

Developing a Digital Synthesizer in C++

Peter Goldsborough
`petergoldsborough@hotmail.com`

December 21, 2014

Contents

1	Sound in the Digital Realm	3
1.1	What is a Music Synthesizer?	3
1.2	From Analog to Digital	3
1.3	Sample Rate	4
1.4	Nyquist Limit	4
1.5	Aliasing	5
1.6	Overflow	7
2	Generating Sound	8
2.1	Simple Waveforms	8
2.2	Complex Waveforms	8
2.2.1	Mathematical Calculation of Complex Waveforms	10
2.2.2	Additive Synthesis	12
2.2.3	Special Waveforms	18
2.3	The Gibbs Phenomenon and the Lanczos Sigma Factor	20
2.4	Real-Time Calculation versus Wavetable Lookup	21
2.4.1	Implementation	21
2.4.2	Interpolation	22
2.4.3	Table Length	22
2.4.4	File Format	23

Chapter 1

Sound in the Digital Realm

1.1 What is a Music Synthesizer?

1.2 From Analog to Digital

Our everyday experience of sound is an entirely analog one. When a physical object emits or reflects a sound wave into space and towards our ears, the signal produced consists of an infinite set of values, spaced apart in infinitesimal intervals. Due to the fact that such a signal has an amplitude value at every single point in time, it is called a continuous signal. (Smith, 1999, p. 11) Figure 1.1 displays the continuous representation of a sine wave.

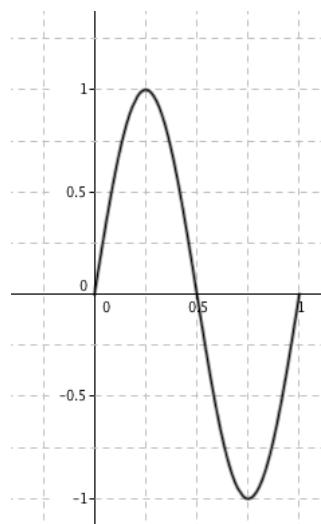


Figure 1.1: The continuous representation of a typical sine wave. In this case, both the signal's frequency f as well as the maximum elongation from the equilibrium a are equal to 1.

While continuous signals and the idea of an infinite, uncountable set of values are easy to model in mathematics and physics — the analog world, computers — in the digital world — effectively have no means by which to represent something that is infinite, since computer memory is a finite resource. Therefore, signals in the digital domain are discrete, meaning they are composed of periodic *samples*. A sample is a discrete recording of a continuous signal's amplitude, taken in a constant time interval (Mitchell, 2008, p. 16). The process by which a continuous signal is converted to a discrete signal is called *quantization*, *digitization* or simply *analog-to-digital-conversion* (Smith, 1999, p. 35-36) (Mitchell, 2008, p. 16). The reverse process of converting discrete samples to a continuous signal is called *digital-to-analog-conversion* (Mitchell, 2008, p. 17). Figure 1.2 shows the discrete representation of a sine wave, the same signal that was previously shown as a continuous signal in Figure 1.1.

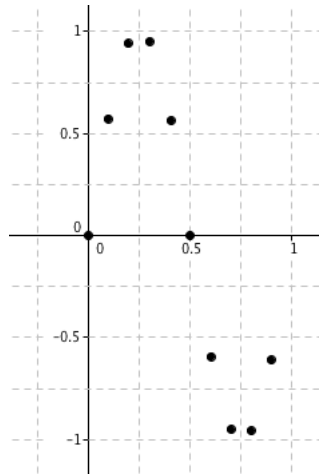


Figure 1.2: The discrete representation of a typical sine wave.

1.3 Sample Rate

The sample rate (often referred to as sampling rate or sampling frequency), commonly denoted by f_s , is the rate at which samples of a continuous signal are taken. The value of the sample rate is given in Hertz (Hz) or samples-per-second. Common values for audio sampling rates are 44.1 kHz, a frequency originally chosen by Sony in 1979 that is still used for Compact Discs, and 48 kHz, the standard audio sampling rate used today (Mitchell, 2008, p. 18) (Colletti, 2013). The reciprocal of the sample rate yields the sampling interval, denoted by T_s and measured in seconds, which is the time period after which a single sample is taken from a continuous signal:

$$T_s = \frac{1}{f_s}$$

The reciprocal of the sample interval again yields the sampling rate:

$$f_s = \frac{1}{T_s}$$

1.4 Nyquist Limit

The sample rate determines the range of frequencies that can be represented by a digital sound system, as only frequencies that are less than or equal to one half of the sampling rate, where

it is possible to take at least one sample above the equilibrium and at least one sample below the equilibrium for every cycle (Mitchell, 2008, p. 18), can be "properly sampled". To sample a signal "properly" means to be able to "reconstruct" a continuous signal, given a set of discrete samples, "exactly", i.e. without any *quantization errors*. The value of one half of the sample rate is called the *Nyquist frequency* or *Nyquist limit*, named after Harry Nyquist, who first described the Nyquist limit and associated phenomena together with Claude Shannon in the 1940s, stating that "a continuous signal can be properly sampled, only if it does not contain frequency components above one half the sampling rate" (Smith, 1999, p. 40). Any frequencies above the Nyquist limit lead to *aliasing*, which is discussed in the next section.

Given the definition of the Nyquist limit and considering the fact that the limit of human hearing is approximately 20 kHz (Cutnell & Johnson, 1998, p. 466), the reason for which the two most common audio sample rates are 40 kHz and above is clear: they were chosen to allow the "proper" representation of the entire audio frequency range, since a sample rate of 40 kHz meets the Nyquist requirement of a sample rate at least twice the maximum frequency component of the signal to be sampled (the Nyquist limit), in this case ca. 20 KHz.

1.5 Aliasing

When a signal's frequency exceeds the Nyquist limit, it is said to produce an *alias*, a new signal with a different frequency that is indistinguishable from the original signal when sampled. This is due to the fact that a signal with a frequency component above the Nyquist limit no longer has one sample taken above and one below the zero level for each cycle, but at arbitrary points of the original signal, which, when reconstructed, yield an entirely different signal. For example, if the sinusoid depicted in Figure 1.3, with a frequency of 4 Hz, is sampled at a sample rate of 5 samples per second, shown in Figure 1.4 meaning the frequency of the continuous signal is higher than the Nyquist limit (here 2.5 Hz), the reconstructed signal, approximated in Figure 1.5, will look completely different from the original sinusoid. "This phenomenon of [signals] changing frequency during sampling is called aliasing, [...] an example of improper sampling" (Smith, 1999, p. 40).

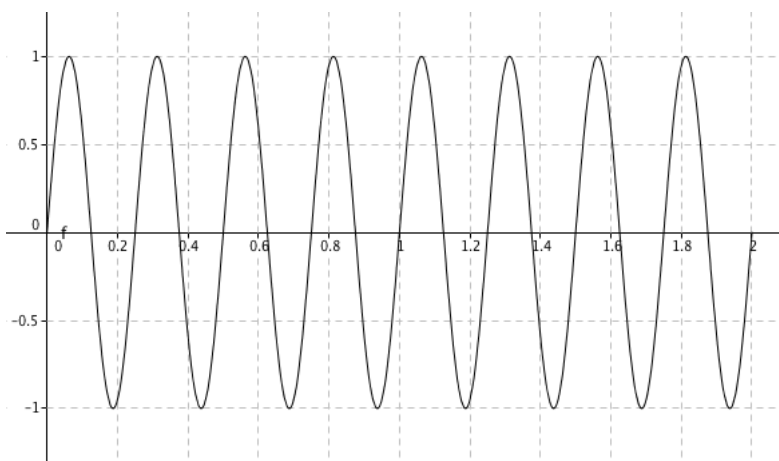


Figure 1.3: A sinusoid with a frequency of 4 Hz.

