

# How to Program a Calculator

## Numerical analysis of common functions

Jake Darby

### **Abstract**

This document will discuss and analyse various numerical methods for computing functions commonly found on calculators. The aim of this paper is to, for each set of functions, compare and contrast several algorithms in regards to their efficiency and accuracy.

# 1 Introduction

For many thousands of years all calculations that a Human might want performing had to be done by hand. For simple calculations such as addition, subtraction and multiplication this was not such an issue, but as society evolved we wanted to know the answer to increasingly hard questions. The greeks sought to find a value for  $\pi$ , and ended up with the bounds that  $\frac{223}{71} < \pi < \frac{22}{7}$ , which while sufficient for their needs is not sufficient for ours in the modern era.

Nowadays we have computers and calculators to calculate functions such as square roots or the trigonometric functions. These devices however must be told how to correctly calculate these values, those methods used will be looked at here. Specifically, for most purposes it does not matter if some small error is present in the results, as long as the value found is close enough to the actual value.

## 1.1 Code and Computers used

During this project I will be discussing the implementation of various algorithms. I will be implementing these algorithms in the C programming language, using the C11 standard.

I chose the C programming language to implement my algorithms in because, once it compiles to binary machine code, the programs produced tend to be very efficient. This is partly due to the low-level of C programming, having relatively close control over direct CPU actions; however this does come at the cost of losing higher functionality that many other programming languages offer. A second reason for the efficiency is due to C's long history, originally being developed in 1969-1970, which has led to several very efficient compilers.

I will be implementing most programs using C's built in primitive types, typically `int`, `unsigned int` and `double`. On a computer an `int` is an integer that can represent both positive and negative bits using twos complement, this gives an `int` using  $n$  bits a minimum value of  $-2^n$  and a maximum value of  $2^n - 1$ . Typically a computer will store an `int` as 32 bits, though some computers may use more or less bits. An `unsigned int` is very similar to an `int`, but does not represent negative values, and thus an `unsigned int` of  $n$  bits has a minimum value of 0 and a maximum value of  $2^n - 1$ .

If an integer of a specific number of bits is needed then the header `stdint.h` may be used which defines `int_N` and `uint_N` which respectively represent `int` of  $N$  bits and `unsigned int` of  $N$  bits; The typical values of  $N$  are 8, 16, 32 and 64.

In C a `double` is a floating point representation of a real value, that typically follows the IEEE 754 standard for double-precision binary floating points. This standard has:

- 1 bit for the sign of the number,  $s$
- 11 bits for the exponent,  $e$
- 52 bits for the significand,  $b = b_0b_1b_2 \dots b_{51}$

- A value that is understood to be:

$$(-1)^s \left( 1 + \sum_{i=1}^{52} b_{52-i} 2^{-i} \right) \times 2^{e-1023}$$

This gives a `double` value a precision of around 15-17 significant decimal digits. While this is good for most applications, there are applications where we may want even more precision than this. To solve this I will be implementing certain algorithms using the GNU Multiple Precision Arithmetic Library (referred to as GMP) as well as GNU MPFR Library(referred to as MPFR), which was built upon GMP to correct and optimise the original. These libraries allow the use of arbitrary precision real values, given enough memory space, as well as integers longer than C's standard integer types can hold.

The downside is that the increased precision does increase computation time for these calculations.

I will be compiling and testing all of my code on a benchmark machine running a light version of Ubuntu `16.04`, using the GNU C Compiler. The specifications of the machine, that may impact performance are:

- An Intel i5-4690K processor running at 4GHz. This processor uses a 64 bit architecture.
- 8Gb of DDR3 RAM
- A modern chipset on the motherboard

**TODO: Expand the introduction and make it more eloquent**

## 2 General Definitions and Theorems

This section will discuss those definitions and theorems that will be applicable and referenced later in the document.

### 2.1 Methods

In this document we will look at various functions, such as root functions, trigonometric functions, among others. Despite the variety of functions being analysed there are several methods that are useful for more than one function, we will discuss these methods below.

#### 2.1.1 Newton-Raphson Method

The Newton-Raphson Method is named after Sir Isaac Newton and Joseph Raphson. It is a method that takes a continuously differentiable function  $f$  and it's derivative  $f'$ , as well as an initial guess  $x_0$ , to create successively more accurate solutions to  $x$  where  $f(x) = 0$ .

The specific definition of the Newton-Raphson method that I will be using in this document is below:

**Definition 2.1.1.1.** Given  $f \in \mathcal{C}(\mathbb{R})$ ,  $f'$  being the derivative of  $f$ , and  $x_0 \in \mathbb{R}$ ; then we define:

$$x_{n+1} := x_n - \frac{f(x_n)}{f'(x_n)} \forall n \in \mathbb{N}$$

The hope is that, if we start with a good initial guess, that the method will converge to some  $x \in \mathbb{R} : f(x) = 0$ .

We derive the above method in the following way. If we have an approximation  $x_n \in \mathbb{R} : f(x_n) \approx 0$  we consider the tangent to  $f(x)$  at  $x_n$ , which is given by  $y = f'(x_n)(x - x_n) + f(x_n)$ . If we set  $y = 0$  and solve for  $x$  we get that

$$x = x_n - \frac{f(x_n)}{f'(x_n)}$$

The iterative formula follows by letting  $x_{n+1} := x$  in the above.

TODO: Re-arrange into a better order

TODO: References

TODO: Possibly include discussions of pros and cons

#### 2.1.2 Taylor Series Expansion

The Taylor Series formulation was created by Brook Taylor in 1715, based off of the work of Scottish mathematician James Gregory. The Taylor Series describes a method of representing a given function by a polynomial, where any finite number of terms will give an approximation to the function itself.

**Definition 2.1.2.1.** Given  $f : \mathbb{R} \rightarrow \mathbb{R}$  which is infinitely differentiable at  $a \in \mathbb{R}$ , we define the Taylor Series of  $f$  at  $a$  to be:

$$\sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x - a)^n$$

We can then, from our definition of a Taylor Series, define a polynomial that will approximate our function  $f$  at  $x \in \mathcal{I} \subset \mathbb{R}$

**Definition 2.1.2.2.** Given  $f : \mathbb{R} \rightarrow \mathbb{R}$  which has a Taylor Series of  $\sum_{n=0}^{\infty} c_n x^n$ , we define the Taylor Polynomial of degree  $N \in \mathbb{N}$  to be

$$p_N(x) := \sum_{n=0}^N c_n x^n = c_0 + c_1 x + c_2 x^2 + \cdots + c_N x^N$$

Some examples of simple Taylor Series are:

$$\frac{1}{1-x} = \sum_{n=0}^{\infty} x^n \quad \forall x \in (-1, 1)$$

TODO: Eloquate this section

TODO: Discuss when a Taylor Series is not equal to it's function

TODO: Describe intervals for which Taylor Series converge

TODO: Create more examples of Taylor Series that are not used later in the document

TODO: Possibly prove Taylor's Theorem

### 2.1.3 CORDIC

TODO: Describe the CORDIC Algorithm in general

TODO: Let other sections define specific changes to CORDIC

## 2.2 Errors

Errors are very useful in this document for discussing convergence and measuring the accuracy of a particular method. There are two distinct types of error that we are interested in:

**Definition 2.2.1.** If we have a value  $v$  and it's approximation  $\tilde{v}$ , then the absolute error is

$$\epsilon := |v - \tilde{v}|$$

The absolute error will be useful in guaranteeing a certain level of accuracy that a given implementation of a method will give. Uses of absolute error in the document will use  $\epsilon$  as their absolute error variable.

**Definition 2.2.2.** If we have a value  $v$  and it's approximation  $\tilde{v}$ , then the relative error is

$$\eta := \frac{\epsilon}{|v|} = \left| 1 - \frac{\tilde{v}}{v} \right|$$

The relative error will also be useful in guaranteeing a certain level of accuracy that a given implementation of a method will give, particularly when the absolute error varies with the size of the input. Uses of relative error in the document will use  $\eta$  as their absolute error variable.

As the previous two errors are hard or impossible to accurately estimate, in a practical manner, then we want an error estimate that we can calculate in an algorithm. Thus we define the iteration error, as the absolute difference of consecutive iterations.

**Definition 2.2.3.** If we have the sequence  $(x_n)_{n \in \mathbb{N}}$ , then the iteration error is defined as:

$$\delta_n := |x_n - x_{n-1}|$$

For our practical purposes we can note that it is almost impossible to calculate  $\epsilon_n$  exactly, while  $\delta_n$  is easy to calculate and tends to be a good enough approximation of  $\epsilon_n$ ; particularly if the convergence is rapid.

## 2.3 Convergence

definition

In this section we consider what it means for a sequence to converge to a limit value, and some results that will be useful in future chapters.

**Definition 2.3.1.** A sequence  $(x_n \in \mathbb{R} : n \in \mathbb{N})$  converges to  $x$  uniformly if

$$\forall \tau \in \mathbb{R}_0^+ \exists N \in \mathbb{N} : \epsilon_n := |x - x_n| < \tau \forall n \in [N, \infty) \cap \mathbb{Z}$$

*Remark 2.3.1.1.* We will typically use the notation that  $\lim_{n \rightarrow \infty} |x_n - x| = 0$ , to denote that  $(x_n : n \in \mathbb{N})$  converges to  $x$ .

**Theorem 2.3.1.**  $(x_n \in \mathbb{R} : n \in \mathbb{N})$  converges to  $x$  uniformly if and only if

$$\forall \tau \in \mathbb{R}_0^+ \exists N \in \mathbb{N} : |x_n - x_m| < \tau \forall m, n \in [N, \infty) \cap \mathbb{Z}$$

*Proof.* For  $\implies$  :

Suppose that  $(x_n : n \in \mathbb{N})$  converges to  $x$  uniformly. Then  $\forall \tau \in \mathbb{R}_0^+ \exists N \in \mathbb{N} : |x_n - x| < \tau \forall n \in [N, \infty) \cap \mathbb{Z}$ .

Thus suppose  $N \in \mathbb{N}$  is such that  $|x_n - x| < \frac{\tau}{2} \forall n \in [N, \infty) \cap \mathbb{Z}$ .

Then if  $n, m \geq N$  we see that

$$|x_n - x_m| \leq |x_n - x| + |x_m - x| \leq \tau$$

For  $\Leftarrow$ :

Omitted for brevity. □

**Definition 2.3.2.** If  $(x_n \in \mathbb{R} : n \in \mathbb{N})$  is a sequence that converges to  $x$ , then it is said to converge:

- Linearly if  $\lambda \in \mathbb{R}^+$  and

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - x|}{|x_n - x|} = \lambda$$

- Quadratically if  $\lambda \in \mathbb{R}^+$  and

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - x|}{|x_n - x|^2} = \lambda$$

- Order  $\alpha \in \mathbb{R}_0^+$  if  $\lambda \in \mathbb{R}^+$  and

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - x|}{|x_n - x|^\alpha} = \lambda$$

For the following proof we will have  $\epsilon_n := |x^* - x_n|$ .

**Theorem 2.3.2.** *Let  $f$  be a twice differentiable function,  $x^*$  be a solution to  $f(x) = 0$  and  $(x_n : n \in \mathbb{N})$  be a sequence produced by the Newton-Raphson Method from some initial point  $x_0$ . If the following are satisfied, then  $(x_n : n \in \mathbb{N}_0)$  converges quadratically to  $x^*$ :*

**NR<sub>1</sub>:**  $f'(x) \neq 0 \forall x \in I := [x^* - r, x^* + r]$ , where  $r \in [|x^* - x_0|, \infty)$

**NR<sub>2</sub>:**  $f''(x)$  is continuous  $\forall x \in I$

**NR<sub>3</sub>:**  $M|\epsilon_0| < 1$  where  $M := \sup \left\{ \left| \frac{f''(x)}{f'(x)} \right| : x \in I \right\}$

*Proof.* By Taylor's Theorem with Lagrange Remainders we have that

$$0 = f(x^*) = f(x_n) + (x^* - x_n)f'(x_n) + \frac{1}{2}(x^* - x_n)^2 f''(y_n)$$

where  $0 < |x^* - y_n| < |x^* - x_n|$ .

Then we get the following derivation:

$$\begin{aligned} f(x_n) + (x^* - x_n)f'(x_n) &= -\frac{1}{2}(x^* - x_n)^2 f''(y_n) \\ \implies \left( \frac{f(x_n)}{f'(x_n)} - x_n \right) + x^* &= -\frac{1}{2} \frac{f''(y_n)}{f'(x_n)} (x^* - x_n)^2 && \text{as NR}_3 \implies f'(x_n) \neq 0 \\ \implies x^* - x_{n+1} &= -\frac{1}{2} \frac{f''(y_n)}{f'(x_n)} (x^* - x_n)^2 \\ \implies \epsilon_{n+1} &= \frac{1}{2} \left| \frac{f''(y_n)}{f'(x_n)} \right| \epsilon_n^2 && \text{by taking absolute values} \end{aligned}$$

As NR<sub>2</sub> holds then  $M$  exists and is positive, and therefore we have:

$$\epsilon_n \leq M \epsilon_{n-1}^2 \leq M^{2^n-1} \epsilon_0^{2^n}$$

We now aim to show that we have convergence, i.e.  $\lim_{n \rightarrow \infty} x_n = x^*$ ; to do this it suffices to show that  $\lim_{n \rightarrow \infty} \epsilon_n = 0$ .

Consider the sequence  $(z_n := M^{2^n-1} \epsilon_0^{2^n} : n \in \mathbb{N}_0)$ . We know that  $0 \leq \epsilon_n \leq z_n \forall n \in \mathbb{N}_0$ , so it then follows that if  $\lim_{n \rightarrow \infty} z_n = 0$ , then  $\lim_{n \rightarrow \infty} \epsilon_n = 0$  by the Squeeze Theorem ??.

Now as  $M\epsilon_0 < 1$  by NR<sub>3</sub>, then we see that:

$$\begin{aligned} \lim_{n \rightarrow \infty} z_n &= \lim_{n \rightarrow \infty} (M\epsilon_0)^{2^n-1} \epsilon_0 \\ &= \epsilon_0 \lim_{n \rightarrow \infty} (M\epsilon_0)^{2^n-1} \\ &= \epsilon_0 \cdot 0 && \text{because } \lim_{n \rightarrow \infty} r^{m_n} = 0 \text{ where } |r| < 1, m_n \geq n \forall n \in \mathbb{N} \\ &= 0 \end{aligned}$$

Now to show that this sequence converges quadratically we see that  $\epsilon_{n+1} = \frac{1}{2} \left| \frac{f''(y_n)}{f'(x_n)} \right| \epsilon_n^2$ , and therefore  $\frac{\epsilon_{n+1}}{\epsilon_n^2} = \frac{1}{2} \left| \frac{f''(y_n)}{f'(x_n)} \right|$ .

Because  $|x^* - y_n| < |x^* - x_n|$  and  $\lim_{n \rightarrow \infty} x_n = x^*$ , then it follows that  $\lim_{n \rightarrow \infty} y_n = x^*$ . Therefore we see that

$$\lim_{n \rightarrow \infty} \frac{\epsilon_{n+1}}{\epsilon_n^2} = \frac{1}{2} \left| \frac{f''(x^*)}{f'(x^*)} \right| \in \mathbb{R}^+$$

Hence as the above limit exists and is positive then the sequence is quadratically convergent.  $\square$

## 2.4 Efficiency and Accuracy Metrics



## **3 Division and Multiplication**

### **3.1 Introduction**

Though the idea of Division and Multiplication can seem fairly simple, particularly from an abstract pure mathematical point of view, these calculations can be computationally difficult. This section will show a few algorithms designed to calculate these values, and discuss their implementation and efficiency.

### **3.2 Multiplication**

#### **3.2.1 Basic Methods**

#### **3.2.2 Advanced Methods**

#### **3.2.3 Analysis**

### **3.3 Division**

#### **3.3.1 Basic Methods**

#### **3.3.2 Advanced Methods**

#### **3.3.3 Analysis**

TODO: Section to be filled out, no work currently done on this section

## 4 Trigonometric Functions

This section will focus on trigonometric functions, which are commonly used cyclic functions. These functions have been studied for hundreds of years, and can be challenging to calculate. We will discuss several methods of calculating them below before comparing methods.

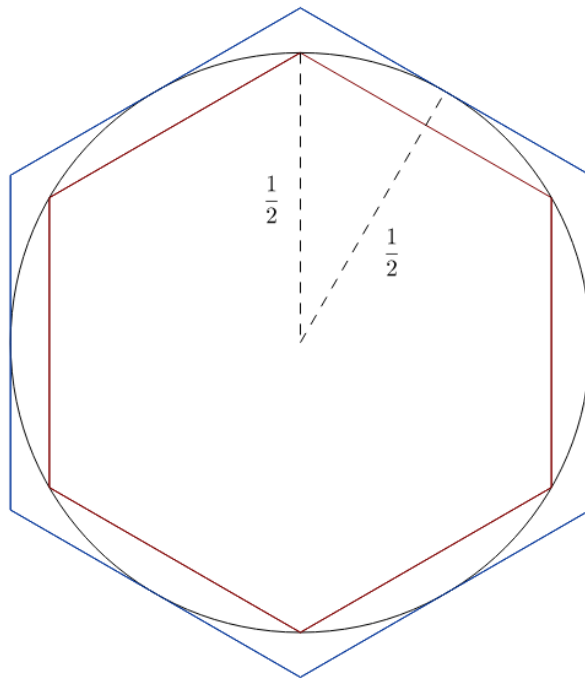
**TODO: Extend and Eloquate introduction**

### 4.1 Calculating $\pi$

Several of the methods in this section require that we already know the value of  $\pi$ , for example when we are applying several trig identities. Here we will briefly discuss several methods for calculating the value of  $\pi$ , so that we may use this value in later subsections.

The first method to consider is the method used by ancient mathematicians, such as the Greeks and Chinese. We know that if the radius of the circle is  $\frac{1}{2}$ , then the circumference of the circle is  $\pi$ , and the value is between the perimeters of the inner and outer polygon perimeters. The internal perimeter is  $p_n = n \sin(\frac{\pi}{n})$  and the external perimeter is  $P_n = n \tan(\frac{\pi}{n})$ .

Figure 4.1.1: Ancient method of calculating  $\pi$



As we know the values of  $\tan(\frac{\pi}{6})$  and  $\sin(\frac{\pi}{6})$ , then we can calculate  $P_6$  and  $p_6$ . It has been shown that  $P_{2n} = \frac{2p_n P_n}{p_n + P_n}$  and  $p_{2n} = \sqrt{p_n P_{2n}}$ , which allows us to create an iterative method to approximate  $\pi$ , by taking the mid-point of the successive polygon perimeters.

Other common historical methods for approximating  $\pi$  are to use infinite series. One such method uses the series expansion of  $\tan^{-1}$ , which is discussed in detail below, where  $\tan^{-1}(1) = \frac{\pi}{4}$ . This gives the following approximation using  $N$  terms:

$$\pi = 4 \sum_{n=0}^N \frac{(-1)^n}{2n+1} = \sum_{n=0}^N \frac{8}{(4n+1)(4n+3)} \quad (4.1.1)$$

This sequence converges very slowly, with sublinear convergence, to the correct value. More modern methods have typically revolved around finding more rapidly converging infinite series, examples include Ramanujan's series:

$$\frac{1}{\pi} = \frac{2\sqrt{2}}{9801} \sum_{n=0}^{\infty} \frac{(4n)!(1103 + 26390n)}{(k!)^n 396^{4n}} \quad (4.1.2)$$

or the Chudnovsky algorithm:

$$\frac{1}{\pi} = 12 \sum_{n=0}^{\infty} \frac{(-1)^n (6n)! (13591409 + 545140134n)}{(3n)! (n!)^3 640320^{3n + \frac{3}{2}}} \quad (4.1.3)$$

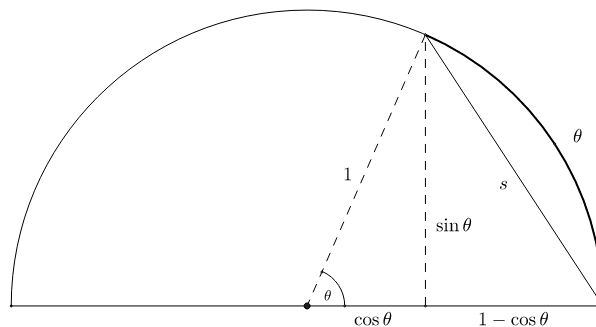
This final series is extremely rapidly convergent to the value of  $\frac{1}{\pi}$ , for example just the first term gives  $\pi$  accurate to 13 decimal places while we can get  $\pi$  accurate to 1000 decimal places with summing just 71 terms. Compared to Equation ?? which takes the summation of 500 terms to achieve the same 1000 digits of accuracy.

To get large degrees of accuracy for  $\pi$  is extremely computer intensive and using the `mpfr` requires the number of bits of precision and number of terms to be set. This makes calculating  $\pi$  to a large number of decimal places, for example 1000000, computationally infeasible on a regular home computer. Therefore for our purposes we will use the precalculated value of  $\pi$  to 1000000 decimal places as listed on [http://www.exploratorium.edu/pi/pi\\_archive/Pi10-6.html](http://www.exploratorium.edu/pi/pi_archive/Pi10-6.html)

## 4.2 Geometric Method

The first method I will be discussing is a method based on geometric properties that are derived on a circle, and we will start by considering values of  $\cos$  in the range  $[0, \frac{\pi}{2}]$ . To do this we will consider the following figure of the unit circle:

Figure 4.2.1: Diagram showing angles to be dealt with



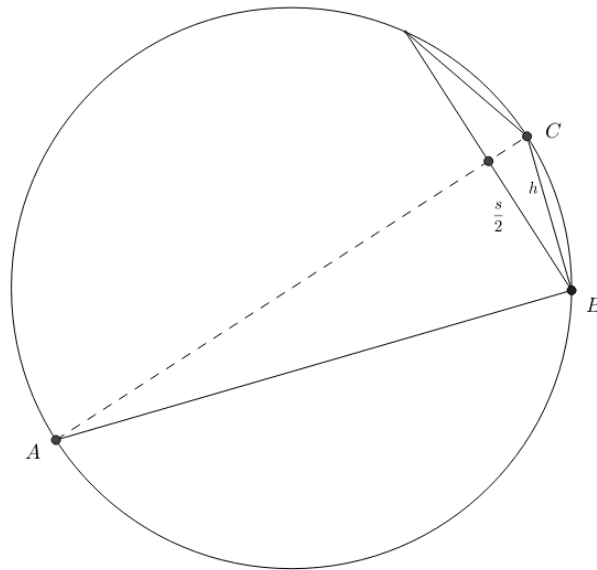
Here theta will be given in radians, and we can note that the labelled arc has length  $\theta$  due to the formula for the circumference of a circle. By using the following derivation we can find a formula for  $\theta$  in terms of  $s$ :

$$\begin{aligned} s^2 &= \sin^2 \theta + (1 - \cos \theta)^2 \\ &= (\sin^2 \theta + \cos^2 \theta) + 1 - 2 \cos \theta \\ &= 2 - 2 \cos \theta \end{aligned} \quad \text{By using } \sin^2 \theta + \cos^2 \theta = 1$$

$$\cos \theta = 1 - \frac{s^2}{2}$$

We will now consider a second diagram which will allow us to calculate an approximate value of  $s$ .

Figure 4.2.2: Diagram detailing how to calculate  $s$

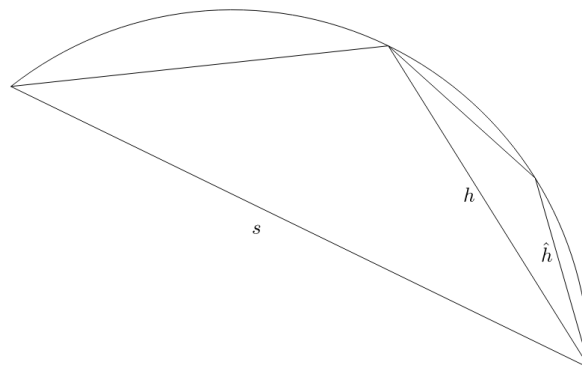


We will first note that by an elementary geometry result we can know that the angle  $ABC$  is a right-angle; also we can consider that  $h$  is an approximation of  $\frac{\theta}{2}$ , which will become relevant later. Now because  $AC$  is a diameter of our circle then it's length is 2 and thus, by utilising Pythagoras' Theorem, we get that the length of  $AB$  is  $\sqrt{AC^2 - BC^2} = \sqrt{4 - h^2}$ .

From here we consider the area of triangle  $ABC$ , which can be calculated as  $\frac{1}{2} \cdot h \cdot \sqrt{4 - h^2}$  and as  $\frac{1}{2} \cdot 2 \cdot \frac{s}{2}$ ; by equating these two, squaring both sides and re-arranging we get that  $s^2 = h^2(4 - h^2)$ . Now we have the basis for a method that will allow us to calculate  $\cos \theta$ .

