match(std::string_view sv) {  
    static_assert(parser<my_grammar, re>::correct);  
    using RE = parser<my_grammar, re>::output_type;  
  
    return match(sv.begin(), sv.begin(), sv.end(), list<RE>{}).success;  
}
```

List? Again?!

Yes!

But first what's the output type of match(...)?

```
template <ForwardIterator It> struct result {  
    ForwardIterator it;  
    bool success;  
    // constexpr constructor etc...  
};
```

Matching a character

```
template <ForwardIterator It, auto C, typename... Ts>
result<It> match(It begin, It it, It end, list<ch<C>, Ts...>) {

    if (it == end) return {it, false};
    if (*it != C) return {it, false};

    return match(begin, std::next(it), end, list<Ts...>{});
}
```


Matching an alpha or a digit class

```
template <ForwardIterator It, typename... Ts>
result<It> match(It begin, It it, It end, list<digit, Ts...>) {
    if (it == end) return {it, false};
    if (!(*it >= '0' && *it <= '9')) return {it, false};
    return match(begin, std::next(it), end, list<Ts...>{});
}
```

```
template <ForwardIterator It, typename... Ts>
result<It> match(It begin, It it, It end, list<alpha, Ts...>) {
    if (it == end) return {it, false};
    if (!(*it >= 'a' && *it <= 'z')) return {it, false};
    return match(begin, std::next(it), end, list<Ts...>{});
}
```

Matching a sequence

```
template <ForwardIterator It, typename... Seq, typename... Ts>
    result<It> match(It begin, It it, It end, list<seq<Seq...>, Ts...>)

    return match(begin, it, end, list<Seq..., Ts...>{});
}
```

Matching an optional

```
template <ForwardIterator It, typename... Opt, typename... Ts>
    result<It> match(It begin, It it, It end, list<opt<Opt...>, Ts...>)

    if (auto out = match(begin, it, end, list<Opt..., Ts...>{})) {
        return out;
    } else {
        // try it without the content of opt<...>
        return match(begin, it, end, list<Ts...>{});
    }
}
```

Matching an alternate

```
template <ForwardIterator It, typename Head, typename... Tail, typename... Ts>
result<It> match(It begin, It it, It end, list<alt<Head, Tail...>, Ts...>) {
    if (auto out = match(begin, it, end, list<Head, Ts...>{})) {
        return out;
    } else {
        // try the next one
        return match(begin, it, end, list<alt<Tail...>, Ts...>{}));
    }
}
```

```
template <ForwardIterator It, typename... Ts>
result<It> match(It begin, It it, It end, list<alt<>, Ts...>) {
    // no option from alternation was successful
    return {it, false};
}
```

Matching a plus cycle

```
template <ForwardIterator It, typename First, typename... Alt, typename... Ts>
result<It> match(It begin, It it, It end, list<plus<Plus...>, Ts...>) {
    for (;;) {
        if (auto inner = match(begin, it, end, list<Plus..., end_of_cycle>{})) {
            if (auto out = match(begin, it, end, list<Ts...>{})) {
                return out;
            } else {
                it = inner.it;
            }
        } else return {it, false};
    }
}
```

```
struct end_of_cycle {};
```

```
template <ForwardIterator It>
result<It> match(It, It it, It, list<end_of_cycle>) {
    return {it, true};
}
```

(The cycle is lazy.)

Matching a star cycle

```
template <ForwardIterator It, typename First, typename... Alt, typename... Ts>
result<It> match(It begin, It it, It end, list<star<Star...>, Ts...>) {
    for (;;) {
        if (auto out = match(begin, it, end, list<Ts...>{})) {
            return out;
        } else if (auto inner = match(begin, it, end, list<Star..., end_of_cycle>{})) {
            if (auto out = match(begin, it, end, list<Ts...>{})) {
                return out;
            } else {
                it = inner.it;
            }
        } else return {it, false};
    }
}
```

(The cycle is lazy.)

How do we know it's the end?

