# High Performance Selection in Athletics
# Michael Whitehead 1078439

**ABSTRACT**

The current selection for athletic persons of a High Performing level against the world stage has issues. There has been an increase in interest to calculate an athlete's ability to perform in competitions for Olympic or World Championship level. Athletics New Zealand currently has a system in place to support and select the best athletes based on previous performances however they are only taking into account the athlete's best ability.

Hidden costs within the sport make it a challenge for a great deal of potentially skilled athletes to continue to the High Performing level based on the current system Athletics New Zealand has in place.

Ideally, based off research done by a PHD student in Auckland University of Technology, the best athletes who perform at the Olympics or World Championships produce performances that are both at a consistent and high achieving level. Looking only at an athlete's best ability is not a good representation of the athlete at all. When comparing athletes abilities the best indication of ability is by the performance 'curve' over their careers.

This process is not only limited to high performing athletes but can be used on lower level athletes to compare ability and determine potential futures for later career success which is greatly sort after by selection committees and sponsors who intend to support these athletes to the High Performing level and beyond.

This report provides a technological solution to the theories and research that Steve Hollings has provided as well as incorporating future designs for the technology and how it can be used for the whole athletic community.

# Contents

# 1. Introduction

## 1.1 Current situation

Athletes throughout New Zealand and around the world have difficulty in competing on the world stage with other High Performing athletes. They all have financial issues that see a great deal of talented athletes leave as they cannot support themselves to travel to important national and international meets or, afford necessary costs such as equipment which needs to be regularly updated, or even physio. An effective way of getting support for their discipline is to be performing on a world level stage, which unless they develop quickly, they will run into financial troubles when getting to the High Performing level. Sponsors and committees have the ability to support these athletes but they too are stuck in a similar situation. How do they figure out who has the potential to be great and is going to go the distance if they get support? The other dilemma that originates from the financial support is how to get to the races of importance. Anybody can compete at athletics events but, without competition athletes have a difficult time developing themselves for the events they have chosen. The national and international events are where these athletes need to be in order to achieve their best results.

Athletics New Zealand currently have a system in place that identifies High Performing athletes within New Zealand and supports them according to their ability using a comparison with High Performing athletes from 1980 - 1990. The program is used only to fund athletes; anything else such as sending athletes to compete at world events is based on an athlete's best performance within a certain timeframe and is compared to the times at the prior competitions. Athletes use a similar method to determine where they sit in comparison to others, comparing best performances done at appropriate times.

## 1.2 Project Focus

This project focuses on collecting, displaying and comparing High Performing athletes with potentially high level athletes. The specific goals include a web scraping tool for data collection and a graphing program that can display and extrapolate results with user data.

The project focuses on Athletics events and works with simple statistics to produce visual data that users can relate to and understand. It takes relevant terms and models used currently in the athletic community to simplify the learning curve for the audience targeted for this program.

The project aims at creating a visual display that people connected to an athlete can use to more correctly determine that athlete's level of ability based on all performances the athlete has achieved, with the added benefit of showing future potential of the athlete based on their performances.

The future intention of the project is to have athletes add performances and compare their abilities not only with high performing athletes but also with athletes of a similar level within their country. Athletes perform on a regular basis so a data base where data can be uploaded and downloaded from will keep athlete programs up-to-date. National level data can also be added in this way giving better comparisons for a wider range of the athletics community.

## 1.3 Steve Hollings work

The methods and ideas for comparing High Performing athletes with athletes of a country were created by Steve Hollings, PHD student at the Auckland University of Technology. Ten years prior Steve had created a similar method of selecting athletes for High Performing races targeting the Olympics and World Championships. He recently reviewed his methods and decided on a better solution to calculate an athlete's potential. From this he wanted a working program that he could add to his PHD. The work and ideas in this project have developed from his base concept idea. They still incorporate the main goal that Steve wanted to produce but also incorporate extra adaptations and methods to display and compare athlete's performances.

## 2. Background

*"Mediocrity acknowledges nothing higher. Talent recognises genius"* – Oscar Wilde.

Countries select athletes to compete for their flag for various reasons. An athlete is the best in their event is generally the most common selection requirement. However, what happens when there are five or ten athletes who are all doing the same performances? This is where other factors come in to figure out who will be the best to support and send to world events.

Various factors include but are not limited to; past performances, comparison with world level athletes, mentality, physical well-being and personal history (i.e. Arrests, drug use). Countries can roughly assess from these factors who to support and send based on their results. However, in recent years countries have been attempting to support younger athletes who they believe have the potential to be the next best athletes for their country, a great idea but it does come with disadvantages. Anything could happen to the athlete before they get to the world stage which could see a loss in return from the support offered. They possibly have hit their peek already or they may decide to give up.

Athletics New Zealand wanted to create a way to support these younger athletes to get them to a high performing level. As a result they created a system to support their athletes as well as a program calculating to make sure they were meeting a certain requirement.

The system has various rules regarding the selection, keeping the sponsorship and amount of support offered. Figure 2.1 on the next page shows the basic idea in the pyramid. There are 4 categories for different levels of athletes. Gold to Green are athletes who have the potential or are in the top 16 in the world for their event. Red is for the younger or lower levelled athletes who are trying to work their way up to the high performing athletes in later years. They are still of a national or higher level.

**Figure 2.1 Athlete Carding Program –Athletics New Zealand**

Their program that they use to calculate what category an athlete should fall works off ten years worth of data and gives a rough estimate off the mean of data from past years performers. It calculates performance levels based off an athlete's best performances. The issue here is that by only looking at a select portion of the data available there is plenty of room for error. The data being sampled is at times outlying data meaning that it does not reflect a person's actual ability, only their best ability. The issue with this is that sending people to compete based off their best performance may not result in their best ability being performed. The other problem that arises is that although there is a good ten years worth of data, athletes abilities are always changing and always improving. This is due to the various methodologies and theories that have been hypothesised, created and trialled to get the most out of the human body. Over the years, someone who won gold at a past Olympics with 10.10 Seconds for the 100 metres would not make the finals for the Olympics of current times.

During the final stages of the graph development interpolation was a useful tool in creating the trend lines which display the athlete's performances over time. This was because the data being used only had 3 points of interest to plot the graph. However when creating a dynamic graph that has zooming capabilities, the ability to re-plot points based on the zoom meant the curve would not lose its shape which is required when comparing lines of performances that are in the tenths of seconds or millimetres in difference.

One of the features thought up towards the end of the project was to be able to predict the future potential of the athlete's development based off the results they had inputted so far. Extrapolation is the process of extending the line off a line of best fit using the equation for that line. This is of use in athletics as various parties are interested in the next one to four years of the athlete's development and whether there is potential for that athlete to improve assuming there are not any major factors that limit them to perform in that time frame.

This project looks at all performances of an athletes career and, creating a $2^{nd}$ order polynomial line of best fit through that data to show their progression over time. This line is then compared with other high performing athlete's lines of their careers to give users an indication of where they are heading in comparison to the world's best. The benefit of this is it can give users an indication of their future performances if they stay on their current path. Whether their lines are rising upwards, levelling out or falling down gives plenty of information on where they are currently performance wise.

The project takes the ideas of Steve Hollings work for comparing High Performing athletes and extends it to be used as a coaching tool as well as producing valuable information around the athlete and their competition or to create a sense of future potential for the athlete to see and work off.

## 3. System Design

The project consists of a web scraper and 3 versions of the program. The first version of the program was created for a client Steve Hollings who was after a quick and simple program that could be used for his work in Auckland. He wanted a simple application that used Microsoft Excel to display High Performing athletes against user data. The second version was a follow on from this idea only it was based in C# and worked with an Excel wrapper to use the Excel graphs and data storage. The final program gets away from using Excel other than using it as data storage and works off its own ability to plot and draw the graphs that Steve Hollings wanted. The three versions show various ways to create the same base idea as well as a development of ideas and concepts from one to the next. The web scraper is a program aside from the main three. This one concerns itself with collecting the data so that the other three versions have data to work with.

The program consists of a section storing the data, a graph displaying the information, a data grid to collect data and a form to navigate to the various event sections. The data storage is a list of data on High Performing athletes that are plotted on the graph. The graph displays both the High Performing athletes as well as the data input by the user to give a comparison in abilities and other valuable information. The data grid is the area for user's data to be entered. The navigation works by collecting user's information such as their name and date of birth, and requests which event they would like to compare their data too. The collected information is then used to calculate from their data that's input to give a line that is comparable to the High Performing athletes collected data.

## 4. Version 1 Excel

This program uses macros and Visual Basic coding to achieve a user friendly environment. It works by opening a form requesting information on the task you want to achieve and then navigating you through the large Excel file to the allocated data set with graph display.

### 4.1 Excel document

Macros can be used to setup, format and calculate various aspects of an Excel document. For this program, the project is already formatted before being given to the user. This is beneficial for two reasons. The first being the processing power on machines. Older or less powerful machines struggle with opening the document as it's a lot larger than a normal Excel document. The other concern with the size of the document and using Excel is that it loads every sheet at runtime making all sheets available for access. However due to the way users operate the program there is no need for all Excel sheets to be created. When operating the program, you can only ever open one Excel Sheet at a time. As a result, having all sheets loaded causes unnecessary load on the computer, making the opening and initial usage of the program slow and resulting in people being less inclined to use it.

The ability to store all the data within sheets works to the client's advantage as the data is stored within the program and is not located in an obscure place that users might accidentally delete; it does not however hold any real benefit to users. In fact, having data so close to the working area of users increases the probability of them accidentally going into it and destroying the program entirely. This was one of the concerns Steve Hollings had and wanted a way of protecting the data so it was not fooled around with.