To complete our method we will consider introducing a new line that is to  $h$  what  $h$  is to  $s$  as shown in the diagram below:

Figure 4.2.3: Detailing the recursive steps



It is easy to see that if we repeat the steps above we get that  $h^2 = \hat{h}^2(4 - \hat{h}^2)$ , and it also follows that  $\hat{h} \approx \frac{\theta}{4}$ . Using this we can take an initial guess of  $h_0 := \frac{\theta}{2^k}$ , for some  $k \in \mathbb{N}$ , and

then calculate  $h_{n+1}^2 = h_n^2(4 - h_n^2)$  where  $n \in [0, k] \cap \mathbb{Z}$ ; finally we calculate  $\cos \theta = 1 - \frac{h_k^2}{2}$ , giving the following algorithm:

Algorithm 4.2.1: Geometric calculation of  $\cos$

```

1  geometric_cos ( $\theta \in [0, \frac{\pi}{2}], k \in \mathbb{N}$ )
2       $h_0 := \frac{\theta}{2^k}$ 
3       $n := 0$ 
4      while  $n < K$ :
5           $h_{n+1}^2 := h_n^2 \cdot (4 - h_n^2)$ 
6           $n \mapsto n + 1$ 
7      return  $1 - \frac{h_k^2}{2}$ 

```

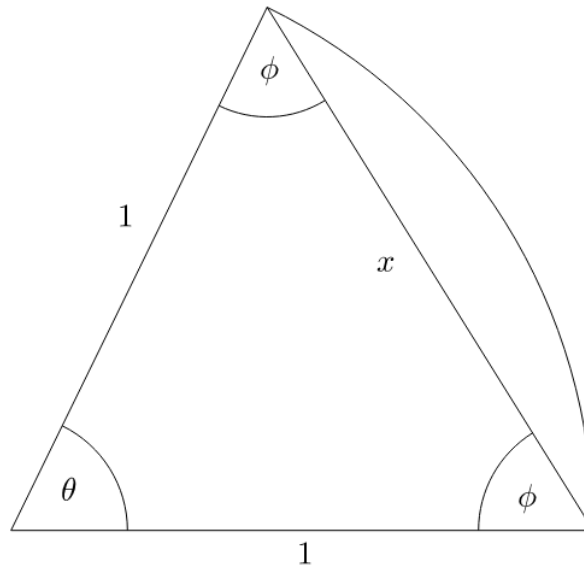
Now we can use the above pseudocode to calculate any trigonometric function value by using various trigonometric identities. First we suppose  $\theta \in \mathbb{R}$ , then we can repeatedly apply the identity  $\cos \theta = \cos(\theta \pm 2\pi)$  to either add or subtract  $2\pi$  until we have a value  $\theta' \in [0, 2\pi)$ . Once we have this value we can utilise the following assignment to calculate  $\cos \theta$ :

$$\cos \theta = \begin{cases} \cos \theta' & : \theta' \in [0, \frac{\pi}{2}] \\ -\cos(\pi - \theta') & : \theta' \in [\frac{\pi}{2}, \pi] \\ -\cos(\theta' - \pi) & : \theta' \in [\pi, \frac{3\pi}{2}] \\ \cos(2\pi - \theta') & : \theta' \in [\frac{3\pi}{2}, 2\pi) \end{cases}$$

Using Algorithm 4.2.1 we can also easily calculate both  $\sin \theta$  and  $\tan \theta$ , by further use of trigonometric identities. In particular we note that  $\sin \theta = \cos(\theta - \frac{\pi}{2})$  and  $\tan \theta = \frac{\sin \theta}{\cos \theta}$ . Hence we can now calculate the trigonometric function value of any angle.

We now wish to analyse the error of our approximation for  $\cos$ , as the other methods have errors that are derivative of the error for approximating  $\cos$ . Now Figure 4.2.4 shows an arc of a circle which creates chord  $x$ , with this we will be able to calculate the exact length of the chord and thus work on the error of our approximations.

Figure 4.2.4: Diagram to find actual arc approximation



To start we will note that  $\phi = \frac{\pi-\theta}{2} = \frac{\pi}{2} - \frac{\theta}{2}$ , and then by using the Sine Law we get

$$\frac{x}{\sin \theta} = \frac{1}{\sin \phi} \implies x = \frac{\sin \theta}{\sin \phi}$$

Now we can recall the double angle formula for sin, which gives  $\sin \theta = 2 \sin \frac{\theta}{2} \cos \frac{\theta}{2}$ , and also  $\sin \phi = \cos \frac{\theta}{2}$ . This allows us to see that

$$x = \frac{2 \sin \frac{\theta}{2} \cos \frac{\theta}{2}}{\cos \frac{\theta}{2}} = 2 \sin \frac{\theta}{2}$$

Therefore we see that  $h_n$  is approximating the chord length associated with angle  $\theta 2^{n-k}$ , and thus  $\epsilon_n = |h_n - 2 \sin(\theta 2^{n-k-1})|$ . Now as  $h_0 = \theta 2^{-k} \approx 2 \sin(\theta 2^{-k-1})$  then it follows that  $\exists \phi$  such that  $h_0 = 2 \sin(\phi 2^{-k-1})$ , from this we can see that  $\phi = 2^{k+1} \sin^{-1}(\theta 2^{-k-1})$ . We will use these facts to prove a couple of propositions.

**Proposition 4.2.1.**  $h_n = 2 \sin(\phi 2^{n-k-1}) \forall n \in [0, k] \cap \mathbb{Z}$  where  $\phi := 2^{k+1} \sin^{-1}(\theta 2^{-k-1})$ .

*Proof.* Proceed by induction on  $n \in [0, k] \cap \mathbb{Z}$ .

$$\mathbf{H}(n): h_n = 2 \sin(\phi 2^{n-k-1})$$

$$\mathbf{H}(0):$$

$$\begin{aligned} 2 \sin(\phi 2^{-k-1}) &= 2 \sin(\sin^{-1}(\phi 2^{-k-1})) \\ &= \theta 2^{-k} \\ &= h_0 \end{aligned} \quad \text{by definition of } h_0$$

$$\mathbf{H}(n) \implies \mathbf{H}(n+1):$$

$$\begin{aligned} h_{n+1} &= h_n \sqrt{4 - h_n^2} \\ &= 2 \sin(\phi 2^{n-k-1}) \sqrt{4 - 4 \sin^2(\phi 2^{n-k-1})} && \text{by } \mathbf{H}(n) \\ &= 4 \sin(\phi 2^{n-k-1}) \cos(\phi 2^{n-k-1}) \\ &= 2 \sin(\phi 2^{n-k}) && \text{by the use of double angle formulas} \end{aligned}$$

□

**Proposition 4.2.2.**  $h_n > 2 \sin(\theta 2^{n-k-1}) \forall n \in [0, k] \cap \mathbb{Z}$

*Proof.* We start by considering the expansion of the exact value of  $h_n$ .

$$\begin{aligned} h_n &= 2 \sin(\phi 2^{n-k-1}) \\ &= 2 \sin(2^{n-k-1} (2^{k+1} \sin^{-1}(\theta 2^{-k-1}))) \\ &= 2 \sin(2^n \sin^{-1}(\theta 2^{-k-1})) \\ &= 2 \sin(\theta 2^{n-k-1} + \frac{1}{6} \theta^3 2^{n-3k-3} + \mathcal{O}(2^{-5k})) \quad \text{Detailed in section ??} \end{aligned}$$

Now as we know that  $n \leq k$ , then it follows that  $\theta 2^{n-k-1} \leq \frac{1}{2} \theta$ .

Also as  $\theta \leq \frac{\pi}{2}$  we know that  $\theta 2^{n-k-1} \leq \frac{\pi}{4}$ .

We can also show that  $\frac{1}{6}\theta^3 2^{n-3k-3} + \mathcal{O}(2^{-5k}) \leq \frac{\pi}{4}$ , though the proof is omitted here for brevity; therefore we see that  $\phi 2^{n-k-1} \leq \frac{\pi}{2}$ , and obviously that  $\phi 2^{n-k-1} > \theta 2^{n-k-1}$ .

Hence, as  $\sin$  is an increasing function in the range  $[0, \frac{\pi}{2}]$ , we conclude that

$$h_n = 2 \sin(\phi 2^{n-k-1}) > 2 \sin(\theta 2^{n-k-1})$$

□

With these two propositions we can now consider the error of our approximation of  $\cos$ . First we will prove the following proposition regarding the error of the approximation of  $s$ :

**Proposition 4.2.3.** *If  $\epsilon_n := |h_n - 2 \sin(\theta 2^{n-k-1})| \forall n \in [0, k] \cap \mathbb{Z}$ , then  $\epsilon_k < 2^k \epsilon_0$ .*

*Proof.*  $\epsilon_n = h_n - 2 \sin(\theta 2^{n-k-1})$  as  $h_n > 2 \sin(\theta 2^{n-k-1})$  by Proposition 4.2.2.

Now we see that:

$$\begin{aligned} \epsilon_{n+1} &= h_{n+1} - 2 \sin(\theta 2^{n-k}) \\ &= h_n \sqrt{4 - h_n^2} - 4 \sin(\theta 2^{n-k-1}) \cos(\theta 2^{n-k-1}) \end{aligned}$$

If we consider the equation  $\alpha\beta - \gamma\delta = (\alpha - \gamma) + \alpha(\beta - 1) - \gamma(\delta - 1)$  and apply it to our current formula we get:

$$\begin{aligned} \epsilon_{n+1} &= (h_n - 2 \sin(\theta 2^{n-k-1})) + h_n(\sqrt{4 - h_n^2} - 1) - 2 \sin(\theta 2^{n-k-1})(2 \cos(\theta 2^{n-k-1}) - 1) \\ &= \epsilon_n + h_n(\sqrt{4 - h_n^2} - 1) - 2 \sin(\theta 2^{n-k-1})(2 \cos(\theta 2^{n-k-1}) - 1) \\ &= 2\epsilon_n + h_n(\sqrt{4 - h_n^2} - 2) - 2 \sin(\theta 2^{n-k-1})(2 \cos(\theta 2^{n-k-1}) - 2) \\ &= 2\epsilon_n + h_n(\sqrt{4 - h_n^2} - 2) + 2 \sin(\theta 2^{n-k-1})(2 - 2 \cos(\theta 2^{n-k-1})) \\ &< 2\epsilon_n + h_n(\sqrt{4 - h_n^2} - 2 \cos(\theta 2^{n-k-1})) \\ &< 2\epsilon_n + h_n(\sqrt{4 - 4 \sin^2(\theta 2^{n-k-1})} - 2 \cos(\theta 2^{n-k-1})) \\ &= 2\epsilon_n + h_n(2 \cos(\theta 2^{n-k-1}) - 2 \cos(\theta 2^{n-k-1})) \\ &= 2\epsilon_n \end{aligned}$$

The inequalities in the above derivation arise from the fact that  $h_n > 2 \sin(\theta 2^{n-k-1})$  by Proposition 4.2.2.

Hence as we now know that  $\epsilon_{n+1} < 2\epsilon_n$ , we then see that  $\epsilon_n < 2^n \epsilon_0$ . Therefore we prove our statement that

$$\epsilon_k < 2^k \epsilon_0$$

□

Obviously  $\epsilon_k = |h_k - s|$ , and we can now use this to find the error of our final answer. First we will start by letting  $\mathcal{C} := 1 - \frac{1}{2}h_k^2$  and note that analytically  $\cos\theta = 1 - \frac{1}{2}s^2$ . Therefore we will now consider  $\epsilon_{\mathcal{C}} = |\mathcal{C} - \cos(\theta)|$ :

$$\begin{aligned}
\epsilon_C &= \left| 1 - \frac{h_k^2}{2} - 1 + \frac{s^2}{2} \right| \\
&= \frac{1}{2} |h_k^2 - s^2| \\
&= \frac{1}{2} |h_k h_k - 2 \sin(\frac{\theta}{2}) 2 \sin(\frac{\theta}{2})| \\
&= \frac{1}{2} (h_k h_k - 2 \sin(\frac{\theta}{2}) 2 \sin(\frac{\theta}{2})) \quad \text{as } 2 \sin(\frac{\theta}{2}) < h_k \\
&= \frac{1}{2} (2\epsilon_k + h_k(h_k - 2) - 2 \sin(\frac{\theta}{2})(2 \sin(\frac{\theta}{2}) - 2)) \\
&< \frac{1}{2} (2\epsilon_k + h_k(h_k - 2 \sin(\frac{\theta}{2}))) \\
&= \frac{1}{2} (2 + h_k) \epsilon_k \\
&= \frac{1}{2} (2 + 2 \sin(\frac{\phi}{2})) \epsilon_k \\
&= (1 + \sin(\frac{\phi}{2})) \epsilon_k \\
&\leq 2\epsilon_k
\end{aligned}$$

As  $\epsilon_C \leq 2\epsilon_k$ , then by Proposition 4.2.3 we see that  $\epsilon_C < 2^{k+1}\epsilon_0$ . Now to consider  $\epsilon_0$  we first observe that  $\epsilon_0 = \theta 2^{-k} - 2 \sin \theta 2^{-k-1}$ , and therefore we can conclude that:

$$\epsilon_C < 2\theta - 2^{k+2} \sin(\theta 2^{-k-1})$$

This looks like an error that may infact grow exponentially large as  $k \rightarrow \infty$ , due to the multiplication by  $2^{k+2}$ . However if we instead consider the series expansion of  $\sin(x)$ , shown in Section 4.3 to be  $\sin(x) = x - \frac{1}{3!}x^3 + \frac{1}{5!}x^5 - \dots$ , and substitute that into our equation we see that:

$$\begin{aligned}
\epsilon_C &< 2\theta - 2^{k+2}(\theta 2^{-k-1} - \frac{1}{3!}\theta^3 2^{-3k-3} + \frac{1}{5!}\theta^5 2^{-5k-5} - \dots) \\
&= 2\theta - 2\theta + \frac{1}{3}\theta^3 2^{-2k-1} - \frac{1}{5!}\theta^5 2^{-4k-3} + \dots \\
&= \frac{1}{3}\theta^3 2^{-2k-1} - \frac{1}{5!}\theta^5 2^{-4k-1} + \dots
\end{aligned}$$

Now obviously the last line tends towards zero as  $k$  tends to infinity, due to it being a formula of order  $\mathcal{O}(2^{-2k-1})$ . Therefore we know that  $\forall \tau \in \mathbb{R}^+ \exists \mathcal{K} \in \mathbb{N} : \epsilon_{C,k} < \tau \forall k \in [\mathcal{K}, \infty) \cap \mathbb{Z}$ . In particular, if we then wish to calculate  $\cos \theta$  accurate to  $N$  decimal places then we are looking to find  $k \in \mathbb{N}$  such that:

$$2\theta - 2^{k+2} \sin(\theta 2^{-k-1}) < 10^{-N} \implies 2^{k+2} \sin(\theta 2^{-k-1}) > 2\theta - 10^{-N}$$

For an example of the above in action we will be taking  $\theta = 0.5$ . The table below shows the minimum  $k \in \mathbb{N}$  to guarantee  $N$  digits of accuracy in the result:



$N$	$k$
5	6
10	14
50	80
100	163
1000	1658

As can be seen the value of  $k$  required to achieve  $N$  digits of accuracy increases roughly linearly when  $\theta = 0.5$ . Testing for other values of  $\theta$  reveals them to have similar required values for  $k$ , at least within the same order of each other.

Another consideration for Algorithm 4.2.1 is that we could "run it in reverse" to attain an algorithm for the inverse cosine function. To start take line 7 which is  $\mathcal{C} = 1 - \frac{1}{2}h_k^2$ , which can be re-arranged to give  $h_k^2 = 2 - 2\mathcal{C}$ , where we know  $\mathcal{C}$  as our initial value.

Line 5 is a little more difficult, but by re-arranging we see that  $h_n^4 - 4h_n^2 + h_{n+1}^2 = 0$ , which can be solved via the quadratic formula to give  $h_n^2 = 2 \pm \sqrt{4 - h_{n+1}^2}$ . Now we can make the observation that if  $x \in \mathbb{R}_0^+$ , then  $\cos^{-1}(-x) = \pi - \cos^{-1}(x)$  and so we can restrict our algorithm to only consider  $x \in [0, 1]$ . With this we know that  $\theta \in [0, \frac{\pi}{2}]$ , and thus  $h_k \leq \sqrt{2}$ . Therefore as  $h_{n+1} > h_n \forall n \in [0, k-1] \cap \mathbb{Z}$  we see that  $h_n^2 \leq 2 \forall n \in [0, k] \cap \mathbb{Z}$ . This allows us to ascertain that to reverse Line 5 we perform  $h_n^2 = 2 - \sqrt{4 - h_{n+1}^2}$ .

Finally line 2 is reversed by returning the value  $2^k h_0$ ; therefore we get the following algorithm for  $\cos^{-1}(x)$  where  $x \in [0, 1]$ :

Algorithm 4.2.2: Geometric calculation of  $\cos^{-1}$

```

1  geometric_aCos ( $x \in [0, 1], k \in \mathbb{N}$ )
2       $h_k := 2 - 2x$ 
3       $n := k - 1$ 
4      while  $n \geq 0$ :
5           $h_n^2 := 2 - \sqrt{4 - h_{n+1}^2}$ 
6           $n \mapsto n - 1$ 
7      return  $2^k h_0$ 

```

Similar to the regular trigonometric functions we can use trigonometric identities to calculate the inverse trigonometric functions from  $\cos^{-1}$ . To start we recall that  $\cos^{-1}(-x) = \pi - \cos^{-1}(x)$  where  $x \in [0, 1]$ , then we can use the identities that  $\sin^{-1}(x) = \frac{\pi}{2} - \cos^{-1}(x)$  and  $\tan^{-1}(x) = \sin^{-1}(\frac{x}{\sqrt{x^2+1}})$ .

If we suppose that all operations in the method are accurately computed then Algorithm 4.2.2 is a computation with high accuracy. This is because there is no initial guess, such as in Algorithm 4.2.1, and so the only introduction of error is assuming that  $2^k h_0 \approx \theta$ . However as we discuss in detail in Section ??, calculating square roots is not a simple task and thus will introduce error to the method in general; therefore the accuracy of the method is roughly as accurate as our method of calculating square roots.

### 4.3 Taylor Series

If we consider our definition of a McClaurin Series from Section ??, we can use this to approximate our Trigonometric Functions. Consider first  $\cos \theta$ , for which we know that  $\frac{d}{d\theta} \cos \theta =$

$-\sin \theta$ ; it then follows that  $\frac{d^2}{d\theta^2} \cos \theta = -\cos \theta$ ,  $\frac{d^3}{d\theta^3} \cos \theta = \sin \theta$  and  $\frac{d^4}{d\theta^4} \cos \theta = \cos \theta$ .

If we let  $f(x) = \cos x$  and use the known values  $\cos(0) = 1$  and  $\sin(0) = 0$ , then we see that:

$$f^{(n)}(0) = \begin{cases} 1 & : 4 \mid n \\ 0 & : 4 \mid n-1 \\ -1 & : 4 \mid n-2 \\ 0 & : 4 \mid n-3 \end{cases}$$

By simplifying this by omitting the 0 coefficient terms we get the following series:

$$\sum_{n=0}^{\infty} \frac{(-1)^n}{(2n)!} x^{2n} = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \dots \quad (4.3.1)$$

By using similar working we can get that the series associated with  $\sin x$ :

$$\sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} x^{2n+1} = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots \quad (4.3.2)$$

Before we go any further we need to consider when Equations 4.3.1 and 4.3.2 converge to their respective functions. To do this we will use the ratio test for series as defined in ??, using Equation 4.3.1 we see that

$$\begin{aligned} L_C &= \lim_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| \\ &= \lim_{n \rightarrow \infty} \left| \frac{\frac{(-1)^{n+1}}{(2n+2)!} x^{2n+2}}{\frac{(-1)^n}{(2n)!} x^{2n}} \right| \\ &= \frac{(2n)!}{(2n+2)!} |x|^2 \\ &= \frac{1}{(2n+2)(2n+1)} |x|^2 \end{aligned}$$

Now it is easy to see that,  $L_C = 0$  for all values of  $x$  as the fractional component decreases as  $n$  increases and  $|x|^2$  is a constant. Therefore we can conclude that Equation 4.3.1 converges to  $\cos(x)$  for all values of  $x$ . We can use a very similar deduction to show that Equation 4.3.2 converges to  $\sin(x)$  for all values of  $x$ .

The above means that  $\cos$  and  $\sin$  can be approximated using Taylor Polynomials, in particular for a given  $N \in \mathbb{N}$ :

$$\cos x \approx \sum_{n=0}^N \frac{(-1)^n}{(2n)!} x^{2n} \quad \text{and} \quad \sin x \approx \sum_{n=0}^N \frac{(-1)^n}{(2n+1)!} x^{2n+1}$$

This allows us to create the following two methods for computing  $\cos x$  and  $\sin x$ :

#### Algorithm 4.3.1: Taylor computation of $\cos$ and $\sin$

```

1  taylor_cos( $x \in \mathbb{R}, N \in \mathbb{N}$ )
2     $\mathcal{C} := 0$ 
3     $n := 0$ 
4    while  $n < N$ :
5       $\mathcal{C} \mapsto \mathcal{C} + (-1)^n \cdot \frac{1}{(2n)!} x^{2n}$ 

```

```

6       $n \mapsto n + 1$ 
7      return  $\mathcal{C}$ 
8
9      taylor_sin( $x \in \mathbb{R}, N \in \mathbb{N}$ )
10      $\mathcal{S} := 0$ 
11      $n := 0$ 
12     while  $n < N$ :
13          $\mathcal{S} \mapsto \mathcal{S} + (-1)^n \cdot \frac{1}{(2n+1)!} x^{2n+1}$ 
14          $n \mapsto n + 1$ 
15     return  $\mathcal{S}$ 

```

As these two methods are obviously very similar and the fact that  $\sin(x) = \cos(x - \frac{\pi}{2})$ , we will continue by examining only the taylor method for approximating  $\cos$ . We will assume that any calculations for  $\sin$  are transformed into a problem of finding a  $\cos$  value.

It should be noted that this  $\cos$  algorithm is particularly inefficient to calculate on a computer implementation; this is primarily due to the way in which the update of  $\mathcal{C}$  is calculated each loop.

In each loop we are calculating  $x^{2n}$ , which has a naive complexity of  $\mathcal{O}(2n)$ . However what we are actually calculating  $x^{2(n-1)} \cdot x^2$  and thus if we store the values of  $x^{2(n-1)}$  and  $x^2$ , the complexity of this step drops to  $\mathcal{O}(1)$ . Similarly we are also calculating  $\frac{1}{(2n)!}$  in each loop which, by the same logic, is  $\frac{1}{2(n-1)!} \cdot \frac{1}{(2n)(2n-1)}$ , and we can use the same storage and update method as for  $x^{2n}$ .

As another step towards optimizing the algorithm we can start with an initial value of  $\mathcal{C} = 1$ , and then perform two updates of  $\mathcal{C}$  each loop until we reach or surpass  $N$ . This saves calculating  $(-1)^n$  each loop, by explicitly performing two different calculations. Implementing all of the above gives us the following two updated methods:

Algorithm 4.3.2: Taylor computation of  $\cos$  optimised

```

1      taylor_cos( $x \in \mathbb{R}, N \in \mathbb{N}$ )
2       $\mathcal{C} := 1$ 
3       $x_2 := x^2$ 
4       $a := 1$ 
5       $b := 1$ 
6       $n := 1$ 
7      while  $n < N$ :
8           $a \mapsto a \cdot \frac{1}{(2n-1)(2n)}$ 
9           $b \mapsto b \cdot x_2$ 
10          $\mathcal{C} \mapsto \mathcal{C} - a \cdot b$ 
11          $a \mapsto a \cdot \frac{1}{(2n+1)(2n+2)}$ 
12          $b \mapsto b \cdot x_2$ 
13          $\mathcal{C} \mapsto \mathcal{C} + a \cdot b$ 
14          $n \mapsto n + 2$ 
15     return  $\mathcal{C}$ 

```

As the next term of the polynomial is known definitively then we can see that it is very easy

to calculate the error of our approximation. We see that

$$\begin{aligned}
\epsilon_N &= |\cos(x) - \text{taylor\_cos}(x, N)| \\
&= \mathcal{O}(|x|^{N'+1}) \quad \text{where } N' \text{ is the smallest} \\
&\quad \text{odd integer such that } N' \geq N \\
&\leq \frac{1}{(2(N'+1))!} |x|^{N'+1} \\
&\leq \frac{1}{(2(N+1))!} |x|^{N+1}
\end{aligned}$$

If we place bounds on the value of  $\cos$  calculated as in Section 4.2, then we know that  $|x| \leq \frac{\pi}{2}$ , and thus we get the following bound for the error of our approximation:

$$\epsilon_N \leq \frac{\pi^{N'+1}}{2^{N'+1}(2(N'+1))!}$$

Thus if we find  $N \in \mathbb{N}$  such that  $\frac{\pi^{N+1}}{2^{N+1}(2(N+1))!} < \tau \in \mathbb{R}^+$  then we know that  $\epsilon_N < \tau$ . If we consider  $\tau = 10^{-k}$ , then we can find  $N \in \mathbb{N}$  such that our approximation is accurate to  $k$  decimal places. Below is a table which details some values of  $k$  and the corresponding minimum  $N$  to guarantee  $k$  decimal places of accuracy:

$k$	$N$
5	4
10	7
50	21
100	36
1000	233

Now for  $\tan x$  we can either calculate both  $\sin x$  and  $\cos x$  using  $\text{taylor\_cos}(x, N)$  and divide the resulting value, or we can calculate  $\tan x$  directly using a Taylor expansion.