```
// and at the end... we need to check for the end :)
template <ForwardIterator It>
    result<It> match(It begin, It it, It end, list<>{}) {
    // if we are at the end input we should accept
    return {it, it == end};
}
```

(Helper for finishing the regex matching.)

In the examples there are **no captures**
or **any advanced RE functionality.**

Benchmarks

How quick or slow is this thing?

Measured code (grep-like utility)

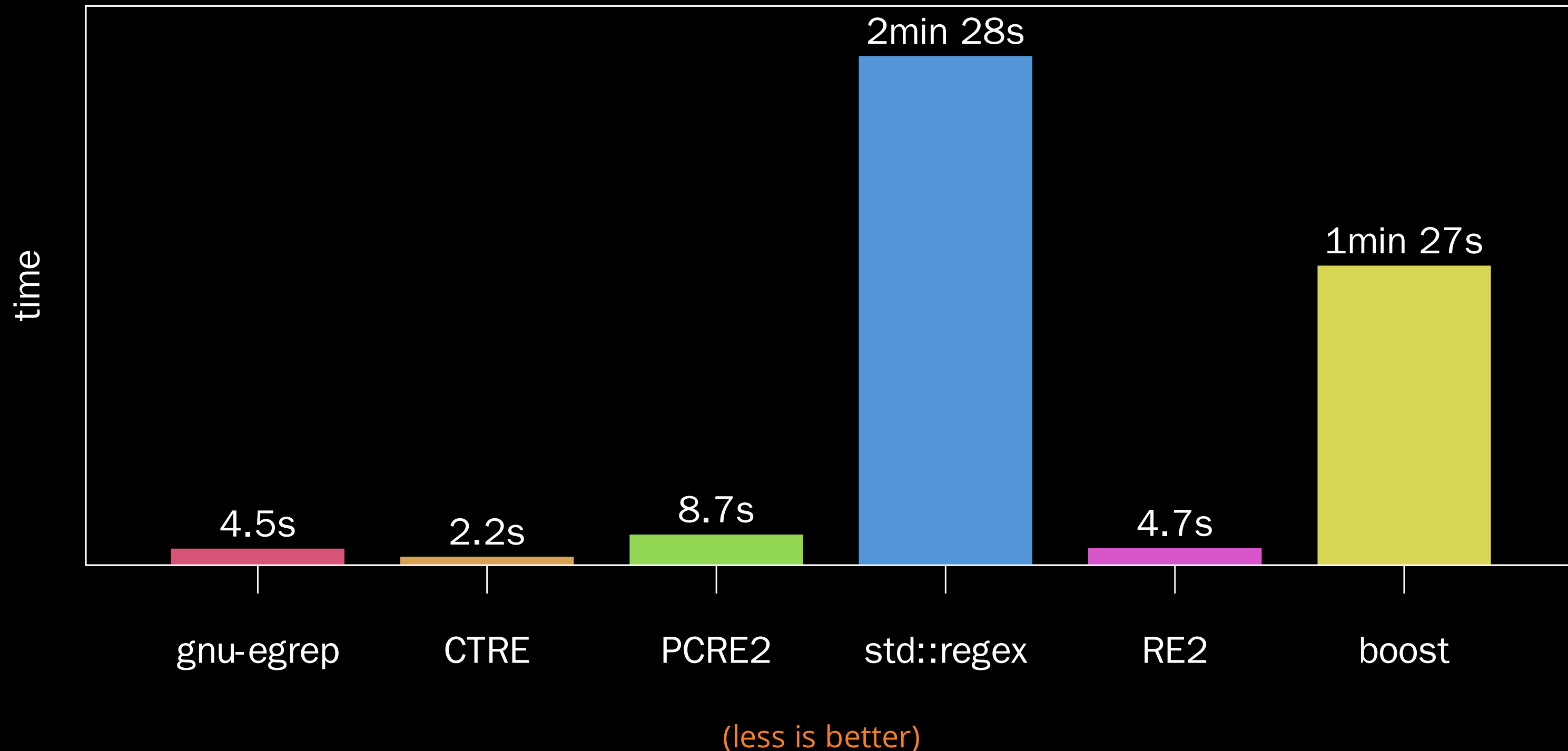
```
int main(int argc, char ** argv) {