### 4.2 Security

Excel has adequate security however it, like most things can be conquered. If a few select people were to be using this software there would not be an issue with security as the few users would likely be trusted with the data meaning they would have admin privileges. However, as this data is used by everyone *it's the equivalent of having all the records stored in a locked room with only a few people having access. If you're sending the file out to people who should only have limited access, that's the equivalent of letting someone into the room and telling them to only look at 1 drawer. They might (and should) do what you said, but there is the chance that they won't. (Luke M. Chandoo.org 2009).*

Security for the data is a necessity as without it, the program becomes useless. The client wanted a way of looking after the data collected. You can lock a file but people need to be let in to use the application. You can lock the sheets but the security on this for Microsoft Excel 2007 and earlier versions are very simple to get past. The most practical solution to get past using software that is very easily broken is to change the working environment.

## 4.3 User Interface

The last issue with creating a program within an Excel document is the fact that there are two programs to work with. The program itself and Excel that the program has been created in. This issue gives too many options and features for users to worry about. An example of this may be that they are trying to print the document, but how do they go about it? Do they look in the program created to find a print function to do the job, or are they required to go through Excels program features to find Excels print function? The Heuristic Evaluation principle Aesthetic and minimalist design for user interface design sums it up nicely. *Dialogues should not contain information which is irrelevant or rarely needed. Every extra unit of information in a dialogue competes with the relevant units of information and diminishes their relative visibility.(2005 by Jakob Nielsen)* So, to get the best results out of people when using the program, simplifying the user interface as much as possible is required and means moving on from Excel applications.
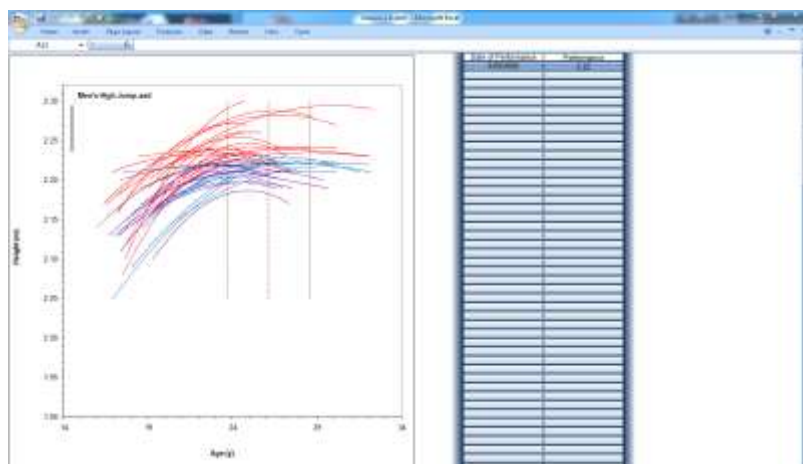


**Figure 4.3.1 Excel Program**

Even though it has its issues, the program itself still achieves the basic functionality required by the client. It has the capability to create graphs showing the data in a quadratic line based on performances and age of performance. It also compares user's data to that of the High Performing athlete's data and gives users valuable information. Overall it achieves the goals it was designed for even if it has some minor flaws.

## 5. Version 2 Excel C#

Version two works with integration of Excel into C#. The plan was to increase the security around the data stored and reduce confusion when navigating through the user interface. Using Excel within the program was to bring across the graph generated by Excel as well as keeping the background calculations associated with it.

### 5.1 Excel Wrapper

The integration of Excel in C# was done using a wrapper that navigated to the Excel file directly. It works by opening the requested document and displays it in the window. Users then had the ability to use the graphing functionality that the Excel document provided while leaving out the Excel Ribbon and titles to the columns and headers of the cells. By itself it proved to do the task that it could do in Excel, however with this version there was no distractions or features being displayed that users were not going to use.

### 5.2 Database

Another way to read and write data to Excel files using C# is by using OLEDB (Object Linking and Embedding, Database). OLEDB is an API designed by Microsoft to allow accessing data from a variety of sources in a uniform manner. This process allows you to extract and input data into Excel files.

OLEDB is much faster at reading a range of cells than the Excel Wrapper. This is because the API directly works with the data in the Excel file where as the Wrapper is a middle ware that gets commands to pass onto the Excel document and gets data from Excel to pass onto the C# application.

The reason this process has not been implemented in this project is because the client wanted to have the Excel graph not only display the information but also when clicking on it there was a hover effect that could show you the name of the athlete whose data you are viewing. The API can only give you snapshots of updated graphs which although shows the same data, there is no interaction with the graph when the user hovers over. As a result it was easier to use the excel wrapper to handle the communication with the Excel to C# integration than using the API.

The reasoning for creating this version was to give it a more professional look. The first version looked unprofessional and poorly organised. The change was intended to give users confidence when using the program that it did do a good job as opposed to giving an impression that it was of poor quality. However, after beginning this version a lot more features that were once not available, now could be implemented easily and be used to make the process more effective. Features such as storing multiple athletes for the same event, displaying multiple athletes at the same time and even displaying information about the project could be done with ease, where once it was not considered as the flexibility of the programming environment limited the features that could be implemented.
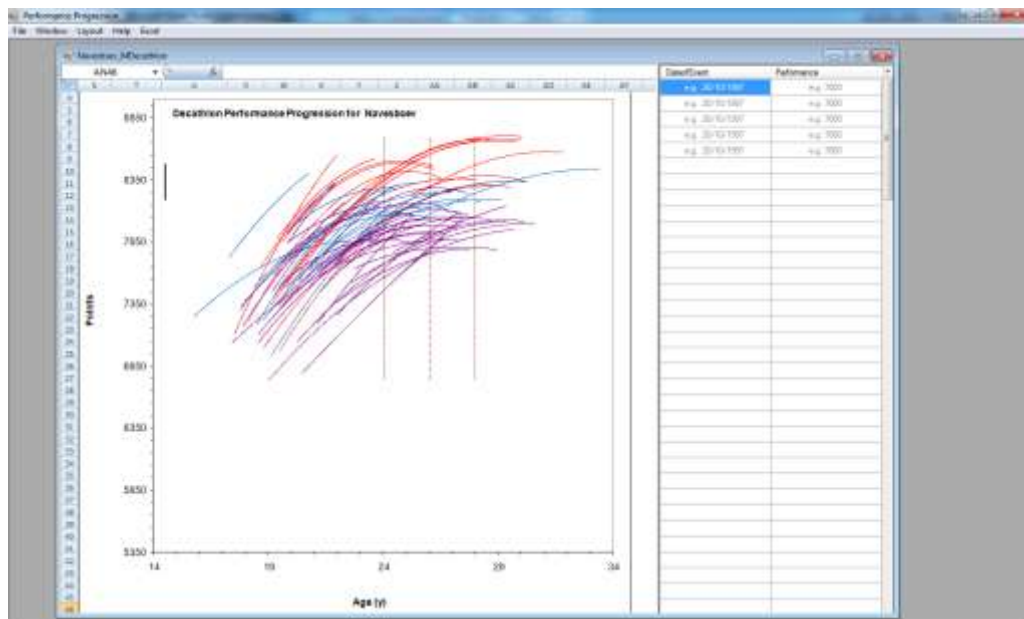


**Figure 5.2.1 C# Program**

One of the huge gains from implementing the project this way was the huge decrease in loading times for the Excel documents. Originally the document was 1.5MB but as the user only ever wants to open up a single page, the new documents that contained only 1 page were 38KB showing a decrease in file size of 97.46%. This meant quicker loading times as well as making it easier on older or smaller computers who do not have much processing power.

## 5.3 Multiple Document Interface (MDI)

MDI proved to be quite challenging. The small feature of opening up multiple athletes at once to compare between athletes proved to be a challenge. The Excel wrapper did the job it needed to so understanding the process of how it connected and communicated to Excel was not a top priority however; when MDI was implemented the program began doing unexpected things.

When implementing MDI while using Excel wrapper to send and receive Excel data to the C# application it began producing issues as to how this communication should work. The Excel wrapper navigates to any Excel document very well however, once navigated to, only the last document to be opened has focus. All the rest are only being displayed. If data is entered into one of the windows then the data is stored only in the last document that was opened. Figure 5.3.2 shows the connection below. The issue was that the wrapper would only link you to the excel file. As soon as another instance was loaded you were still linked to the file but you could no longer send or receive data to that file. The link to send or receive any new data would be set to the new excel document that was opened.
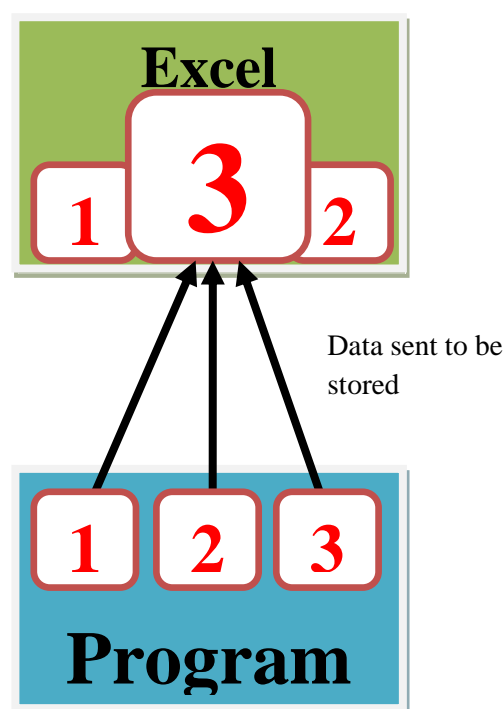


**Figure 5.3.2 Excel Wrapper connection**

To fix this bug the old windows needed to become the main centre of focus once more by forcing Excel to make the document being worked on to become the main point of focus and put the last one to be displayed to run in the back.