In calculating the McClaurin series for  $\tan x$  we start by letting  $\tan x = \sum_{n=0}^{\infty} a_n x^n$ , and then noting that as  $\tan x$  is an odd series then it's McClaurin series only contains non-zero coefficients for odd powers of  $x$ ; therefore we get that  $\tan x = \sum_{n=0}^{\infty} a_{2n+1} x^{2n+1} = a_1 x + a_3 x^3 + a_5 x^5 + \dots$ .

Next we consider that  $\frac{d}{dx} \tan x = 1 + \tan^2 x$ , and knowing the McClaurin series form of  $\tan x$  we get the following:

$$\begin{aligned}
\sum_{n=0}^{\infty} (2n+1) a_{2n+1} x^{2n} &= 1 + \left( \sum_{n=0}^{\infty} a_{2n+1} x^{2n+1} \right)^2 \\
&= 1 + a_1^2 x^2 + (2a_1 a_3) x^4 + (2a_1 a_5 + a_3^2) x^6 + \dots
\end{aligned}$$

Considering the co-efficients of powers on the right hand side of the above equation we see that  $2a_1 a_3 = a_1 a_3 + a_3 a_1 = a_1 a_{4-1} + a_3 a_{4-3}$  and  $2a_1 a_5 + a_3^2 = a_1 a_5 + a_3 a_3 + a_5 a_1 = a_1 a_{6-1} + a_3 a_{6-3} + a_5 a_{6-5}$ . This indicates that our general form for the co-efficient of  $2n$  on the right hand side is  $\sum_{k=1}^n a_{2k-1} a_{2n-2k+1}$ , and thus returning to our equation we get

$$a_1 + \sum_{n=1}^{\infty} (2n+1) a_{2n+1} x^{2n} = 1 + \sum_{n=1}^{\infty} \left( \sum_{k=1}^n a_{2k-1} a_{2n-2k+1} \right) x^{2n}$$

Using this we conclude that  $a_1 = 1$  and  $a_{2n+1} = \frac{1}{2n+1} \sum_{k=1}^n a_{2k-1} a_{2n-2k+1} \forall n \in \mathbb{N}$ . We can note immediately that the calculation of any previous co-efficients will provide no help in calculating later co-efficients and so the entire sum must be calculated each loop, while also storing each co-efficient already calculated.

This means that the complexity to calculate co-efficient  $a_{2n+1}$  is  $\mathcal{O}(n)$  and will be the  $n^{\text{th}}$  such calculation, making the complexity of calculating  $n$  co-efficients to be  $\mathcal{O}(n^2)$ . Comparing this to the `taylor_cos` method we see that to calculate up to  $n$  co-efficients of both `cos` and `sin` has complexity  $\mathcal{O}(n)$ . Therefore it is more efficient to calculate `tan` by calculating both `cos` and `sin` using Algorithm ??, and performing division than directly using Taylor Polynomial approximation.

We would also like to be able to calculate the inverse trigonometric functions using this method, which means we need to find our McClaurin series of the inverse trigonometric functions. The simplest of these is  $\tan^{-1}$ , where we start by recalling that  $\frac{d}{dx} \tan^{-1} x = \frac{1}{1+x^2}$  and then by intergrating both sides we get:

$$\begin{aligned} \tan^{-1} x &= \int \frac{1}{1+x^2} dx \\ &= \int (1 - (-x^2))^{-1} dx \\ &= \int \sum_{n=0}^{\infty} (-x^2)^n dx && \text{by Equation ??} \\ &= \int \sum_{n=0}^{\infty} (-1)^n x^{2n} dx \\ &= c + \sum_{n=0}^{\infty} \frac{(-1)^n}{2n+1} x^{2n+1} \end{aligned}$$

As  $\tan^{-1}(0) = 0$  then we see that  $c = 0$  and thus gives us the following formula for  $\tan^{-1}$ :

$$\tan^{-1} x = \sum_{n=0}^{\infty} \frac{(-1)^n}{2n+1} x^{2n+1}$$

Now due to the restrictions from Equation ?? the above is only valid for  $x \in [-1, 1]$ , but we know that the domain of  $\tan^{-1}$  is  $x \in \mathbb{R}$ . To fix this we will first recognise that  $\tan^{-1}(-x) = -\tan^{-1}(x)$ , so we can restric our problem to  $x \in \mathbb{R}_0^+$ . Now if we take the double angle formula for `tan`:

$$\tan(\alpha + \beta) = \frac{\tan(\alpha) + \tan(\beta)}{1 - \tan(\alpha) \tan(\beta)}$$

By substituting  $\alpha = \tan^{-1}(x)$  and  $\beta = \tan^{-1}(y)$  into the above then we get

$$\tan^{-1}(x) + \tan^{-1}(y) = \tan^{-1} \left( \frac{x+y}{1-xy} \right)$$

Using this, suppose we are looking for  $\tan^{-1}(z)$  where  $z \in (1, \infty)$  and let  $y = 1$ , then  $\tan^{-1}(y) = \frac{\pi}{4}$ . We can then re-arrange the equation  $z = \frac{x+1}{1-x}$  to get  $x = \frac{z-1}{z+1}$ ; finally as  $z > 1$ , then  $0 < x < 1$ . This allows us to calculate:

$$\tan^{-1}(z) = \frac{\pi}{4} + \tan^{-1}\left(\frac{z-1}{z+1}\right)$$

In the above the calculated value is in the range  $[0, 1]$  and so it is valid to use a Taylor polynomial using our McClaurin series above. This gives the following method

Algorithm 4.3.3: Taylor Method for  $\tan^{-1}$

```

1  taylor_aTan (x ∈ [0, 1], N ∈ ℕ)
2      T := 0
3      x2 := x2
4      y := x
5      n := 0
6      while n < N:
7          T ↦ T +  $\frac{1}{2n+1}$ y
8          y ↦ y · x2
9          T ↦ T -  $\frac{1}{2n+2}$ y
10         y ↦ y · x2
11         n ↦ n + 2
12     return T

```

Similar to Algorithm ?? the error of Algorithm 4.3.3 is easy to calculate. We see that

$$\begin{aligned}
 \epsilon_N &= |\tan^{-1}(x) - \text{taylor\_aTan}(x, N)| \\
 &\leq \frac{1}{2N+3} |x|^{2N+3} \\
 &\leq \frac{1}{2N+3} \quad \text{as } x \leq 1
 \end{aligned}$$

The next function we will consider is  $\sin^{-1}$ , which starts it's derivation in much the same way as  $\tan^{-1}$ . First we start by recalling that  $\frac{d}{dx} \sin^{-1}(x) = (1 - x^2)^{-\frac{1}{2}}$ , then by taking integrals of both sides we get the following derivation:

$$\begin{aligned}
\sin^{-1}(x) &= \int (1 - x^2)^{-\frac{1}{2}} dx \\
&= \int \sum_{n=0}^{\infty} \binom{-\frac{1}{2}}{n} (-x^2)^n \\
&= c + \sum_{n=0}^{\infty} (-1)^n \left( \prod_{k=1}^n \frac{-\frac{1}{2} - k + 1}{k} \right) \frac{x^{2n+1}}{2n+1} \\
&= c + \sum_{n=0}^{\infty} \frac{(-1)^n}{n!(2n+1)} \left( \prod_{k=1}^n \frac{\frac{1}{2} - k}{k} \right) x^{2n+1} \\
&= c + \sum_{n=0}^{\infty} \frac{(-1)^{2n}}{n!(2n+1)} \left( \prod_{k=1}^n \frac{2k-1}{2} \right) x^{2n+1} \\
&= c + \sum_{n=0}^{\infty} \frac{1}{n!(2n+1)2^n} \left( \prod_{k=1}^n 2k-1 \right) x^{2n+1} \\
&= c + \sum_{n=0}^{\infty} \frac{1}{n!(2n+1)2^n} (1 \times 3 \times 5 \times \cdots \times (2n-1)) x^{2n+1} \\
&= c + \sum_{n=0}^{\infty} \frac{1}{n!(2n+1)2^n} \times \frac{1 \times 2 \times 3 \times \cdots \times (2n)}{2 \times 4 \times \cdots \times (2n)} x^{2n+1} \\
&= c + \sum_{n=0}^{\infty} \frac{(2n)!}{(n!)^2 (2n+1) 4^n} x^{2n+1}
\end{aligned}$$

As  $\sin^{-1}(0) = 0$  then we see that  $c = 0$ . Because the above is valid for  $x \in (-1, 1)$ , and we know the values of  $\sin^{-1}(-1)$  and  $\sin^{-1}(1)$ , then we can have the following method for evaluating  $\sin^{-1}$ :

**Algorithm 4.3.4: Taylor Method for  $\sin^{-1}$**

```

1  taylor_aSin( $x \in [-1, 1], N \in \mathbb{N}$ )
2      if  $x = 1$ :
3          return  $\frac{\pi}{2}$ 
4      if  $x = -1$ :
5          return  $-\frac{\pi}{2}$ 
6       $\mathcal{S} := x$ 
7       $x_2 := x^2$ 
8       $y := x$ 
9       $a := 1$ 
10      $b := 1$ 
11      $c := 1$ 
12      $n := 1$ 
13     while  $n < N$ :
14          $a \mapsto 2n \cdot (2n-1) \cdot a$ 
15          $b \mapsto n^2 \cdot b$ 
16          $c \mapsto 4 \cdot c$ 
17          $y \mapsto x_2 \cdot y$ 
18          $\mathcal{S} \mapsto \mathcal{S} + \frac{a}{b \cdot c \cdot (2n+1)} \cdot y$ 
19          $n \mapsto n + 1$ 

```

The error for this method is similar to the  $\tan^{-1}$  method, in that  $\epsilon_N \leq \frac{(2(N+1))!}{((N+1)!(2N+1)4^{N+1})}$ . Finally we note that  $\cos^{-1}(x) = \frac{\pi}{2} - \sin^{-1}(x)$ , and thus can be calculated from a value calculated with Algorithm 4.3.4.

## 4.4 CORDIC

CORDIC is an algorithm that stands for COordinate Rotation DIgital Computer and can be used to calculate many functions, including Trigonometric Values. The CORDIC algorithm works by utilising Matrix Rotations of unit vectors. This algorithm is less accurate than some other methods but has the advantage of being able to be implemented for fixed point real numbers in efficient ways using only addition and bitshifting.

CORDIC works by taking an initial value of  $\mathbf{x}_0 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$  which can be rotated through an anti-clockwise angle of  $\gamma$  by the matrix

$$\begin{pmatrix} \cos \gamma & -\sin \gamma \\ \sin \gamma & \cos \gamma \end{pmatrix} = \frac{1}{\sqrt{1 + \tan^2 \gamma}} \begin{pmatrix} 1 & -\tan \gamma \\ \tan \gamma & 1 \end{pmatrix}$$

By taking taking smaller and smaller values of  $\gamma$  we can create an iterative process to find  $\mathbf{x}_n$  which converges, for a given  $\beta \in (-\frac{\pi}{2}, \frac{\pi}{2})$ , to

$$\begin{pmatrix} \cos \beta \\ \sin \beta \end{pmatrix}$$

To do this we repeatedly add and subtract our values for  $\gamma$  from  $\beta$  to bring it as close to 0 as possible. For our purposes we wish to have a sequence  $(\gamma_k : k \in [0, n] \cap \mathbb{Z})$  which will allow us to construct all angles in the range  $(-\frac{\pi}{2}, \frac{\pi}{2})$  to within a known level of accuracy. There are many possible choices here, but we wish to consider  $(\gamma_k : k \in [0, n] \cap \mathbb{Z})$  such that  $\tan \gamma_k = 2^{-k} \forall k \in [0, n] \cap \mathbb{Z}$ .

We can note that the powers of 2 have a useful property, in that if  $m > n \in \mathbb{N}$  we see that  $\sum_{k=n}^{m-1} 2^k = 2^m - 2^n$ . We wish to show that our choice for  $\gamma_k$  have a similar property which will be usefull in showing that they are a good choice for our CORIC algorithm.

**Proposition 4.4.1.** *If  $m \in \mathbb{Z}_0^+$  and  $n \in \mathbb{Z}^+$  such that  $m > n$  and  $\gamma_k = \tan^{-1}(2^{-k}) \forall k \in \mathbb{Z}_0^+$ , then  $\gamma_m < \gamma_n + \sum_{k=m+1}^n \gamma_k$ .*

*Proof.* We know that  $2^{-m} = 2^{-n} + \sum_{k=m+1}^n 2^{-k}$ , and thus by applying  $\tan^{-1}$  to both sides we get:

$$\tan^{-1} 2^{-m} = \gamma_m = \tan^{-1}(2^{-m-1} + 2^{-m-2} + \dots + 2^{-n} + 2^{-n})$$

Let  $a := 2^{-m-1} + 2^{-m-2} + \dots + 2^{-n} + 2^{-n}$  and  $b := 2^{-m-2} + \dots + 2^{-n} + 2^{-n}$ . Obviously  $a < b$  and further we know that  $\tan^{-1}$  is continuous on  $[a, b]$  and differentiable on  $(a, b)$ . Therefore we can apply the Mean Value Theorem from calculus to find that

$$\exists c \in (a, b) : \frac{1}{c^2 + 1} = \frac{\tan^{-1}(b) - \tan^{-1}(a)}{b - a}$$



By re-arranging we see that

$$\begin{aligned}\tan^{-1}(b) &= \frac{2^{-m-1}}{c^2 + 1} + \tan^{-1}(a) \\ &< \frac{2^{-m-1}}{2^{-2m-2} + 1} + \tan^{-1}(a)\end{aligned}$$

It can be shown, by considering the series expansion of  $\tan^{-1}(2^{-m-1})$ , that  $\frac{2^{-m-1}}{2^{-2m-2}+1} < \tan^{-1}(2^{-m-1}) \forall m \in \mathbb{Z}_0^+$ ; therefore we get that:

$$\tan^{-1}(b) < \tan^{-1}(2^{-m-1}) + \tan^{-1}(a)$$

Following this and using the assumed value of  $\gamma_{m+1}$ , we see that:

$$\gamma_m < \gamma_{m+1} + \tan^{-1}(2^{-m-2} + \dots + 2^{-n} + 2^{-n})$$

By repeating the above process we eventually see that:

$$\gamma_m < \sum_{k=m+1}^{n-1} \gamma_k + \tan^{-1}(2^{-n} + 2^{-n})$$

In a similar manner we can repeat the above process with  $a := \tan^{-1}(2^{-n})$  and  $b := \tan^{-1}(2^{-n} + 2^{-n})$ . This will show that:

$$\gamma_m < \gamma_n + \sum_{k=m+1}^n \gamma_k$$

□

Using the previous proposition we can then show that our  $\gamma_k$  have the property that every angle in  $(-\frac{\pi}{2}, \frac{\pi}{2})$  can be approximated by either adding or subtracting successive  $\gamma_k$  to within a tolerance of  $\gamma_n$ .

**Proposition 4.4.2.** *If  $\gamma_k = \tan^{-1}(2^{-k}) \forall k \in \mathbb{Z}$ , then for any  $n \in \mathbb{N}$*

$$\exists (c_k \in \{-1, 1\} : k \in [0, n] \cap \mathbb{Z}) : |\beta - \sum_{k=0}^n c_k \gamma_k| \leq \gamma_n \quad \forall \beta \in (-\frac{\pi}{2}, \frac{\pi}{2})$$

*Proof.* We let  $\beta \in (-\frac{\pi}{2}, \frac{\pi}{2})$  and then will proceed by induction on  $n \in \mathbb{N}$ .

$$\mathbf{H}(n): \exists (c_k \in \{-1, 1\} : k \in [0, n] \cap \mathbb{Z}) : |\beta - \sum_{k=0}^n c_k \gamma_k| \leq \gamma_n$$

$\mathbf{H}(0)$ : We have 4 cases to consider:

**Case  $\beta \in [0, \frac{\pi}{4}]$ :** In this case  $-\frac{\pi}{4} \leq \beta - \gamma_0 < 0$   
Therefore  $|\beta - \gamma_0| \leq \gamma_0$ .

**Case  $\beta \in [\frac{\pi}{4}, \frac{\pi}{2})$ :** In this case  $0 \leq \beta - \gamma_0 < \frac{\pi}{4}$   
Therefore  $|\beta - \gamma_0| \leq \gamma_0$ .

**Case  $\beta \in (-\frac{\pi}{4}, 0)$ :** In this case  $0 < \beta + \gamma_0 < \frac{\pi}{4}$

Therefore  $|\beta - \gamma_0| < \gamma_0$ .

**Case  $\beta \in (-\frac{\pi}{2}, -\frac{\pi}{4}]$ :** In this case  $-\frac{\pi}{4} < \beta - \gamma_0 \leq 0$

Therefore  $|\beta - \gamma_0| < \gamma_0$ .

Therefore we see that  $H(0)$  holds true.

**$H(n) \implies H(n+1)$ :**

By  $H(n) \exists (c_k \in -1, 1 : k \in [0, n] \cap \mathbb{Z}) : |\beta - \sum_{k=0}^n c_k \gamma_k| \leq \gamma_n$ ; so let  $\beta_n := \beta - \sum_{k=0}^n c_k \gamma_k$ .

By Proposition 4.4.1 we know that  $\gamma_n < 2\gamma_{n+1}$ , and so we can proceed by case analysis:

**Case  $\beta_n \in [0, \gamma_{n+1})$ :**

$-\gamma_{n+1} \leq \beta_n - \gamma_{n+1} < 0 \implies |\beta - \sum_{k=0}^{n+1} c_k \gamma_k| \leq \gamma_{n+1}$  where  $c_{n+1} = -1$ .

**Case  $\beta_n \in [\gamma_{n+1}, \gamma_n)$ :**

$0 \leq \beta_n - \gamma_{n+1} < \gamma_{n+1} \implies |\beta - \sum_{k=0}^{n+1} c_k \gamma_k| \leq \gamma_{n+1}$  where  $c_{n+1} = -1$ .

**Case  $\beta_n \in [-\gamma_{n+1}, 0)$ :**

$0 \leq \beta_n + \gamma_{n+1} < \gamma_{n+1} \implies |\beta - \sum_{k=0}^{n+1} c_k \gamma_k| \leq \gamma_{n+1}$  where  $c_{n+1} = 1$ .

**Case  $\beta_n \in (-\gamma_n, -\gamma_{n+1})$ :**

$-\gamma_{n+1} < \beta_n + \gamma_{n+1} < 0 \implies |\beta - \sum_{k=0}^{n+1} c_k \gamma_k| \leq \gamma_{n+1}$  where  $c_{n+1} = 1$ .

Therefore as we have found a suitable  $c_n$  in all cases then we have shown that  $H(n) \implies H(n+1)$ .  $\square$

With this proposition we see that our choice for  $\gamma_k$  is a good choice to use for the CORDIC algorithm as it covers the entire range of  $(-\frac{\pi}{2}, \frac{\pi}{2})$ .

Now, as stated before, the basis of our algorithm is to calculate  $\begin{pmatrix} \cos \beta \\ \sin \beta \end{pmatrix}$  by using rotations of a unit vector. By putting our values for  $\gamma_k$  into our rotation matrix we get the following:

$$\begin{pmatrix} \cos \gamma_k & -\sin \gamma_k \\ \sin \gamma_k & \cos \gamma_k \end{pmatrix} = \frac{1}{\sqrt{1+2^{-2k}}} \begin{pmatrix} 1 & -2^{-k} \\ 2^{-k} & 1 \end{pmatrix}$$

Then if we take a current estimate of  $\begin{pmatrix} \cos \beta \\ \sin \beta \end{pmatrix}$  at step  $k$  to be  $\begin{pmatrix} x_n \\ y_n \end{pmatrix}$ , we see that

$$\begin{pmatrix} \cos \gamma_k & -\sin \gamma_k \\ \sin \gamma_k & \cos \gamma_k \end{pmatrix} \begin{pmatrix} x_k \\ y_k \end{pmatrix} = \frac{1}{\sqrt{1+2^{-2k}}} \begin{pmatrix} x_k - 2^{-k} y_k \\ y_k + 2^{-k} x_k \end{pmatrix}$$

This gives a very simple formula for the update of  $x_k$  and  $y_k$ , which can be used as the basis of the CORDIC Algorithm.

As seen in our proof of Proposition 4.4.2, we can approximate our desire angle at step  $n$  by keeping a track of  $\beta_n := \beta - \sum_{k=0}^{n-1} c_k \gamma_k$ . At step  $n$  we then have  $\beta_{n+1} = \beta_n - \gamma_n$  if  $\beta_{n+1} \geq 0$ , and  $\beta_{n+1} = \beta_n + \gamma_n$  otherwise. This leads us to the general implementation of CORDIC for Trigonometric Functions:

#### Algorithm 4.4.1: General Cordic

1	$\text{CORIC}(\beta \in (-\frac{\pi}{2}, \frac{\pi}{2}), n \in \mathbb{N}) :$
2	$x := 1$

```

3 |     y := 0
4 |     k := 0
5 |     while k < n:
6 |         if β ≥ 0:
7 |             t := x
8 |             x ↦  $\frac{1}{\sqrt{1+2^{-2k}}}(x - 2^{-k}y)$ 
9 |             y ↦  $\frac{1}{\sqrt{1+2^{-2k}}}(y + 2^{-k}t)$ 
10 |            β ↦ β - tan-1(2-k)
11 |         else:
12 |             t := x
13 |             x ↦  $\frac{1}{\sqrt{1+2^{-2k}}}(x + 2^{-k}y)$ 
14 |             y ↦  $\frac{1}{\sqrt{1+2^{-2k}}}(y - 2^{-k}t)$ 
15 |            β ↦ β + tan-1(2-k)
16 |         k ↦ k + 1
17 |     return (x, y)T

```

There are few improvements we can make on the general algorithm, however if we start to consider implementations of the algorithm we can find several ways to make our algorithm more efficient.

First we consider the representation of our values in the program, and while in many of the previous algorithms a floating point `double` value, as described in Section ??, we will see here that we wish to use a fixed point representation. If we have a fixed point representation of our values, then we are using an  $N$  bit integer to represent the value in question, with a fixed number of bits set aside for the integer part and the remainder for the fractional part. In this case the process of addition, subtraction as well as multiplication and division by powers of 2 is the same as that for integers.

In particular as our values never exceed the range of  $(-2, 2)$ , then we can use  $N - 2$  bits of our  $N$  bit integer to be the fractional part; this gives us a maximum precision of  $2^{2-N}$ . Further as we are only performing multiplication and division by two, this operation can be performed by bitshifting the values, which is much quicker than actual integer multiplication.

Second we can precalculate all of the values needed for the algorithm to trade storage space for a reduction in computational complexity. The values which we need to pre-calculate are  $\gamma_k = \tan^{-1}(2^{-k})$  and  $\frac{1}{\sqrt{1+2^{-2k}}}$  for  $k \in [0, n) \cap \mathbb{Z}$ . The first thing to note about this is that instead of calculating the multiplication  $\frac{1}{\sqrt{1+2^{-2k}}}$  at each stage we can actually take this value out of the loops and pre-calculate  $\prod_{k=0}^n \frac{1}{\sqrt{1+2^{-2k}}}$  for  $k \in [0, n) \cap \mathbb{Z}$ . Using these precalculated products we can then replace  $x := 1$  with  $x := \prod_{k=0}^n \frac{1}{\sqrt{1+2^{-2k}}}$  in the initialisation stage.

