

## 4 Assignment 2 [Assignment ID: `cpp_compile_time`]

### 4.1 Preamble (Please Read Carefully)

Before starting work on this assignment, it is **critically important** that you **carefully** read Section 1 (titled “General Information”) which starts on page 1-1 of this document.

### 4.2 Topics Covered

This assignment covers material primarily related to the following: compile-time computation, `constexpr`, literal types.

### 4.3 Problems — Part A — Nonprogramming Exercises

- 8.1 [lvalues/rvalues]
- 8.9 [temporary objects]
- 8.12 [moving vs. copying]
- 8.28 [data structures]

### 4.4 Problems — Part B — Compile-Time Fractal Computation

B.1 *Constexpr basic string class template* (`cexpr_basic_string`). In this exercise, a class template for representing an arbitrary sequence of characters (i.e., a string) that can be used in `constexpr` contexts will be developed. This class template is called `cexpr_basic_string`. The `cexpr_basic_string` class template has two template parameters:

- (a) T. The type of each character in the string (e.g., **char**, **unsigned char**, **wchar\_t**).
- (b) M. The maximum number of characters that can be held in the string (which does not include the dummy null character, discussed later).

The class template uses a C-style array (embedded in the `cexpr_basic_string` object itself) for storing the character data for a string, which is completely arbitrary. The `cexpr_basic_string` class template has the interface given in Listing 1.

The `cexpr_basic_string` class template always stores an additional dummy null-character (i.e., a character with the value `value_type(0)`) immediately following the last character in the string buffer (or at the start of the string buffer if the string is empty). This dummy character is always present, even in the case of an empty string. This dummy character is not considered one of the characters in the string so this dummy character is not included in the count returned by the `size` member function. It is always the case that, for any `cexpr_basic_string` object `s`, `s[s.size()]` refers to the dummy null character (which immediately follows the last character of the string). The above null-termination is employed so that the pointer returned by the data member function can be used in contexts where a null-terminated character array is required.

Listing 1: Interface for class template `cexpr_basic_string`

```

1 namespace ra::cexpr {
2
3     // A basic string class template for use in constexpr contexts.
4     template <class T, std::size_t M>
5     class cexpr_basic_string
6     {
7     public:
8
9         // An unsigned integral type used to represent sizes.
10        using size_type = std::size_t;
11
12        // The type of each character in the string (i.e., an alias for