## 5.4 Excel Hierarchy

Of the three main Microsoft Office tools (Excel, Word and PowerPoint) Excel opens and works with multiple documents under the same Excel program. The others open up new versions of themselves for each document. The hierarchical structure of excel is shown in Figure 5.4.1 on the following page.
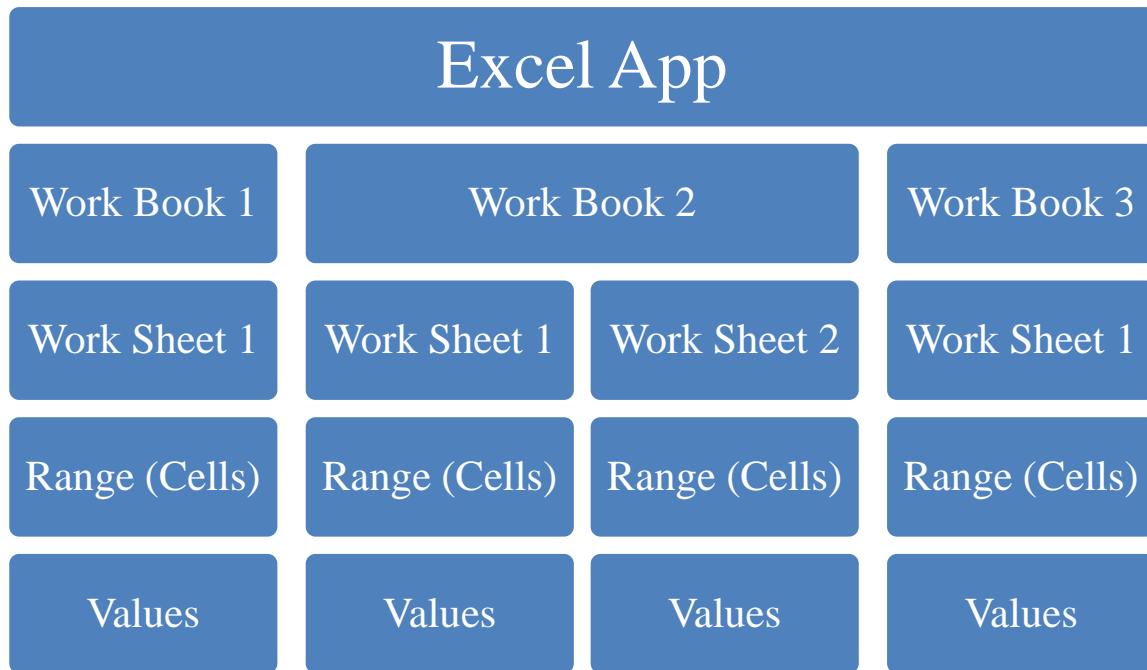
**Figure 5.4.1 Excel Hierarchy**

The Excel Application runs for 1 to n workbooks. Each workbook contains 1 to n worksheets. Each worksheet contains a range of cells which each contain some value. This structure of having every workbook open under the same application was the reason for the Excel wrappers having difficulty in communicating with their individual workbook. In Excel, if multiple workbooks have been opened, only the most recent workbook worked on has focus. If Word had been used instead, the focus for the Excel wrapper to the word document (Word wrapper) would have always had the right focus as the documents are always opened separately. The solution was to force focus on the workbook that the Excel Wrapper was suppose to be connected to.

When closing the Excel documents in the program the process of closing was also different. To close an Excel document meant closing the workbook first (for saving purposes) and then closing the Excel application. Having multiple workbooks worked in a similar fashion. Close the workbook or workbooks you wanted closed before closing Excel. However, the reason for all this was that if you had multiple workbooks opened and you closed the application then all the workbooks were left in an open state. This meant that if you tried opening these documents later on, they would say that a different application already had them open and the only way to free them up again was to kill the Excel process completely.

## 6. Version 3 C# with Excel

This version was created to work past the Excel boundaries. The main reason for this was to get more flexibility with the graph. Originally the graph was static and unmoving. If users had data that lay outside the X and Y limits then they were not receiving the full extent of what the program had to offer. The other main reason for the change was the elements that are to be introduced at a later date. Keeping the original excel wrapper would have limited the features and potential displays that the program could produce.
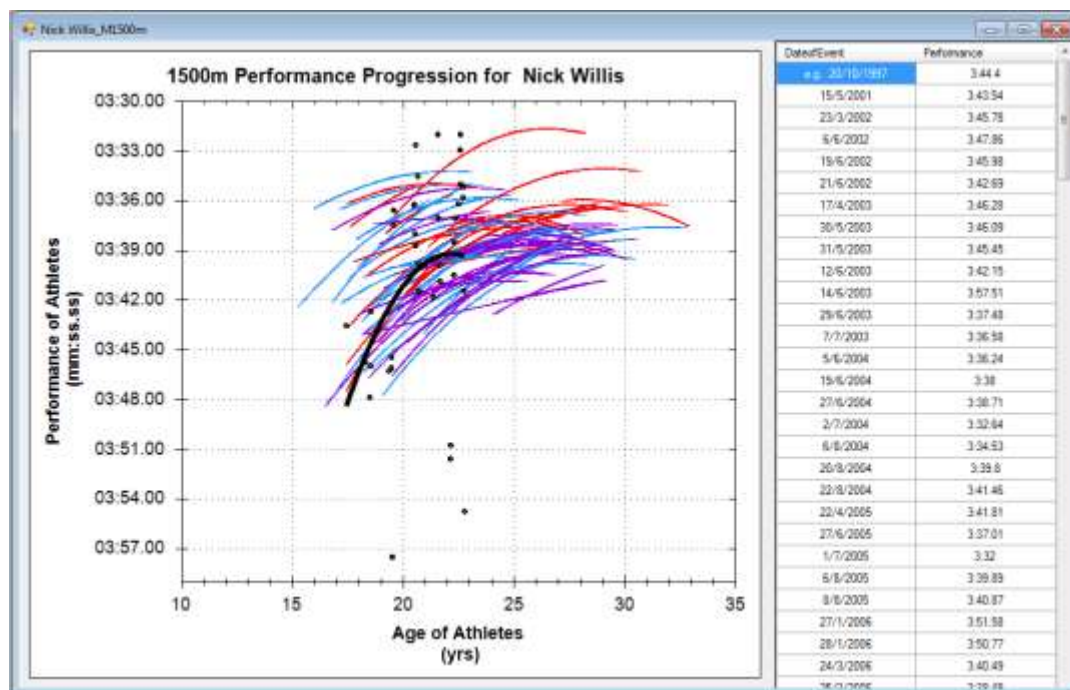


**Figure 6.0.1 C# Program Child Form**

By removing the Excel wrapper a number of improvements became apparent. Firstly, the new graph that is displayed did not have any diminishing results. You are able to move and resize the graph with little hassle. When the window is resized, the graph is also resized which benefits users who want to compare multiple charts greatly. When using the Excel graph the graph would not resize with the form so comparing athletes meant playing around with the size of the forms to get them just right, failing after more than two forms were being used.

When closing the program if multiple documents were left open in the program, the wrappers were not closing the links they had made to the Excel files. The process was still running even after closing the program. The solution was simple enough. Close the excel file when you close a MDI child. However, MDI forms have a particular way to close child and parent forms.

When closing the program if there were MDI children forms open the program would fire the closing trigger on every child form before firing the closing trigger on the parent form, then it fires all the closed triggers for every child before firing the closed event for the parent form. The issue here was how the closing was done. The closing of multiple child forms needed to be closed with a certain piece of code so that the wrappers would close their connections to the Excel documents. This was placed in the closing trigger of the parent form. However, by the time the parent form triggers a firing event, every child event has already triggered their closing events once. If no save dialogue was included this issue never would have arose. The solution was to ignore the second closing trigger if it was triggered so that only one request of asking to save was initiated.

## 6.1 Milestones

Milestones are used in many areas of life to set out and plan goals that want to be achieved by the user. In running, milestones are highly implemented in coaching and management aspects of the athlete's career. Milestones in athletics generally target specific times or events that an athlete sees themselves doing two to four years down the track. This connects with the program through younger athletes using the program. In the beginning the program was designed for the high end athletes who may be able to compare their times with the world's best. However, through the development of the program, it became apparent that younger or less skilled athletes may also want to use the program. The downside here is even though they want to use it, unless they can compare themselves to someone, the programs usefulness to the user becomes limited.
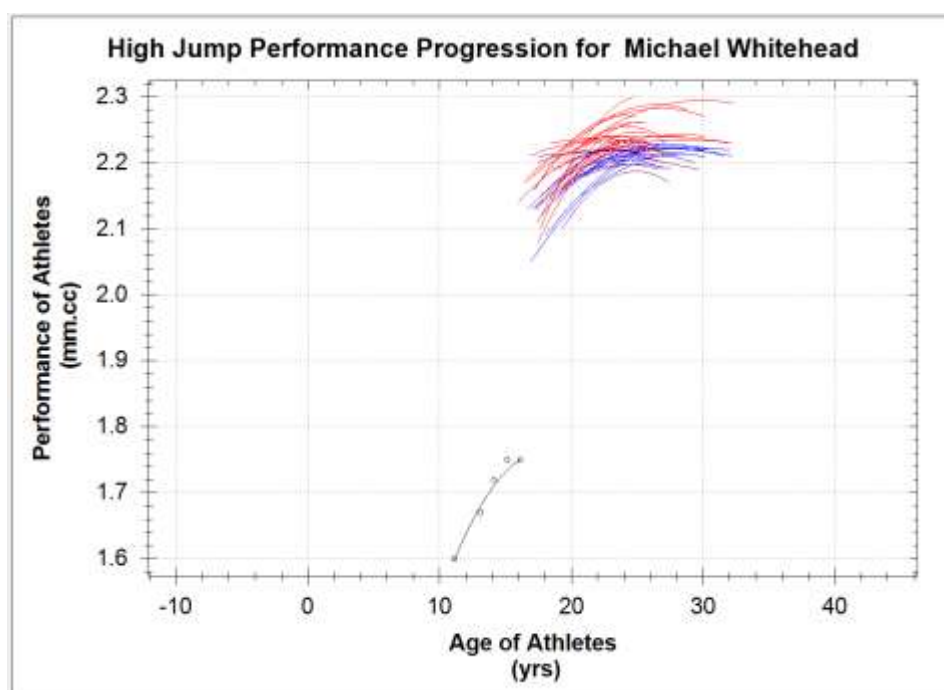