Now to consider an actual implementation, suppose we are using the 16 bit integer `int16_t` to represent our values; which will have the leading two bits represent the integer part and the remaining 14 bits represent the fractional part. In this case the level of precision is  $2^{-14} = 0.00006103515625$  and further we can show that as  $\gamma_{14} = \tan^{-1}(2^{-14}) \approx 2^{-14}$ ; therefore the largest we will choose  $n := 14$  to ensure the maximum possible accuracy, without performing excessive calculations

This means we can simplify our algorithm further by calculating only  $\prod_{k=0}^{14} \frac{1}{\sqrt{1+2^{-2k}}}$  and  $\tan^{-1}(2^{-k}) \forall k \in [0, 14] \cap \mathbb{Z}$ . One further note is that these values then need to be converted to approximations in our 16 bit fixed point representation. The first value is:

$$\begin{aligned} \prod_{k=0}^{14} \frac{1}{\sqrt{1+2^{-2k}}} &= 0.60725293651701023412897124207973889082\dots \\ &\approx 00.10011011011101_2 \\ &= 26dd_{16} \end{aligned}$$

Below is a table of all the angles in the relevant formats

$\gamma_k$	Exact Form	Binary	Hexadecimal
$\gamma_0$	0.7853981633...	00.11001001000011 <sub>2</sub>	3243 <sub>16</sub>
$\gamma_1$	0.4636476090...	00.01110110101100 <sub>2</sub>	1dac <sub>16</sub>
$\gamma_2$	0.2449786631...	00.00111110101101 <sub>2</sub>	0fad <sub>16</sub>
$\gamma_3$	0.1243549945...	00.00011111110101 <sub>2</sub>	07f5 <sub>16</sub>
$\gamma_4$	0.0624188099...	00.00001111111110 <sub>2</sub>	03fe <sub>16</sub>
$\gamma_5$	0.0312398334...	00.00000111111111 <sub>2</sub>	01ff <sub>16</sub>
$\gamma_6$	0.0156237286...	00.00000100000000 <sub>2</sub>	0100 <sub>16</sub>
$\gamma_7$	0.0078123410...	00.00000010000000 <sub>2</sub>	0080 <sub>16</sub>
$\gamma_8$	0.0039062301...	00.00000001000000 <sub>2</sub>	0040 <sub>16</sub>
$\gamma_9$	0.0019531225...	00.00000000100000 <sub>2</sub>	0020 <sub>16</sub>
$\gamma_{10}$	0.0009765621...	00.00000000010000 <sub>2</sub>	0010 <sub>16</sub>
$\gamma_{11}$	0.0004882812...	00.00000000001000 <sub>2</sub>	0008 <sub>16</sub>
$\gamma_{12}$	0.0002441406...	00.00000000000100 <sub>2</sub>	0004 <sub>16</sub>
$\gamma_{13}$	0.0001220703...	00.00000000000010 <sub>2</sub>	0002 <sub>16</sub>
$\gamma_{14}$	0.0000610351...	00.00000000000001 <sub>2</sub>	0001 <sub>16</sub>

This allows us to then write the following method in C to calculate both  $\cos \beta$  and  $\sin \beta$ , provided  $\beta \in [-\frac{\pi}{2}, \frac{\pi}{2}]$  is given in 16 bit fixed point representation:

```
int16_t *cordic_16(int16_t beta)
{
    const int16_t GAMMA = {0x3243, 0x1dac, 0x0fad, 0x07f5, 0x03fe,
                           0x01ff, 0x0100, 0x0080, 0x0040, 0x0020,
                           0x0010, 0x0008, 0x0004, 0x0002, 0x0001};

    int16_t x = 0x26dd, y = 0x0000, t, result;

    for(int k = 0; k <= 14; ++k)
    {
        t = x;
        if(beta >= 0)
        {
            beta -= GAMMA[k];
            x = x - (y >> k);
            y = y + (t >> k);
        }
        else
    }
```

```

        {
            beta += GAMMA[k];
            x = x + (y >> k);
            y = y - (t >> k);
        }
    }

    //This line is required by C to allow the value to be returned
    result = malloc(2 * sizeof(int16_t));

    result[0] = x;
    result[1] = y;
    return result;
}

```

As can easily be seen in the algorithm the number of calculations each iteration is constant, and the number of iterations is fixed at 15. This means that the algorithm is an  $\mathcal{O}(1)$  algorithm, and guarantees an answer accurate to 4 decimal places as  $2^{-14} < 10^{-4}$ . Further as the only calculations are integer addition, subtraction and bitshifting this method executes extremely quickly.

Similar methods exist for other fixed length formats such as using `int32_t` or `int64_t`. To examine in more detail how the method converges we will consider an implementation using `int64_t`, which will be approximating  $\cos(0.5)$ . The code used is included in the Appendix ?? and can perform the calculations with  $n \leq 63$ . Below are some of the functions approximations for different values of  $n$ :

$n$	Output with bold accurate digits
1	<b>0.70710678118654757273731</b>
2	<b>0.94868329805051376801827</b>
3	<b>0.84366148773210747346951</b>
4	<b>0.90373783889353875853345</b>
5	<b>0.87527458786899225984257</b>
6	<b>0.88995346811933362385360</b>
...	...
19	<b>0.87758301847694786257392</b>
20	<b>0.87758210404530012649360</b>
21	<b>0.87758256126152311971111</b>
22	<b>0.87758278986933524468128</b>
...	...
53	<b>0.87758256189037275873943</b>
54	<b>0.87758256189037264771712</b>
55	<b>0.87758256189037275873943</b>
56	<b>0.87758256189037275873943</b>
...	...
63	<b>0.87758256189037275873943</b>

This table shows us several interesting features of the algorithm, the first being that while there are points at which a certain number of decimal places are guaranteed; before that point the number of decimal places of accuracy can vary, such as in the first few iterations. As we know that the error after  $n$  iterations is at most  $\gamma_n = \tan^{-1}(2^{-n})$ , then we can guarantee

that we have at least  $d$  decimal places of accuracy if we use at least  $\log_2(\cot(10^{-d}))$  iterations.

Second there are some values of  $n$  which have uncharacteristically close approximations of the actual value, such as the case when 21 iterations are used. This arises due to the algorithm finding a good approximation for  $\beta$ , but then successive numbers of iterations move away from this value, thus once more decreasing the number of decimal digits of accuracy.

Finally at the end of the table we see that from 55 iterations onwards, the results do not get any more accurate. It turns out this is due to the program converting the `int64_t` fixed point values into `double` values, which typically have a precision of around  $2^{-55}$ . If we instead modify the program to use a more precise floating point representation we see that the 53 to 56 section of the table becomes:

$n$	Output with bold accurate digits
53	<b>0.8775825618903727</b> 3965747
54	<b>0.877582561890372</b> 68653156
55	<b>0.87758256189037271</b> 298609
56	<b>0.8775825618903727</b> 2621336

This is much more inline with what we would expect to see from the known error of the algorithm.

Now another use of CORIC is to effectively run it in reverse, which will allow us to calculate the Inverse Trigonometric functions. To do this we will start by considering the method for calculating  $\tan^{-1}$ , and then use trigonometric identities to calculate both  $\cos^{-1}$  and  $\sin^{-1}$ .

To accomplish this we will be fixing some initial values for  $\sin \theta$  and  $\cos \theta$ , and then running the CORDIC algorithm to move the approximation of  $\sin \theta$  towards zero. In doing this we will effectively run our algorithm in reverse, and if we keep track of the angles that we rotate through we can find  $\tan^{-1}$ .

We know that  $\tan \theta = \frac{\sin \theta}{\cos \theta}$ , which means that if we have a current approximation  $\begin{pmatrix} x_k \\ y_k \end{pmatrix}$  then  $\frac{y_k}{x_k} \approx \tan \theta$ . Using this, if we have an input of  $\tan \theta = z$  then we can take our initial values to be  $x_0 := \frac{1}{z}$  and  $y_0 := z$ . This has the desired property that  $\frac{y_0}{x_0} = z$ , and if we have  $y_n$  tending to 0 then the angle we approximate in the process will be  $\theta$ .

If we again consider a 16 bit fixed point implementation for our algorithm we can implement it as follows:

```
int16_t *cordic_atan_16(int16_t z)
{
    const int16_t GAMMA = {0x3243, 0x1dac, 0x0fad, 0x07f5, 0x03fe,
                          0x01ff, 0x0100, 0x0080, 0x0040, 0x0020,
                          0x0010, 0x0008, 0x0004, 0x0002, 0x0001};

    int16_t x = 0x2000, y = z >> 1, t, theta;

    for(int k = 0; k <= 14; ++k)
    {
        t = x;
        if(y < 0)
```

```

        {
            theta -= GAMMA[k];
            x = x - (y >> k);
            y = y + (t >> k);
        }
        else
        {
            theta += GAMMA[k];
            x = x + (y >> k);
            y = y - (t >> k);
        }
    }

    return theta;
}

```

Similar to our considerations when dealing with the Taylor method of calculating  $\tan^{-1}$ , we need to ensure that the input value is not too large, and so can perform the same transformations to the value to ensure we are always calculating a value in the range  $[0, 1)$ . Using this we can then use the identities  $\sin^{-1}(z) = \tan^{-1}(\frac{z}{\sqrt{1-z^2}})$  and  $\cos^{-1}(z) = \tan^{-1}(\frac{\sqrt{1-z^2}}{z})$ .

Obviously there are basic exceptional values that need to be checked for, in particular  $\cos^{-1}(0) = \frac{\pi}{2}$ , and  $\sin^{-1}(\pm 1) = \pm \frac{\pi}{2}$ . If these values are checked before hand then we are never dividing by 0,  $z \in [-1, 1] \cap \mathbb{Z}$ , and thus we have a complete algorithm, that calculates the inverse Trigonometric Functions.

This method, like the CORDIC method for the regular Trigonometric Functions, has an approximation that is accurate to within  $\gamma_n$ . Thus for our 16 bit implementation, the output will be accurate to within an error of  $2^{-14} = 0.00006103515625$ , in particular guaranteeing at least 4 decimal places of accuracy. A final note is that the Inverse Trigonometric Functions, again much like the regular CORDIC algorithm, is an  $\mathcal{O}(1)$  algorithm with simple calculations, making the algorithm extremely efficient.

## 4.5 Comparrison of Methods

We have observed three different methods for calculating the Trigonometric Functions, as well as their inverses and so should compare their efficiency and accuracy properties.

First we will compare how quickly each algorithm approaches the correct value for different inputs of  $n$ , and using  $\theta = 0.5$ . The comparrison will use `double` values for computation, so that all three methods may be equally compared. The table below compares the convergence of  $\cos \theta$ , with the bold digits being the correct digits found:

$n$	geometric_Cos(0.5, $n$ )	taylor_Cos(0.5, $n$ )	CORDIC(0.5, $n$ )
1	<b>0.87</b> 6953125000000000	<b>1.</b> 000000000000000000	<b>0.707</b> 106781186547572
2	<b>0.877</b> 426177263259887	<b>0.877</b> 604166666666629	<b>0.948</b> 683298050513768
3	<b>0.8775</b> 43526076081437	<b>0.877</b> 604166666666629	<b>0.843</b> 661487732107473
4	<b>0.8775</b> 72806699400187	<b>0.87758256</b> 2158978118	<b>0.903</b> 737838893538758
5	<b>0.87758</b> 0123327654892	<b>0.87758256</b> 2158978118	<b>0.875</b> 274587868992259
6	<b>0.87758</b> 1952264380182	<b>0.87758256189037</b> 3424	<b>0.889</b> 953468119333623
7	<b>0.877582</b> 409484792491	<b>0.87758256189037</b> 3424	<b>0.882</b> 719918613777410
8	<b>0.8775825</b> 23789035007	<b>0.8775825618903727</b> 58	<b>0.879</b> 022003513595939
9	<b>0.8775825</b> 52365041901	<b>0.8775825618903727</b> 58	<b>0.877</b> 152884812089639
10	<b>0.8775825</b> 59509040183	<b>0.8775825618903727</b> 58	<b>0.878</b> 089122532394572

This table demonstrates that `taylor_Cos` has the fastest convergence, and also demonstrates the staggered increase in accuracy as each step of the algorithm calculates two updates to  $\cos \theta$ , and thus the output only gets more accurate every other value of  $n$ . The `geometric_Cos` method has the second best convergence, while the CORDIC algorithm lags behind, having inconsistent convergence as measured in correct digits.

Next we will note that all algorithms can guarantee 10 digits of accuracy in a fixed number of steps. In particular we can guarantee 10 digits of accuracy for `geometric_Cos` when  $n \geq 16$ , `taylor_Cos` when  $n \geq 8$  and CORDIC when  $n \geq 34$ . Using the lower bounds of each of these values for  $n$  we can directly compare the speed of the methods.

To compare the methods we will be testing 1000 random values in the range  $[0, \frac{\pi}{2})$  for which we will calculate the cosine of with each method 100000 times. This will then also be compared to the standard C implementation of the `cos` function, available in `math.h`. The results of my personal testing follow, where the given times are for individual values, not individual method execution times:

	geometric_cos	taylor_cos	cordic_cos	builtin_cos
Total time:	16.029s	7.937s	21.471s	0.243s
Average time:	0.016s	0.007s	0.021s	0.000s
Minimum time:	0.015s	0.007s	0.020s	0.000s
Maximum time:	0.022s	0.013s	0.030s	0.000s

These values show that the fastest algorithm that we have discussed is Algorithm ?? (`taylor_Cos`), while the slowest is the CORDIC algorithm. However all of our Algorithms are much less efficient than the built-in `cos` function of C. It turns out this discrepancy is due to inefficient implementation as the `cos` function also uses a Taylor approximation, but is implemented in a much lower-level method that optimises the execution of the code.

TODO: Ref the C code

TODO: [https://sourceware.org/git/?p=glibc.git;a=blob;f=sysdeps/ieee754/dbl-64/s\\_sin.c;hb=HEAD1281](https://sourceware.org/git/?p=glibc.git;a=blob;f=sysdeps/ieee754/dbl-64/s_sin.c;hb=HEAD1281)

Next we will compare our methods for the Inverse Trigonometric Functions, starting with how they converge to the correct value, as detailed in the following table:



$n$	geometric_aCos(0.5, $n$ )	taylor_aCos(0.5, $n$ )	CORDIC(0.5, $n$ )
1	<b>2.35</b> 1425307918200591	<b>2.27</b> 0796326794896735	<b>2.35</b> 6194490192344837
2	<b>2.34</b> 7503635391542609	<b>2.32</b> 7962993461563101	<b>1.89</b> 2546881191538687
3	<b>2.346</b> 521397812842746	<b>2.34</b> 0568243461563113	<b>2.13</b> 7525544318402914
4	<b>2.3462</b> 75724597314926	<b>2.344</b> 244774711563117	<b>2.26</b> 1880538865164602
5	<b>2.34621</b> 4299177873829	<b>2.345</b> 470795757570225	<b>2.32</b> 4299348861121661
6	<b>2.34619</b> 8942378459939	<b>2.345</b> 913166442261221	<b>2.35</b> 5539182291389810
7	<b>2.346195</b> 103149716576	<b>2.346</b> 081295659538934	<b>2.33</b> 9915453670913247
8	<b>2.3461941</b> 43336564508	<b>2.3461</b> 47594614218956	<b>2.34</b> 7727794731014228
9	<b>2.3461939</b> 03386887935	<b>2.34617</b> 4467628018511	<b>2.34</b> 3821564599047224
10	<b>2.3461938</b> 43452078375	<b>2.34618</b> 5594784405026	<b>2.34</b> 5774687115525836

Here we see for the inverse trigonometric functions the convergence speed has been altered with the Geometric method now having the fastest convergence, the Taylor Method converges much slower and the CORDIC method is more stable. One interesting behaviour that emerges for larger values of  $n$  in the geometric\_aCos is demonstrated in the following table:

$n$	geometric_aCos(0.5, $n$ )
13	<b>2.34619382</b> 2083380897
14	<b>2.34619381</b> 2716280469
15	<b>2.34619373</b> 7779483257
...	...
22	<b>2.3460</b> 97524754926944
23	<b>2.3412</b> 02123910687049
24	<b>2.3510</b> 23238547698124

This behaviour arises due to the use of `double` to calculate values of very small magnitude, this causes the value to become effectively 0 and thus lead to the inaccuracies seen. If we use a higher precision representation for the calculations we get the following table instead:

$n$	geometric_aCos(0.5, $n$ )
13	<b>2.346193823</b> 718087586
14	<b>2.3461938234</b> 83759158
15	<b>2.34619382342</b> 5177051
...	...
22	<b>2.3461938234056</b> 50874
23	<b>2.346193823405649</b> 980
24	<b>2.3461938234056497</b> 57

With this we see that Algorithm ?? continues in the same pattern as before and is actually correct. So we may again look to time our functions to test their efficiency as compared to each other. To do this we will again use 1000 random values, this time in the range  $(-1, 1)$ , each of which we will calculate  $\cos^{-1}$  using each method 100000 times. We note that the algorithms can guarantee 10 decimal places of accuracy for different values of  $n$ , in particular geometric\_aCos when  $n \geq 18$ , taylor\_aCos when  $n \geq 30$  and CORDIC when  $n \geq 50$ .

	geometric_cos	taylor_cos	cordic_cos	builtin_cos
Total time:	27.273s	14.358s	29.142s	2.143s
Average time:	0.027s	0.014s	0.029s	0.002s
Minimum time:	0.026s	0.014s	0.028s	0.001s
Maximum time:	0.033s	0.018s	0.032s	0.006s

Again this table shows that the Taylor method is the quickest of those analysed and the CORDIC method is the slowest, however they also both are much slower than the built-in methods. One thing to note is that the inverse trigonometric functions are simply less efficient to calculate, as can be seen in the execution time of the built-in method, which appears to be two orders of magnitude greater than the corresponding trigonometric method.

We conclude that for most implementations the Taylor method is the most appropriate method to use to ensure a high accuracy quickly. However the CORDIC algorithm is of use when more advanced features such as floating point type values, or hardware multipliers are not present; further it is possible to create hardware implementations of the CORDIC algorithm which can even further speed up the calculations.

## 5 Root Functions

In this section of the document we will consider several methods for approximating root functions. For our purposes we are only going to consider roots of  $N \in \mathbb{R}_0^+$ , this is because if  $N \in \mathbb{R}^-$  then it follows that  $\sqrt{N} = i\sqrt{|N|}$ .

### 5.1 Digit by Digit Method

The first method we will examine is an old method, that has been observed in Babylonian Mathematics over 2000 years ago, which is used to accurately generate the square root of numbers one digit at a time. This method differs from others discussed as it generates each digit of the root with perfect accuracy, one at a time, thus in a theoretical sense this algorithm is the most accurate of the methods we will view; we will see however that this method is slow.

Now suppose we are looking for  $\sqrt{N}$ , then we know that  $\sqrt{N} = a_0 10^n + a_1 10^{n-1} + a_2 10^{n-2} + \dots$  for some  $n \in \mathbb{Z}$ ; it then follows that  $N = (a_0 10^n + a_1 10^{n-1} + a_2 10^{n-2} + \dots)^2$ . By expanding the quadratic value we get that

$$N = a_0^2 10^{2n} + (2a_0 a_1) 10^{2n-2} + (2(a_0 a_2 + a_1^2)) 10^{2n-4} + \dots + (2 \sum_{i=0}^{k-1} a_i 10^{k-i-1} + a_k^2) 10^{2n-2k}$$

An observation should be made regarding the value of  $n$  that we use for the theorem. We could of course try different values of  $n$ , in some structured procedure, that will find the largest  $n$  such that  $10^n \leq N$ . However we can note that  $\log_{10}(\sqrt{N}) = \frac{1}{2} \log_{10}(N)$ , thus  $10^{\lfloor \frac{1}{2} \log_{10}(N) \rfloor} = \sqrt{N}$ . Using this information, and the fact that  $n \in \mathbb{Z}$ , we can have  $n := \lfloor \frac{1}{2} \log_{10}(N) \rfloor$ .

This allows us to get successive approximations of  $N$  where  $N_0 = a_0^2 10^{2n}$ ,  $N_1 = N_0 + (2a_0 a_1) 10^{2n-2}$ ,  $N_2 = N_1 + (2(a_0 a_2 + a_1^2)) 10^{2n-4}$ . This will allow us to create an algorithm that will give successive approximations of  $\sqrt{N} = a_0 10^n + a_1 10^{n-1} + \dots$ , more importantly each approximation will give us the exact next digit in the decimal representation of  $\sqrt{N}$ .

Thus we can have an iterative method to solve the problem, where at each stage we are trying to find the largest digit which satisfies the inequality  $(2 \sum_{i=0}^{k-1} a_i 10^{k-i-1} + a_k) a_k 10^{2n-2k} \leq N - N_{k-1}$ . Thus we get the following pseudocode, which outputs two sequences, one indicating the digits before the decimal point and one afterwards. I will use set notation to indicate the sequences, but in this case order is important and repetition is allowed.

Algorithm 5.1.1: Exact Digit by Digits Square Root

```

1  exactRootDigits( $N \in \mathbb{R}_0^+, d \in \mathbb{N}$ ):
2       $Digits_a := \emptyset$ 
3       $Digits_b := \emptyset$ 
4       $k := 0$ 
5       $n := \lfloor \frac{1}{2} \log_{10}(N) \rfloor$ 
6      while  $k < d$ :
7           $a_k := \max \left\{ t \in [0, 9] \cap \mathbb{Z} : \left( 2 \sum_{i=0}^{k-1} a_i 10^{k-i-1} + t \right) t 10^{2n-2k} \leq N \right\}$ 
8           $N \mapsto N - \left( 2 \sum_{i=0}^{k-1} a_i 10^{k-i-1} + a_k \right) a_k 10^{2n-2k}$ 
9          if  $n - k < 0$ :
```

```

10          $Digits_b \mapsto Digits_b \cup \{a_k\}$ 
11     else :
12          $Digits_a \mapsto Digits_a \cup \{a_k\}$ 
13      $k \mapsto k + 1$ 
14     if  $Digits_a = \emptyset$ :
15          $Digits_a := \{0\}$ 
16     if  $Digits_b = \emptyset$ :
17          $Digits_b := \{0\}$ 
18     return  $(Digits_a, Digits_b)$ 

```

This method has a computational complexity of  $\mathcal{O}(d^2)$ , as each loop requires the operations of summing  $k$  elements, and the loop is repeated for  $k = 0 \rightarrow d$ . We will see that by considering some changes to the algorithm we can change the complexity class to be  $\mathcal{O}(d)$ .

First we will note that line 5 is not an issue, as if we only care about the first significant digit of  $\frac{1}{2}\log_{10}(N)$ , then this is  $\mathcal{O}(|\log(N)|)$ . This can be seen as if we start from  $n = 0$  we can either count up or down until we find  $10^{2n}$  at most or at least  $N$ , respectively. This obviously takes at most  $|\log_{10}(N)|$  steps, giving us our stated complexity. We will also assume that  $\mathcal{O}(|\log(N)|) \leq \mathcal{O}(d)$ , as we have already seen that we can manipulate our input  $N$  to be within a reasonable range.

Second we note that on line 7 we calculate  $\sum_{i=0}^{k-1} a_i 10^{k-i-1}$  for each value of  $t$ ; we can reduce the complexity of this line by pre-calculating this value. However we can do even better if we consider that at step  $k + 1$  we are calculating  $\sum_{i=0}^k a_i 10^{k-i} = a_k + 10 \sum_{i=0}^{k-1} a_i 10^{k-i-1}$ . Thus if we introduce  $P_0 := 0$ , and for each  $k$  we calculate  $P_{k+1} := 10P_k + a_k$ , then we can reduce the complexity from  $\mathcal{O}(k)$  to  $\mathcal{O}(1)$ .

This calculation of  $P_k$ , then carries over to reduce the complexity of line 8 to be  $\mathcal{O}(1)$  instead of  $\mathcal{O}(k)$ . Combining this we can create the modified algorithm below:

#### Algorithm 5.1.2: Exact Digit by Digits Square Root version 2

```

1  exactRootDigits_v2 ( $N \in \mathbb{R}_0^+, d \in \mathbb{N}$ ):
2       $Digits_a := \emptyset$ 
3       $Digits_b := \emptyset$ 
4       $k := 0$ 
5       $n := \lfloor \frac{1}{2}\log_{10}(N) \rfloor$ 
6       $P_0 := 0$ 
7      while  $k < d$ :
8           $a_k := \max \{t \in [0, 9] \cap \mathbb{Z} : (20P_k + t) t 10^{2n-2k} \leq N\}$ 
9           $N \mapsto N - (20P_k + a_k) a_k 10^{2n-2k}$ 
10          $P_{k+1} := 10P_k + a_k$ 
11         if  $n - k < 0$ :
12              $Digits_b \mapsto Digits_b \cup \{a_k\}$ 
13         else :
14              $Digits_a \mapsto Digits_a \cup \{a_k\}$ 
15          $k \mapsto k + 1$ 
16     if  $Digits_a = \emptyset$ :
17          $Digits_a := \{0\}$ 
18     if  $Digits_b = \emptyset$ :
19          $Digits_b := \{0\}$ 

```

This method is usefull, but can be difficult to implement as it requires high precision for the representation of the real value of  $N$ . In my implementation using C, I utilised the MPFR library to utilise high precision integers, but still encountered issues regarding loss of precision.

As an example the table below shows the number of digits of accuracy I was able to calculate for  $\sqrt{2}$  using the above algorithm, compared to the number of bits of precision used in the calculations.

Bits of Precision	Maximum Accuracy
8	2
16	5
32	9
64	18
128	39
256	77
512	154
1024	308
2048	615
4096	1234
8192	2466

This data is highly structured and so we can hope to create a simple function that would allow us to calculate how much precision would be needed for a given number of digits of accuracy, at least for single digit inputs for  $N$ . We can see that the average ratio of Precision to Accuracy is 3.41259..., which ranges from 3.31928... to 4.0. From this we can draw a general trend that Digits of Accuracy  $\approx 3.4 \times$  Bits of Precision; thus if we take the more generous assumption that Digits of Accuracy  $4 \times$  Bits of Precision, we can use this to pre-determine the accuracy needed.

It should be noted that to ensure accuracy we should over-estimate the required precision, however if we overestimate the precision, then our calculations will be performed using unnecessarily large data structures and thus computation time will increase.

One particular use of this technique is to find an approximation of a squar root to it's integer part, calculated in base 2. This algorithm is of note as we will see that it has a computation time of (1).

The algorithm uses the same basis as the base 10 version, for it's calculations, but due to the nature of being in binary several changes can be made for computational efficiency. To do this we will view the problem as follows: if we know some  $r \in \mathbb{Z}_0^+$  which is our current approximation of our root, we are looking for some  $e \in \mathbb{Z}_0^+$  such that  $(r + e)^2 \leq N$ . Expanding this out we get  $r^2 + 2re + e^2 \leq N$ , and if we keep track of  $M = N - r^2$ , we can test if  $2re + e^2 \leq M$ .