    auto re = PREPARE("PATTERN");
    std::ifstream stream{argv[1], std::ifstream::in};
    std::string line;
    while (std::getline(stream, line)) {
        if (re.MATCH(line)) {
            std::cout << line << '\n';
        }
    }
}
```

Measurement methodology

- CSV file (1.3 GiB) with 6.5 MLoC
- Median of 3 measurements
- MacBook Pro 13" 2016 i7 3.3Ghz

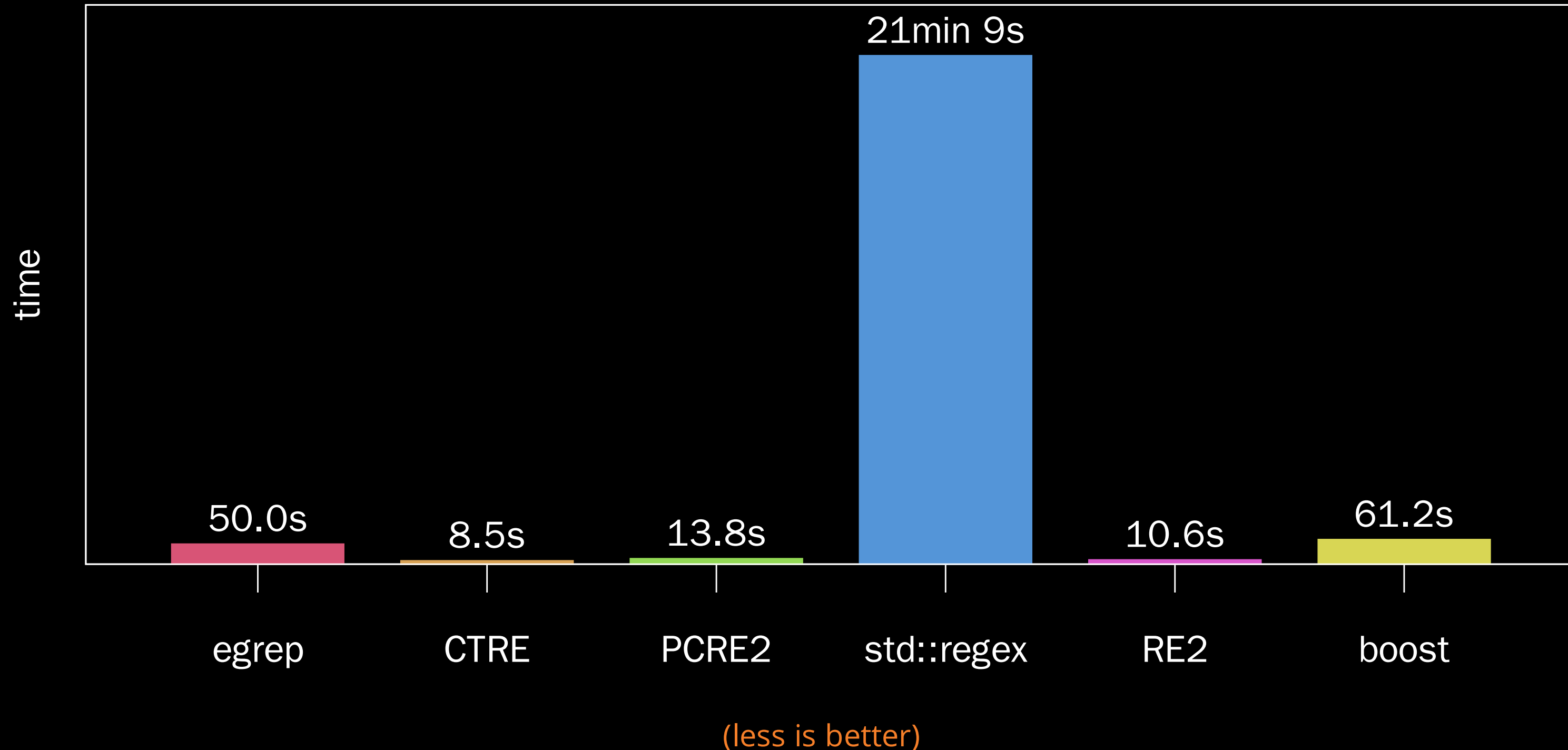
Runtime Matching (GCC, Linux): `ABCD | DEF'GH | EFGHI | A{4, }`