```

```
13     // the template parameter T).
14     using value_type = T;
15
16     // The type of a mutating pointer to each character in the string.
17     using pointer = T*;
18
19     // The type of a non-mutating pointer to each character in the
20     // string.
21     using const_pointer = const T*;
22
23     // The type of a mutating reference to a character in the string.
24     using reference = T&;
25
26     // The type of a non-mutating reference to a character in the
27     // string.
28     using const_reference = const T&;
29
30     // A mutating iterator type for the elements in the string.
31     using iterator = pointer;
32
33     // A non-mutating iterator type for the elements in the string.
34     using const_iterator = const_pointer;
35
36     // Default construct a string.
37     // Creates an empty string (i.e., a string containing no
38     // characters).
39     //
40     // Time complexity:
41     // Linear in M.
42     constexpr cexpr_basic_string();
43
44     // Copy construct a string.
45     //
46     // Time complexity:
47     // Linear in M.
48     constexpr cexpr_basic_string(const cexpr_basic_string&) =
49         default;
50
51     // Copy assign a string.
52     //
53     // Time complexity:
54     // Linear in M.
55     constexpr cexpr_basic_string& operator=(
56         const cexpr_basic_string&) = default;
57
58     // Destroy a string.
59     //
60     // Time complexity:
61     // Constant.
62     ~cexpr_basic_string() = default;
63
64     // Creates a string with the contents given by the
65     // null-terminated character array pointed to by s.
66     // If the string does not have sufficient capacity to hold
67     // the character data provided, an exception of type
68     // std::runtime_error is thrown.
69     //
```

```
70     // Time complexity:
71     // Linear in the length of the string s.
72     constexpr cexpr_basic_string(const value_type* s);
73
74     // Creates a string with the contents specified by the characters
75     // in the iterator range [first, last).
76     // If the string does not have sufficient capacity to hold
77     // the character data provided, an exception of type
78     // std::runtime_error is thrown.
79     //
80     // Time complexity:
81     // Linear in the size of the range [first, last).
82     constexpr cexpr_basic_string(const_iterator first,
83                               const_iterator last);
84
85     // Returns the maximum number of characters that can be held by a
86     // string of this type.
87     // The value returned is the template parameter M.
88     //
89     // Time complexity:
90     // Constant.
91     static constexpr size_type max_size();
92
93     // Returns the maximum number of characters that the string can
94     // hold. The value returned is always the template parameter M.
95     //
96     // Time complexity:
97     // Constant.
98     constexpr size_type capacity() const;
99
100    // Returns the number of characters in the string (excluding the
101    // dummy null character).
102    //
103    // Time complexity:
104    // Constant.
105    constexpr size_type size() const;
106
107    // Returns a pointer to the first character in the string.
108    // The pointer that is returned is guaranteed to point to a
109    // null-terminated character array.
110    // The user of this class shall not alter the dummy null
111    // character stored at data() + size().
112    //
113    // Time complexity:
114    // Constant.
115    value_type* data();
116    const value_type* data() const;
117
118    // Returns an iterator referring to the first character in the
119    // string.
120    //
121    // Time complexity:
122    // Constant.
123    constexpr iterator begin();
124    constexpr const_iterator begin() const;
125
126    // Returns an iterator referring to the fictitious
```

```

127     // one-past-the-end character in the string.
128     //
129     // Time complexity:
130     // Constant.
131     constexpr iterator end();
132     constexpr const_iterator end() const;
133
134     // Returns a reference to the i-th character in the string if i
135     // is less than the string size; and returns a reference to the
136     // dummy null character if i equals the string size.
137     // Precondition: The index i is such that i >= 0 and i <= size().
138     //
139     // Time complexity:
140     // Constant.
141     constexpr reference operator[](size_type i);
142     constexpr const_reference operator[](size_type i) const;
143
144     // Appends (i.e., adds to the end) a single character to the
145     // string. If the size of the string is equal to the capacity,
146     // the string is not modified and an exception of type
147     // std::runtime_error is thrown.
148     //
149     // Time complexity:
150     // Constant.
151     constexpr void push_back(const T& x);
152
153     // Erases the last character in the string.
154     // If the string is empty, an exception of type std::runtime_error
155     // is thrown.
156     //
157     // Time complexity:
158     // Constant.
159     constexpr void pop_back();
160
161     // Appends (i.e., adds to the end) to the string the
162     // null-terminated string pointed to by s.
163     // Precondition: The pointer s must be non-null.
164     // If the string has insufficient capacity to hold the new value
165     // resulting from the append operation, the string is not modified
166     // and an exception of type std::runtime_error is thrown.
167     //
168     // Time complexity:
169     // Linear in the length of the string s.
170     constexpr cexpr_basic_string& append(const value_type* s);
171
172     // Appends (i.e., adds to the end) to the string another
173     // cexpr_basic_string with the same character type (but
174     // possibly a different maximum size).
175     // If the string has insufficient capacity to hold the new value
176     // resulting from the append operation, the string is not modified
177     // and an exception of type std::runtime_error is thrown.
178     //
179     // Time complexity:
180     // Linear in other.size().
181     template <size_type OtherM>
182     constexpr cexpr_basic_string& append(
183         const cexpr_basic_string<value_type, OtherM>& other);

```

```

184
185     // Erases all of the characters in the string, yielding an empty
186     // string.
187     //
188     // Time complexity:
189     // Constant.
190     constexpr void clear();
191
192 };
193
194 }

```

For information on how to throw an exception, refer to Section 4.6.

Since the `char` type is commonly used as the character type for strings, the following alias template is provided for convenience:

```

namespace ra::cexpr {
    template <std::size_t M>
        using cexpr_string = cexpr_basic_string<char, M>;
}

```

One additional helper function is provided as follows:

```

namespace ra::cexpr {
    constexpr std::size_t to_string(std::size_t n, char* buffer,
        std::size_t size, char** end);
}

```

The `to_string` function converts the integer `n` to its equivalent (decimal) null-terminated string representation. The buffer to be used to store the result starts at the location pointed to by `buffer` and has a size of `size` characters. The resulting string produced by the function is null-terminated. The number of characters written to the buffer, excluding the null character, is returned. If `end` is non-null, `*end` is set to point to the null character at the end of the converted string. If the buffer provided does not have sufficient capacity to hold the string resulting from the conversion process, an exception of type `std::runtime_error` is thrown.

The code for the `cexpr_basic_string` class template and other helper code should be placed in the file `include/ra/cexpr_basic_string.hpp`. Note that the above identifiers (e.g., `cexpr_basic_string`, `cexpr_string`, and `to_string`) are all in the `ra::cexpr` namespace.

Write a program called `test_cexpr_basic_string` to test the code for the `cexpr_basic_string` class template and its associated helper code. The source code for the test program should be placed in the file `app/test_cexpr_basic_string.cpp`.