Now we can consider our choice of  $e$ , the most practical method is to test succeasive  $e_m := 2^m$ , where  $m$  is descending starting with  $m = \max m \in \mathbb{Z}_0^+ : 4^m \leq N$ . We can use an iterative

formula to build up the integer square root, where we start with  $r = 0, M = N$  and have  $rr + e_m$  whenever  $2re_m + e_m^2 \leq M$ , stopping when  $m < 0$ . This is then implemented as follows:

#### Algorithm 5.1.3: Integer Square Root Algorithm

```

1  integerSquareRoot( $N \in \mathbb{Z}_0^+$ ):
2       $M := N$ 
3       $m := \max m \in \mathbb{Z}_0^+ : 4^m \leq M$ 
4       $r := 0$ 
5      while  $m \geq 0$ :
6          if  $2r(2^m) + 4^m \leq M$ :
7               $M \mapsto M - 2r(2^m) + 4^m$ 
8               $r \mapsto r + 2^m$ 
9               $m \mapsto m - 1$ 
10     return  $r$ 

```

If we now consider an implementation of the above algorithm using an unsigned integer system with  $K$  bits, where  $2|K$ . We will use `res` to represent  $2re_m$ , which means at the start of the algorithm we will have `res = 0`; similarly we can use `bit` to represent  $e_m^2$ . As we know that  $K$  bits are used and  $2|K$ , it then follows that the largest power of 4 less than the maximum representable value ( $2^K - 1$  is  $2^{K-2}$ , which can be calculated as `bit = 1 << (K - 2)` using bitshift operations. Finally we will use `num` to represent  $M$ .

Now that we have discussed the setup we can consider how to implement some of the steps above. First to implement line 3 we can simply keep dividing `bit` by 4 while `bit > num`, which can be efficiently implemented as `bit >> 2` by using bitshifts in place of division by powers of 2. The same technique can be used in place of line 9, which leads us to re-evaluating our usage of line 5. As we are using bitshifting and a bitshift that would take a number past 0 instead results in 0, we also know that  $2|K$  and so eventually we will reach `bit == 1`, which represents  $m = 0$ ; therefore we can use `bit > 0` as our stopping criteria on line 5.

Line 6 is easy to convert, given our definitions of `res`, `bit` and `num`, as is line 7. All that remains is to consider how to update `res`, which has two different ways of being updated depending on whether `res + bit <= num`. If it is false that `res + bit <= num`, then we wish for `res` to represent  $2re_{m-1}$ ; this is easily achieved if we consider that  $2re_{m-1} = \frac{1}{2}(2re_m)$ , which prompts the update `res = res >> 1`. For the second case, when `res + bit <= num` is true, we want `res` to represent  $2(r + e_m)e_{m-1}$ ; to implement this we consider the following derivation:

$$\begin{aligned}
 2(r + e_m)e_{m-1} &= \frac{1}{2} \cdot 2(r + e_m)e_m \\
 &= \frac{1}{2} \cdot 2(re_m + e_m^2) \\
 &= \frac{1}{2}(2re_m) + e_m^2
 \end{aligned}$$

Using this above derivation we see that we can calculate this as `res = (res >> 1) + bit`. Below is a simple implementation of this in C using the unsigned 32 bit integer type `uint32_t`. A more commented and slightly modified version can be found in Appendix ??, File ??.

```

|| uint32_t int_sqrt(uint32_t num)

```

```

{
    uint32_t res = 0, bit = (1 << 30);

    while (bit > num)
        bit = bit >> 2;

    while (bit > 0)
    {
        if (res + bit <= num)
        {
            num = num - (res + bit);
            res = (res >> 1) + bit;
        }
        else
            res = res >> 1;

        bit = bit >> 2;
    }

    return res;
}

```

We should consider the final step of the loop, when `bit == 1`. In this case when `res` is updated we have `res` represent either  $2(r+e_0)e_{-1} = r+e_0$ , or  $2re_{-1} = r$ ; thus the algorithm exits with the correct value.

Now that the algorithm is correctly constructed using simple unsigned integer addition, subtraction and bitshifting (which we can assume all have computational time of  $\mathcal{O}(1)$ ), we can look at the worst case complexity of the algorithm:

- The complexity of the set up of variables is constant time.
- The worst case complexity would be to have `bit <= num` at the start.
- The loop would execute 16 times for our 32 bit integers, and contains a single operation which is  $\mathcal{O}(1)$  complexity.
  - The worst case within the loop is to have `res + bit <= num` for each iteration.
  - Within the first `if` branch there are a constant 4 operations.
  - Each loop has an additional operation operation to update `bit`.
  - This makes 5 operations per loop, giving  $\mathcal{O}(1)$  complexity within the loops.

Therefore we see that the algorithm has  $\mathcal{O}(1)$  time complexity, and even has the same in storage complexity. In particular our 32 bit example requires 163 operations, including assignments, comparisons and calculations. This means that the integer square root of any number up to 4294967295 can be calculated extremely quickly.

## 5.2 Bisection Method

The Bisection Method is a general method for approximating the zero,  $\alpha$ , of a function,  $f$ , on a bounded interval,  $I := [a, b]$ , where  $f$  has the property  $f(x)f(y) < 0 \forall (x, y) \in [a, \alpha) \times (\alpha, b]$ ; we may assume, without loss of generality, that  $f(x) < 0 \forall x \in [a, \alpha]$ .

The bisection method starts with initial bounds  $a_0 = a, b_0 = b$ , where the initial approximation for the root is  $x_0 = \frac{1}{2}(a + b)$ . We will consider pseudocode of the iteration process, that uses  $b_n - a_n < \tau$  or  $f(x_n) = 0$  as exit criteria. Here  $\tau$  is a tolerance threshold, and if the exit criteria is met it means that  $|x_n - \alpha| \leq \frac{\tau}{2}$ , while the other exit criteria means we have reached an exact solution.

#### Algorithm 5.2.1: General Bisection Method

```

bisectionMethod ( $a \in \mathbb{R}, b \in (a, \infty), f \in \mathcal{C}[a, b], \tau \in \mathbb{R}^+$ )
   $a_0 := a$ 
   $b_0 := b$ 
   $x_0 := \frac{1}{2}(a + b)$ 
   $n := 0$ 
  while  $f(x_n) \neq 0$  AND  $b_n - a_n > \tau$ :
    if  $f(x_n) < 0$ :
       $a_{n+1} := x_n$ 
       $b_{n+1} := b_n$ 
    else:
       $a_{n+1} := a_n$ 
       $b_{n+1} := x_n$ 
     $n \mapsto n + 1$ 
     $x_n := \frac{1}{2}(a_n + b_n)$ 
  return  $x_n$ 

```

For our purposes we are trying to find the zero of  $f(x) = x^2 - N$ , which is a strictly increasing function on  $\mathbb{R}_0^+$ . If  $N \geq 1$ , then  $\sqrt{N} \in [0, N]$ , while  $N < 1 \implies \sqrt{N} \in [0, 1]$ . It is obvious that our function has the required property, and thus we get the following method for finding the square root of  $N$ :

#### Algorithm 5.2.2: Bisection Method for Square Roots

```

bisectionSquareRoot ( $N \in \mathbb{R}_0^+, \tau \in \mathbb{R}^+$ )
   $a_0 := 0$ 
   $b_0 := \max 1, N$ 
   $x_0 := \frac{1}{2}(a_0 + b_0)$ 
   $n := 0$ 
  while  $x_n^2 - N \neq 0$  AND  $b_n - a_n > \tau$ :
    if  $x_n^2 - N < 0$ :
       $a_{n+1} := x_n$ 
       $b_{n+1} := b_n$ 
    else:
       $a_{n+1} := a_n$ 
       $b_{n+1} := x_n$ 
     $n \mapsto n + 1$ 
     $x_n := \frac{1}{2}(a_n + b_n)$ 
  return  $x_n$ 

```

The implementation of this method is efficiently achieved in C using only addition, subtraction and multiplication by a constant. Before this method is implemented, however, we must first consider if and when it converges to the correct answer. From an intuitive standpoint we



would assume that if there is only one root in the interval, it would follow that we would converge to the root.

**Proposition 5.2.1.**  $\lim_{n \rightarrow \infty} x_n = \sqrt{N}$  for Algorithm 5.2.2

*Proof.* To prove this statement it suffices to prove that  $\sqrt{N} \in [a_n, b_n] \forall n \in \mathbb{N}$  and  $\lim_{n \rightarrow \infty} |x_n - \sqrt{N}| = 0$ .

*Claim 1:*  $\sqrt{N} \in [a_n, b_n] \forall n \in \mathbb{N}$

*Proof.*  $a_0 := 0 \implies a_0 \leq \sqrt{N}$   
 $b_0 := \max\{1, N\} \implies b_0 \geq \sqrt{N}$

Therefore it is obvious that  $\sqrt{N} \in [a_0, b_0]$

Now suppose  $\sqrt{N} \in [a_n, b_n]$  for some  $n \in \mathbb{N}$

It should be noted that  $a_n, b_n, x_n \in \mathbb{R}_0^+ \forall n \in \mathbb{N}$  as  $a_0, b_0 \in \mathbb{R}_0^+$  and all the subsequent values are derived from these using only addition and multiplication by positive factors.

We then see that  $x_n := \frac{1}{2}(a_n + b_n)$ , and we consider the two cases that  $x_n^2 - N \leq 0$  or  $x_n^2 - N \geq 0$ .

**Case**  $x_n^2 - N \leq 0$ :  $a_{n+1} := x_n, b_{n+1} := b_n$

It is therefore obvious that  $\sqrt{N} \leq b_{n+1}$ .

Now we see that  $x_n^2 - N \leq 0 \implies x_n^2 \leq N \implies x_n \leq \sqrt{N}$  as all the values are non-negative.

Thus  $\sqrt{N} \in [a_{n+1}, b_{n+1}]$ .

**Case**  $x_n^2 - N \geq 0$ :  $a_{n+1} := a_n, b_{n+1} := x_n$

It is therefore obvious that  $\sqrt{N} \geq a_{n+1}$ .

Now we see that  $x_n^2 - N \geq 0 \implies x_n^2 \geq N \implies x_n \geq \sqrt{N}$  as all the values are non-negative.

Thus  $\sqrt{N} \in [a_{n+1}, b_{n+1}]$ .

Hence  $\sqrt{N} \in [a_n, b_n] \implies \sqrt{N} \in [a_{n+1}, b_{n+1}] \forall n \in \mathbb{N}$

As  $\text{sqr}t N \in [a_0, b_0]$  then we see that  $\sqrt{N} \in [a_n, b_n] \forall n \in \mathbb{N}$  ■

*Claim 2:*  $\lim_{n \rightarrow \infty} |x_n - \sqrt{N}| = 0$

*Proof.* Let  $n \in \mathbb{N}$  be arbitrary.

As  $x_n := \frac{1}{2}(a_n + b_n)$  then we see that  $|a_n - x_n| = |b_n - x_n| = \frac{1}{2}(b_n - a_n)$ .

Now as  $\sqrt{N} \in [a_n, b_n]$  it follows that  $|\sqrt{N} - x_n| \leq \frac{1}{2}(b_n - a_n)$ .

As the modulus function is a mapping from  $\mathbb{R}$  to  $\mathbb{R}_0^+$ , it is clear that  $|\sqrt{N} - x_n|$  is bounded below by 0.

Now as for each  $n \in \mathbb{N}$ , either  $a_{n+1} = x_n$  or  $b_{n+1} = x_n$ , we see that  $b_{n+1} - a_{n+1} = \frac{1}{2}(b_n - a_n)$ .

Further we can see that  $b_n - a_n \geq 0 \forall n \in \mathbb{N}$  because  $b_n \geq a_n$ .

Therefore the sequence of  $\frac{1}{2}(b_n - a_n)$  is a strictly decreasing sequence that is bounded below, by 0. Thus  $\lim_{n \rightarrow \infty} \frac{1}{2}(b_n - a_n) = 0$

Therefore  $\lim_{n \rightarrow \infty} |x_n - \sqrt{N}| = \lim_{n \rightarrow \infty} \frac{1}{2}(b_n - a_n) = 0$  ■

By using our two claims above we see that  $\lim_{n \rightarrow \infty} x_n = \sqrt{N}$ . □

The algorithm can be generalised to search for  $\sqrt{k}N$ , where  $k \in [2, \infty) \cap \mathbb{Z}$ . We can do this by using the integer power function discussed previously in section ?? . This gives the following algorithm:

Algorithm 5.2.3: Bisection Method for General Roots

```

kRootBisectionMethod( $N \in \mathbb{R}_0^+, k \in [2, \infty) \cap \mathbb{Z}, \tau \in \mathbb{R}^+$ )
   $a_0 := 0$ 
   $b_0 := \max 1, N$ 
   $x_0 := \frac{1}{2}(a_0 + b_0)$ 
   $n := 0$ 
  while  $\text{intPow}(x_n, k) - N \neq 0$  AND  $b_n - a_n > \tau$ :
    if  $\text{intPow}(x_n, k) - N < 0$ :
       $a_{n+1} := x_n$ 
       $b_{n+1} := b_n$ 
    else:
       $a_{n+1} := a_n$ 
       $b_{n+1} := x_n$ 
     $n \mapsto n + 1$ 
     $x_n := \frac{1}{2}(a_n + b_n)$ 
  return  $x_n$ 

```

The proof that this converges to the correct value is very similar to the proof for square roots.

We can now consider the accuracy that can be achieved by our algorithm, for our purposes we will be considering  $\sqrt{N}$ , though the same applies for  $\sqrt{k}N$ . We know that  $\sqrt{N} \in [a_n, b_n] \forall n \in \mathbb{N}$ , and in particular we know that either  $\sqrt{N} \in [a_n, x_n]$  or  $\sqrt{N} \in [x_n, b_n] \forall n \in \mathbb{N}$ ; therefore we know that  $\epsilon_n := |x_n - \sqrt{N}| \leq \frac{1}{2}(b_n - a_n) \forall n \in \mathbb{N}$ . Then as we know that  $b_{n+1} - a_{n+1} = \frac{1}{2}(b_n - a_n)$ , we know that  $\epsilon_n \leq \frac{1}{2^n}(b_0 - a_0)$ .

We can consider that  $\forall N \in \mathbb{R}_0^+ \exists (r, k) \in [\frac{1}{4}, 1) \times \mathbb{Z} : N = r \cdot 2^{2k}$ ; using this we know that  $\sqrt{N} = \sqrt{r} \cdot 2^k$ . As we have the fixed initial bounds of  $a_0 = 0$  and  $b_0 = 1$ , then if we are finding  $\sqrt{r}$  we know that  $\epsilon_n \leq \frac{1}{2^n} \forall n \in \mathbb{N}$ . Hence we can calculate the precision of our current estimate beforehand for any  $n \in \mathbb{N}$ , and thus we can guarantee  $d$  significant digits of accuracy for  $r \in [\frac{1}{4}, 1)$ .

To get this accuracy must find  $n \in \mathbb{N}$  such that  $\epsilon_n \leq \frac{1}{10^d}$ , to achieve this we must find  $n \in \mathbb{N}$  such that  $2^n \geq 10^d$ . For example the following table indicates the required  $n$ , required for certain significant digits of accuracy.

$d$	$n : 2^n \geq 10^d$
1	0
5	15
10	30
20	64
50	163
100	329

Now usually finding  $r$  and  $k$  as above would be as hard as calculating the logarithm of  $N$ ; however due to the way that C stores real numbers as either `double` or in the MPFR

library, finding these values are actually fairly trivial. Both provide a functionality to find  $(a, b) \in [\frac{1}{2}, 1) \times \mathbb{Z} : N = a \cdot 2^b$ , and from this we merely require a simple comparison and division by 2 if  $b$  is not even. This leads to the following algorithm, which has the above maximum number of iterations for a required accuracy:

Algorithm 5.2.4: Bisection Method for Square Roots with fixed bounds

```

bisectionSquareRoot ( $N \in \mathbb{R}_0^+, \tau \in \mathbb{R}^+$ )
  Let  $(r, e) \in [\frac{1}{2}, 1) : N = r \cdot 2^e$ 
  if  $2 \nmid e$ :
     $r \mapsto \frac{r}{2}$ 
     $e \mapsto e - 1$ 
   $a_0 := 0$ 
   $b_0 := 1$ 
   $x_0 := \frac{1}{2}(a_0 + b_0)$ 
   $n := 0$ 
  while  $x_n^2 - N \neq 0$  AND  $b_n - a_n > \tau$ :
    if  $x_n^2 - N < 0$ :
       $a_{n+1} := x_n$ 
       $b_{n+1} := b_n$ 
    else:
       $a_{n+1} := a_n$ 
       $b_{n+1} := x_n$ 
     $n \mapsto n + 1$ 
     $x_n := \frac{1}{2}(a_n + b_n)$ 
  return  $x_n \cdot 2^{\frac{e}{2}}$ 

```

### 5.3 Newton's Method for Square Roots

If we consider  $f(x) = x^2 - N$  then if  $x^*$  is a solution to  $f(x) = 0$  we see that  $x^* = \sqrt{N}$ . As  $f'(x) = 2x$ , then the Newton's Method, will give  $x_{n+1} = x_n - \frac{x_n^2 - N}{2x_n}$ , where  $x_0$  is a given initial guess.

We can see that, in C, each iteration will calculate  $x = x - (x*x - N) / (2*x)$ , which requires 5 operations; however if we re-arrange our equation, we instead get  $x_{n+1} = \frac{1}{2}x_n + \frac{N}{x_n}$ . Implementing our new iterative formula we get  $x = 0.5 * (x + N/x)$ , which now uses only 3 operations.

We can then use the following pseudocode as the basis of our implementations of the Newton-Raphson Method for Square Roots:

Algorithm 5.3.1: Basic Newton Method for Square Root

```

NewtonSquareRoot ( $N \in \mathbb{R}, x_0 \in \mathbb{R}, \tau \in (0, 1)$ ):
   $n := 0$ 
  loop:
     $x_{n+1} := \frac{1}{2}(x_n + \frac{N}{x_n})$ 
     $\delta_n := |x_{n+1} - x_n|$ 
    if  $\delta_n \leq \tau$ :
      return  $x_{n+1}$ 

```

Next we want to consider our initial estimate  $x_0$ ; it is prudent to first consider when our initial estimate will converge to the correct root. By looking at a graph of the function, and in particular the tangents to the curve, it would seem reasonable to wonder if  $\lim_{n \rightarrow \infty} x_n = \sqrt{N}$ .

**Proposition 5.3.1.** *If  $x_0 \in \sqrt{N}, \infty$  and  $\{x_n : n \in \mathbb{N}\}$  is a sequence of approximations of  $\sqrt{N}$  found via the Newton-Raphson Method, as detailed above, then:*

$$\lim_{n \rightarrow \infty} x_n = \sqrt{N}$$

*Proof.* Suppose  $x_n > \sqrt{N}$ , then

$$\begin{aligned} x_{n+1} &= \frac{1}{2} \left( x_n + \frac{N}{x_n} \right) \\ &< \frac{1}{2} \left( x_n + \frac{N}{\sqrt{N}} \right) && \text{as } \sqrt{N} < x_n \implies \frac{1}{x_n} < \frac{1}{\sqrt{N}} \\ &= \frac{1}{2} (x_n + \sqrt{N}) \\ &< \frac{1}{2} (2x_n) \\ &= x_n \end{aligned}$$

Therefore we see that  $\{x_k : k \in [n, \infty) \cap \mathbb{Z}\}$  is a strictly decreasing sequence.

Now suppose that  $x_n \geq \sqrt{N}$  and then, for a contradiction, assume that  $x_{n+1} < \sqrt{N}$ . We then see that:

$$\begin{aligned} \frac{1}{2} \left( x_n + \frac{N}{x_n} \right) &< \sqrt{N} \\ \implies x_n + \frac{N}{x_n} &< 2\sqrt{N} \\ \implies x_n^2 + N &< 2\sqrt{N}x_n \\ \implies x_n^2 - 2\sqrt{N}x_n + N &< 0 \\ \implies (x_n - \sqrt{N})^2 &< 0 \end{aligned}$$

This is a contradiction as  $x_n, \sqrt{N} \in \mathbb{R} \implies (x_n - \sqrt{N})^2 \geq 0$ .

Therefore  $x_n \geq \sqrt{N} \implies x_{n+1} \geq \sqrt{N}$ .

Hence if  $x_0 > \sqrt{N}$ , then it follows that  $\{x_n : n \in \mathbb{N}\}$  is a strictly decreasing sequence that is bounded below. Therefore by an elementary result from limit theory, we see that  $\lim_{n \rightarrow \infty} x_n = \inf\{x_n : n \in \mathbb{N}\}$ .  $\square$

The most obvious choice for  $x_0$  would be  $N$ , but we see that  $N \in (0, 1)$ , then  $N < \sqrt{N}$ . In this case, we could choose  $x_0 = 1$  for the case that  $N \in (0, 1)$ . Therefore we can choose

$$x_0 := \begin{cases} N & : N \in (1, \infty) \\ 1 & : N \in (0, 1) \end{cases}$$

In our choice of  $x_0$ , we have so far left out the cases where  $N \in \{0, 1\}$ . In both of these cases we already know the correct answer, namely  $\sqrt{N} = N$  provided  $N \in \{0, 1\}$ . Therefore we can exclude them from our calculations, as we can pre-assess the value of  $N$ , simply returning the correct answer if one of these cases is encountered.

This then leads to an updated version of the above pseudocode:

**Algorithm 5.3.2: Basic Newton Method for Square Root**

```
NewtonSquareRoot ( $N \in \mathbb{R}_0^+, \tau \in (0, 1)$ ) :
```

```
  if  $N \in \{0, 1\}$ :
```

```
    return  $N$ 
```

```
  if  $N > 1$ :
```

```
     $x_0 := N$ 
```

```
  else:
```

```
     $x_0 := 1$ 
```

```
   $n := 0$ 
```

```
  loop:
```

```
     $x_{n+1} := \frac{1}{2}(x_n + \frac{N}{x_n})$ 
```

```
     $\delta_n := |x_{n+1} - x_n|$ 
```

```
    if  $\delta_n \leq \tau$ :
```

```
      return  $x_{n+1}$ 
```

```
     $n \mapsto n + 1$ 
```

**TODO: Write up examination and implementation of this pseudocode**

An alternative would be to use the integer square root method discussed in Section 5.1 to improve our initial choice of  $x_0$ . We will start by showing, that for intervals  $I \subset \mathbb{R}^+$ , the first two criteria for quadratic convergence of the Newton Raphson method are met.

**Proposition 5.3.2.** *If  $I \subset \mathbb{R}^+$  then  $NR_1$  and  $NR_2$  are satisfied for  $f(x) = x^2 - N$*

*Proof.*  $f(x) = x^2 - N \implies f'(x) = 2x \implies f''(x) = 2$

Now as  $x \in \mathbb{R}^+ \forall x \in I$ , then it is obvious that  $f'(x) > 0$

Therefore  $f'(x) \neq 0 \forall x \in I$ , and so  $NR_1$  is satisfied.

As  $f''(x)$  is a constant function, then it is continuous on all of  $\mathbb{R}$ .

Hence  $f''(x)$  is continuous  $\forall x \in I$  and so  $NR_2$  is satisfied. □

Now the integer square root function will always produce a root that is at most a distance of 1 from  $\sqrt{N}$ ; therefore we can consider  $I = [\sqrt{N} - 1, \sqrt{N} + 1]$ . Now if  $N \leq 1$ , then  $I \subset \mathbb{R}^+$  and so we cannot guarantee the satisfaction of  $NR_1$ . Therefore we can proceed with our analysis of the case that  $N > 1$ .

If  $N > 1$  we need to find when we can satisfy  $NR_3$ . First, we remember that  $M := \sup \left| \frac{f''(x)}{f'(x)} \right| : x \in I$  and  $\epsilon_0 := |x_0 - \sqrt{N}|$ . Then to satisfy  $NR_3$ , we must have that  $M\epsilon_0 < 1$ .

We can guarantee that  $\epsilon_0 \leq 1$  because  $x_0 \in I$  from the integer square root algorithm; therefore it suffices to find the situation where  $M < 1$ . As both  $f'$  and  $f''$  are continuous and non-zero

on  $I$  it follows that  $M = \sup x^{-1} : x \in I = (\sqrt{N} - 1)^{-1}$ . We then see that:

$$\begin{aligned} M < 1 &\iff \sqrt{N} - 1 > 1 \\ &\iff \sqrt{N} > 2 \\ &\iff N > 4 \end{aligned}$$

Therefore we can get the following new choice for  $x_0$ , and thus new pseudocode:

$$x_0 := \begin{cases} 1 & : N \in (0, 1) \\ N & : N \in (1, 4] \\ \text{intSqrt}(N) & : N \in (4, \infty) \end{cases}$$

#### Algorithm 5.3.3: Basic Newton Method for Square Root

NewtonSquareRoot( $N \in \mathbb{R}_0^+, \tau \in (0, 1)$ ):

```

  if  $N \in \{0, 1\}$ :
    return  $N$ 
  if  $N < 1$ :
     $x_0 := 1$ 
  else:
    if  $N \leq 4$ :
       $x_0 := N$ 
    else:
       $x_0 := \text{IntSqrt}(N)$ 
   $n := 0$ 
  loop:
     $x_{n+1} := \frac{1}{2}(x_n + \frac{N}{x_n})$ 
     $\delta_n := |x_{n+1} - x_n|$ 
    if  $\delta_n \leq \tau$ :
      return  $x_{n+1}$ 
     $n \mapsto n + 1$ 

```

**TODO:** Write up examination of different versions tried, such as using  $x_0 = N$ , etc...

If we consider any  $N \in \mathbb{R}_0^+$ , then  $\exists a \in [\frac{1}{2}, 1), b \in \mathbb{Z} : N = a \times 2^b$ . Finding this value would be as hard as finding the logarithm of  $N$  base 2, but due to the representation of numbers within C, both standard C and MPFR have functions that allow us to extract these two values with minimal computational expenditure.