Matching pattern against Big CSV file (1.3 GiB, 6.5 MLoC).



Runtime Matching (Clang, Mac): `ABCD | DEFGH | EFGHI | A{4, }`

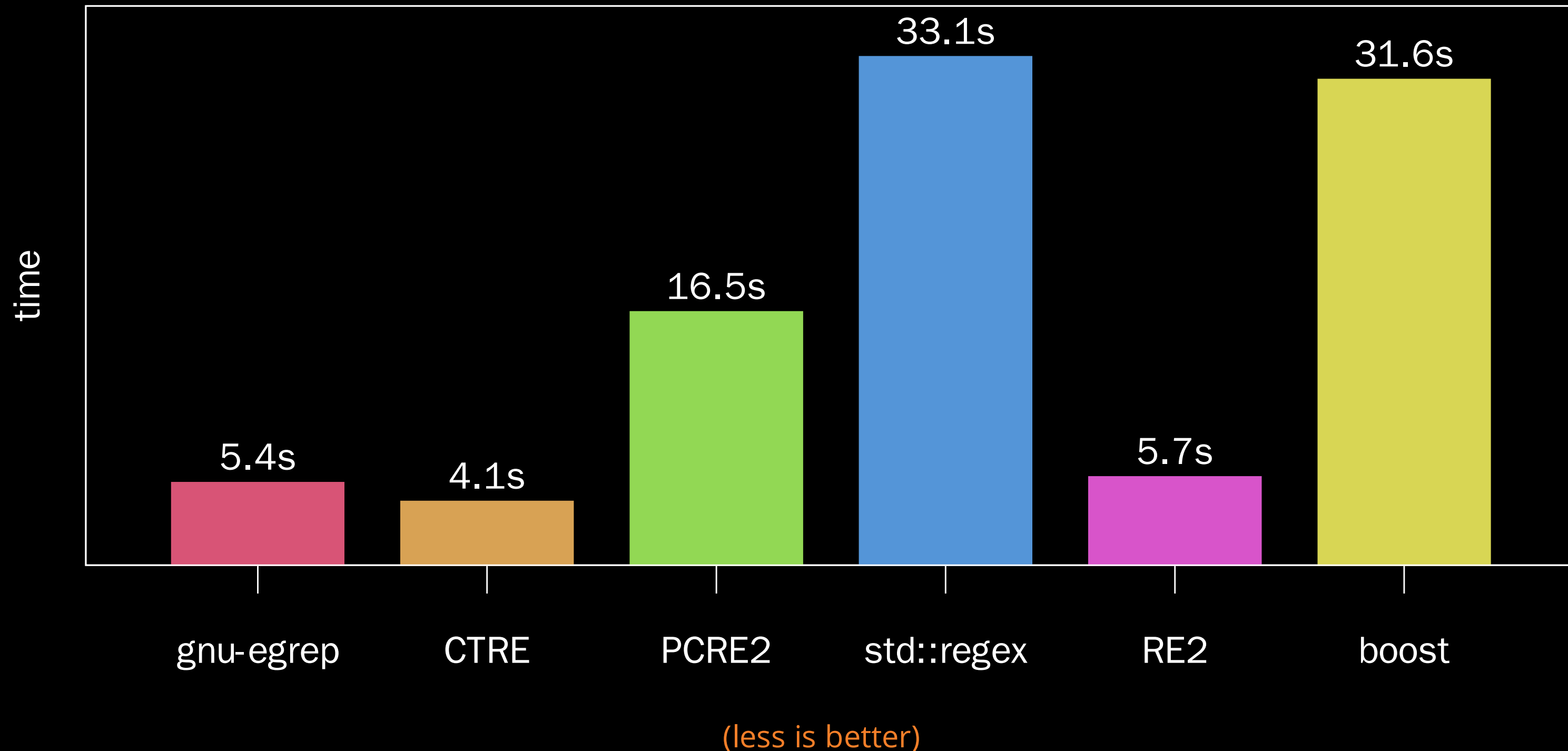
Matching pattern against Big CSV file (1.3 GiB, 6.5 MLoC).



Looks like `std::regex`
in `libc++` is slow...

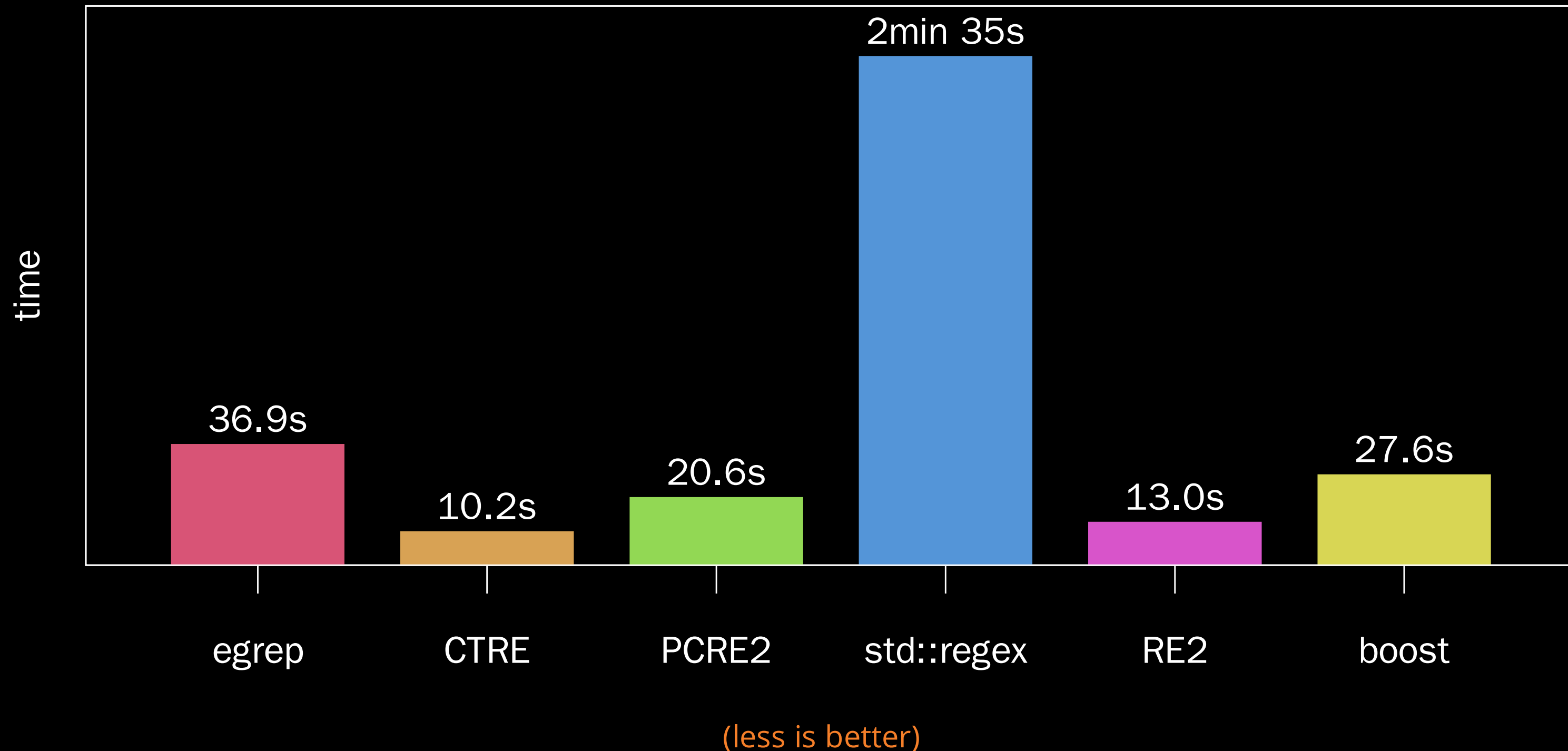
Runtime Matching (GCC, Linux): `[0-9a-fA-F]{8,16}`

Matching pattern against Big CSV file (1.3 GiB, 6.5 MLoC).



Runtime Matching (Clang, Mac): `[0-9a-fA-F]{8,16}`

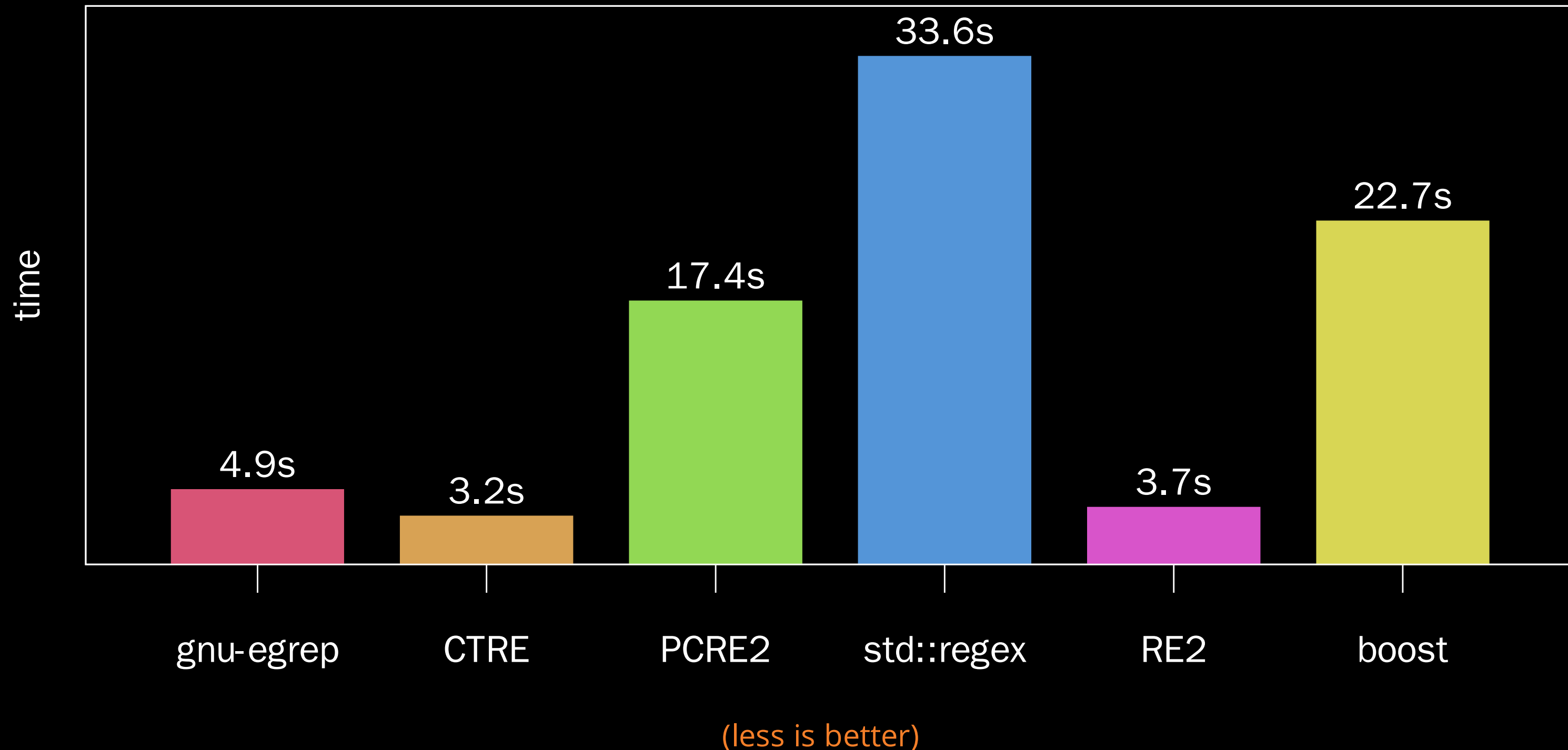