**B.2 Mandelbrot variable template** (`mandelbrot`). In this exercise, code is developed that can be used to compute an image representation of the Mandelbrot set at compile time. The computed image is made available through a variable template called `mandelbrot`, which is of type `cexpr_string`. This string variable holds the image encoded in the text-based bitmap PNM format.

Let  $\mathbb{C}$  denote the set of complex numbers. The Mandelbrot set  $S$  is the set of all  $c \in \mathbb{C}$  such that the sequence  $z_0, z_1, z_2, \dots$  does not tend toward infinity, where

$$z_n = \begin{cases} z_{n-1}^2 + c & n \geq 1 \\ c & n = 0. \end{cases} \quad (2)$$

As it turns out, the boundary of  $S$  is a fractal curve. More information about the Mandelbrot set can be found at:

[https://en.wikipedia.org/wiki/Mandelbrot\\_set](https://en.wikipedia.org/wiki/Mandelbrot_set)

The Mandelbrot set  $S$  can be represented in the form of a binary image as follows. Define the function  $\chi_S$  that maps  $\mathbb{C}$  to  $\{0, 1\}$  as

$$\chi_S(z) = \begin{cases} 1 & z \in S \\ 0 & \text{otherwise} \end{cases}$$

(i.e.,  $\chi_S$  is effectively a boolean predicate that tests if a complex number is a member of  $S$ ). Let  $F$  denote a binary image function defined on points in  $\Lambda = \{0, 1, \dots, W-1\} \times \{0, 1, \dots, H-1\}$  (i.e., a rectangular grid of width  $W$  and height  $H$ ). Define a sampling function  $\lambda$  that maps a point  $(k, \ell) \in \Lambda$  to a point in the rectangular region  $[a_0, b_0] \times [a_1, b_1]$  of  $\mathbb{C}$  as given by

$$\lambda[(k, \ell)] = \left( a_0 + k \left( \frac{b_0 - a_0}{W-1} \right), a_1 + (H-1-\ell) \left( \frac{b_1 - a_1}{H-1} \right) \right),$$

where  $(a_0, a_1) = (-1.6, -1.1)$  and  $(b_0, b_1) = (0.6, 1.1)$ . The function  $F$  is then given by

$$F[(k, \ell)] = \chi_S(\lambda[(k, \ell)]).$$

The function  $F$  is effectively the Mandelbrot set represented in the form of an image.

The function  $\chi_S$  is computed from the definition of the Mandelbrot set given by (2). In order to ensure some consistency in the results obtained by different implementations, the function  $\chi_S$  must be computed as follows. To determine the value of  $\chi_S(c)$ , the implementation must use (2) to compute  $z_i$  for successively larger values of  $i$  (starting from 0). This iteration stops when either of the following two conditions is met: 1)  $|z_i| > 2$  or 2)  $i = 16$ . If iteration stops due to the first condition, the implementation should assume that the sequence  $z_0, z_1, z_2, \dots$  grows without bound, in which case  $\chi_S(c) = 0$ . If iteration stops due to the second condition, the implementation should assume that the sequence  $z_0, z_1, z_2, \dots$  remains bounded for all  $i$ , in which case  $\chi_S(c) = 1$ .

The `mandelbrot` variable template has two template parameters:

- (a) `W`. The width of the image (in samples).
- (b) `H`. The height of the image (in samples).

The interface for this variable template is as described in Listing 2. In order to minimize the amount of work required for this exercise, the `cexpr_string` class template and `to_string` function developed in Exercise B.1 should be used. The `mandelbrot` variable template must be of type `cexpr_string<M>` for some value of `M`. The string must be encoded in the text-based PNM format for bitmap images. This particular format is described shortly. The particular value of `M` to be used is at the discretion of the implementation. Clearly, however, `M` must be chosen sufficiently large that the `cexpr_string` object can hold the complete PNM-encoded character sequence for an image of width `W` and height `H`. The source for the `mandelbrot` variable template and all of its supporting code should be placed in the file `include/ra/mandelbrot.hpp`.

Listing 2: Interface for the `mandelbrot` variable template

```

1 namespace ra::fractal {
2
3     // A variable template for a string that represents an image depicting
4     // the Mandelbrot set. The image has width W and height H.
5     // This object must be of type cexpr_string<M> for some appropriate M.
6     // The string is a binary image encoded in the text-based bitmap PNM
7     // format.
8     // The values of W and H must be such that W >= 2 and H >= 2.
9     template <std::size_t W, std::size_t H>
10    constexpr auto mandelbrot = implementation-defined;
11
12 }
```

The text-based PNM format for bitmap images is very simple. This format represents a bitmap image of width  $W$  and height  $H$  using a character sequence that consists of the following (in order):

- (a) A signature, which consists of the character “P” followed by the character “1”.
- (b) A space character.
- (c) The width  $W$  of the image, which consists of a sequence of one or more decimal digits.
- (d) A space character.
- (e) The height  $H$  of the image, which consists of a sequence of one or more decimal digits.
- (f) A newline character.
- (g) The (binary) image samples. For each of the  $H$  rows in the image starting with the top row, a string of  $W$  characters that are either “0” or “1” characters, followed by a newline character.

