

# COMP390 Evolving a Sorting Algorithm with SNGP

Robin Lockyer

*Student ID: 201148882*

*Primary Supervisor: Dr. David Jackson*

*Secondary Supervisor: Dr. Valentina Tamma*

2018/2019



### **Abstract**

Genetic programming is a technique for creating programmes not by writing them by hand, but instead by creating a population of random programmes and modifying them using an evolutionary algorithm. The desired result is that after several generations a programme that performs well at a given task is generated. GP has previously been used to successfully evolve sorting algorithms.

Single node genetic programming is a variation on GP invented by Dr Jackson which structures the population of programmes in a manner that allows the use of dynamic programming when computing the result of the programmes in an effort to more efficiently generate a working solution.

This project aims to compare the effectiveness of the two methods in evolving a sorting algorithm.

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Background</b>	<b>4</b>
2.1	Standard Genetic Programming . . . . .	4
2.2	Single Node Genetic Programming . . . . .	6
<b>3</b>	<b>Design and Realisation</b>	<b>9</b>
3.1	GP Implementation . . . . .	9
3.2	SNGP Implementation . . . . .	12
3.3	Fitness Function . . . . .	14
3.4	Test Data . . . . .	16
<b>4</b>	<b>Results</b>	<b>17</b>
<b>5</b>	<b>Evaluation and Learning Points</b>	<b>22</b>
<b>6</b>	<b>Professional Issues and Data Required</b>	<b>23</b>
<b>7</b>	<b>Bibliography</b>	<b>24</b>
<b>8</b>	<b>Appendices</b>	<b>24</b>

# 1 Introduction

This project is was done for my project supervisor Dr David Jackson. The aim of this project is to attempt to evolve a sorting algorithm using node genetic programming (SNGP) and, if successful, compare the effectiveness of evolving sorting algorithms using standard genetic programming (GP) to evolving sorts with SNGP.

The purpose of GP is to automate the creation of algorithms and programmes. This is done by applying a genetic algorithm to a population of random programmes so that successive generations of programmes improve at the desired characteristics until a functional programme is created. The standard approach to GP requires evaluating hundreds of programmes per generation over potentially thousands of generations and as such GP can take up a large amount of processing time. Several variations of GP have been created that try to reduce the amount of processing, including Linear Genetic Programming and Parallel Distributed GP [5].

SNGP is one such variation devised by Dr Jackson in *A New, Node-Focused Model for Genetic Programming* [1]. This variation makes use of a form of dynamic programming to re-use results of previously evaluated programmes. It has been shown that SNGP tends to perform better than standard GP in terms of processing time, solution rate, and solution size [1]. SNGP has reduced efficiency when dealing with problems with side-effects because this prevents re-use of evaluations, although it still performs better than GP at some problems with side effects [2].

A sorting algorithm reads and manipulates an array of integers as it executes, and so must make use of side effects. Regular genetic programming has been shown to be capable of evolving a working sorting algorithm [3, 4]. This makes evolving a sort a good problem to evaluate SNGP on as comparisons can be made with previous research.

## 2 Background

### 2.1 Standard Genetic Programming

Genetic programming allows the user to automate finding solutions to problems without the need to know much about the solutions themselves. This is done by a stochastic process where successive generations of programmes are altered to create the next generation, with the goal of each generation being an improvement on the last until a desired result is found.

In standard the standard form of GP, described in *A Field Guide to Genetic programming* [5], programmes are encoded as a tree of primitive functions and terminals. Each function has a number of child nodes equal to it's arity, and each individual child represents a single operand of the parent function. Nodes with no children represent terminal functions or constants which require no input. Figure 2.1 shows an example of a programme encoded as a tree.

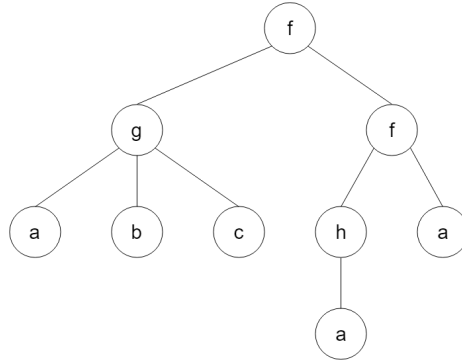


Figure 2.1: This tree encodes the programme  $f( g( a, b, c ), f( h( a ), a ) )$ , where  $f$ ,  $g$ , and  $h$  are functions and  $a$ ,  $b$ , and  $c$  are terminals.

An initial population of random programme trees is generated according to certain parameters that can vary depending on implementation. Each member of the population is executed in turn. The results of these executions are evaluated and each member of the population is given a fitness score so that better programmes have higher scores.

Members of the population are then selected to produce offspring programs based on their fitness. Simply selecting the best individuals to reproduce can lead to low diversity in the resulting population, so a probabilistic is often used to prevent a single programme's offspring from dominating subsequent generations. Common methods for doing this include fitness proportionate selection, where the probability of selection is proportional to the fitness of the individual relative to the fitness of the whole population, and tournament selection, where a small subset of programmes are chosen with equal probability and the fittest of the subset is selected for reproduction.

Genetic operators are then applied to the chosen programmes to create a new generation. The operator selected to create each member is chosen with a pre selected probability. The three most commonly used operators, which are also the three operators Koza used to successfully evolve a sort [3], are:

- **Reproduction:**

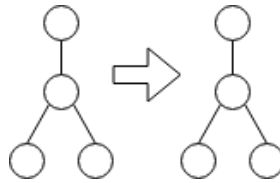


Figure 2.2: Example of reproduction operator

The selected programme is copied into the new generation without modi-

fication. Reproduction tends to have a low selection probability to ensure the following generation is sufficiently different.

- **Crossover:**

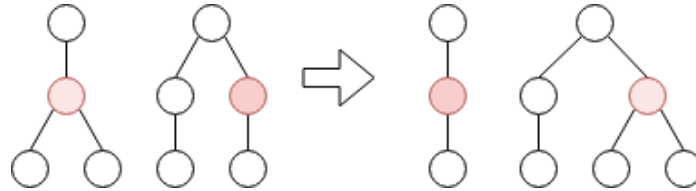


Figure 2.3: Example of crossover operation. The highlighted nodes are the crossover points

A random node is selected as the crossover point in each of the two chosen programmes and the subtrees rooted at the selected nodes are swapped.

- **Mutation:**

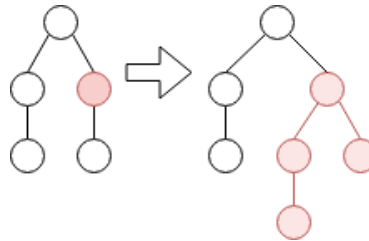


Figure 2.4: Example of mutation operator. The highlighted node is replaced by the highlighted subtree.

A random node is selected in the chosen programme. A new, random subtree is generated to replace the subtree rooted at the selected node.

The process of creating new generation is continued until a certain condition is reached, usually after a certain number of generations or a programme exceeds a fitness specified by the user.

Genetic programming has been previously shown to evolve a sorting algorithm by Kenneth E Kinnear [4, 3]. His work was not focused on evolving an optimal sort, but instead attempting to see if GP was capable of evolving any kind of sort at all. Kinnear considered sorting to be an interesting problem to solve using GP as there are many possible ways to implement a sorting algorithm, and the domain of the problem is infinite.

## 2.2 Single Node Genetic Programming

Single Node Genetic Programming is a variant of GP that uses a form of dynamic programming to speed up the evaluation of GP programmes.

The population is organised a single directed acyclic graph. This allows the nodes to be organised into a topological ordering such that any given node only uses nodes later in the ordering as operands. The nodes are stored in an array in this order, with all the terminals occupying the lowest values, as shown in Figure 2.5.

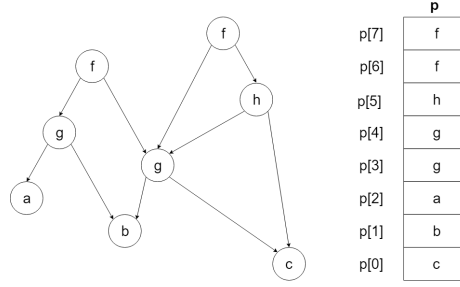


Figure 2.5: This SNGP graph is stored in topological order in an array. Terminals a,b, and c occupy the lowest positions of the array.

Every subtree within the graph represents a single programme. For example the graph in Figure 2.5 contains the following programmes:

- a
- b
- c
- g( a, b )
- h( g( b, c ), c )
- f( g( a, b ), g( b, c ) )
- f( g( b, c ), h( g( b, c ), c ) )

This allows many different programmes to be represented in a single graph as each node is itself a programme, while in regular GP there are many nodes per programme. The structure also allows for sharing of subtrees between programmes so that if there exists a programme with a particularly high fitness, other programmes can make use of the same subtrees to improve their fitness.

Evaluation of the population starts from the subtree rooted at the lowest element of the array and works upwards. The results of each evaluation are saved and can be re-used whenever the subtree needs to be executed as part of a larger programme, making the evaluation of an SNGP population very efficient.

SNGP uses only one genetic operator called *successor mutate*. This operator picks a random function node in the graph and changes an operand with another random node with a lower index in the population array. Instead of the whole population being re-evaluated only the the modified node and programmes that make use of that node need to be evaluated. These nodes can be determined by keeping track of the parents of each node and using this information backtrack up the graph from the modified node. When backtracking to determine which

nodes to update it is important to visit nodes in array order to ensure that subtrees are updated before their results are re-used.