**Figure 6.1.1 Graph representing an instance of a younger athlete being plotted.**

As you can see in Figure 6.1.1 there is a huge gap between the area that young athlete's data is plotted in comparison to High Performing athletes. The solution is to include National level athlete's data that fall roughly in the gap between the young athlete and the High Performing athletes. Having this extra data in the graph a step pattern begins to emerge that looks similar to a milestone chart that most athletes are very familiar with. Figure 6.1.2 below is not an exact representation of the milestones in action however, the idea of filling the gap with national level athletes will have a similar look to it but with overlapping of lines between the national and High Performance areas as some national athletes compete on the world stage.
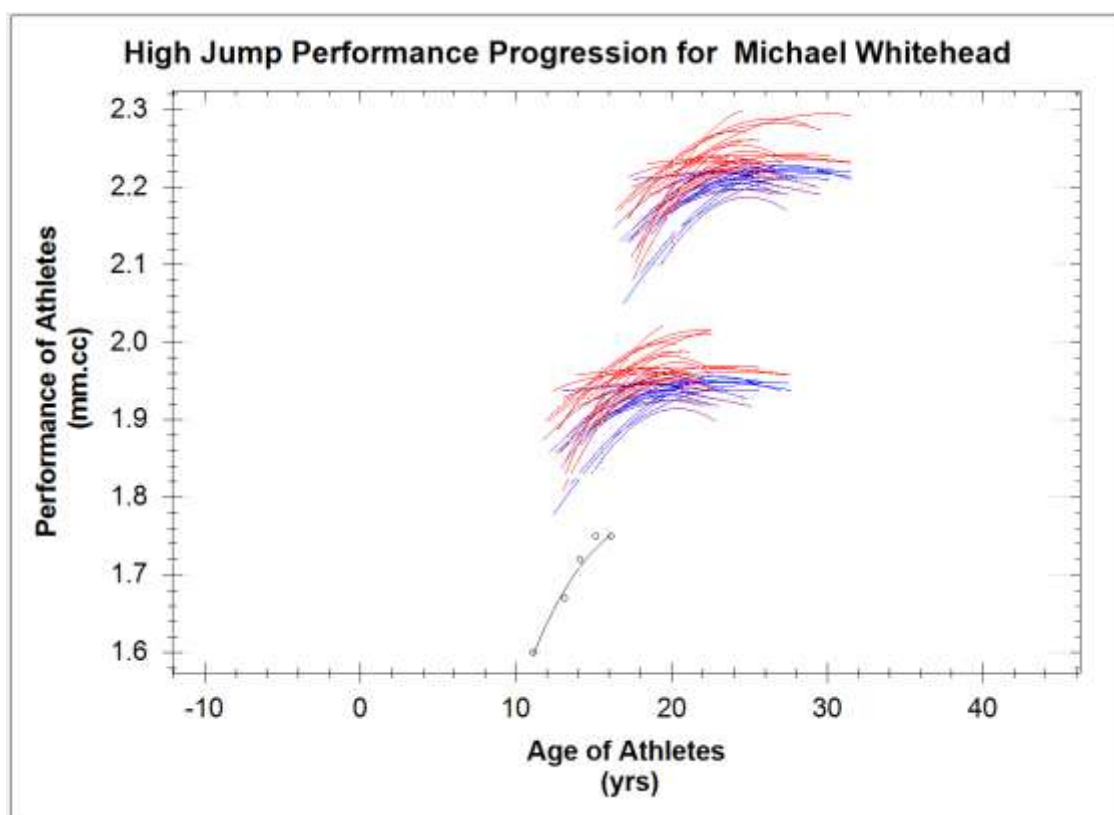


**Figure 6.1.2 Milestone Example for Young athletes plotted performances.**

# 7. HCI DESIGN

When designing the program and how it was to be used, it became apparent that the client and user testing became quite valuable to the design process. To get it to look and work the best for user's five independent people gave feedback in the design of the project.

## 7.1 Testing

The testing of the program was not a formal process. The program was distributed to parties who either were closely connected or had an investment of some sort in the program. Having the view points of independent people on the project meant a more useable system that people could use. The process of finding flaws was very similar to Heuristic Evaluation.

As all the testers were independent of one another, different issues and suggestions would arise to better the usability of the program. The general way the testing was conducted was taken from one of the main principles in the Agile Method Lean which was "Delivering as fast as possible". This ideology works off the idea that you should send your project out into the world to be operated with by users so that you can fix any issues that they find, as opposed to leaving everything to the end and finding out too late that what you have created does not work well.



**Figure 7.1.1 Curve showing the proportion of usability problems in an interface found by heuristic evaluation using various numbers of evaluators.**

## 7.2 Lean

The lean method used in the creation of this program also gave another useful principle "Decide as late as possible". This principle meant that while under creation of the program, instead of guessing or assuming the feature was right as soon as it was created, feedback was returned quickly when a feature of the program was turning down the wrong road. Issues and

bugs could be fixed early on which meant that towards the end of the program, very few things needed changing or removing.

An example of this was around the graph. The client wanted to change the axis minimums which would have taken a considerable amount of time and if they were not to his liking then it would have been a great deal of time wasted. The solution was to create 2 mock up graphs and show the various options the client had available. This proved to be a good decision as the client immediately disliked the look of the new graphs. The outcome was time spent on features the user wanted rather than wasting time on areas of the program that the user was not going to be happy with.

The reason for choosing lean to work on this program was due to past experiences but also because of this story a friend had said in passing. "A client does not know what they want. They want a painting and they have a rough idea of how it is to look but they cannot do it themselves. It's up to the programmer who has all the paints and brushes to create this painting for them. However, no matter how good the client is at describing their painting it will never look like what they wanted by the end. Instead, the programmer needs to continuously show the painting to the client as much as they can while they're painting so they can get feedback and so they get it to resemble as closely as possible to what the client wants." Lean works by getting the product to the client as quickly as possible to get feedback quickly so that big decisions can be delayed until all the facts are there and it is exactly what the client wants without any uncertainty.

## 7.3 Gaming Interface

One of the big changes to the design of the program was the way users were navigated around the program. Originally the concept was to open one mini form after another that gave you possible areas that a user may want to visit. This would have done the job that the client wanted however, usability wise it was not very friendly and caused some confusion. The comments and specifications that the client and the other testers came back with pointed towards a design very similar to that of a gaming introductory menu. Users want to be lead through an application the first time they use it and maybe a couple more times afterwards, basically having their hand held as they progress through. This is to create an easy learning curve for users using the program for the first time. An example of the intro screen can be seen on the following page in Figures 7.3.1 and 7.3.2.

Flash games on the internet and mobile devices all share a consistent design to them. You could have created the best game in the world but made a poor main menu screen. If users struggle with your main menu then the likely hood of them playing your game is very low. The same rule applies for books. An Eye-Tracking study showed that participants spent an average of 20 seconds on each website. (*1995-2012 ScienceDaily LLC)* Bookstore browsers take ten to twenty seconds on average to decide whether to purchase or re-shelve a book. (2012, Meredith Corporation). Although different sources of information, both have shown through studies that it takes around about 20 seconds for users to decide whether they want to work with a product or not. The same applies for programs. With the game menu style design, users can decide what they want to do with the options available and decide if they want to use the program or not. The idea is to offer as much information as the user requires to assess whether they will use the software or not.
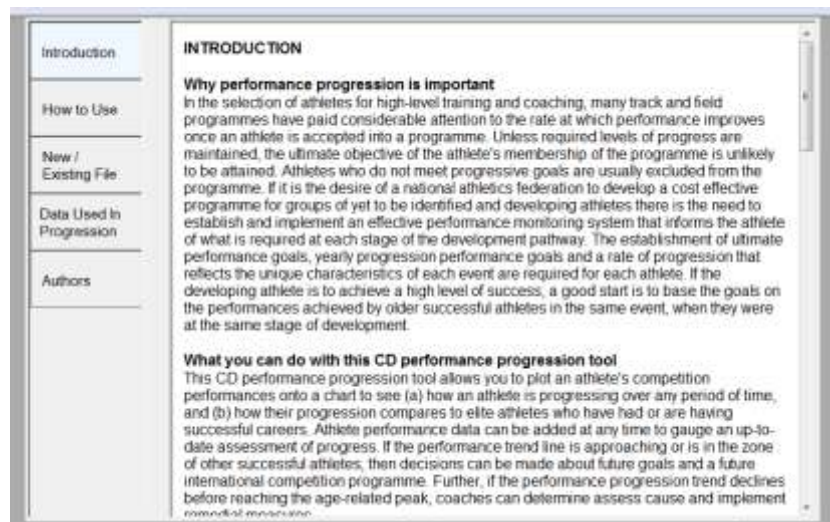