To further illustrate this image format, we now consider a simple example. Consider an image of width 8 and height 4 whose contents resemble a crude letter “V”. This image would be encoded using the following string:

```
P1 8 4
10000001
01000010
00100100
00011000
```

Since the amount of complex arithmetic in this exercise is minimal, it is probably easiest to perform the necessary computations without the formality of using a complex-number class. Unfortunately (at least, as of C++17), the class template `std::complex` does not provide `constexpr` versions of all of the functions that are likely to be needed for the complex arithmetic in this exercise. Consequently, the `std::complex` class template is somewhat less helpful. In any case, regardless of whether a class such as `std::complex` is used, all real-arithmetic should be performed using double-precision (i.e., **double**) arithmetic.

A program called `test_mandelbrot` should be used to test the `mandelbrot` variable template. The source for this program should be placed in the file `app/test_mandelbrot.cpp`. An example of the source for a very trivial test program is shown in Listing 3. This example is intended only for illustrative purposes, however. It is probably advisable to perform more thorough testing than what is achieved by this trivial program. The image generated by this trivial test program is shown in Figure 1.

Listing 3: Source listing for a very trivial example of a `test_mandelbrot` program.

```
1 #include <iostream>
2 #include "ra/mandelbrot.hpp"
3
4 int main()
5 {
6     // Force the image (in PNM format) to be computed at compile time.
7     constexpr auto s = ra::fractal::mandelbrot<256, 256>;
8
9     // Output the image (in PNM format).
10    std::cout << s.begin() << '\n';
11 }
```

Admittedly, this exercise is somewhat of an abuse of compile-time computation in the sense that one would probably not normally perform this much computation at compile time. Nevertheless, this exercise serves to clearly demonstrate the power of compile-time computation. In the `CMakeLists` file, it may be necessary to add a compiler flag to increase the maximum allowable amount of compile-time computation for the compilation of the source file `app/test_mandelbrot.cpp`. (In particular, adding such a flag is necessary if the compilation of this source file fails due to the compile-time computation limits being exceeded.) In the case of GCC, the maximum allowable amount of compile-time computation can be controlled with the `-fconstexpr-loop-limit` option. This option takes a parameter that is an integer value specifying the loop limit in `constexpr` computation. A value of 10000 should be sufficient to compute the value of `mandelbrot<512, 512>`. In the case of Clang, the maximum allowable amount of compile-time computation can be controlled with the `-fconstexpr-steps` option. This option takes a parameter that is an integer value specifying the maximum number of steps in `constexpr`

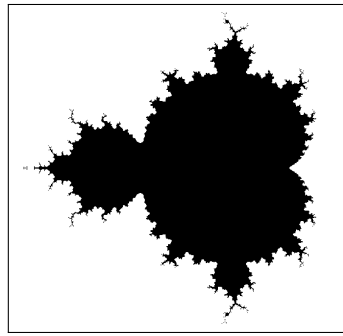


Figure 1: Mandelbrot set image.

computation. A value of 100000000 should be sufficient to compute the value of `mandelbrot<512, 512>`. Note that the option for specifying the maximum amount of compile-time computation should only be used for the compilation of the file `app/test_mandelbrot.cpp` (not for any other files). With CMake, adding a compile flag to only a single source file can be easily accomplished with the `set_source_files_properties` command. This might be accomplished with a command like the following in the `CMakeLists.txt` file:

```
# Set the variable EXTRA_COMPILE_FLAGS to the desired compile flags
# to add for the particular compiler being used (as explained above).
# For GCC:
#   set(EXTRA_COMPILE_FLAGS "-fconstexpr-loop-limit=10000")
# For Clang:
#   set(EXTRA_COMPILE_FLAGS "-fconstexpr-steps=100000000")
# Add the desired flags to the compilation of the
# app/test_mandelbrot.cpp file.
set_source_files_properties(app/test_mandelbrot.cpp PROPERTIES
    COMPILER_FLAGS ${EXTRA_COMPILE_FLAGS})
```

As part of the testing of the code for the `mandelbrot` variable template, confirm that the generated image does, in fact, resemble the Mandelbrot set. This can be done by printing the (string) value of the `mandelbrot` variable template to a file and then displaying this file with an image viewer that supports the PNM format. For example, the following command can be used to display a file called `mandelbrot.pnm` (in PNM format) with the ImageMagick software:

```
display mandelbrot.pnm
```

## 4.5 Problems — Part C — Compile-Time Filter Design

C.1 *Constexpr math constants and functions (i.e., `pi`, `sin`, `cos`, `tan`, `sqrt`, and others)*. In this exercise, numerous variable and function templates are developed that provide support for math constants (such as  $\pi$ ) and math functions (such as `sin` and square root) that can be used in `constexpr` contexts.