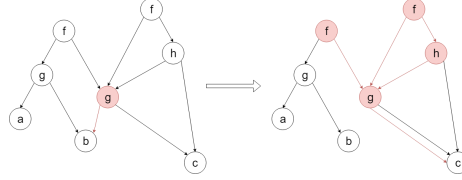


Figure 2.6: A successor mutate operation is applied to the highlighted node on the left graph, resulting in the right graph. Highlighted nodes on the left graph should be re-evaluated

Successor mutate allows for the modification of many programmes with a single genetic operator, and the modification itself is very efficient. Evaluation of the new population tend to be more efficient than in GP as SNGP can make use of the prior evaluation from the unchanged subtrees from the previous generation.

After application of the successor mutate operator the population as a whole is given a fitness value. If the application has improved the fitness value then the change is kept, otherwise the change is reverted. There are variation of SNGP that differ in how to assign a fitness to the population:

- **SNGP/A:** The population fitness is the average fitness of the whole population.
- **SNGP/B:** The population fitness is the fitness of the best individual within the population.

Like GP, this process is repeated for a certain number of generations or a programme exceeds a fitness specified by the user.

SNGP has been shown to produce a working solutions after fewer generations than GP, as well as producing many more working solutions in the same number of generations [1]. The programmes produced tend to be smaller than those in GP, most likely due to the maximum size of any programme being capped by the size of the population, and the re-use of subtrees within a single programme.

SNGP cannot make use of dynamic programming when dealing with primitives with side effects. The result of a subtree may differ depending on what is executed beforehand, so the subtree must be executed again for each other programme that makes use of it. However, even with this disadvantage SNGP/B still requires fewer evaluations to find a solution than GP in some cases, although this is dependant on the specific problem being solved [2].



## 3 Design and Realisation

### 3.1 GP Implementation

The GP implementation was based example code given by Dr Jackson and also on an implementation of the TinyGP system found in the book *A Field Guide to Genetic Programming* [5]. It was written in C as the example code given to me was also in C. The parameters are the same as described in Koza's work [4].

Each programme is stored as a string of primitive functions and terminals. The string is in prefix notation, and as all primitives have fixed arity no bracketing is needed. Programmes also store their fitness and their size. The information is stored as a struct, detailed in Figure 3.1.

Prog	
Primitive[]	code
double	fitness
int	proglen

Figure 3.1: Data structure to store a programme. The *code* array stores a string representing the prefix notation of the programme tree, *fitness* stores the programmes fitness, and *proglen* the number of nodes in the programme tree.

The primitives used for this implementation are the same nine Kinnear used:

- **INDEX:** A terminal that gives the current index of the ITERATE function.
- **LENGTH:** A terminal that gives the length of the current array.
- **ITERATE(start, end, function):** Sets index to the value of start and increments index by one until index equals end or index equals length -1. Each iteration will execute function a single time. This primitive will return the minimum of end and LENGTH. To prevent excessive execution time, programmes are capped at 10000 total iterations.
- **SWAP(x,y):** Swaps the elements of the array at positions x and y. Returns x.
- **SMALLEST(x,y):** Compares the values of the array at x and y. Returns the index of the smaller.
- **LARGEST(x,y):** Compares the values of the array at x and y. Returns the index of the larger.
- **SUBTRACT(x,y):** Returns the value of  $x - y$ .
- **INCREMENT(x):** Returns the value of  $x+1$ .
- **DECREMENT(x):** Returns the value of  $x-1$ .

The primitives are defined in an enumeration, and a table is kept to store the names and arities of the primitives. An additional *DUMMY* primitive is also included to occupy the 0 value of the enumeration and make all actual primitives

have arity  $\geq 0$ . This is useful when debugging and in certain algorithms for identifying when a primitive has not been assigned. These data structures are shown in Figure 3.2.

$$\text{enum Primitive} = \{p1, p2, \dots, pn - 1, pn\}$$

TableEntry	
Primitive	primitive
int	arity
char[]	name

$$\text{TableEntry[] arityTable} = \{\{p1, 0, "PrimitiveName"\}, \dots, \{pn, k, "PrimitiveName"\}\}$$

Figure 3.2: Data structures used to store information about primitives

The each member of the population is initialised by using the grow method. A maximum tree depth is specified and a random primitive is selected as the root. Random primitives are selected as the children for primitives that have been previously selected. The primitives can be either functions or terminals, unless selecting a function would cause the tree to exceed the maximum depth. The initial maximum depth is chosen to be 6 to match Koza's work. The pseudocode for the algorithm is shown in algorithm 1.

---

**Algorithm 1:** Algorithm to initialise a programme as a random tree with a given depth

---

```

Function createTree (maxDepth):
    if maxTreeDepth = 1 then
        | Output: randomTerminal()
    else
        | elseif-block
    end
    root ← randomPrimitive()
    for i ← 1 to arityTable[root].arity do
        | root.operand[i] ← createTree (maxDepth - 1)
    end
    Output: root
end

```

---

Three genetic operators used are the three mentioned in section 2.1. Crossover is used to generate 80% of the next generation, and 20% is generated by reproduction. There is also a 10% change that mutation will be applied to each member of the new population. The pseudocode for this is shown in algorithm 2.

---

**Algorithm 2:** Genetic Programming Algorithm

---

```
initialisePopulation()
evaluatePopulation(0)
for generation  $\leftarrow$  1 to NUM_GENERATIONS do
    for j  $\leftarrow$  to POPULATION_SIZE do
        if randInt(1,10)  $\leq$  8 then
            | newPopulation[j]  $\leftarrow$  crossover(selectProg(),selectProg())
        else
            | newPopulation[j]  $\leftarrow$  reproduce(selectProg())
        end
    end
    foreach programme  $\in$  newPopulation do
        if randInt(1,10) = 1 then
            | mutate(programme)
        end
    end
    population  $\leftarrow$  newPopulation
    evaluatePopulation(generation)
end
Output: best(population)
```

---

The three genetic operators use the following algorithms:

---

**Algorithm 3:** Reproduction Operator

---

**Input:** parentProg  
**Output:** parentProg

---

---

**Algorithm 4:** Crossover Operator

---

**Input:** parentProg1  
**Input:** parentProg2  
crossoverPoint1  $\leftarrow$  randInt(1,parent1.proglen) crossoverPoint2  $\leftarrow$   
randInt(1,parent2.proglen)  
childProg1  $\leftarrow$   
replaceSubtree(parentProg1,crossoverPoint1,parentProg2.prog[crossoverPoint2])  
childProg2  $\leftarrow$   
replaceSubtree(parentProg2,crossoverPoint2,parentProg1.prog[crossoverPoint1])  
**Output:** childProg1  
**Output:** childProg2

---

---

**Algorithm 5:** Mutation Operator

---

**Input:** parentProg  
randomNode  $\leftarrow$  randInt(1,parentProg.proglen)  
childProg  $\leftarrow$   
replaceSubtree(parentProg,randomNode,randomTree())  
**Output:** childProg

---

### 3.2 SNGP Implementation

The SNGP implementation was solely based on Dr Jackson example code. Much of it was re-used from my GP implementation, including the primitives and arity table.

Each node in the SNGP population stores its current fitness, fitness from th previous generation, an array of operands, an array of predecessors, and the number of nodes in the programme tree. This is stored in the data structure shown in figure 3.3. An arity table is also stored in the same was as in GP.

Node	
Primitive	primitive
double	fitness
double	oldFitness
int[]	operands
iny[]	predecessors
int	proglen

Figure 3.3: Data structure to store information about each node in the SNGP population

The array of predecessors is a fixed length, the same as the population array. It stores the indicies of each predecessor in a linked list, where *predecessor*[0] contains the value of the predecessor with the smallest index,  $p_1$ , and then *predecessor*[ $p_1$ ] stores the value of the next largest predecessor,  $p_2$ , and so on until *predecessor*[ $p_n$ ] = 0. This method allows for efficient traversal, insertion, and deletion, while using a fixed amount of memory.

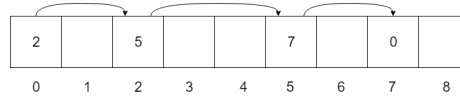


Figure 3.4: A predecessor array containing the indices 2, 5, and 7

Figure 3.5: fig:predecessor

This structure is used again to determine which nodes should be updated after performing a successor mutate operation.

---

**Algorithm 6:** Create list of nodes to be updated

---

```
Function buildUpdateList (mutatedNodeIndex):  
  addToUpdateList(mutatedNodeIndex)  
  predecessors  $\leftarrow$  population[mutatedNodeIndex].predecessors  
  nextPredecessorIndex  $\leftarrow$  predecessors[0]  
  while  $\neg$ (nextPredecessorIndex = 0) do  
    inList  $\leftarrow$  addToUpdateList(nextPredecessorIndex)  
    if inList then  
      buildUpdateList(nextPredecessorIndex)  
    end  
    nextPredecessorIndex  $\leftarrow$  predecessors[nextPredecessorIndex]  
  end  
end  
Function addToUpdateList (nodeIndex):  
  nextValue  $\leftarrow$  0  
  while  $\neg$ (updateList[nextValue] = 0)  $\wedge$  updateList[nextValue] <  
    nodeIndex do  
    nextValue  $\leftarrow$  updateList[nextValue]  
  end  
  if  $\neg$ (updateList[nextValue] = nodeIndex) then  
    updateList[nodeIndex]  $\leftarrow$  updateList[nextValue]  
    updateList[nextValue]  $\leftarrow$  nodeIndex  
    Output: False  
  else  
    Output: True  
  end  
end
```

---

The pseudocode for the overall SNGP algorithm is:

---

**Algorithm 7:** SNGP Algorithm

---

```
initialisePopulation() fitness  $\leftarrow$  evaluatePopulation()  
oldFitness  $\leftarrow$  infinity  
for generation  $\leftarrow$  1 to NUM_GENERATIONS do  
  randomNodeIndex  $\leftarrow$  randInt(1, POPULATION_SIZE)  
  successorMutate(randomNodeIndex)  
  buildUpdatelist(randomNodeIndex)  
  oldFitness  $\leftarrow$  fitness  
  fitness  $\leftarrow$  evaluatePopulation()  
  if fitness  $\leq$  oldFitness then  
    undoSuccessorMutate() fitness  $\leftarrow$  oldFitness  
  end  
end  
Output: bestProgramme()
```

---

The two methods of calculating the populations overall fitness are SNGP/A and SNGP/B, which are described in the following algorithms:

---

**Algorithm 8:** Calculate the SNGP/A fitness

---

**Input:** population  
 $totalFitness \leftarrow 0$   
**foreach**  $programme \in population$  **do**  
     $totalFitness \leftarrow totalFitness + evaluate(programme)$   
**end**  
**Output:**  $totalFitness / size(population)$

---



---

**Algorithm 9:** Calculate the SNGP/B fitness

---

**Input:** population  
 $bestFitness \leftarrow 0$   
**foreach**  $programme \in population$  **do**  
     $totalFitness \leftarrow \max(bestFitness, evaluate(programme))$   
**end**  
**Output:**  $bestFitness$

---

### 3.3 Fitness Function

While many different fitness functions were used, they all counted the number of inversions in a test array before and after executing a programme. To efficiently count inversions I used a modified merge sort algorithms, detailed in algorithm 10. The algorithm is  $O(n \log n)$ , and merging is done using two preallocated buffers to avoid unnecessary memory allocations while the GP algorithm is running.