Matching pattern against Big CSV file (1.3 GiB, 6.5 MLoC).



So the `std::regex` in `libc++`
is consistently bad.

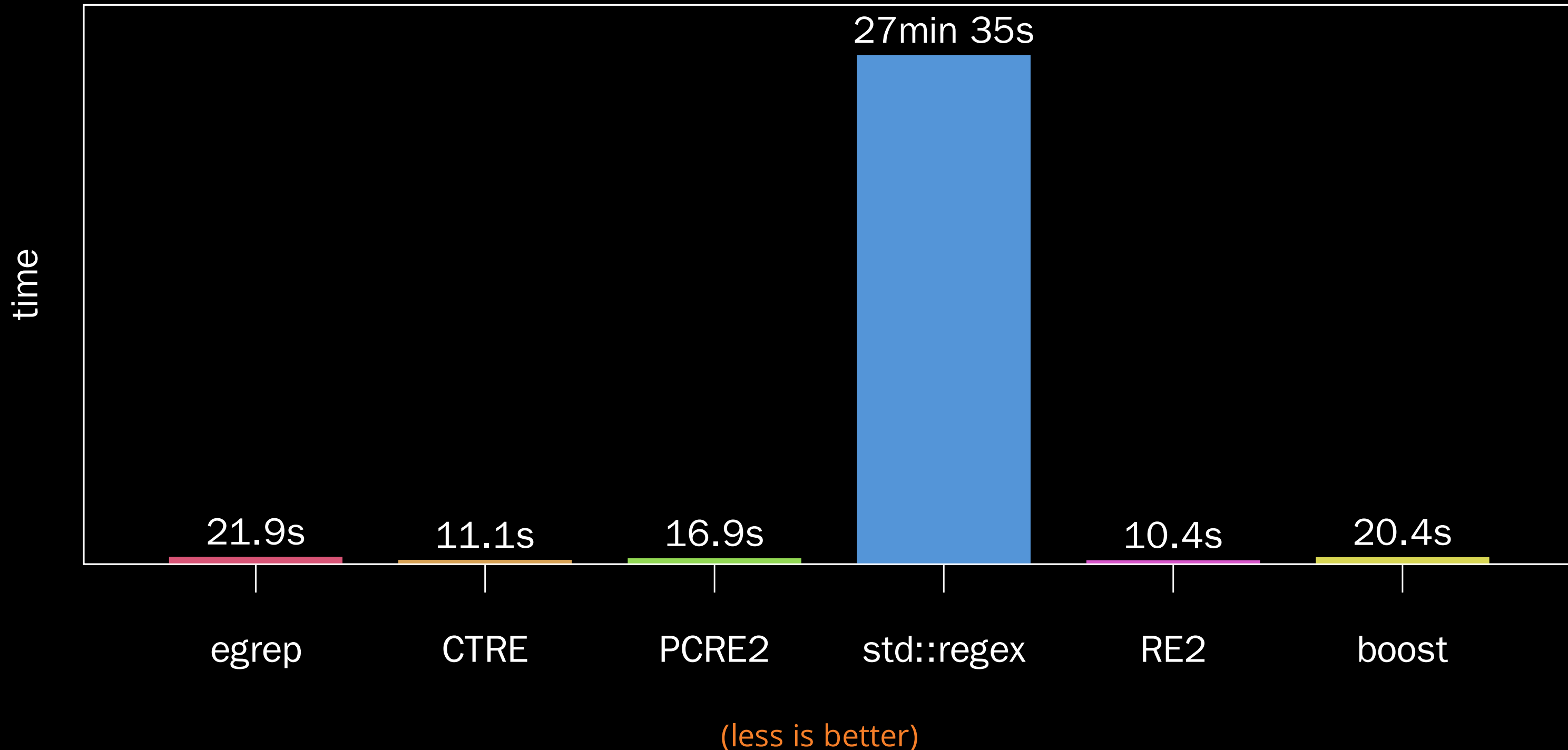
Runtime Matching (GCC, Linux): $([0-9]\{4,16\})?[aA]$

Matching pattern against Big CSV file (1.3 GiB, 6.5 MLoC).



Runtime Matching (Clang, Mac): `([0-9]{4,16})?[aA]`

Matching pattern against Big CSV file (1.3 GiB, 6.5 MLoC).



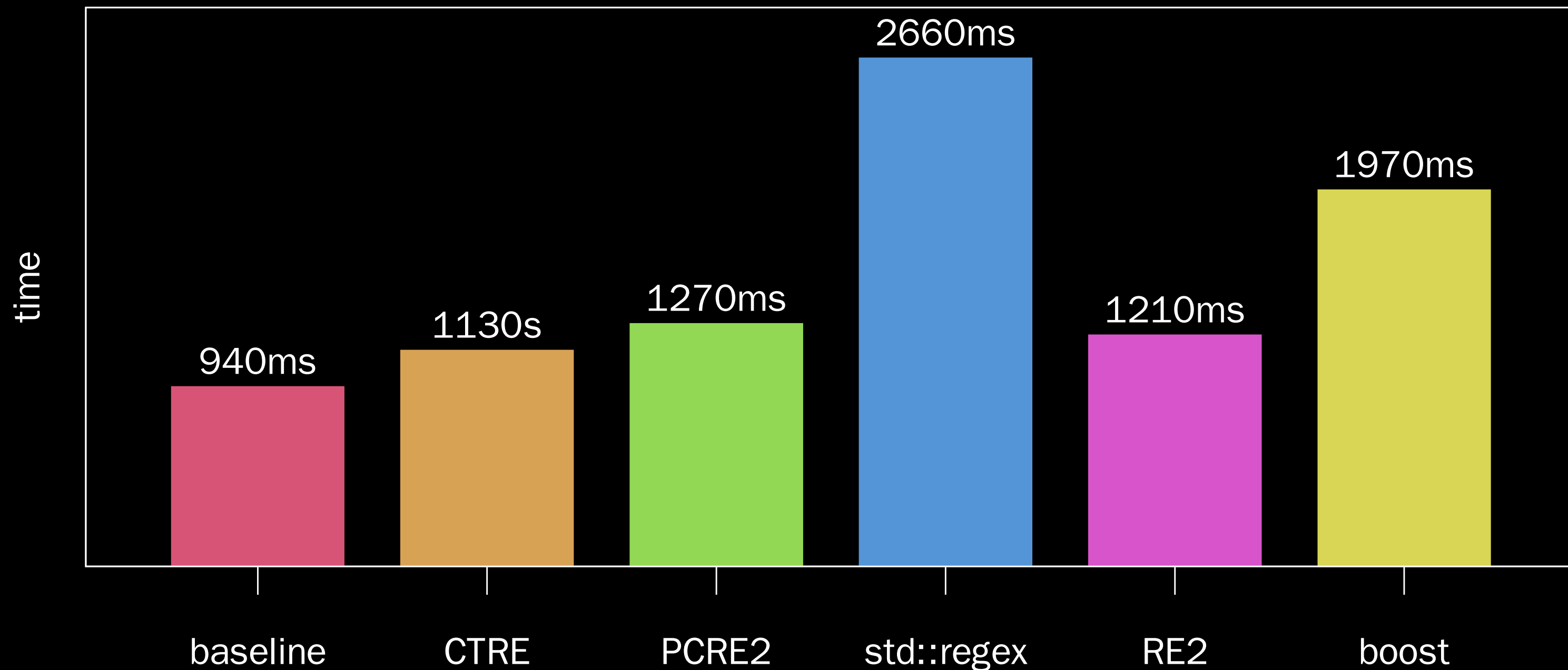
Something **interesting**...

Compile Time Benchmark