Each of these templates has a template parameter `T` that corresponds to the floating-point type to be used to represent real numbers. Only floating-point types (i.e., `float`, `double`, and `long double`) can be used for `T`. The interface provided by the above variable and function templates is given in Listing 4. (Note that all of these templates are in the namespace `ra::constexpr_math`.) The source for these variable and function templates is to be placed in the file `include/ra/constexpr_math.hpp`.

Listing 4: Interface for `constexpr_math` functions

```
1 namespace ra::constexpr_math {
2
3     // The math constant pi.
```



```
4 // The type T is a floating-point type.
5 template <class T>
6 constexpr T pi = std::numbers::pi_v<T>;
7
8 // Returns the absolute value of x.
9 // The type T is a floating-point type.
10 template <class T>
11 constexpr T abs(T x);
12
13 // Returns the square of x.
14 // The type T is a floating-point type.
15 template <class T>
16 constexpr T sqr(T x);
17
18 // Returns the cube of x.
19 // The type T is a floating-point type.
20 template <class T>
21 constexpr T cube(T x);
22
23 // Returns the remainder after division when x is divided by y.
24 // In particular, x - n y is returned where n is the result obtained by
25 // dividing x by y and then rounding (to an integer value) toward zero.
26 // If y is zero, an exception of type std::overflow_error is thrown.
27 // The type T is a floating-point type.
28 template <class T>
29 constexpr T mod(T x, T y);
30
31 // Returns the sine of x.
32 // Note that a particular algorithm must be used to implement this
33 // function.
34 // The type T is a floating-point type.
35 template <class T>
36 constexpr T sin(T x);
37
38 // Returns the cosine of x.
39 // Note that a particular algorithm must be used to implement this
40 // function.
41 // The type T is a floating-point type.
42 template <class T>
43 constexpr T cos(T x);
44
45 // Returns the tangent of x.
46 // Note that a particular algorithm must be used to implement this
47 // function.
48 // If the tangent of x is infinite, an exception of type
49 // std::overflow_error is thrown.
50 // The type T is a floating-point type.
51 template <class T>
52 constexpr T tan(T x);
53
54 // Returns the square root of x.
55 // If x is negative, an exception of type std::domain_error is thrown.
56 // Note that a particular algorithm must be used to implement this
57 // function.
58 // The type T is a floating-point type.
59 template <class T>
60 constexpr T sqrt(T x);
```

61  
62 }

For information on how to throw an exception, refer to Section 4.6.

As mentioned in the interface description, the implementation technique for some of the function templates must follow a particular approach. In what follows, these restrictions on the implementation strategies are discussed.

The `sin` function must be implemented by using the triple-angle formula and small-angle approximation for `sin`. The triple-angle formula and small-angle approximation are respectively given by

$$\sin 3x = 3 \sin x - 4 \sin^3 x \quad \text{and} \quad (3)$$

$$\sin x \approx x \quad \text{for small nonnegative } x. \quad (4)$$

By rearranging (3), we can obtain the following recursive equation for `sin`:

$$\sin x = 3 \sin(x/3) - 4 \sin^3(x/3). \quad (5)$$

To implement the `sin` function, (5) should be used recursively for the computation of `sin` with the base case for the recursion corresponding to  $x \leq 10^{-6}$ . In this base case, `sin` is simply assumed to be equal to `x`, with this assumption being valid for nonnegative `x` due to (4). (Computing `sin` for negative `x` can be handled by using the fact that `sin` is an odd function.) The maximum recursion depth for this algorithm could potentially become large if `sin` is computed for large magnitude `x`. To eliminate this problem, the algorithm must use the fact that `sin` is  $2\pi$ -periodic in order to reduce the problem of computing `sin` for arbitrary `x` to that of computing `sin` for  $x \in [0, 2\pi)$ . (That is, do not apply (5) to the problem of computing `sin` until first ensuring that  $x \in [0, 2\pi)$ .) To reduce the problem in this way, the `mod` function will likely be helpful.

The `cos` function must be implemented by using the `sin` function and the fact that `cos` is identical to `sin` except for a translation (i.e., shift).