---

**Algorithm 10:** Algorithm that sorts array and counts number of inversions

---

```

Function countInversions (arr):
  if length(arr) ≤ 1 then
    | inversions ← 0
  else
    sizeA ← ⌊length(arr)/2⌋
    sizeB ← length(arr) - sizeA
    A ← [First sizeA elements of arr]
    B ← [Last sizeB elements of arr]
    inversions ← countInversions(A) + countInversions(B)
    C ← Empty Array
    while length(A) > 0 and length(B) > 0 do
      if B[1] < A[1] then
        | Append B[1] to C
        | inversions ← inversions + length(A)
        | Remove first element of B
      else
        | Append A[1] to C
        | Remove first element of A
      end
    end
    if length(A) = 0 then
      | Append remaining elements of B to C
    else if length(B) = 0 then
      | Append remaining elements of A to C
    end
    arr ← C
  end
  Output: inversions
end

```

---

The fitness function I used as a starting point for both my GP and SNGP implementations was the same Kinnear used in [3, 4]. The function is defined as the following:

$$\begin{aligned}
fitness(prog) &= \frac{adjusted(prog)}{\sum_{p \in population} adjusted(p)} \\
adjusted(prog) &= \frac{1}{1 + raw(prog)} \\
raw(prog) &= praw(prog) - \min_{p \in population} praw(p) \\
praw(prog) &= \left( \sum_{t=1}^{NumberOfTests} res(t) \right) \cdot of + size(prog) \cdot sf \\
res(t) &= rdis(t) + pdis(t) \\
pdis(t) &= \max(rdis(t) - idis(t), 0) \cdot 100
\end{aligned}$$

Figure 3.6: Kinnear's Fitness Function

Where  $rdis(t)$  is the remaining number of inversions in a test array  $t$  after running a member of the population,  $idis(t)$  is the initial number of inversions in a test array  $t$ , and the order factor  $of$  and size factor  $sf$  are constant weights used to adjust the resulting fitness.

For SNGP, I also used other fitness function I created myself. The vast majority of these were not effective at all, and so will not be included. The ones that seemed to work best were some variation on the following:

$$fitness(x) = \begin{cases} 1 & \text{if } rdis(t) = 0 \\ -rdis(t) & \text{if } rdis(t) > idis(t) \wedge rdis(t) \neq 0 \\ 0 & \text{if } rdis(t) = idis(t) \wedge rdis(t) \neq 0 \\ 1 - rdis(t)/idis(t) & \text{otherwise} \end{cases}$$

Figure 3.7: Fitness Function 1

### 3.4 Test Data

Each programme generated during a GP run or an SNGP run is executed on a set of test arrays. This data comes from pre-generated file containing a series of random arrays loaded into the programme at the beginning of execution. This has the advantage of allowing for reproduction of results, easier testing of the programme, and prevents random number generation from increasing the execution time.

The arrays are organised into fixed size groupings. Each programme in the same GP generation is executed on the same group of arrays. This was done so that each programme was subjected to the same conditions in each generation, while still preventing over fitting to a fixed set of data. This method also allows



each group to be customized to meet certain requirements (e.g. contains a sorted array, a reversed array, or a very short array), but this was found to be unnecessary.

The file is a standard text file where the first two values on each line are the length and number of inversions of the array, followed by the array itself. Each value is separated by a space. A blank line indicates the end of the current group of arrays.

When loaded into the programme, the test arrays are stored as the following struct:

Array	
int	size
int	inversions
int[]	arr

Figure 3.8: Data structure to store each test array

The files were generated by a python script, and the inversions counted using algorithm 10. This allowed me to easily change the parameters of the test files and experiment with what worked.

## 4 Results

I was successful in replicating Kinnear’s work in evolving a sorting algorithm using standard GP, and found his methods able to evolve a sort very easily with no modification. Figure 4.1 shows his method very easily able to evolve a sort within a couple of generations, and quickly plateaued to around 500 solutions per generation.

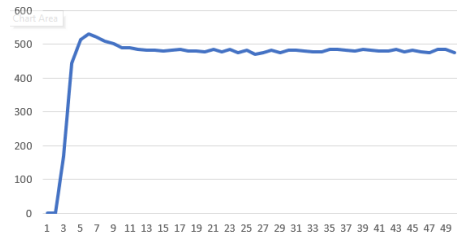


Figure 4.1: .Graph of number of solutions generated by GP (vertical) against generation (horizontal). Population size 1000, 50 generations. Size factor = 5, order factor = 5

My SNGP implementation was far less successful, however. I was not able to create an implementation that could reliably generate a working solution. I do not believe this to be an error in my code as by seeding the initial population with a handmade working solution my implementation was able to identify it

as a working sorting algorithm and re-use some of the structures within it to great success. I also attempted to seed the population with a partial solution containing useful structures, and the implementation was able to use it to create a working solution in approximately 6.7% of runs, using a population of 1000, 50 generations, and my own fitness function mentioned in figure 3.7.

I found little difference in results no matter the parameters used for SNGP. Using Kinnear's fitness function always resulted in plateauing between a fitness of 0.005 and 0.01 after about 15 generations. My own fitness functions faired no better, reaching its maximum of 0.4 by 20 generations.

All tests were over 20 runs. The custom fitness function refers to the function in figure 3.7. Parameters tested include:

- SNGP/A or SNGP/B
- Population size
- Number of generations
- Fitness function
- Size factor and order factor in Kinnear's function
- Strict or non strict successor mutate. Strict successor mutate only keeps changes that result in strictly greater fitness and reverts the rest, while non strict keeps changes that have greater than or the same fitness.

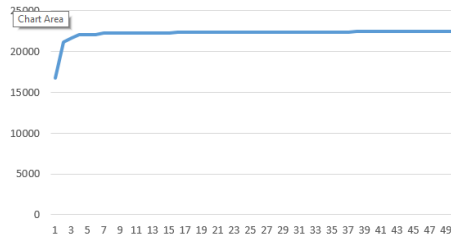


Figure 4.2: Graph of fitness against generation. SNGP/A, Population size 1000, 50 generations. Size factor = 5, order factor = 5. Kinnear's fitness function.

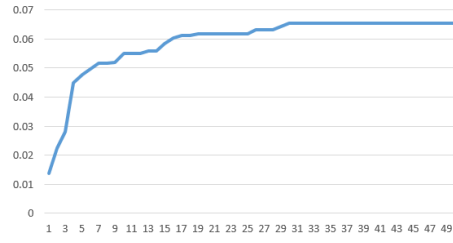


Figure 4.3: Graph of fitness against generation. SNGP/A, Population size 1000, 50 generations. Size factor = 5, order factor = 5. Kinnear's fitness function.

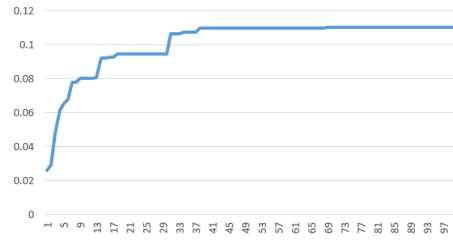


Figure 4.4: Graph of fitness against generation. SNGP/B, Population size 500, 100 generations. Size factor = 5, order factor = 5. Kinnear's fitness function.

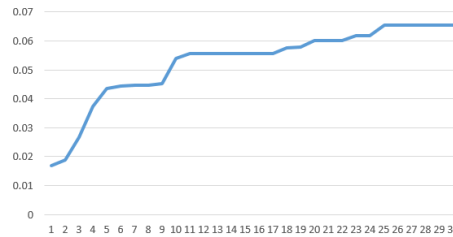


Figure 4.5: Graph of fitness against generation. SNGP/B, Population size 1500, 30 generations. Size factor = 5, order factor = 5. Kinnear's fitness function.

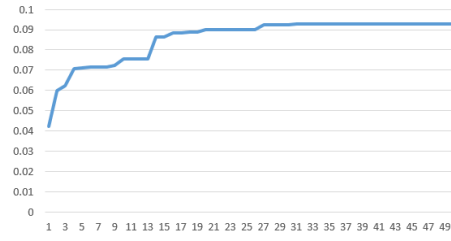


Figure 4.6: Graph of fitness against generation. SNGP/B, Population size 1000, 50 generations. Size factor = 0, order factor = 5. Kinnear's fitness function.