- All tests include **the same set of headers**.
- I measured **compilation time only**, without linking.

Optimized **Compile Time** of The Benchmark code.

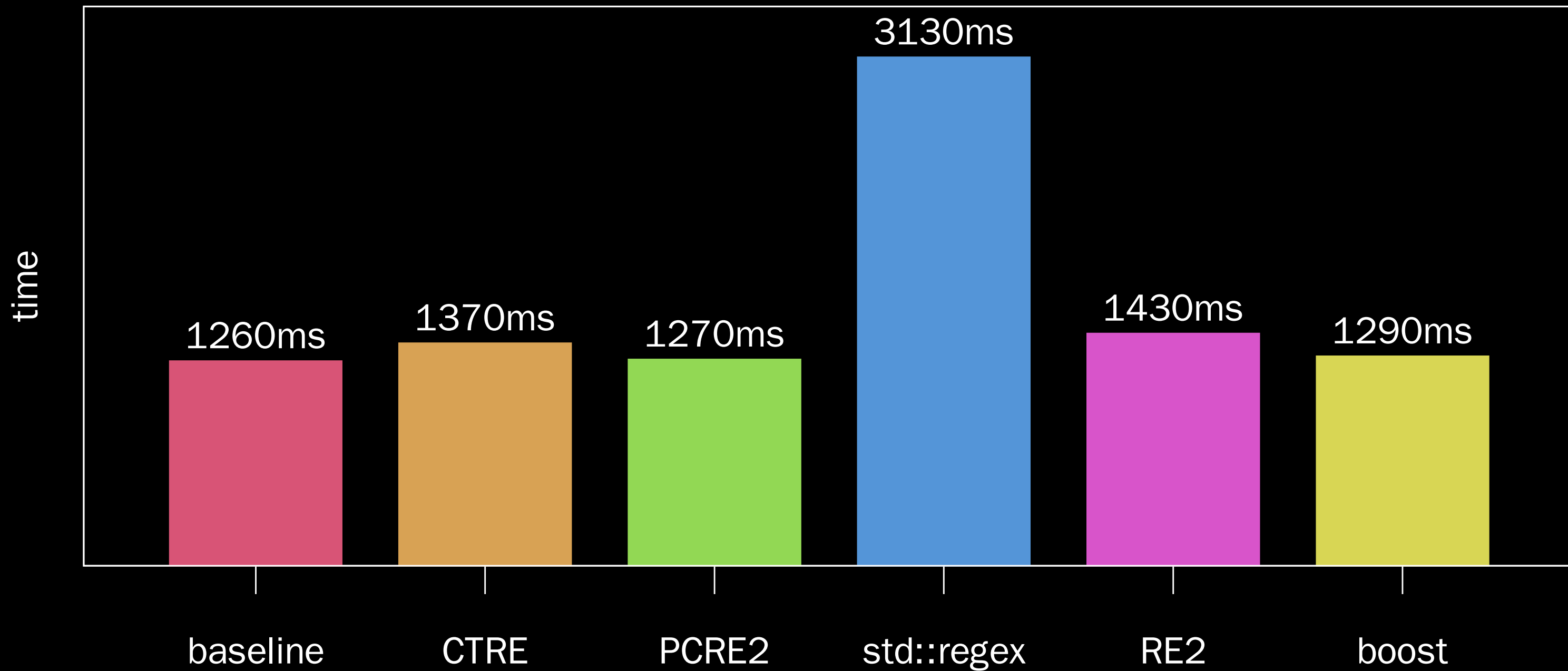
```
clang -c -O3
```



(less is better)

Optimized **Compile Time** of The Benchmark code.

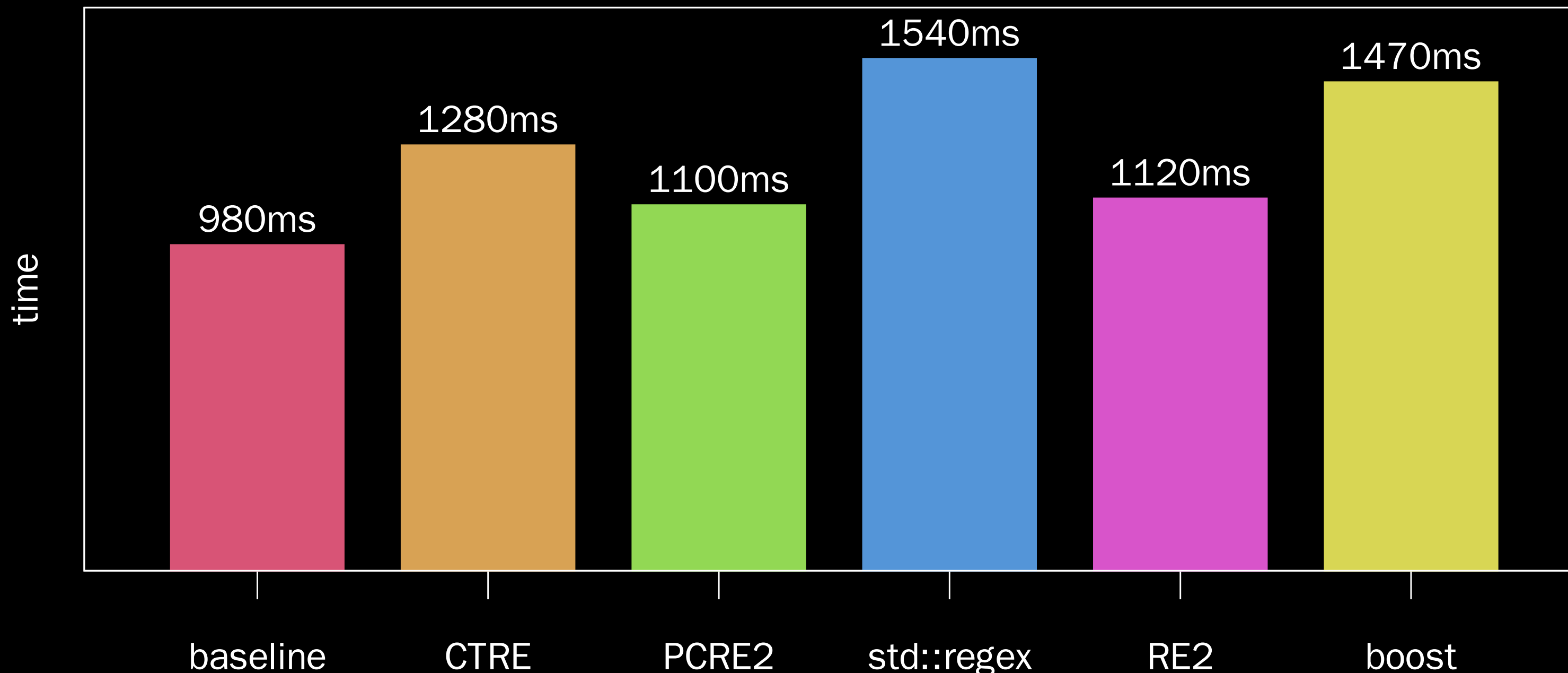
`gcc -c -O3`



(less is better)

Compile Time of The Benchmark code.

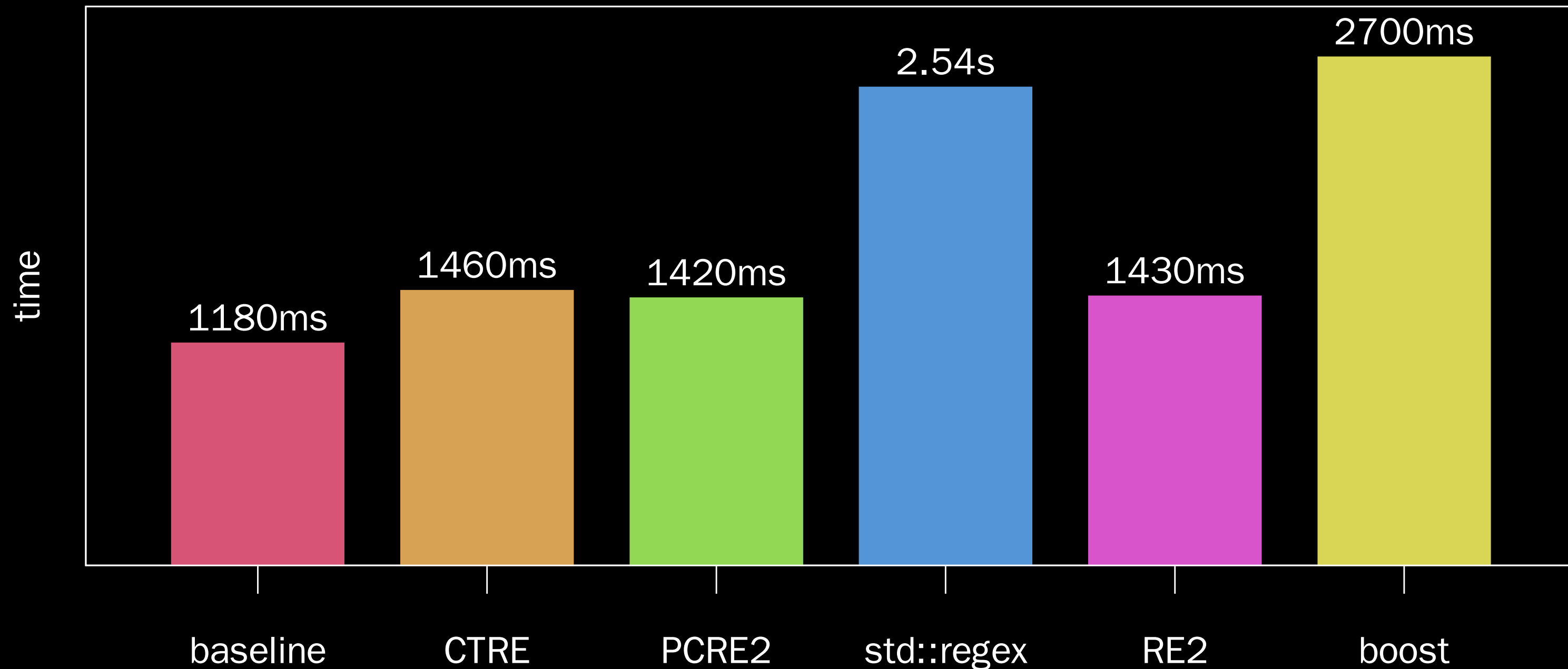
```
clang -c -O0
```



(less is better)

Compile Time of The Benchmark code.