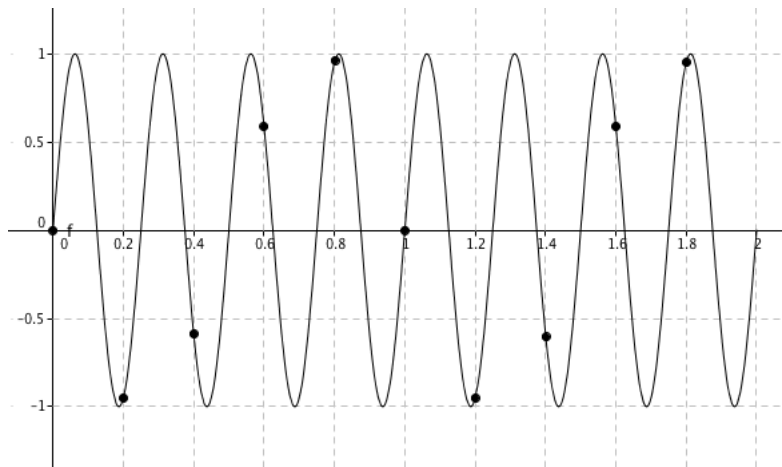


Figure 1.4: A sinusoid with a frequency of 4 Hz, sampled at a sample rate of 5 Hz. According to the Nyquist Theorem, this signal is sampled improperly, as the frequency of the continuous signal is not less than or equal to one half of the sample rate, in this case 2.5 Hz.

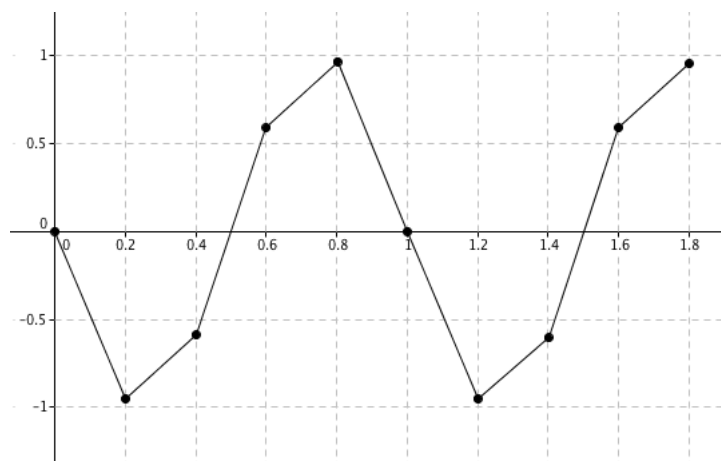


Figure 1.5: An approximated representation of the signal created from the sampling process shown in Figure 1.4. Visually as well as acoustically, the new signal is completely different from the original signal. However, in terms of sampling, they are indistinguishable from each other due to improper sampling. These signals are said to be *aliases* of each other.

1.6 Overflow

Another interesting property of digital sound, which is not encountered in the analog world, is that it can overflow. When we attempt to increase the loudness of something in the analog world, e.g. by hitting a drum more intensely, the expected result is a louder sound. In the digital realm however, it may occur that attempting to increase a signal's amplitude does not result in an increased loudness, but in distortion. The cause of this phenomenon lies in the way digital audio is stored. Since computer memory is a finite resource, each sample has a dedicated portion of computer memory allocated to it. For example, the Waveform Audio File Format (WAVE), a common computer file format for audio data, stores each sample of an audio track as a 16-bit signed integer. A 16-bit signed integer gives a possible range of -2^{16-1} to 2^{16-1} (16 – 1 because the most significant bit is used as the sign bit in two's complement representation). This means that a signal with an amplitude of 1 will be stored as 32767, an amplitude of 0.5 as 16384, an amplitude of -1 as -32768 and so on. If one tries to increase the amplitude of a signal whose value has already saturated the available range and space allocated to it, in this case 32767 on the positive end and -32768 on the negative end, the result is that the integer with which the sample is stored *overflows*. Because WAVE files (and many other storage media) store samples as *signed* integers, overflow always results in a change of sign:

$$32767_{10} = 0111111111111111_2$$

$$0111111111111111_2 + 1 = 1000000000000000_2 = -32768_{10}$$

$$1000000000000000_2 - 1 = 0111111111111111_2 = 32767_{10}$$

A visualization of the result of increasing the amplitude of a signal with a saturated value (an amplitude of 1), shown in Figure 1.6, is given in Figure 1.7.

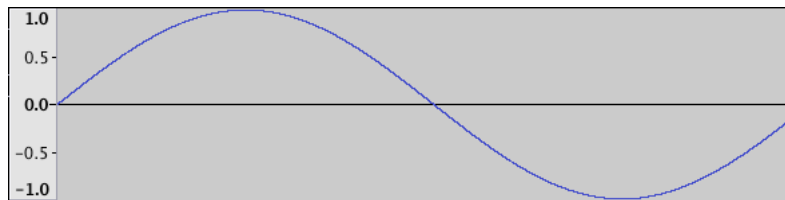


Figure 1.6: A typical sinusoidal signal with an amplitude of 1. The integer range provided by the allocated memory for the top-most sample is saturated, meaning it is equal to 32767_{10} or 0111111111111111_2 .

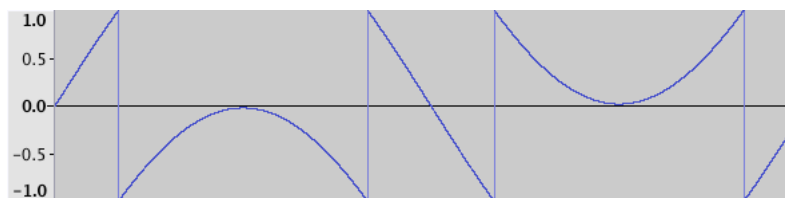


Figure 1.7: What happens when the amplitude of the signal from Figure 1.6 is increased by a factor of 2. Multiple samples have overflowed and thus changed their sign. Because of the way two's-complement representation is implemented, the signal continues its path as if no overflow had ever occurred. The only difference being, of course, that the sign has changed mid-way.

Chapter 2

Generating Sound

The following sections will outline how digital sound can be generated in theory and implemented in practice, using the C++ programming language.