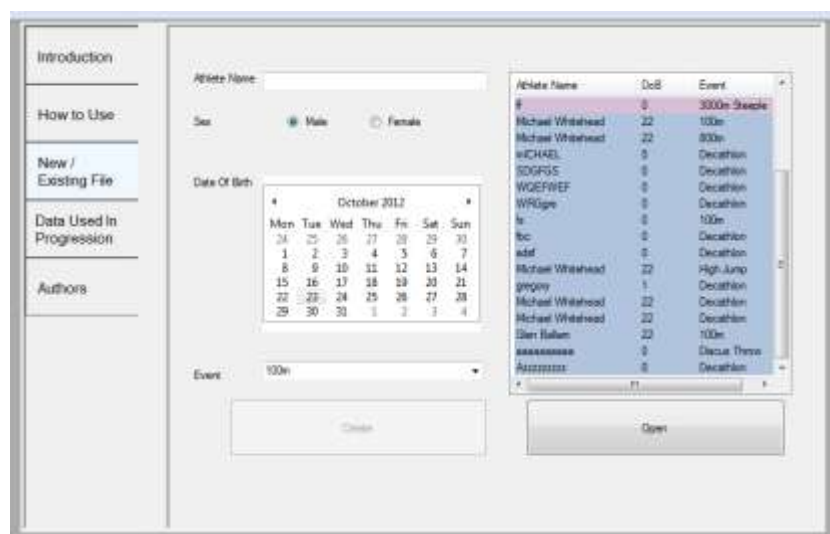Figure 7.3.1 Menu Opened when no other form is opened



Figure 7.3.2 Menu Opened to Create and Open sessions

**8. Data Collection**

The data used in this project consists of an Athletes name, where they positioned at Olympics or World champs, their age at the day of performance and the performance they achieved. This data is then processed down from n amount of data into three x, y coordinates that are used to draw the trend line of that athletes data. An example of the data used is given in Figure 8.0.1.

| Athlete1 | Status1 | Age1 | Perf1 | Athlete2 | Status2 | Age2 | Perf2 |
|---|---|---|---|---|---|---|---|
| Assefa ... | Medal | 17 | 27:49.9 | Abderrahim... | Final | 25.1 | 27:49.6 |
| Assefa ... | Medal | 20.6 | 27:36.4 | Abderrahim... | Final | 27.3 | 27:47.9 |
| Assefa ... | Medal | 24.2 | 27:30.7 | Abderrahim... | Final | 29.6 | 27:49.2 |
| Charles ... | Medal | 21.3 | 27:36.4 | Abdi ... | Final | 21.3 | 28:10.4 |
| Charles ... | Medal | 23.8 | 27:30.5 | Abdi ... | Final | 25.3 | 27:54.0 |
| Charles ... | Medal | 26.3 | 27:28.2 | Abdi ... | Final | 29.3 | 27:47.2 |
| Haile ... | Medal | 19.4 | 27:36.1 | Abdullah ... | Final | 20.1 | 27:38.9 |
| Haile ... | Medal | 25 | 27:06.2 | Abdullah ... | Final | 22.6 | 27:35.7 |
| Haile ... | Medal | 30.5 | 26:54.5 | Abdullah ... | Final | 25.1 | 27:36.2 |

**Figure 8.0.1 Data stored in the Excel Documents**

The reason to simplify the data from a large list down to three points for users is because on average there are 83 performances per athlete that has been collected. The trend lines are the important part of the data that was being collected. Having every point is meaningless when the graph is only concerned with displaying the lines. Later on, when more data on high performing athletes are collected there will be more lines to draw. The only points that need to be displayed on the graph are that of the users. The three points represent the beginning and end of the trend line, where the middle point is the level of curve given to the line.

## 8.1 Manual Data Collection

Originally, when collecting the data the process was done manually. This meant that athletes who had been caught cheating i.e. drug use were not collected as it was not a true indication of their performances over time. It also gave the opportunity to work through the data to make sure that the data on the database was legitimate and did not have any bad data. From the process, although a good deal of time was taken, it did remove a lot of bad data. Also it was a way to familiarise with the data and get a good understanding of what to expect when working with it in the programs.

## 8.2 Process of Collection

The process of collecting the data was performed in two parts. The first was to collect the results from the Olympic and World Championships for athletics and collect the top 16 place getters for all the track and field events for both men and woman. The track events consisted of 10 events ranging from the 100metres to the 10,000 metres. The field events consisted of 9 events including the decathlon/ heptathlon. In total, 2,017 competitors and 168,576 performances were collected over all 38 events between the years 2000 to 2010.

The first process of finding out who had competed and achieved top 16 was the most difficult. It involved going to the International Association of Athletics Federations (IAAF) website and working through their links to the various events that were held. The IAAF is the international governing body for the sports of athletics. The Olympics and World champs are held by the country selected by the IAAF and as a result, each Olympics or World Championships has its own website associated with the country and results performed there. The IAAF keep a record of the links to each site so interested parties can find an event regardless of what country they are from.

After collecting the names the next step was to collect the data on the athlete's performances over their careers. This was done by using the Tilastopaja website which is database containing information on any athlete who performed in an international event. Athletes who perform at an official event understand that they have given permission for their results and name to be used by the athletic community. Data on this site is open to the public so there are no privacy issues that arise. The site is also kept up to date by its members so data from all over the world is regularly added.

Collecting the data was a slow and tedious process that would not make for a good approach to collecting the data in the future. The process consisted on locating an athlete who, either changed their name (Marriage), name that was listed may have been spelt different somewhere else, changed their country or whose date of birth was not properly recorded and resulted in multiple birthdates. Also there are many people with the same name making the process no easier. Once the athlete was found the next step was to find what years they performed and collect data for every year located on different pages. If they were found with doping offences they were immediately removed from the study as it would not be a fair comparison between athletes up and coming with the High Performing athletes.

Finally once all the data for an event was collected, the data was put through a program that returned 3 points for an athlete's data to create a trend line. The data also was given one of three titles to differentiate the "class" of the athlete. $1^{st}$ to $3^{rd}$ was 'Medal' status, $4^{th}$ to $8^{th}$ was 'Final' and $9^{th}$ to $16^{th}$ was 'Other'.

## 8.3 Web Scrapper

The process of collecting the data manually is very difficult and time consuming. The solution to making this process quicker for the future was to create a web scraper that went through collecting all the data on a single athlete. The program that collects the data currently only works for the Tilastopaja website as out of all the sources on the web that data is collected, their formatting is consistent within their site as well as they have the most data to collect. This is a program outside the main graph display as it seemed unnecessary for users to have to download and process raw data. The intention is to have a database that users can connect to get the latest data. The data being an athlete with 3 points for a trend line to be created from.

Tilastopaja does not have an API to connect you to their database so web scraping becomes the easiest way to collect the data. Also, every athlete on their site is given a particular ID so that no two athletes can ever be the same. This means that the scraper cannot simply go onto a page based on a name. Therefore, using the sites search engine and other GUI to navigate around the site was used.

The athlete name is used in their search engine which then searches through the database and returns one of two screens. Either the athletes found after only using a name or, you are sent to a results page that shows athletes that have similarities to the name you searched or, there are no results for the athlete you have chosen. To find the correct athlete out of the list refining the search is needed to locate the appropriate one. This is done by searching through the list for names of athletes who share the same country and date of birth. If these filters turn up tails the next step is to ask for a different athlete. The country is not as important of a filter as the date of birth as the country of the athlete may be regularly changed which is not regularly updated. However the date of birth is normally consistent over the course of their careers. There is a potential loop hole that arises often if the search has returned no exact match after the very first search. The database does contain some phantom athletes that have been added and have not yet been removed by the moderators. In the event of this happening, going onto this page profile and searching for results on the athlete for a particular event will

normally tell whether this is the correct athlete or not. Finally, for extreme precautions, the athlete will have as an absolute minimum of 2 years worth of results. The assumption can be made as the athletes that are of interest have competed at a World or Olympic event, meaning that to be selected or even noticed in the IAAF or countries eyes, an athlete must perform highly for a minimum of 3 years prior.

If the athlete is located and is the correct athlete based off the search then the data can be collected. The profiles work by displaying all the athletes career highlights and on the same page the selected year (by default this is the current year) of all their performances for that year. To collect every year that they performed means navigating the page to every year's data. The amount of years stored begin at 1970 meaning a great deal of pages for more recent athletes do not have any data on them.

To collect each year for an athlete the design originally consisted of using Threading to collect the data. This was to download the data quickly so that athletes waiting to be collected would not have to wait long before being processed. However, as the collection process is a onetime event that is done by one person as opposed to every user using the software, time is not an issue.

## 8.4 Future Data
A problem will begin to arise over the years as more athletes become world class. The graph will begin to fill with so many lines that the information that the graph gives will have diminishing returns. At this stage the plan is to display a limited amount of athletes on the graph at one time. This could be done in various ways however; the envisioned proposal is to display athletes who achieved World or Olympic status over the previous 10 years. This is so the data is regularly updated for the athletes of those generations. Displaying athletes of 20 years ago will hold little meaning to people unless they were record holders. Another possibility that may become an option for users is to display the recent 10 years of athletes and record holding athletes. Athletes are not always satisfied with only being the best, most are after the title of World Record holder, putting them into a new league of their own. As this is highly sort after, it would be useful to include it in the graph for those athletes wanting the next level.