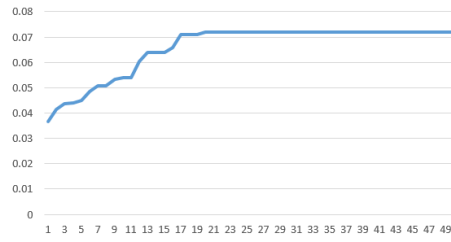


Figure 4.7: Graph of fitness against generation. SNGP/B, Population size 1000, 50 generations. Size factor = 2, order factor = 5. Kinnear's fitness function.

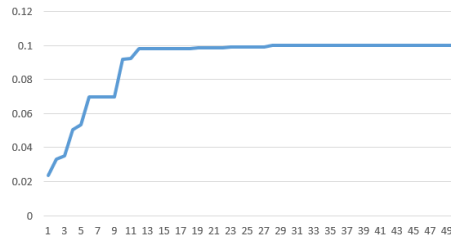


Figure 4.8: Graph of fitness against generation. SNGP/B, Population size 1000, 50 generations. Size factor = 8, order factor = 5. Kinnear's fitness function.

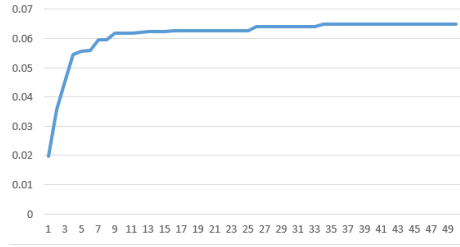


Figure 4.9: Graph of fitness against generation. SNGP/B strict, Population size 1000, 50 generations. Size factor = 5, order factor = 5. Kinnear's fitness function.

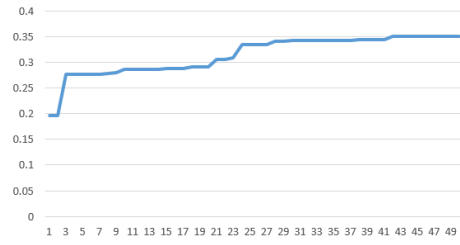


Figure 4.10: Graph of fitness against generation. SNGP/B, Population size 1000, 50 generations. Custom fitness function.

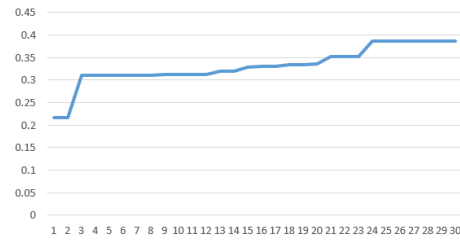


Figure 4.11: Graph of fitness against generation. SNGP/B, Population size 1500, 30 generations. Custom fitness function.

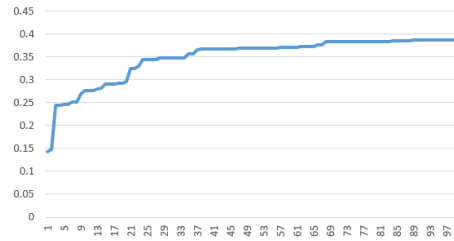


Figure 4.12: Graph of fitness against generation. SNGP/B, Population size 500, 100 generations. Custom fitness function.

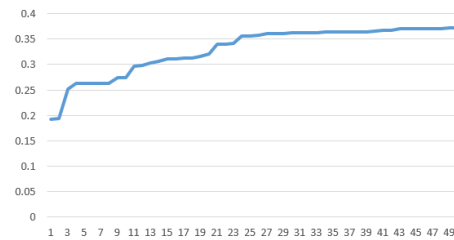


Figure 4.13: Graph of fitness against generation. SNGP/B strict, Population size 500, 100 generations. Custom fitness function.

In conclusion, I believe that SNGP is not effective at evolving a sorting algorithm.

## 5 Evaluation and Learning Points

While I was not successful in evolving a sort using SNGP, I achieved my goal of assessing the performance of SNGP in evolving a sort. There is room for further analysis of why SNGP is not effective at this task, such as by experimenting with other parameters or finding other problems that aren't suitable to be solved by SNGP and comparing them.

Throughout the course of this project I have learnt how to specify a substantial problem and how to put together a plan to complete the problem. I have improved my skills at locating relevant information and making use of it within the project, and I can use it to help evaluate and contextualise my own work.

I now have a much wider knowledge of current research topics within the field, and a much deeper understanding of machine learning, the areas where it can be applied, and the strengths and weaknesses of the techniques used in this domain.

This project was also another chance for me to improve my ability to produce a presentation and use it to convey information in a clear and concise manner. Similarly, I have also been able to practice my writing skills when making my

specification and design documents and this report. My C programming skills have certainly improved, and I now feel much more confident in my ability to produce a working piece of software in C.

## 6 Professional Issues and Data Required

This project follows the four sections of the British Computer Society Code of Conduct:

- **Public Interest:** This section concerns third parties and the general public. This project does not contain any public components, or interact with any third parties or members of the public. Any documents or code used in the course of this project were freely available to me through the University of Liverpool and cited properly in the bibliography. As such this project does not affect the public health, privacy, or security of others, and it does not breach the right of third parties, discriminate against anybody in any way.
- **Professional Competence and Integrity:** This section is about not claiming a level of competence that is not possessed. This project was specifically created by Dr Jackson for someone with the level of competence a third year computer science student should have. As such I believe that this project was well within my level of confidence when I requested to have this project assigned to me, and as I have successfully completed the project I still believe this to be the case.
- **Duty to Relevant Authority:** This section focuses on the responsibilities I am held to by a relevant authority. In this case the relevant authority is the University of Liverpool. I have not disclosed any confidential information, misrepresented any of my work, or caused a conflict of interest. I have carried out my responsibilities and accept professional responsibility for my work.
- **Duty to the Profession:** This section relates to the reputation of the computer science profession. By following the previous section I believe I have avoided any action that could bring the profession into disrepute. I have also not committed any crime or become bankrupt for the duration of this project, so I have no need to notify the BCS.

In regard to the data required for this project, I believe there is no need for proof of consent or ethical source of data as there was no human participation or human data used in this project.

## 7 Bibliography

### References

- [1] JACKSON, D. A new, node-focused model for genetic programming. In *Genetic Programming* (2012), A. Moraglio, S. Silva, K. Krawiec, P. Machado, and C. Cotta, Eds., Lecture Notes in Computer Science, Springer Berlin Heidelberg, pp. 49–60.
- [2] JACKSON, D. Single node genetic programming on problems with side effects. In *Parallel Problem Solving from Nature - PPSN XII* (2012), C. A. C. Coello, V. Cutello, K. Deb, S. Forrest, G. Nicosia, and M. Pavone, Eds., Lecture Notes in Computer Science, Springer Berlin Heidelberg, pp. 327–336.
- [3] KINNEAR, K. E. Evolving a sort: Lessons in genetic programming. In *in Proceedings of the 1993 International Conference on Neural Networks* (1993), IEEE Press, pp. 881–888.
- [4] KINNEAR, K. E. Generality and difficulty in genetic programming: Evolving a sort. In *ICGA* (1993).
- [5] POLI, R., LANGDON, W. B., AND MCPHEE, N. F. *A Field Guide to Genetic Programming*. Lulu Enterprises, UK Ltd.

## 8 Appendices

Code listing for GP\_Sort.c

```
#include <stdio.h>
#include <stdlib.h>
#include <time.h>

#define RANDOMSEED time(NULL)//1894

#define NUM_TERMINALS 2
#define NUM_FUNCTIONS 7
#define NUM_PRIMITIVES NUM_TERMINALS+NUM_FUNCTIONS
#define MAX_ARITY 3

#define POPULATION_SIZE 1000 //1000 in final
#define MAX_PROG_SIZE 3000
#define INITIAL_MAX_DEPTH 6
#define NUM_GENERATIONS 50 //50 in final
#define NUM_TESTS 15
#define MAX_RUNS 20 //20 in final
```



```

#define NUM_TEST_SETS 3000

#define SF 5
#define OF 5

#define MAX_MUTATE_DEPTH_INCREASE 1.15

#define OUT_FILE "results\\gp-success-trace-2.csv"
FILE* successTrace = NULL;

int numWorkingProgs = 0;
int workingProgramme = -1;

typedef enum {
DUMMY,
INDEX,
LENGTH,
ITERATE,
SWAP,
SMALLEST,
LARGEST,
SUB,
INC,
DEC
} Primitive;

typedef struct {
Primitive primitive;
int arity;
char* name;
} TableEntry;

TableEntry primitiveTable[NUM_PRIMITIVES+1] = {
{DUMMY, -1, "DUMMY"},
{INDEX, 0, "INDEX"},
{LENGTH, 0, "LENGTH"},
{ITERATE, 3, "ITERATE"},
{SWAP, 2, "SWAP"},
{SMALLEST, 2, "SMALLEST"},
{LARGEST, 2, "LARGEST"},
{SUB, 2, "SUB"},
{INC, 1, "INC"},
{DEC, 1, "DEC"}
};

typedef struct {

```

```

char code[MAX_PROG_SIZE+1]; // Includes null terminator
double fitness;
int progLen;
} Prog;

char* progNode = NULL;

Prog popBuffer1[POPULATION_SIZE];
Prog popBuffer2[POPULATION_SIZE];

Prog* population = &popBuffer1[0];
Prog* newPopulation = &popBuffer2[0];

typedef struct{
int size;
int inversions;
int arr[];
} Array;

#define arrayMem(x) (sizeof(Array) + sizeof(int) * (x))

Array* tests[NUM_GENERATIONS][NUM_TESTS];
int maxTestSize;

Array* results = NULL;
int* mergeBuffer1 = NULL;
int* mergeBuffer2 = NULL;

int index = 0;
int progIterations = 0;
#define MAX_PROG_ITERATIONS 10000

//Returns random integer in interval [min,max] inclusive
int randRange(int min, int max){

return min + rand()/(RAND_MAX/(max-min+1)+1);

}

int initialiseTestData(char* path){

FILE* file = fopen(path,"r");

if(!file) return 1;
int setNum = 0;
int testNum = 0;