2.1 Simple Waveforms

The simplest possible waveform is the sine wave. As a function of time, it can be mathematically represented by Equation 2.1, where A is the maximum amplitude of the signal, f the frequency in Hertz and ϕ an initial phase offset in radians:

$$f_s(t) = A \sin(2\pi ft + \phi) \quad (2.1)$$

A computer program to compute the values of a sine wave with a variable duration, implemented in C++, is shown in Table 2.1. Another waveform similar to the sine wave is the cosine wave, which differs only in a phase offset of 90° or $\frac{\pi}{2}$ radians:

$$f_c(t) = A \cos(2\pi ft + \phi) = A \sin(2\pi ft + \phi + 90) \quad (2.2)$$

Therefore, the program from Table 2.1 could be modified to compute a cosine wave by changing line 22 from:

```
22 | double phase = 0;  
  
to  
22 | double phase = pi/2.0;
```

2.2 Complex Waveforms

Now that the process of creating simple sine and cosine waves has been discussed, the generation of more complex waveforms can be examined. Generally, there are two methods by which complex waveforms can be created in a digital synthesis system: mathematical calculation or additive synthesis.


```

1  #include <cmath>
2
3  int main(int argc, char * argv[])
4  {
5      // The sample rate, 48 kHz.
6      const unsigned short samplerate = 48000;
7
8      // The duration of the generated sine wave, in seconds.
9      const unsigned long duration = 1;
10
11     // The number of samples that will be generated, derived
12     // from the number of samples per second (the sample rate)
13     // and the number of seconds to be generated for.
14     const unsigned long numberOfSamples = duration * samplerate;
15
16     const double pi = 3.141592653589793;
17
18     const double twoPi = 6.28318530717958;
19
20     // The frequency of the sine wave
21     double frequency = 1;
22
23     // The phase counter. This variable can be seen as phi.
24     double phase = 0;
25
26     // The amount by which the phase is incremented for each
27     // sample. Since one period of a sine wave has 2 pi radians,
28     // dividing that value by the sample rate yields the amount
29     // of radians by which the phase needs to be incremented to
30     // reach a full 2 pi radians.
31     double phaseIncrement = frequency * twoPi / samplerate;
32
33     // The maximum amplitude of the signal, should not exceed 1.
34     double maxAmplitude = 0.8;
35
36     // The buffer in which the samples will be stored.
37     double * buffer = new double[numberOfSamples];
38
39     // For every sample.
40     for (unsigned long n = 0; n < numberOfSamples; ++n)
41     {
42         // Calculate the sample.
43         buffer[n] = maxAmplitude * sin(phase);
44
45         // Increment the phase by the appropriate
46         // amount of radians.
47         phase += phaseIncrement;
48
49         // Check if two pi have been reached and
50         // reset if so.
51         if (phase >= twoPi)
52         {
53             phase -= twoPi;
54         }
55     }
56
57     // Further processing ...
58
59     // Free the buffer memory.
60     delete [] buffer;
61 }

```

Table 2.1: C++ implementation of a complete sine wave generator.

2.2.1 Mathematical Calculation of Complex Waveforms

In the first case — mathematical calculation, waveforms are computed according to certain mathematical formulae and thus yield *perfect* or *exact* waveforms, such as a square wave that is equal to the maximum amplitude exactly one half of a period and equal to the minimum amplitude for the rest of the period. While these waveforms produce a very crisp and clear sound, they are rarely found in nature due to their degree of perfection and are consequently rather undesirable for a music synthesizer. Nevertheless, they are considerably useful for modulating other signals, as tiny acoustical imperfections such as those found in additively synthesized waveforms can result in unwanted distortion which is not encountered when using mathematically calculated waveforms. Therefore, exact waveforms are the best choice for modulation sources such as Low Frequency Oscillators (LFOs), which are discussed in later chapters. (Mitchell, 2008, p. 71)

The following paragraphs will analyze how four of the most common waveforms found in digital synthesizers, the square, the sawtooth, the ramp and the triangle wave, can be generated via mathematical calculation.

Note: There have found to be disparities in literature over which waveform is a sawtooth and which a ramp wave. This thesis will consider a sawtooth wave as descending from maximum to minimum amplitude with time and a ramp wave as ascending from minimum to maximum amplitude with time.

Square Waves

Ideally, a square wave is equal to its maximum amplitude for exactly one half of a period and equal to its minimum amplitude for the other half of the same period. A single period of a square wave can be calculated as shown in Equation 2.3, where the independent variable t as well as the period T can be either in samples or in seconds. A mathematical Equation for a full, periodic square wave function is given by Equation 2.4, where t is time in seconds and the frequency f in Hertz. An equivalent C++ computer program is shown in Table 2.2.

$$f(t) = \begin{cases} 1, & \text{if } 0 \leq t < \frac{T}{2} \\ -1, & \text{if } \frac{T}{2} \leq t < T \end{cases} \quad (2.3) \quad f(t) = \begin{cases} 1, & \text{if } \sin(2\pi ft) > 0 \\ -1, & \text{otherwise} \end{cases} \quad (2.4)$$

Sawtooth Waves

An ideal sawtooth wave descends from its maximum amplitude to its minimum amplitude linearly before jumping back to the maximum amplitude at the beginning of the next period. A mathematical Equation for a single period of such a sawtooth wave function, calculated directly from the phase, is given by Equation 2.5 (Mitchell, 2008, p. 68). Alternatively, the function can depend on time or on samples, as shown by Equation 2.6, where T is the period. A computer program to compute one period of a sawtooth wave is given in Table 2.3.

$$f(\phi) = \begin{cases} -\frac{\phi}{\pi} + 1, & \text{if } 0 \leq \phi < 2\pi \end{cases} \quad (2.5) \quad f(t) = \begin{cases} -\frac{2t}{T} + 1, & \text{if } 0 \leq t < 1 \end{cases} \quad (2.6)$$

```

1 double* square(const unsigned int period)
2 {
3     // the sample buffer
4     double * buffer = new double [period + 1];
5
6     // time for one sample
7     double sampleTime = 1.0 / period;
8
9     // the midpoint of the period
10    double mid = 0.5;
11
12    double value = 0;
13
14    // fill the sample buffer
15    for (int n = 0; n < period; n++)
16    {
17        buffer[n] = (value < mid) ? -1 : 1;
18
19        value += sampleTime;
20    }
21
22    return buffer;
23 }

```

Table 2.2: C++ code to generate and return one period of a square wave, where `period` is the period duration in samples. Note that this function increments in sample time, measured in seconds, rather than actual samples. This prevents a one-sample quantization error at the mid-point, since time can always be halved whereas a sample is a fixed entity and cannot be broken down any further.

```

1 double* sawtooth(const unsigned int period)
2 {
3     // the sample buffer
4     double * buffer = new double [period];
5
6     // how much we must decrement the
7     // index by at each iteration
8     double incr = -2.0 / period;
9
10    double value = 1;
11
12    for (int n = 0; n < period; n++)
13    {
14        buffer[n] = value;
15
16        value += incr;
17    }
18
19    return buffer;
20 }

```

Table 2.3: C++ code to generate one period of a sawtooth wave function, where `period` is the period duration in samples.