## 8.5 Milestones

Milestones were mentioned earlier as a potential addition to the graph. As the graph is displaying a comparison for athletes with High Performing athletes, a milestone feature is very beneficial for athletes to track their progress. The data used to create this effect will be obtained from the Athletics New Zealand Rankings website. Currently national level athletes who's results meet a certain level have their times recorded on a fairly new ranking page within the Athletics New Zealand website. At the moment, the page is in its infant stage with very little data however, in 3-5 years time there will be plenty of data to collect that will give better trend lines for athletes to compare against if they are at a lower or similar level.

The data is still being stored in Excel. This is done so because the initial design of the project was created in Excel by the client's request. Even when creating the second program using C# there was no need to change it as the program was using the graph drawing feature that Excel offered. For the third design most of the focus was around the graph so just displaying the correct information was priority. As there is not a lot of information stored at this stage, storing the data in Excel was the simplest solution for now.

For future development there will be a lot of information being collected and sent out. Having all this data spread out over 38 Excel files could become quite disjointed. Keeping the information in one location and updating people's data storage would be the ideal solution. The data storage would be looked after by using mySQL and Structured Query Language (SQL) for the communication. This is done because it will be easier to send requests direct to the database for a specific set of information and update the local database rather than download a whole Excel document for data to update the program. This will work well if users begin sharing data and request either specific sets of athletes or requesting an individual's dataset to be added into their data.

# 9. GRAPH DEVELOPMENT

The first version of the graph was done in Excel. The way of displaying the data of the athletes and user data was to use second order polynomials. This was because of Steve's research into athlete performances over their careers. For athletes, the general process of their careers works with 2 major spikes in their lives, one while being a teenager, the second at around but not limited to their mid twenties. The intention for Steve's work was to work around the second spike. The main reason for this was his research found that younger athletes who were the best in their level would generally leave for various reasons. The focus was around the athletes who were going to hit the High Performing level for this program. The other reason for the quadratic line was to show the athletes general performance ability around a certain age to give a good estimate of their ability at a particular age.

The second version that used C# Excel integration used an Excel wrapper to link to the Excel spreadsheet where the graph was displayed. Unfortunately there were some issues regarding the window display. Basically the window would occasionally not send the entire user interface through like the scroll bar which made it difficult to navigate around the graph. Also there were varying degrees of change in the appearance of graphs on different people's computers. What looked correct on one display looked semi assembled on another.



Athletic Level

Version 1 — High Performance / International

Version 2 — National; Youth National / Regional
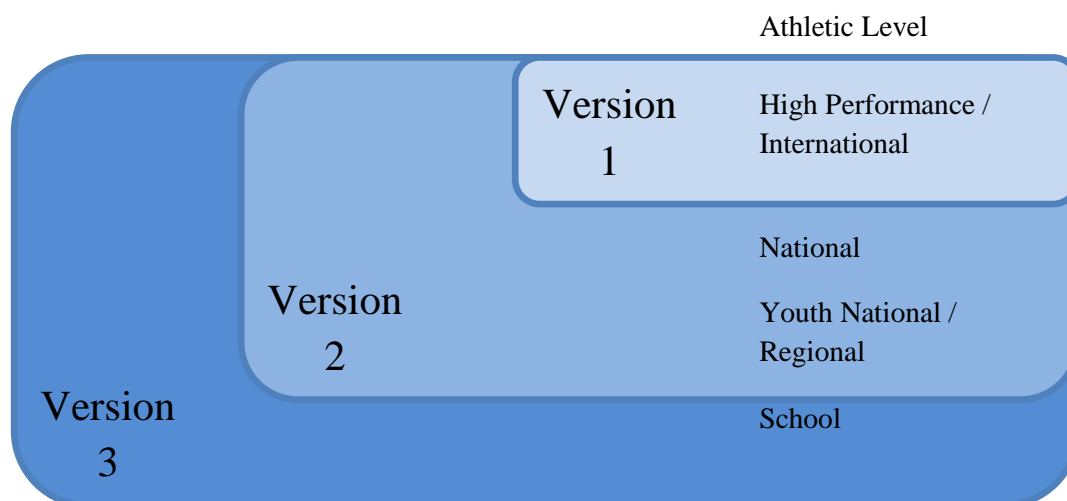
Version 3 — School

Figure 9.0.1 Progress of 3 Programs and their target audiences

The second program used the Excel graph to display the data. In this version it was also upgraded to include more athlete performances. Originally the program was focusing on High Performing athletes in New Zealand however; it managed to miss the National level athletes who have the potential to become high performing athletes. The way it missed these athletes was the way the graph was formatted. The graph was designed to display the data in the best

way possible so that people using the system could get decent information out of it. The curves 'look' and difference from another curve was the target here. Having them clustered together meant no useful information could be obtained. The issue was that although High performing athlete's data points would always show on this graph, it was not suited for athletes who were up and coming who had data points that fell just outside the x-y axis of the graph.

The solution was to change the x-y ranges of the axis so that more information could be displayed. There was two potential fixes. The first was to manually change the graph to static values, the second was to create a dynamic graph that could be panned and zoomed. The client was not too enthusiastic with the dynamic graph so the alteration of the x-y axis maximum values was done. The reason for the dynamic graph is for the precision that it offers through zooming and panning functionality. An example of this is shown below in Figure 9
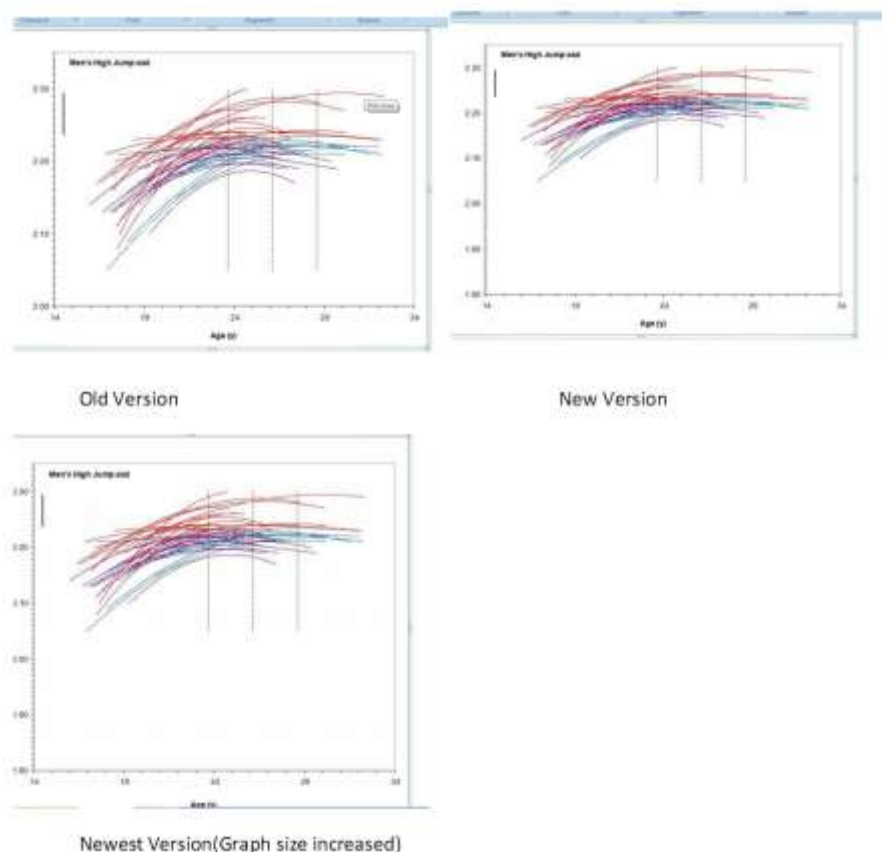


Old Version                                    New Version



Newest Version(Graph size increased)

**Figure 9.0.1 Axis Alterations**

The graph needed to include athlete's data which was of a national standard but also needed to keep the aspect ratio to keep the curves of the data to give useful information. The solution was to gather the New Zealand ranking lists minimum standards and get a middle value

between the old graph maximums and the ranking minimums. Unfortunately, even doing this made the graphs visually return different information. The end solution was to resize the graphs as well to get nice looking curves but also cause there was a lot of room in the Excel wrapper that was going to waste now that the data was being added to a data grid in C# as opposed to putting it straight into Excel.

There are disadvantages to using Excel to display the data in a graph. The main issue being that the graph is static and does not easily change its maximum and minimum values for its axis. Although it was not a must for the client to have a dynamic graph I decided to proceed with it to see if there were advantages to creating this. There happens to be a huge advantage to showing data on a dynamic graph that was otherwise unable to be shown on a static graph.

The first big advantage is the ability to map any data points. This means being able to include data for any person and display it next to the data on the graph without compromising the curve of the pre-existing data. The graph can be zoomed in and out to show all data or just a selection. It also has functionality to pan around to get the area of data that you want to display. The excel version does not give any of this freedom to work with the graph, it only shows all the data and assumes that fitting all the data in, or leaving data outside the graph is what the user wants from it.

Another advantage that's given from this dynamic graph is the ability to show milestones. With the Excel version there is the issue of it either showing you everything or only a section of all the data. If it shows all the data then the user cannot see the data up close and they lose some of the valuable information by having the data zoomed in. With the dynamic graph, users can zoom out to show the whole data as an overview before zooming in and panning around the various data sets to demonstrate milestones.

Zedgraph had basic calculations to create standard graphs however; if you wanted to create more complex looking graphs then you would need to add in your own calculations to create them. A second order polynomial line was one of these calculations. Normally when working with a line of best fit through a set of points on a graph the line is straight. For this version however Steve wanted to show a rise and fall of performances of an athlete's career which is why a second order polynomial was used. Had a straight line been used the data being displayed would not have given a good representation of an athlete's career. Some lines would rise, some would fall and most would flatten out giving no real insight as to the athlete's peak performances during their career.