```

```

maxTestSize = 0;

//Iterate through test data file until we have a test set for each g
while(setNum < NUMGENERATIONS){

char firstC = getc(file);

//If this line is blank, then we have reached the end of this test s
if(firstC == '\n'){
++setNum;
testNum = 0;
continue;

//if we have exeeded the number of tests per set, skip the remainder
}else if(testNum >= NUM_TESTS){
++setNum;
testNum = 0;

do{
fscanf(file,"%*[^\\n]",NULL);
getc(file);
firstC = getc(file);
}while(firstC != '\\n');
continue;
}
//If end of file has been reached, exit loop
else if(firstC == EOF){

break;

}

ungetc(firstC, file);

int arrSize = 0;
int inversions = 0;

fscanf(file, "%d %d", &arrSize, &inversions);

tests[setNum][testNum] = malloc( arrayMem(arrSize) );

if(arrSize > maxTestSize) maxTestSize = arrSize;

Array* test = tests[setNum][testNum];

```

```

test->size = arrSize;
test->inversions = inversions;

for(int i = 0; i < arrSize; i++){

int n;

fscanf( file , "%d", &n);

test->arr[i] = n;
}
getc( file );
++testNum;

}
fclose( file );
return 0;
}

//progNode must be set to Prog.code before this is called
//progIterations must be set to 0 before this is called
int execute(){

switch(*progNode++){

case INDEX:{

return index;

}break;

case LENGTH:{

return results->size;

}break;

case ITERATE:{

int len = results->size;
int start = execute();
int end = execute();
char* function = progNode;

int oldIndex = index;

```

```

for(index = start; index <= end && index < len && progIterations < M
progNode = function;

execute();
}

index = oldIndex;

return (end < len) ? end : len ;

}break;

case SWAP:{

int x = execute();
int y = execute();

//If x or y is not a valid index, return 0
if(x<0 || y<0 || x>=results->size || y>=results->size) return 0;

int t = results->arr[x];
results->arr[x] = results->arr[y];
results->arr[y] = t;

return x;

}break;

case SMALLEST:{

int x = execute();
int y = execute();

if(results->arr[x] < results->arr[y]){
return x;
}
else{

return y;

}

}break;

case LARGEST:{

```

```

int x = execute();
int y = execute();

if(results->arr[x] > results->arr[y]){
return x;
}
else{

return y;

}

}break;

case SUB:{

int x = execute();
int y = execute();

return x-y;

}break;

case INC:{

int x = execute();

return x+1;

}break;

case DEC:{

int x = execute();

return x-1;

}break;

default:
printf("\nINVALID PRIMITIVE (%d)\n", *(progNode-1));
return -1;
}

}

```

```

int countInversionsRec(int* arr, int* working, int size, int offset)

if(size <= 1){

return 0;

}

int sizeA = size / 2;
int sizeB = size - sizeA;

int* A = working + offset;
int* B = &working[sizeA] + offset;

int inversions = countInversionsRec(working, arr, sizeA, offset) + countInversionsRec(working, &working[sizeA], sizeB, offset);

int aCounter = 0;
int bCounter = 0;
int cCounter = offset;

while( aCounter < sizeA && bCounter < sizeB ){

if( A[aCounter] <= B[bCounter] ){

arr[cCounter] = A[aCounter];
aCounter++;

} else {

arr[cCounter] = B[bCounter];
inversions += (sizeA - aCounter);
bCounter++;

}

cCounter++;

}

if(aCounter < sizeA){

memcpy( &arr[cCounter], &A[aCounter], (sizeA - aCounter) * sizeof(int) );
aCounter = sizeA;

} else if(bCounter < sizeB){

```

```

memcpy( &arr[cCounter], &B[bCounter], (sizeB - bCounter) * sizeof(int)
}

return inversions;

}

int countInversions(Array* arr){
    if(mergeBuffer1 == NULL || mergeBuffer2 == NULL) return -1;

    memcpy(mergeBuffer1, arr->arr, arr->size*sizeof(int));
    memcpy(mergeBuffer2, arr->arr, arr->size*sizeof(int));

    arr->inversions = countInversionsRec(mergeBuffer1, mergeBuffer2, arr);

    return arr->inversions;
}

//returns the result of each test by calculating the remaining disorder
//the test and adding a penalty disorder pDis.
int res(int popIndex, int testSet, int testNum){

    Array* test = tests[testSet][testNum];

    //If inversions not counted, count inversions
    if(test->inversions == -1){

        memcpy(results, test->arr, test->size);
        countInversions(test);

    }

    index = 0;

    memcpy(results, test, arrayMem(test->size));

    progIterations = 0;
    progNode = population[popIndex].code;
    execute(popIndex);

    int iDis = test->inversions;
    int rDis = countInversions(results);

```



```

int pDis = (rDis > iDis) ? (rDis - iDis)*100 : 0;

return rDis + pDis;

}

//Calculates the provisional raw fitness of a programme
int praw(int testSet, int popIndex){

int resSum = 0;

for(int testNum = 0; testNum < NUMTESTS; testNum++){
resSum += res(popIndex, testSet, testNum);
}

if(resSum == 0){
numWorkingProgs += 1;
workingProgramme = popIndex;
}

return (resSum * OF) + (population[popIndex].progLen * SF);

}

//Calculates the fitness for each member of the population
void evaluatePopulation(int testSet){

numWorkingProgs = 0;

int prawTable[POPULATION_SIZE];
int raw[POPULATION_SIZE];
float adj[POPULATION_SIZE];
float adjSum = 0;

prawTable[0] = praw(testSet, 0);

int minpraw = prawTable[0];

//Calculate the provisional raw fitness and minimum raw fitness
for(int popIndex = 1; popIndex < POPULATION_SIZE; ++popIndex){

prawTable[popIndex] = praw(testSet, popIndex);

if(prawTable[popIndex] < minpraw) minpraw = prawTable[popIndex];
}

```

```

}

//Calculate the adjusted fitness for each member of the population
for(int popIndex = 0; popIndex < POPULATION_SIZE; ++popIndex){

int raw = prawTable[popIndex] - minpraw;

adj[popIndex] = 1.0/(1.0+raw);

adjSum += adj[popIndex];

}

//Set the fitness for each member of the population as the normalised
//fitness
for(int popIndex = 0; popIndex < POPULATION_SIZE; ++popIndex){

population[popIndex].fitness = adj[popIndex]/adjSum;

}

fprintf(successTrace,"%d",numWorkingProgs);

}

char* createTree(char* tree, int depth, int maxNodes, int full){
char* treeEnd = tree;

if (depth == 1 || maxNodes == 1){
*treeEnd = randRange(1,NUM_TERMINALS);
treeEnd++;
}

else{

char primitive;
int arity;
do{

primitive = full ? randRange(NUM_TERMINALS+1, NUM_PRIMITIVES): randR
arity = primitiveTable[primitive].arity;

}while(arity > maxNodes-1);

*tree = primitive;

```

```

treeEnd++;

for(int argument = 0; argument < arity; ++argument){
    treeEnd = createTree(treeEnd, depth-1, maxNodes - (treeEnd - tree) -
    }
    }
    return(treeEnd);
}

int fitnessProportionalSelection(){
    float randNum = (float)rand() / (float)RAND_MAX;
    for(int popIndex = 0; popIndex < POPULATION_SIZE; ++popIndex ){
        if(population[popIndex%POPULATION_SIZE].fitness >= randNum) return popIndex;
        else randNum -= population[popIndex].fitness;
    }

    //If this line of code is reached it is due to errors in floating point arithmetic
    return POPULATION_SIZE-1;
}

int subtreeLength(char* start){
    int remaining = primitiveTable[*start].arity;
    int length = 1;

    while(remaining > 0){
        remaining--;
        start++;
        remaining += primitiveTable[*start].arity;
        length++;
    }

    return length;
}

```

```

int subTreeDepth(char* start){

int maxDepth = 0;

int arity = primitiveTable[*start].arity;

start++;

for(int i = 0; i<arity; i++){

int depth = subTreeDepth(start);

if(depth > maxDepth) maxDepth = depth;

start += subtreeLength(start);

}

return maxDepth+1;

}

//Creates a new programme at newProg by copying baseProg and replace
//with a copy of newSubtree
void replaceSubtree(char* baseProg, int baseLen, int oldSubtree, char* newSubtree)

memcpy(newProg, baseProg, oldSubtree);
int subtreelen = subtreeLength(baseProg + oldSubtree);
int newSubtreeLen = subtreeLength(newSubtree);
memcpy(newProg + oldSubtree, newSubtree, newSubtreeLen);
memcpy(newProg + oldSubtree + newSubtreeLen, baseProg + oldSubtree +
newProg[baseLen-subtreelen+newSubtreeLen] = '\0';

}

void reproduction(int popIndex){

newPopulation[popIndex] = population[fitnessProportionalSelection()]

}

void mutate(Prog* baseProg, int popIndex){

//int selectedProg = fitnessProportionalSelection();

int baseDepth = subTreeDepth(baseProg->code);