Ramp Waves

A ramp wave is, quite simply, an inverted sawtooth wave. It ascends from its minimum amplitude to its maximum amplitude linearly, after which it jumps back down to the minimum amplitude. Consequently, Equation 2.7 and 2.8 differ from sawtooth Equation 2.5 and 2.6 only in their sign and offset. The equivalent C++ implementation shown in Table 2.4 also reflects these differences.

$$f(\phi) = \begin{cases} \frac{\phi}{\pi} - 1, & \text{if } 0 \leq \phi < 2\pi \end{cases} \quad (2.7) \quad f(t) = \begin{cases} \frac{2t}{T} - 1, & \text{if } 0 \leq t < 1 \end{cases} \quad (2.8)$$

```

1 double* ramp(const unsigned int period)
2 {
3     // the sample buffer
4     double * buffer = new double [period];
5
6     // index incrementor
7     double incr = 2.0 / period;
8
9     double value = -1;
10
11     for (int n = 0; n < period; n++)
12     {
13         buffer[n] = value;
14
15         value += incr;
16     }
17
18     return buffer;
19 }

```

Table 2.4: C++ ramp wave generator. This code differs from the program shown in Table 2.3 solely in the amplitude offset (-1 instead of 1) and the increment, which is now positive.

Triangle waves

A triangle wave can be seen as a combination of a ramp wave and a sawtooth wave, or as a linear, "edgy", sine wave. It increments from its minimum amplitude to its maximum amplitude linearly one half of a period and decrements back to the minimum during the other half. Simply put, "[a] triangle wave is a linear increment or decrement that switches direction every π radians" (Mitchell, 2008, p. 69). A mathematical definition for one period of a triangle wave is given by Equation 2.9, where ϕ is the phase in radians. If ϕ is kept in the range of $[-\pi; \pi]$ rather than the usual range of $[0; 2\pi]$, the subtraction of π can be eliminated, yielding Equation 2.10. If the dependent variable is time, in seconds, or samples, Equation 2.11 can be used for a range of $[0; T]$, where T is the period, and Equation 2.12 for a range of $[-\frac{T}{2}; \frac{T}{2}]$. A C++ implementation is shown in Table 2.5.

$$f(\phi) = \begin{cases} 1 - \frac{2|\phi - \pi|}{\pi}, & \text{if } 0 \leq \phi < 2\pi \end{cases} \quad (2.9)$$

$$f(\phi) = \begin{cases} 1 - \frac{2|\phi|}{\pi}, & \text{if } -\pi \leq \phi < \pi \end{cases} \quad (2.10)$$

$$f(t) = \begin{cases} 1 - \frac{4|t - \frac{T}{2}|}{T}, & \text{if } 0 \leq t < T \end{cases} \quad (2.11)$$

$$f(t) = \begin{cases} 1 - \frac{4|t|}{T}, & \text{if } -\frac{T}{2} \leq t < \frac{T}{2} \end{cases} \quad (2.12)$$

2.2.2 Additive Synthesis

The second method of generating complex waveforms, additive synthesis, produces waveforms that, despite not being mathematically perfect, are closer to the waveforms found naturally. This method involves the summation of a theoretically infinite, practically finite set of sine and/or cosine waves with varying parameters and is often called Fourier Synthesis, after the 18th century French scientist, Joseph Fourier, who first described the process and associated phenomena of summing sine and cosine waves to produce complex waveforms. This calculation of a complex, periodic waveform from a sum of sine and cosine functions is also called a

```

1 double* triangle(const unsigned int period) const
2 {
3     double* buffer = new double[period];
4
5     double value = -1;
6
7     // 4.0 because we're incrementing/decrementing
8     // half the period and the range is 2, so it's
9     // actually 2 / period / 2.
10    double incr = 4.0 / period;
11
12    // Boolean to indicate direction
13    bool reachedMid = false;
14
15    for (unsigned int n = 0; n < period; n++)
16    {
17        wt[n] = value;
18
19        // Increment or decrement depending
20        // on the current direction
21        value += (reachedMid) ? -incr : incr;
22
23        // Change direction every time
24        // the value hits a maximum
25        if (value >= 1 || value <= -1)
26        { reachedMid = !reachedMid; }
27    }
28
29    return buffer;
30 }

```

Table 2.5: C++ program to compute one period of a triangle wave.

Fourier Transform or a Fourier Series, both part of the Fourier Theorem. In a Fourier Series, a single sine/cosine component is either called a harmonic, an overtone or a partial. All three name the same idea of a waveform with a frequency that is an *integer multiple* of some fundamental pitch or frequency. (Mitchell, 2008, p. 64) Throughout this thesis the term *partial* will be preferred.

Equation 2.13 gives the general definition of a discrete Fourier Transform. Equation 2.14 shows a simplified version of Equation 2.13. Table 2.6 presents a C++ struct to represent a single partial and Tables 2.7 and 2.8 a piece of C++ code to compute one period of any Fourier Series.

$$f(t) = \frac{a_0}{2} + \sum_{n=1}^{\infty} (a_n \cos(\omega n t) + b_n \sin(\omega n t))$$

Equation 2.13: Formula to calculate an infinite Fourier series, where $\frac{a_n}{2}$ is the center amplitude, a_n and b_n the partial amplitudes and ω the angular frequency, which is equal to $2\pi f$.

$$f(t) = \sum_{n=1}^N a_n \sin(\omega n t + \phi_n)$$

Equation 2.14: Simplification of Equation 2.13. Note the change from a computationally impossible infinite series to a more practical finite series. Because a cosine wave is a sine wave shifted by 90° or $\frac{\pi}{2}$ radians, the \cos function can be eliminated and replaced by an appropriate \sin function with a phase shift ϕ_n .

```

1 struct Partial
2 {
3     Partial(unsigned short number, double ampl, double phsOffs = 0)
4         : num(number), amp(ampl), phaseOffs(phsOffs)
5     { }
6
7     /*! The Partial's number, stays const. */
8     const unsigned short num;
9
10    /*! The amplitude value. */
11    double amp;
12
13    /*! A phase offset */
14    double phaseOffs;
15 };

```

Table 2.6: C++ code to represent a single partial in a Fourier Series.

```

1 template <class PartItr>
2 double* additive(PartItr start,
3                 PartItr end,
4                 unsigned long length,
5                 double masterAmp = 1,
6                 bool sigmaAprox = false,
7                 unsigned int bitWidth = 16)
8 {
9     static const double pi = 3.141592653589793;
10
11    static const double twoPi = 6.28318530717958;
12
13    // calculate number of partials
14    unsigned long partNum = end - start;
15
16    double * buffer = new double [length];
17
18    double * amp = new double [partNum]; // the amplitudes
19    double * phase = new double [partNum]; // the current phase
20    double * phaseIncr = new double [partNum]; // the phase increments
21
22    // constant sigma constant part
23    double sigmaK = pi / partNum;
24
25    // variable part
26    double sigmaV;
27
28    // convert the bit number to decimal
29    bitWidth = pow(2, bitWidth);
30
31    // the fundamental increment of one period
32    // in radians
33    static double fundIncr = twoPi / length;
34
35    // fill the arrays with the respective partial values
36    for (unsigned long p = 0; start != end; ++p, ++start)
37    {
38        // initial phase
39        phase[p] = start->phaseOffs;
40
41        // fundIncr is two Ï€ / tablelength
42        phaseIncr[p] = fundIncr * start->num;
43
44        // reduce amplitude if necessary
45        amp[p] = start->amp * masterAmp;
46
47        // apply sigma approximation conditionally
48        if (sigmaAprox)
49        {
50            // following the formula
51            sigmaV = sigmaK * start->num;
52
53            amp[p] *= sin(sigmaV) / sigmaV;
54        }
55    }

```

Table 2.7

```

56 // fill the wavetable
57 for (unsigned int n = 0; n < length; n++)
58 {
59     double value = 0;
60
61     // do additive magic
62     for (unsigned short p = 0; p < partNum; p++)
63     {
64         value += sin(phase[p]) * amp[p];
65
66         phase[p] += phaseIncr[p];
67
68         if (phase[p] >= twoPi)
69         { phase[p] -= twoPi; }
70
71         // round if necessary
72         if (bitWidth < 65536)
73         {
74             Util::round(value, bitWidth);
75         }
76
77         buffer[n] = value;
78     }
79
80     delete [] phase;
81     delete [] phaseIncr;
82     delete [] amp;
83
84     return buffer;
85 }

```

Table 2.8: *Continuation of Table 2.7.* C++ program to produce one period of an additively synthesized complex waveform, given a start and end iterator to a container of partials, a buffer length, a maximum, "master", amplitude, a boolean whether or not to apply sigma approximation and lastly a maximum bit width parameter.