## 9.1 2$^{\text{nd}}$ Order Polynomial

To calculate the data to get the best fitting curve f(x) has the least square error,

$$\Pi = \sum_{i=1}^{n} [y_i - f(x_i)]^2 = \sum_{i=1}^{n} [y_i - (a + bx_i + cx_i^2)]^2 = \min$$

Note that $a$, $b$ and $c$ are coefficients while all $x_i$ and $y_i$ are given values; $x_i$ is all performances and $y_i$ is the age at the performance. To obtain the least squares error, the unknown coefficients $a$, $b$ and $c$ must yield zero first derivatives.

$$\frac{\partial \Pi}{\partial a} = 2 \sum_{i=1}^{n} [y_i - (a + bx_i + cx_i^2)] = 0$$

$$\frac{\partial \Pi}{\partial b} = 2 \sum_{i=1}^{n} x_i [y_i - (a + bx_i + cx_i^2)] = 0$$

$$\frac{\partial \Pi}{\partial c} = 2 \sum_{i=1}^{n} x_i^2 [y_i - (a + bx_i + cx_i^2)] = 0$$

Expanding the above equations we get,

$$\sum_{i=1}^{n} y_i = a \sum_{i=1}^{n} 1 + b \sum_{i=1}^{n} x_i + c \sum_{i=1}^{n} x_i^2$$

$$\sum_{i=1}^{n} x_i y_i = a \sum_{i=1}^{n} x_i + b \sum_{i=1}^{n} x_i^2 + c \sum_{i=1}^{n} x_i^3$$

$$\sum_{i=1}^{n} x_i^2 y_i = a \sum_{i=1}^{n} x_i^2 + b \sum_{i=1}^{n} x_i^3 + c \sum_{i=1}^{n} x_i^4$$

Once the 3x3 matrix (A) and 1x3matrix (B) are obtained, a simple calculation is performed to obtain $a$, $b$, $c$ matrix (C)

$$A^{-1} \times B = C$$

Once $a$, $b$, $c$ are obtained, you can find the trend line for any amount of data and draw the line using the equation

$$y = ax^2 + bx + c$$

| Age | Performance | Predicted Performance |
|-----|-------------|-----------------------|
| 14.6 | 1.70 | 1.61 |
| 15 | 1.50 | 1.71 |
| 15.2 | 1.83 | 1.76 |
| 15.8 | 2.00 | 1.88 |
| 16.2 | 1.90 | 1.95 |
| 16.5 | 1.93 | 1.99 |
| 17.1 | 2.05 | 2.04 |
| 17.4 | 2.12 | 2.06 |
| 17.6 | 2.07 | 2.06 |
| 17.9 | 2.02 | 2.06 |

MATRIX 3x3 (A)

| 10 | 163.3 | 2678.87 |
|-----|-------|---------|
| 163.3 | 2678.87 | 41030.1 |
| 2678.8 | 41030.1 | 634491 |

INVERSE MATRIX 3x3 (A)

| -0.193708 | -0.0751258 | 0.00567 |
|-----------|------------|---------|
| -0.075125 | 0.00991158 | -0.0003 |
| 0.0056759 | -0.0003237 | -1.5E-06 |

MATRIX 1x3 (B)

| 19.12 |
|-------|
| 313.89 |
| 5175.75 |

MATRIX 1x3 (C)

| a | -0.0432872 |
|---|------------|
| b | 1.54499912 |
| c | -11.721756 |

**Figure 9.1.1 Tables showing the calculations over a dataset**

The difference between the above calculations to draw the new trend lines in comparison to how Excel calculates the trend lines are very minute in difference. The graph below in Figure 9.1.2 shows an example of a plotted graph with random values. The green line shows a trend line drawn off the calculations that Excel has built in. The red line demonstrates the result from the before mentioned calculation. As you can see there is a slight difference in the dip of the graph as well as where the starting and ending points are drawn. This can be explained by the rounding of values Excel does. However the difference is so minute that it will not have an effect on the way the program functions at all.
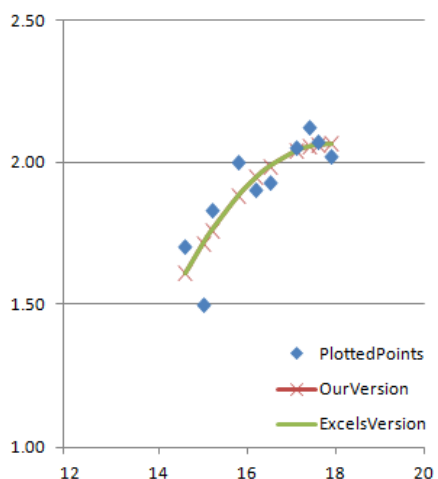


**Figure 9.1.2 Plotted data from Figure 9.1.1 calculation**

The program works by giving you an estimate of the athlete's potential and is used for prediction; there is no need for the graph to be 100% precise, 99% works just as well. The prediction has a standard deviation from the line for error. The other reason is that the future version of the program will include an error funnel technique which will not be affected by the slight difference in the curve.
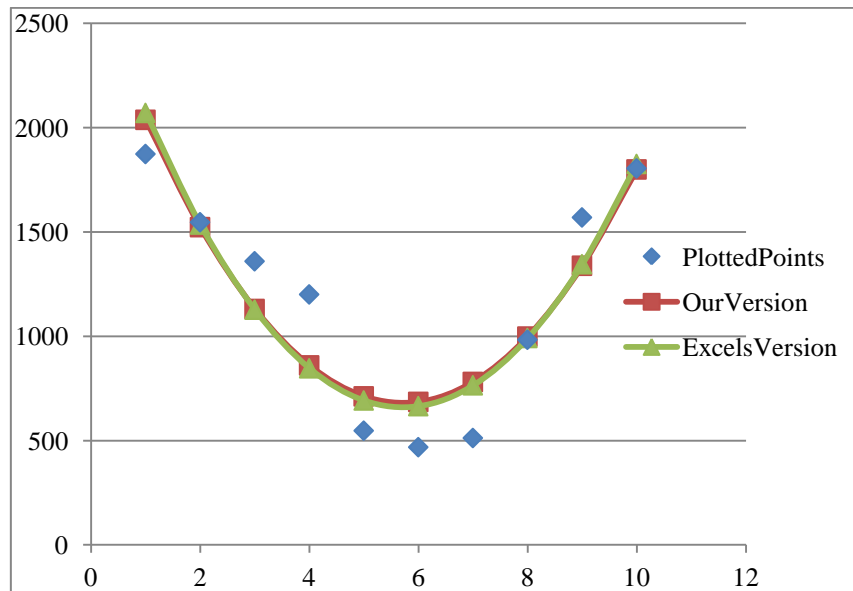


**Figure 9.1.3 Plotted Data comparing the calculation against Excels calculations**

## 9.2 Error Funnel

Excel does very well in storing and displaying large sums of data however, there is a limit to how much it can display and hold. C# on the other hand is only limited by the computer's ability to compute. One of the features that were sort after was an error funnel.

The idea of the error funnel is to give users an indication of the potential future an athlete has if they continue to progress on the path they have chosen.

There are two parts to the error funnel for this graph; one element is to do re-sampling to create a funnel similar to Figure 9.2.1 below representing the 95% confidence interval of the data. The other element is to extrapolate the data to give a potential future of what the athlete may potentially do.
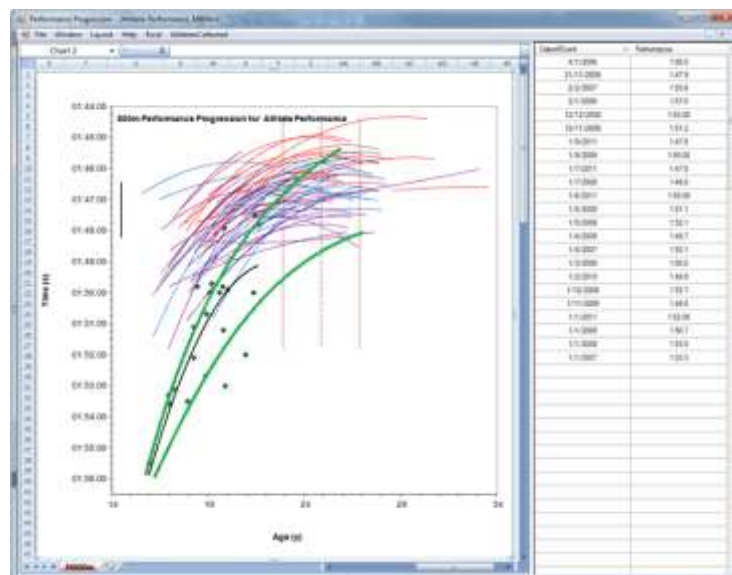


**Figure 9.2.1 Error Funnel**

There are various ways to predict off the data. The chosen method is to use re-sampling and extrapolation. The idea was to work off the data that was already collected off all the high performing athletes.

How this would work is the end of the data that the user had put in would be the starting point of the future development. The program would then look at all the athletes from that age onwards and get an average quadratic line. That line would then be added to the end of the user's athlete to give a prediction of what they may potentially do from that age onwards. This method to predict the athlete's potential future had its positives, however, as the program started heading towards having High Performing down to a College level of athletes using it, the prediction would not work for all users. The younger athlete's predictions would

be highly inaccurate due to how young athletes develop. So the process would only be useful for the high performing level athletes. The other way of predicting athletes performances would be more individualised to the athlete, focussing on their performances rather than what other athletes have done.