```

```

int maxDepth = baseDepth * MAX_MUTATE_DEPTH_INCREASE;

int mutatedNode = randRange(0, baseProg->progLen - 1);

int subTreeDep = subTreeDepth(&baseProg->code[mutatedNode]);
int subTreeLen = subtreeLength(&baseProg->code[mutatedNode]);

char newTree[MAX_PROG_SIZE];
createTree(newTree, maxDepth - baseDepth + subTreeDep, MAX_PROG_SIZE);

replaceSubtree(baseProg->code, baseProg->progLen, mutatedNode, newTree);

newPopulation[popIndex].progLen = strlen(newPopulation[popIndex].code);
}

void crossover(int popIndex1, int popIndex2){

int parIndex1 = fitnessProportionalSelection();
int parIndex2 = fitnessProportionalSelection();
while(parIndex1 == parIndex2) parIndex2 = fitnessProportionalSelection();

Prog* parent1 = &population[parIndex1];
Prog* parent2 = &population[parIndex2];

int crossoverPoint1 = randRange(0, parent1->progLen - 1);
int crossoverPoint2 = randRange(0, parent2->progLen - 1);

replaceSubtree(parent1->code, parent1->progLen, crossoverPoint1, parent2->code);
newPopulation[popIndex1].progLen = strlen(newPopulation[popIndex1].code);

if(popIndex2 < POPULATION_SIZE){

replaceSubtree(parent2->code, parent2->progLen, crossoverPoint2, parent1->code);
newPopulation[popIndex2].progLen = strlen(newPopulation[popIndex2].code);
}

}

void initialisePopulation(){

for(int popIndex = 0; popIndex < POPULATION_SIZE; ++popIndex){

Prog* prog = &population[popIndex];

char* code = prog->code;

```

```

createTree(code , 6, MAX_PROG_SIZE,0);

prog->progLen = strlen(code);

}

}

int init(char* path){
if(path == NULL) return 1;
if(initialiseTestData(path)) return 1;
srand(RANDOMSEED);
results = malloc( arrayMem(maxTestSize) );
mergeBuffer1 = malloc( sizeof(int) * maxTestSize );
mergeBuffer2 = malloc( sizeof(int) * maxTestSize );
return 0;

}

#include "GP_Sort_Debug.c"
int main(int argc , char* argv[]){

printf("Start\n\n");
printf(argv[1]);
printf("\n\n");

printPrimitiveTable();

if(init(argv[1])){
printf("File Not Found");
return 1;
}else{

printf("\nTest file loaded\n\n");

}

successTrace = fopen(OUT_FILE,"w");

for(int run = 0; run<MAXRUNS; ++run){

//if(numWorkingProgs>0) break;

printf("\n\nRun %d\n\n",run);

```

```

initialisePopulation();

//evaluate the initial population (generation 0)
evaluatePopulation(0);

for(int generation = 1; generation < NUMGENERATIONS; ++generation){
// if (numWorkingProgs>0) break;

for(int popIndex = 0; popIndex < POPULATION_SIZE; popIndex++){

if(randRange(0,9) <= 8){
crossover(popIndex,(popIndex+1)%POPULATION_SIZE);
++popIndex;
}
else reproduction(popIndex);

}

for(int popIndex = 0; popIndex < POPULATION_SIZE; popIndex++){
if(randRange(0,9) == 0){
Prog p = newPopulation[popIndex];
mutate(&p,popIndex);
}

Prog* temp = newPopulation;
newPopulation = population;
population = temp;

evaluatePopulation(generation);

if(generation%100 == 0){

printf("\n%d", generation);

}

}

fprintf(successTrace,"\n");
}

```

```

//printPopulation();

if (numWorkingProgs>0){
printf("\n\nSuccess!\n\n");
printNode(workingProgramme);
}

fclose(successTrace);

return 0;

}

```

Code listing for GP\_Sort\_Debug.c

```

#ifndef GP_SORT_DEBUG_C
#define GP_SORT_DEBUG_C

#include<stdio.h>

void printIntArray(int* a, int n){
for(int i = 0; i<n; i++){
printf("%d ", a[i]);
}
}

void printPrimitiveTable(){
for(int i = 0; i < NUM_PRIMITIVES; i++){
TableEntry entry = primitiveTable[i];
printf("i: %d, e: %d, a: %d n: %s\n", i, entry.primitive, entry.arity);
}

printf("\n");
}

void printTestData(){
for(int i = 0; i < NUM_GENERATIONS; i++){

```



```

printf("Set %d\n\n", i);

for(int j = 0; j < NUM_TESTS; j++){
    Array* test = tests[i][j];

    printf("Test: %d-%d Size: %d Inversions: %d\n", i, j, test->size, te

    for(int k = 0; k < test->size; k++){
        printf("%d ", test->arr[k]);
    }

    printf("\n");
}

printf("\n");

printf("\n");
}

printf("\n");
}

void printNode(int popIndex){

    Prog prog = population[popIndex];
    printf(
        "Index: %d ProgLen: %d Fitness: %f \nProg: ",
        popIndex,
        prog.progLen,
        prog.fitness
    );

    for(int i = 0; i < prog.progLen; ++i){

        printf("%s ", primitiveTable[prog.code[i]].name);
    }
}

```

```

void printPopulation(){

for(int popIndex = 0; popIndex < POPULATION_SIZE; popIndex++){
printNode(popIndex);

printf("\n\n");
}

printf("\n");
}

void setExampleProgramme(Prog* prog){

Primitive code[18] = {ITERATE,SUB,LENGTH,LENGTH,DEC, LENGTH ,ITERATE

for(int i = 0; i<18; i++){

prog->code[i] = code[i];

}

prog->progLen = 18;

}

void testExecution(){

setExampleProgramme(&population[0]);

res(0,0,0);

printIntArray(results->arr , results->size);

}

#endif

```

Code listing for SNGPSort.c

```

#include <stdio.h>
#include <stdlib.h>
#include <limits.h>
#include <time.h>

```

```

#define RANDOMSEED time(NULL)//2928

#define NUMTERMINALS 2
#define NUMFUNCTIONS 7
#define NUM_PRIMITIVES NUMTERMINALS+NUMFUNCTIONS
#define MAXARITY 3

//Macro to define which version of SNGP is used
#define evaluatePopulation(updateList,testSet) evaluatePopulationSNGP

//The maximum number of times we apply the successor mutate operation
#define MAX_OPS 50
#define NUMGENERATIONS MAX_OPS+1
#define POPULATION_SIZE 1000
#define NUM_TESTS 15
#define MAX_RUNS 20
#define NUM_TEST_SETS 30000
#define BETTER_THAN >=

#define OUTPUT_INTERVAL 5

#define OUT_FILE "results\\sngpA-fitness-trace-ff1.csv"
FILE* fitnessTrace = NULL;

#define SF 5
#define OF 5

typedef struct{
int size;
int inversions;
int arr[];
} Array;

#define arrayMem(x) (sizeof(Array) + sizeof(int) * (x))

Array* tests[NUMGENERATIONS][NUM_TESTS];
int maxTestSize;

Array* results = NULL;
int* mergeBuffer1 = NULL;
int* mergeBuffer2 = NULL;
int updateList[POPULATION_SIZE];

int index = 0;
int progIterations = 0;
#define MAX_PROG_ITERATIONS 2000

```

```

int success = 0;
int workingProgramme = -1;

//For SNGP/A
float totalNodeFitness = 0;

typedef enum {
INDEX,
LENGTH,
ITERATE,
SWAP,
SMALLEST,
LARGEST,
SUB,
INC,
DEC
} Primitive;

typedef struct {
Primitive primitive;
int arity;
char* name;
} TableEntry;

TableEntry primitiveTable[NUM_PRIMITIVES] = {
{INDEX, 0, "INDEX"},
{LENGTH, 0, "LENGTH"},
{ITERATE, 3, "ITERATE"},
{SWAP, 2, "SWAP"},
{SMALLEST, 2, "SMALLEST"},
{LARGEST, 2, "LARGEST"},
{SUB, 2, "SUB"},
{INC, 1, "INC"},
{DEC, 1, "DEC"}
};

typedef struct {
Primitive primitive;
double fitness;
double oldFitness;
int operands[MAX_ARITY];
int predecessors[POPULATION_SIZE];
int progLen;
} Node;

```

```

Node population[POPULATION_SIZE];

void addPredecessor(int popIndex, int newPredIndex){
    //If the node is a terminal, we don't need to maintain its predecess
    if(popIndex < NUM_TERMINALS) return;

    Node* node = &population[popIndex];
    int* predArray = node->predecessors;
    int nextValue = 0;
    while(predArray[nextValue] != 0 && predArray[nextValue] < newPredInd
    nextValue = predArray[nextValue];
}

if(predArray[nextValue] != newPredIndex){
    predArray[newPredIndex] = predArray[nextValue];
    predArray[nextValue] = newPredIndex;
}

}

int removePredecessor(int popIndex, int newPredIndex){
    //If the node is a terminal, we don't need to maintain its predecess
    if(popIndex < NUM_TERMINALS) return 0;

    Node* node = &population[popIndex];
    int* predArray = node->predecessors;
    int nextValue = 0;
    while(predArray[nextValue] != 0 && predArray[nextValue] != newPredIn

```

```

nextValue = predArray[nextValue];

}

//if predIndex not found in predArray, return error
if(predArray[nextValue] == 0) return 1;

predArray[nextValue] = predArray[newPredIndex];

predArray[newPredIndex] = 0;

return 0;
}

//Returns random number in interval [min,max] inclusive
int randRange(int min, int max){

return min + rand()/(RAND_MAX/(max-min+1)+1);

}

void initialisePopulation(){

for(int terminalIndex = 0; terminalIndex < NUM_TERMINALS; terminalIndex++)
{
Node* node = &population[terminalIndex];

node->primitive = terminalIndex;

node->fitness = -1;
node->oldFitness = -1;
node->progLen = 1;

}

for(int functionIndex = NUM_TERMINALS; functionIndex < POPULATION_SIZE; functionIndex++)
{
Node* node = &population[functionIndex];

Primitive primitive = randRange(NUM_TERMINALS, NUM_PRIMITIVES-1);

node->primitive = primitive;
node->fitness = -1;
node->oldFitness = -1;
node->progLen = 1;
}
}