The `tan` function must be implemented by using the `sin` and `cos` functions and the fact that  $\tan x = \frac{\sin x}{\cos x}$  (for  $\cos x \neq 0$ ). The implementation must not divide by zero, under any circumstances. If a circumstance arises that would lead to a division by zero, an exception should be thrown (as explained in the interface definition).

The `sqrt` function must be implemented by using the Newton-Raphson (root-finding) method. In the Newton-Raphson method, to solve for a root of the equation  $f(x) = 0$ , we select an initial estimate  $x_0$  of the root. Then, we apply the following iterative process to improve the accuracy of our initial estimate:

$$x_{n+1} = x_n - f(x_n)/f'(x_n), \quad (6)$$

where  $f'$  denotes the first derivative of  $f$ . To find the square root of  $c$ , we can simply solve for the (real non-negative) root of the equation  $f(x) = x^2 - c$ . The initial estimate  $x_0$  of the root should be chosen as  $c$ . The iteration in (6) should continue until  $|x_{n+1} - x_n| \leq \epsilon$  (i.e., the root estimate does not change by more than  $\epsilon$  from one iteration to the next) or  $n > m$ , where the tolerance  $\epsilon$  and the maximum iteration count  $m$  can be chosen as `std::numeric_limits<T>::epsilon()` and `std::numeric_limits<T>::max_exponent`, respectively. (Note that these choices of  $\epsilon$  and  $m$  are overkill, but at least this eliminates the need to worry that the computed square root will fail to be accurate enough for our application.) If the number whose square root is to be computed is negative, an exception should be thrown (as explained in the interface definition).

It is absolutely critical that the exact algorithms specified above be employed for `sin`, `cos`, `tan`, and `sqrt`. Failure to do so will likely result in implementations of these functions with accuracy properties that are very substantially different from what is required, which could lead to your code failing many test cases (due to results that are not sufficiently accurate).

A program called `test_cexpr_math` is to be developed to test the variable and function templates from above. The source for this test program is to be placed in the file `app/test_cexpr_math.cpp`.

**C.2 Biquad filter design functions.** In this exercise, code is developed that can be used to design several types of discrete-time biquad filters. This code consists of a class template `biquad_filter_coefs` that is used to represent the coefficients of a biquad filter as well as several function templates that can be used to design various types of filters. These filter-design function templates can be employed in `constexpr` contexts.

In audio and music processing applications, biquad filters are often employed. A (real-coefficient) discrete-time biquad filter has a transfer function  $H$  of the form

$$H(z) = \frac{a_0 + a_1z^{-1} + a_2z^{-2}}{b_0 + b_1z^{-1} + b_2z^{-2}}, \quad (7)$$

where  $a_0, a_1, a_2, b_0, b_1,$  and  $b_2$  are real coefficients. Such a transfer function can always be normalized such that the constant term in the denominator is 1 (provided  $b_0 \neq 0$ , which will always be the case here). That is, such a transfer function can always be expressed in the form

$$H(z) = \frac{a'_0 + a'_1z^{-1} + a'_2z^{-2}}{1 + b'_1z^{-1} + b'_2z^{-2}}. \quad (8)$$

This can be accomplished by simply dividing each of the numerator and denominator of  $H(z)$  by  $b_0$ . For various reasons, it is often more convenient to express biquad transfer functions in this normalized form.

In discrete-time signal processing, the sampling frequency is an important quantity. Often, rather than specifying frequencies directly, it is more convenient to specify them relative to the sampling frequency, leading to the notion of normalized frequency. The normalized frequency  $f$  that corresponds to the actual frequency  $f'$  and sampling frequency  $f_s$  is given by  $f = 2f'/f_s$ . That is, the normalized frequency 1 corresponds to the Nyquist frequency (i.e.,  $f_s/2$ ). Note that the normalized frequency is always a value in  $[0, 1]$ . In what follows, the term “normalized frequency” refers to this definition just introduced.

Although many types of biquad filters are possible, we only consider the design of the following basic types: 1) lowpass, 2) highpass, 3) bandpass, 4) low-frequency shelving boost filter. and 5) low-frequency shelving cut filter. A lowpass biquad filter with normalized cutoff frequency  $f$  and Q factor  $Q$  has the filter coefficients given by

$$\begin{aligned} \Omega &= \tan\left(\frac{\pi}{2}f\right), \\ a_0 &= \Omega^2, \quad a_1 = 2\Omega^2, \quad a_2 = \Omega^2, \\ b_0 &= \Omega^2 + \Omega/Q + 1, \quad b_1 = 2(\Omega^2 - 1), \quad \text{and} \quad b_2 = \Omega^2 - \Omega/Q + 1. \end{aligned}$$