The following paragraphs will examine how the four waveforms presented in Section 2.2.1, the square, the sawtooth, the ramp and the triangle wave, can be synthesized additively.

Square Waves

When speaking of additive synthesis, a square wave is the result of summing all odd-numbered partials (3rd, 5th, 7th etc.) at a respective amplitude equal to the reciprocal of their partial number ($\frac{1}{3}$, $\frac{1}{5}$, $\frac{1}{7}$ etc.). The amplitude of each partial must decrease with increasing partial numbers to prevent amplitude overflow. A mathematical equation for such a square wave with N partials is given by Equation 2.13, where $2n - 1$ makes the series use only odd partials. A good maximum number of partials N for near-perfect but still naturally sounding waveforms is 64, a value determined empirically. Higher numbers have not been found to produce significant improvements in sound quality. Table 2.9 displays the C++ code needed to produce one period of a square wave in conjunction with the `additive` function from Tables 2.7 and 2.8. Figure 2.1 shows the result of summing 2, 4, 8, 16, 32 and finally 64 partials.

$$f(t) = \sum_{n=1}^N \frac{1}{2n-1} \sin(\omega(2n-1)t) \quad (2.13)$$

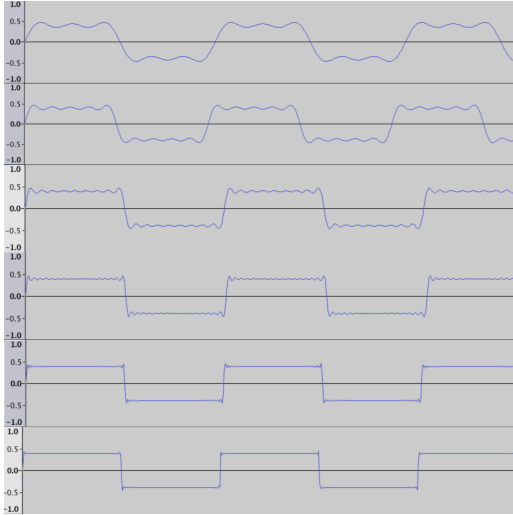


Figure 2.1: Square waves with 2, 4, 8, 16, 32 and 64 partials.

Table 2.9: C++ code for a square wave with 64 partials.

```

1 | std::vector<Partial> vec;
2 |
3 | for (int i = 1; i <= 128; i += 2)
4 | {
5 |     vec.push_back(Partial(i, 1.0/i));
6 | }
7 |
8 | double* buffer = additive(vec.begin(),
9 |                           vec.end(),
10 |                          48000)

```

Sawtooth Waves

A sawtooth wave is slightly simpler to create through additive synthesis, as it requires the summation of every partial rather than only the odd-numbered ones. The respective amplitude is again the reciprocal of the partial number. Equation 2.14 gives a mathematical definition for a sawtooth wave, Figure 2.2 displays sawtooth functions with various partial numbers and Table 2.10 shows C++ code to generate such functions.

$$f(t) = \sum_{n=1}^N \frac{1}{n} \sin(\omega n t) \quad (2.14)$$

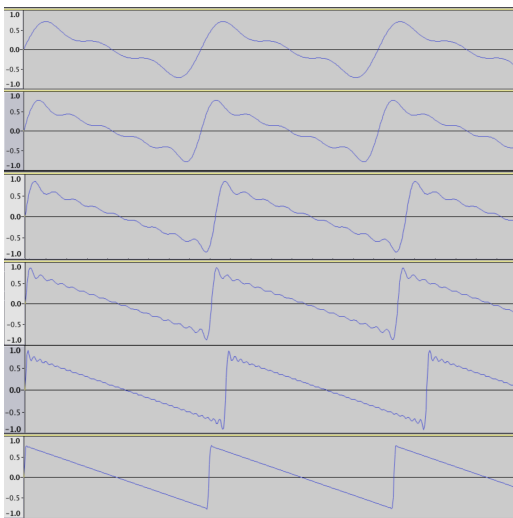


Figure 2.2: Sawtooth waves with 2, 4, 8, 16, 32 and 64 partials.

Table 2.10: C++ code for a sawtooth wave with 64 partials.

```

1 | std::vector<Partial> vec;
2 |
3 | for (int i = 1; i < 64; ++i)
4 | {
5 |     vec.push_back(Partial(i, 1.0/i));
6 | }
7 |
8 | double* buffer = additive(vec.begin(),
9 |                           vec.end(),
10 |                          48000)

```


Ramp Waves

Ramp waves are essentially the inverse of sawtooth waves. Therefore, we can simply sum sinusoids as we did for a sawtooth function but change the sign of each partial's amplitude to negative instead of positive. Equation 2.15 gives the mathematical definition, Figure 2.3 displays a set of ramp waveforms with different partial numbers and Table 2.11 again shows the accompanying C++ code.

$$f(t) = \sum_{n=1}^N -\frac{1}{n} \sin(\omega n t) \quad (2.15)$$

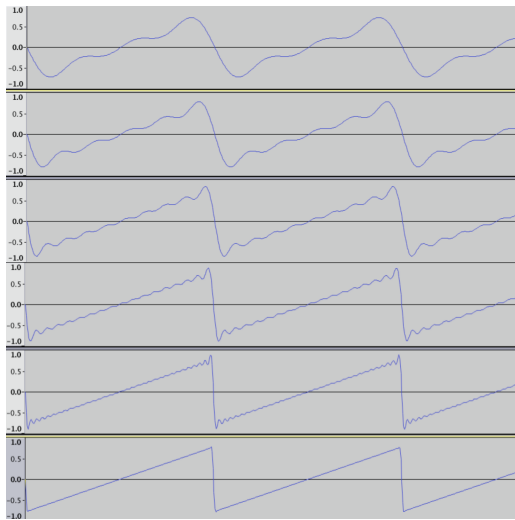


Figure 2.3: Ramp waves with 2, 4, 8, 16, 32 and 64 partials.