```

```

node->predecessors[0] = 0;

for(int operandIndex = 0; operandIndex < MAX_ARITY; operandIndex++){

    if( operandIndex < primitiveTable[primitive].arity ){

        int randomOperand = randRange(0,functionIndex-1);

        node->operands[operandIndex] = randomOperand;
        node->progLen += population[randomOperand].progLen;

        addPredecessor(randomOperand, functionIndex);

    } else{

        node->operands[operandIndex] = -1;

    }

}

}

}

}

}

int initialiseTestData(char* path){

FILE* file = fopen(path,"r");

if(!file) return 1;
int setNum = 0;
int testNum = 0;
maxTestSize = 0;

//Iterate through test data file until we have a test set for each g
while(setNum < NUM_GENERATIONS){

char firstC = getc(file);

//If this line is blank, then we have reached the end of this test s
if(firstC == '\n'){
++setNum;
testNum = 0;
continue;

```

```

//if we have exeeded the number of tests per set, skip the remainder
}else if(testNum >= NUM_TESTS){
++setNum;
testNum = 0;

do{
fscanf(file,"%*[^\\n]",NULL);
getc(file);
firstC = getc(file);
}while(firstC != '\\n');
continue;
}
//If end of file has been reached, exit loop
else if(firstC == EOF){

break;

}

ungetc(firstC, file);

int arrSize = 0;
int inversions = 0;

fscanf(file, "%d %d", &arrSize, &inversions);

tests[setNum][testNum] = malloc( arrayMem(arrSize) );

if(arrSize > maxTestSize) maxTestSize = arrSize;

Array* test = tests[setNum][testNum];
test->size = arrSize;
test->inversions = inversions;

for(int i = 0; i < arrSize; i++){

int n;

fscanf(file, "%d", &n);

test->arr[i] = n;
}
getc(file);
++testNum;

```



```

}
fclose ( file );
return 0;
}

int execute (int popIndex){

Node* node = &population [popIndex];

switch (node->primitive){

case INDEX:{

return index;

}break;

case LENGTH:{

return results->size;

}break;

case ITERATE:{

int len = results->size;
int start = execute (node->operands [0]);
int end = execute (node->operands [1]);
int functionIndex = node->operands [2];

int oldIndex = index;

for (index = start; index <= end && index < len && progIterations < 2

execute (functionIndex);
}

index = oldIndex;

return (end < len) ? end : len ;

}break;

case SWAP:{

int x = execute (node->operands [0]);

```

```

int y = execute(node->operands[1]);

//If x or y is not a valid index, return 0
if(x<0 || y<0 || x>=results->size || y>=results->size) return 0;

int t = results->arr[x];
results->arr[x] = results->arr[y];
results->arr[y] = t;

return x;

}break;

case SMALLEST:{

int x = execute(node->operands[0]);
int y = execute(node->operands[1]);

if(results->arr[x] < results->arr[y]){
return x;
}
else{

return y;

}

}break;

case LARGEST:{

int x = execute(node->operands[0]);
int y = execute(node->operands[1]);

if(results->arr[x] > results->arr[y]){
return x;
}
else{

return y;

}

}break;

case SUB:{

```

```

int x = execute(node->operands[0]);
int y = execute(node->operands[1]);

return x-y;

}break;

case INC:{

int x = execute(node->operands[0]);

return x+1;

}break;

case DEC:{

int x = execute(node->operands[0]);

return x-1;

}break;

default:
printf("\nINVALID PRIMITIVE (%d)\n", node->primitive);
return -1;
}

}

int countInversionsRec(int* arr, int* working, int size, int offset)

if(size <= 1){

return 0;

}

int sizeA = size / 2;
int sizeB = size - sizeA;

int* A = working + offset;
int* B = &working[sizeA] + offset;

int inversions = countInversionsRec(working, arr, sizeA, offset) + countInversionsRec(working, arr, sizeB, offset);

```

```

int aCounter = 0;
int bCounter = 0;
int cCounter = offset;

while( aCounter < sizeA && bCounter < sizeB ){
    if( A[aCounter] <= B[bCounter] ){
        arr[cCounter] = A[aCounter];
        aCounter++;
    } else {
        arr[cCounter] = B[bCounter];
        inversions += (sizeA - aCounter);
        bCounter++;
    }
    cCounter++;
}

if(aCounter < sizeA){
    memcpy( &arr[cCounter], &A[aCounter], (sizeA - aCounter) * sizeof(int) );
} else if(bCounter < sizeB){
    memcpy( &arr[cCounter], &B[bCounter], (sizeB - bCounter) * sizeof(int) );
}

return inversions;
}

int countInversions(Array* arr){
    if(mergeBuffer1 == NULL || mergeBuffer2 == NULL) return -1;

    memcpy(mergeBuffer1, arr->arr, arr->size*sizeof(int));
    memcpy(mergeBuffer2, arr->arr, arr->size*sizeof(int));
}

```

```

arr->inversions = countInversionsRec(mergeBuffer1, mergeBuffer2, arr);

return arr->inversions;

}

float testNode(int popIndex, int testSet, int testNum){

Array* test = tests[testSet][testNum];

//If inversions not counted, count inversions
if(test->inversions == -1){

memcpy(results, test->arr, test->size);
countInversions(test);

}

index = 0;

memcpy(results, test, arrayMem(test->size));

progIterations = 0;
execute(popIndex);


int inversions = countInversions(results);

float fitness;// = test->inversions - inversions;

if(inversions == test->inversions && inversions!=0) fitness = 0;
else if(test->inversions!=0) fitness = 1 - inversions/(float)test->i
else if(inversions == 0) fitness = 1;
else fitness = -inversions;

return fitness;
}

float evaluateNode(int popIndex, int testSet){

float nodeTotalFitness = 0;

for(int testNum = 0; testNum < NUM_TESTS; testNum++){

nodeTotalFitness += testNode(popIndex, testSet, testNum);

```

```

}

population[popIndex].oldFitness = population[popIndex].fitness;
population[popIndex].fitness = nodeTotalFitness / NUM_TESTS;
if(population[popIndex].fitness > 0.8) success = 1;
return population[popIndex].fitness;
}

void updateProgLen(int* updateList){
    if (updateList == NULL){
        for(int nextUpdateNode = 0; nextUpdateNode < POPULATION_SIZE; nextUpdateNode++)
            Node* node = &population[nextUpdateNode];
        node->progLen = 1;
        for(int arg = 0; arg < primitiveTable[node->primitive].arity; arg++)
            node->progLen += population[node->operands[arg]].progLen;
    }

}

}
else{
    for(int nextUpdateNode = updateList[0]; nextUpdateNode != 0; nextUpdateNode++)
        Node* node = &population[nextUpdateNode];
    node->progLen = 1;
    for(int arg = 0; arg < primitiveTable[node->primitive].arity; arg++)
        node->progLen += population[node->operands[arg]].progLen;
}
}

```

```

    }
    }

}

void normaliseFitness(){

    int prawTable[POPULATION_SIZE];
    int raw[POPULATION_SIZE];
    float adj[POPULATION_SIZE];
    float adjSum = 0;

    prawTable[0] = population[0].fitness;

    int minpraw = prawTable[0];

    for(int popIndex = 1; popIndex < POPULATION_SIZE; ++popIndex){
        prawTable[popIndex] = population[popIndex].fitness;
        if(prawTable[popIndex] < minpraw) minpraw = prawTable[popIndex];
    }

    for(int popIndex = 0; popIndex < POPULATION_SIZE; ++popIndex){
        int raw = prawTable[popIndex] - minpraw;
        adj[popIndex] = 1.0/(1.0+raw);
        adjSum += adj[popIndex];
    }

    for(int popIndex = 0; popIndex < POPULATION_SIZE; ++popIndex){
        population[popIndex].fitness = adj[popIndex]/adjSum;
    }
}

float evaluatePopulationSNGP_A(int* updateList, int testSet){

```

```

if (updateList == NULL){
totalNodeFitness = 0;
for (int popIndex = 0; popIndex < POPULATION_SIZE; popIndex++){
totalNodeFitness += evaluateNode (popIndex, testSet);
}
} else{
for (int nextUpdateNode = updateList[0]; nextUpdateNode != 0; nextUpdateNode++){
totalNodeFitness += evaluateNode (nextUpdateNode, testSet);
totalNodeFitness -= population [nextUpdateNode].oldFitness;
}
}
return totalNodeFitness/POPULATION_SIZE;
}

float evaluatePopulationSNGP_B (int* updateList, int testSet){
float bestNodeFitness = 0;
if (updateList == NULL){
bestNodeFitness = evaluateNode (0, testSet);
for (int popIndex = 1; popIndex < POPULATION_SIZE; popIndex++){
float nodeFitness = evaluateNode (popIndex, testSet);
if (bestNodeFitness BETTER_THAN nodeFitness) bestNodeFitness = nodeFitness;
}
} else{
for (int nextUpdateNode = updateList[0]; nextUpdateNode != 0; nextUpdateNode++){
evaluateNode (nextUpdateNode, testSet);
}
}
}

```



```

bestNodeFitness = population[0].fitness;

for(int i = 1; i<POPULATION_SIZE; ++i){

if(population[i].fitness BETTER_THAN bestNodeFitness) bestNodeFitness

}
}

return bestNodeFitness;
}

void successorMutate(int popIndex, int randomOperandIndex, int newOperandValue)
{
Node* node = &population[popIndex];
removePredecessor(node->operands[randomOperandIndex], popIndex);
node->operands[randomOperandIndex] = newOperandValue;
addPredecessor(newOperandValue, popIndex);
}

int addToUpdateList(int popIndex){
int nextValue = 0;

while(updateList[nextValue] != 0 && updateList[nextValue] < popIndex)
nextValue = updateList[nextValue];
}

if(updateList[nextValue] != popIndex){
updateList[popIndex] = updateList[nextValue];
updateList[nextValue] = popIndex;

//popIndex is not already in update list
return 0;

} else{

//popIndex is already in update list
return 1;
}