A highpass filter with normalized cutoff frequency  $f$  and Q factor  $Q$  has the filter coefficients given by

$$\begin{aligned} \Omega &= \tan\left(\frac{\pi}{2}f\right), \\ a_0 &= 1, \quad a_1 = -2, \quad a_2 = 1, \\ b_0 &= \Omega^2 + \Omega/Q + 1, \quad b_1 = 2(\Omega^2 - 1), \quad \text{and} \quad b_2 = \Omega^2 - \Omega/Q + 1. \end{aligned}$$

A bandpass filter with normalized center frequency  $f$  and Q factor  $Q$  has the filter coefficients given by

$$\begin{aligned} \Omega &= \tan\left(\frac{\pi}{2}f\right), \\ a_0 &= \Omega/Q, \quad a_1 = 0, \quad a_2 = -\Omega/Q, \\ b_0 &= \Omega^2 + \Omega/Q + 1, \quad b_1 = 2(\Omega^2 - 1), \quad \text{and} \quad b_2 = \Omega^2 - \Omega/Q + 1. \end{aligned}$$

A low-frequency shelving boost filter with normalized cutoff frequency  $f$  and gain-control parameter  $A$  has the filter coefficients given by

$$\begin{aligned} \Omega &= \tan\left(\frac{\pi}{2}f\right), \\ a_0 &= A\Omega^2 + \sqrt{2A}\Omega + 1, \quad a_1 = 2(A\Omega^2 - 1), \quad a_2 = A\Omega^2 - \sqrt{2A}\Omega + 1, \\ b_0 &= \Omega^2 + \sqrt{2}\Omega + 1, \quad b_1 = 2(\Omega^2 - 1), \quad b_2 = \Omega^2 - \sqrt{2}\Omega + 1. \end{aligned}$$

A low-frequency shelving cut filter with normalized cutoff frequency  $f$  and gain-control parameter  $A$  has the filter coefficients given by

$$\begin{aligned}\Omega &= \tan\left(\frac{\pi}{2}f\right), \\ a_0 &= \Omega^2 + \sqrt{2}\Omega + 1, \quad a_1 = 2(\Omega^2 - 1), \quad a_2 = \Omega^2 - \sqrt{2}\Omega + 1, \\ b_0 &= A\Omega^2 + \sqrt{2A}\Omega + 1, \quad b_1 = 2(A\Omega^2 - 1), \quad b_2 = A\Omega^2 - \sqrt{2A}\Omega + 1.\end{aligned}$$

To provide a convenient container for holding the coefficients of a biquad filter, a class template called `biquad_filter_coefs` is used. This class template has a single template parameter `Real`, which specifies the real-number type used to represent the filter coefficients. This template can be instantiated with `Real` being any floating-point type (i.e., **float**, **double**, or **long double**). The `biquad_filter_coefs` class template has the very simple interface provided in Listing 5. This class template only has public members and all such members are identified in the interface description.

Listing 5: Interface for `biquad_filter_coefs` class template

```

1 namespace ra::biquad {
2
3     // Biquad filter coefficients class.
4     template <class Real>
5     struct biquad_filter_coefs
6     {
7         // The real number type used to represent the filter coefficients.
8         using real = Real;
9
10        // Creates a set of filter coefficients where the coefficients
11        // a0, a1, a2, b0, b1, and b2 are initialized to a0_, a1_, a2_,
12        // b0_, b1_, and b2_, respectively.
13        constexpr biquad_filter_coefs(real a0_, real a1_, real a2_, real b0_,
14        real b1_, real b2_);
15
16        // Creates a set of filter coefficients by copying from another set.
17        // Since Real and OtherReal need not be the same, this constructor
18        // can be used to convert between filter coefficients of different
19        // types.
20        template <class OtherReal>
21        constexpr biquad_filter_coefs(
22        const biquad_filter_coefs<OtherReal>& coefs);
23
24        // The filter coefficients a0, a1, a2, b0, b1, and b2.
25        real a0;
26        real a1;
27        real a2;
28        real b0;
29        real b1;
30        real b2;
31    };
32
33 }
```

Several function templates are provided for designing various types of biquad filters. These function templates have the interface shown in Listing 6. Each of these functions has a template parameter `Real`, which controls the real-number type used for computing filter coefficients. Depending on how accurate the computed filter coefficients need to be, the template could be instantiated with `Real` being any one of the floating-point types (i.e., **float**, **double**, or **long double**).