Table 2.11: C++ code for a ramp wave with 64 partials.

```

1  std::vector<Partial> vec;
2
3  for (int i = 1; i < 64; ++i)
4  {
5      vec.push_back(Partial(i, -1.0/i));
6  }
7
8  double* buffer = additive(vec.begin(),
9                             vec.end(),
10                            48000)

```

Triangle Waves

The process of generating triangle waves additively differs from previous waveforms. The amplitude of each partial is no longer the reciprocal of the partial number, $\frac{1}{n}$, but now the inverse of the partial number squared: $\frac{1}{n^2}$. Moreover, the sign of the amplitude alternates for each partial in the series. As for square waves, only odd-numbered partials are used. Mathematically, such a triangle wave is defined as shown in Equation 2.16 or, more concisely, in Equation 2.17. Figure 2.4 displays such a triangle wave with various partial numbers (2,4,8,16,32 and 64) and Table 2.12 implements C++ code to compute a triangle wave.

$$f(t) = \sum_{n=1}^{\frac{N}{2}} \frac{1}{(4n-3)^2} \sin(\omega(4n-3)t) - \frac{1}{(4n-1)^2} \sin(\omega(4n-1)t) \quad (2.16)$$

$$f(t) = \sum_{n=0}^N \frac{(-1)^n}{(2n+1)^2} \sin(\omega(2n+1)t) \quad (2.17)$$

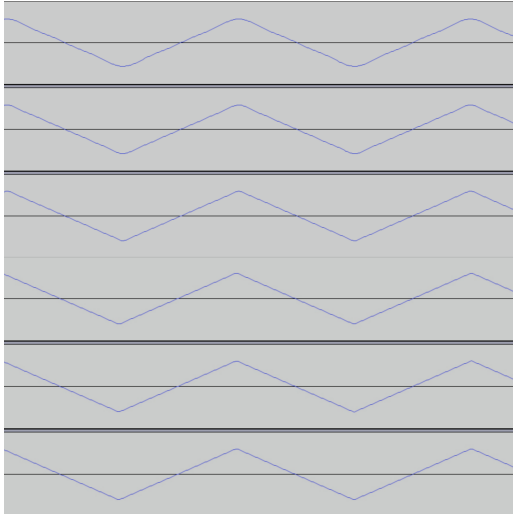


Figure 2.4: Triangle waves with 2, 4, 8, 16, 32 and 64 partials. Note that already 2 partials produce a very good approximation of a triangle wave.

Table 2.12: C++ code for a triangle wave with 64 partials.

```

1  std::vector<Partial> vec;
2
3  double amp = -1;
4
5  for(int i = 1; i <= 128; i += 2)
6  {
7      amp = (amp > 0) ? (-1.0/(i*i)) : (1.0/(i*i));
8
9      vec.push_back(Partial(i,amp));
10 }
11
12 double* buffer = additive(vec.begin(),
13                             vec.end(),
14                             48000);

```

2.2.3 Special Waveforms

One noticeable problem with square, sawtooth and ramp waves is that they have very sharp transitions at certain points in their waveforms, such as the square wave at the midpoint, $\frac{T}{2}$, where the amplitude jumps from its maximum to its minimum. While these transitions and jumps in amplitude contribute to the characteristic sound of these waves when used for music, they can pose a problem when they are used for modulation, as sharp transitions in amplitude translate to clicking sounds acoustically. One possible solution to reduce the impact of such transitions would be to overlay an audio envelope directly on the waveform, to bend or smooth the amplitude at the transition points. However, a more effective and direct solution would be to mathematically calculate a waveform that is already smoothed out at these critical points. The following two paragraphs propose a set of empirically determined functions to calculate "smooth" square and sawtooth waves.

Smooth Square Waves

Smooth square waves can be calculated mathematically as a set of four power functions with a very large exponent that change their sign after the midpoint. The exponent chosen, 50, was determined entirely experimentally. Equations 2.18 to 2.21 display the four power functions and their respective range of definition. Figure 2.5 shows the resulting visualization, where each individual power function is highlighted in a different colour as a visual aid.

$$a(t) = \begin{cases} (t-1)^{50} - 1, & \text{if } 0 \leq t < \frac{T}{4} \end{cases} \quad (2.18)$$

$$b(t) = \begin{cases} (t+0.5)^{50} - 1, & \text{if } \frac{T}{4} \leq t < \frac{T}{2} \end{cases} \quad (2.19)$$

$$c(t) = \begin{cases} -(t-1.5)^{50} + 1, & \text{if } \frac{T}{2} \leq t < \frac{3T}{4} \end{cases} \quad (2.20)$$

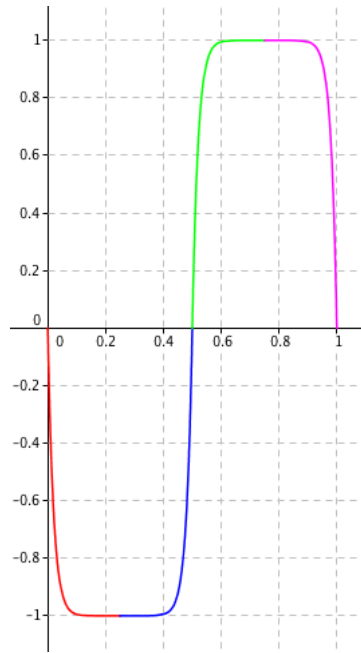


Figure 2.5: One period of a smooth square wave as the sum of four piecewise power functions.

$$d(t) = \begin{cases} -t^{50} + 1, & \text{if } \frac{3T}{4} \leq t < T \end{cases} \quad (2.21)$$

Smooth Sawtooth Waves

For sawtooth waves, the only critical area of transition is the point at which one period ends and the next begins, as the amplitude jumps back to the maximum from its minimum value. To make this transition smoother, one can use two very steep quadratic functions, one positive and one negative, that intersect at the midpoint of the transition. If one allocates nine tenths of one period for the normal descent of the sawtooth function from its maximum to its minimum, the transition can take place in the last tenth of the period. Equation 2.22 shows an altered version of the sawtooth wave function introduced in Section 2.2.1 and Equations 2.23 and 2.24 the two aforementioned quadratic functions. Figure 2.6 visualizes the resulting waveform.

$$a(t) = \begin{cases} \frac{-2t}{0.9T} + 1 = \frac{-20t}{9T} + 1, & \text{if } 0 \leq t < \frac{9T}{10} \end{cases} \quad (2.22)$$

$$b(t) = \begin{cases} 400(t - 0.9)^2 - 1, & \text{if } \frac{9T}{10} \leq t < \frac{95T}{100} \end{cases} \quad (2.23)$$

$$c(t) = \begin{cases} -400(t - 1)^2 + 1, & \text{if } \frac{95T}{100} \leq t < T \end{cases} \quad (2.24)$$

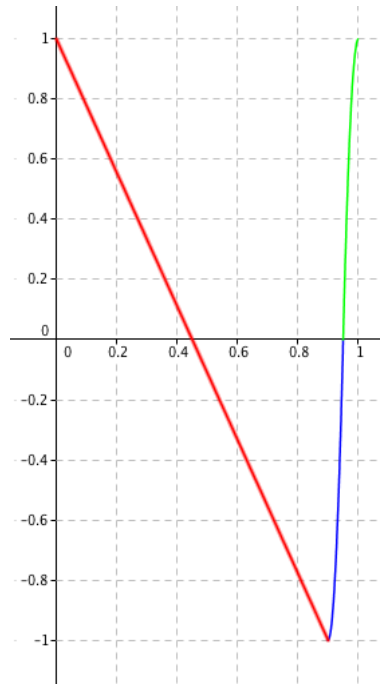


Figure 2.6: One period of a smooth sawtooth function as a sum of three piecewise functions.