### 9.21 Re-sampling

Re-sampling is achieving one of three methods. For this issue the method of estimating the precision of sample statistics by using subsets of available data or drawing randomly with replacement from a set of data points (bootstrapping). This will give a funnel which will contain a 95% confidence interval of all the potential trend lines of the user's data to show where they could potentially go with their performances so far.

### 9.22 Extrapolation

Extrapolation is the other technique used to get a sense of future potential for the user. This technique does come with risks. The further from the actual data the extrapolation is the less reliable it becomes. This is not a huge issue though. Selectors, sponsors, athletes and coaches may be interested in only up to 4 years later but more likely interested in a year or under. This will mean that the prediction will be fairly accurate for the area that users are interested in. Also, as most are interested in 4years or less, the prediction will only go to 4 years to lessen the potential error. The reason for 4 years is this is the normal planning timeframe that coaches use to determine the athletes training. They also use milestones of 4 years as Olympics are held every 4 years. World championships are every year except for the year that the Olympics are held.

### 9.23 Bootstrapping

Bootstrapping works by sampling data from all the data on a group. Get the sample data and create a new trend line off it. Do this around 10,000 times. The graphs on the following page in Figure 9.23.1 have not been taken from the athletes data, nor are they using $2^{nd}$ order polynomials however, they display the process of bootstrapping that will be incorperated into the program. The top left graph represents a collection of data with a line drawn using the least squares method and a $n^{th}$ order polynomial. The graph on the right has used the bootstrapping method to create 1,000 lines using a sample of the total data. The result is 1,000 $n^{th}$ order polynomial lines.
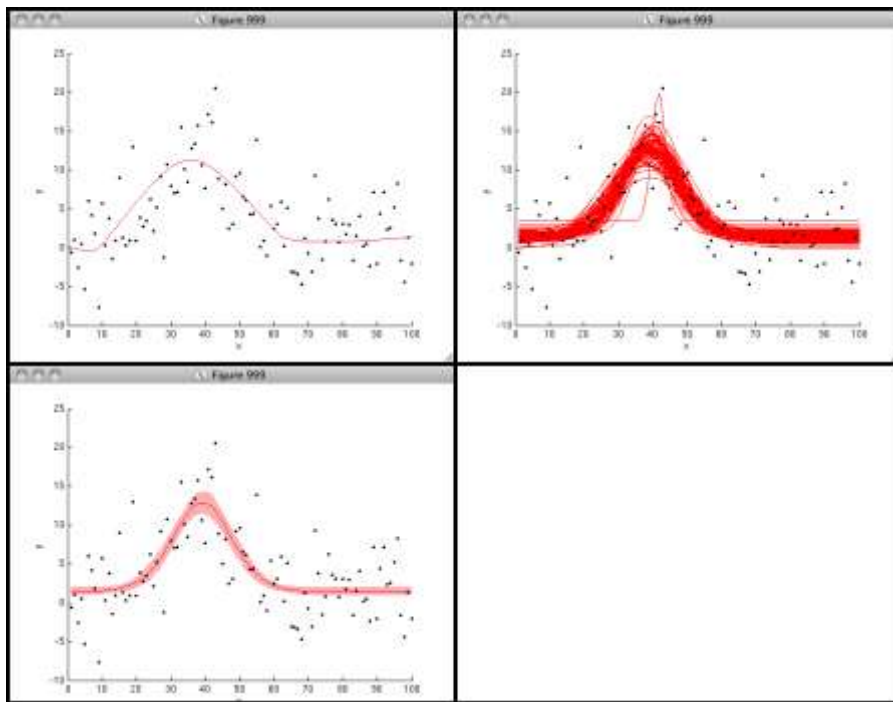
**Figure 9.23.1 graphs showing Bootstrapping method**

The graph in the bottom shows a solid red line representing the mean of the 1000 sample trendlines and the light red funnel around it is the standard deviation of the 1000 samples. The intention is to do a similar process with a users data, only show them the final funnel view rather than the process of views leading up to it.

To calculate the maximum and minimum for the funnels bounds, the program goes through a standard interval along the x axis and finds the maximum and minimum points from the 10,000 trend lines. This will return the x, y coordinates for the maximum and minimum trend lines of the funnel. From the example there are some outliers that do not fit well with the rest of the lines. A standard deviation would be the solution to removing these outliers and giving more precision to the funnel. The way to do this would be to go 2.5% of the lines in from the top and bottom of the available trend lines to get the standard deviation of all the trend lines. An example of the standard deviation of the lines is shown in Figure 9.23.1 on the bottom left.

The last step to show a potential future and to change the look of the funnel from a cylinder shape to a funnel shape is to implement extrapolation. The way to calculate this is to use the coefficients a, b, c from the equations that were created for the maximum, main and minimum trend lines. Then extending the lines by the before mentioned 4 years to give a 'funnel look' demonstrating where the data that the user has supplied is potentially heading.
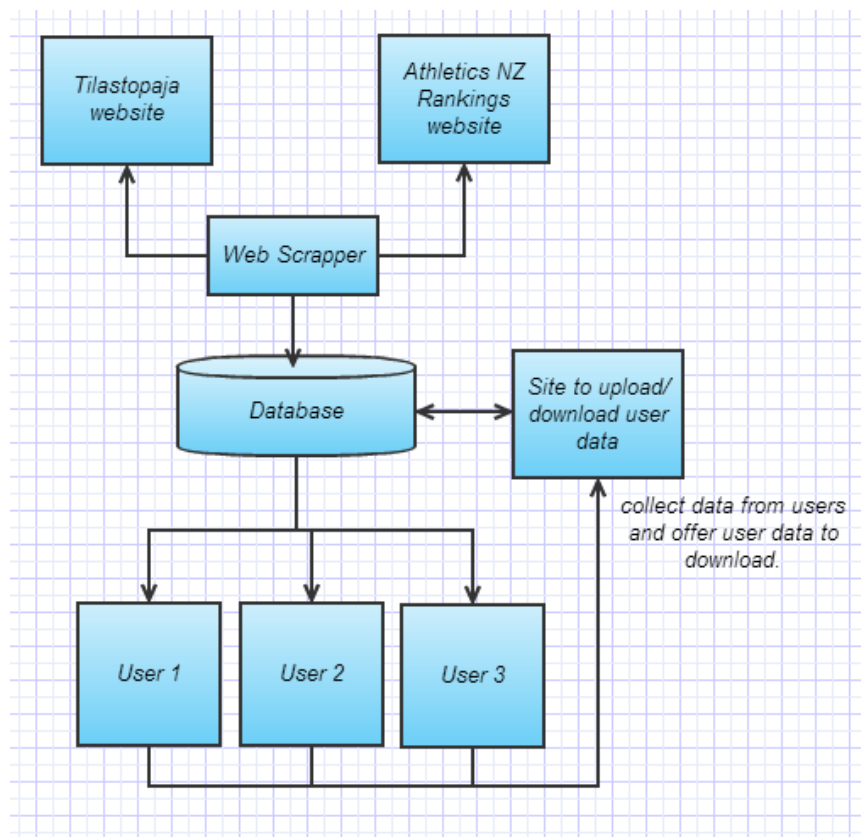
36

## 10. Reflection



**Figure 10.0.1**

## 10.1 Future Work

The future of the project structure is to combine the techniques already discussed with connecting users with one another and sharing data on athletes if they choose to. How it will work is; there is a web scraper that downloads data of High Performing athletes from the Tilastopaja website. It also downloads data from the Athletics NZ rankings website to get National level athletes data. The scrapper stores this information in a database. The main program is distributed to users who get yearly updates of the data. They can then create their own trend lines of athlete's performances which can be uploaded to a site for other users to download.

The program only needs yearly or season updates as athletics like many sports is a seasonal activity. During season anyone using this software is going to get little gain from seeing daily or weekly updated trend lines for athletes for two reasons.

The first is that athletes have already achieved the build up that they are going into the season with. The program does not give you information to work off that will help you each time you put new data in. It works by giving an overview of all your data. Adding one entry is important but generally will not alter the graph alone, the calculations work better with more data.

The second reason is people who use this software are using it to get a picture of where an athlete stacks up with the rest of the competition at the same age. The data is used as a predictor to calculate whether an athlete is matching up to their fellow competitors and whether they are going to rise or fall in the future based on the results from the past.

As far as data privacy is concerned, users who upload data have the right to remove the data if they so wish. The data that they share and the data collected on High Performing and National level athletes are public source anyway. Any person who performs in a club level or higher event agrees for their results to be displayed to everyone when they perform otherwise it is not counted as a legitimate performance.

Athletes who are not on a world or national level may want to see what their level is currently doing and compare themselves with their competition. This will generally affect younger athletes of a regional and national level and their coaches. There will be the potential for users to upload fake or misleading data. As this is a possibility it will be up to users to decide who they download from. A rating and comment section could be added to help users download from the best up-loaders for better data.

## 10.2 Evaluation

Athletics New Zealand has bought the program and has begun implementing it into their selection process.

# References

(Luke M. Chandoo.org 2009)

http://chandoo.org/forums/topic/excel-2007-security-best-practices#post-39583

(2005 by Jakob Nielsen)

http://www.useit.com/papers/heuristic/heuristic_list.html

*(1995-2012 ScienceDaily LLC)*

http://www.sciencedaily.com/releases/2012/02/120216094726.htm

(2012, Meredith Corporation)

http://www.divinecaroline.com/122444/100663-judging-book-its-cover-sells.

(2012 eFunda, Inc.)

http://www.efunda.com/math/leastsquares/lstsqr2dcurve.cfm

(OLE DB)

http://msdn.microsoft.com/en-us/library/windows/desktop/ms722784(v=vs.85).aspx