Listing 6: Interface for biquad filter design functions

```

1 namespace ra::biquad {
2
3     // Returns the coefficients for a biquad lowpass filter with normalized
4     // cutoff frequency f and Q factor q, where f in [0,1] and q > 0.
5     // The filter coefficients are always normalized such that the
6     // coefficient b0 is 1.
7     // The type Real is a floating-point type.
8     // All real arithmetic should be performed with the Real type.
9     template <class Real>
10    constexpr biquad_filter_coefs<Real> lowpass(Real f, Real q);
11
12    // Returns the coefficients for a biquad highpass filter with
13    // normalized cutoff frequency f and Q factor q, where f in [0,1]
14    // and q > 0.
15    // The filter coefficients are always normalized such that the
16    // coefficient b0 is 1.
17    // The type Real is a floating-point type.
18    // All real arithmetic should be performed with the Real type.
19    template <class Real>
20    constexpr biquad_filter_coefs<Real> highpass(Real f, Real q);
21
22    // Returns the coefficients for a biquad bandpass filter with
23    // normalized cutoff frequency f and Q factor q, where f in [0,1]
24    // and q > 0.
25    // The filter coefficients are always normalized such that the
26    // coefficient b0 is 1.
27    // The type Real is a floating-point type.
28    // All real arithmetic should be performed with the Real type.
29    template <class Real>
30    constexpr biquad_filter_coefs<Real> bandpass(Real f, Real q);
31
32    // Returns the coefficients for a biquad low-frequency shelving
33    // boost filter with normalized cutoff frequency f and gain-control
34    // parameter a, where f in [0,1] and a >= 0.
35    // For a gain in dB of G (where G > 0), choose a = 10 ^ (G / 20).
36    // The filter coefficients are always normalized such that the
37    // coefficient b0 is 1.
38    // The type Real is a floating-point type.
39    // All real arithmetic should be performed with the Real type.
40    template <class Real>
41    constexpr biquad_filter_coefs<Real> lowshelf_boost(Real f, Real a);
42
43    // Returns the coefficients for a biquad low-frequency shelving
44    // cut filter with normalized cutoff frequency f and gain-control
45    // parameter a, where f in [0,1] and a >= 0.
46    // For a gain in dB of G (where G < 0), choose a = 10 ^ (-G / 20).
47    // The filter coefficients are always normalized such that the
48    // coefficient b0 is 1.
49    // The type Real is a floating-point type.
50    // All real arithmetic should be performed with the Real type.
51    template <class Real>
52    constexpr biquad_filter_coefs<Real> lowshelf_cut(Real f, Real a);
53
54 }

```

The code for the `biquad_filter_coefs` class template and the filter-design function templates should be placed

in the file `include/ra/biquad_filter.hpp`.

A program called `test_biquad_filter` should be used to test the filter-design code. The source for this program should be placed in the file `app/test_biquad_filter.cpp`. An online biquad filter coefficient calculator, which may prove useful for testing purposes, can be found at:

<http://www.earlevel.com/main/2021/09/02/biquad-calculator-v3>

(In the case of low-frequency shelving filters, this online calculator selects the boost and cut cases if the gain  $G$  in dB is positive and negative, respectively.) For example, using the above online calculator, one can confirm that the filter coefficients obtained for a lowpass filter with a normalized cutoff frequency of  $\frac{1}{4}$  and a Q factor of  $\frac{1}{\sqrt{2}}$  are as follows:

$$a_0 \approx 0.09763076084032914, \quad a_1 \approx 0.19526152168065827, \quad a_2 \approx 0.09763076084032914, \\ b_0 \approx 1.0, \quad b_1 \approx -0.9428060277021066, \quad \text{and} \quad b_2 \approx 0.3333290710634233.$$

## 4.6 Comments on Exceptions

An exception is an object that describes an error condition that occurs during the execution of a program. An error condition can be signalled by throwing an exception, which is accomplished with a `throw` statement. The standard library defines numerous exception types in the header file `stdexcept` (e.g., `std::runtime_error`, `std::overflow_error`, and `std::domain_error`). Each standard library exception type provides a constructor that takes a pointer to a character array (i.e., string). This string provides some additional information about the error that the exception is signalling. For example, if one is asked to take the square root of a negative number, one might throw a domain error (i.e., `std::domain_error`) as in the code example given below.

```
1  #include <stdexcept>
2
3  double sqrt(double x)
4  {
5      double result = 0.0;
6      if (x < 0) {
7          // We are attempting to take the square root of a negative number.
8          throw std::domain_error("square root of negative number");
9      }
10     // ... (initialize result)
11     return result;
12 }
```

More code examples that throw exceptions can be found in the “Exceptions” section of the lecture slides. If additional information related to throwing exceptions is desired, the following URLs might also be helpful:

- <http://en.cppreference.com/w/cpp/error/exception>
- <http://en.cppreference.com/w/cpp/header/stdexcept>