This helps as we can then narrow our problem, to only finding  $\sqrt{a} : a \in [\frac{1}{2}, 1)$ , and then calculating

$$\sqrt{N} = \sqrt{a} \times 2^{\lfloor \frac{b}{2} \rfloor} \times \alpha \text{ where } \alpha = \begin{cases} 1 & : b \in 2\mathbb{Z} \\ \sqrt{2} & : b \in \mathbb{Z}^+ \setminus 2\mathbb{Z} \\ \frac{1}{\sqrt{2}} & : b \in \mathbb{Z}^- \setminus 2\mathbb{Z} \end{cases}$$

We then get the following algorithm, which implements this:

#### Algorithm 5.3.4: Newton Method for Square Root v3

NewtonSquareRoot( $N \in \mathbb{R}_0^+, \tau \in (0, 1)$ ):

Let  $(a, b) \in [\frac{1}{2}, 1) \times \mathbb{Z}$  s.t.  $N = a * 2^b$

```

 $x_0 := 1$ 
if  $b \equiv 0 \pmod{2}$ :
     $\alpha := 1$ 
else:
    if  $b > 0$ :
         $\alpha := \sqrt{2}$ 
    else:
         $\alpha := \frac{1}{\sqrt{2}}$ 
 $n := 0$ 
loop:
     $x_{n+1} := \frac{1}{2}(x_n + \frac{a}{x_n})$ 
     $\delta_n := |x_{n+1} - x_n|$ 
    if  $\delta_n \leq \tau$ :
        return  $\alpha \cdot x_{n+1} \cdot 2^{\lfloor \frac{b}{2} \rfloor}$ 
     $n \mapsto n + 1$ 

```

We must first consider the fact that the algorithm requires the pre-calculation of both  $\sqrt{2}$  and  $\frac{1}{\sqrt{2}}$ , to be able to calculate all values. However, it turns out we can use the algorithm itself to generate these values as  $2 = \frac{1}{2} \cdot 2^2$ , and as the exponent of 2 is even then the algorithm does not require  $\sqrt{2}$  for this computation. Similarly  $\frac{1}{2} = \frac{1}{2} \cdot 2^0$ , which again is an even exponent. We can thus run our algorithm to find an arbitrarily accurate values for  $\sqrt{2}$  and  $\frac{1}{\sqrt{2}}$  to allow us to run the algorithm for other values.

With this observation can then consider  $N \in [\frac{1}{2}, 1)$ . As this is a small range and, as per our previous algorithm, we use an initial guess of  $x_0 = 1$ , then we can prove that our algorithm will converge quadratically to  $\sqrt{N}$ .

**Proposition 5.3.3.** *Algorithm 5.3.4, satisfies the criteria of Theorem 2.3.2, and thus has quadratic convergence to  $\sqrt{N}$ .*

*Proof.* To fulfill the criteria of Theorem 2.3.2, we must find an interval  $I := [\sqrt{N} - r, \sqrt{N} + r]$  for some  $r \geq \epsilon_0$ .

Consider  $\epsilon_0 = |\sqrt{N} - x_0| = 1 - \sqrt{N}$ . We see that as  $N \geq \frac{1}{2}$  then  $\sqrt{N} \geq \sqrt{2}^{-1}$ , and thus  $\epsilon_0 \leq 1 - \sqrt{2}^{-1}$ . Let us have  $r := 1 - \frac{1}{\sqrt{2}}$ , and  $I$  as defined above.

If we look at the lower bound of  $I$ , then we see that:

$$\begin{aligned}
 \sqrt{N} - r &\geq \frac{1}{\sqrt{2}} - (1 - \frac{1}{\sqrt{2}}) \\
 &= \frac{2}{\sqrt{2}} - 1 \\
 &= \sqrt{2} - 1 \\
 &> 0
 \end{aligned}$$

Therefore we see that  $I \subset \mathbb{R}^+$ , and so by Proposition 5.3.2 we get that  $\text{NR}_1$  and  $\text{NR}_2$  are satisfied. It then remains to show that  $\text{NR}_3$  is satisfied on  $I$ .

Now by the definition in Theorem 2.3.2, we have that  $M = \sup \left\{ \frac{1}{2} \left| \frac{f''(x)}{f'(y)} \right| : x, y \in I \right\}$ . We know that  $I$  is bounded,  $f''(x) = 2$  and  $f'(x) = 2x$  meaning that  $\frac{1}{2} \left| \frac{f''(x)}{f'(y)} \right| = \frac{1}{f'(x)}$  as  $x \in \mathbb{R}^+$ .

Therefore our problem is reduced to finding  $\max \left\{ \frac{1}{2x} : x \in I \right\}$ , which is equivalent to finding  $\min \{x : x \in I\} = \sqrt{N} - r$ . Therefore by passing this information back up the chain we get that

$$M = \frac{1}{2(\sqrt{N} - r)}$$

Then we see that:

$$\begin{aligned} M_{\epsilon_0} &= \frac{1 - \sqrt{N}}{2(\sqrt{N} - r)} \\ &\leq \frac{1 - \frac{1}{\sqrt{2}}}{2(\sqrt{N} - r)} \quad \text{as } \sqrt{N} \geq \frac{1}{\sqrt{2}} \\ &\leq \frac{1 - \frac{1}{\sqrt{2}}}{2(\frac{1}{\sqrt{2}} - r)} \quad \text{as } \sqrt{N} \geq \frac{1}{\sqrt{2}} \\ &= \frac{1 - \frac{1}{\sqrt{2}}}{2(\frac{2}{\sqrt{2}} - 1)} \\ &= \frac{1 - \frac{1}{\sqrt{2}}}{2\sqrt{2}(1 - \frac{1}{\sqrt{2}})} \\ &= \frac{1}{2\sqrt{2}} \\ &< 1 \quad \text{as } 2\sqrt{2} > 1 \end{aligned}$$

As we have confirmed that  $M_{\epsilon_0} < 1$ , then we have confirmed that  $\text{NR}_3$  is satisfied on  $I$ , and so the algorithm converges quadratically to the desired root.  $\square$

Using the previous proposition we can, similar to our previous methods, consider how many iterations would be needed to reach a required tolerance. To start we consider that, as mentioned in the proof of Theorem 2.3.2, that  $\epsilon_n \leq (M_{\epsilon_0})^{2^{n-1}} \epsilon_0$ .

We know that  $M_{\epsilon_0} \leq \frac{1}{2\sqrt{2}}$  and that  $\epsilon_0 \leq 1 - \frac{1}{\sqrt{2}}$ , giving:

$$\epsilon_n \leq \left( \frac{1}{2\sqrt{2}} \right)^{2^{n-1}} \left( 1 - \frac{1}{\sqrt{2}} \right)$$

Thus if we want to achieve a tolerance of  $\epsilon_n \leq \tau$ , then it suffices to find  $n \in \mathbb{N}_0$  such that:

$$\left( \frac{1}{2\sqrt{2}} \right)^{2^{n-1}} \leq \tau$$

Then,

$$(2^n - 1) \log \left( \frac{1}{2\sqrt{2}} \right) \leq \log \left( \frac{\tau}{1 - \frac{1}{\sqrt{2}}} \right)$$



By noting that  $\log(\frac{1}{a}) = -\log(a)$ , then we get

$$(1 - 2^n) \log(2\sqrt{2}) \leq \log \left( \frac{\tau}{1 - \frac{1}{\sqrt{2}}} \right)$$

Once this is rearranged we get the following inequality:

$$2^n \geq \frac{\log \left( \frac{2(\sqrt{2}-1)}{\tau} \right)}{\log(2\sqrt{2})}$$

By taking logarithms again and re-arranging we get that

$$n \geq \frac{\log \left( \frac{\log \left( \frac{2(\sqrt{2}-1)}{\tau} \right)}{\log(2\sqrt{2})} \right)}{\log(2)} = \log_2 \left( \log_{2\sqrt{2}} \left( 2 \frac{\sqrt{2}-1}{\tau} \right) \right)$$

Now for an example, suppose we want to know how many iterations we need to perform to find  $\sqrt{N}$  to within 10 decimal places, i.e.  $\tau = 10^{-10} = 0.0000000001$ . We remember that  $\sqrt{N} \in [\frac{1}{2}, 1)$ , and then we will apply transformations to this value afterwards, therefore this is equivalent to finding 10 significant digits of accuracy for our square root (ignoring any loss of accuracy that may arise from multiplications afterwards).

Now in this case we want to find  $n \in \mathbb{N}$  such that  $n \geq \log_2(\log_{2\sqrt{2}}(2 \cdot 10^{10}(\sqrt{2}-1)))$ . Using Wolfram Alpha to calculate this value we get that we need  $n \geq 4.457144...$  and so we can take  $n = 5$ . This means that we could modify our algorithm and implementation to do 5 fixed iterations of Newton's Method to guarantee at least 10 decimal places of accuracy.

In terms of efficiency versus accuracy tradeoff modifying the problem thus would improve it's efficiency by removing, now unnecessary, calculation and comparison of  $\delta_n$  at each stage. However this does need a fixed guaranteed accuracy, and therefore such a program would no longer be suitable if we needed to calculate a square root accurate to 15 decimal places.

Below is a table that lists the minimum  $n \in \mathbb{N}$  such that  $n$  satisfies our inequality, where our tolerance is  $10^k$  for some  $k \in \mathbb{N}$ . This will give us the maximum number of iterations that must be performed for the required accuracy.

$k : \tau = 10^k$	$n$
5	4
10	5
100	8
1,000	12
1,000,000	22

## 5.4 Newton's Inverse Square Root Method

As discussed in Section ??, computers are more efficient at multiplication over division. We would therefore prefer to find a way of utilising Neton's Method without having to perform any costly division operations.

If we consider  $f(x) = N - \frac{1}{x^2}$  then if  $x^*$  is a solution to  $f(x) = 0$  we see that  $x^* = \frac{1}{\sqrt{N}}$ . As  $f'(x) = \frac{2}{x^3}$ , then the Newton's Method, will give

$$x_{n+1} = x_n - \frac{N - \frac{1}{x_n^2}}{\frac{2}{x_n^3}} = x_n \left( \frac{3}{2} - \frac{N}{2} x_n^2 \right)$$

where  $x_0$  is a given initial guess. As can be seen this algorithm requires no division if we multiply by real constants rather than the division implied above.

We can then consider that, similar to Algorithm 5.3.4, any  $N$  can be represented as  $a \cdot 2^b$  where  $a \in [\frac{1}{2}, 1)$ . This will, again allow us to narrow our problem to a known range of values, by using the following transformations.

$$\begin{aligned} N = a \cdot 2^b &\implies \frac{1}{N} = \frac{1}{a} \cdot 2^{-b} \\ &\implies \frac{1}{\sqrt{N}} = \frac{1}{a} \cdot 2^{\lfloor \frac{-b}{2} \rfloor} \cdot \alpha \quad \alpha := \begin{cases} 1 & : b \equiv 0 \pmod{2} \\ \sqrt{2} & : b \equiv 1 \pmod{2}, b \in \mathbb{Z}^- \\ \frac{1}{\sqrt{2}} & : b \equiv 1 \pmod{2}, b \in \mathbb{Z}^+ \end{cases} \\ &\implies \sqrt{N} = N \cdot \frac{1}{\sqrt{a}} \cdot 2^{\lfloor \frac{-b}{2} \rfloor} \cdot \alpha \end{aligned}$$

Therefore we only need to calculate inverse square roots for values of  $N$  in the range  $[\frac{1}{2}, 1)$ . Thus giving us the following algorithm:

#### Algorithm 5.4.1: Newton Inverse Square Root Method

```

NewtonInvSquareRoot( $N \in \mathbb{R}_0^+, \tau \in (0, 1)$ ):
  Let  $(a, b) \in [\frac{1}{2}, 1) \times \mathbb{Z}$  s.t.  $N = a \cdot 2^b$ 
   $x_0 := 1$ 
  if  $b \equiv 0 \pmod{2}$ :
     $\alpha := 1$ 
  else:
    if  $b > 0$ :
       $\alpha := \frac{1}{\sqrt{2}}$ 
    else:
       $\alpha := \sqrt{2}$ 
   $n := 0$ 
  loop:
     $x_{n+1} := x_n(\frac{3}{2} + \frac{a}{2}x_n^2)$ 
     $\delta_n := |x_{n+1} - x_n|$ 
    if  $\delta_n \leq \tau$ :
      return  $N \cdot \alpha \cdot x_{n+1} \cdot 2^{\lfloor \frac{-b}{2} \rfloor}$ 
     $n \mapsto n + 1$ 

```

With this method we can once again consider it's convergence properties, in particular does it satisfy the criteria for quadratic convergence in Theorem 2.3.2.

**Proposition 5.4.1.** *Algorithm 5.4.1 satisfies the criteria of Theorem 2.3.2, and thus has quadratic convergence to  $\sqrt{N}$ .*

*Proof.* We know that we only need to consider  $N \in [\frac{1}{2}, 1)$ , and therefore  $\sqrt{N^{-1}} \in (1, \sqrt{2}]$ . Also  $x_0 = 1$  and so we see that

$$\epsilon_0 = |x_0 - \sqrt{N^{-1}}| = \sqrt{N^{-1}} - x_0 \leq \sqrt{2} - 1$$

Now let  $r := \epsilon_0 = \sqrt{N} - 1$  and  $I := [\sqrt{N^{-1}} - r, \sqrt{N^{-1}}]$ . If we consider the lower bound of  $I$  we see that  $\sqrt{N^{-1}} - (\sqrt{N^{-1}} - 1) = 1$ , and in particular  $0 \notin I$ .

Next we know that  $f(x) = N - x^{-2}$ , and therefore we get  $f'(x) = 2x^{-3}$ ,  $f''(x) = -6x^{-4}$ . It is obvious that  $\nexists x \in \mathbb{R} : f'(x) = 0$ , which means that  $f'(x) \neq 0 \forall x \in I$  and so  $\text{NR}_1$  is satisfied. Also as  $f''$  is only discontinuous at  $x = 0$  and  $0 \notin I$ , then  $f''(x)$  is continuous  $\forall x \in I$ , meaning this satisfies  $\text{NR}_2$ .

Now  $M = \sup \left\{ \frac{1}{2} \left| \frac{2x^3}{6y^4} \right| : x, y \in I \right\}$ , we can simplify the function we are trying to minimise to get  $\frac{1}{6} \frac{x^3}{y^4}$ . It is obvious that in order to maximise this function we should find the largest possible  $x$  and smallest possible  $y$ , as both are positive. Hence by taking  $x = \sqrt{N^{-1}} + r$  and  $y = 1$ , then  $M = \frac{1}{6}(2\sqrt{N^{-1}} - 1)^3 \leq \frac{1}{6}(2\sqrt{2} - 1)^3$ .

Now we consider  $M_{\epsilon_0}$ :

$$\begin{aligned} M_{\epsilon_0} &= \frac{1}{6}(2\sqrt{N^{-1}} - 1)^3(\sqrt{N} - 1) \\ &\leq \frac{1}{6}(2\sqrt{2} - 1)^3(\sqrt{2} - 1) \\ &\approx 0.42199376 \dots \\ &< 1 \end{aligned}$$

Therefore as  $M_{\epsilon_0} < 1$  we have satisfied  $\text{NR}_3$ , and as such we have quadratic convergence of our method to  $\sqrt{N^{-1}}$ .  $\square$

## 5.5 Comparrison of Methods

We have observed several methods that can be used to calculate Square Roots, and so now we will see how the methods compare to each other in practice. The exact root method that we first discussed is the hardest to compare to the other methods as it works in a very different manner to the others. For now we will merely observe that it is an inefficient method that can be will be shown to tae longer than the others.

Second we need to compare the different methods discussed for the Newton Square Root method. As they work in the same general method, we really only need to test the computation time of the different methods. To do this we will be testing 1000 values in the range  $(0, 1000)$  and will calculate each of these values 100000 times, accurate to within a tolerance of  $10^{-1}$ , for each method to give the most accurate results. The table below gives the calculated results:

	mpfr_newton_sqrt_v1	mpfr_newton_sqrt_v2	mpfr_newton_sqrt_v3
Total time:	10.507s	12.707s	8.188s
Average time:	0.010s	0.012s	0.008s
Minimum time:	0.003s	0.004s	0.005s
Maximum time:	0.016s	0.021s	0.016s

Here we see that our third method, as expected, is the fastest of the proposed methods and so we will use this method going forwards. One unexpected result is that the second method is actually slower than the first, which is likely due to the extra conversions, comparisons and method calls; this slows down the execution more than it is sped up by reduction in number of iterations required.

Now for the comparison of methods we will be comparing modified versions of Algorithms 5.2.4, 5.3.4 and 5.4.1, which will execute for a given number of steps, rather than testing for the approximate error. To do this we need to consider how many iterations each method needs to reach a particular number of decimal places of accuracy.

We have seen the required number of iterations for a tolerance  $\tau = 10^{-k} : k \in \mathbb{N}$ , for both the bisection and basic newton square root methods, and similar to the basic newton method, we can show that for the inverse newton method we are looking for  $n \in \mathbb{N}$  that satisfies the following inequality:

$$n > \log_2 \left( \log_{\frac{1}{6}(\sqrt{2}-1)(2\sqrt{2}-1)^3} \left( \frac{\tau}{\sqrt{2}-1} \right) \right) - 1$$

This gives the following table:

$k : \tau = 10^k$	Bisection Method	Newton Method	Inverse Newton
5	16	4	4
10	33	5	5
100	332	8	9
1,000	3321	12	12
1,000,000	3219280	22	22

To show the above in action we have the table below which shows the convergence of all 3 methods to  $\sqrt{0.75} \approx 0.86602540378$ , for different numbers of iterations  $n$  with the bold digits being those correct:

$n$	bisectSquareRoot(0.75, $n$ )	NewtonSquareRoot(0.75, $n$ )	NewtonInvSquareRoot(0.75, $n$ )
0	<b>0</b> .50000000000000000000	<b>1</b> .00000000000000000000	<b>0</b> .75000000000000000000
1	<b>0</b> .75000000000000000000	<b>0</b> .87500000000000000000	<b>0</b> .84375000000000000000
2	<b>0</b> .87500000000000000000	<b>0</b> .866071428571428603	<b>0</b> .865173339843750000
3	<b>0</b> .81250000000000000000	<b>0</b> .866025405007363691	<b>0</b> .866024146705512976
4	<b>0</b> .84375000000000000000	<b>0</b> .866025403784438596	<b>0</b> .866025403781701674
5	<b>0</b> .85937500000000000000	<b>0</b> .866025403784438596	<b>0</b> .866025403784438596
6	<b>0</b> .86718750000000000000	<b>0</b> .866025403784438596	<b>0</b> .866025403784438596

If we compare the methods so that they guarantee an accuracy of 10 decimal places, then we will be able to see their relative efficiency. In particular we will again be testing the three

methods using 1000 values in the range (0, 1000), and calculating the square root of each of these values 10000 times for each method; further we will be including the the digit by digit method and the built-in C `sqrt` function. The results calculated are present in the following table:

	<code>root_digits_precise</code>	<code>bisect_sqrt</code>	<code>newton_sqrt</code>	<code>newton_inv_sqrt</code>	<code>builtin_sqrt</code>
Total time:	227.620s	2.520s	1.028s	0.646s	0.000s
Average time:	0.227s	0.002s	0.001s	0.000s	0.000s
Minimum time:	0.160s	0.002s	0.000s	0.000s	0.000s
Maximum time:	0.429s	0.004s	0.004s	0.001s	0.000s

Here we see the expected result that the digit by digit method is the least efficient method, taking two orders of magnitude more time than the second least efficient. We also see that while the two different newton methods are similar in time, and that even though they each performed the same number of iterations, the inverse square root method is the faster; this is due to the method having no division operations to perform. The quickest is of course the built-in `sqrt` function from C, this is due to an implementation that uses several low-level features of the C language to achieve the displayed level of performance.

In conclusion we can say that the best method that we have considered is Algorithm ?? which has rapid convergence to the sought square root, while also having fast execution. However if we are in a situation where we require large numbers of digits of accuracy, and yet do not have a suitable floating point types large enough to store these values, then the digit by digit method can be used to get an arbitrary number of digits accuracy.

## 6 Logarithms and Exponentials

Exponentiation is the operation of calculating  $x^y$  where  $x$  and  $y$  are members of some field, for the purposes of this document we will be considering  $x, y \in \mathbb{R}$ . This operation is widely used by many different branches of mathematics and industry, for example many real world phenomena can be modelled by exponentials; we would therefore like to be able to calculate  $x^y$  quickly and efficiently.

The first thing we consider is that  $x^y$  when  $x \in \mathbb{R}$  and  $y \in \mathbb{R} \setminus \mathbb{Z}$  is not well-defined on  $\mathbb{R}$ , and requires consideration of the function on the complex plane. Due to this we will not be considering negative numbers to non-integer bases; in particular, unless stated otherwise, we will be assuming that  $x \in \mathbb{R}_0^+$ .

Now we also know that  $x^{-y} = \frac{1}{x^y}$  when  $y \in \mathbb{R}$ , and as such we will also be restricting this section to the assumption that  $y \in \mathbb{R}_0^+$ . Further we consider the following facts:

$$x^0 = 1 \forall x \in \mathbb{R}_0^+$$

$$0^y = 0 \forall y \in \mathbb{R}^+$$

If we take out these known trivial cases then we can restrict this section to considering only  $(x, y) \in \mathbb{R}^{+2}$ .

Now if we have  $y \in \mathbb{R}^+$  then it follows that  $\exists(a, b) \in \mathbb{Z}_0^+ \times [0, 1)$  such that  $y = a + b$ . This allows us to use the identity that  $x^{m+n} = x^m x^n$  to consider the following two cases separately:

$$x^a : a \in \mathbb{Z}_0^+ \tag{6.0.1}$$

$$x^b : b \in [0, 1) \tag{6.0.2}$$

### 6.1 Calculating $x^a$

As we know that  $a \in \mathbb{Z}_0^+$ , then we know that  $x^a = \underbrace{x \times \dots \times x}_a$ ; i.e. the problem is equivalent

to finding  $x$  multiplied with itself  $a$  times. As we are only dealing with  $a \in \mathbb{Z}_0^+$ , then we will be considering  $x \in \mathbb{R}$  as we can calculate exponentials of negative numbers.

The naive way to go about calculating  $x^a$  is to simply perform the multiplication of  $x$  by itself  $a$  times. The algorithm for that can be seen below:

Algorithm 6.1.1: Naive integer exponentiation

```

1  naive_int_exp( $x \in \mathbb{R}, a \in \mathbb{Z}_0^+$ ):
2       $n := 0$ 
3       $z := 1$ 
4      while  $n < a$ :
5           $z \mapsto x \cdot z$ 
6      return  $z$ 
```

This algorithm is very simple and has complexity of  $\mathcal{O}(a)$ , which is a reasonable complexity, but still has the chance to grow large as  $a$  grows. Instead we can consider a more informed approach, in particular we know that either  $2 \mid a$  or  $2 \nmid a$ , which then gives us the following:

$$x^a = \begin{cases} (x^2)^{\frac{a}{2}} & : 2 \mid a \\ x \cdot (x^2)^{\frac{a-1}{2}} & : 2 \nmid a \end{cases}$$

We can use this fact to build a recursive method of calculating  $x^a$ , where we repeatedly call the method from within itself. To ensure the method ends correctly we need to identify a base case for the recursion, i.e. where the process stops and returns the correct value. We can see that eventually the above will reach the point where  $a = 0$ , in which case we know that  $x^0 = 1$ ; this will be the base case of our recursion.

We want to ensure that the algorithm will terminate, which we can do by seeing that it terminates when  $a = 0$  and then considering  $a \in \mathbb{Z}^+$ . Now if  $2 \mid a$  then  $\frac{a}{2} \in \mathbb{Z}^+$  and also  $\frac{a}{2} < a$ , similarly if  $2 \nmid a$  then  $\frac{a-1}{2} \in \mathbb{Z}_0^+$  because  $a \geq 1$  and also  $\frac{a-1}{2} < a$ . Thus we see that the sequence produced by  $a \in \mathbb{Z}^+$  is a strictly decreasing sequence that is bounded below by 0 and thus we must eventually reach 0, meaning the algorithm terminates.