```
gcc -c -O0
```



(less is better)

What's the behaviour of the library?

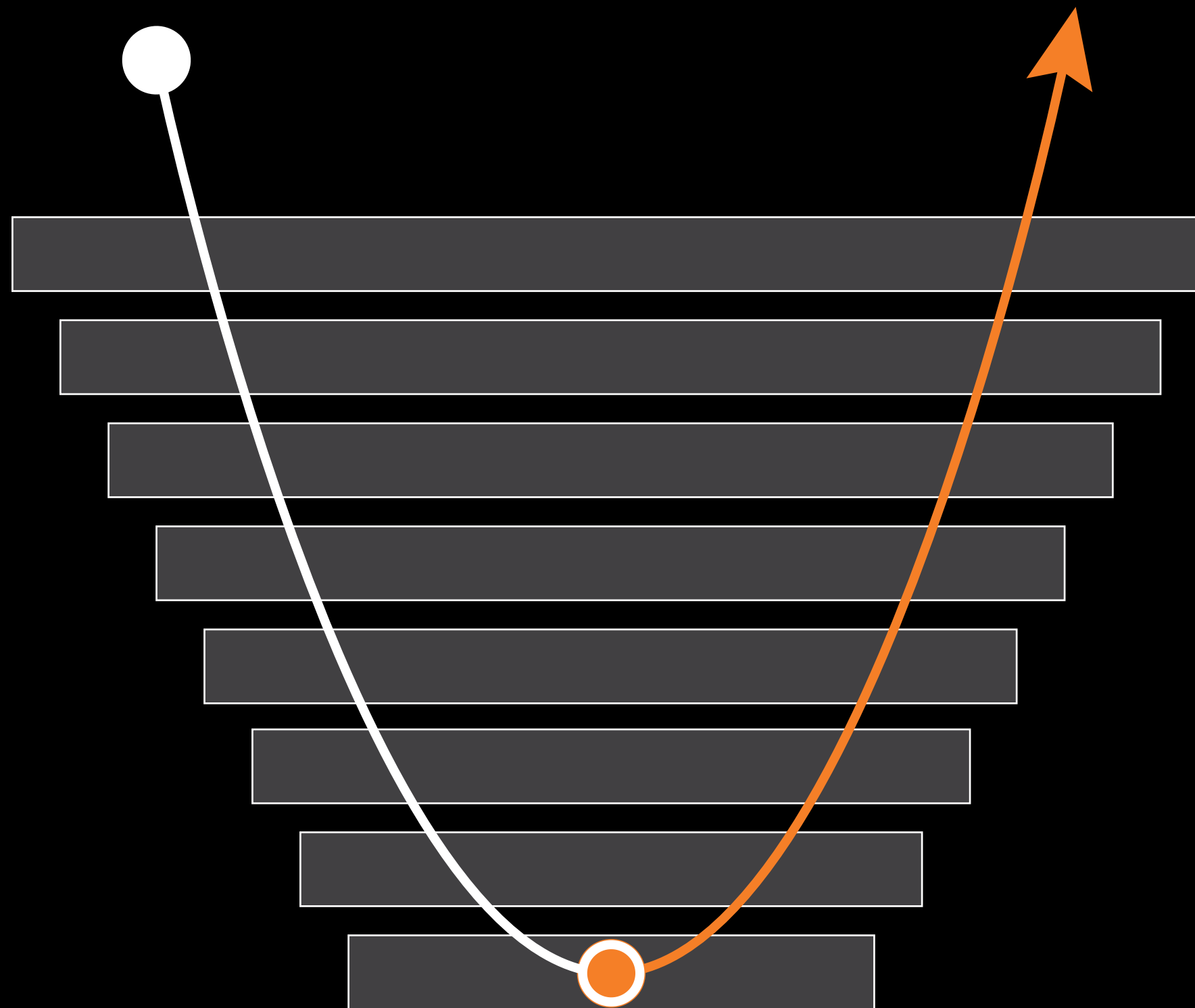
- Multiple different REs cost **~200ms** each.
- Repeating the same REs cost **virtually nothing** (due to caching).
- Cost of one RE parsing is **linear** with its complexity.

What did we learn?

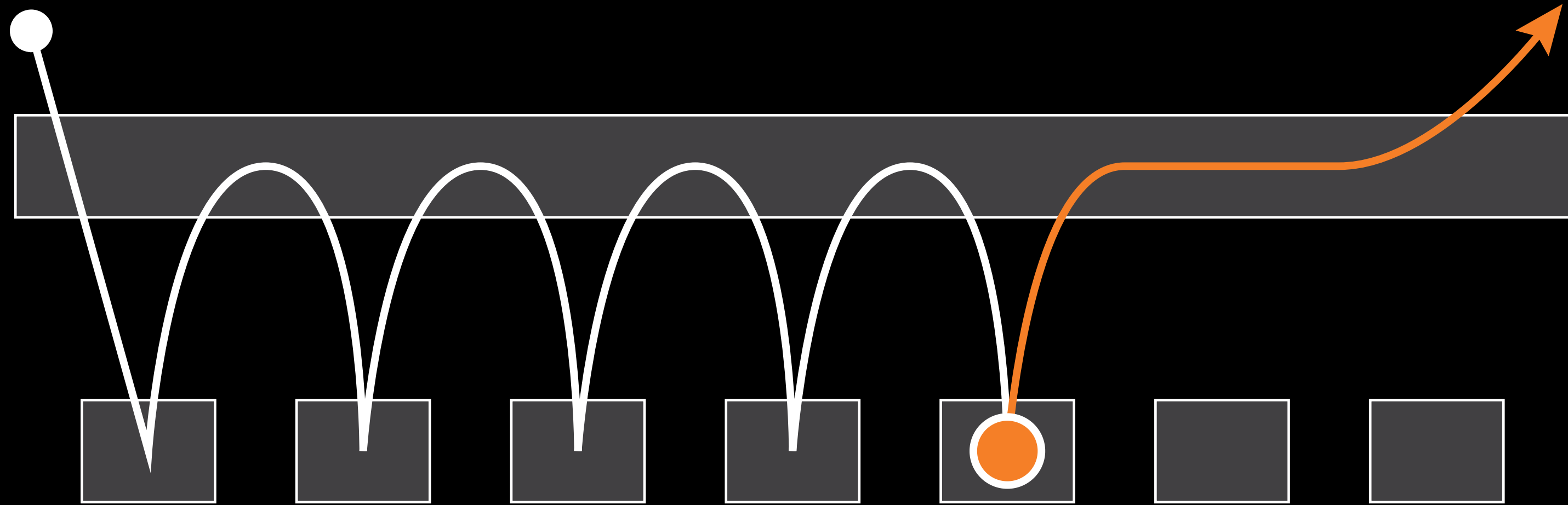
- **Recursion** and **function overloading resolution** are powerful tools.
- Pure functional programming keeps your code **simple**.
- **Trust your compiler** (but check its output).
- **Type-list** is the most **useful** compile-time data structure.

One more thing...

1024



**But there is a no cycle
usable with type-list!**



- divide algorithm into pieces char by char
- store & restore state of parser

Store the parser state