2.3 The Gibbs Phenomenon and the Lanczos Sigma Factor

Examining Figure 2.7, which displays an additively synthesized square wave function with 64 partials, one may observe that additive synthesis, the summation of sine waves to produce complex waveforms, produces an overshoot — slight ripples or "horns" — at the ends of each peak and trough of a waveform. This is known as the Gibbs Phenomenon, named after the American scientist Josiah Willard Gibbs who first described it in 1898, and is "the result of summing a finite series of partials rather than the infinite series of partials specified by the Fourier transform" (Mitchell, 2008, p. 67). Acoustically, this ringing artifact has little influence on the sound of the produced waveform. However, for modulation, this imperfection may render the resulting sound unsatisfactory. A common way to reduce the Gibbs Phenomenon is to apply the Lanczos Sigma (σ) Factor to each partial of a Fourier Series. This is often called *sigma-approximation*. The definition of the Lanczos Sigma Factor is given in Equation 2.25, where n is the current partial number and M the total number of partials in a Fourier Summation. In the `additive` function shown in Tables 2.7 and 2.8, the Lanczos Sigma Factor is implemented in lines 48 to 54. Figure 2.8 shows the same square wave from Figure 2.7 after sigma-approximation.

$$\sigma = \text{sinc}\left(\frac{n\pi}{M}\right) = \frac{\sin\left(\frac{n\pi}{M}\right)}{\frac{n\pi}{M}} \quad (2.25)$$

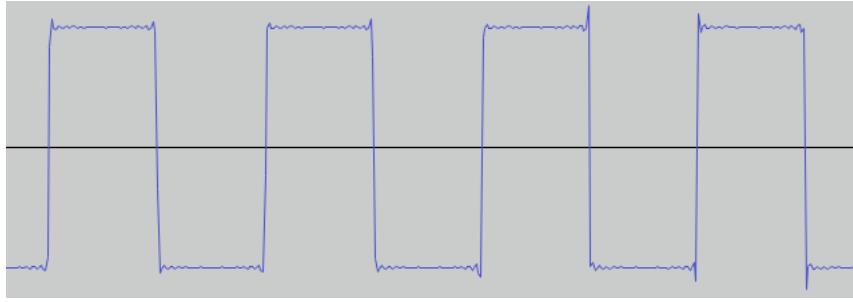


Figure 2.7: An additively synthesized square wave with 64 partials before sigma-approximation.

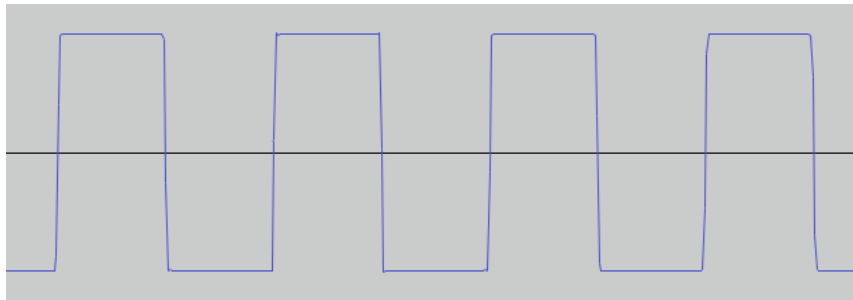


Figure 2.8: An additively synthesized square wave with 64 partials after sigma-approximation.

2.4 Real-Time Calculation versus Wavetable Lookup

Following the discussion of the creation of complex waveform, the two options for playing back waveforms in a digital synthesis system must be examined: continuous real-time calculation of waveform samples or lookup from a table that has been calculated once and then written to disk — a so-called Wavetable. To keep matters short, the second method was found to be computationally more efficient and thus the better choice, as memory is a more easily expended resource than computational speed.

2.4.1 Implementation

A Wavetable is a table (practically speaking an array) in which one stores pre-calculated waveform amplitude values for lookup. This way, the computation of individual samples as shown in Sections 2.2.1 and 2.2.2 can be replaced by the incrementing of a table index variable and retrieval of samples by dereferencing of the array at the current index. If a table holds sample values for one period of a waveform, at a frequency of $1Hz$, the frequency of the waveform can be adjusted during playback by multiplying the increment value by some factor other than one. "For example, if we increment by two [instead of one], we scan the table in half the time and produce a frequency twice the original frequency. Other increment values will yield proportionate frequencies." (Mitchell, 2008, p. 80-81) The fundamental increment for a table index, i_{fund} , the value by which a table index must be incremented after each sample to traverse a waveform at a frequency of $1Hz$, given the table length L and the samplerate f_s , is shown in Equation 2.26. To alter the frequency f of the played-back waveform, Equation 2.27 should be used to calculate the appropriate increment value i .

$$i_{fund} = \frac{L}{f_s} \quad (2.26)$$

```
1 | sample = table[integral] + ((table[integral + 1] - table[integral]) * fractional)
```

Table 2.13: An interpolation algorithm in pseudo-code.

```
1 | template<class T>
2 | double interpolate(T table [], double index)
3 | {
4 |     long integral = static_cast<long>(index); // The truncated integral part
5 |     double fractional = index - integral; // The remaining fractional part
6 |
7 |     // grab the two items in-between which the actual value lies
8 |     T value1 = table[integral];
9 |     T value2 = table[integral+1];
10 |
11 |     // interpolate: integer part + (fractional part * difference between value2 and value1)
12 |     double final = value1 + ((value2 - value1) * fractional);
13 |
14 |     return final;
15 | }
```

Table 2.14: Full C++ template function to interpolate a value from a table, given a fractional index.

$$i = i_{fund} \times f = \frac{L}{f_s} \times f \quad (2.27)$$

2.4.2 Interpolation

For almost all configurations of frequencies, table lengths and sample rates, the table index i produced by Equation 2.27 will not be an integer. Since using a floating point number as an index for an array is a syntactically illegal operation in C++, there are two options. The first is to truncate or round the table index to an integer, thus "introducing a quantization error into the signal [...]". This is obviously a suboptimal solution which would result in a change of phase and consequently distortion. (Mitchell, 2008, p. 84) The second option, interpolation, tries to approximate the true value from the sample at the current index and at the subsequent one. Interpolation is achieved by summation of the sample value at the floored, integral part of the table index, $\lfloor i \rfloor$, with the difference between this sample and the sample value at the next table index, $\lfloor i \rfloor + 1$, multiplied by the fractional part of the table index, $i - \lfloor i \rfloor$. Table 2.13 displays the calculation of a sample by means of interpolation in pseudo-code and Table 2.14 in C++. (based on pseudo-code, Mitchell, 2008, p. 85)

2.4.3 Table Length

The length of the Wavetable must be chosen carefully to be both memory efficient and at the same time provide a decently accurate waveform resolution. An equation to calculate the size of a single wavetable in Kilobytes is given by Equation 2.28, where L is the table length and N the number of bytes provided by the precision of the data type used for samples, e.g. 8 bytes for the double precision floating point data type `double`.

$$Size = \frac{L \times N}{1024} \quad (2.28)$$

From: Jari Kleimola jari.kleimola@aalto.fi
 Subject: RE: Question about license
 Date: 10 Mar 2014 23:30 PM
 To: Peter Goldsborough petergoldsbrough@hotmail.com

Hi Peter,

You need to handle WT[i+1] case carefully in order not to go out of bounds. One solution is to append the first entry of the WT to the end of the WT. Note that direct computation is sometimes faster than wavetables (especially if WT is too big to fit inside processor cache). As a rule of thumb, do not use a bigger wavetable than the cache.