Instead of a recursive algorithm that calls itself the algorithm below is an iterative version which performs the same function:

Algorithm 6.1.2: Exponentiation by squaring

```

1  exp_by_squaring( $x \in \mathbb{R}, a \in \mathbb{Z}_0^+$ ):
2       $n := a$ 
3       $z := 1$ 
4       $\hat{x} := x$ 
5      while  $n > 0$ :
6          if  $2 \nmid n$ :
7               $z \mapsto \hat{x} \times z$ 
8               $n \mapsto n - 1$ 
9               $\hat{x} \mapsto \hat{x}^2$ 
10              $n \mapsto \frac{n}{2}$ 
11     return  $z$ 

```

This algorithm is much more efficient than Algorithm 6.1.1 due to the number of times the inner loop is executed. The inner loop drives  $a$  towards 0 by dividing by 2 each step, this means that as  $a = \mathcal{O}(2^{\log_2(a)})$ , then this goal is achieved in only  $\log_2(a)$  loops. Therefore the complexity of this algorithm is  $\mathcal{O}(\log_2(a))$ , which is an improvement upon the previous algorithm's complexity of  $\mathcal{O}(a)$ .

To see this difference in efficiency in action the following table shows the times taken for each method when comparing 1000 different pairs of values  $(x, a) \in [0, 10] \times ([0, 100] \cap \mathbb{Z})$ . With these values we calculated  $x^a$  using both methods 100000 times to get the following results:

	naive_int_exp	squaring_int_exp
Total time:	16.800s	2.593s
Average time:	0.016s	0.002s
Minimum time:	0.000s	0.000s
Maximum time:	0.037s	0.004s

With this we will move on to further subsections as there are few improvements that can be made on an  $\mathcal{O}(\log_2(a))$  algorithm, particularly in this instance.

## 6.2 Calculating $x^b$

If we have  $b \in (0, 1)$ , then we obviously can't use the our previous subsubsection for calculating  $x^y$ . The most common way of calculating such exponentiation is by considering that  $x = e^{\ln(x)}$  and thus  $x^b = (e^{\ln(x)})^b = e^{b \ln(x)}$ ; however this now raises the problem of how to calculate both  $e^\alpha$  and  $\ln(\beta)$ . The following will deal with how to calculate these values and thus use them in conjunction to calculate  $x^b$ .

## 6.3 TBC

The mathematical constant  $e$  has been known since the early 1600s and was originally calculated by Jacob Bernoulli, and was studied by Leonhard Euler, where it appeared in Euler's *Mechanica* in 1736. While several possible equivalent definitions of  $e$  exist the most common such definition is that  $e := \lim_{n \rightarrow \infty} (1 + \frac{1}{n})^n$ .

If we now consider the definition of  $e$  and also consider  $e^x$ , then we can show that  $e^\alpha = \lim_{n \rightarrow \infty} (1 + \frac{x}{n})^n$ . This gives us our first basic method of how to calculate  $e^x$ :

Algorithm 6.3.1: Basic Method for calculating  $e^\alpha$

1  
2

```
basic_exp( $x \in \mathbb{R}, n \in \mathbb{N}$ )  
    return  $(1 + \frac{x}{n})^n$ 
```

If we consider  $(1 + \frac{x}{n})^n$  as a function of a continuous  $n$  then we can find the following derivation:

$$\begin{aligned} \frac{d}{dn} \left[ \left(1 + \frac{x}{n}\right)^n \right] &= \left(1 + \frac{x}{n}\right)^n \frac{d}{dn} \left[ n \ln \left(1 + \frac{x}{n}\right) \right] \\ &= \left(1 + \frac{x}{n}\right)^n \left( \frac{d}{dn} [n] \ln \left(1 + \frac{x}{n}\right) + n \frac{d}{dn} \left[ \ln \left(1 + \frac{x}{n}\right) \right] \right) \\ &= \left(1 + \frac{x}{n}\right)^n \left( \ln \left(1 + \frac{x}{n}\right) + \frac{n}{1 + \frac{x}{n}} \frac{d}{dn} \left(1 + \frac{x}{n}\right) \right) \\ &= \left(1 + \frac{x}{n}\right)^n \left( \ln \left(1 + \frac{x}{n}\right) - \frac{x}{n + x} \right) \\ &= \frac{\left(1 + \frac{x}{n}\right)^n}{x + n} \left( (x + n) \ln \left(1 + \frac{x}{n}\right) - x \right) \end{aligned}$$

By the last line of this we can see that because  $(x, n) \in \mathbb{R}^{+2}$  then  $\ln(1 + \frac{x}{n}) > 0$  and thus we conclude that  $(x + n) \ln(1 + \frac{x}{n}) - x > 0$ . Therefore we see that  $\frac{d}{dn} \left[ \left(1 + \frac{x}{n}\right)^n \right] > 0$  for all  $(x, n) \in \mathbb{R}^{+2}$ , and in particular this means that  $(1 + \frac{x}{n})^n < (1 + \frac{x}{n+1})^{n+1} \forall n \in \mathbb{N}$ .

One consequence of this is that  $(1 + \frac{x}{n})^n < e^x \forall n \in \mathbb{N}$ , therefore we can define the error of algorithm 6.3.1 as  $\epsilon_N := |e^x - (1 + \frac{x}{n})^n| = e^x - (1 + \frac{x}{n})^n$ . Now as  $\lim_{n \rightarrow \infty} (1 + \frac{x}{n})^n = e^x$  then we see that  $\lim_{n \rightarrow \infty} \epsilon_n = 0$ , and thus our algorithm is correct and valid for approximating  $e^x$ .

Next we see that this method, while simple, approximates  $e^x$  very poorly. In particular the table below shows the approximation of  $e^{0.75}$  for different values of  $n$ , where the bold digits are the correctly approximated digits.



$n$	Approximation of $e^{0.75}$
1	1.80000000000000000044
10	<b>2.158924997272786787</b>
100	<b>2.218468215957572747</b>
1000	<b>2.224829248807374831</b>
10000	<b>2.225469716120127850</b>
100000	<b>2.225533806810873500</b>
1000000	<b>2.225540216319864358</b>
10000000	<b>2.225540857275162929</b>
100000000	<b>2.225540921370736781</b>
1000000000	<b>2.225540927780294606</b>

With this table we see that the method very poorly approximates  $e^x$ , requiring a very large  $n$  to get just a few digits of accuracy. While this does not require more calculations from the method, requiring this large a value of  $n$  can lead to inaccuracies in the implementation of the algorithm using `double` data types in C.

In general there are better methods of approximating  $e^x$  and also  $\ln(x)$ , which while requiring more calculations are much more accurate than the most basic method presented here.

## 6.4 Taylor Series Method

If we take the elementary result from calculus that  $\frac{d}{dx}e^x = e^x$ , then we can calculate the McClaurin series of  $e^x$ . By the definition of a McClaurin series we know that the series expansion of  $e^x$  about 0 is

$$\sum_{k=0}^{\infty} \frac{d^{(k)}}{dx^k}[e^x](0) \frac{x^k}{k!}$$

As  $\frac{d^{(k)}}{dx^k}[e^x] = e^x \forall k \in \mathbb{Z}_0^+$  and  $e^0 = 1$  then we see that the series becomes

$$\sum_{k=0}^{\infty} \frac{x^k}{k!}$$

Using this we see that  $e^x \approx \sum_{k=0}^n \frac{x^k}{k!}$ , which gives the following method for approximating  $e^x$ :

Algorithm 6.4.1: Taylor Method for calculating  $e^x$

```

1  taylor_exp( $x \in \mathbb{R}, n \in \mathbb{Z}_0^+$ )
2       $t = 1$ 
3       $z = 1$ 
4       $k = 1$ 
5      while  $k < n$ :
6           $t \mapsto \frac{t \cdot x}{n}$ 
7           $z \mapsto z + t$ 
8           $k \mapsto k + 1$ 
9      return  $z$ 
```

This allows us to calculate  $e^x$  more efficiently as we can note that the error of the approximation is easy to approximate. We know that  $\epsilon_n := |e^x - \sum_{k=0}^n \frac{x^k}{k!}| \leq \frac{|x|^{n+1}}{(n+1)!}$  for all  $n \in \mathbb{Z}_0^+$ .

While we can't guarantee the size of  $x$  in general we will consider  $x \in (0, 1)$  for the purposes of analysing this function.

As  $x \in (0, 1)$  then it follows that  $x < 1$  and thus we can see that  $\epsilon_n < \frac{1}{n!} \forall n \in \mathbb{Z}_0^+$ . Using this we can see that if we wish to use our method such that the error is at most  $\tau_d := 10^{-d}$ , then we need to find  $n \in \mathbb{Z}_0^+ : \frac{1}{n!} < \tau_d$ . The table below shows some values for  $(n, d)$  pairs such that  $n$  is the smallest positive integer such that  $\frac{1}{n!} < \tau_d$ :

$d \in \mathbb{N}$	$\arg \min \{n \in \mathbb{N} : n! > 10^d\}$
1	4
10	14
100	70
1000	450

Thus we see that that we can guarantee 100 digits of accuracy with an input of  $n \geq 70$  and 1000 digits of accuracy with  $n \geq 450$ , this is much less than our previous method where an input of  $n = 1000$  only gave 2 decimal places of accuracy.

The inverse of the function  $z = e^x$  is the logarithm function  $\ln(z) = x$ , which we can again consider for Taylor Series expansion. First we will show the result from calculus that  $\frac{d}{dx}[\ln(x)] = \frac{1}{x}$ :

**Proposition 6.4.1.**

$$\frac{d}{dx}[\ln(x)] = \frac{1}{x}$$

*Proof.* We will prove this from the first principles using the definition that  $\frac{d}{dx}[f(x)] = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$

$$\begin{aligned} \frac{d}{dx}[\ln(x)] &= \lim_{h \rightarrow 0} \frac{\ln(x+h) - \ln(x)}{h} \\ &= \lim_{h \rightarrow 0} \frac{\ln(1 + \frac{h}{x})}{h} \\ &= \lim_{h \rightarrow 0} \ln \left( \left(1 + \frac{h}{x}\right)^{\frac{1}{h}} \right) \end{aligned}$$

If we let  $u := \frac{h}{x}$ , then we get that  $ux = h$  and  $\frac{1}{h} = \frac{1}{ux}$ . Also  $\lim_{h \rightarrow 0} u = 0$ , and so we get the following:

$$\begin{aligned} \frac{d}{dx}[\ln(x)] &= \lim_{u \rightarrow 0} \ln((1+u)^{\frac{1}{ux}}) \\ &= \frac{1}{x} \lim_{u \rightarrow 0} \ln((1+u)^{1/u}) \end{aligned}$$

If we now let  $n := \frac{1}{u}$  and consider that  $\lim_{u \rightarrow 0} n = \infty$ , then our derivative becomes:

$$\begin{aligned}
\frac{d}{dx} \ln(x) &= \frac{1}{x} \lim_{n \rightarrow \infty} \ln\left(1 + \frac{1}{n}\right)^n \\
&= \frac{1}{x} \ln\left(\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n\right) \\
&= \frac{1}{x} \ln(e) && \text{by the definition of } e \\
&= \frac{1}{x}
\end{aligned}$$

□

Now we know that  $\frac{d^k}{dx^k} \left[\frac{1}{x}\right] = (-1)^k k! x^{-k-1}$ , and thus we can build up a Taylor Series expansion. In this case, rather than centering the series about  $x = 0$  for a McClaurin series we can instead center the series around  $x = 1$  which gives the following series expansion for  $\ln(x)$ :

$$\begin{aligned}
\sum_{k=0}^{\infty} \frac{\frac{d^k}{dx^k} [\ln(x)](1)}{k!} (x-1)^k &= \ln(1) + \sum_{k=1}^{\infty} \frac{\frac{d^{k-1}}{dx^{k-1}} [x^{-1}](1)}{k!} (x-1)^k \\
&= \sum_{k=1}^{\infty} \frac{[(-1)^{k-1} (k-1)! x^{-k}](1)}{k!} (x-1)^k \\
&= \sum_{k=1}^{\infty} \frac{(-1)^{k-1}}{k} (x-1)^k \\
&= - \sum_{k=1}^{\infty} \frac{(1-x)^k}{k}
\end{aligned}$$

We know that  $\ln(x) = - \sum_{k=1}^{\infty} \frac{(1-x)^k}{k}$  when the series  $\sum_{k=1}^{\infty} \frac{(1-x)^k}{k}$  converges. We thus need to know when the sum converges.

**Proposition 6.4.2.** *The series  $\sum_{k=1}^{\infty} \frac{(1-x)^k}{k}$  converges when  $x \in (0, 2)$ .*

*Proof.* Let  $a_k := \frac{(1-x)^k}{k}$ . We will proceed by using the ratio test to show when the series converges absolutely. The test states that the series converges when  $\lim_{k \rightarrow \infty} \left| \frac{a_{k+1}}{a_k} \right| < 1$ .

Now we can consider the following derivation:

$$\begin{aligned}
\lim_{k \rightarrow \infty} \left| \frac{a_{k+1}}{a_k} \right| &= \lim_{k \rightarrow \infty} \left| \frac{\frac{1}{k+1} (1-x)^{k+1}}{\frac{1}{k} (1-x)^k} \right| \\
&= \lim_{k \rightarrow \infty} \left| \frac{k}{k+1} (1-x) \right| \\
&= |1-x| \lim_{k \rightarrow \infty} \left| \frac{k}{k+1} \right| \\
&= |1-x|
\end{aligned}$$

Therefore our series converges when:

$$\begin{aligned} |1 - x| < 1 &\iff -1 < 1 - x < 1 \\ &\iff -1 < x - 1 < 1 \\ &\iff 0 < x < 2 \end{aligned}$$

Hence  $\sum_{k=1}^{\infty} \frac{(1-x)^k}{k}$  converges when  $x \in (0, 2)$ . □

Now as we can't know if  $x \in (0, 2)$  then we can consider that  $\forall x \in \mathbb{R}^+ \exists (a, b) \in [\frac{1}{2}, 1) \times \mathbb{Z} : x = a \cdot 2^b$ ; thus we see that  $\ln(x) = \ln(a \cdot 2^b) = b \ln(2) + \ln(a)$ . As previously noted in Section ?? this operation, while theoretically complex, is simple to calculate for most computers by how they represent floating point values.

With this we can then use the following method to approximate  $\ln(x)$  by the Taylor polynomial  $-\sum_{k=1}^n \frac{(1-x)^k}{k}$ :

Algorithm 6.4.2: Taylor Method for calculating  $\ln(x)$

```

1  taylor_nat_log( $x \in \mathbb{R}^+, n \in \mathbb{N}$ ):
2      Find  $(a, b) \in [\frac{1}{2}, 1) \times \mathbb{Z}$  such that  $x = a \cdot 2^b$ 
3       $y := 1 - a$ 
4       $t := y$ 
5       $z := y$ 
6       $k := 1$ 
7      while  $k < n$ :
8           $t \mapsto t \cdot y$ 
9           $z \mapsto z + \frac{t}{k}$ 
10          $k \mapsto k + 1$ 
11     return  $b \ln(2) - z$ 

```

The first thing to consider for the above method is how to calculate  $\ln(2)$ . It is not possible to directly calculate  $\ln(2)$  using the above algorithm as  $2 = \frac{1}{2} \cdot 2^2$ , however  $\frac{1}{2} = \frac{1}{2} \cdot 2^0$  and so we do not need to know  $\ln(2)$  to calculate  $\ln(\frac{1}{2})$ . We can see that  $\ln(2) = -\ln(\frac{1}{2})$ , and so we can calculate our constant value  $\ln(2)$  to be used in the algorithm by using the algorithm itself.

Now similar to previous Taylor approximations the final error of our approximation using the above method is  $\epsilon_n := |\ln(x) - \text{taylor\_log}(x, n)|$ . As the next term of the approximation would be  $\frac{(1-x)^n}{n}$ , then we know that  $\epsilon_n \leq \left| \frac{(1-a)^n}{n} \right|$ ; further we know that  $a \in [\frac{1}{2}, 1)$  and thus  $\epsilon_n < \frac{1}{2^n n}$ .

Using this approximation we can see that if we wish to guarantee  $d$  decimal places of accuracy then it suffices to find  $n \in \mathbb{N}$  such that  $\frac{1}{2^n n} < 10^{-d} \implies 2^n n > 10^d$ . As  $n \in \mathbb{N}$  then  $2^n < 2^n n$  and so we merely need to find  $n \in \mathbb{N}$  such that  $2^n > 10^d$  to guarantee  $d$  decimal places of accuracy. Some example values are included in the table below:

$d \in \mathbb{N}$	$\arg \min \{n \in \mathbb{N} : 2^n > 10^d\}$
1	4
10	34
100	333
1000	3322

As we now have Taylor methods for approximating both  $e^x$  and  $\ln(x)$ , then we can use the two to derive a Taylor method of calculating  $x^y$  and  $\log_x(y)$ . To start we will consider  $x^y = e^{y \ln(x)}$  and  $x = a \cdot 2^b$ , giving the solution as  $x^y = e^{y(b \ln(2) + \ln(a))}$ . Similarly we note that  $\log_x(y) = \frac{\ln(y)}{\ln(x)}$ , and if we consider that  $x = a \cdot 2^b$  and  $y = c \cdot 2^d$ , then we see that  $\log_x(y) = \frac{d \ln(2) + \ln(c)}{b \ln(2) + \ln(a)}$ . Below are the Taylor methods for approximating these functions:

Algorithm 6.4.3: Taylor Method for calculating  $x^y$  and  $\log_x(y)$

```

1  taylor_log( $x \in \mathbb{R}^+, y \in \mathbb{R}^+, n \in \mathbb{N}$ ):
2       $a := \text{taylor\_nat\_log}(y, n)$ 
3       $b := \text{taylor\_nat\_log}(x, n)$ 
4      return  $\frac{a}{b}$ 
5
6  taylor_pow( $x \in \mathbb{R}^+, y \in \mathbb{R}, n \in \mathbb{N}$ ):
7       $a := \text{taylor\_nat\_log}(x, n)$ 
8       $a \mapsto y \cdot a$ 
9      return taylor_exp( $a, n$ )

```

To test the convergence of the Taylor methods above we are going to test calculations of  $7.3^{4.8}$ ,  $7.3^{-4.8}$ ,  $0.21^{4.8}$ ,  $7.3^{0.21}$ ,  $\log_{7.3}(4.8)$ ,  $\log_{0.21}(4.8)$  and  $\log_{7.3}(0.21)$ . These values are calculated for several different values of  $n$  with the bold digits representing the correct values in the tables below:

$n$	$7.3^{4.8}$	$7.3^{-4.8}$	$0.21^{4.8}$	$7.3^{0.21}$
1	<b>1.0000000000</b>	<b>1.0000000000</b>	<b>1.0000000000</b>	<b>1.0000000000</b>
2	<b>10.561319400</b>	<b>-8.561319400</b>	<b>-6.422212933</b>	<b>1.4183077237</b>
3	<b>56.076838311</b>	<b>36.990949511</b>	<b>21.518877680</b>	<b>1.5046585363</b>
4	<b>200.85920964</b>	<b>-107.8118783</b>	<b>-48.47602784</b>	<b>1.5167171202</b>
5	<b>546.24576990</b>	<b>237.58122696</b>	<b>82.710783892</b>	<b>1.5179778747</b>
6	<b>1205.3726532</b>	<b>-421.5471761</b>	<b>-113.8668463</b>	<b>1.5180831956</b>
7	<b>2253.5829747</b>	<b>626.66342673</b>	<b>131.57101558</b>	<b>1.5180905223</b>
8	<b>3682.4131809</b>	<b>-802.1668232</b>	<b>-131.0877429</b>	<b>1.5180909589</b>
9	<b>5386.6141612</b>	<b>902.03416303</b>	<b>114.86315726</b>	<b>1.5180909816</b>
10	<b>7193.4074522</b>	<b>-904.7591286</b>	<b>-89.85299062</b>	<b>1.5180909827</b>
...	...	...	...	...
20	<b>13901.238666</b>	<b>-11.00988984</b>	<b>-0.092958315</b>	<b>1.5180909827</b>
...	...	...	...	...
40	<b>13929.955484</b>	<b>0.0000717862</b>	<b>0.0005580236</b>	<b>1.5180909827</b>
...	...	...	...	...
80	<b>13929.955484</b>	<b>0.0000717877</b>	<b>0.0005580236</b>	<b>1.5180909827</b>

As we can see in the table the `taylor_pow` does not converge perfectly, and may even diverge from the correct value for small values of  $n$ ; however we see that the methods do converge for large values of  $n$ . This behaviour is due to the values being outside the restrictions used in the analysis of the functions.

$n$	$\log_{7.3}(4.8)$	$\log_{0.21}(4.8)$	$\log_{7.3}(0.21)$
1	<b>0.8431178860</b>	<b>-1.086107266</b>	<b>-0.776274970</b>
2	<b>0.8431178860</b>	<b>-1.086107266</b>	<b>-0.776274970</b>
3	<b>0.8045021618</b>	<b>-1.025878600</b>	<b>-0.784207957</b>
4	<b>0.7938608884</b>	<b>-1.011309817</b>	<b>-0.784982875</b>
5	<b>0.7906472231</b>	<b>-1.007102721</b>	<b>-0.785071082</b>
6	<b>0.7896173849</b>	<b>-1.005776909</b>	<b>-0.785082036</b>
7	<b>0.7892739993</b>	<b>-1.005337682</b>	<b>-0.785083473</b>
8	<b>0.7891562591</b>	<b>-1.005187460</b>	<b>-0.785083668</b>
9	<b>0.7891150494</b>	<b>-1.005134935</b>	<b>-0.785083695</b>
10	<b>0.7891003970</b>	<b>-1.005116266</b>	<b>-0.785083699</b>
...	...	...	...
50	<b>0.7890920869</b>	<b>-1.005105681</b>	<b>-0.785083699</b>

This shows that `taylor.log` converges better than `taylor.exp`, however part of this is due to the values tested having magnitudes close to 1. Answers with a larger or smaller magnitudes tend to converge slower, which can be seen in the table for `taylor.exp` where the value that had best convergence had an answer of about 1.5 and all other values tested had answers that were several orders of magnitude different from 1.

## 6.5 Hyperbolic Series Method

There are more efficient series which can be used to find  $\ln$ , which converge quicker than the Taylor approximation; one such method is to consider the Hyperbolic Trigonometric function  $\tanh$ . We start by considering the definition that  $\tanh(x) := \frac{e^x - e^{-x}}{e^x + e^{-x}}$ , and then find a formula for  $\tanh^{-1}(x)$ :

$$\begin{aligned}
 z = \frac{e^x - e^{-x}}{e^x + e^{-x}} &\implies z = \frac{e^{2x} - 1}{e^{2x} + 1} \\
 &\implies ze^{2x} + z = e^{2x} - 1 \\
 &\implies e^{2x}(1 - z) = 1 + z \\
 &\implies e^{2x} = \frac{1 + z}{1 - z} \\
 &\implies e^x = \left( \frac{1 + z}{1 - z} \right)^{\frac{1}{2}} \\
 &\implies x = \frac{1}{2} \ln \left( \frac{1 + z}{1 - z} \right)
 \end{aligned}$$

Using this we can see that  $2 \tanh^{-1} \left( \frac{x-1}{x+1} \right) = \ln(x)$ , and we can use the Taylor Expansion of  $\tanh^{-1}$  to approximate  $\ln$ .

Now to attain the Taylor series for  $\tanh^{-1}(x)$  we can use the same method as when we calculated the Taylor series for  $\ln$ . The exact calculations are omitted, but the end result is that we get that:

$$\tanh^{-1}(x) = \sum_{n=0}^{\infty} \frac{x^{2n+1}}{2n+1} \forall x \in \mathbb{R}^+$$

And thus by using this series we get the result that:

$$\ln(x) = 2 \sum_{n=0}^{\infty} \frac{1}{2n+1} \left( \frac{x-1}{x+1} \right)^{2n+1} \quad \forall x \in \mathbb{R}^+$$

The implementation of this is similar to previous implementations of series approximations of a function and is detailed below:

Algorithm 6.5.1: Hyperbolic series method for  $\ln$

```

1  hyperbolic_nat_log( $x \in \mathbb{R}^+, n \in \mathbb{Z}_0^+$ ):
2       $a := \frac{x-1}{x+1}$ 
3       $b := y^2$ 
4       $c := a$ 
5       $k := 0$ 
6      while  $k \leq n$ :
7           $a \mapsto a \cdot b$ 
8           $c \mapsto c + \frac{a}{2^{k+1}}$ 
9           $k \mapsto k + 1$ 
10     return  $2 \cdot c$ 

```

Using this we see that if we have  $\epsilon_n := |\ln(x) - \text{hyperbolic\_nat\_log}(x, n)|$ , then we know that  $\epsilon_n \leq \frac{1}{2n+3} \left| \frac{x-1}{x+1} \right|^{2n+3}$ . If we consider restricting our calculations to  $x \in [\frac{1}{2}, 1)$  by using the same calculations as shown for Algorithm ??, then we can see that  $|x-1| \leq \frac{1}{2}$  and  $|x+1| \geq \frac{3}{2}$ ; therefore  $\epsilon_n \leq \frac{1}{3^{2n+3}(2n+3)}$ .