```

```

}

}

//updateList[0] must be set to 0 before this function is called
void buildUpdateList(int popIndex){

addToUpdateList(popIndex);

int* predArray = population[popIndex].predecessors;

int nextPredecessor = predArray[0];

while(nextPredecessor != 0){

int inList = addToUpdateList(nextPredecessor);
if(!inList) buildUpdateList(nextPredecessor);
nextPredecessor = predArray[nextPredecessor];

}

}

void restoreFitnessValues(int* list){

int nextNode = list[0];

while(nextNode != 0 && nextNode < POPULATION_SIZE){

population[nextNode].fitness = population[nextNode].oldFitness;

nextNode = list[nextNode];

}

}

int init(char* path){
if(path == NULL) return 1;
if(initialiseTestData(path)) return 1;
srand(RANDOMSEED);
results = malloc( arrayMem(maxTestSize) );
mergeBuffer1 = malloc( sizeof(int) * maxTestSize );
mergeBuffer2 = malloc( sizeof(int) * maxTestSize );

```

```

return 0;

}

#include "SNGP_Sort_Debug.c"

int main(int argc, char* argv[]) {

printf("Start\n\n");
printf(argv[1]);
printf("\n\n");

printPrimitiveTable();

if (init(argv[1])) {
printf("File Not Found");
return 1;
} else {

printf("\nTest file loaded\n\n");

}
/*
initialisePopulation();
initialiseExamplePopulation();
updateProgLen(NULL);
evaluatePopulation(NULL, 0);
printPopulation();*/

fitnessTrace = fopen(OUT_FILE, "w");

for (int run = 0; run < MAX_RUNS; ++run) {

if (success == 1) break;

printf("\n\nRun %d\n\n", run);

initialisePopulation();
//initialisePartialSolution();
//updateProgLen(NULL);
float oldFitness = (1 BETTER_THAN 0)? INT_MAX: INT_MIN;

//evaluate the initial population (generation 0)
float fitness = evaluatePopulation(NULL, 0);

```

```

for(int generation = 1; generation < NUMGENERATIONS; ++generation){

if(success==1) break;

if(generation % OUTPUT.INTERVAL ==0)printf("\nGeneration %d",generation);

int randomNodeIndex = randRange(NUM_TERMINALS, NUM_PRIMITIVES-1);
Node* randomNode = &population[randomNodeIndex];
int randomOperandIndex = randRange(0, primitiveTable[randomNode->primitiveIndex]);
int oldOperandValue = randomNode->operands[randomOperandIndex];
int randomOperandValue = randRange(0,randomNodeIndex-1);

successorMutate(randomNodeIndex, randomOperandIndex, randomOperandValue);

updateList[0] = 0;

buildUpdateList(randomNodeIndex);

updateProgLen(updateList);

oldFitness = fitness;

fitness = evaluatePopulation(updateList, generation % NUM.TEST.SETS);

if(!(fitness BETTER_THAN oldFitness)){

restoreFitnessValues(updateList);

fitness = oldFitness;

removePredecessor(randomOperandValue, randomNodeIndex);

randomNode->operands[randomOperandIndex] = oldOperandValue;

addPredecessor(oldOperandValue, randomNodeIndex);

}

fprintf(fitnessTrace,"%f",fitness);
if(generation % OUTPUT.INTERVAL ==0)printf(" Fitness: %f\n", fitness);

}

fprintf(fitnessTrace,"\n");
}
printPopulation();

```

```

fclose(fitnessTrace);

if(success == 1){
printf("\n\nSuccess!\n\n");
printNode(workingProgramme);
} else{
printf("\n\nNo working programme\n\n");
}

return 0;
}

```

Code listing for SNGP\_Sort\_Debug.c

```

#ifndef SNGP_SORT_DEBUG_C
#define SNGP_SORT_DEBUG_C

#include<stdio.h>

void printIntArray(int* a, int n){
for(int i = 0; i<n; i++){
printf("%d ", a[i]);
}
}

void printPrimitiveTable(){
for(int i = 0; i < NUM_PRIMITIVES; i++){
TableEntry entry = primitiveTable[i];
printf("i: %d, e: %d, a: %d n: %s\n", i, entry.primitive, entry.arity, entry.name);
}

printf("\n");
}

```

```

void printTestData(){
    for(int i = 0; i < NUMGENERATIONS; i++){
        printf("Set %d\n\n", i);
        for(int j = 0; j < NUMTESTS; j++){
            Array* test = tests[i][j];
            printf("Test: %d-%d Size: %d Inversions: %d\n", i, j, test->size, te
            for(int k = 0; k < test->size; k++){
                printf("%d ", test->arr[k]);
            }
            printf("\n");
        }
        printf("\n");
    }
    printf("\n");
}

void printNode(int popIndex){
    Node node = population[popIndex];

    printf(
        "Index: %d primitive: %s Arity: %d\nFitness: %f OldFitness: %f ProgLen
        popIndex,
        primitiveTable[node.primitive].name,
        primitiveTable[node.primitive].arity,
        node.fitness,
        node.oldFitness,
        node.progLen
    );

    for(int j = 0; j < primitiveTable[node.primitive].arity; j++){

```

```

printf("%d ", node.operands[j]);
}

printf("\nPredecessors: ");

/*for(int j = 0; j < POPULATION_SIZE; ++j){
if(node.predecessors[j]!=0) printf("%d ", node.predecessors[j]);
}*/

for(int i = node.predecessors[0]; i != 0; i = node.predecessors[i]){

printf("%d ", i);
}
}

void printPopulation(){

for(int popIndex = 0; popIndex < POPULATION_SIZE; popIndex++){
printNode(popIndex);

printf("\n\n");
}

printf("\n");
}

void initialiseExamplePopulation(){

population[0].primitive = INDEX;

population[1].primitive = LENGTH;

population[2].primitive = INC;
population[2].operands[0] = 0;

population[3].primitive = SMALLEST;
population[3].operands[0] = 2;
population[3].operands[1] = 0;

```

```

population[4].primitive = SWAP;
population[4].operands[0] = 3;
population[4].operands[1] = 0;

population[5].primitive = SUB;
population[5].operands[0] = 1;
population[5].operands[1] = 1;

population[6].primitive = DEC;
population[6].operands[0] = 1;

population[7].primitive = ITERATE;
population[7].operands[0] = 5;
population[7].operands[1] = 6;
population[7].operands[2] = 4;

population[8].primitive = ITERATE;
population[8].operands[0] = 5;
population[8].operands[1] = 6;
population[8].operands[2] = 7;

population[9].primitive = SWAP;
population[9].operands[0] = 6;
population[9].operands[1] = 0;

}

void initialisePartialSolution(){

population[12].primitive = INC;
population[12].operands[0] = 0;

population[20].primitive = SMALLEST;
population[20].operands[0] = 2;
population[20].operands[1] = 0;

population[25].primitive = SWAP;
population[25].operands[0] = 3;
population[25].operands[1] = 0;

population[33].primitive = SUB;
population[33].operands[0] = 1;
population[33].operands[1] = 1;

population[40].primitive = DEC;

```



```

population[40].operands[0] = 1;

population[42].primitive = ITERATE;
population[42].operands[0] = 33;
population[42].operands[1] = 40;
population[42].operands[2] = 25;

}

void testExecution(){

initialiseExamplePopulation();
printPopulation();

results = malloc(arrayMem(maxTestSize));

printf("Test data initialised\n");

int testResult = testNode(8,2,1);

for(int i = 0; i < results->size; i++){

printf("%d ", results->arr[i]);

}

}

#endif

```

Code listing for testGenerator.py

```

from random import randint
from sys import argv

#length20, inversions = 86
testArr1 = [82,160,286,169,65,275,35,376,288,87,398,57,66,80,41,195,

#length = 17, inversions = 82

testArr2 = [206, 344, 377, 270, 117, 14, 304, 125, 178, 134, 222, 18

minIntVal = 0
maxIntVal = 400

minArrSize = 2
maxArrSize = 30

```

```

numTests = 30
maxOps = 30000

def countInversions(arr):

    if len(arr) <= 1:
        return 0

    mid = len(arr)//2
    a = arr[:mid]
    b = arr[mid:]

    inversions = countInversions(a) + countInversions(b)

    aCounter = 0
    bCounter = 0
    cCounter = 0

    while aCounter < len(a) and bCounter < len(b):
        if a[aCounter] <= b[bCounter]:
            arr[cCounter] = a[aCounter]
            aCounter+=1
        else:
            arr[cCounter] = b[bCounter]
            bCounter+=1
            inversions+=(len(a)-aCounter)

        cCounter+=1

    if aCounter < len(a):
        for i in range(aCounter, len(a)):
            arr[cCounter] = a[i]
            cCounter+=1
    else:
        for i in range(bCounter, len(b)):
            arr[cCounter] = b[i]
            cCounter+=1

    return inversions

def randArr():
    arr = []
    for i in range(randint(minArrSize, maxArrSize+1)):
        arr.append(randint(minIntVal, maxIntVal+1))
    return arr

```

```

def main():
    #print(countInversions(testArr2))
    file = open(argv[1], "w+")
    outString = ""
    for testSet in range(maxOps):
        for test in range(numTests):
            arr = randArr()
            outString += str(len(arr)) + " " + str(countInversions(arr.copy())) +
            outString += "\n"
            file.write(outString)

if __name__ == "__main__":
    main()

```