Figure 2.9: An excerpt of an E-Mail exchange with Jari Kleimola.

Daniel R. Mitchell advises that the length of the table be a power of two for maximum efficiency. (Mitchell, 2008, p. 86) Moreover, during an E-Mail exchange with Jari Kleimola, author of the 2005 master's thesis titled "Design and Implementation of a Software Sound Synthesizer", it was discovered that "as a rule of thumb", the table length should not exceed the processor cache size. A relevant excerpt of this e-mail exchange is depicted in Figure 2.9. Considering both pieces of advice as well as the fact that my computer, a 2.5 GHz Intel Core i5, 13-inch, Mid-2012 MacBook Pro, has an L2-Cache size of 256 KB, it was decided that a Wavetable size of 4096, which translates to 32 KB of memory, would be suitable.

One important fact to mention, also discussed by Jari Kleimola as seen in Figure 2.9, is that because the interpolation algorithm from Table 2.14 must have access to the sample value at index $i + 1$, i being the current index, an additional sample must be appended to the Wavetable to avoid a BAD_ACCESS error. This added sample usually has the same value as the first sample in the table to avoid a discontinuity. Therefore, the period of the waveform actually only fills 4095 of the 4096 indices of the Wavetable, as the 4096th sample is equal to the first.

2.4.4 File Format

For maximum efficiency, the Wavetables are not created at program startup but read from a binary file. To prevent the synthesizer program from accidentally reading faulty files, some simple identification string must be added to the file. Therefore, Wavetable files contain a 6-char ID string with the name of the synthesizer, Anthem, after which the 32 KB of Wavetable data follow. Additionally, Wavetable files end with a .wavetable file extension. Table 2.15 displays a function to read such a Wavetable file and Figure 2.10 shows the first few bytes of a Wavetable file when opened as plain-text. The total size of a Wavetable file is exactly 32774 bytes, 32768 bytes (32KB) of Wavetable data and 6 bytes for the ID string.

```
ANTHEM{tÁâ#Y?ô`)~ä#i?`æ Ä•/r?·z05É#y?÷?Ë°\l?À°≠μð/Ç?ìμ7$~Ö?±@^-d#â?ôÄÖVfGâ?~E" lè?%~≠ä;Hë?
BtoYd/1?@iâ?äli?áq#." .i?ÄððYÄêó?g"wË"ô?ÎYÄË~Yô?êÄ†Gú?cñÜ°"Yù?@VK°-kü?Æä'-ð~†?μk:ôG°?¶4μ~ó
ç?â5ñÎiYç?b,Ç~äçÉ?
```

Figure 2.10: The first few bytes of a Wavetable file.

```
1  #include <fstream>
2  #include <stdexcept>
3
4  double* readWavetable(const std::string &fname)
5  {
6      std::ifstream file(fname);
7
8      if (! file.is_open())
9      { throw std::runtime_error("Could not find wavetable file: " + fname); }
10
11     if (! file.good())
12     { throw std::runtime_error("Error opening wavetable: " + fname); }
13
14     char signature[6];
15
16     file.read(signature, 6);
17
18     if (strncmp(signature, "ANTHEM", 6))
19     { throw std::runtime_error("Invalid signature for Anthem file!"); }
20
21     int len = 4096;
22     int size = len * sizeof(double);
23
24     double * wt = new double [len];
25
26     file.read(reinterpret_cast<char*>(wt), size);
27
28     return wt;
29 }
```

Table 2.15: C++ code to read a Wavetable file.

Bibliography

- [1] Daniel R. Mitchell,
BasicSynth: Creating a Music Synthesizer in Software.
Publisher: Author.
1st Edition,
2008.
- [2] Steven W. Smith,
The Scientist and Engineer's Guide to Digital Signal Processing.
California Technical Publishing,
San Diego, California,
2nd Edition,
1999.
- [3] Gordon Reid,
Synth Secrets, Part 12: An Introduction To Frequency Modulation.
<http://www.soundonsound.com/sos/apr00/articles/synthsecrets.htm>
Accessed: 8 October 2014.
- [4] Justin Colletti,
The Science of Sample Rates (When Higher Is Better – And When It Isn't)
2013.
<http://www.trustmeimascientist.com/2013/02/04/the-science-of-sample-rates-when-higher-is-better-and-when-it-isnt/>
Accessed: 17 December 2014.
- [5] John D. Cutnell and Kenneth W. Johnson,
Physics.
Wiley,
New York,
4th Edition,
1998.

List of Figures

1.1	The continuous representation of a typical sine wave. In this case, both the signal's frequency f as well as the maximum elongation from the equilibrium a are equal to 1.	3
1.2	The discrete representation of a typical sine wave.	4
1.3	A sinusoid with a frequency of 4 Hz.	5
1.4	A sinusoid with a frequency of 4 Hz, sampled at a sample rate of 5 Hz. According to the Nyquist Theorem, this signal is sampled improperly, as the frequency of the continuous signal is not less than or equal to one half of the sample rate, in this case 2.5 Hz.	6
1.5	An approximated representation of the signal created from the sampling process shown in Figure 1.4. Visually as well as acoustically, the new signal is completely different from the original signal. However, in terms of sampling, they are indistinguishable from each other due to improper sampling. These signals are said to be <i>aliases</i> of each other.	6
1.6	A typical sinusoidal signal with an amplitude of 1. The integer range provided by the allocated memory for the top-most sample is saturated, meaning it is equal to 32767_{10} or 011111111111111_2	7
1.7	What happens when the amplitude of the signal from Figure 1.6 is increased by a factor of 2. Multiple samples have overflowed and thus changed their sign. Because of the way two's-complement representation is implemented, the signal continues it's path as if no overflow had ever occurred. The only difference being, of course, that the sign has changed mid-way.	7
2.1	Square waves with 2, 4, 8, 16, 32 and 64 partials.	16
2.2	Sawtooth waves with 2, 4, 8, 16, 32 and 64 partials.	16
2.3	Ramp waves with 2, 4, 8, 16, 32 and 64 partials.	17
2.4	Triangle waves with 2, 4, 8, 16, 32 and 64 partials. Note that already 2 partials produce a very good approximation of a triangle wave.	18
2.5	One period of a smooth square wave as the sum of four piecewise power functions.	19
2.6	One period of a smooth sawtooth function as a sum of three piecewise functions.	20
2.7	An additively synthesized square wave with 64 partials before sigma-approximation.	21
2.8	An additively synthesized square wave with 64 partials after sigma-approximation.	21
2.9	An excerpt of an E-Mail exchange with Jari Kleimola.	23
2.10	The first few bytes of a Wavetable file.	23