By considering the final simplification that  $\epsilon_n < \frac{1}{3^{2n+3}}$ , then if we wish to have  $\epsilon_n < \tau \in \mathbb{R}^+$  it suffices to find  $n \in \mathbb{N}$  such that  $\frac{1}{3^{2n+3}} < \tau$ . In particular we consider when  $\tau = 10^{-d}$  which will guarantee  $d$  decimal places of accuracy, below is a table showing the smallest  $n \in \mathbb{N}$  that guarantees  $d$  decimal places of accuracy:

$d \in \mathbb{N}$	$\arg \min \{n \in \mathbb{N} : 3^{2n+3} > 10^d\}$
1	1
10	8
100	104
1000	1047

As can be seen in the table, fewer iterations are needed to approximate  $\ln(x)$  to the same degree of accuracy using hyperbolic series as when using the Taylor series. Further the calculations performed each iteration are very similar in complexity, both being  $\mathcal{O}(1)$ , and thus we can expect that Algorithm 6.5.1 will execute faster than 6.4.2.

## 6.6 Continued fractions

Another method for evaluating  $e^x$  is the use of continued fractions, which are a way of approximating real functions by a rational number with a recursive structure. Such fractions have been studied for many years and can be used to rationally approximate functions. Some examples of continued fractions for real numbers are:

$$e = 2 + \frac{1}{1 + \frac{1}{2 + \frac{1}{1 + \frac{1}{1 + \frac{1}{4 + \ddots}}}}} \quad \pi = 3 + \frac{1}{7 + \frac{1}{15 + \frac{1}{1 + \frac{1}{292 + \frac{1}{1 + \ddots}}}}}$$

In general a continued fraction for a number  $x \in \mathbb{R}$  has the form:

$$b_0 + \frac{a_1}{b_1 + \frac{a_2}{b_2 + \frac{a_3}{b_3 + \ddots}}} \quad (6.6.1)$$

As writing of continued fractions in the above manner takes up a lot of room and has a degree of ambiguity we will use the following notation:

$$\mathbf{K}_{n=1}^{\infty} \frac{a_n}{b_n} := \frac{a_1}{b_1 + \frac{a_2}{b_2 + \frac{a_3}{b_3 + \frac{a_4}{b_4 + \ddots}}}} \quad (6.6.2)$$

Therefore we can re-write Equation 6.6.1 as  $b_0 + \mathbf{K}_{n=1}^{\infty} \frac{a_n}{b_n}$ .

One of the most useful formulas regarding continued fractions was formulated by Leonhard Euler, and deals with the sum  $a_0 + a_0 a_1 + a_0 a_1 a_2 + \cdots + (a_0 \cdots a_n) = \sum_{i=0}^n (\prod_{j=0}^i a_j)$ . The formula derived by Euler is known as Euler's Continued Fraction Formula and is as follows:

$$\sum_{i=0}^n \left( \prod_{j=0}^i a_j \right) = \mathbf{K}_{i=0}^n \frac{\alpha_i}{\beta_i} \text{ where } \alpha_i := \begin{cases} a_0 & : i = 0 \\ -a_i & : i \in [1, n] \cap \mathbb{Z} \end{cases} \quad (6.6.3)$$

$$\beta_i := \begin{cases} 1 & : i = 0 \\ 1 + a_i & : i \in [1, n] \cap \mathbb{Z} \end{cases}$$

Many Taylor series have a structure that is compatible with equation 6.6.3 and so can be approximated by a continued fraction in this way. In particular we are looking at  $e^x = \sum_{k=0}^n \frac{x^k}{k!}$  where we note that we can re-write the series as  $e^x = 1 + \sum_{i=1}^n (\prod_{j=1}^i \frac{x}{j})$  and therefore by using Euler's Continued Fraction Formula we see that:



$$\begin{aligned}
e^x &= 1 + \frac{x}{1 - \frac{\frac{1}{2}x}{1 + \frac{1}{2}x - \frac{\frac{1}{3}x}{1 + \frac{1}{2}x - \frac{\frac{1}{n-1}x}{1 + \frac{1}{n-1}x - \frac{\frac{1}{n}x}{1 + \frac{1}{n}x}}}}} \\
&= 1 + \mathbf{K}_{i=1}^n \frac{\alpha_i}{\beta_i}
\end{aligned} \tag{6.6.4}$$

$$\begin{aligned}
\text{where } \alpha_i &:= \begin{cases} x & : i = 1 \\ -\frac{1}{i}x & : i \in [2, n] \cap \mathbb{Z} \end{cases} \\
\beta_i &:= \begin{cases} 1 & : i = 1 \\ 1 + \frac{1}{i}x & : i \in [2, n] \cap \mathbb{Z} \end{cases}
\end{aligned}$$

We would like to simplify the the above equation to remove the fractional co-efficients. If we consider multiplying  $\alpha_1$  by some constant  $c_1$ , then to have an equivalent fraction we would have to multiply it's denominator by  $c_1$ ; in practice this means multiplying  $\beta_1$  and  $\alpha_2$  by  $c_1$ . We can suppose that we could continue in a similar manner for constants  $c_2, c_3, \dots$  multiplying  $\alpha_2, \alpha_3, \dots$

**Proposition 6.6.1.** *If we have a continued fraction  $b_0 + \mathbf{K}_{i=1}^n \frac{a_i}{b_i}$  and constants  $(c_i : i \in [1, n] \cap \mathbb{Z})$ , then:*

$$b_0 + \mathbf{K}_{i=1}^n \frac{a_i}{b_i} = b_0 + \mathbf{K}_{i=1}^n \frac{c_{i-1}c_i a_i}{c_i b_i}$$

where  $c_0 := 1$ , for any  $n \in \mathbb{N}$ .

*Proof.* We will proceed by induction on  $n \in \mathbb{N}$ .

$$\mathbf{H}(n): b_0 + \mathbf{K}_{i=1}^n \frac{a_i}{b_i} = b_0 + \mathbf{K}_{i=1}^n \frac{c_{i-1}c_i a_i}{c_i b_i}$$

$\mathbf{H}(1):$

$$\begin{aligned}
b_0 + \frac{c_0 c_1 a_1}{c_1 b_1} &= b_0 + \frac{c_1 a_1}{c_1 b_1} && \text{as } c_0 = 1 \\
&= b_0 + \frac{a_1}{b_1} && \text{as required}
\end{aligned}$$

$\mathbf{H}(n) \implies \mathbf{H}(n+1):$

$$\begin{aligned}
b_0 + \mathbf{K}_{i=1}^{n+1} \frac{c_{i-1}c_i a_i}{c_i b_i} &= b_0 + \left( \mathbf{K}_{i=1}^n \frac{c_{i-1}c_i a_i}{c_i b_i} \right)_+ \frac{c_n c_{n+1} a_{n+1}}{c_{n+1} b_{n+1}} \\
&= b_0 + \left( \mathbf{K}_{i=1}^n \frac{c_{i-1}c_i a_i}{c_i b_i} \right)_+ c_n \frac{a_{n+1}}{b_{n+1}}
\end{aligned}$$

Now let us define  $b'_i$  as:

$$\begin{cases} b_n + \frac{a_{n+1}}{b_{n+1}} & : i = n \\ b_i & : i \neq n \end{cases}$$

Therefore we can continue and see that:

$$\begin{aligned}
b_0 + \mathbf{K}_{i=1}^{n+1} \frac{c_{i-1}c_i a_i}{c_i b_i} &= b_0 + \mathbf{K}_{i=1}^n \frac{c_{i-1}c_i a_i}{c_i b'_i} \\
&= b_0 + \mathbf{K}_{i=1}^n \frac{a_i}{b'_i} \quad \text{by H}(n) \\
&= b_0 + \mathbf{K}_{i=1}^{n+1} \frac{a_i}{b_i}
\end{aligned}$$

□

Using this proposition we can see that if we have the sequence  $(c_1, c_2, \dots, c_n)$  defined as  $c_i = i$  and apply it to our sequence for  $e^x$  we get that:

$$\begin{aligned}
e^x &= 1 + \frac{x}{1 - \frac{x}{2 + x - \frac{2x}{\ddots - \frac{(n-1)x}{n-1 + x - \frac{x}{n+x}}}}} \quad (6.6.5) \\
&= 1 + \mathbf{K}_{i=1}^n \frac{\alpha_i}{\beta_i} \quad \text{where } \alpha_i := \begin{cases} x & : i = 1 \\ -(i-1)x & : i \in [2, n] \cap \mathbb{Z} \end{cases} \\
&\quad \beta_i := \begin{cases} 1 & : i = 1 \\ x + i & : i \in [2, n] \cap \mathbb{Z} \end{cases}
\end{aligned}$$

This is a much simpler continued fraction, but evaluating it would still be expensive due to the repeated division operations; to get around this we can consider what are known as the convergents of a continued fraction. It is obvious that if we use a continued fraction to approximate some value  $z$  by the continued fraction  $b_0 + \mathbf{K}_{i=1}^n \frac{a_i}{b_i}$ , then there are some  $A_n, B_n \in \mathbb{N}$  such that  $z = \frac{A_n}{B_n}$ .

To start we will define  $A_{-1} := 1$  and  $B_{-1} := 0$ , and consider the case when  $n = 0$ ; for this case  $z = b_0$ , which means that  $A_0 = b_0$  and  $B_0 = 1$ . For the case when  $n = 1$  we have  $z = b_0 + \frac{a_1}{b_1}$ , which when rearranged is  $z = \frac{b_0 b_1 + a_1}{b_1}$ . This means that  $A_1 = b_0 b_1 + a_1 = b_1 A_0 + a_1 A_{-1}$  and  $B_1 = b_1 = b_1 B_0 + a_1 B_{-1}$ , and for the case when  $n = 2$  we get the similar result that  $A_2 = b_0 b_1 b_2 + a_2 b_0 + a_1 b_2 = b_2 A_1 + a_2 A_0$  and  $B_2 = b_1 b_2 + a_2 = b_2 B_1 + a_2 B_0$ .

It is actually true that this relationship continues for all  $n \in \mathbb{N}$ , and thus we get what are known as the Fundamental Recurrence Formulas for continued fractions:

$$\begin{aligned}
A_{-1} &= 1 & B_{-1} &= 0 \\
A_0 &= b_0 & B_0 &= 1 \\
A_{n+1} &= b_{n+1} A_n + a_{n+1} A_{n-1} & B_{n+1} &= b_{n+1} B_n + a_{n+1} B_{n-1} \quad \forall n \in \mathbb{Z}_0^+
\end{aligned}$$

Using this and our simplified continued fraction for  $e^x$  we can use the following method to approximate  $e^x$  by using a continued fraction up to  $a_n, b_n$  where  $n \geq 2$ :

Algorithm 6.6.1: Continued fraction for  $e^x$

1	<code>cont_frac_exp</code> ( $x \in \mathbb{R}, n \in \mathbb{N}$ ):
2	$A_1 := x + 1$

```

3      B1 := 1
4      A2 := x2 + 2x + 2
5      B2 := 2
6      a := -x
7      b := 2 + x
8      k := 2
9      while k < n:
10         a ↦ a - x
11         b ↦ b + 1
12         Ak+1 := bAk + aAk-1
13         Bk+1 := bBk + abk-1
14         k ↦ k + 1
15     return  $\frac{A_k}{B_K}$ 

```

One observation of the above algorithm, when implemented on a computer, is that if we pre-generate  $b_i$  and  $a_i$  for  $i \in [2, n] \cap \mathbb{Z}$  then the calculations of  $A_i$  and  $B_i$  are independent. This means that, if supported by the computer, both  $A_i$  and  $B_i$  could be computed in parallel. This may allow an implementation of the algorithm to be more efficient than one that computes the function in sequence.

Now while continued fractions are useful for approximating functions, it is difficult to evaluate the error of their output analytically. One result is that if  $a_n = 1 \forall n \in \mathbb{N}$  when approximating some value  $z$ , then  $|z - \frac{A_n}{B_n}| \leq \frac{1}{|B_{n+1}B_n|}$ . If we transform the continued fraction of  $e^x$  into this form by using Proposition ??, then we get that:

$$e^x = 1 + \frac{1}{(-\frac{1}{x}) + \frac{1}{(-\frac{2+x}{2x}) + \frac{1}{(-\frac{3+x}{3x}) + \ddots}}}$$

By using a computer to implement the calculations for a test value of  $x = 1$ , we see that  $\frac{1}{B_5B_6} = 0.009131261889664\dots$  and  $\frac{1}{B_{10}B_{11}} = 0.000041307209877\dots$ ; thus we can guarantee two decimal place of accuracy with `cont_frac_exp(1, 5)` and 4 with `cont_frac_exp(1, 10)`. However if we instead have  $x = 14$  then  $\frac{1}{B_{10}B_{11}} = 0.314711263190806\dots$  and convergence is similarly poor for negative values.

Further computations show that convergence of  $x \in (0, 1)$  is better than the convergence when  $x = 1$ , and thus we can use the identities and conversions to ensure good convergence. In particular if  $x \in \mathbb{Z}$  then we can calculate the reciprocal of `cont_frac_exp(-x, n)` and if  $x \in (1, \infty)$  we use the identity that  $x = a \cdot 2^b$ ; with this we see that  $e^x = e^{2^b} e^a$  and  $2^b \in \mathbb{Z}^+$ . If we evaluate  $e^{2^b}$  separately using one of the previous methods, such as algorithm ??, then we can restrict `cont_frac_exp` to  $x \in (0, 1)$ .

With this restriction in place we know that algorithm 6.6.1 converges at least as quickly as it does for  $x = 1$ , and thus we can use it's convergence to guarantee the convergence of our method. Below is a table that shows the minimum  $n$  needed to acheive the associated  $d$  decimal places of accuracy:

$d$	Minumum $n$ to guarantee $d$ decimal places of accuracy
1	2
10	22
100	235
1000	2386

While this is useful to consider, the actual implementations of the algorithm will use the iterative error to test for convergence. This is because due to the speed of convergence when the difference between successive estimates is small enough then the latest estimate is also close enough to the correct answer.

An alternative continued fraction for  $e^x$  that arises from the generalized hypergeometric series is:

$$e^x = \frac{1}{1 - \frac{x}{1 + \frac{x}{2 - \frac{x}{3 + \frac{2x}{4 - \frac{2x}{5 + \frac{3x}{6 - \ddots}}}}}}} \quad (6.6.6)$$

$$= \mathbf{K}_{i=1}^n \frac{\alpha_i}{\beta_i} \quad \text{where } \alpha_i := \begin{cases} 1 & : i = 1 \\ -x & : i = 2 \\ (-1)^{i-1} \lfloor \frac{i-1}{2} \rfloor x & : i \in [3, \infty) \cap \mathbb{Z} \end{cases}$$

$$\beta_i := \begin{cases} 1 & : i = 1 \\ i - 1 & : i \in [2, \infty) \cap \mathbb{Z} \end{cases}$$

Due to the  $(-1)^{i-1} \lfloor \frac{i-1}{2} \rfloor$  factor in the definition of  $\alpha_i$  it is more efficient to perform two updates each step rather than one. Below is the implementation of this method:

Algorithm 6.6.2: Continued fraction for  $e^x$  version 2

```

1  cont_frac_exp_v2( $x \in \mathbb{R}, n \in \mathbb{N}$ ):
2     $A_1 := 1$ 
3     $B_1 := 1$ 
4     $A_2 := 1$ 
5     $B_2 := 1 - x$ 
6     $a := 1$ 
7     $b := 1$ 
8     $k := 2$ 
9    while  $k < n$ :
10      $a \mapsto xa$ 
11      $b \mapsto b + 1$ 
12      $A_{k+1} := bA_k + aA_{k-1}$ 
13      $B_{k+1} := bB_k + aB_{k-1}$ 
14      $k \mapsto k + 1$ 

```

```

15 |         b ↦ b + 1
16 |         Ak+1 := bAk - aAk-1
17 |         Bk+1 := bBk - abk-1
18 |         k ↦ k + 1
19 |     return  $\frac{A_k}{B_K}$ 

```

The fraction needed to analyse this method is again found by using proposition ?? and is:

$$\begin{array}{c}
 \frac{1}{1 + \frac{1}{-\frac{1}{x} + \frac{1}{-2 + \frac{1}{\frac{3}{x} + \frac{1}{2 + \frac{1}{-\frac{5}{x} + \frac{1}{-2 + \ddots}}}}}}}
 \end{array}$$

By implementing this with a computer we get similar results to above, particularly that there is rapid convergence for  $x \in (0, 1)$ . Further the convergence of values in  $x \in (0, 1)$  is more rapid than  $x = 1$  and so we can use the convergence of  $x = 1$  as an upper bound of our method. Below is the table showing the smallest  $n \in \mathbb{N}$  needed to ensure  $d$  decimal places of accuracy for some  $d \in \mathbb{N}$ :

$d$	Minumum $n$ to guarantee $d$ decimal places of accuracy
1	4
10	12
100	61
1000	405

As can be seen the convergence of 6.6.6 appears to be significantly faster than that of 6.6.5 and one might be satisfied by this, however an even better solution exists.

As the fraction 6.6.6 can be shown to converge for all values of  $x$  to  $e^x$  then we can consider the even and odd convergents. The even convergents of the sequence are  $\frac{A_0}{B_0}, \frac{A_2}{B_2}, \frac{A_4}{B_4}, \dots$ , while the odd convergents are  $\frac{A_1}{B_1}, \frac{A_3}{B_3}, \frac{A_5}{B_5}, \dots$ . As  $\lim_{n \rightarrow \infty} \frac{A_n}{B_n} = e^x$ , then  $\lim_{n \rightarrow \infty} \frac{A_{2n}}{B_{2n}} = \lim_{n \rightarrow \infty} \frac{A_{2n+1}}{B_{2n+1}} = e^x$ ; the following proposition gives an explicit form for the odd and even convergents.

**Proposition 6.6.2.** *If  $z = K_{i=1}^{\infty} \frac{a_i}{1}$ , then the limit of the odd convergent of  $z$  is:*

$$\begin{array}{c}
 x_{odd} = a_1 - \frac{a_1 a_2}{1 + a_2 + a_3 - \frac{a_3 a_4}{1 + a_4 + a_5 - \frac{a_5 a_6}{1 + a_6 + a_7 - \ddots}}}
 \end{array}$$

while the limit of the even convergent is:

$$x_{even} = \frac{a_1}{1 + a_2 - \frac{a_2 a_3}{1 + a_3 + a_4 - \frac{a_4 a_5}{1 + a_5 + a_6 - \ddots}}}$$

*Proof.* Ommitted □

If we apply proposition ?? to 6.6.6, to acheive the form  $\mathbf{K}_{i=1}^{\infty} \frac{a_i}{1}$  then we end up with the following fraction:

$$e^x = \frac{1}{1 + \frac{-x}{1 + \frac{\frac{1}{2}x}{1 + \frac{-\frac{1}{6}x}{1 + \frac{\frac{1}{6}x}{1 + \frac{-\frac{1}{10}x}{1 + \frac{\frac{1}{10}x}{1 + \frac{-\frac{1}{14}x}{1 + \ddots}}}}}}}} \quad (6.6.7)$$

Now if we apply proposition 6.6.2 to the above fraction we see that:

$$e^x = 1 + \frac{x}{1 - x + \frac{1}{2}x + \frac{\frac{1}{12}x^2}{1 - \frac{1}{6}x + \frac{1}{6}x + \frac{\frac{1}{60}x^2}{1 - \frac{1}{10}x + \frac{1}{10}x + \frac{\frac{1}{140}x^2}{\ddots}}}} \quad (6.6.8)$$

Finally by simplifying and applying proposition ?? one more time we reach the following continued fraction for  $e^x$ :

$$e^x = 1 + \frac{2x}{2 - x + \frac{x^2}{6 + \frac{x^2}{10 + \frac{x^2}{14 + \ddots}}}} \quad (6.6.9)$$

$$= 1 + \mathbf{K}_{i=1}^{\infty} \frac{\alpha_i}{\beta_i}$$

$$\text{where } \alpha_i := \begin{cases} 2x & : i = 1 \\ x^2 & : i \in [2, \infty) \cap \mathbb{Z} \end{cases}$$

$$\beta_i := \begin{cases} 2 - x & : i = 1 \\ 4i - 2 & : i \in [2, \infty) \cap \mathbb{Z} \end{cases}$$

If we implement this method by using the Fundamental Recurrence Formula then we get the following:

### Algorithm 6.6.3: Continued fraction for $e^x$ version 3

```

1  cont_frac_exp_v3( $x \in \mathbb{R}, n \in \mathbb{N}$ ):
2       $A_0 := 1$ 
3       $B_0 := 1$ 
4       $A_1 := 2 + x$ 
5       $B_1 := 2 - x$ 
6       $a := x^2$ 
7       $b := 2$ 
8       $k := 1$ 
9      while  $1 < n$ :
10          $b \mapsto b + 4$ 
11          $A_{k+1} := bA_k + aA_{k-1}$ 
12          $B_{k+1} := bB_k + aB_{k-1}$ 
13          $k \mapsto k + 1$ 
14     return  $\frac{A_k}{B_k}$ 

```

As with the previous two continued fraction methods of approximating  $e^x$  we can apply proposition ?? to 6.6.9 to find the following equivalent continued fraction:

$$1 + \frac{1}{\frac{1}{x} - \frac{1}{2} + \frac{1}{\frac{12}{x} + \frac{1}{\frac{5}{x} + \frac{1}{\frac{28}{x} + \dots}}}}$$

Again a computer was used to evaluate  $B_k$  of the above fraction, which gave the expected results of quick convergence for  $x \in (0, 1)$  and more rapid convergence for  $x \in (0, 1)$  than  $x = 1$ . Using this the table below was generated to show the minimum  $n \in \mathbb{N}$  that guarantees  $d$  digits of accuracy:

$d$	Minimum $n$ to guarantee $d$ decimal places of accuracy
1	2
10	6
100	30
1000	202

This has the fastest theoretical convergence of the three methods, and thus is expected to perform the best.

## 6.7 Comparison of Methods

We have introduced several methods for calculating both logarithms and exponentials in this chapter, and considered their theoretical convergence; we now look at a direct comparison of the methods as implemented in C.

The first consideration is which values to use while comparing methods. While all the methods converge for all values, or can be made to by using transformations of the inputs and outputs, most methods converge best for small values. Therefore values being tested will typically be

in the range of  $[0.5, 1)$ .

The first methods to be compared here are the versions of the continued fraction method discussed previously. Below we have the outputs of different versions of the program for different values of  $n$ , with the bold digits being the correctly approximated digits.

$n$	cont_frac_exp_v1	cont_frac_exp_v2	cont_frac_exp_v3
1	1.94499999999999	3.33333333333333	2.0769230769230
2	2.00216666666666	3.33333333333333	2.0132689987937
3	2.01217083333333	2.0769230769230	2.0137543842848
4	2.01357141666666	2.0054200542005	2.0137527042253
5	2.0137348180555	2.0132689987937	2.0137527074744
6	2.0137511581944	2.0137906192914	2.0137527074704
7	2.0137525879565	2.0137543842848	2.0137527074704
8	2.0137526991603	2.0137526161232	2.0137527074704
9	2.0137527069445	2.0137527042253	2.0137527074704
10	2.0137527074399	2.0137527076056	2.0137527074704

As can be seen here the first two methods have similar convergence, however despite having a very poor theoretical convergence the first method converges better than the second version. However it is obvious that the third method has the fastest convergence, and thus should be the one to use in further comparisons.

Now we can compare the speed of the Taylor and continued fraction methods of calculating exponential values. For this we will use 1000 values in the range  $[\frac{1}{2}, 1)$  and calculate each 100000 times to compare the speed of the method. We will be using values of  $n$  which guarantee 10 decimal places of accuracy, in particular  $n = 14$  for `taylor_exp` and  $n = 6$  for `cont_frac_exp_v3`.

The results of the tests run on my computer are included in the table below alongside those for the built in `exp` function in `math.h`:

	Total time:	Average time:	Minimum time:	Maximum time:
<code>taylor_exp</code>	12.430s	0.012s	0.012s	0.019s
<code>cont_frac_exp</code>	4.741s	0.004s	0.004s	0.007s
<code>builtin_exp</code>	2.608s	0.002s	0.002s	0.004s

This shows that the continued fractions method of evaluating exponential functions is almost three times as efficient as the standard Taylor series method. However both fall short of the built in method, despite the hyperbolic series method being a close second. This is likely due to a lower level implementation of the exponential function with various highly efficient programming practices implemented to optimize the code execution speed.

The two methods discussed to evaluate  $\ln$  have their convergence for different values of  $n$  shown below, where they are approximating the value 0.7, with the bold digits representing the correctly approximated digits:



$n$	taylor_nat_log	hyperbolic_nat_log
1	-0.300000000000000	-0.356604925707
2	-0.300000000000000	-0.356673383305
3	-0.345000000000000	-0.356674906089
4	-0.354000000000000	-0.356674942973
5	-0.356025000000000	-0.356674943913
6	-0.356511000000000	-0.356674943938
7	-0.356632500000000	-0.356674943938
8	-0.356663742857	-0.356674943938
9	-0.356671944107	-0.356674943938
10	-0.356674131107	-0.356674943938

## 7 Preliminary References

<http://math.exeter.edu/rparris/peanut/cordic.pdf>

Inside your Calculator by Gerald R Rising

Wolfram Alpha