```
enum class decision {
    undecided,
    correct,
    wrong
};

// move to next character
template <size_t Pos, typename Stack, typename T>
static constexpr auto move(pop_input, Stack stack, T subject) {
    // return parse<Pos+1>(stack, subject);
    return state<Pos+1, Stack, T, decision::undecided>();
}
```

Accept and Reject

```
// instead of returning pair with subject and boolean flag

template <size_t Pos, typename Stack, typename T>
    static constexpr auto move(reject, Stack s, T subject) {
        return state<Pos, Stack, T, decision::correct>();
    }

template <size_t Pos, typename Stack, typename T>
    static constexpr auto move(accept, Stack s, T subject) {
        return state<Pos, Stack, T, decision::wrong>();
    }
}
```

Restore the parser state

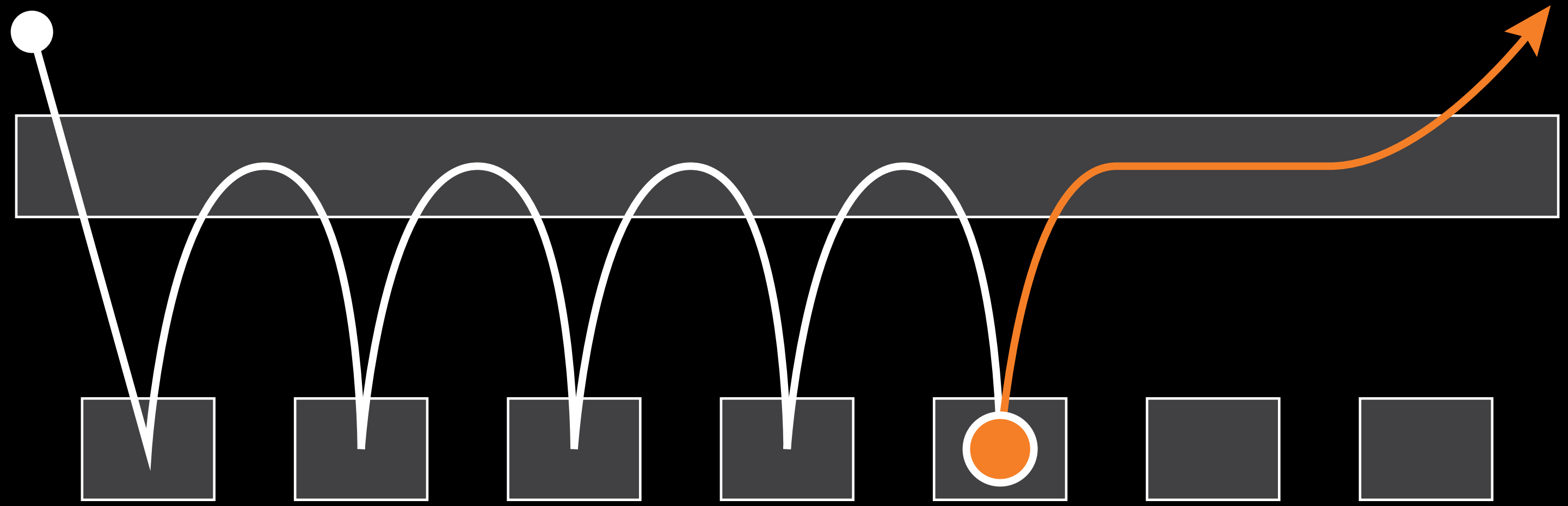
```
struct placeholder {};  
  
template <size_t Pos, typename Stack, typename Subject, decision Decision>  
struct state {  
    auto operator+(placeholder) {  
        if (Decision == decision::undecided) {  
            return parse<Pos>(Stack(), Subject());  
        } else {  
            return *this;  
        }  
    }  
};
```

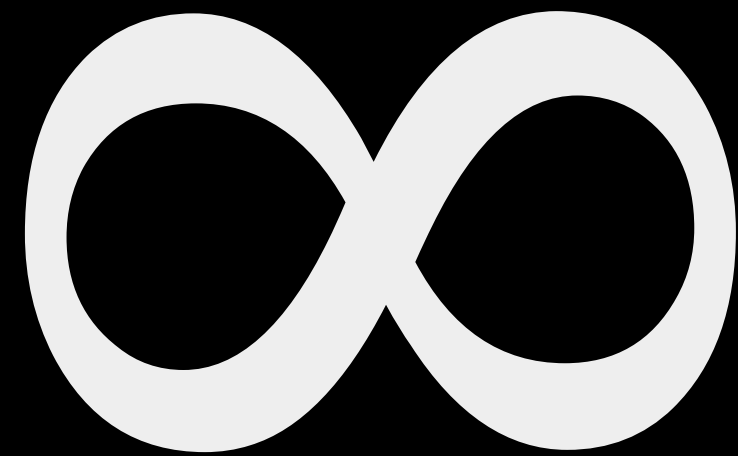
How to make it work?

```
template <size_t Idx> using index_placeholder = placeholder;

template <typename Subject> auto tr_parse(Subject s) {
    return tr_parse(s, std::make_index_sequence<input.length()>());
}

template <typename Subject, size_t... Idx>
    auto tr_parse(Subject s, std::index_sequence<Idx...>) {
        return pair(parse<0>(s) + ... + index_placeholder<Idx>());
    }
```





Thank You!

You can find the slides and sources at:
compile-time.re